J. Jake Nichol
Candidate
Computer Science
Department
This dissertation is approved, and it is acceptable in quality and form for publication: Approved by the Dissertation Committee:
Melanie E. Moses, Chair
G. Matthew Fricke, Member
Abdullah Mueen, Member
Tobias P. Fischer, Member
Laura P. Swiler, Member

Seeking Structure in Complex Systems: From Feature Analysis to Space-Time Causal Discovery with Earth Science Applications

BY

J. Jake Nichol

Bachelor of Science, Computer Science, University of New Mexico, 2016

Master of Business Administration, University of New Mexico, 2017

DISSERTATION

Submitted in Partial Fulfillment of the Requirements for the Degree of

Doctor of Philosophy

Computer Science

The University of New Mexico
Albuquerque, New Mexico

July, 2025

DEDICATION

To my wife, Whitney, for your unwavering love, support, and stubborn encouragement. You bring
peace to my restless mind.
To my parents, Carolyn and Jeff, for setting me on the path to become who I am. I learned more
from you than you know.
To my brother, Frank, for being my first partner and confidant. I learned more from you than you
should know.
"Whatsoever is contrary to nature is contrary to reason, and whatsoever is contrary to reason is
absurd."
—Baruch Spinoza

ACKNOWLEDGMENTS

I want to thank ...

...I could not have achieved this without them.

Seeking Structure in Complex Systems: From Feature Analysis to Space-Time Causal Discovery with Earth Science Applications

by

J. Jake Nichol

Bachelor of Science, Computer Science, University of New Mexico, 2016Master of Business Administration, University of New Mexico, 2017Doctor of Philosophy, Computer Science, University of New Mexico, 2025

ABSTRACT

Complex systems are difficult to study because of their many interacting parts, emergent phenomena, and feedback loops. These systems underpin all life on Earth. We need improved tools for seeking an understanding of them. This body of research presents my investigations into data-driven methods for understanding complex systems, including my invention of a novel causal discovery meta-algorithm for space-time gridded data. I demonstrated machine learning feature importance and causal discovery capabilities for comparing simulated and observed climate data. I developed a new benchmark for modeling space-time dynamics of locally driven phenomena and examined a prominent causal discovery algorithm. Finding that contemporary causal discovery struggles with the high-dimensionality of space-time gridded data, I developed Causal Space-Time Stencil Learning (CaStLe), a causal discovery meta-algorithm for recovering the space-time evolution of advective phenomena. Finally, I extended CaStLe to recover multivariate space-time dynamics. This research enhances scientists' capabilities to explore and understand complex systems in our universe.

Contents

1 Introduction				
	1.1	The Pursuit of Causal Discovery	3	
	1.2	Statistical Learning	6	
		1.2.1 Explainability in Machine Learning	7	
		1.2.2 Bayesian Networks	9	
	1.3	Causal Network Learning	11	
		1.3.1 Definitions, Notations, and Key Causal Assumptions	11	
		1.3.2 Consistency	16	
		1.3.3 Validation and Falsifiability	17	
		1.3.4 Time Series Causal Discovery	19	
	1.4	Earth Science Challenges	20	
		1.4.1 Earth Science Data	22	
2	Rela	ated Work	24	
	2.1	Causal Discovery	25	
		2.1.1 Causal Network Learning	27	
		2.1.2 Structural Causal Models	30	

	2.2	Attribution in Climate Science	31
	2.3	Causal Discovery for Earth Systems Science	32
		2.3.1 Specific Application Challenges	34
		2.3.2 Recent Efforts to Overcome Application Challenges	38
	2.4	Applications of Causal Network Discovery for Climate Science	40
]	Part I Foundations of Structure Learning for Earth Systems	46
3	Ma	chine Learning Feature Analysis Illuminates Disparity Between	
	E35	SM Climate Models and Observed Climate Change	47
	3.1	Publication Notes	47
	3.2	Abstract	48
	3.3	Introduction	48
	3.4	Related Work	51
	3.5	Data and Methods	52
		3.5.1 Data	53
		3.5.2 Random Forests	56
		3.5.3 Pre-Processing	58
		3.5.4 Model Training and Hyper-Parameter Tuning	59
		3.5.5 Feature Importance Measurement	60
		3.5.6 Model Evaluation	61
	3.6	Results	62
	3 7	Discussion	65

	3.8 Conclusions	68
4	Learning Why: Data-Driven Causal Evaluations of Climate Models	7 1
	4.1 Publication Notes	71
	4.2 Abstract	72
	4.3 Introduction	72
	4.4 Data	76
	4.5 Approach	78
	4.5.1 The PCMCI method	80
	4.5.2 Comparing and evaluating causal models	80
	4.6 Anticipated Contributions	81
5	Causal Evaluations for Identifying Differences between Observations	
	and Earth System Models	83
	5.1 Publication Notes	83
	Part II Local Causal Discovery in High-Dimensional Gridded Data	99
6	Benchmarking the PCMCI Causal Discovery Algorithm for Spatiotem-	
	poral Systems	100
	6.1 Publication Notes	100
7	Space-Time Causal Discovery in Earth System Science: A Local Sten-	
	cil Learning Approach	157

7.1	Publication Notes					
7.2	Abstract					
7.3	Introduction					
7.4	Background: Causal Discovery and Formal Mathematical Scope	167				
	7.4.1 Related Work: Causal Structure Learning	169				
	7.4.2 PDE-Like Systems	176				
	7.4.3 Causal Discovery of Physical Dynamics: Dynamical Con-					
	straints	178				
7.5	Data: The 1991 Mt. Pinatubo Eruption	181				
	7.5.1 Held-Suarez-Williamson-Volcanic	183				
	7.5.2 Mt. Pinatubo in E3SMv2-SPA	185				
7.6	Methodology: Causal Discovery with CaStLe	185				
	7.6.1 Notation	185				
	7.6.2 Causal Space-Time Stencil Learning	186				
	7.6.3 The CaStLe Meta-Algorithm	187				
	7.6.4 Theoretical Properties	191				
	7.6.5 Methodological Limitations	193				
	7.6.6 Strategies for Addressing Limitations	194				
7.7	Results: Discovering Atmospheric Dynamics in Global Climate Mod-					
	els	196				
	7.7.1 Validation with HSW-V	197				
	7.7.2 Extending to More Complexity: E3SMv2-SPA Modeled Aeroso	ds20′				

7.8	Validation and Benchmarking					
	7.8.1	Evaluating CaStLe: A Comparative Analysis	205			
7.9	Discu	ssion	209			
A	Unde	rstanding Assumptions	221			
	A. 1	CaStLe Assumptions	221			
	A.2	Causal Discovery Assumptions	222			
	A.3	Relationship Between Assumption Sets	223			
	A.4	Potential Violations and Their Manifestations	224			
В	Statis	tical and Time Complexity	226			
	B.1	Time Complexity	227			
	B.2	Statistical Consistency	230			
C	Asym	nptotic Consistency	231			
D	Application to Non-Linear Dynamics: Continuous Systems via Burg-					
	ers' Equation					
	D.1	Burgers' Equation: Model and Parameters	238			
	D.2	Advection Angle Estimation	239			
	D.3	Experimental Setup	240			
E	Propo	osed Modification of Statistical Methods for CaStLed Data	242			
F	Limit	rations of Dimensionality Reduction for Space-Time Causal Dis-				
	cover	y	243			
G	Additional experimental details for Section 7.7					
Н	Analy	vsis of Spatial Blocking	252			

	I	Analy	sis of Assumption Violation Examples	256
		I.1	Time Resolution is Too Coarse (Assumption T1)	256
		I.2	Time Interval is Too Long (Assumption T2)	256
		I.3	Grid Resolution is Too Coarse (Assumption S1)	259
		I.4	Block Sizes are Too Large (Assumption S2)	260
	J	Addit	tional GCM Results	263
	K	Addit	tional VAR Results	264
	L	PC-St	table-Single	267
8	M-(CaStL	e: Uncovering Local Causal Structures in Multivariate Space	-
	Tim	ne Grio	dded Data	270
	8.1	Public	cation Notes	270
	8.2	Abstr	ract	270
	8.3	Introd	duction	272
		8.3.1	Background and Motivation	275
		8.3.2	Foundations of the CaStLe Framework	276
		8.3.3	Theoretical Properties and Empirical Validation of CaStLe	280
		8.3.4	Research Gap and Motivation for Multivariate Extension	283
		8.3.5	Contributions	285
		8.3.6	Paper Organization	286
	8.4	Metho	ods	286
		8.4.1	Phase 1: The Locally Encoded Neighborhood Structure (LENS))287
		8.4.2	Phase 2: The Parent-Identification Phase (PIP)	288

		8.4.3	Interpretability: Decomposing the Multivariate Stencil	289
	8.5	Bench	nmarking Multivariate Causal Space-Time Stencil Learning (M-	
		CaStL	Le) with vector autoregression models (VARs)	290
		8.5.1	Background: Univariate Space-Time vector autoregression mod-	-
			els (VARs)	290
		8.5.2	Multivariate Space-Time vector autoregression models (VARs)	292
	8.6	Resul	ts	293
		8.6.1	Metrics	294
		8.6.2	Data Generation	295
		8.6.3	Multivariate Performance	296
		8.6.4	Comparison to the PC Algorithm	297
		8.6.5	Exploring Recall	298
	8.7	Discu	ssion	300
	A	Comp	oleted Data Generation Parameters	305
	В	Addit	ional vector autoregression model (VAR) Results	307
		B.1	PC Comparison Results	307
9	Cor	ıclusio	n	309
			Synthesis of Foundations Work	
			Machine Learning Feature Importance for Climate Models	
			Causal Discovery for Climate Model Evaluation	
	0.2		·	
	9.2		I: Discovery of Local Dynamics	
		9.2.1	Grid-Level Benchmarking of PCMCI	314

	9.2.2 CaStLe: Grid-Level Causal Discovery	315
	9.2.3 M-CaStLe: Multivariate Grid-Level Causal Discovery	319
9.3	Connections and Research Frontiers	320

List of Figures

3.1	Comparison of observed, pan-Arctic mean September sea ice ex-	
	tent with predictions from Energy Exascale Earth System Model	
	(E3SM)'s historical ensembles 1-5. The mean of Energy Exascale	
	Earth System Model (E3SM) simulations is shown with 95% confi-	
	dence interval (shaded)	56
3.2	June feature importance. Standard box-and-whisker plot (McGill	
	et al., 1978) of values for 13 predictions generated by 385 models.	
	The average R^2 , anomaly correlation coefficient (ACC), and mean	
	absolute error (MAE) are displayed in the gray boxes. The blue line	
	in each dataset is the mean importance of a random variable in each	
	feature set	62
3.3	July feature importance. Standard box-and-whisker plot (McGill	
	et al., 1978) of values for 13 predictions generated by 385 models.	
	The average R^2 , anomaly correlation coefficient (ACC), and mean	
	absolute error (MAE) are displayed in the gray boxes. The blue line	
	in each dataset is the mean importance of a random variable in each	
	feature set	63

3.4	August feature importance. Standard box-and-whisker plot (McGill	
	et al., 1978) of values for 13 predictions generated by 385 models.	
	The average R^2 , anomaly correlation coefficient (ACC), and mean	
	absolute error (MAE) are displayed in the gray boxes. The blue line	
	in each dataset is the mean importance of a random variable in each	
	feature set	65
4.1	Comparison of observed, pan-Arctic mean September sea ice ex-	
	tent with predictions from Energy Exascale Earth System Model	
	(E3SM)'s historical ensembles 1-5. The mean of E3SM simulations	
	is shown with 95% confidence interval (shaded)	77
4.2	Diagram showing correlated relationships between variables in June	
	from the observed dataset between 1979 to 2014. Green indicates a	
	positive correlation and orange indicates a negative correlation. The	
	correlation threshold is ± 0.6	79

7.1	Schematic overview of the key elements of CaStLe and the process
	followed in its application to Mount Pinatubo's eruption of strato-
	spheric aerosols. Beginning with Earth system model output, Step
	1. is to collect stratospheric wind and aerosol data. Step 2. is to
	apply our novel CaStLe meta-algorithm to the aerosol data to obtain
	a causal graph describing the space-time evolution of the aerosols.
	Finally, we use the wind fields to help validate the causal graph re-
	sults in Step 3

1.2	Illustration of CaStLe (Algorithm 1) as applied to space-time data	
	on a 4×4 grid. Step A (§7.6.3): for every interior grid cell, its 3×3	
	(Moore) neighborhood is selected. (Note, all four 4×4 grids in the	
	second panel are identical.) Step B (§7.6.3): Data are represented in	
	a reduced coordinate space obtained by appending time series from	
	each neighborhood according to its position relative to the neighbor-	
	hood's center. Step C (§7.6.3): during the Parent Identification Phase	
	(PIP), a causal discovery algorithm is used to estimate the parents of	
	the center time series; the resulting graph forms the causal stencil.	
	Step D (§7.6.3): the estimated stencil is expanded to its equivalent	
	representation in the original space. Note that each time chunk (col-	
	ored intervals in the center panel) in the reduced space corresponds	
	to an interior grid cell of the original data, and that each edge in the	
	final causal graph reflects to a stencil edge learned during the PIP.	
	See §7.6.3 for details	187

7.3 Application of CaStLe-PC-Stable to HSW-V simulation of the 1991 Mt. Pinatubo eruption. The stencils estimated by CaStLe (white) capture the underlying high-altitude wind fields (green) using only satellite-measured AOD, with near perfect accuracy in high aerosol regions (red-orange). Autodependencies are shown with black nodes where grid cells cause themselves, and gray nodes where there is no autodependence. All links represent a six hour time lag, the time resolution of the HSW-V dataset. On longer horizons (bottom row), CaStLe is able to recover equatorial wind currents as far away as South America, half-way around the world from Mt. Pinatubo (white triangle). CaStLe accurately identifies the prevailing westerly atmospheric winds because it was able to identify the spacetime dependence between neighboring grid cells. Additional details 215 7.5 Application of CaStLe-PC-Stable to E3SMv2-SPA simulation of the 1991 Mt. Pinatubo eruption. The stencils estimated by CaStLe (white) capture the underlying high-altitude wind fields (green) using only total aerosol optical depth (AOD). Autodependencies are shown with black nodes where grid cells cause themselves, and gray nodes where there is no autodependence. All links represent a one day time lag, the time resolution of the E3SMv2-SPA dataset. The heatmap depicts AOD from any source at 50 hPa. The top panel depicts learning from the first 20 days after eruption, which began on day 15. The bottom panel depicts learning approx 6 months after the eruption over a 20-day time period. In the more challenging setting of the fully-coupled E3SMv2-SPA model, our results in the first weeks are still generally consistent with those in HSW-V presented in Section 7.7.1, showing that CaStLe is largely robust to greater complexity. In the bottom panel, the aerosols and winds are in a different regime. CaStLe stencils are still consistent in the tropics and now begin to recover dynamics pushing aerosols northwards above central Asia and southwards through western North America. A more complex model and smaller block sizes illustrate more nuanced dynamics, and there is more to learn from these, however, we leave deeper atmospheric dynamics analysis to future work. 217 7.6 Comparison of CaStLed and non-CaStLed causal discovery approaches on linear-Gaussian dynamics, including Granger causality or FullCI (orange), PC (green), PCMCI (red), and DYNOTEARS (purple), as well as a statistical model of the data generating process (blue) presented with both MCC and F_1 metrics. In the low-sample size regime (T=10, left) CaStLed approaches can accurately recover the underlying causal graph, with performance increasing on larger grid sizes (solid lines); by contrast, non-CaStLed approaches are unable to perform better than mere chance (dashed lines). Even a model based on the underlying data generating process (Sparse VAR, blue) is significantly outperformed by its CaStLed counterpart. In the high-sample size regime (T=150, right), non-CaStLe approaches have improved performance but still compare unfavorably with their CaStLed coun-218 D1 Application of CaStLe-PC to advection estimation from non-linear PDE dynamics. In the left panel, the first three columns depict realizations of Burgers' equation under different advection-to-diffusion regimes; the fourth column depicts the causal stencil identified by CaStLe-PC; and the final column compares the estimated advection angle with the true advection angle. The right panel depicts the accuracy of CaStLe-PC under various signal-to-noise conditions. Each combination of advection and diffusion rates were tested with 500 angles sampled uniformly from $[0^{\circ}, 360^{\circ})$. In low-diffusion (high SNR) scenarios, CaStLe-PC can identify the underlying advection clearly (top row of left panel and yellow-green columns in right panel). By contrast, in low-advection (low SNR) scenarios, CaStLe-PC struggles to accurately identify the underlying advective dynamics (bottom row of left panel and blue bars in right panel). Even in highly diffusive scenarios, CaStLe-PC is able to accurately estimate the underlying advection when it is sufficiently strong (around $M/c \ge 20$) as shown in the middle row of the left panel. Additional

F1	PCA study of Burgers' equation solution ($\theta = 45^{\circ}$, $M = 6$, $c =$	
	0.05). Four empirical orthogonal functions (EOFs) capture \approx 91% of	
	variance, with spatial patterns (left) and temporal evolution (right).	
	The bottom panels show explained variance distribution and PCMCI	
	causal graph, which fails to accurately represent the known direc-	
	tional advection process in the underlying PDE, highlighting limita-	
	tions of this approach for local causal structures in space-time systems	.246
F2	PCA-Varimax study of Burgers' equation solution ($\theta = 45^{\circ}$, $M =$	
	6, $c = 0.05$). Four empirical orthogonal functions (EOFs) capture	
	\approx 91% of variance, with spatial patterns (left) and temporal evolution	
	(right). The bottom panels show explained variance distribution and	
	PCMCI causal graph, which fails to accurately represent the known	
	directional advection process in the underlying PDE, highlighting	
	limitations of this approach for local causal structures in space-time	
	avetame.	247

F3	PCA study of the HSW-V dataset, in the time interval 21 days post-	
	eruption. Four empirical orthogonal functions (EOFs) capture $\approx 85\%$	
	of variance, with spatial patterns (left) and temporal evolution (right).	
	The bottom panels show explained variance distribution and PCMCI	
	causal graph, which fails to accurately represent the known direc-	
	tional advection process in the underlying system, highlighting limi-	
	tations of this approach for local causal structures in space-time sys-	
	tems	248
F4	PCA-Varimax study of the HSW-V dataset, in the time interval 21	
	days post-eruption. Four empirical orthogonal functions (EOFs)	
	capture \approx 85% of variance, with spatial patterns (left) and tempo-	
	ral evolution (right). Since varimax rotation does not preserve the	
	explained variance ordering, we reordered EOFs according to the	
	identified centroid's longitude. The bottom panels show explained	
	variance distribution and PCMCI causal graph, which fails to accu-	
	rately represent the known directional advection process in the un-	
	derlying system, highlighting limitations of this approach for local	
	causal structures in space-time systems	249

FS	PCA study of the E3SMv2-SPA dataset, in the time interval of days	
	15-35. Nine empirical orthogonal functions (EOFs) capture \approx 87%	
	of variance, with spatial patterns (left) and temporal evolution (right).	
	The bottom panels show explained variance distribution and PCMCI	
	causal graph, which fails to accurately represent the known direc-	
	tional advection process in the underlying system, highlighting limi-	
	tations of this approach for local causal structures in space-time sys-	
	tems	250
F6	PCA-Varimax study of the E3SMv2-SPA dataset, in the time interval	
	of days 15-35. Nine empirical orthogonal functions (EOFs) capture	
	\approx 87% of variance, with spatial patterns (left) and temporal evolution	
	(right). Since varimax rotation does not preserve the explained vari-	
	ance ordering, we reordered EOFs according to the identified cen-	
	troid's longitude. The bottom panels show explained variance distri-	
	bution and PCMCI causal graph, which fails to accurately represent	
	the known directional advection process in the underlying system,	
	highlighting limitations of this approach for local causal structures	
	in space-time systems	251

H1	Results of CaStLe applied to HSW-V 21 days after the Mt. Pinatubo	
	eruption with three different block sizes, $12^{\circ} \times 12^{\circ}$, $20^{\circ} \times 20^{\circ}$, and	
	$60^{\circ} \times 60^{\circ}$. We find that results are generally consistent over the same	
	area for each block size, with smaller block sizes allowing for addi-	
	tional nuance in some areas. Note that the $20^{\circ}\times20^{\circ}$ block panel is	
	similar to the results shown in Figure 3, but more longitudes were	
	added to get a space factorable by more integers, such as 12, 20, and	
	60	254
H2	The PC algorithm and CaStLe applied to E3SMv2-SPA in the 15 $^{\circ}$ \times	
	15° block between 15.00° to $30.00^{\circ}N$ and 75° to $90^{\circ}E$. from the	
	day of the eruption to 20 days later. PC struggles to estimate an	
	interpretable and physically meaningful graph of the dependence	
	structure in this area. In contrast, CaStLe is able to identify an in-	
	terpretable dependence structure that represents the local dynamics	
	within the space.	255
I1	Results of using a coarsened temporal resolution (two-daily) in the	
	E3SMv2-SPA study. CaStLe finds many fewer links in this setting.	
	It is clear that when time is too coarse, causal structures fail to be	
	detected. However, the remaining links that are found are largely	
	true positives, suggesting that CaStLe is relatively robust to coarser	
	time sampling	257

12	Results of applying CaStLe to a longer time interval from day 15 to
	65. CaStLe identifies more links, indicating it is learning too many
	causal structures in the data, but still finds many of the true positives
	we found in our initial study. This indicates that many of the blocks
	in this interval have temporal causal stationarity, leading CaStLe to
	perform adequately
I3	Results of applying CaStLe to a time interval that is too long and
	contains too many causal structures, day 15 to 200. We see that
	CaStLe identifies many links in each block. Comparing them to
	the winds is ineffective because the wind arrows are averages over
	the whole period rather than reflections of how they change in time,
	which CaStLe is learning from. With such a density of links, it is
	further challenging to know which are correct and which are spurious. 258
I4	Results of using a coarse grid (9°) in the E3SMv2-SPA study. We
	find that CaStLe performs very well overall. There are few false
	positives and it clearly captures the overall advection dynamics of
	the system

I5	Results of using a coarse grid (18°) in the E3SMv2-SPA study. CaS-	
	tLe performs well in the early time interval, clearly identifying the	
	east-to-west advection pattern. However, in the later time interval,	
	it finds no spatial structures apart from autodependencies in each	
	block. This is likely because the east-to-west advection is weaker in	
	this period and the grid is too coarse to capture the narrower bands	
	of northward advection that dominates the interval	261
I6	Results of using block sizes too large in the E3SMv2-SPA study. We	
	see that many true positives are found, but many false positives as	
	well. CaStLe seems to identify multiple contradictory causal struc-	
	tures within many cells, which may lead to more spurious links dis-	
	covered. Even where links appear correct, they are largely uninter-	
	pretable in the presence of contradictions	262

J1	Application of CaStLe-DYNOTEARS to HSW-V simulation of the	
	1991 Mt. Pinatubo eruption. The stencils estimated by CaStLe	
	(white) capture the underlying high-altitude wind fields (green) us-	
	ing only satellite-measured AOD, with near perfect accuracy in high	
	aerosol regions (red-orange). On longer horizons (bottom row), CaS-	
	tLe is able to recover equatorial wind currents as far away as South	
	America, half-way around the world from Mt. Pinatubo (white tri-	
	angle). CaStLe accurately identifies the prevailing westerly atmo-	
	spheric winds because it was able to identify the space-time depen-	
	dence between neighboring grid cells	263
K 1	Matthews correlation coefficient (MCC) comparison between CaS-	
	tLed and non-CaStLed causal discovery approaches on 2D VAR dy-	
	namics for each sparsity level, including Granger causality (orange),	
	PC (green), PC-Stable-Single (cyan), PCMCI (red), DYNOTEARS	
	(purple), and a statistical model of the data generating process (blue).	
	See Section 7.8.1 for experimental details	265
K2	F ₁ score comparison between CaStLed and non-CaStLed causal dis-	
	covery approaches on 2D VAR dynamics for each sparsity level, in-	
	cluding Granger causality (orange), PC (green), PC-Stable-Single	
	(cyan), PCMCI (red), DYNOTEARS (purple), and a statistical model	
	of the data generating process (blue). See Section 7.8.1 for experi-	
	mental details.	266

8.1	A conceptual diagram of the Locally Encoded Neighborhood Struc-	
	ture (LENS) that CaStLe constructs for learning underlying local	
	causal dynamics in gridded data. This encoding transforms the orig-	
	inal grid space into a local neighborhood structure without marginal-	
	ization, preserving all of the local relationships in the gridded time	
	series data.	280
8.2	A demonstration of the full CaStLe process to produce a causal sten-	
	cil graph on an example input 4×4 gridded space-time system. In	
	the Locally Encoded Neighborhood Structure (LENS) phase, neigh-	
	borhood information is collected from each of the interior grid cells,	
	which are then concatenated to form the Locally Encoded Neighbor-	
	hood Structure (LENS). Finally, the PIP phase applies an adapted	
	time series causal discovery algorithm to learn the space-time par-	
	ents of the center node. The learned stencil depicts the underlying	
	space-time structure of each grid cell in the original data	281

8.3	A schematic diagram of the input, computational phases, and output	
	of Multivariate Causal Space-Time Stencil Learning (M-CaStLe).	
	Similar to CaStLe's procedure (c.f. Figure 8.2), the first phase col-	
	lects local neighborhood information into the Locally Encoded Neigh-	
	borhood Structure (LENS), which now collects information for each	
	variable's time series in each grid cell. The second phase applies the	
	Parent-Identification Phase (PIP) to every variable at every position	
	in the Locally Encoded Neighborhood Structure (LENS) to deter-	
	mine which variables cause the center variables from each location	
	in the Locally Encoded Neighborhood Structure (LENS). Finally,	
	the resulting multivariate stencil graph can be decomposed into the	
	spatial graph and reaction graph for improved interpretability and	
	potential analysis	285
8.4	Showing precision and recall alongside predicted positive rate, a	
	measure of how often a positive is predicted among all other predic-	
	tions. As variables increase, the predicted positive rate decreases,	
	which diminishes recell	204

8.5	A comparison between Multivariate Causal Space-Time Stencil Learn-	
	ing (M-CaStLe)-PC and PC considering the F_1 score for $V=4$ as	
	the number of links increases on a 4×4 grid. Multivariate Causal	
	Space-Time Stencil Learning (M-CaStLe)-PC outperforms PC in ev-	
	ery case because PC struggles with the very high dimensionality of	
	the systems since it is naive to the spatial and variable structures	298
8.6	In simple chains of multivariate stencils, even with an extremely	
	large number of variables, recall can be captured perfectly if the sig-	
	nal strength is large enough.	299
A7	Parameter ranges used in our experimental design, showing the link	
	count distribution for each grid size and variable count combination.	
	Each horizontal line represents the span of network links tested, with	
	each parameter combination having at least 30 replicate experiments	
	(n values shown). Our experiments covered grid sizes from 4×4 to	
	10×10 and 1-6 variables per grid. All experiments used 1000 time	
	samples and coefficient values between 0.1 and 1.0. The network	
	density, d , defined as the ratio of actual links, L , to maximum possi-	
	ble links $d = \frac{L}{(3 \times 3 \times V^2)}$, where $d \in (0, \dots 0.5]$. Not all density values	
	produced 30 stable systems within our computational constraints,	
	particularly at higher densities. This visualization shows which pa-	
	rameter combinations successfully generated sufficient replicates for	
	statistical analysis.	305

A8	The relationship between link coefficients and the number of links	
	present. As the number of links increases, maximum (blue) and min-	
	imum (green) link coefficients show a clear decreasing trend, with	
	their distribution becoming narrower and centered around lower val-	
	ues. This reveals that networks with more links have weaker signals,	
	suggesting that highly interconnected systems cannot be stable with	
	large dependencies	06
В9	Comparisons between Multivariate Causal Space-Time Stencil Learn-	
	ing (M-CaStLe)-PC and PC considering the F ₁ score, precision, and	
	recall for all V as the number of links increases on a 4×4 grid. Mul-	
	tivariate Causal Space-Time Stencil Learning (M-CaStLe)-PC out-	
	performs PC in every case because PC struggles with the very high	
	dimensionality of the systems since it is naive to the spatial and vari-	
	able structures	08

List of Tables

3.1	Training Features and June Data Excerpt: total cloud cover per-	
	centage (CLT), downward longwave flux at surface (FLWS), pres-	
	sure at the surface (PS), sea ice extent (SIE), sea ice volume (SIV),	
	near-surface specific humidity (SSH), sea surface temperature (SST),	
	temperature at the surface (TS), wind u component/zonal (uwind),	
	and wind v component/meridional (vwind). Values listed are means	
	over the pan-Arctic grid for each day of the month, rounded to two-	
	decimal places for display only.	55
1	Capabilities of CaStLe for Earth science applications. This table	
	summarizes the key methodological advantages of CaStLe and their	
	relevance to specific Earth science phenomena, highlighting appli-	
	cations where grid-level causal discovery enables analyses that were	
	previously infeasible with prior causal discovery approaches	220

1 Introduction

The principal function of science is to explore and explain our universe. To fulfill this charge, scientists seek to answer the questions of 'how?' and 'why?' In this pursuit, we strive to expand human knowledge, improve the well-being of all life, and develop practical applications that transform our world. Complex systems are fundamental to science because they represent the intricate reality of our world. By their nature, complex systems are difficult to study because of their many interacting parts, emergent phenomena, feedback loops, and tipping points. While many complex systems underpin life on Earth, our tools for studying them are limited.

This dissertation investigates the state of the art in data-driven *structure learning* methodologies for explaining and understanding complex systems, particularly for space-time Earth systems. As I use it in this work, structure learning is a class of methods that identify underlying dynamics, or structure, from data. In Part I, I outline the basics of the structure learning task and study how machine learning feature importance and causal discovery can be used to estimate structure in the Earth system.

Finding that the state-of-the-art primarily tackles global-scale emergent structures, Part II focuses on identifying local-scale structures in gridded datasets. This

work begins with benchmarking causal discovery algorithms for learning grid-level space-time dynamics. It corroborates that causal discovery algorithms struggle with datasets containing hundreds of thousands of grid cells, each with several orders of magnitude fewer observations in time. This imbalance is one aspect of the *curse of dimensionality* (Bellman, 1957; Bühlmann and Geer, 2011), where many variables relative to sample size limits conventional statistical methods and renders many forms of inference, including causal discovery, unreliable without dimensionality reduction.

To resolve that challenge, I developed a novel method, Causal Space-Time Stencil Learning (CaStLe), that significantly improves the performance and efficiency of causal discovery in local space-time dynamics. It does so via two parts: (i) the Locally Encoded Neighborhood Structure (LENS) reorganizes the given data such that the high-dimensionality of gridded data is eliminated and the sample complexity of the underlying grid-level structure is maximized; and (ii) the Parent-Identification Phase (PIP), which selectively applies causal discovery to minimize the search space while side-stepping spatial confounding. The initial implementation of CaStLe was univariate, in that it could only identify the space-time structure of a single quantity of interest. This work concludes with extending CaStLe to Multivariate Causal Space-Time Stencil Learning (M-CaStLe), which adapts the LENS and PIP to capture space-time structure between multiple quantities of interest.

1.1 The Pursuit of Causal Discovery

The scientific method provides consistent rigor to answer the 'how?' and 'why?' questions. With it, we design experiments, collect data on what we observe, and determine what we can learn from those data. Causal inference is the process of answering these questions and determining when such an answer is attainable. Pearl and Mackenzie (2018a) suggest that causal inference is conducted via three operations, which he calls the *Ladder of Causation*:

rung one: seeing (observing and collecting information)

rung two: doing (intervention and experimentation)

rung three: imagining alternatives (counterfactual analysis)

Causal discovery is an algorithmic methodology for finding causal hypotheses and eliminating spurious correlations in data, grounded in strict assumptions that represent domain expertise. Machine learning is typically classified as rung one, seeing; it produces observational distributions from which predictions of future states can be made. Causal graph discovery is rung two, doing; it produces interventional distributions in the form of causal models. These can be used to reason about the effects of intervention. Finally, structural causal models and digital twins are examples of rung three, because they enable one to reason about the implications of alternative scenarios. (Peters et al., 2017)

Statistical and machine learning are standard toolsets to quantify and predict

relationships when only observational (non-manipulated) data is available. Statistics can describe data and inform us of the underlying distribution, but it generally defers further inference (Pearl and Mackenzie, 2018a). Correlated relationships between variables are bidirectional and often ambiguous. Since correlation does not imply causation, one can only make stronger inferences with stronger assumptions.

Machine learning models capture patterns rather than learn to understand underlying mechanisms by computing statistics and fitting functions that separate data. The algorithms learn functions that map input to output, predicting a probabilistic distribution. Its primary goal is to model the given data to predict the classification or future values, i.e., regression. Machine learning has proven to be an informative and useful tool, but prediction is only correlation and, thus, also does not imply causation. The nascent field of explainable machine learning is bearing fruit in some domains. However, it is also limited to descriptive statistics and correlated information. Using explainable machine learning for elucidating the dynamics in a system may be a promising starting point towards finding causality when ground truth is nebulous. Later, in Chapter 3, I will discuss an analysis with random forest feature importance (Nichol et al., 2021).

The most reliable, though still imperfect, method of estimating causal relationships is with the randomized control trial (RCT) framework. In conducting an RCT, scientists make tacit assumptions called identifiability conditions: the causal Markov condition, ignorability/exchangeability, positivity/overlap, no in-

terference, and consistency. Ideal RCTs meet these assumptions by design; however, errors or biases, such as selection bias, may break identifiability. Hernán and Robins (2020) explain the remaining assumptions for causal inference in detail. I define the causal assumptions important for causal discovery in Section 1.3.1.

RCTs are a powerful tool, but they are not feasible in many cases, such as when randomizing treatment is unethical, impossible, or too expensive or inconvenient. One such example is the Earth science domain. In geophysics, many natural events are impossible to conduct ourselves, i.e., we cannot make an earthquake occur. In other fields, such as atmospheric science, we cannot ethically intervene randomly without fully understanding the downstream impacts of each intervention, e.g., stratospheric aerosol injection for solar radiation management. We have one Earth, and we cannot afford to disrupt it carelessly.

In some cases where RCTs are infeasible, we can conduct observational studies with frameworks like the target trial (Rubin, 1974; Robins, 1986; Dorn, 1953; Feinstein, 1970; Dawid, 2000). However, this relies upon enough sampling to measure a representative distribution of possible outcomes, posing another challenge for causal inference in Earth sciences: we can only observe one instance of the possible outcomes of the Earth's dynamics. One solution may lie in simulations, and numerical Earth system models (ESMs) are an ongoing research area. However, their complexity makes models imperfect, computationally expensive, and challenging to evaluate.

Founded on principles from path analysis (Wright, 1921), contemporary causal

discovery is developing into a rigorous mathematical framework, primarily due to work by Rubin (1974); Spirtes, Glymour, and Scheines (1993); Pearl (1995a); Peters, Janzing, and Schlkopf (2017). This framework can mathematically describe the causal questions asked, counterfactuals, interventions, relevant variables to measure, and potential answers to the causal questions. In the past two decades, algorithms have been designed to leverage this framework for reconstructing causal graphs or, interchangeably, causal networks. We can compute statistical relationships and make strict assumptions with observational data and the true underlying causal structure to reconstruct the causal structure that generated the observed data. These assumptions are also known as the identifiability conditions in causal inference. Algorithmically reconstructing causal graphs is called causal discovery, causal network learning, or causal learning.

1.2 Statistical Learning

Peters et al. (2017, p.46) write that "formally, learning causal models is substantially different from the [statistical] learning scenario because it aims at inferring a model that describes the behavior of the system under interventions and not just observations taken from the same distribution. Therefore, there is no straightforward way to adopt arguments from statistical learning theory, to obtain a learning theory for causal relations." Statistical machine learning generally aims to learn a function that fits given data, and we hope it can extrapolate from unseen data. Explainability tools, either derived directly from the model (e.g., decision trees and

random forest Gini importance) or many models trained on permuted data, fundamentally describe the models *alone*, rather than the true underlying dynamics in the data.

Tautologically, if the goal is to identify and learn about the dynamics in a system, then causality is fundamentally the only way to reason about those dynamics. As Pearl and Mackenzie (2018a) state, contemporary machine learning fundamentally cannot consider the causality in a system because it lacks a language for causality, i.e., counterfactuals and interventions. While we hope a trained model has learned some true underlying function in the data's generating process, it is causally unverifiable. Showing that a model consistently handles new data well increases the confidence that the model has generalized the true causal process, but the error in a model is merely a correlated observation; it does not verify causality.

1.2.1 Explainability in Machine Learning

The black-box nature of most machine learning models poses a big challenge for interpreting and validating their results. *Trustworthy machine learning* and *fairness in machine learning* efforts have turned to uncertainty quantification and explainability methods to validate further, and to understand how and why a particular model has been fit. Some machine learning methods have an inherent explainability, such as decision trees and random forests (Breiman, 2001; Nembrini et al., 2018). Because these models are built iteratively, Gini impurity, the probability of misclassifying an observation, is computed for every node split in the trees. Gini

impurities can be aggregated after learning to produce a Gini importance, or feature importance, for each feature. These importance values measure how much each feature contributed to reducing the model's error on average.

Other machine learning models must use ad hoc and post hoc methods to measure the importance of features for model learning. Examples include Shapley values (Lundberg and Lee, 2017), Locally Interpretable Model-Agnostic Explanations (LIME) (Krishnapuram et al., 2016), SHapley Additive exPlanations (SHAP) (Lundberg and Lee, 2017), and DeepLIFT (Shrikumar et al., 2017). Shapley values, LIME, and SHAP, are model agnostic methods, so they can be applied to support vector machines, random forests, neural networks, etc. DeepLIFT is a member of a class of methods specifically for neural networks. All of these function by measuring the contribution of each feature to the model or a specific prediction. They train many models and vary whether each feature is included by permuting each feature.

Explainability may illuminate causality with respect to the model, but it cannot illuminate causality within the studied system. That is evident because explainability methods make no assumptions about the system itself, nor the data observed. The methods and the models have no way of knowing whether the data and features are representative of the system. They are only aware of the models and data given. Because of this, inferences from these will always fail to rise above making purely associational observations of the given data.

In general, it is acceptable that explainability methods fail to elucidate causality

within a system because they make no claims beyond a rigorous attempt at explaining the given model. To these methods, the generating process is not what created the input data but the model itself. They fundamentally address a different question from causal discovery.

1.2.2 Bayesian Networks

Judea Pearl wrote in his book, *The Book of Why*, that he initially made the same mistake as many philosophers and economists, and that I would suggest is made by many in machine learning now: putting probability first and causality second (Pearl and Mackenzie, 2018a, p.50). He thought that uncertainty was the most important thing missing from artificial intelligence and insisted that uncertainty be represented by probabilities. With that in mind, he developed Bayesian networks to reason under uncertainty.

Bayesian networks encode conditional probabilities between events. Given that we observe certain probabilities of events, Bayesian networks can compute the likelihood of other events or whether certain facts are true or false. This computation is called *belief propagation*.

Pearl says that while Bayesian networks are still popular for reasoning under uncertainty, they fail to accomplish what he was after: identifying and quantifying causality. Bayesian networks fail to climb beyond rung one in his Ladder of Causation. He says, "Bayesian networks inhabit a world where all questions are reducible to probabilities, or degrees of association between variables..." (Pearl and

Mackenzie, 2018a, p.51). Pearl solved this problem after putting aside Bayesian networks to develop structural causal models (SCM) and the *Do* notation which provide a mathematical language for writing down *what we know* and *what we want to know*. His *Do-Calculus* (Pearl, 2012) enables us to compute counterfactual and interventional distributions from observational data, as opposed to probabilistic distributions alone.

Bayesian networks are quite similar to causal networks, however. Pearl (1995b) and Pearl and Mackenzie (2018a) write about how to transition from a Bayesian network to a causal network in. Bayesian networks' probabilistic and belief propagation properties are still valid in causal networks. The main difference is in how they are constructed. Bayesian networks are a graphical form of conditional probability tables.

A causal network changes the language of the relationships between nodes; the meaning of their construction and interpretation change. Rather than a relationship between nodes indicating that they probabilistically coincide, in a causal network, it indicates which node another node "'listens' to before choosing its value," (Pearl and Mackenzie, 2018a, p.129). The *listening* analogy describes the causal assumptions, i.e., the knowledge one has of the system. A missing link between nodes denotes that the two are independent in both Bayesian networks and causal networks. Though, in a causal network, a missing link may also indicate two nodes are indirectly independent. As Pearl notes, this implies that causal assumptions cannot be made-up and can be falsified against the observed data. Pearl's transition from

Bayesian networks to causal networks coincided with the work of Spirtes and Glymour's (1991) development of causal discovery, which is the reconstruction of the causal network from observational data.

1.3 Causal Network Learning

In this work, I will focus on conditional independence-based causal network learning¹ (Spirtes et al., 1993; Runge et al., 2019a) for reconstructing causal graphs. Time series adaptations are well suited for the stochastic, highly autocorrelated, and high-dimensional data in Earth science (Runge, 2018a; Runge et al., 2019a). Other forms of causal discovery include nonlinear state-space methods (Arnhold et al., 1999; Sugihara et al., 2012), and structural causal models (Spirtes and Zhang, 2016a; Peters et al., 2017).

1.3.1 Definitions, Notations, and Key Causal Assumptions

Causal Graphs

For a multivariate time series \mathbf{X} , X^i denotes the time series of the i^{th} variable, $X^i_{t-\tau}$ denotes the time series lagged by τ time steps, and $\mathbf{X}^-_t = (\mathbf{X}_{t-1}, \mathbf{X}_{t-2}, ...)$ are lagged time series of \mathbf{X} , representing its temporal parents.

A causal graph, or a causal network, is a directed acyclic graph (DAG) or partially-directed acyclic graph (PDAG) that encodes the causal structure between variables in a system. Using DAGs to represent causal relationships is credited to

¹Causal network learning is also known as "causal discovery," "causal graph discovery," and "structure learning," and I will use these terms interchangeably throughout this dissertation.

Pearl (1995a, 1998). A causal time series graph adapts the causal DAG to incorporate time lags. Each variable has a node for the original, present time t, and every time lag, $t - \tau$. This is theoretically an infinite graph, but in practice, we truncate the graph to a maximum time lag, τ_{max} .

A link between variables in a causal graph, G, marks a dependence between two variables. Variables $X_{t-\tau}^i$ and X_t^j are connected by a lag-specific directed link, $X_{t-\tau}^i \to X_t^j \in G$ for $\tau > 0$, if and only if

where $\not\perp$ denotes statistical dependence (\perp would denote independence). Thus, Equation 1.1 can be read as " X_t^j is dependent on $X_{t-\tau}^i$, conditional on $[X_t^-, excluding the set <math>\{X_{t-\tau}^i\}$]." Autodependencies are links where i=j. Links from X_t^i to X_t^j are called contemporaneous links. Some algorithms represent these with an undirected edge in the graph, others can use collider rules to possibly orient these (Runge, 2020; Spirtes et al., 1993).

The parents of a node, X_t^i , in G, are mathematically written as

$$\mathscr{P}(X_t^i) = \{X_{t-\tau}^k : X^k \in \mathbf{X}, \tau > 0, X_{t-\tau}^k \to X_t^i\}. \tag{1.2}$$

D-separation

Independence between nodes within a graph is called *d-separation*, for directed-separation, or sometimes just *separation*. It tells us where and when association

can *flow*, or be measured, between two nodes. If two nodes are not d-separated, then their data will be correlated. This is an important property for interpreting graphs, but with the assumptions detailed in the following section, we can infer the graph from measured dependencies in data.

To explain d-separation, we first need to explain how association flows between variables in a causal graph. The rules of d-separation operate on the three main components in a causal DAG: the *chain*, $(X \to Y \to Z \to ...)$ and $(... \leftarrow X \leftarrow Y \leftarrow Z)$; the *fork*, $(X \leftarrow Y \to Z)$; and the *collider* or *V-structure*, $(X \to Y \leftarrow Z)$. In chains and forks, association flows between all variables. For the chains/fork example above, $X \not\perp\!\!\!\!\!\perp Y \not\perp\!\!\!\!\!\perp Z$. Note that the two chains and the fork all have the same independence relationships. This set of independence relationships represent a *Markov equivalence class* of causal graphs.

Dependence is transitive, so we also have that $X \not\perp\!\!\!\perp Z$. In colliders, association flows only between the parents (i.e., X and Z here) and the child (i.e., Y) node. Thus, in the collider example above, $X \perp\!\!\!\!\perp Z$, but $X \not\perp\!\!\!\!\perp Z$ and $Z \not\perp\!\!\!\!\perp Y$.

When we condition on variables, we say they are blocking variables because they may block the flow of association. When a variable is conditioned on, or blocked, in chains and forks, they no longer allow the flow of association between the variables. In this way, we can "close" chains and forks. In the case of the chains and fork above, if we condition on Y, then Y is blocked, and we get the dependence relation $X \perp\!\!\!\perp Z \mid Y$. On the other hand, when the child node in a collider is conditioned upon, we have the opposite; colliders "open," and association

flows between parents. In the example above, when we condition on Y, we get the relationship $X \not\perp\!\!\!\!\perp Z \mid Y$.

From this, the definition of d-separation is as follows:

Nodes X and Y are d-separated given a conditioning set S, with $X,Y \notin S$, if and only if all paths between X and Y are blocked, denoted

$$X \bowtie Y \mid S, \tag{1.3}$$

where S may be empty. D-separation applies to the children of nodes as well. If Z in the collider above had a child node, W, then Z and W would be d-separated just as X and Z are d-separated.

Causal Assumptions

Like many statistical machine learning approaches, causal discovery has specific assumptions, some that depend on the algorithm and the data. In addition, there are three untestable assumptions and require domain expertise to safely assume: the causal Markov condition, faithfulness, and causal sufficiency. These assumptions represent the domain expertise required to infer beyond mere statistical inference to answer causal questions. They are summarized below, and are detailed further in Runge (2018a), which includes clear examples for each assumption that illustrate how algorithms can infer incorrect links when assumptions are not met.

The **causal Markov condition** is necessary for all independence-based methods. It states that if and only if the joint distribution of a time series process, **X**,

with the true causal graph G,

$$X_t^- \backslash \mathscr{P}_{Y_t} \bowtie Y_t \mid \mathscr{P}_{Y_t} \Longrightarrow \mathbf{X}_t^- \backslash \mathscr{P}_{Y_t} \perp \!\!\!\perp Y_t \mid \mathscr{P}_{Y_t},$$
 (1.4)

for all $Y_t \in \mathbf{X}_t$, with parents \mathscr{P}_{Y_t} . Essentially, this states that d-separation in the graph implies independence in the data. The contraposition is implied:

$$\mathbf{X}_{t}^{-} \backslash \mathscr{P}_{Y_{t}} \underline{\times} Y_{t} \mid \mathscr{P}_{Y_{t}} \Longrightarrow X_{t}^{-} \backslash \mathscr{P}_{Y_{t}} \bowtie Y_{t} \mid \mathscr{P}_{Y_{t}}$$

$$\tag{1.5}$$

The **faithfulness** assumption guarantees that the graph contains all conditional independence relationships that the causal Markov condition implies. A causal graph is faithful if and only if for the joint distribution of a time series process, X, with the true causal graph G, for all disjoint subsets of nodes $Y, Z, S \subset G$

$$X_Y \perp \!\!\! \perp X_Z \mid X_S \Longrightarrow Y \bowtie Z \mid S.$$
 (1.6)

This states that d-separation is implied by independence. The contraposition is also implied,

$$Y \bowtie Z \mid S \Longrightarrow X_Y \not\perp \!\!\!\perp X_Z \mid X_S.$$
 (1.7)

Causal sufficiency is often the more difficult to assume in open and complex systems. It assumes that all common causes of two or more variables are included in the analysis. Formally, a set of variables, S, is causally sufficient for a process, X, if and only if every common cause, or parent, of any two or more variables in W, is included in W, or has some value for all units in the population.

In this work, we are primarily interested in time series data and time-lagged relationships, and these methods require the **time-order assumption**: that the past causes the future, causality cannot travel faster than the speed of light, and that the future cannot cause effects in the past. Depending on the algorithm, assumptions for **stationarity** and **dependency type** are necessary. Glymour et al. (2019) argue that assuming nonstationarity may be allowed in some cases and could even be leveraged as more information. However, as Runge (2018a) notes, stationarity may be indicative of a confounding variable that violates causal sufficiency.

1.3.2 Consistency

Consistency is an important trait of a causal discovery algorithm. If an algorithm is consistent, it has been proven to converge to the true causal graph in the limit of infinite sample sizes. Each algorithm will be defined in part by a set of causal assumptions that are integral to the proof. A common set of those assumptions are described in Section 1.3.1.

Some algorithms, such as conditional independence-based approaches, require additional statistical assumptions. For example, conditional independence-based algorithms testing with a non-parametric regression independence test will need to assume that the function estimator converges correctly, that the noise in the model is additive and independent, and that the unconditional independence test of the residuals converges.

Universal consistency is defined for iterative causal algorithms (Runge, 2018a):

Denoted by \hat{G}_n , the estimated graph of some causal estimator from a sample of distribution P, with sample size n, and by the true causal graph G. A causal estimator is said to be universally consistent if \hat{G}_n converges in probability to G for every distribution P,

$$\lim_{n \to \infty} \Pr(\hat{G}_n \neq G) = 0. \tag{1.8}$$

This says that the probability of misestimating the true graph becomes arbitrarily small for large sample sizes for any distribution P.

Universal consistency is weaker than *uniform consistency*, which "bounds the error probability as a function of the sample size, giving a rate of convergence" (Runge, 2018a). For a merely universally consistent algorithm, the sample size required for a given error threshold will be different for each distribution, *P.* Runge (2018a) notes that uniformly consistent conditional independence-based algorithms can only exist under additional assumptions.

1.3.3 Validation and Falsifiability

As discussed by Runge et al. (2019b), method development in causal discovery requires benchmark datasets with ground truth causal structures. CauseMe.net is a website the authors have made for collecting benchmarking datasets for validating causal discovery algorithms. Ground truth for those data sets must come from expert knowledge or randomized experiments. Observational causal networks can be falsified with experimental results. Unfortunately, much of the motivation to use causal discovery is in cases where experimental results do not exist, when random-

ized control trials are infeasible. When expert knowledge of a causal structure and experimental results do not exist, causal models must be validated by validating each of the causal assumptions made by the algorithm. Since causal discovery algorithms can be proven to be *consistent*, as defined in Section 1.3.2, validating the assumptions can show that the resulting causal network is asymptotically correct to infinitely large sample sizes.

Peters et al. (2017, p.120) also discuss the falsifiability of causal models. They state that traditional machine learning algorithms build probabilistic models, structural causal models can be used for counterfactual models, and causal graphical models can be used for interventional models. They write that two models are equivalent if their corresponding predictions agree. Likewise, we can falsify a probabilistic or interventional model if the corresponding distributions disagree with the observed data. In the case of traditional machine learning, this is commonly computed with validation datasets to ensure that prediction distributions agree with unseen data. In the case of interventional, causal graphical models, if a model correctly predicts the observational distribution but fails to predict the interventional distribution, from a randomized trial, for example, then the model is falsified. Peters et al. (2017, p.120) state that falsifying counterfactual models is difficult in general.

1.3.4 Time Series Causal Discovery

Temporal information is critical to inferring the Earth system's dynamics because the Earth system is a temporal process. Many causal discovery methods imply the inherent temporal aspects of causality without representing it explicitly. Peters et al. (2017, p.10) note that although it is sometimes said that causality discussions must account for time, usually time is not necessary to discuss the effect of interventions. They write that both statistical learning and causal learning can be thought of as "abstractions of an underlying more accurate physical model that describes reality more fully." This is quite obviously true for numerical Earth system models in which differential equations define the dynamics of hundreds of quantities around the globe. It is even more so for the natural system that Earth system models attempt to estimate.

Peters et al. (2017, p.26) note that "an event can only influence events lying in its light cone, since no signal can travel faster than the speed of light in a vacuum." That is, physics explicitly excludes causation from the future to the past. They explain that although this is true, it is widely believed that microscopic and quantum mechanical systems are invertible. They say that the asymmetry of time-order is less critical for describing a causal relationship than the asymmetry of the information carried causal function between events. This is why time is not included in descriptions of physical laws, such as $F = m \times a$. However, the time-order asymmetry is sometimes essential for inferring the direction of causal dependence from

data alone (Runge, 2018a).

The consequence of discarding time-order asymmetries in data is that temporal information for interpreting dependence relationships is lost and cannot inform inference. If no other asymmetry is captured in the data, then we need temporal information to elucidate. Some systems, such as climatological processes, are often best summarized in data by time series. Rather than a set of independent samples, a summary in time is necessary to describe Earth system dynamics accurately. Conditional independence-based causal discovery is flexible enough to be adapted for time series input (Runge et al., 2019a).

Many causal discovery algorithms are designed for independent and identically distributed samples. The causal graph can include temporally lagged variables to capture temporal relationships between variables. Each node is multiplied into nodes for each time step. Theoretically, this creates an infinitely large *time series* graph, which each variable, X, is represented as many nodes, $\{X_t, X_{t-1}, X_{t-2}, X_{t-3}, ...\}$. In practice, we limit the number of lags to a time step that is large enough to capture the theoretical temporal dependence between the variables of interest.

1.4 Earth Science Challenges

A critical problem in Earth science is identifying the causal pathways from an intervening Earth system event, such as a wildfire, volcanic eruption, or atmospheric injection, to the many impacts on climate, weather, ecology, and human livelihoods in various places on Earth. Causal pathways are paths through a graph of nodes

representing various climate impacts or quantities of interest. There is a critical need for analyses that trace the causal path from an intervention, through mediating effects, to impacts that affect life, economic systems, natural resources, and more.

Climate interventions of interest include anthropogenic climate change and natural and artificial stratospheric aerosol injection (SAI). Volcanoes are an occasional source of natural interventions in the climate, injecting teragrams of gases into the stratosphere (Guo et al., 2004a); though eruptions of that magnitude are rare, only occurring every 50 to 100 years. Artificial SAI events are manufactured efforts to change climate regionally or globally. Examples include geoengineering ideas, such as reducing global mean temperatures with sulfuric gas injection. A related example is weather interventions, such as China's rain-making effort, Sky River (Gimeno et al., 2014; Wang et al., 2018), which attempts to bring more rain to a historically arid region. Understanding the downstream impacts of these interventions is vital for evaluating the risks of geoengineering and predicting the impact on neighboring regions.

Reconstructing the causal space-time pathways from intervention to impact will provide critical insights to understand intentional and unintentional interventions. If the global community decides to attempt geoengineering to mitigate climate change impacts, experiments may start small and localized. We need tools to understand the effects of the experiments. If another country decides to implement geoengineering for itself, perhaps at the expense of its neighbors' moisture, then

causal analysis will be critical for understanding those impacts.

Many Earth science problems, particularly those considering a relatively short time window, are very data-sparse. Measurement frequency can vary depending on the variable, quality needs, and equipment. Sometimes daily or sub-daily observations are available, but not for very far into the past, often weekly or monthly data is most abundant. The dataset may contain hundreds of variables on millions of grid cells. Frequently, one may want to understand the interdependence of a few variables in several hundred positions with an order of magnitude fewer observations per variable/position pair. This presents a high-dimensional problem, posing poor statistical power and high sample complexity for statistical methods.

1.4.1 Earth Science Data

Earth system data is obtained from several different sources, such as station measurements, satellites, data-fused reanalysis products, and Earth system model output. The data is multimodal and can have a large variety of spatial and temporal resolutions. Station measurements can poll a quantity very often, but only provides data for a point in space. Satellites cover large strips of space over the globe, but measurements can be less frequent, particularly in a specific area of interest, and still often have missing data due to cloud cover. Reanalysis products combine station measurements, satellite data, and weather or climate modeling to produce a hybrid of fused, interpolated data that generally covers all space on the globe.

Reanalysis products and Earth system model output are most convenient be-

cause they are spatially complete and temporally consistent, but come with more assumptions than raw measurements. Spatially, the data from these sources is generally arranged on a discrete 3D grid. Grids can take many forms, most common are cubed latitude-longitude grids. Geodesic grids are also used in order to achieve better geometric regularity between cells. (Ebert-Uphoff and Deng, 2014). Earth system model output is frequently analyzed on a per-run basis, a per-model basis with ensembles of runs, or with Coupled Model Intercomparison Project (CMIP) output. CMIP is a collaboration project that combines output from over 100 models, sourced from over 50 modeling centers.

The research in this dissertation addresses many of the ideas and challenges above. It examines the capabilities and limitations of contemporary data-driven modeling. After identifying a key research gap in grid-level causal discovery, this work introduces two novel methodologies for causal discovery of local grid-level dynamics, CaStLe and M-CaStLe, that advance the state-of-the-art in performance and efficiency. I demonstrate these advances with new benchmarking approaches and realistic applications in the Earth sciences. The following chapters detail the path from explainable machine learning for Earth system model evaluation to causal discovery of Earth system dynamics to novel causal discovery approaches for gridded space-time data. With these advances, this work contributes to toolsets for further scientific discovery.

2 Related Work

The RCT was the first innovation to measure causal effects in experiments directly. Ronald A. Fisher is credited with first using randomization for experiments in 1925 (Fisher, 1925; Hall, 2007). Around the same time, Wright (1921) wrote about using what he called path analysis to evaluate and represent directed statistical dependencies. According to Pearl and Mackenzie (2018a), path analysis is a direct ancestor to modern causal inference techniques, though it was not recognized as such until the 1950s. Splawa-Neyman et al. (1923) were the first to publish on a potential outcomes framework, providing a notation for causal effects in a randomized setting (Rubin, 2005).

In the 1970s, Donald Rubin's potential outcomes framework opened the door to causal inference in non-randomized observational studies (Rubin, 1974). Potential outcomes try to address the *fundamental problem of causal inference*: once treatment is given to an individual, we can no longer observe what could have occurred had the individual not received treatment. More specifically, as Holland (1986) writes, "it is impossible to observe the value of $Y_t(u)$ and $Y_c(u)$ on the same unit and, therefore, it is impossible to observe the effect of t on u" for potential outcomes, Y, of treatment, t, and control, c, on the individual unit, u (Holland, 1986). While these quantities cannot be observed or computed, this framework

allows us to compute other causal quantities based on certain assumptions in nonrandomized studies.

Pearl (2012) added to Rubin's potential outcomes notation with the *do-calculus*, a way of clarifying the notation for describing the change in probability distributions of a given quantity from *doing* an intervention on that quantity. In 2000, Pearl presented the structural causal model (SCM), which is a nonparametric form of structural equation models (SEM) (Pearl, 2000, 2001). Economists and sociologists have used SEMs for decades, and they trace their conceptions to Spearson (Tarka, 2018).

Among many other contributions, Robins (1986) introduced a graphical approach to causal inference, the finest fully randomized, causally interpretable structure tree graph. Pearl (1995b) improved on this approach by introducing directed acyclic graphs (DAGs) from computer science and graph theory to causal inference. In that work, Pearl shows how independencies can be described in a Bayesian network graph and how we can similarly represent causal relationships.

2.1 Causal Discovery

Causal discovery, or causal learning, is the pursuit of computing the causal structure from observational data. It intends to outline when an association is causal or merely correlated (Peters et al., 2017). Many algorithms do this by detecting spurious correlations in data, and after making strict assumptions, the causal structure can be found (Runge et al., 2019b). The necessary assumptions are derived

from decades of previous causal inference literature. They can be used to prove *consistency*, which is the property that an algorithm converges to the true causal graph in the limit of infinite sample data (Runge, 2018a).

Wiener (1956) published the idea that a variable could be considered causal to another if the ability to predict the second is improved by including information about the first. Granger (1969) later published a practical method for computing on this notion, now known as *Granger causality*. Typically, Granger causality refers to linear bivariate analysis using linear regression models (Peters et al., 2017) or vector autoregressive models (Runge, 2018a). Granger causality has several limitations, outlined in Peters et al. (2017), including an inability to detect indirect causes, failure in the presence of deterministic dependencies, a limitation to only detecting lagged dependencies, and problems with sub-sampled time series (Runge et al., 2019b).

A nonlinear, multivariate approach to Granger causality is called transfer entropy (Peters et al., 2017; Runge, 2018a; Runge et al., 2019b). Peters et al. (2017) state that transfer entropy fails in many of the same scenarios as Granger causality. However, they write, "we emphasize that the qualitative statement about presence or absence of causal inference in the case of two causally sufficient time series only fails for a rather artificial scenario, while quantifying the causal influence via transfer entropy can be problematic also in less artificial scenarios," (Peters et al., 2017, p.207).

2.1.1 Causal Network Learning

In the 1990s, Peter Spirtes, Clark Glymour, and Richard Scheines developed graphical causal discovery, also known as causal network learning (Spirtes et al., 1993). Spirtes and Glymour invented the PC algorithm, named for their first names (Spirtes and Glymour, 1991). This algorithmically attempts to reconstruct the causal structure from observational data. The main underlying idea stems from Reichenbach's Common Cause Principle (Reichenbach, 1956): that if two variables are statistically dependent, there must be a causal relationship between the two or a third common driver of the two.

The full description and pseudocode for PC can be found in Spirtes et al. (1993), and I will provide a brief outline here. It begins with a fully connected graph in which each node is assigned a variable. To leverage Reichenbach's principle, PC iteratively tests each pair of variables, X and Y, for independence, conditioned on a set of one or more variables, X denoted $X \perp \!\!\!\perp Y \mid X$, while dependence would be denoted $X \perp \!\!\!\perp Y$. If two variables are conditionally independent, their link is removed. This first phase results in an undirected *skeleton* graph. In short, the second and third phases use rule sets to orient edges based on principles of how association flows between nodes in a graph. See *d-separation*, detailed here in Section 1.3.1 and in Spirtes et al. (1993). To accurately estimate causal effects, PC relies on strict assumptions, including faithfulness, the causal Markov condition, and causal sufficiency.

Causal sufficiency is one of the more challenging and commonly violated assumptions in causal inference. Spirtes, Glymour, and Scheines' fast causal inference algorithm (FCI) does not require the causal sufficiency assumption (Spirtes et al., 1993). This algorithm does not require the causal sufficiency assumption and, as a consequence, will only produce a Markov equivalence class of partially directed acyclic graphs. The consistency of PC and FCI is shown in Spirtes et al. (1993).

Runge et al. (2019a) published an adaptation to the PC algorithm called PC momentary conditional independence (PCMCI). PCMCI is specifically written for reconstructing lagged-causal time-series graphs (Runge, 2018a). This two-phase algorithm first uses a modified PC algorithm adapted for time series, called PC₁, which attempts to construct a sparse partially directed graph. In the second phase, momentary conditional independence (MCI) is computed for each connected variable pair to reduce the graph further to converge on the estimated causal graph. MCI conditions on both the parents of a given variable, X, as well as the lagged, or time-shifted, parents of X.

Each phase of PCMCI serves a specific purpose in identifying the causal structure. PC₁ removes irrelevant conditions of each variable via iterative conditional independence tests. PC₁ tests only the condition subset with the largest association instead of testing all possible combinations like PC (Runge, 2018a). The MCI phase then controls the relatively high false-positive rate for highly interdependent time series. Conditioning on lagged parents of each variable controls for highly

autocorrelated time series data and makes MCI an estimator of causal strength. Both PC₁ and MCI can be implemented with any conditional independence test. Tests for linear models, nonlinear additive noise models, and nonparametric models exist (Peters et al., 2017; Runge, 2018a; Runge et al., 2019b).

Runge (2018a); Runge et al. (2019a) show empirical results from tests on synthetic data to benchmark PCMCI against several other algorithms, including PC, FCI, convergent cross-mapping, LiNGAM (Shimizu et al., 2006), and Granger-causality. They show that PCMCI performs best or above average in terms of high true positive rates and low false positive rates on time series data in several tests with dynamical noise, autocorrelation, and high dimensionality. After identifying the graph, PCMCI was also able to compute true causal effects well (Runge et al., 2019a).

The PC, FCI, and PCMCI algorithms are examples of causal discovery's conditional independence (CI) based causal network learning pillar. These are highly adaptable algorithms because they can be implemented with any conditional independence test. Choosing the correct one depends on specific assumptions about the data and the functional form of the dependencies within. These range from the linear partial correlation test, nonparametric residual-based tests for nonlinear dependencies with additive Gaussian noise (Ramsey, 2014; Runge et al., 2019a), kernel-based approaches (Zhang et al., 2011a), information-theoretic conditional mutual information (Runge, 2018b), and neural networks (Sen et al., 2017).

2.1.2 Structural Causal Models

Because Granger causality and many causal network learning algorithms require a time delay between cause and effect, they cannot easily determine contemporaneous dependencies (Runge et al., 2019b). Contemporaneous dependencies primarily exist when causation occurs faster than the available time-sampling interval. SCMs typically ignore the time-order of causal dynamics; instead, they operate on the assumption that the past is already coded into covariates (Peters et al., 2017). They can estimate contemporaneous effects because they make additional assumptions about the functional forms between dependencies (Runge et al., 2019b).

As SEM's causal-descendant, SCMs are used to model nonlinear causal relationships and require added assumptions for correct estimation (Peters et al., 2017). These allow for the estimation of direct and indirect causal effect, a quantitative estimate of causal strength, without further assumptions on the functional forms interdependencies or distribution of error terms in the data (Tarka, 2018). Peters et al. (2017) overview SCMs in the bivariate and multivariate cases. They describe SCMs' uses for causal discovery and applications to machine learning. Despite their advantages, SCMs have not yet been applied to Earth system sciences (Runge et al., 2019b).

2.2 Attribution in Climate Science

While the climate science literature does not broadly use causal discovery or causal inference techniques explicitly, a primary interest in climate science is detecting and attributing changes in our climate. *Detection* and *attribution* have precise definitions in climate science. Detecting a signal change requires demonstrating that the observed signal differs in a statistically significant way from natural variability. Detection does not imply an attribution of that change. Attribution requires (1) showing that an observed signal is unlikely in natural variability, (2) consistent with estimated changes to the signal given anthropogenic and natural forcing, and (3) not consistent with alternative, plausible explanations of the observed signal (Houghton et al., 2001).

In 1996, Klaus Hasselmann published one of the first attempts to quantitatively attribute climate changes (Hasselmann, 1997). Until then, there was mounting evidence that global warming could be attributed to anthropogenic forcing, but it was largely qualitative or circumstantial. He provides a multi-pattern fingerprinting framework for statistically attributing climate signals.

Hasselmann states that for the attribution problem, further hypotheses regarding the cause of a detected change need to be considered. This demonstrates the counterfactual theory required for causal inference (Pearl and Mackenzie, 2018a). He further writes that an obstacle for quantitative signal-to-noise analyses is that they require information on the space-time structure of the predicted climate signal

and the climate variability. This implies the same expert provided causal structural knowledge that Runge, Pearl, Peters, and others suggest is critical for effective causal inference (Runge et al., 2019b; Pearl and Mackenzie, 2018a; Peters et al., 2017). He then describes an idea similar to causal sufficiency: "A discrimination between competing forcing mechanisms can clearly be meaningfully attempted only if all candidate mechanisms and their associated climate change signals are specified."

Finally, because of the finite nature of real data, Hasselmann states that it can never be ruled out that there may be other overlooked forcing mechanisms that would generate the observed signal. The consequences of this fact are "unequivocal attribution is achieved only in the hypothetical infinite-sequence limit ... We must, therefore, restrict ourselves in principal to a statistical definition of attribution that applies only in the limited sense of establishing a ranking within a given finite set of candidate forcing mechanisms." This essentially iterates the same limitations of finite data in causal discovery (see *consistency* in Section 2.1), detailed by Runge (2018a) and Peters et al. (2017) and described in Section 2.3.1.

2.3 Causal Discovery for Earth Systems Science

Causal discovery has been applied to Earth systems science several times recently. Runge et al. (2019b) cite several papers in which Granger causality, causal network learning algorithms, and nonlinear state-space methods have been applied to climate science problems. Causal network learning applications are relatively

recent and primarily focused on climate science (Ebert-Uphoff and Deng, 2012; Kaufman et al., 2020; Kretschmer et al., 2016; Nowack et al., 2020a; Runge et al., 2014). These will be detailed further in Section 2.4.

Runge et al. (2015a) present a framework for identifying gateways, mediators, and causal effects in Earth systems. First, they use varimax-rotated principal component analysis on gridded sea-level pressure data to identify localized areas of variability, such as the El Niño Southern Oscillation and the Quasi-biennial Oscillation climate modes, as described in Vejmelka et al. (2014). With those, they can project the original data onto the selected components to create a time series signal for the given quantity in several regions. They then use the regions as nodes in their time series causal discovery algorithm, which identifies the causal relationship between nodes and removes spurious associations found in the data. With that, they are able to identify teleconnections between climate modes and sea level pressure components. Beyond that, they use their established causal networks to compute causal effect metrics for how much a component impacts others in the space-time system.

Runge et al. (2019b) give an overview of causal discovery methods for Earth systems science problems. They identify several classes of causal discovery methods suited for several classes of problems. The classes of problems they list are causal hypothesis testing, complex network analysis, analysis of the causes of extreme events, and causal model comparisons. They also provide examples of these methods used to solve various space-time problems, including an Arctic climate

problem, an ecology problem, and a cardiology problem. They discuss the many challenges in applying causal discovery to Earth systems science, from methodological to data to computational and statistical challenges. These are discussed in detail in 2.3.1. Finally, they present future research directions for causal discovery and call for more scientists to work on using causal discovery to solve the challenges in Earth systems science.

Eyring et al. (2019) published a perspective paper on climate model evaluation tools. In it, they say that better tools are required to effectively evaluate the quality of climate models. Climate models are our primary means of studying and experimenting with climate dynamics, and understanding how well they perform is critical to that research. They say, "other promising diagnostic developments on the horizon that should be further advanced include studies that assess responses to perturbations rather than mean climate, and the application of innovative data science methods in Earth system science such as neural networks, machine learning-based anomaly detection techniques, graphical models and causal discovery."

2.3.1 Specific Application Challenges

Runge et al. (2019b) overview the process, data, and computational and statistical challenges faced in applying causal discovery to Earth sciences. The following is a recapitulation of the relevant challenges in that overview.

Process Challenges

The time-dependent processes in Earth systems give rise to strong **autocorrelation** and **time delays** for processes to act on one another. **Nonlinearities** and **state-dependencies and synergies** make selecting an estimation method critical. The wrong method may struggle to disentangle the autocorrelation and internal dynamics and thus fail to achieve the correct causal structure. Various geoscience time series may be acting on **different time scales**, which can be separated to incorporate and interpret different relationships. Many statistical methods make assumptions about the **noise distribution** in the data. Many methods assume additive Gaussian noise, but nonlinear and model-free solutions exist (Peters et al., 2017; Runge, 2018a; Runge et al., 2019b). Processes with heavy tails and extreme outliers may violate linearity and normality assumptions.

Data Challenges

Climate data is space-time, meaning it is measured and computed on a 3-dimensional grid over the Earth's land, oceans, and atmosphere. Hundreds of individual quantities can be collected, leading to a very high-dimensional problem. Extracting features from this data can be a big challenge.

Observational data is incomplete; some processes cannot be adequately measured and quantified. It comes from satellite and station measurements and can include several forms of **measurement error**, such as measurement noise, instrumental biases, and missing data. Often, observational data comes in the form of

reanalysis data. Reanalysis is a data assimilation effort to fill data gaps and meaningfully represent quantities of interest via observed data and model output. Finally, satellite measurements only date to 1979, so observational time series are often short. If problems with the observational data are directly related to the processes of interest, then selection bias may be a problem.

Simulation data is vast. Although its spatial resolution is generally smaller than observational data, it is typically 0.5 degrees to 1.0 degrees latitude and longitude. The temporal resolution and timescales of simulations are often higher than observational datasets. They can span hundreds of years and include hourly data.

Because of the high-dimensional and complex data, **variable extraction** is difficult. Time series variables need to be extracted from space-time data; sometimes, feature construction techniques are necessary to form causally relevant features. To do this, fingerprinting (Hasselmann, 1997) and dimensionality reduction techniques (Vejmelka et al., 2014), such as empirical orthogonal functions (EOF)¹ (Hannachi et al., 2007) and varimax-rotated principal component analysis (PCA) (Hannachi et al., 2007), are often necessary. Additionally, these features should be interpretable, representing physical processes in the system.

Often, causal drivers cannot be measured, which leads to **latent**, **or unobserved**, **variables** in the analysis. The absence of common causes, or a variable that causes two or more other variables, can lead to spurious links detected in the causal discovery algorithm. Runge notes that failing to account for important

¹The climate community refers to principal components as EOFs.

drivers, such as anthropogenic climate forcings, may render time series stationary.

Like latent variables, **subsampling** is when a time series is too infrequently sampled. If the causal mechanism acts on a smaller time scale than measured, the mechanism may not be detectable. On the other hand, **Time-aggregation** may reduce the data size and algorithm's computational complexity, but it can make relationships appear contemporaneous or cyclic.

Computational and Statistical Challenges

Sample size and dimensionality is an issue for the scalability and time complexity of many causal discovery methods. While many methods are proven to be correct in the limit of unlimited data by consistency, they are typically relatively slow, some polynomially and some cubically (Runge et al., 2019b). The opposite problem is more likely in observational climate science because, as mentioned earlier, the observed record is still short. When sample sizes are too small, causal relationships may not be reliably estimated. In the case of PC and related methods, conditional independence tests may produce incorrect results, and orientation rules may contradict each other if sample sizes are too small. If dimensionality is high and the sample size is small, conditional independence tests may be underpowered. Lastly, uncertainty quantification, which includes statistical test uncertainties and data measurement uncertainties, is an ongoing research challenge for causal inference.

Rejoinder to the Challenges

Most of these challenges discussed are also challenges for traditional correlation, regression, and machine learning methods. However, interpretation of those remains nebulous and often leads to incorrect conclusions. The assumptions made by causal inference and causal discovery merely require subject matter expertise; they encode the domain knowledge to infer causal dependencies and reject spurious association from observational data. Likewise, it is a mistake to embark on traditional statistical and machine learning endeavors without subject matter expertise because of the propensity to mishandle data and make spurious inferences. Because of these factors, Runge et al. (2019b) note that there is "no strong reason to avoid adoption and exploration of modern causal inference techniques."

It seems clear that climate attribution, described in Section 2.2, and causal discovery are fundamentally equivalent endeavors, from intent to results and limitations. Given that both are approached correctly, they are equally valid in asserting the causal dependence between climatological processes. This further iterates Runge's assertion that there is no reason to avoid the exploration of modern causal inference for learning about the Earth's climate.

2.3.2 Recent Efforts to Overcome Application Challenges

As described above, one of the challenges in statistical and causal inference in climate science is the amount of data available. It is common for observational and simulated datasets to be available on a coarse temporal resolution, such as monthly.

When we seek to discover causal dependencies that occur on a finer resolution than measured, we may only find contemporaneous or undirected dependencies. In fact, one of the basic assumptions of the PC, FCI, and PCMCI algorithms is *no instantaneous effects* (Spirtes et al., 1993; Peters et al., 2017; Runge, 2018a). That is, no two variables may act on one another instantly or, practically speaking, within one observed timestep.

To detect contemporaneous links, rather than assume they do not exist, Runge published an adapted version of his PCMCI algorithm, which he calls PCMCI+ (Runge, 2020). Runge notes that autocorrelation is key to increasing contemporaneous link orientation recall. PCMCI+ also "improves the reliability of CI tests by optimizing the choice of conditioning sets and yields much higher recall, well-controlled false positives, and faster runtime than the original PC algorithm for highly autocorrelated time series." Empirically, it maintains performance for time series with low autocorrelation.

Similar to FCI, Runge's Latent PCMCI (LPCMCI) is an implementation of PCMCI to handle the case in which causal sufficiency cannot be assumed, when latent variables exist (Gerhardus and Runge, 2020). This algorithm is tolerant of latent variables while possibly illuminating their existence. Tolerating latent confounding is critical in many open systems in which it is impossible to observe and account for all confounding. The downside of these methods is that they can only estimate the causal structure up to a Markov equivalence class.

2.4 Applications of Causal Network Discovery for Climate Science

Ebert-Uphoff and Deng (2012) may have been the first to apply causal networks to climate science in 2012. They cite inspiration from seminal papers from Tsonis and Roebber (2004) and Tsonis et al. (2006) for their initial work on correlated climate teleconnections, and from Pearl and Mackenzie (2018a) and Spirtes et al. (1993), for their causal discovery work. Ebert-Uphoff and Deng apply the PC algorithm to 500 millibars geopotential height at individual grid cell locations. Geopotential height is the height above sea level of a specific pressure level in a specific location, adjusted for the variations in gravity due to changes in latitude.

Ebert-Uphoff and Deng's work is similar to previous work identifying correlated teleconnections by creating a network of dependencies between grid cells on the globe of one variable. Their contribution is to apply causal inference to those teleconnections, removing spurious relationships and identifying a causal network. There are a couple of limitations to their approach. Without including a time series implementation of the PC algorithm, their method treats each day's observation as an independent sample rather than a time-dependent process. They also use neighboring grid cells in the network, possibly violating independence assumptions in the conditional independence tests. Major modes of climate variability may not be adequately captured in single grid cells either, so a weaker signal may lead to undetected links. Still, grid cell level nodes may increase the total captured spatial

variability because spatial aggregation and dimensionality reduction techniques can reduce variance.

Kretschmer et al. (2016) applied causal discovery to detect causal effects in Arctic midlatitude winter circulation. They apply a version of the PC algorithm adapted for time series and use seven different variables in the Arctic. They are regional ice, ocean, and atmospheric quantities. They aggregated daily data into monthly means because they were specifically looking for processes acting on a monthly scale. Finally, they used weighted spatial averaging to convert the data into 1-dimensional time series. They validated their findings by careful analysis of the variable selections. They selected variables from work conducted in previous Arctic climatological studies and included proxies for some unmeasurable complex phenomena.

Nowack et al. (2020a) used PCMCI to evaluate how similar climate model runs were to observed dynamics. Specifically, they developed graphs depicting how sea level pressure in 50 regions on the globe relates to each other region. In the correlation setting, a relation between variables across space on the globe is called a teleconnection. They discovered causal graphs for 20 models in the Coupled Model Intercomparison Project Phase 5 (CMIP5). Each model was represented by several simulation runs, each used to generate their own graph. They used the F_1 score to measure the similarity between graphs.

Using spatial sea level pressure data, Nowack et al. (2020a) detected 50 regions of interest using a common technique in climate science. First, PCA identified

the first 100 orthogonal components. Then, they use the varimax rotation algorithm, which has been found to increase the interpretability of components and localize them in space. They note, "principal components without rotation consecutively maximize variance and therefore often mix contributions of physically defined modes such as the El Niño Southern Oscillation, Pacific Decadal Oscillation, or the North Atlantic Oscillation, whose time-behavior is not orthogonal, making patterns more difficult to interpret."

Finally, they select 50 of the 100 components based on spatial separability and frequency spectra. Resulting are 50 discrete regions with high variability and independent patterns. They used the 50 components for each node in the causal discovery analysis. Lastly, they note that "the selection of components defining the network nodes will typically be guided by expert knowledge in conjunction with dimension reduction techniques."

Tibau et al. (2022) built on the dimensionality reduction approach, augmenting it to output grid-cell-level networks. They specifically delineate *mode-level* (dimensionality reduction or cell aggregation) and grid-level causal discovery. Their augmentation is called Mapped-PCMCI, which first applies dimensionality reduction, then computes a mode-level causal network with PCMCI, and finally maps the grid cells within the modes to each other using the network previously constructed. Their resulting network consists of edges between grid cells, but the method assumes that cells within modes are fully connected, i.e., each cell is dependent on all of its neighbors. In contrast, our work specifically seeks inter-cell

spatial relationships. Finally, they also describe the failure of a traditional causal discovery approach for grid-cell-level data, "[if] we apply PCMCI directly at the grid-level, the low power of this high-dimensional and redundant estimation problem (see Section 2.2.2) leads to most links being missing."

Recently, a new tradition, causal representation learning, developed out of machine learning to leverage causal reasoning for their models (Schölkopf et al., 2021). While still a developing field, it shows particular promise for estimating relationships in the presence of latent confounding. Boussard et al. (2023) and Brouillard et al. (2024) developed the Causal Discovery with Single-parent Decoding (CDSD) algorithm within the causal representation learning framework and applied it to the climate science field. CDSD performs well in high-dimensional data settings but through a different mechanism. It performs dimensionality reduction by learning latent variables and enforcing a "single-parent" constraint where each grid cell belongs to exactly one latent factor. This naturally clusters grid cells into coherent, often contiguous regions and enables the discovery of causal relationships between these larger-scale patterns. In contrast to grid-level structure learning, CDSD identifies broader teleconnection pathways between regional climate modes. Thus, CDSD abstracts to a higher level by mapping the native grid space to an identifiable latent representation before performing causal discovery.

Several studies have addressed local-scale phenomena. Pfleiderer et al. (2020) applied causal discovery to identify precursors to seasonal hurricane frequency. They utilized the precursors to inform a predictive model. Polkova et al. (2021)

identified local drivers of marine cold-air outbreaks in the Barents Sea. These demonstrate that existing causal discovery approaches can be valuable for seasonal and sub-seasonal phenomena. However, both marginalized large regions prior to analysis, reducing the space's dimensionality, and did not evaluate the space-time evolution of phenomena nor grid-level dynamics.

There are some examples of causal discovery algorithms leveraging spatial information. Zhu et al. (2016) developed pg-Causality that applies space-time pattern mining and a Gaussian Bayesian Network to seek local dependencies in the space-time propagation of air quality data. Sheth et al. (2022) developed STCD for understanding hydrological systems. They constrained the discovery of spatial structures by only allowing higher elevation nodes to be parents of lower elevation nodes because water follows the gravity gradient. While both cleverly use mined or known spatial structure to inform their causal discovery, they are both limited to use in sparse point-measured data from static base stations rather than gridded data. Further, these methods enforce constraints as filtering mechanisms. Neither address the scalability challenges in high-dimensional gridded data.

Parallel Approaches in Neuroscience: Causal Discovery for High-Dimensional Spatial-Temporal Data

Other scientific domains face similar challenges with high-dimensional space-time data. Neuroscience, for example, needs to study mechanisms in brain interactions, and fMRI images may contain thousands to millions of pixels. The anatomy of the brain also exhibits locality constraints. Ramsey (2014) made computational

optimizations to the Greedy Equivalence Search algorithm, including sparsity constraints and limiting the distance of potential parents, to recover graphs with millions of nodes. Saetia et al. (2021) marginalized regions of interest in the brain using spatial averaging and then applied the PCMCI algorithm to construct causal graphs. There is a common interest in recovering graphs of high-dimensional gridlevel data throughout the sciences. Developing more tools that enhance the estimation and interpretability of causal graphs in these spaces will help advance our understanding of space-time structures across the sciences.

What is clear from prior work is that grid-level analyses are challenging, both statistically and computationally, due to how many grid cell dependencies need to be estimated, the enormous number of observations needed, and the redundant information content of nearby cells.

Part I

Foundations of Structure Learning for Earth Systems

3 Machine Learning Feature Analysis Illuminates Disparity Between E3SM Climate Models and Observed Climate Change

3.1 Publication Notes

Citation: Nichol, J. Jake, et al. "Machine learning feature analysis illuminates disparity between Energy Exascale Earth System Model (E3SM) (E3SM Project, 2018) climate models and observed climate change." Journal of Computational and Applied Mathematics, vol. 395, 2021, p. 113451.

Publication date: October 2021

Publisher: Journal of Computational and Applied Mathematics

Formatting: The original published text has been preserved as much as possible while still adhering to the formatting requirements of this dissertation.

Data and Software Availability: The paper is available at https://www.osti.gov/biblio/1782577.

Funding: This work is supported by Sandia Earth Science Investment Area Laboratory Directed Research and Development funding. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and

Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.

3.2 Abstract

In September of 2020, Arctic sea ice extent was the second-lowest on record. State of the art climate prediction uses Earth system models (ESMs), driven by systems of differential equations representing the laws of physics. Previously, these models have tended to underestimate Arctic sea ice loss. The issue is grave because accurate modeling is critical for economic, ecological, and geopolitical planning. We use machine learning techniques, including random forest regression and Gini importance, to show that the Energy Exascale Earth System Model (E3SM) (E3SM Project, 2018) relies too heavily on just one of the ten chosen climatological quantities to predict September sea ice averages. Furthermore, E3SM gives too much importance to six of those quantities when compared to observed data. Identifying the features that climate models incorrectly rely on should allow climatologists to improve prediction accuracy.

3.3 Introduction

We have observed dramatic declines in Arctic sea ice since the advent of satellite imaging (Stroeve and Notz, 2018). This change is of critical importance to global economic, social, political, and ecological landscapes, not least because of the

opening of new navigable sea routes and the impact on wildlife (Arc, 2019; Smith and Stephenson, 2013). As an essential component of the Earth's climate, sea ice loss drives the positive feedback between surface albedo and Arctic warming and may contribute to changes in ocean circulation and mid-latitude weather (Goosse et al., 2018; Sevellec et al., 2017; Cohen et al., 2018; Cvijanovic et al., 2017).

Earth system models (ESMs) provide state of the art simulations of the global climate. They include general circulation and thermodynamic models for ocean and atmosphere, and models for land, sea ice, and land ice processes. Collecting an ensemble of parameterized ESM runs produces a distribution of forecasts that provide bounds on predictions. Simulations of Arctic sea ice in these models include complex interactions between the ice, ocean, and atmosphere. However, limitations in ESMs, such as the inability to resolve critical small-scale processes, can lead to biases when compared to observations. It is, therefore, critical to identify sources of bias.

Previous generations of ESMs have, on average, underestimated the rate of sea ice loss in the Arctic (Rosenblum and Eisenman, 2017). This is apparent in data from the Coupled Model Intercomparison Project (CMIP), which includes simulation results from a broad array of ESMs from modeling centers around the globe. CMIP *phases* mark improvements in the state of the art. The extent of sea ice loss has been a consistent problem, first identified in phase 3 (Meehl et al., 2007; Stroeve et al., 2007). By phase 5 (CMIP5), overall model bias had improved (Taylor Karl E., Stouffer Ronald J., 2012). However, Rosenblum and Eisenman

(Rosenblum and Eisenman, 2017), in an analysis of 118 simulation runs from 40 CMIP5 simulations, found that 89% of CMIP5 model runs underpredicted the rate at which sea ice extent is lost (km²/decade) by more than a standard deviation; and 2014 loss by an average of 2 million km². The disagreement with observation may imply that ESMs' parameters are not well-tuned. Stroeve et al. (Stroeve et al., 2007) suggest this discrepancy is due to missing key causal mechanisms or represent a misunderstanding of underlying physical processes.

The Energy Exascale Earth System Model (E3SM) (E3SM Project, 2018), developed by the United States Department of Energy (DOE), is included in phase 6 (CMIP6) (Eyring et al., 2016) (March 2019). E3SM is a new state of the science climate modeling and prediction project. In CMIP5 and E3SM, the rates of pan-Arctic sea ice change are similar to observation before 1996 but deviate from observation afterward. In CMIP5's case, the rate of loss is less than observed (Rosenblum and Eisenman, 2017), while E3SM's is greater than observed (Section 3.5.1: Data). These differences in sea ice loss rates lead to inaccurate long term predictions about absolute sea ice extent in the Arctic. To our knowledge, our work is the first mechanistic analysis of E3SM accuracy.

We use random forest regression (RFR) (Breiman, 2001) and Gini importance (Nembrini et al., 2018) to determine which E3SM features drive climate predictions. We perform an identical study of historical observations to identify the features that are most influential on prediction of actual sea ice loss. By comparing the two, we determined that E3SM relies too heavily on some features, to the detri-

ment of others, resulting in a divergence from observation. This work elucidates differences in sea ice response between observational data and E3SM simulations and can help improve sea ice prediction.

3.4 Related Work

Stroeve et al. (Stroeve et al., 2012) analyze the agreement between simulated Arctic models, CMIP3 and CMIP5, and observed data. They report that while phase 5 models are an improvement over phase 3 they consistently overestimate forecasted ice extent in the Arctic. The authors suggest that modeling may be improved by including more complex mechanisms such as sea ice albedo parameterization, thickness distributions, and melt ponds.

Rosenblum and Eisenman (Rosenblum and Eisenman, 2017) examined CMIP5's sea ice extent predictions in the Arctic and found overprediction of sea ice extent. Correcting the models required an increase in warming well above observed rates, leading the authors to conclude that the current methods were systematically flawed.

Ionita et al. (Ionita et al., 2018) presented a method for using multiple linear regression to predict the September sea ice extent minimums in the pan-Arctic region and the East Siberian Sea. Notably, they used step-wise regression because it may highlight the underlying coupled physical mechanisms between factors. For the pan-Arctic region, their model was able to predict sea ice extent anomalies for May, June, and July fairly accurately (reporting r-values between 0.84 and 0.9).

Although they found a "skillful" model could be built from their list of Arctic features, they did not analyze the relative importance of those features for their models.

Reid and Tarantino used support vector regression (SVR) to predict the Arctic sea ice extent (Reid and Tarantino, 2014). SVRs were able to construct predictive models, but they only considered sea ice extent as a predictor and could not analyze any other features for their importance. They chose SVRs because they are successful in predicting complex dynamical systems such as climate. The authors reported the comparative results of tuning the SVR, and compared them to CMIP5 ensembles but not to observation.

3.5 Data and Methods

Our methods were able to account for discrepancies in climate simulations and observations. Like multiple linear regression and its associated term-weights, random forests are a machine learning method that is wholly transparent (Breiman, 2001), unlike many other so-called "black box" methods, such as SVRs. We used RFRs and their corresponding Gini importance measure to determine how much influence each input feature has on E3SM predictions. With those tools, we analyzed each feature's impact on historical sea ice extent and used that information to highlight discrepancies with E3SM.

3.5.1 Data

Our machine learning (ML) models used monthly averages of June, July, and August data from the atmosphere, ocean, and sea ice to predict September sea ice extent for a given year. Results from observational and reanalysis data products are then compared against results from five ensemble members of the E3SM *historical* dataset. The features our ML models are trained on are a subset of physical quantities simulated by E3SM in the Arctic. We chose these features because they match observable features in nature and that we hypothesized would be good predictors of sea ice loss. Each feature of each dataset is a time series beginning with the start of the satellite era in 1979 and ending with the last year of available E3SM output, 2014.

The observational data included monthly sea ice extent computed from gridded, daily, passive-microwave satellite observations of sea ice concentration provided by the National Snow & Ice Data Center (NSIDC) (Peng et al., 2013). Sea ice concentration is a percentage value of ice in each grid cell, and sea ice extent (SIE) is computed as the total area of cells containing more than 15% ice. Sea ice volume reanalysis data were provided by the Pan-Arctic Ice Ocean Modeling and Assimilation System (PIOMAS) (Schweiger et al., 2011). Atmospheric data (total cloud cover percentage (CLT), downward longwave flux at surface (FLWS), pressure at the surface (PS), near-surface specific humidity (SSH), temperature at the surface (TS), wind u component/zonal (uwind), and wind v component/meridional

(vwind)) were from an atmosphere reanalysis provided by the National Centers for Environmental Prediction (NCEP) (NOAA et al., 2019a). Sea surface temperature (SST) was provided by the National Oceanic and Atmospheric Administration (NOAA) (NOAA et al., 2019b). For each of the atmospheric data variables, as well as SST, monthly Arctic area averages were computed from the global gridded fields.

We used the DOE's E3SM for climate simulation data in this work (E3SM Project, 2018; Golaz et al., 2019). E3SM version 1 was a fork of the community Earth system model (Kay et al., 2015), which was a part of the CMIP5 collection analyzed by Rosenblum and Eisenman (Rosenblum and Eisenman, 2017). E3SM is a global model comprised of submodels for land, atmosphere, land ice, sea ice, oceans, and rivers. Specifically, we used data from E3SM's *historical* ensembles 1-5 at one-degree global resolution.

E3SM published five historical ensemble runs to offer a distribution of forecasts. The runs were initialized from different years of a 500-year pre-industrial control simulation. The historical runs start in 1850, running for 165 years to 2014. The final 36 years, 1979 to 2014, were used in our analysis to match the years of observed data. Small differences in each run's initial conditions can significantly impact long-term results, though average behavior between runs is expected to be consistent.

Table 3.1 summarizes the observed features we collected; an excerpt of June values is included. Each feature is a time series of the feature's mean in a given

Table 3.1: **Training Features and June Data Excerpt**: total cloud cover percentage (CLT), downward longwave flux at surface (FLWS), pressure at the surface (PS), sea ice extent (SIE), sea ice volume (SIV), near-surface specific humidity (SSH), sea surface temperature (SST), temperature at the surface (TS), wind u component/zonal (uwind), and wind v component/meridional (vwind). Values listed are means over the pan-Arctic grid for each day of the month, rounded to two-decimal places for display only.

					June						Sept.
	CLT	FLWS	PS	SIE	SIV	SSH	SST	TS	uwind	uwind	SIE
Year	(%)	(W/m^2)	(Pa)	$(10^6 \mathrm{km}^2)$	$(10^6 \mathrm{km}^3)$	(mg/kg)	(°C)	(°C)	(m/s)	(m/s)	(10^6km^2)
1979	42.08	256.56	97930.00	12.53	29.79	4.31	0.56	273.46	0.94	0.48	5.90
1980	40.89	259.51	97901.00	12.20	29.15	4.44	0.68	274.67	0.99	0.47	6.83
1981	40.47	258.13	98098	12.43	26.82	4.27	0.65	274.27	0.06	0.06	6.40
÷	:	:	:	:	:	:	÷	÷	÷	÷	:
2012	40.36	271.60	98 105.00	10.67	16.00	5.12	1.39	277.28	-0.03	-0.06	3.55
2013	40.66	266.93	97989.00	11.36	17.54	4.98	1.26	276.50	0.93	0.42	5.27
2014	39.84	263.94	98.19	11.03	17.68	4.72	1.47	275.67	0.00	0.04	5.38

month from 1979 to 2014. Values in the time series are an area-sum over the pan-Arctic oceanic region. Each feature's monthly data is a mean of every Arctic sample in the given month, resulting in a single value per month. Generally, the observational and reanalysis datasets have similar magnitudes to the simulation data. However, for CLT, the NCEP reanalysis is significantly lower than the E3SM data. This is a known bias in the NCEP reanalysis data, and future work could investigate feature analyses of alternative reanalysis datasets (Zib et al., 2012).

The data used in this work is publicly available on the E3SM website. The five historical ensemble runs were retrieved from the v1 one-degree data CMIP6 release. To disambiguate them from our machine learning models and observed data, we will refer to E3SM's *historical ensembles 1-5* as *simulations 1-5*, simulation runs, or simply E3SM runs for the remainder of this paper. Figure 3.1 shows a comparison of the observed and simulation datasets evaluated in this work.

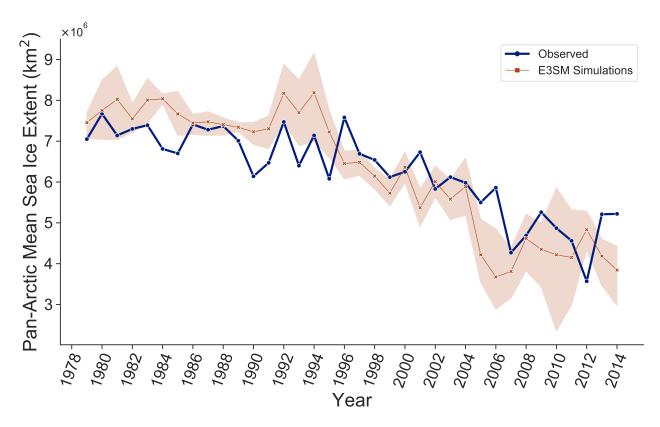


Figure 3.1: Comparison of observed, pan-Arctic mean September sea ice extent with predictions from E3SM's historical ensembles 1-5. The mean of E3SM simulations is shown with 95% confidence interval (shaded).

3.5.2 Random Forests

We found that linear models performed poorly on our data. For this work, we used RFR models because they are relatively simple, intuitive models that can learn nonlinear relationships between features. As a part of their training, the decision trees in random forests generate Gini impurity measures. These measures are aggregated after training to determine the Gini importance of each feature. In our case, we computed importance as the total reduction in mean absolute error (MAE) caused by each feature.

RFR is an ensemble learning technique, similar to a combination of bootstrap

aggregation (bagging (Breiman, 1996)) and decision tree regression. Bagging is a method to combine the knowledge of many naive estimators, or trees in our case, by providing a subset of the full sample set to each estimator. The result is the average of many noisy, but unbiased, estimators, reducing overall variance. Random forests improve the bagging method by choosing random subsets of the feature set for each node split in every tree (Banfield et al., 2007). The number of random features each node considers, and when to split are tuned hyper-parameters. The final forest's estimate is the average prediction from the random trees.

For N trees, $T_1, ..., T_N$, random forest regression prediction is computed as follows:

$$RF(N) = \frac{1}{N} \sum_{n=1}^{N} T_n(x)$$

given the training sample, x.

The random forest implementation we used was the random forest regressor from Python's sci-kit learn package (Pedregosa et al., 2011). The implementation uses a perturb and combine technique (Breiman, 1998a) made for tree regressors. Perturb and combine reduces test set error by introducing a diverse set of regressors via randomized regressor construction. For the rest of the data analysis, we used Python's Numpy package (Van Der Walt et al., 2011). We utilized Python's Seaborn package (Waskom and the seaborn development team, 2020) for data visualization.

3.5.3 Pre-Processing

To prepare the data for training, we split it into training and testing years. Our goal was not to develop predictive models for next year's sea ice extent. We were more interested in finding models that have learned the data well that we then used for feature analysis. Thus, we split the training and testing data randomly.

Because some years are easier to forecast than others, we should model every combination of training and testing years. For 36 total years and 18 testing years, we computed $\binom{36}{18} = 9075135300.00$ total combinations of training and testing years. Since it is infeasible to train that many models and evaluate each feature's importance, we used this standard method to compute a sample size:

$$\frac{(z\text{-}score)^2 \times \sigma \times (1-\sigma)}{e^2}$$

with a *z-score* computed with 95% confidence, e = 5% margin of error, and standard deviation σ , which yielded 385 sample sets on which to train and test our models. We illustrate with 18 testing years because it is the maximum value of $\binom{36}{X}, X \in [1,36]$.

Decision trees, and thus random forests, are scale-invariant (Breiman et al., 1984). This means that although our data varies greatly in scale between, for example, sea ice extent, in millions of km², and wind speeds, less than 1.00 m/s, the models' accuracy is unaffected. This is an advantage over many other ML models, and we can leave the data generally untouched. However, random forests extrap-

olate poorly for data outside of their training's minimum and maximum values (Hengl et al., 2018). This presented a problem for our analysis of the dataset because, as shown in Figure 3.1, the latter third of the data has values generally lower than any in the first two thirds. We detrended training and testing data separately to mitigate that problem by forcing the data to have a zero mean. After training and fitting our models, we retrended the data and the model's predictions to evaluate their error.

3.5.4 Model Training and Hyper-Parameter Tuning

Finally, we trained RFR models on the data the training splits provided. Note that the trees in our forests were allowed to grow until all leaves were pure, even if they contained a single sample. Decision trees are often pruned to reduce overfitting, but Breiman (Breiman, 1998b) suggests letting trees grow fully in random forests to boost accuracy and increase ensemble diversity. Banfield et al. (Banfield et al., 2003, 2007) also discuss ensemble size in random forests and conclude that many more trees are necessary than are typically used. Ensemble size is an important hyper-parameter to tune because the number of trees in the forest directly impacts the possible feature sets the forest can explore, and too many trees can reduce a random forest's performance while also sacrificing run-time. Our forests comprised 250 decision trees. The number of trees was determined empirically. Forests of size 10, 50, 100, 250, 500, and 1000 trees were evaluated and their performance was measured on the basis of the test $\overline{R^2}$ (average R^2) and average test anomaly

correlation coefficient (ACC), which are detailed in Section 3.5.6. We found that 250 tree models maximized $\overline{R^2}$ and \overline{ACC} . Lastly, the trees in each forest used mean squared error as their nodes' splitting criterion.

3.5.5 Feature Importance Measurement

We used Gini importance because of the non-linearities in climate data; in particular, Gini importance is not susceptible to data multicollinearities. Given that all of our features come from the same complex system, it would be difficult to eliminate features by simple correlation measures. In standard usage, Gini importance is normalized to compare relative importance within a single dataset. We chose to preserve the absolute importance values, letting us compare across datasets.

We also considered drop-column and permutation importance methods (Breiman, 2001). However, we found them to be unsuitable because they are highly susceptible to multicollinearity. Because many physical processes are directly acting on each other, Arctic features are inherently correlated, and any leave-one-out importance method will highlight that correlation. We found that the correlation leads these methods to attribute more importance to the least correlated feature, and it becomes difficult to glean meaningful insights.

3.5.6 Model Evaluation

We used the R^2 (coefficient of determination) from the Nash-Sutcliffe efficiency definition, given by:

$$R^{2}(\hat{y}, y) = 1 - \frac{\sum (y - \hat{y})^{2}}{\sum (y - \overline{y})^{2}},$$

where y are the true values, \hat{y} are the predicted values, and \bar{y} is the mean of y. This definition has a range of $(-\inf, 1]$ where 1 is the best possible score.

In addition to $\overline{R^2}$, we evaluated model performance with average MAE (\overline{MAE}) and \overline{ACC} . Again, average here means the mean value measured in 385 models with random training and testing year splits. Since \overline{MAE} is in millions of km², we took the Sea Ice Outlook's 2019 season report (Bhatt et al., 2020) as a baseline. This report includes several different types of data-driven models and presents one-year forecasts. These should have less error than ours, given how many more years we forecasted at once. With the exception of a few outliers between 2008 and 2019, sea ice forecast error was between -0.4 and 0.6 million km².

ACC is the Pearson's correlation coefficient (r-value) of sea ice extent anomalies. A time series' anomaly is a measure of the data's deviation from its *climatology*. In our case, the climatology is the mean value of the true values the models are attempting to forecast. This function is defined by:

$$ACC(\hat{y}, y) = \frac{\sum [(\hat{y} - \overline{y})(y - \overline{y})]}{M \times \sigma_{\hat{y}} \times \sigma_{y}}$$

where y are the true values, \hat{y} are the predicted values, M is the number of samples

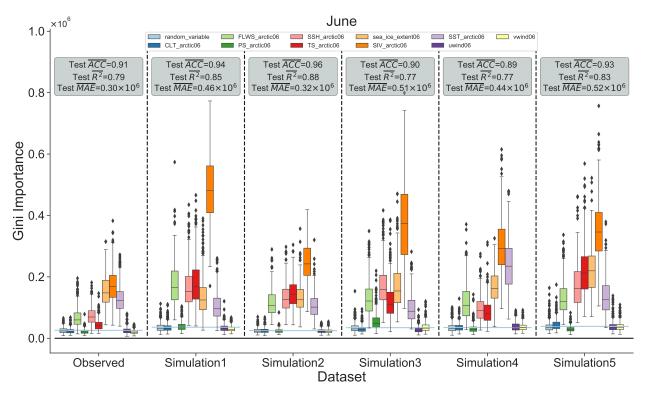


Figure 3.2: June feature importance. Standard box-and-whisker plot (McGill et al., 1978) of values for 13 predictions generated by 385 models. The average R^2 , anomaly correlation coefficient (ACC), and mean absolute error (MAE) are displayed in the gray boxes. The blue line in each dataset is the mean importance of a random variable in each feature set.

in y and \hat{y} , \bar{y} is the mean or climatology of y, $\sigma_{\hat{y}}$ is the standard deviation of the predicted values, and σ_{y} is the standard deviation of the true values.

3.6 Results

Our goal is to learn the importance of climate features on the predictions made by E3SM and compare that to the actual importance of those features on observed sea ice extent. We found that was best accomplished by training RFRs on 23 uniformly randomly chosen years and testing with the remaining 13. Our performance measure was based on the mean of $\overline{R^2}$ scores among datasets for the June

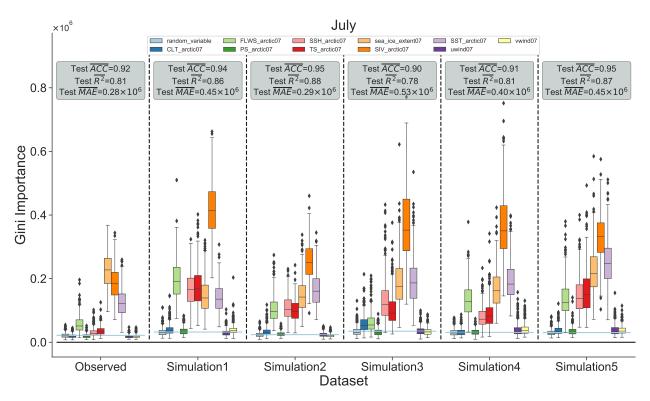


Figure 3.3: July feature importance. Standard box-and-whisker plot (McGill et al., 1978) of values for 13 predictions generated by 385 models. The average R^2 , anomaly correlation coefficient (ACC), and mean absolute error (MAE) are displayed in the gray boxes. The blue line in each dataset is the mean importance of a random variable in each feature set.

input data. This train-test-split resulted in maximum and minimum $\overline{R^2}$ scores of 0.88 and 0.77, respectively, yielding a measure of 0.83. $\overline{R^2}$ denotes the average R^2 of the 385 models.

We replicated our analysis for each month between June and August, predicting September SIE. Each subsequent month generates less error. Within each dataset, each feature's relative importance changes. Some features' importance is correlated with the progression of months, while others appear to change randomly.

Figure 3.2 shows June's feature importance values. The average train and test error values indicate that the models generally learn the data well. The blue line shows the mean feature importance of a random variable included in each model's

feature set. The random variable indicates a lower bound on importance; any feature with an importance value near this line has virtually no importance. We found that adding a random variable decreases individual model performance, but the effect is minimized when taking the mean over every model.

There are some similarities between each dataset. They share the same list of six important features, though their order and magnitudes differ. sea ice volume (SIV) is consistently the most important, though the degree of absolute importance varies. SIV, TS, SSH, SIE, FLWS, and SST are important in each dataset. The datasets, except for simulation 3, share the same list of unimportant features as well. These are CLT, PS, uwind, and vwind. One apparent exception is June's PS in Figure 3.2: simulation 3; however, excluding PS from the training data, results in a negligible difference in $\overline{R^2}$ (0.7681 vs. 0.7682).

July features, shown in Figure 3.3, predicted as well or better than June in each of our error metrics; simulation 3 had the lowest $\overline{R^2}$, 0.78, and simulation 2 had the highest, 0.88. The same features were important in July as in June, but the relative importance values changed. June's sea ice extent became more important in the observed dataset, surpassing the importance of SIV. SSH became less important in the observed dataset, too, settling just above the random variable. SSH remained as important in the simulation datasets.

The most dramatic change in importance occurs in August. These results are in Figure 3.4. Error was significantly better with simulations 3 and 4 having the minimum $\overline{R^2}$, 0.87, and simulations 1 and 2 having the maximum, 0.91. In August, sea

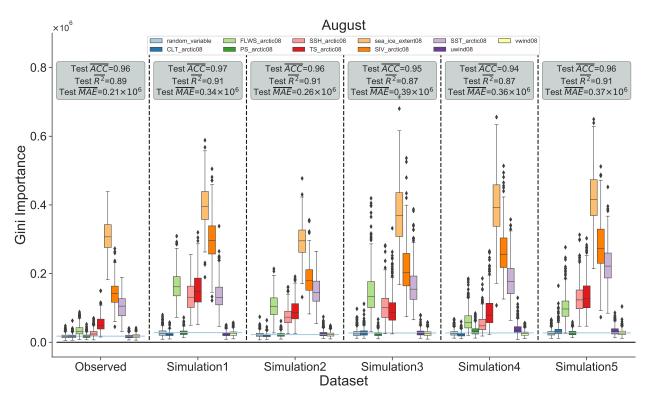


Figure 3.4: August feature importance. Standard box-and-whisker plot (McGill et al., 1978) of values for 13 predictions generated by 385 models. The average R^2 , anomaly correlation coefficient (ACC), and mean absolute error (MAE) are displayed in the gray boxes. The blue line in each dataset is the mean importance of a random variable in each feature set.

ice extent was always the most important. The importance values of the remaining features generally changed very little throughout datasets.

3.7 Discussion

We found that our RFR ML models were able to accurately learn each of the datasets. After examining the Gini importances computed within each model, we discovered some key differences in how each dataset relates to September pan-Arctic sea ice extent.

A problem with our dataset is that the satellite record only goes back to 1979.

One solution is to adapt the models to forecast sea ice extent continuously throughout each year. This is in line with Reid and Tarantino's approach (Reid and Tarantino, 2014) (see Section 3.4), but with random forests instead of support vector machines and including many features instead of only sea ice extent. The models would train on the full year of data and see 432 data points rather than 36 in the time series. Several observed features are measured more frequently than monthly, some every few hours of every day, so a means to incorporate inconsistent sampling resolutions of features should be investigated to leverage all of the data available. Another solution could be to use a surrogate model to generate more data that is similar to the first 15 years of observed data, which have a much flatter trend. The surrogate model would let the new data agree with what the model learns about input features.

The combined error metrics and general consistency of results between each dataset suggests that our models have learned the data well, and the feature analysis can identify key patterns. It is meaningful that the same six features are considered important across datasets and input-months. Since our analysis is of the pan-Arctic region, it is possible that the set of unimportant features would be more important in specific subregions of the Arctic.

Though the most important feature in June and August is consistent between simulation and observation, the absolute importance differs markedly. One clear pattern is that June shows an acute reliance on sea ice volume for both observations and simulations. By August the reliance is traded for sea ice extent. This

finding is consistent with earlier studies evaluating sea ice predictability using lagcorrelation analyses with ESM ensemble data (Ordonez et al., 2018; Blanchard-Wrigglesworth et al., 2011).

Although the observed and simulated data share patterns, there is a clear difference between them. In July, simulations and observed data do not agree on the most important feature. In June, July, and August, simulated data relies too heavily on almost all the important features. In each dataset, importance values diminish for the remaining features in June and July, and their distributions overlap more than they did in June, but the observed dataset still shows the least importance in FLWS, SSH, TS, and SST.

Interestingly, simulations 1 and 2 forecasted with the highest $\overline{R^2}$ each input month, and simulations 3 and 4 had the lowest $\overline{R^2}$ in each input month. Simulations 1 and 2 have the lowest \overline{MAE} and highest \overline{ACC} among the simulation runs, and 3 and 4 have the highest \overline{MAE} and lowest \overline{ACC} among the simulations runs. Although the differences are small, these consistencies may indicate some commonality between these simulation runs.

Our ML models performed better on the observed data than on the simulations as measured by \overline{MAE} and \overline{ACC} , but is not reflected in $\overline{R^2}$. That suggests that the mean value, or the trend after retrending, was very predictable, but its intervariability, which R^2 explains, was less predictable. The likely explanation is in the difference in the complexity of the systems. Observed features of the continuous Earth system are artificially discretized. In any complex system, intervariability is

difficult to forecast. However, because we chose largely relevant features as predictors, we could capture the macro-level patterns, as evidenced by the macro-level error measures: \overline{MAE} and \overline{ACC} .

3.8 Conclusions

We demonstrated that random forest regression and the associated Gini importance measure can provide insight into why ESMs incorrectly estimate sea ice extent in recent decades. We found a discrepancy in the feature importance between observed and simulation datasets. In particular, the discrepancy between E3SM and observation appear to be due to an over-reliance on June sea ice extent and August sea ice volume. The order of feature importance was also different between E3SM and observation, and the ordering was not consistent within E3SM ensemble members. In all cases, E3SM over-relies on six features compared to observed data. Machine learning allows us to fill the gaps in the underlying physics of ESMs, providing a metric for Stroeve et al.'s (Stroeve et al., 2012) hypothesis that ESMs are missing complex relations and causal mechanisms.

In the future, we can evaluate more features that can be measured or constructed in each dataset. An analysis, including all months of the year in each model will be elucidating as well. Sea ice extent is measured daily via satellite imagery. We can understand how each dataset explains sea ice extent at a higher resolution every month of the year.

We can repeat our analysis on other regions, including Antarctica, where there

are also problematic disagreements with observations (Rosenblum and Eisenman, 2017). An analysis like this of other climate models could be insightful too. It would be particularly interesting to compare simulations in which there few to no correlated features. That would allow for variations on the analysis, such as more modeling approaches, which require linearly independent features, and more feature analysis methods, such as drop-column importance, which would otherwise struggle with multicollinearities.

Further insight could be gained by repeating our analysis with a machine learning method other than RFR, however the following methods have their own challenges. Most neural network models would need more observed data than is available to converge. We found that multiple linear regression cannot learn the data well because the relationships between features are nonlinear. Reid and Tarantino (Reid and Tarantino, 2014) found that SVR can forecast the data well, but it is unclear what the best feature analysis method would be.

Given the discoveries in this paper, we can run experiments with E3SM to determine how reducing feature disagreements between the observed and simulation datasets impact E3SM's forecasts. That process may not yield results for several reasons, including that E3SM's real feature set is large and complex, focusing analysis on the Arctic region is too restricting to estimate the effects of the global Earth model, or our ML models are too limited by small datasets. Despite these challenges, our results can potentially guide climate modelers as they develop the next generation of ESMs.

Acknowledgments

This work is supported by Sandia Earth Science Investment Area Laboratory Directed Research and Development funding. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.

This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

4 Learning Why: Data-Driven Causal Evaluations of Climate Models

4.1 Publication Notes

Citation: Nichol, J. Jake, et al. "Learning Why: Data-Driven Causal Evaluations of Climate Models." ICML 2021 Workshop Tackling Climate Change with Machine Learning, 2021.

Publication date: 2021

Conference: International Conference on Machine Learning 2021 Workshop Tackling Climate Change with Machine Learning

Formatting: The original published text has been preserved as much as possible while still adhering to the formatting requirements of this dissertation.

Data and Software Availability: The paper is available at https://www.osti.gov/servlets/purl/1888471.

Funding: This work is supported by Sandia Earth Science Investment Area Laboratory Directed Research and Development funding. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security

Administration under contract DE-NA-0003525.

4.2 Abstract

We plan to use nascent data-driven causal discovery methods to find and compare causal relationships in observed data and climate model output. We will look at ten different features in the Arctic climate collected from public databases and from the Energy Exascale Earth System Model (E3SM). In identifying and comparing the resulting causal networks, we hope to find important differences between observed causal relationships and those in climate models. With these, climate modeling experts will be able to improve the coupling and parameterization of E3SM and other climate models.

4.3 Introduction

Climate models are critical to our understanding of climate change. We believe there is an opportunity to apply causal inference methods to these models to improve predictions. We can understand the quality of a model by comparing it with observations of the natural phenomena being simulated. From there, we can make the necessary improvements to the model, but where to start? Currently models are developed using a trial and error approach, in which a model is designed and parameterized and the resulting accuracy is observed. For computationally expensive models this approach quickly becomes inefficient. We propose to investigate the causal relationships between features and their weights to better target reparam-

eterization and feature selection efforts. We propose to focus on the pan-Arctic region because we previously studied Earth system model (ESM) prediction discrepancies there (Nichol et al., 2021). The Arctic climate, though important in itself, also has global climate implications.

In Runge et al. (2019c), a recent review of causal methods, they argue that causal discovery is well-suited to improving climate models. Nowack et al. (2020b) provide an example analysis of a global climate model. This work proposes to build these publications, by extending this nascent field to Energy Exascale Earth System Model (E3SM) (E3SM Project, 2018) and a including multiple feature analysis.

In contrast to methods based in statistical correlations, causal inference tells us why systems behave the way they do. Discovering the underlying causal structure in data and then comparing those structures from observed and simulated datasets will give us a richer understanding of the differences between the data sources.

Commonly, causal effects are determined and quantified by interventionist experiments, usually in randomized trials. Because of the magnitude, complexity, and uniqueness of the Earth's climate, there are significant feasibility and ethical problems with controlling and intervening in the climate for experimentation. For this reason, climate science is largely studied with coupled numerical models. Each model encapsulates subsystems and subprocesses that work together to determine the long-term climate.

The status-quo in Earth system model evaluation is based on simple descriptive

statistics, like mean, variance, climatologies, and spectral properties of model output derived from correlation and regression methods (Runge et al., 2019c). These methods can be simple to implement and interpret but are often ambiguous or misleading; resulting associations can be spurious and the directions of effects is fundamentally unknown.

In recent decades, a rigorous mathematical framework has been developed for observational causal inference by Pearl, Spirtes, Glymour, Scheines, and others (Spirtes et al., 1993; Pearl, 2009; Spirtes and Zhang, 2016b). The framework is largely based on Reichenbach's (Reichenbach, 1956) Common Cause Principle: that if two variables are dependent, there must be a causal relationship between the two or a third common driver of the two. Most importantly, causal methods identify the direction of observed effects between variables and detect spurious correlations.

The model we are interested in for this work is the United States Department of Energy (DOE) Energy Exascale Earth System Model (E3SM) (E3SM Project, 2018). This model is a coupling of atmospheric, ocean, river, land, land ice, and sea ice numerical models. Its goal is to use exascale computing to output high-resolution simulations of natural and anthropogenic effects in the climate.

The Arctic climate has significant direct and indirect impacts on global climate, ecology, geopolitics, and economics (Assessment, 2004; Arc, 2019; Smith and Stephenson, 2013). In particular, the volume and extent of Arctic sea ice are important indicators for the current state and projections of global climate change

(Goosse et al., 2018; Sevellec et al., 2017; Runge et al., 2015b; Cvijanovic et al., 2017). Because of this, effectively understanding the causal drivers in the Arctic climate system is requisite for understanding the future of our climate and how we can mitigate or intervene in climate change.

Climate models are in active development and the Coupled Model Intercomparison Project (CMIP) is a group that collects and curates modern climate models for world-wide collaboration. Researchers have found that models in phases 3 and 5 of CMIP underestimate the rate of Arctic sea ice loss on average (Rosenblum and Eisenman, 2017; Taylor Karl E., Stouffer Ronald J., 2012; Stroeve et al., 2007). Figure 4.1 shows the difference between observed sea ice extent and E3SM's modeled prediction.

In previous work, we used random forest feature analysis to determine which summer-time features in the Arctic are most predictive of yearly sea ice extent minimums in September (Nichol et al., 2021). We then compared results from observed data and simulation output data. This approach allowed us to discover and compare nonlinear relationships in the climate systems. Random forest feature importance values are correlations and direction can only be inferred from each feature to the single predictand. Therefore, inter-feature relationships in the model cannot be interpreted causally. Finding differences between in causal relationships between climate models and observed data will identify clear, actionable problems with the models.

4.4 Data

We selected time series data for ten features in the Arctic consisting of monthly mean values for each year of available data. Empirical data was collected from observational and reanalysis data products, and simulated data were taken from five ensemble members of the E3SM *historical* dataset (E3SM Project, 2018; Golaz et al., 2019). The selected features are a subset of physical quantities simulated by E3SM in the Arctic and are the same ones used in our previous work with random forests, (Nichol et al., 2021). We originally chose these features because they match observable features in nature and we hypothesized they would be good predictors of sea ice loss. Through feature analysis, we discovered that some inputs were far more predictive than others, but we did not have a causal inference framework to explain why. Each feature of the observed dataset is a time series beginning with the start of the satellite era in 1979 to 2018. The E3SM *historical ensembles* span 1850 to 2014.

The observational data includes monthly sea ice extent computed from gridded, daily, passive-microwave satellite observations of sea ice concentration provided by the National Snow & Ice Data Center (NSIDC) (Peng et al., 2013). Sea ice concentration is a percentage value of ice in each grid cell, and sea ice extent (SIE) is computed as the total area of cells containing more than 15% ice. Sea ice volume (SIV) reanalysis data were provided by the Pan-Arctic Ice Ocean Modeling and Assimilation System (PIOMAS) (Schweiger et al., 2011). Atmospheric data, total

cloud cover percentage (CLT), downward longwave flux at surface (FLWS), pressure at the surface (PS), near-surface specific humidity (SSH), temperature at the surface (TS), wind u component/zonal (uwind), and wind v component/meridional (vwind)) were from an atmosphere reanalysis provided by the National Centers for Environmental Prediction (NCEP) (NOAA et al., 2019a). Sea surface temperature (SST) was provided by the National Oceanic and Atmospheric Administration (NOAA) (NOAA et al., 2019b). For each of the atmospheric data variables, as well as SST, monthly Arctic area averages were computed from the global gridded fields. Simulated data features were selected to match the observation dataset.

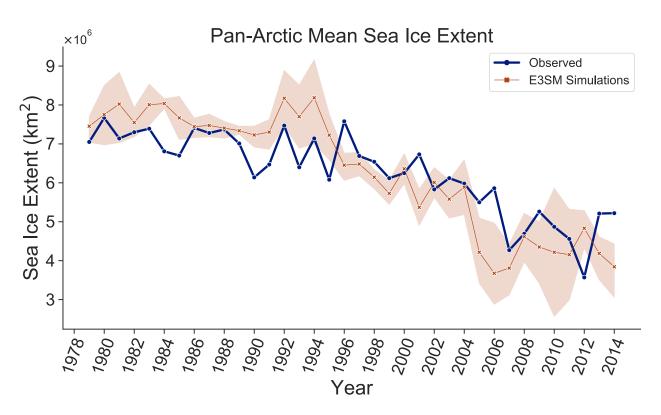


Figure 4.1: Comparison of observed, pan-Arctic mean September sea ice extent with predictions from E3SM's historical ensembles 1-5. The mean of E3SM simulations is shown with 95% confidence interval (shaded).

Figure 4.1 shows the difference between observed and E3SM's simulated sea ice extent in September each year between 1979 and 2014. September is when sea ice extent is at its minimum. The model generally predicts the same trend but fails to determine critical lows in yearly sea ice extent. While the simulations generally predict sea ice extent well, there are significant departures (fall outside the 95% CI) in particular years. For example, in 2012 there was a reversal between simulation, which predicted a year-over-year increase in sea ice, but instead a record low was observed. Since sea ice extent has a non-linear effect on the global climate, providing a causal explanation for these departures is critical.

4.5 Approach

Causal inference is a mathematical framework for answering questions about why phenomena occur. Causal modeling is an effort to discover, describe, and analyze the relationships between cause and effect (Pearl, 2009; Spirtes and Zhang, 2016b). The calculus of causation is defined in two languages: a causal diagram, expressing what we know, and a symbolic language, expressing what we want to know (Pearl and Mackenzie, 2018b). The methods we propose derive a causal diagram from the given data.

A causal diagram is a directed graph where arcs represent the causal relationships between variables. Figure 4.2 is a diagram depicting correlations between variables in the observed dataset from our previous work. Only mean values from June in each year between 1979 and 2014 were included. For example, the PC

algorithm (Spirtes et al., 1993) could take a diagram such as the one in Figure 4.2 as input and iteratively remove spurious correlations and determine the causal direction between the remaining links.

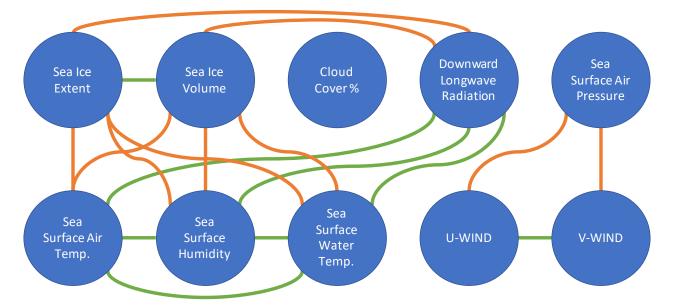


Figure 4.2: Diagram showing correlated relationships between variables in June from the observed dataset between 1979 to 2014. Green indicates a positive correlation and orange indicates a negative correlation. The correlation threshold is ± 0.6 .

There are multiple methods for constructing causal networks that are candidates for investigation in this work. These include causal network learning algorithms, such as the Peter-Clark (PC) algorithm, structural causal model frameworks, such as LiNGAM, and the fast causal inference (FCI) algorithm. Each of these require sets of assumptions about the given data describing the system. We will need to determine which assumptions we can meet with the available data. Due to the nonlinear, stochastic, high-dimensional nature of the climate system, it is likely that causal network learning algorithms and structural causal models will be more effective.

4.5.1 The PCMCI method

We plan to attempt our analysis with PCMCI (Runge et al., 2019d) first. PCMCI extends the PC-algorithm by adding momentary conditional independence (MCI) tests. These remove false-positives left by the PC algorithm and conditions on each variable's causal parent and its time-shifted parents as well. Thus, the algorithm is designed to remove spurious relationships and identify concurrent and time-lagged causal relationships. PCMCI was specifically designed for highly interdependent time series such as climate data.

In (Nowack et al., 2020b), the authors used time series data for sea level pressure data collected at 50 locations around the globe. The authors then examined the relationship between precipitation and the causal network skill scores for sea level pressure to demonstrate that this method can help identify dynamic coupling mechanisms arising from underlying physical processes. The Nowack et al. study is one of the first causal network inference studies using large-scale spatiotemporal data and provides a proof-of-concept that such methods are viable for analyzing climate systems. They looked at a single variable in various regions. In contrast, we plan to use PCMCI to analyze several different quantities in the same region.

4.5.2 Comparing and evaluating causal models

An obvious first approach for comparing causal diagrams is with standard graph comparison metrics such as global properties and summary statistics: edge density, global clustering coefficient, degree distribution, counts of subgraphs, hamming distance, etc. However, these are defined by correlation and do not address the causal nature of the networks.

Other metrics grounded in information theory, such as information flow, are more appropriate for causal networks but possibly more difficult to interpret holistically. In (Runge, 2015), the authors present a framework for determining information flow from multivariate causal diagrams.

A different approach is to consider the resulting models' performance. This includes metrics such as true positive rate (TP), false positive rate (FP), accuracy, positive predictive value, false omission rate, the S-score, and the G-measure and F1-score (metrics combining TP and FP). These require a baseline model, such as the causal diagram of the observed dataset, to measure the performance of a test model. These are easier to interpret than information flow but are relative measures and cannot be assessed independently.

4.6 Anticipated Contributions

The contributions of this work will bring climate modeling experts a step closer to understanding *why* E3SM does not model certain Arctic quantities well, such as sea ice extent. In our previous work, random forests were able to elucidate which features were more or less important for model predictability in observed and E3SM data. This work should support those results and help explain the causal drivers behind observed and E3SM results. Future research after this work could include: considering more features in the Arctic; other regions with known mod-

eling biases, such as the Antarctic; and other climate modeling problems, such as determining the effects and sources of major climate events. Clear examples are volcanic eruptions and anthropogenic climate change and intervention. Developing more informative analytics for climate models will hasten their improvement and better inform policy decisions to mitigate and combat global climate change.

5 Causal Evaluations for Identifying Differences between Observations and Earth System Models

5.1 Publication Notes

Citation: Nichol, J. Jake, et al. "Causal Evaluations for Identifying Differences between Observations and Earth System Models." [Report No. 1820528]. 2021. U.S. Department of Energy, Office of Scientific and Technical Information.

Publication date: 2021

Publisher: U.S. Department of Energy, Office of Scientific and Technical Information.

Formatting: The original published text has been preserved as published.

Data and Software Availability: The paper is available at https://www.osti.gov/biblio/1820528.

Funding: This work is supported by Sandia Earth Science Investment Area Laboratory Directed Research and Development funding. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell

International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.

SAND2021-11449 R

LDRD PROJECT NUMBER: 224486

LDRD PROJECT TITLE: Causal Evaluations for Identifying Differences

between Observations and Earth System Models

PROJECT TEAM MEMBERS: PI: Matt Peterson (1461), PM: Susan Altman

(8140), Jake Nichol (1461), Kara Peterson (1442)

ABSTRACT:

We use a nascent data-driven causal discovery method to find and compare causal relationships in observed data and climate model output. We consider ten different features in the Arctic climate collected from public databases on observational and Energy Exascale Earth System Model (E3SM) data. In identifying and analyzing the resulting causal networks, we make meaningful comparisons between observed and climate model interdependencies. This work demonstrates our ability to apply the PCMCI causal discovery algorithm to Arctic climate data, that there are noticeable similarities between observed and simulated Arctic climate dynamics, and that further work is needed to identify specific areas for improvement to better align models with natural observations.

INTRODUCTION AND EXECUTIVE SUMMARY OF RESULTS:

The Arctic is changing rapidly and feedbacks between the ocean, atmosphere, and sea ice may be accelerating that change [12]. Accurate predictions of the future sea ice extent in the Arctic depend on understanding the impacts of greenhouse gas forcing and the superimposed internal variability of the complex Earth system. In particular, sea ice loss in the Arctic has been shown to have a linear relationship with global average surface temperature in both observational data and simulation data, with most predictions indicating that the Arctic will be seasonally ice free by mid-century [12,13]. The correlation is generally explained by a common dependency of temperature and sea ice concentration on greenhouse gas concentration, but causality has not typically been assessed. Other studies have found that internal variability in the climate system can accelerate or impede sea ice loss and there is currently no consensus on the dominant processes in the ocean and atmosphere that have the largest impact [14, 15, 16].

Earth system models (ESMs) are critical to our understanding of climate change, but the complex nature of the interactions between atmosphere, ocean, ice, and land can obscure causal relationships. Here, we investigate the causal relationships between Arctic climate features to better understand the complex feedbacks that result in rapid Arctic change and sea ice loss. This effort extends our feature analysis that identified features important for predicting yearly minimum sea ice concentration and compared feature importance between simulations and observations [1].

In [2], a recent review of causal discovery methods for complex systems, they argue that causal discovery is well-suited to improving climate models. In [3], authors provide an example analysis of a global climate model, though focus on a single feature in many separate regions of the globe. This work Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.









builds on these publications by extending this nascent field to the U.S. Department of Energy's Energy Exascale Earth System Model (E3SM) [4] and including a multiple feature analysis within one common region. E3SM is a coupling of atmospheric, ocean, river, land, land ice, and sea ice numerical models. Its stated goal is to use exascale computing to output high-resolution simulations of natural and anthropogenic effects in the climate.

Commonly, causality is determined and quantified by interventionist experiments, usually in randomized trials. Because of the magnitude, complexity, and uniqueness of the Earth's climate, there are significant feasibility and ethical problems with controlling and intervening in the climate for experimentation. For this reason, climate science is largely studied with ESMs, which are coupled numerical models. Each model encapsulates subsystems and subprocesses coupled together to approximate the long-term climate.

The status-quo in ESM evaluation is based on descriptive statistics, like mean, variance, climatologies, and spectral properties of model output derived from correlation and regression methods [2]. These methods can be simple to implement and interpret but are often ambiguous or misleading; resulting associations can be spurious and the directions of effects is fundamentally unknown.

In recent decades, a rigorous mathematical framework has been developed for observational causal inference by Spirtes, Glymour, Scheines, Pearl, Rubin, and others [5, 6, 7, 8]. The framework for causal discovery is largely based on Reichenbach's [9] Common Cause Principle: that if two variables are statistically dependent, there must be a causal relationship between the two, or a third common driver of the two. Most importantly, causal discovery methods attempt to identify the direction of observed effects between variables and detect spurious correlations. Effectively understanding the causal drivers in the Arctic climate system is requisite for understanding the future of our climate and how we can mitigate or intervene in climate change.

In previous work, we used a random forest feature analysis to determine which summertime features in the Arctic are most predictive of yearly sea ice extent minimums in September [1]. We then compared results from observed data and simulation output data. This approach allowed us to discover and compare nonlinear relationships in the climate systems. Random forest feature importance values are correlations and direction can only be inferred from each feature to the single predictand. Therefore, inter-feature relationships in the model cannot be interpreted causally. This research expands on our previous work by identifying causal relationships in the data and comparing causal networks from historical simulations and observations.

Causal discovery of observational data is notoriously difficult because spurious correlations and incomplete data leads to spurious inferences. In this work we use conditional independence-based causal discovery, which relies on several assumptions for estimating causal links. One of which is causal sufficiency, that all confounding variables are observed. Because the complex dynamics of the Arctic system are actively researched, and there is no strong consensus on the dominant processes in the Arctic climate, we cannot validate causal sufficiency. We chose our variable set because of their strong

correlation with sea ice extent and their success in predicting sea ice extent [17, 18, 1], and they serve as a good hypothesis for a sufficient set.

In our analysis, we were able to fit a network depicting conditional dependencies between features to each of six data sets, observed and five simulated. We then applied a similarity score to evaluate how well the simulated datasets agree with the observed data and each other. Finally, we discuss the next steps for this work and how to derive meaningful differences between the networks.

DETAILED DESCRIPTION OF RESEARCH AND DEVELOPMENT AND METHODOLOGY:

Data

We collected ten features of the Arctic climate. Each was a timeseries of monthly mean values, averaged spatially over the region above 60 degrees North latitude. The observed dataset consisted of natural observations and output from reanalysis products. Simulated data was from the five members, or runs, of the E3SM *historical* ensemble [4]. The *historical* ensemble is a set of runs simulating the Earth system from 1850 to 2014. These runs were initialized by a 500-year-long pre-industrial control simulations, named *piControl*. The selected features are a subset of the quantities E3SM models and were chosen to match observable natural quantities and have been shown in previous work to have strong correlations with sea ice extent [17, 18, 1]. Resulting are six separate datasets, one observational and five E3SM simulation datasets.

The specific quantities we used were mostly the same as outlined in our plan (as seen in Addendum A). We did choose to change a few details. Rather than limit each variable to the same temporal range, 1979-2014, we instead included all the data available for each. We used the entire 150-year span of the E3SM data. The observational timeseries' date range varied by each feature, though they all start in 1979 and continue at least through 2017. Additionally, the full 150-year surface zonal and meridional wind timeseries were not readily available, so we opted to use surface wind magnitude, SWind, in their place, which does not include a directional component. Lastly, we included monthly precipitation rate data from E3SM and from the National Centers for Environmental Prediction for the observational dataset. Full data details are in Addendum C.

Preprocessing

The method detailed below, PCMCI, assumes the data is statistically stationary, i.e., its summary statistics do not change in time. First, we tested each timeseries for stationarity. This consisted of using the Kwiatkowski–Phillips–Schmidt–Shin (KPSS) and augmented Dickey-Fuller (ADF) hypothesis tests. KPSS tests the null hypothesis that a timeseries is stationary around a deterministic trend while ADF tests the null hypothesis that a timeseries is nonstationary around a deterministic trend. If KPSS fails to reject the null hypothesis, and the ADF test rejects, then we considered a timeseries stationary. We used an alpha value of 0.05 to determine significance and found that most features were nonstationary. To



keep dependencies and inferences consistent, we applied a 12-month differencing transform to every timeseries. A 12-month difference transform is the process of subtracting a timeseries by itself lagged 12 months. Resulting is a timeseries of the original's change from one year to the next. Differencing removes trend and choosing 12-months will remove yearly seasonality in the data.

Causal network learning

Causal discovery is the process of reconstructing the causal structure from purely observational data [10]. Traditional causality research to determine the causal effect, inferences about the strength of effects between variables, is done when the causal structure is already known. Causal discovery is used when the causal structure is mostly unknown. The causal structure discovered is often represented as a directed acyclic graph in which the nodes represent observed variables, and the edges represent causal relationships.

Causal discovery generally makes four major assumptions: (1) the causal Markov assumption, that if two nodes, X and Y, are d-separated in a graph G, given a conditioning set Z, then X and Y are conditionally independent in their joint probability distribution, given Z; (2) the faithfulness assumption, that if two variables, X and Y, are conditionally independent, given a set of variables, Z, then their nodes in a graph, G, must be d-separated, given Z; (3) causal sufficiency, that there are not any unobserved confounding variables of any variables in the graph; and (4) acyclicity, that there are no cycles in the graph.

In this work, we applied the PCMCI algorithm [11]. PCMCI is an extension to the PC causal network learning algorithm [5], named for its authors Peter Spirtes and Clark Glymour. PC is known for a relatively high false positive rate and struggles with high dimensional, autocorrelated data [11]. In [11], Runge et al. adapted PC to use its skeleton discovery phase for condition selection and then utilize a momentary conditional independence (MCI) phase. PCMCI estimates the causal links between all variable pairs, including their temporal lags.

The first important determination in applying PCMCI is to choose a conditional independence test. The authors have implemented three, the partial correlation, a linear parametric test, gaussian process regression and distance correlation, a nonlinear parametric test, and conditional mutual information with a k-nearest-neighbors estimator, a nonlinear nonparametric test. Generally, the functional form of the dependencies in the feature set needs to be assumed and the appropriate test is chosen. In our case though, we knew it was likely that nonlinear dependencies existed in the data but could not assume if they remained after the data was transformed.

To estimate the dependencies' functional form, we plotted each feature with another one in a scatter plot. The resulting plot depicts how each feature varies with the other. With this, linearities and nonlinearities can be found by eye. Applying this process to the untransformed data, we indeed found several nonlinearities of various forms as well as linear dependencies. Applying it to the transformed data revealed no clear nonlinearities, and multiple clearly linear relationships. With this discovery, we selected the partial correlation parametric linear conditional independence test.

PCMCI has two primary hyperparameters for tuning. The first is the maximum lag, τ_{max} , the maximum lag to evaluate for each variable. τ_{max} is an estimate of the maximum time that every variable may have an effect on the others. The estimation of τ_{max} may come from prior knowledge or by analyzing the linear dependence of each variable with every other variable at a range of lags. The second parameter to estimate is the alpha significance threshold for edges in the graph. Every pairwise dependence is determined with conditional independence tests and has an associated p-value for its significance. Alpha is the threshold for whether the p-value of each link is small enough to be included in the final graph.

To estimate τ_{max} , we plotted the cross-dependencies between each variable at lags between 0 and 24 months and looked for dependence to reach zero for every graph. See Figure 1 for an example from the observed dataset. We repeated this process for each dataset and found that $\tau_{max}=12$ months was adequate for each variable pair. To estimate alpha, we followed the procedure in [3], which selects from the list $\{0.0001, 0.001, 0.01, 0.05, 0.1, 0.2, 0.3, 0.4, 0.5\}$ by computing the Akaike information criterion (AIC) of the models fit by each value in the list. That list is slightly more extensive than in [3] because we found each graph was selecting 0.05 and wanted to be sure it was not just selecting the smallest available value.

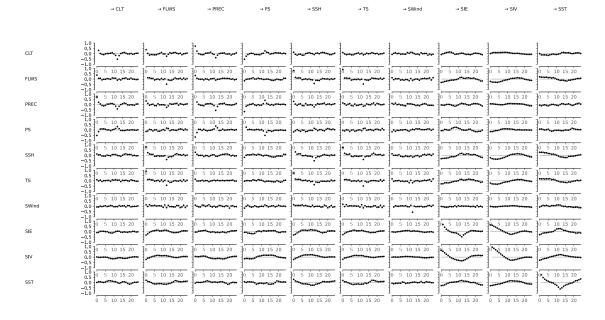


Figure 1: Plots of each feature as a function of each other feature's lags. The vertical axes denote linear dependence, and the horizontal axes denote the number of lags in months.

Causal network comparison

We utilized the F_1 score used in [3] to compare each pair of graphs. The F_1 score is a graph similarity metric with bounds [0,1], with 0 indicating no similarity and 1 indicating perfect similarity. The metric is computed from the precision, P, and recall, R, of a graph in comparison to a reference graph. Precisely, these values are computed as:

$$F_1 = \frac{2 \cdot P \cdot R}{P + R}$$

where

$$P = \frac{TP}{TP + FP}$$
$$R = \frac{TP}{TP + FN}$$

and TP is the number of true positives, FP is the number of false positives, and FN is the number of false negatives. These terms often assume a ground truth, although because the observed graph is an estimated causal graph and not ground truth, it is important to consider this metric as a relative score and not absolute.

RESULTS AND DISCUSSION:

Before analyzing the results, we filtered links from each network with less than 0.001 significance. For each dataset, PCMCI independently selected pc-alpha value to be 0.05 via AIC. PCMCI evaluated lags between 0 and 12 months for each feature. The simplified graphs in Figure 2 and Figure 3 hide the nodes of each features' lags and only presents a single node per variable. The full timeseries graphs inferred by PCMCI include nodes for each feature's lags up to the maximum lag of 12 months. Because the date ranges on simulated and observed data are not the same, we present results from networks learned from the fully available date ranges, as well as from a homogenous range, 1979 to 2014. Although the algorithm has less data to learn from, this may be a fairer comparison to observed dynamics.

Simplified graphs label links with a list of the lags with significant dependency in order of magnitude. Node color depicts a feature's auto-dependency, how dependent a feature is on its lags. Edge color depicts cross dependency, how dependent a feature is on another feature. Negative, or blue, cross dependency indicates that as the parent's value increases or decreases, the child's value changes inversely. Positive dependence indicates parent and child values increase and decrease together. Because we used a linear conditional independence test, these relationships are linear. Since these colors span many lags, the color chosen for the simplified graphs is the maximum absolute link between two features or a feature and itself.

Figure 2 is a simplified causal network estimation, trained from the full range of observed data. Resulting is relatively sparse partially directed acyclic graphs, with only 5.3% of all possible links existing

in the graph. Directed links represent discovered dependencies between features. Undirected links represent contemporaneous dependencies.

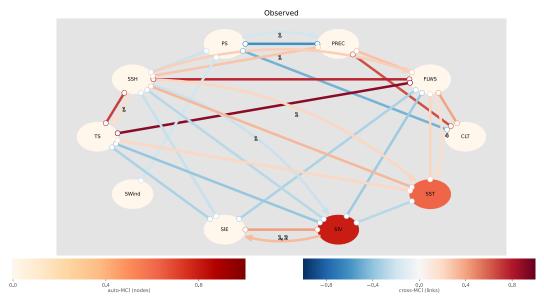


Figure 2: Simplified graph resulting from applying PCMCI with the partial correlation test on observational data in the fully available date range. The pc-alpha parameter was selected by AIC to be 0.05, the links are defined by a significance threshold of 0.001.

Figure 3 is the simplified graph fit by simulation 1 of the E3SM *historical* ensemble in the fully available date range. Although many similar links exist in this graph, it contains many more than the observed data graph. The remaining simulation graphs can be found in the Addendum C. They all differ but are more alike than the observed data graph and contain more links. An average of 8.6% of all possible links exist in the simulation graphs.

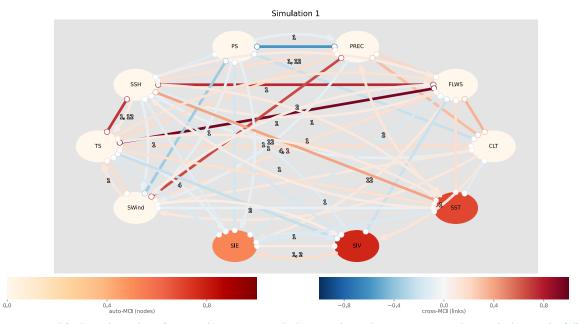


Figure 3: Simplified graph resulting from applying PCMCI with the partial correlation test on simulation 1's data in the fully available date range. The pc-alpha parameter was selected by AIC to be 0.05, the links are defined by a significance threshold of 0.001.

To better quantify the similarity between each graph, we computed the F_1 score of each pair of graphs. For this analysis, we included the fully detailed networks. These include a node for each lag of each feature. Figure 4 shows these results for the fully available date range graphs. The simulation networks are the most similar with each other, while the observed network is the most different from all other networks. The average simulation to simulation F_1 score is 0.83. The average simulation to observed F_1 score is 0.7.

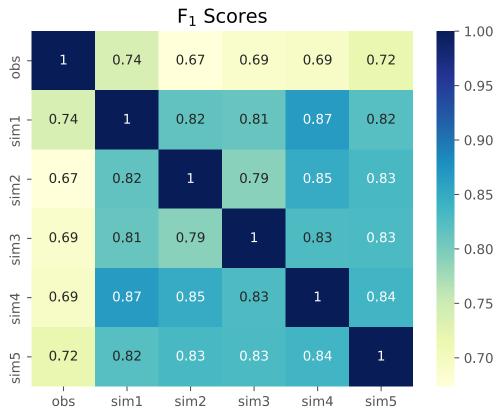


Figure 4: Matrix of F_1 similarity scores of each pair of graphs for the fully available date range graphs.

The homogenous date range changed the observed graph minimally but altered the simulated graphs noticeably. 5.5% of all possible links exist in the observed graph, while an average of only 4.8% exist in the simulation graphs for this date range. Figure 5 shows the F_1 similarity scores for the homogenous date range, 1979-2014. In this, the simulation networks lose some similarity, dropping to an average value of 0.71 simulation to simulation. The average similarity to the observed network increases slightly though, to 0.73. It is intuitive that the simulations would diverge in later years, after having been initialized equivalently, and eliminating the early years makes this apparent in their similarity scores.

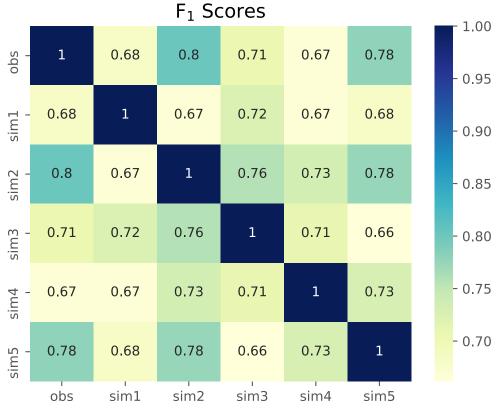


Figure 5: Matrix of F_1 similarity scores of each pair of graphs for the homogenous date range, 1979-2014.

In this work, we developed a strong foundation for applying conditional independence-based causal discovery algorithms. The differencing transforms we applied to the data were important for removing seasonality and trend, which removes the unobserved confounders driving them. We have found that we can apply causal discovery algorithms to Arctic climate data and find strong consistencies between observed and simulated timeseries. Although we cannot validate the causal sufficiency assumption with certainty, we can see that discovered conditional dependencies are similar in each dataset. In future work, we can develop and apply node-to-node similarity metrics to find which nodes are most responsibility for dissimilarity between graphs.

It is important to remember that each feature was transformed to create stationary timeseries. The 12-month differencing transform means that each timeseries is a series each month's deltas from that month's previous year. This means that a directed link from feature X to feature Y would be interpreted as the change in Y from year to year is dependent on the change in X from year to year.

The primary limitation of our findings is the inability to justify the causal sufficiency assumption. The remaining assumptions can be considered satisfied as they assume that an underlying causal structure exists in the data, and that cause and effect does not occur instantaneously. That is assured by the

physical and temporal nature of these quantities. The challenge of causal sufficiency exists in any open complex system. We plan to apply causal discovery algorithms that do not rely on the causal sufficiency assumption, such as the Fast Causal Inference algorithm [5] or Latent PCMCI (LPCMCI) [19]. LPCMCI augments PCMCI to discover causal links in the presence of latent, or unobserved, features.

ANTICIPATED OUTCOMES AND IMPACTS:

During this project we presented our findings to the International Conference on Machine Learning (ICML) in the form of a workshop paper (as seen in Addendum A) and an online poster presentation. We also gave another presentation internally to the Validation and Verification of Machine Learning Models discussion group (as discussed in Addendum B). Later this fall there will be presentation at the Chesapeake Large-Scale Analytics Conference (CLSAC) about this work. These presentations allowed us to network with other groups around the labs and externally; organizations include 5493, 1463, 0515, and professors at the University of New Mexico (as discussed in Addendum B).

The major lesson learned in this project was that ground truth for artic climate dynamics is an ongoing research problem, which this work depends on for validating our results are causal. Currently we are relying heavily on climate experts to validate our causal models, but to fully develop metrics for comparing our models we need a concrete understanding of arctic climate dynamics as well as global dynamics. Once these climate dynamics are sufficiently validated, we can utilize these causal models to help us improve our simulated models.

This work will continue in the *CLimate impact: Determining Etiology thRough pAthways (CLDERA) Grand Challenge* project starting in FY22. We plan on improving and adding metrics for comparing similarities and differences between causal models. We are also looking into determining how well a given model fits the data used for training. Some other research areas we want to explore include incorporating spatial data features into our analysis. The work done in this project used averaged values over the entire arctic. We could have divided the data into subregions of the arctic, but with this being a Late-Start LDRD with limited time we decided it was best to simplify the problem space. This will be important for CLDERA because we will be working with data on a global scale and averaging values over the whole globe would not work as easily.

CONCLUSION: (400 word limit)

In this work, we found strong similarities between conditional dependencies discovered in observed and simulated climate dynamics. If the assumptions of causal discovery were to hold, we would find that E3SM climate simulation runs are causally similar to each other and, importantly, causally similar to observations. Although we cannot validate the causal sufficiency assumption, there is evidence that our feature set is a good hypothesis. The largest remaining sources of confounding may be from remaining seasonality and trend from external forcing such as periodic-natural and anthropogenic climate changes. A clear next step is to apply a causal discovery algorithm that does not require causal sufficiency and then compare results.

REFERENCES:

- [1] J. Jake Nichol, Matthew G. Peterson, Kara J. Peterson, G. Matthew Fricke, and Melanie E. Moses. 2021. Machine learning feature analysis illuminates disparity between E3SM climate models and observed climate change. *J Comput Appl Math* 395, (October 2021), 113451. DOI:https://doi.org/10.1016/j.cam.2021.113451
- [2] Jakob Runge, Sebastian Bathiany, Erik Bollt, Gustau Camps-Valls, Dim Coumou, Ethan Deyle, Clark Glymour, Marlene Kretschmer, Miguel D. Mahecha, Jordi Munoz-Mari, Egbert H. van Nes, Jonas Peters, Rick Quax, Markus Reichstein, Marten Scheffer, Bernhard Scholkopf, Peter Spirtes, George Sugihara, Jie Sun, Kun Zhang, and Jakob Zscheischler. 2019. Inferring causation from time series in Earth system sciences. *Nat Commun* 10, 1 (2019). DOI:https://doi.org/10.1038/s41467-019-10105-3
- [3] Peer Nowack, Jakob Runge, Veronika Eyring, and Joanna D. Haigh. 2020. Causal networks for climate model evaluation and constrained projections. *Nat Commun* 11, 1 (2020), 1--11. DOI:https://doi.org/10.1038/s41467-020-15195-y
- [4] Jean-christophe Golaz, Peter M Caldwell, Luke P Van Roekel, Mark R Petersen, Qi Tang, Jonathan D Wolfe, Guta Abeshu, Valentine Anantharaj, Xylar S Asay-davis, David C Bader, Sterling A Baldwin, Gautam Bisht, Peter A Bogenschutz, Marcia Branstetter, Michael A Brunke, Steven R Brus, Susannah M Burrows, Philip J Cameron-smith, Aaron S Donahue, Michael Deakin, Richard C Easter, Katherine J Evans, Yan Feng, Mark Flanner, James G Foucar, Jeremy G Fyke, Elizabeth C Hunke, Robert L Jacob, Douglas W Jacobsen, Nicole Jeffery, Philip W Jones, Noel D Keen, Stephen A Klein, Vincent E Larson, L Ruby Leung, Hong-yi Li, Wuyin Lin, William H Lipscomb, Po-lun Ma, Renata B Mccoy, Richard B Neale, Stephen F Price, Yun Qian, Philip J Rasch, J E Jack Reeves Eyre, William J Riley, Todd D Ringler, Andrew F Roberts, Erika L Roesler, Andrew G Salinger, Zeshawn Shaheen, Xiaoying Shi, Balwinder Singh, Milena Veneziani, Hui Wan, Hailong Wang, Shanlin Wang, and Dean N Williams. 2019. The DOE E3SM Coupled Model Version 1: Overview and Evaluation at Standard Resolution. *J Adv Model Earth Sy* (March 2019). DOI:https://doi.org/10.1029/2018ms001603
- [5] Peter Spirtes, Clark N Glymour, Richard Scheines, and David Heckerman. n.d. *Causation, prediction, and search*.
- [6] Judea Pearl. 2009. Causal inference in statistics: An overview. *Statistics Surv* 3, September (2009), 96--146. DOI:https://doi.org/10.1214/09-ss057

- [7] Peter Spirtes and Kun Zhang. 2016. Causal discovery and inference: concepts and recent methodological advances. *Applied Informatics* 3, 1 (2016), 3. DOI:https://doi.org/10.1186/s40535-016-0018-x
- [8] Donald B. Rubin. 1974. Estimating causal effects of treatments in randomized and nonrandomized studies. *J Educ Psychol* 66, 5 (1974), 688–701. DOI:https://doi.org/10.1037/h0037350
- [9] Hans Reichenbach. n.d. *The direction of time*. Univ of California Press.
- [10] Clark Glymour, Kun Zhang, and Peter Spirtes. 2019. Review of Causal Discovery Methods Based on Graphical Models. *Frontiers Genetics* 10, (2019), 524. DOI:https://doi.org/10.3389/fgene.2019.00524
- [11] Jakob Runge, Peer Nowack, Marlene Kretschmer, Seth Flaxman, and Dino Sejdinovic. 2019. *Detecting and quantifying causal associations in large nonlinear time series datasets*. Retrieved from http://advances.sciencemag.org/
- [12] Dirk Notz and Julienne Stroeve. The trajectory towards a seasonally ice-free Arctic ocean. *Current Climate Change Reports*, 4:407–416, 2018.
- [13] Julienne Stroeve and Dirk Notz. Changing state of Arctic sea ice across all seasons. *Environ-mental Research Letters*, 13:103001, 2018.
- [14] Dirk Olonscheck, Thorsten Mauritsen, and Dirk Notz. Arctic sea-ice variability is primarily driven by atmospheric temperature fluctuations. *Nature geoscience*, 12:430–434, 2019.
- [15] Ana C. Ordonez, Cecilia M. Bitz, and Edward Blanchard-Wrigglesworth. Processes controlling Arctic and Antarctic sea ice predictability in the Community Earth System Model. *Journal of Climate*, 31:9771–9786, 2018.
- [16] Qinghua Ding, Axel Schweiger, Michelle L'Heureux, Eric J. Steig, David S. Battisti, Nathaniel C. Johnson, Eduardo Blanchard-Wrigglesworth, Stephen Po-Chedley, Qin Zhang, Kirstin Harnos, Mitchell Bushuk, Bradley Markle, and Ian Baxter. Fingerprints of internal drivers of Arctic sea ice loss in observations and model simulations. *Nature Geoscience*, 12(1):28–33, 2019.



[17] Jiwon Kim, Kwangjin Kim, Jaeil Cho, Yong Q. Kang, Hong-Joo Yoon, and Yang-Won Lee. Satellite-based prediction of Arctic sea ice concentration using a deep neural network with multi-model ensemble. *Remote Sensing*, 11(19), 2019.

- [18] Monica Ionita, Klaus Grosfeld, Patrick Scholz, Renate Treffeisen, and Gerrit Lohmann. September Arctic sea ice minimum prediction a new skillful statistical approach. *Earth System Dynamics*, 10:189–203, 2019.
- [19] Andreas Gerhardus and Jakob Runge. n.d. High-recall causal discovery for autocorrelated time series with latent confounders. In *Advances in Neural Information Processing Systems*, Curran Associates, Inc., 12615–12625. DOI:https://doi.org/10.5194/egusphere-egu21-8259

Part II

Local Causal Discovery in High-Dimensional Gridded Data

6 Benchmarking the PCMCI Causal Discovery Algorithm for Spatiotemporal Systems

6.1 Publication Notes

Citation: Nichol, J. Jake, et al. "Benchmarking the PCMCI Causal Discovery Algorithm for Spatiotemporal Systems." [Report No. 1991387]. 2023. U.S. Department of Energy, Office of Scientific and Technical Information.

Publication date: 2023

Publisher: U.S. Department of Energy, Office of Scientific and Technical Information.

Formatting: The original published text has been preserved as published.

Data and Software Availability: The paper is available at https://www.osti.gov/biblio/1991387.

Funding: This work is supported by Sandia Earth Science Investment Area Laboratory Directed Research and Development funding. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.

SANDIA REPORT

SAND2023-05141 Printed June, 2023



Benchmarking the PCMCI Causal Discovery Algorithm for Spatiotemporal Systems

J. Jake Nichol, Michael Weylandt, Mark Smith, Laura P. Swiler

Prepared by Sandia National Laboratories Albuquerque, New Mexico 87185 Livermore, California 94550 Issued by Sandia National Laboratories, operated for the United States Department of Energy by National Technology & Engineering Solutions of Sandia, LLC.

NOTICE: This report was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government, nor any agency thereof, nor any of their employees, nor any of their contractors, subcontractors, or their employees, make any warranty, express or implied, or assume any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represent that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government, any agency thereof, or any of their contractors or subcontractors. The views and opinions expressed herein do not necessarily state or reflect those of the United States Government, any agency thereof, or any of their contractors.

Printed in the United States of America. This report has been reproduced directly from the best available copy.

Available to DOE and DOE contractors from

U.S. Department of Energy Office of Scientific and Technical Information P.O. Box 62 Oak Ridge, TN 37831

Telephone: (865) 576-8401 Facsimile: (865) 576-5728 E-Mail: reports@osti.gov

Online ordering: http://www.osti.gov/scitech

Available to the public from

U.S. Department of Commerce National Technical Information Service 5301 Shawnee Road Alexandria, VA 22312

Telephone: (800) 553-6847 Facsimile: (703) 605-6900 E-Mail: orders@ntis.gov

Online order: https://classic.ntis.gov/help/order-methods



ABSTRACT

Causal discovery algorithms construct hypothesized *causal graphs* that depict *causal dependencies* among variables in observational data. While powerful, the accuracy of these algorithms is highly sensitive to the underlying dynamics of the system in ways that have not been fully characterized in the literature. In this report, we benchmark the PCMCI causal discovery algorithm in its application to gridded spatiotemporal systems. Effectively computing grid-level causal graphs on large grids will enable analysis of the causal impacts of transient and mobile spatial phenomena in large systems, such as the Earth's climate. We evaluate the performance of PCMCI with a set of structural causal models, using simulated spatial vector autoregressive processes in one- and two-dimensions. We develop computational and analytical tools for characterizing these processes and their associated causal graphs.

Our findings suggest that direct application of PCMCI is not suitable for the analysis of dynamical spatiotemporal gridded systems, such as climatological data, without significant preprocessing and down-scaling of the data. PCMCI requires unrealistic sample sizes to achieve acceptable performance on even modestly sized problems and suffers from a notable curse of dimensionality. This work suggests that, even under generous structural assumptions, significant additional algorithmic improvements are needed before causal discovery algorithms can be reliably applied to grid-level outputs of earth system models.

ACKNOWLEDGMENTS

We thank members of the CLDERA Grand Challenge LDRD project team for helpful discussions and comments on an early draft of this manuscript. JJN also acknowledges support from his Ph. D. advisors, Dr. Matthew Fricke and Dr. Melanie Moses of the Department of Computer Science at the University of New Mexico.

JJN and MS developed the 1D model and performed and analyzed relevant simulations. JJN and MW developed the 2D model, characterized its VAR dynamics, and performed and analyzed relevant simulations. JJN and MW wrote and edited the manuscript. LPS supervised all research and edited the manuscript.

CONTENTS

1.	Intro	oduction	11			
	1.1.	Background and Related Work	12			
		1.1.1. Structural Causal Modelling				
		1.1.2. Causal Discovery & the PCMCI Algorithm				
	1.2.	Contributions				
2.	Methods					
	2.1.	Spatiotemporal Data Generation Models	16			
		2.1.1. Model Definition: One Spatial Dimension				
		2.1.2. Model Definition: Two Spatial Dimensions				
	2.2.	PCMCI Algorithm: Tuning Parameters				
3.	Simulation Design					
	3.1.	Simulation Design: One Spatial Dimension	27			
	3.2.					
4.	Res	ults	29			
	4.1.	Performance Measures	29			
	4.2.	One-dimensional model	32			
	4.3.	Two-dimensional model	33			
5.	Disc	cussion	42			
Bil	bliog	raphy	44			
Αp	pend	dices	47			
Α.	A. Additional Simulation Results: Two-Dimensional Model					

LIST OF FIGURES

A time series graph representation of the SCM in Equation (1.1). By associating each variable with a node for each time lag, it is possible to fully capture relationship between variables and their temporal ancestors.	14
Causal graphs of variables X,Y,Z at grid cells A,B,C,D , for the SCM defined by Equation (2.4). Here, each variable exhibits temporal autocorrelation at each grid cell (orange arrows), while we observe spatial dependence among X and cross-variable dependence $X \to Y \to Z$. All dependencies occur after a single lag	19
Spatial Updates in the Two-Dimensional Model (Section 2.1.2). The 3×3 NDM is expanded to a $N^2 \times N^2$ matrix which fully characterizes the action of the NDM and can be used to analyze the behavior of the resulting system. The sparsity pattern of this matrix is reflected in the time series causal network for this process	24
Dynamics matrix for the 3×3 NDM (a b cd e fg h i) as applied on a 4×4 lattice. Note the "nested circulant" structure of this matrix, where each colored block has a circulant structure, as well as the block circulant structure of the dynamics matrix as a whole.	25
A causal graph for the one-dimensional simulation model. The five variables V,W,X,Y,Z are each observed on 10 grid cells. Each variable exhibits temporal autocorrelation (orange), while only V and Y exhibit spatial (left/right) dependencies. Cross-variable dependencies exist at every grid cell according to the causal structure $V \to W \to X \to Y \to Z$. Both the cross-variable and left-to-right dependencies occur at one lag (red), while the right-to-left dependencies occur at two lags (green). This graph has $50 = 5 \times 10$ nodes and 130 edges: the time series causal graph would have $100 = 50 \times (\max \log 2)$ nodes	28
F_1 scores for the One-Dimensional model of Sections 3.1 and 3.1. Coefficients a and c represent autocorrelation and cross-correlation dependence coefficients, respectively, where cross-correlation relates to both variable-to-variable and cell-to-cell dependencies. Only stable coefficient combinations are shown. PCMCI performs better with more time samples, larger a values, and larger c values. When a or c are sufficiently large, the system becomes unstable.	35
F_1 score results from the one-dimensional spatial example with varying autocorrelation, a , and constant cross-correlation, c . Each data point includes all possible T time	20
F_1 score results from the one-dimensional spatial example with varying cross-correlation, c , and constant autocorrelation, a . Each data point includes all possible T time sam-	36
	each variable with a node for each time lag, it is possible to fully capture relationship between variables and their temporal ancestors

Figure 4-4. F_1 score results from the one-dimensional spatial example with varying T time sample Each data point includes all possible a and c dependence coefficients	
Figure 4-5. Effect of grid size (N) on PCMCI F_1 and MCC scores. Both metrics decrease relative slowly in N . Other simulation parameters are fixed to $\sigma = 1.0$, $NDD = \frac{3}{9}$, and $T = 1000$.	ly
Figure 4-6. Effect of increasing sample size (T) on PCMCI performance (MCC). Performance increases sublinearly in T , with $T > 575$ being necessary to obtain acceptable performance (MCC > 0.7). Box labels report median MCC across replicates. Other simulation parameters fixed as $N = 10$, $\sigma = 1.0$, and NDD $= \frac{6}{9}$	ce er- er
Figure 4-7. Effect of sample size, (T) grid size, (N) , and neighborhood dependence density of PCMCI performance (MCC). For sufficiently large sample sizes, PCMCI is able to consistently recover the true graph structure; the effect of grid size and NDD are lepronounced than T . Values shown are mean performance over 30 replicates	on to ess
Figure 4-8. Effect of sample size, (T) grid size, (N) , and neighborhood dependence density of PCMCI performance (MCC). For sufficiently large sample sizes, PCMCI is about to consistently recover the true graph structure; the effect of grid size and ND are limited. Values shown are mean performance over 30 replicates. $\sigma = 1$ for a simulations.	on lle D
Figure 4-9. Probability of PCMCI Success as a function of grid size N and sample size T , with success defined as MCC above a user-defined threshold. Results are empirical probabilities over 30 replicates: σ and neighborhood density are fixed to 1.0 and $\frac{6}{9}$ respectively. Lines depict a simple linear model of grid size on success probability, with shaded regions depicting (non-multiplicity adjusted) confidence intervals	th b- e- th
Figure 4-10. Effect of Innovation Magnitude (σ) on PCMCI performance (MCC). Changing appears to have no systematic effect on PCMCI performance	σ
Figure A-1. False Discovery Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits low FDR for $T > 50$. FDR decreases with the number of cause effects (density) and with increasing time samples.	al
Figure A-2. True Positive Rate of PCMCI under the scenarios described in Section 3.2. PCMC consistently exhibits low true positive rates for $T < 350$. TPR decreases with the number of causal effects and with increasing grid sizes.	CI ne
Figure A-3. False Negative Rate of PCMCI under the scenarios described in Section 3.2. PCMC consistently exhibits relatively high false negative rates in all scenarios, indicating lo statistical power. FNR generally increases with the number of causal effects and with increasing grid sizes.	CI w th 50
Figure A-4. True Negative Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits near perfect true negative rates in all scenarios. To the extent varies, TNR decreases with the number of causal effects and with decreasing grid states.	it
Figure A-5. False Positive Rate of PCMCI under the scenarios described in Section 3.2. PCMC consistently exhibits near perfect false positive rates in all scenarios. To the extent varies, FPR increases with the number of causal effects and with decreasing grid size	CI it

LIST OF TABLES

Table 4-1.	The stable autocorrelation (a) and cross-correlation (c) dependence coefficients identi-	
	fied for the one-dimensional model.	32

NOMENCLATURE

Abbreviation	Definition			
ANM	Additive Noise Model			
CI	Conditional Independence			
DAG	Directed Acyclic Graph			
DOE	Department of Energy			
ENSO	El Niño Southern Oscillation			
FCI	Fast Causal Inference [algorithm]			
LiNGAM	Linear Non-Gaussian Acyclic Model			
LPCMCI	Latent-PCMCI			
MCC	Matthews Correlation Coefficient			
NDD	Neighborhood Dependency Density			
NDM	Neighborhood Dynamics Matrix			
OCE	Optimal Causal Entropy			
PC	Peter-Clark [algorithm]			
PCA	Principle Component Analysis			
PCMCI	PC -Momentary Conditional Independence [algorithm]			
PDAG	Partially Directed Acyclic Graph			
PDE	Partial Differential Equation			
SCM	Structural Causal Model			
SEM	Structural Equation Model			
VAR	Vector Autoregressive [model]			

1. INTRODUCTION

Automated causal structure discovery is an exciting frontier of data-driven science and domain-informed machine learning, but techniques for causal discovery are still rather untested in complex domains. As part of a larger investigation of causal discovery and attribution in climate systems, we investigate the performance of a state-of-the-art algorithm for causal discovery from climate data. The algorithm returns a *causal graphical model* of the given variables. Causal graphical models are usually directed acyclic graphs (DAGs) that relate the causal dependence (graph edges) between variables (graph nodes). Due to the scientific, computational, and statistical difficulties of characterizing climate systems, we instead draw upon well-established techniques for the *benchmarking* of machine learning algorithms for the evaluation of causal discovery. Our results highlight the limitations of modern causal discovery approaches and demonstrate the unreliable performance of these algorithms, even in the most amenable scenarios.

To create the benchmark test cases and perform the various studies we show in this report, we rely on the ideas of benchmarking. According to Olson et al. [1], "the term benchmarking is used in machine learning to refer to the evaluation and comparison of ML methods regarding their ability to learn patterns in 'benchmark' datasets that have been applied as 'standards'. Benchmarking could be thought of simply as a sanity check to confirm that a new method successfully runs as expected and can reliably find simple patterns that existing methods are known to identify." There are many benchmark datasets available: readers may be familiar with the ImageNet database which is commonly used for image classification test problems [2]. Recently, there has been a growth in scientific machine learning benchmarks as well, see Thiyagalingam et al. [3, 4]. The benchmarking approach typically involves a few main steps: identification of training datasets which provide the benchmark data or "gold standard" data, identification of the algorithm or method being tested and associated algorithm choices that might be examined (e.g. number of layers in a neural network, activation function used, optimization algorithm to determine hyperparameters, etc.), and a set of performance metrics with which to evaluate the algorithm. Depending on the extent and focus of the benchmark exercise, the ML algorithm can be run with many algorithm choices and the best" choices can be identified, according to the performance metrics which typically involve "goodness" of fit" with respect to predicting the benchmark data but which also may include time to train, time to make a prediction or inference, amount of computing power needed, etc.

We note that causal discovery does not necessarily fall into the machine learning category: it involves aspects of statistical modeling and network inference. However, we feel the benchmark terminology as defined above represents the goal of our efforts well. We also have leveraged verification and validation concepts from the computational science community which focuses on PDE solutions for physical systems, with the goal of improving the credibility of computational models and assessing their predictive capability [5–8]. There are some aspects of verification, specifically solution verification, in the work presented in this report. In the subsequent sections, however, we use the benchmarking terminology.

Benchmarking becomes more challenging for structure-learning algorithms (such as causal discovery), because they require a complete ground-truth graph to evaluate correctness, rather than additional obser-

vations as traditional machine learning requires. This typically limits structure-learning benchmarking to high-fidelity simulation output or hypothesized ground-truth, developed from randomized control trials. While there are a number of metrics that measure the performance of a machine learning model (such as cross-validation error, leave-one-out error, *etc.*), they typically only apply to models predicting additional data points from observational probability distributions, rather than intervention distributions¹, because they capture the ability of the model to represent the training and/or testing data. They do not address other questions such as the correct implementation of the algorithm or the properties and performance it exhibits on various classes of problems. For causal modeling and causal discovery algorithms, there has been limited work specifically seeking to address the issue of "is the inferred graph or causal structure that the algorithm produces correct?" though the works of Runge [10], Runge et al. [11] provide limited, but promising, initial results in this space. In this work, we seek to partially address this important lacuna.

In this work, we report the results of an extensive benchmarking exercise for the PCMCI algorithm of Runge et al. [11]. We specifically focus on the performance of this algorithm as applied to data with spatial and temporal dependence. Our results rely upon a simulation framework inspired by statistical models for time series and by the spatial dynamics of cellular automata. While limited benchmarking of PCMCI has previously been performed, ours is distinguished by a thorough analysis of the effect of spatial structure on performance.

1.1. Background and Related Work

The philosophical and statistical aspects of causal inference and causal discovery are subtle but powerful and our discussion here is necessarily informal. For a further discussion of these issues, we refer the reader to the books by Peters et al. [9] and by Pearl and Mackenzie [12], as well as the many references therein.

1.1.1. Structural Causal Modelling

Causal network discovery, or causal structure learning, is the process of estimating a causal graph² of an underlying *structural causal model* (SCM) from observational data³ and subject matter expertise⁴. An SCM is a semi-mechanistic model, which augments a classical statistical model with a notion of causal structure.⁵ While exact estimation of the SCM is typically impossible, it is often possible to accurately estimate the causal network associated with that SCM. A causal network is a DAG representation of the SCM, where variables represent different aspects of the data and directed edges connect "cause" to "effect."

¹Intervention distributions are what causal graphs predict. We omit discussion of that topic and refer the reader to Peters et al. [9, p. 120-121]

²Also known as a *causal network*.

³Observational data is characterized as non-experimental data; it contains no planned interventions or controls.

⁴Subject matter expertise is represented by critical causal assumptions, which causal discovery algorithms leverage to reason about the statistical properties found in observational data.

⁵Classical probabilistic statistical models do not naturally incorporate causal structure, instead representing data as a simultaneous draw from an underlying probability distribution. Any temporal object, such as the sample path of Brownian motion, is a draw of a single time-indexed object from an underlying space, rather than a system obeying causal laws.

Many SCMs imply the same causal network, but, under reasonable assumptions,⁶ there is a unique DAG for any SCM. When considering SCMs of temporal data, there exist multiple ways of depicting the causal network; see the works by Eichler [13] and Peters et al. [9, p. 198] for details.

As a simple example, consider the following SCM:

$$W_{t} := 0.9W_{t-1} + \eta_{t}^{W}$$

$$X_{t} := 0.8X_{t-1} + 0.4W_{t-1} + 0.2Z_{t-3} + \eta_{t}^{X}$$

$$Y_{t} := 0.5Y_{t-1} + 0.2X_{t-2} + \eta_{t}^{Y}$$

$$Z_{t} := 0.6Z_{t-1} + 0.3Y_{t-1} + \eta_{t}^{Z}$$

$$(1.1)$$

where each $\eta \sim \mathcal{N}(0,1)$ is IID Gaussian noise. These relations form a SCM for simulated realizations of this process.

Figure 1-1 is a causal graph for the SCM in Equation (1.1). Specifically, it is a *time series graph* [10], which captures the temporal dependencies of each node. Each node is a temporally lagged instantiation of each variable. Notice that each variable is autocorrelated in Equation (1.1), with a link between itself and its past self, over 1 lag. The 2 and 3 lag dependencies in $X \to Y$ and $Z \to X$, respectively, are also depicted passing over their respective lag lengths. Without the lagged representation, time-delayed feedbacks⁷ would be illustrated as cycles, which violates an important assumption of causal graphs: acyclicity.

While an SCM maps to a DAG, causal network discovery algorithms often output partially-directed acyclic graphs (PDAGs) [14], in which some edges are undirected. Undirected edges indicate a dependence was identified, but not the direction of dependence. Edges sometimes fail to be oriented because of violated assumptions or too little data, but most causal discovery algorithms can only estimate up to the correct Markov equivalence class of graphs, even when assumptions are met and sampling is sufficient. See Peters et al. [9, p. 102] for more on the Markov equivalence of graphs.

Estimated graphs can be annotated with more information indicating the strength of dependence between nodes, causal effect size, causal susceptibility, *etc.* [15, 16], but in this work, we are only concerned with estimating the topology (edge structure) of the time series causal network.

Algorithms for reconstructing causal networks from data generated by an SCM are discussed in the next section.

1.1.2. Causal Discovery & the PCMCI Algorithm

Many algorithms for causal discovery have been proposed in the previous 30 years, most notably the PC algorithm [17], named for its authors Peter Spirtes and Clark Glymour, the Fast Causal Inference (FCI) [17], and the Linear Non-Gaussian Acyclic Model (LiNGAM) [18]. While these general-purpose

⁶These assumptions include causal faithfulness, the causal Markov condition, and causal sufficiency. Put simply, the faithfulness assumption states that separation of two nodes in the causal network is implied by independence, the causal Markov condition states that separation in the graph implies independence in the data, and causal sufficiency states that we have included all common causes of two or more variables in the analysis. Again, for a more detailed discussion, see the books of Peters et al. [9] and Pearl and Mackenzie [12].

⁷Such as that from $X \to Y \to Z \to X$ over 2, 1, and 3 lags, respectively.

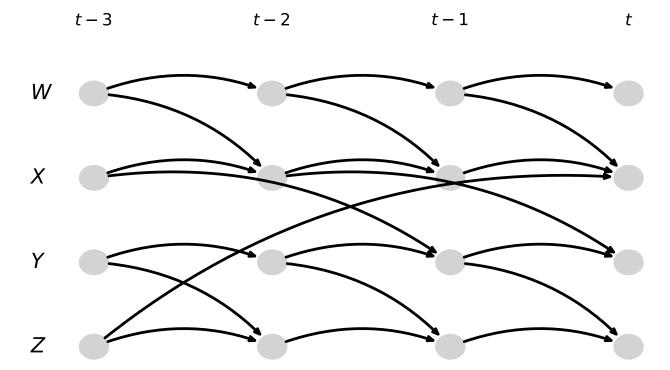


Figure 1-1. A time series graph representation of the SCM in Equation (1.1). By associating each variable with a node for each time lag, it is possible to fully capture relationship between variables and their temporal ancestors.

algorithms are primarly designed for non-temporal data, temporally-aware variants of these algorithms exist [16] as well as novel approaches specific to time series, such as the Optimal Causal Entropy (OCE) algorithm [19, 20]. In this work, we consider the PC-Momentary Conditional Independence (PCMCI) algorithm of Runge et al. [16]. We focus on PCMCI because it was specifically designed to deal with the complex temporal structure of climate data and it has found wide use among the causal climate community [15, 21–25].

PCMCI modifies the classical PC algorithm [17] by adding so-called "Momentary Conditional Independence" tests. These tests take advantage of the temporal structure of the data to greatly reduce the number of potential causal effects, thereby decreasing the space of possible causal networks and improving inferential performance. Like the PC algorithm, the output of PCMCI is a PDAG, however, the time-order of lagged dependencies helps PCMCI orient more edges than it would without temporal information.

The standard variant of PCMCI assumes all causal relationships work on a lag and that there are no contemporaneous dependencies in the data. While we focus on the standard PCMCI algorithm, our simulation study could easily be applied to PCMCI variants, including the Latent-PCMCI of Gerhardus and Runge [26], which allows for unobserved confounders, and PCMCI+ of Runge [27], which allows for contemporaneous dependencies.

Runge et al. [11] detail PCMCI thoroughly and provide an open-source implementation of the approach⁸. PCMCI is a two-phase algorithm: the first phase uses a modified version of the PC algorithm to construct a sparse causal PDAG; this modified algorithm, which they call PC₁, performs a series of iterative conditional

⁸https://jakobrunge.github.io/tigramite/

independence (CI) tests in a search for the causal parents of each variable. PC₁ modifies this search to only condition on the potential confounders with the largest correlations to the variables in question. While this significantly increases computational performance, the full impact of this heuristic modification has not yet been fully characterized.

The second phase of the PCMCI algorithm uses MCI tests to prune this graph in an attempt to eliminate temporally-induced spurious causality. MCI tests extend traditional conditional independence tests by conditioning on lagged (time-shifted) observations of variables. In doing so, they specifically examine whether apparent causal dependencies are artifacts of autocorrelation and prune these spurious graph edges and reduce the false positive rate of PC_1 .

As with the original PC algorithm, both the PC₁ and MCI steps of PCMCI can be used with arbitrary conditional independence tests. Test with the conditional Pearson correlation, the *partial correlation*, are easily implemented and widely used, but their performance is only guaranteed for (jointly) Gaussian data. Peters et al. [9] discuss alternative independence tests; see also the discussion by Runge [10].

Finally, we note that while PCMCI is commonly used for climate data, it does not take advantage of the spatial structure typically present in such data. Rather than dealing with spatial structure explicitly, common practice is to summarize data into non-spatial components before applying PCMCI. This summarization is typically done with a statistical technique such as Principal Components Analysis (PCA) or variants thereof or by using external climate knowledge to divide spatial data into pre-defined regions or modes, which are assumed to have no further spatial dependencies [15, 22, 25, 28, 29]. While powerful, these approaches have several drawbacks: PCA-type approaches construct features that are composed of all of the features of the underlying data, so the implied causal relationships are often of an "all-to-all" nature; *a priori* knowledge is useful for well-studied climate phenomena but is difficult to apply to novel studies. In this work, we consider working with unaggregated spatial data observed on a regular grid, such as the output of a large-scale earth system model or geo-referenced observational data. As we will see below, this approach poses novel difficulties in simulation and estimation.

1.2. Contributions

In this paper, we perform an extensive simulation study to benchmark the performance of PCMCI on a set of spatially-inspired SCMs. By using data generated from a known SCM, we are able to accurately quantify the performance of PCMCI on a variety of metrics. In addition to the analysis of PCMCI, our data simulation procedures may be of independent interest. Our findings inform the feasibility of causal discovery from real and simulated climate data and identify several challenges that must be addressed before applying these algorithms at scale.

Section 2 introduces the mathematical framework used to generate spatiotemporal data generation studies, while section 3 describes the specific parameter values used in our simulations. The results of our simulation studies are shown in Section 4, along with a detailed discussion of their implications for causal discovery practice. Finally, Section 5 summarizes our results and discusses potential directions of future research.

2. METHODS

2.1. Spatiotemporal Data Generation Models

Causal dependencies in multivariate data are often expressed as SCMs, *e.g.*, SCM (1.1). If there exists a direct causal dependence from X to Y, which we denote $X \to Y$, then we posit a relationship of the form:

$$Y := f_Y(X) + \eta_Y \qquad (\eta_Y \perp \!\!\! \perp X) \tag{2.1}$$

where f_Y is a (measurable) function relating the cause variable X to the effect variable Y and η_Y is additive noise. If X is random, then we assume X and η are independent ($\eta_Y \perp \!\!\! \perp X$), though this assumption may be relaxed in some circumstances. In the common case where $f_Y(\cdot)$ is a linear function of X, we recover the well studied class of linear structural equation models (SEM) [12]. As Peters et al. [9] discuss, the assumption of additive noise in Equation (2.1) is not essential, but it is standard in the field and we will use it throughout our analysis.

The SCM (2.1) is an additive noise $model^2$ (ANM) [9, p. 50], a restriction on the class of SCMs that is also useful for identifying variables which do not exhibit a causal effect on Y. Suppose that

$$Y = f(X,Z) + \eta_Y$$

for some function f. It can be shown that Z is not a parent of Y if there exists some function g(X) such that f(X,Z) = g(X) for all (X,Z) or equivalently $Y = g(X) + \eta_Y$.

When modeling temporal data, the ANM (2.1) must be modified to allow for a variable to depend on its previous values. Let X_t be the state of a system of interest at time t; we make two standard assumptions on the behavior of X_t :

- **T1)** Lagged dependence: $X_{i,t} \not\to X_{j,t-\tau}^3$ for any (i,j) and any $\tau \ge 0$.
- **T2)** Temporal Causal Stationarity: the dynamics governing the evolution of X_t do not change over time.

These assumptions are essentially unavoidable in causal analysis of temporal data: Assumption T1 states that causal dependencies follow the "arrow of time" while Assumption T2 implies that there is a fixed causal structure that we are seeking to estimate. If T2 did not hold, then it is unclear what our target of

¹We also assume that $Y \nsubseteq X$, *i.e.*, that Y does not appear on both sides of equation (2.1): this is essentially equivalent to the common assumption that the causal graph of the system is a DAG.

²Following the notation in Peters et al. [9], we will hereafter use assignment (:=) when describing SCM definitions, and equivalence (=) when specifying ANMs and, later, autoregressive models. In this work, the ANMs and autoregressive models are generative models, so they are no less causal.

³For our purposes, \rightarrow indicates no direct dependence between variables.

estimation actually is.⁴ Under these assumptions, the ANM for a system with only a single temporal lag⁵ becomes:

$$\boldsymbol{X}_t = f(\boldsymbol{X}_{t-1}) + \boldsymbol{\eta} \tag{2.2}$$

where, as before, η is an independent noise variable. In the temporal context, where the effect of the randomly sampled $\eta_{i,t}$ terms persists over time, we will typically refer to the $\eta_{i,t}$ terms as *innovations* rather than *error* or *noise* to emphasize that they are not measurement error, but rather are the fundamental driving element of the system.

In simulation settings, $f(\cdot)$ often represents one step of a (explicit) PDE solver [9]. If $f(\cdot)$ is a linear function, then Equation (2.2) is a Vector Autoregressive (VAR) model [10, 11] and can be written as

$$\boldsymbol{X}_t = \boldsymbol{A}\boldsymbol{X}_{t-1} + \boldsymbol{\eta}$$

where A is a fixed matrix encoding the causal dynamics of the system. Specifically, we note that the sparsity pattern of A exactly captures the causal structure of the system:

$$X_{i,t-1} \to X_{j,t} \Leftrightarrow A_{ij} \neq 0$$

As we will observe in the sequel, this property of VARs is particularly useful when simulating from and estimating causal structure in temporal data.

So far, our development has not posited any spatial structure to X_t , only the temporal lagged-dependence structure of Equation (2.2). We next introduce two spatial causal assumptions that parallel our temporal assumptions:

- **S1)** Neighborhood dependence: if (i, j) are not neighbors (in a problem specific sense) then $X_i \not\to X_j$.
- **S2)** Spatial Causal Stationarity: the dynamics governing the evolution of X_t do not change over space.

Assumption S1 attempts to capture a sense of "locality" and to disallow "action at a distance." When applying this assumption to physical systems, this implies a certain relationship between the temporal and spatial discretizations used: at sufficiently low observation rates, it is possible for a causal effect to exist beyond immediate neighbors. We do not explore the details of that relationship here, but we do note that similar concerns are well-studied in the design of numerical differential equation solvers where spatial and temporal discretizations must be chosen in a suitably consistent manner. Like Assumption T2, Assumption S2 ensures that PCMCI is learning the same causal structure throughout the space. Assumption S2 is not essential in this application and can be easily relaxed. These dynamics are similar to rule-based cellular automata (CA), where the state of each cell is dependent on its immediate neighbors and the update rules are fixed across all cells and time steps.

Under these assumptions, we obtain the single-lag spatiotemporal ANM:

$$X_{i,t} = f(X_{i,t-1}, \{X_{j,t-1}\}_{j \in \mathcal{N}(i)}) + \eta_{i,t}$$

⁴Assumption T2 can be weakened to only require the *causal* structure of the dynamics, and not the full dynamics, to remain constant over time, but we do not pursue this relaxation.

⁵ For higher order lags, we have $\mathbf{X}_t = \sum_{\tau=1}^T f_{\tau}(\mathbf{X}_{t-\tau}) + \boldsymbol{\eta}$, but we omit higher lags for simplicity of exposition unless noted otherwise. See Peters et al. [9, p. 208] for additional discussion.

⁶For example, consider a simple system in which $X_{i+1,t+1} = X_{i,t} + \eta_{i,t}$ for all (i,t). If i is interpreted as a spatial coordinate in a single dimension, this system satisfies S1. If we reduce our sampling and can only observe $\mathbf{Y}_t = \mathbf{X}_t, \mathbf{Y}_{t+1} = \mathbf{X}_{t+2}, \dots, \mathbf{Y}_{t+\tau} = \mathbf{X}_{t+2\tau}$, we instead have the causal relationship $Y_{i+2,t} = Y_{i,t}$ which appears to violate S1.

where $\mathcal{N}(i)$ denotes the neighborhood of i, . If f is further assumed to be linear, then we have

$$X_{i,t} = \alpha X_{i,t-1} + \sum_{j \in \mathcal{N}(i)} \beta_j X_{j,t-1} + \eta_{i,t}$$
 (2.3)

The sparsity of the α and β_j coefficients dictates the causal structure of X_t . We will occasionally refer to α as a temporal autocorrelation coefficient and β_j as a cross-dependence coefficient, though they are not numerically equal to the actual autocorrelation function of the process X_t .

It is clear that Equation (2.3) can be again expressed as a linear VAR system, with the spatial assumptions S1 and S2 posing additional constraints on the structure of the dynamics (coefficient) matrix. In the next two sections, we characterize these constraints for one- and two-dimensional systems, leaving higher-dimensional systems to the reader.

Specifically, we consider two spatial cases to evaluate different kinds of spatiotemporal dynamics. In Section 2.1.1, we consider a multivariate, multi-lagged model supported on a one-(spatial)-dimensional array. In Section 2.1.2, we consider a univariate single-lag model supported on a two-(spatial)-dimensional array. For both models, we assume the underlying space has a toroidal topology, with the leftmost and rightmost elements of the one-dimensional space being neighbors, and similarly for the topmost and bottommost elements in the two-dimensional case. In one-dimension, the torus is a circle, while the two-dimensional torus is a "donut" shape. We note that this topology differs from that of the surface of a sphere, in that moving far north does not have the same effect as moving far to the west and that there is no analogue of a pole where all cells coincide, but our results can be extended to that setting. Under these two settings, we design an extensive simulation study to characterize the performance of causal discovery algorithms on spatial data.

⁷More informally, we simulate dynamics in a world which "wraps" like the classic arcade game PACMAN.

2.1.1. Model Definition: One Spatial Dimension

We first consider simulating causal dynamics on a one-dimensional spatial lattice of size *N*. Under our assumption S2, we note that each cell can only causally depend on itself and its immediate left and right neighbors, suitably lagged. We further consider a "multivariate" setting in which multiple variables are observed for each cell, and where the causal structure for different variables may not coincide.

We describe the structure of our one-dimensional model in some detail, noting that most of the intuition transfers to the two-dimensional case we consider in the following section. On a lattice of size N=4, we observe three variables, X,Y,Z. Within a single variable, only X exhibits spatial dependencies, such that each cell depends on the neighbor to its left. The causal structure between variables is $X \to Y \to Z$. This sort of model is suitable for simplified modeling of atmospheric aerosol advection and their interaction with radiation and atmospheric temperatures: for some aerosol species, wind can advect aerosols to spatially neighboring regions, while the causal structure $X \to Y \to Z$ reflects the aerosol particles' radiation absorption and subsequent temperature impact, e.g., $H_2SO_4 \to radiative$ flux \to atmospheric temperature. See Figure 2-1a for a spatial illustration of this structure, and Figure 2-1b for a time series graph of the same example. Figure 2-1a is an example of a *summary graph* [9, p. 199].

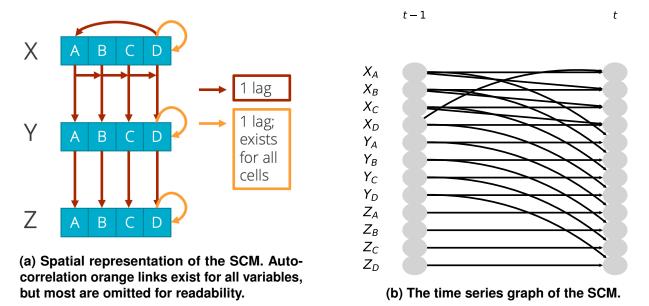


Figure 2-1. Causal graphs of variables X,Y,Z at grid cells A,B,C,D, for the SCM defined by Equation (2.4). Here, each variable exhibits temporal autocorrelation at each grid cell (orange arrows), while we observe spatial dependence among X and cross-variable dependence $X \to Y \to Z$. All dependencies occur after a single lag.

If we assume linear dynamics for this system, we obtain the SCM:

$$X_{A,t} := \alpha_{X,A} X_{A,t-1} + \beta_{X,A} X_{D,t-1} + \eta_{X,A,t}$$

$$X_{B,t} := \alpha_{X,B} X_{B,t-1} + \beta_{X,B} X_{A,t-1} + \eta_{X,B,t}$$

$$X_{C,t} := \alpha_{X,C} X_{C,t-1} + \beta_{X,C} X_{B,t-1} + \eta_{X,C,t}$$

$$X_{D,t} := \alpha_{X,D} X_{D,t-1} + \beta_{X,D} X_{C,t-1} + \eta_{X,D,t}$$

$$Y_{A,t} := \alpha_{Y,A} Y_{A,t-1} + \gamma_{X \to Y,A} X_{A,t-1} + \eta_{Y,A,t}$$

$$Y_{B,t} := \alpha_{Y,B} Y_{B,t-1} + \gamma_{X \to Y,B} X_{B,t-1} + \eta_{Y,B,t}$$

$$Y_{C,t} := \alpha_{Y,C} Y_{C,t-1} + \gamma_{X \to Y,C} X_{C,t-1} + \eta_{Y,C,t}$$

$$Y_{D,t} := \alpha_{Y,D} Y_{D,t-1} + \gamma_{X \to Y,D} X_{D,t-1} + \eta_{Y,D,t}$$

$$Z_{A,t} := \alpha_{Z,A} Z_{A,t-1} + \gamma_{Y \to Z,A} Y_{A,t-1} + \eta_{Z,A,t}$$

$$Z_{B,t} := \alpha_{Z,B} Z_{B,t-1} + \gamma_{Y \to Z,B} Y_{B,t-1} + \eta_{Z,B,t}$$

$$Z_{C,t} := \alpha_{Z,C} Z_{C,t-1} + \gamma_{Y \to Z,C} Y_{C,t-1} + \eta_{Z,C,t}$$

$$Z_{D,t} := \alpha_{Z,D} Z_{D,t-1} + \gamma_{Y \to Z,D} Y_{D,t-1} + \eta_{Z,D,t}$$

Because this system is linear, we have an equivalent vector autoregressive (VAR) process representation, $\chi = \Gamma + \eta$:

Here, the α parameters control the temporal autocorrelation of each cell-variable series with itself, the β parameters control the spatial dependence within a variable, and the γ parameters capture cross-variable dependencies. In this scenario, we assume only variable X has spatial dependencies within the same variable, while variables Y and Z exhibit only autocorrelation and the cross-variable structure $X \to Y \to Z$. If we further assume causal stationarity for this model (Assumption S2), these dynamics simplify further to

That is:

 $ilde{oldsymbol{\chi}} = ilde{oldsymbol{\Gamma}} + oldsymbol{\eta}$:

- $\alpha_v = \alpha_{v,\ell}$ for all variables v and spatial locations ℓ ;
- $\beta = \beta_{X,\ell}$ for all spatial locations ℓ ;
- $\gamma_{v \to w} = \gamma_{v,w,\ell}$ for all variables v,w and all spatial locations ℓ

Further examination of this matrix reveals several sub-blocks with circulant structure, including an α_X , β block, a $\gamma_{X \to Y}$ block, a $\gamma_{Y \to Z}$ block, and α_Y and α_Z blocks: we will return to this observation in the next section.

The specific values of α_v , β , and $\gamma_{v\to w}$ determine whether the resulting stochastic process has spatiotemporal statistical stationarity, which we will call "stability" for brevity. PCMCI assumes the given time series are statistically stationary, so we need to filter the coefficients that constitute a stable process. To do that, we constructed a *companion matrix* [30, p. 259], which is of the general form:

$$m{F} = egin{bmatrix} m{ ilde{\Gamma}}_{t-1} & m{ ilde{\Gamma}}_{t-2} & \dots & m{ ilde{\Gamma}}_{t- au} \ m{I} & m{0} & \dots & m{0} \ m{0} & \ddots & m{0} & dots \ m{0} & m{0} & m{I} & m{0} \end{bmatrix}$$

for τ lags in the model. The companion matrix is a matrix composed of the $\tilde{\Gamma}$ coefficient matrices (defined above), and the identity matrices and zero matrices that match the size of $\tilde{\Gamma}$. If all eigenvalues of the companion matrix are less than one, then the chosen coefficients will constitute a stable system [30, p. 259]. In Section 3.1, we describe a two-lag system used for experiments, and the companion matrix we used for determining stability is given by:

$$m{F_1} = egin{bmatrix} ilde{m{\Gamma}}_{t-1} & ilde{m{\Gamma}}_{t-2} \ ilde{m{I}} & m{0} \end{bmatrix}$$

In Section 3.1 we give specifics of the various model parameters used in our simulations. Because our spatiotemporal model thus reduces to a standard VAR process, for which the PCMCI causal discovery algorithm has previously been found to be effective, we note that our results complement and extend what has previously been shown for the PCMCI algorithm [11].

2.1.2. Model Definition: Two Spatial Dimensions

We next consider simulating causal dynamics on a two-dimensional finite lattice of dimension *N*. As before, we require that the simulated system has VAR-type dynamics and satisfies assumptions S1-2 and T1-2.

In two spatial dimensions, Assumption S2 implies that each cell has eight neighbors in its so-called "Moore neighborhood", yielding a total of nine potential causal parents (eight neighboring cells and the dependent cell's own previous value). As such, the causal dynamics of the system are dictated by a 3-by-3 matrix, which we term the *neighborhood dependence matrix* (NDM). To simulate dynamics from the NDM, we

⁸In the study of cellular automata, the Moore neighborhood of a cell includes both orthogonal and diagonal neighbors, while the von Neumann neighborhood includes only orthogonal (up, down, left, right) neighbors.

update each element of X_t by taking the inner product of the NDM and the immediate neighborhood of a grid cell: that is,

$$X_{ij,t} = \langle X_{\mathcal{N}(ij),t-1}, \mathbf{NDM} \rangle + \eta_{ij,t} = \operatorname{Trace}(X_{\mathcal{N}(ij),t-1}^{\top} \mathbf{NDM}) + \eta_{ij,t}$$

where $X_{\mathcal{N}(ij)}$ is the submatrix of X consisting of the $(i,j)^{\text{th}}$ element and its immediate neighbors. The NDM defines an invariant "update kernel" which is applied separately to each grid cell in order to simulate its expected value at the next time step. As such, the NDM update dynamics are a sliding dot product of the NDM and the spatial grid, defined by X_t :

$$\boldsymbol{X}_{t} = \mathbf{NDM} \star \boldsymbol{X}_{t-1} + \boldsymbol{\eta}_{t} \tag{2.5}$$

For two matrices $\mathbf{A} \in \mathbf{R}^{n \times n}$ and $\mathbf{B} \in \mathbf{R}^{N \times N}$, we define their sliding dot product $\mathbf{C} \in \mathbf{R}^{N \times N}$ to be the matrix with $(k, l)^{\text{th}}$ element given by

$$C_{kl} = \sum_{i=-\lceil n/2 \rceil}^{\lceil n/2 \rceil} \sum_{j=-\lceil n/2 \rceil}^{\lceil n/2 \rceil} A(k+i \mod N, l+j \mod N) B(2i+1, 2j+j). \tag{2.6}$$

where the mod operator is used to enforce wrapping at the boundaries of our lattice. In our context, the dimension of the sliding dot product kernel $\mathbf{A} = \mathbf{NDM}$ is fixed as n = 3, reflecting the size of the local neighborhood of each cell; the dimension of the state variable $\mathbf{B} = \mathbf{X}_t$ varies with the size of the lattice.

While it is possible to simulate dynamics according to Equation (2.5) for any NDM, the resulting multivariate time series is not statistically stationary without additional assumptions on **NDM**. In order to guarantee stationarity, we seek to represent Equation (2.5) as a (linear) VAR model and apply standard stationarity requirements [30]. In particular, we know that if we have VAR dynamics of the form

$$\mathbf{Y}_t = \mathbf{A}\mathbf{Y}_{t-1} + \mathbf{\eta}_t$$

the time series $\{Y_t\}$ is stationary if $\|A\|_{op} < 1$, where $\|\cdot\|_{op}$ denotes the operator or spectral norm of a matrix, *i.e.*, the magnitude of its largest (possibly complex) eigenvalue. Hence, for a given NDM A, it suffices to find a matrix $\widetilde{A} \in \mathbb{R}^{N^2 \times N^2}$ such that

$$\operatorname{vec}(\boldsymbol{X}_{t}) = \widetilde{\boldsymbol{A}}\operatorname{vec}(\boldsymbol{X}_{t-1}) + \operatorname{vec}(\boldsymbol{\eta}_{t})$$
(2.7)

Figure 2-2 demonstrates how the NDM, \mathbf{A} , can be used to form an equivalent VAR coefficient matrix, \mathbf{A} . For each grid cell, a suitably padded and shifted version of the NDM is constructed and then multiplied with the previous length N^2 state vector, $\text{vec}(\mathbf{X}_{t-1})$. Repeating this process for all N^2 grid cells creates the N^2 -by- N^2 coefficient matrix for the VAR representation. We do not seek to fully characterize the algebraic properties of this matrix here, but we do note that it exhibits a block convolutional structure, as shown in Figure 2-3; that is, it has the form of a N-by-N circulant matrix where each element is itself an N-by-N (sub)block matrix. Because the sliding dot product is closely related to a convolution, this circulant block structure is not unexpected.

⁹Denoted by ★; also known as a cross-correlation in signal processing, or a flipped convolution à la convolutional neural networks.

With this representation in hand, we are now able to characterize NDMs that give statistically stationary spatiotemporal data (which for brevity we will call "stable NDMs"): a 3-by-3 NDM, \boldsymbol{A} yields stable dynamics if its equivalent N-by-N VAR coefficient matrix $\boldsymbol{\tilde{A}}$ satisfies $\|\boldsymbol{\tilde{A}}\|_{op} < 1$.

In our simulations below, we leverage this characterization as the basis of an Accept-Reject sampling scheme for statistically stationary NDM matrices from the asymmetric Gaussian ensemble. See Algorithm 1. While the efficiency of Algorithm 1 was more than sufficient for this study, more work is needed to efficiently sample stationary NDMs on larger grids. We note that, though natural, this characterization of stationary NDMs does not appear to have been previously considered in the literature and the VAR representation appears to be novel. Previous simulation studies of PCMCI, such as that of Runge [10] and Runge et al. [11], do not sample from the space of stable NDMs and instead explicitly construct a selection of SCMs with small coefficients whose stationarity is then verified empirically through simulation.

Algorithm 1 Sampling Stable Gaussian NDMs: Accept/Reject Algorithm

• **Output: A** sampled from $\mathbf{A} \sim \mathcal{N}(\mathbb{R}^{3\times 3})|\mathbf{A}|$ is stationary

• Repeat:

- 1. Sample $\mathbf{A} \in \mathbb{R}^{3 \times 3}$ from the 9-dimensional standard Gaussian distribution
- 2. Construct \widetilde{A} according to the process of Figure 2-2
- 3. If $\|\mathbf{A}\|_{op} < 1$ return \mathbf{A}

In our two-dimensional simulation studies below, we only consider the single-lag single-variable VAR defined by Equations (2.5) and (2.7). Extensions to more complex models are straight-forward. For our model, the multilag extension of Equation 2.5 is given by

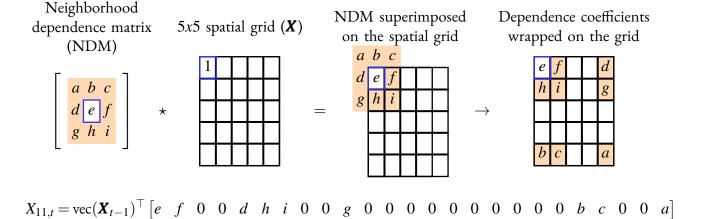
$$\boldsymbol{X_t} = \sum_{\ell=1}^{L} \boldsymbol{A}_{\ell} \star \boldsymbol{X}_{t-\ell} + \boldsymbol{\eta}_t$$
 (2.8)

for L lags, while the multilag, the *multivariate* extension of Equation 2.5 is given by

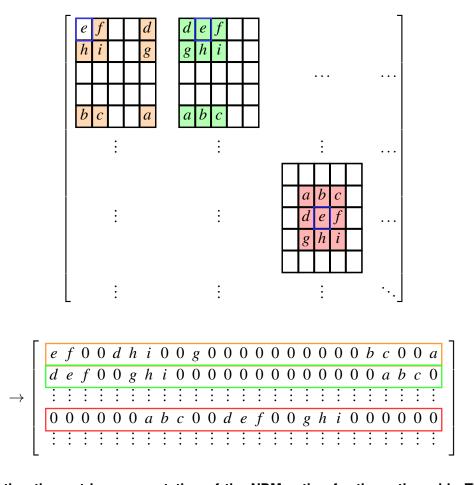
$$\boldsymbol{X}_{t}^{(J)} = \sum_{\ell=1}^{L} \sum_{j=1}^{\mathfrak{J}} \boldsymbol{A}_{\ell}^{(j \to J)} \star \boldsymbol{X}_{t-\ell}^{(j)} + \boldsymbol{\eta}_{t} \qquad \text{for } J = 1, \dots, \mathfrak{J}$$
(2.9)

for L lags and $\mathfrak J$ variables. Here $\pmb A_\ell$ denotes the lag- ℓ NDM while $\pmb A_\ell^{(j\to J)}$ denotes the multivariate dependence NDM of J on j at lag ℓ .

Finally we note that the single variable VAR(1) here represents the easiest case for causal discovery algorithms. The introduction of more lags, more variables, or non-linear dependencies would only increase the difficulty of causal discovery. As such, the experiments we show below represent an *upper bound* on the performance of PCMCI as applied in more realistic scenarios.



(a) Mapping the action of a neighborhood dependence matrix (NDM) on a single grid cell to a matrix representation. As the NDM is applied to the top left grid cell of the 5×5 spatial grid, the update incorporates all 8 neighbors, which wrap both vertically and horizontally around the edge of our 2D torus. The action of the NDM on a particular grid cell is represented by the top right matrix, which can easily be seen to be equivalent to the vector-matrix product formulation shown below.



(b) Constructing the matrix representation of the NDM action for the entire grid. The process described in Figure is repeated for each grid cell in the 5×5 lattice, which produces a 5×5 matrix, each element of which is a 5×5 matrix reflecting the NDM on a particular cell. Vectorizing these matrices yields the full 25×25 -update matrix shown in the final row.

Figure 2-2. Spatial Updates in the Two-Dimension $\frac{1}{2}$ Model (Section 2.1.2). The 3×3 NDM is expanded to a $N^2 \times N^2$ matrix which fully characterizes the action of the NDM and can be used to analyze the behavior of the resulting system. The sparsity pattern of this matrix is reflected in the time series causal network for this process.

```
e f 0 d h i 0 g 0 0 0 0 b c 0 a
d e f 0 g h i 0 0 0 0 0 0 a b c 0
0 d e f 0 g h i 0 0 0 0 0 0 a b c
f 0 d e i 0 g h 0 0 0 0 0 c 0 a b
b c 0 a e f 0 d h i 0 g 0 0 0 0 0
a b c 0 d e f 0 g h i 0 0 0 0 0
0 a b c 0 d e f 0 g h i 0 0 0 0 0
c 0 a b f 0 d e i 0 g h 0 0 0 0 0
0 0 0 0 b c 0 a e f 0 d h i 0 g
0 0 0 0 a b c 0 d e f 0 g h i 0
0 0 0 0 a b c 0 d e f 0 g h i 0
0 0 0 0 a b c 0 d e f 0 g h i 0
0 0 0 0 a b c 0 d e f 0 g h i 0
0 0 0 0 a b c 0 d e f 0 g h i
0 0 0 0 0 a b c 0 d e f 0 g h i
0 0 0 0 0 a b c 0 d e f 0 g h i
0 0 0 0 0 a b c 0 d e f 0 g h i
0 0 0 0 0 a b c 0 d e f 0 g h
h i 0 g 0 0 0 0 b c 0 a e f 0 d
g h i 0 0 0 0 a b c 0 d e f
i 0 g h 0 0 0 0 c 0 a b f 0 d e
```

Figure 2-3. Dynamics matrix for the 3×3 NDM $\begin{pmatrix} a & b & c \\ d & e & f \\ g & h & i \end{pmatrix}$ as applied on a 4×4 lattice. Note the "nested circulant" structure of this matrix, where each colored block has a circulant structure, as well as the block circulant structure of the dynamics matrix as a whole.

2.2. PCMCI Algorithm: Tuning Parameters

The PCMCI algorithm has two tuning parameters which must be set by the analyst:

- au_{max} , the maximum dependence lag
- $lpha_{PC}$, the significance threshold used for each conditional independence test

 au_{max} can be chosen based on expert knowledge of the system to determine the maximum hypothetical time for causality to propagate. In general, setting au_{max} too low will significantly distort the estimated causal structure, while setting au_{max} too high will slightly increase the runtime and the false positive rate of PCMCI; as such, users should err on the high side of possible values of au_{max} when the optimal value is unknown.

The PCMCI algorithm uses the α_{PC} parameter for pruning links in the PC Condition Selection phase of the algorithm. During this phase, the (classical) PC Condition Selection algorithm is used for Markov blanket discovery, where it proceeds by running a series of conditional independence tests and removes the link between two variables if the associated test has a *p*-value less than α_{PC} . As Runge et al. [11] notes, PCMCI does not account for dependencies among the various independence tests or for multiple testing and α_{PC} is better interpreted as a regularization parameter than a statistical significance level, as the false

positive rate of the PCMCI algorithm is not controlled. *Ceteris paribus*, decreasing α_{PC} will result in a sparser estimated causal graph.

Other free parameters include a *minimum* lag τ_{min} , autocorrelation control parameters p_X and p_Y , and a final threshold level α_G which is applied as a heuristic multiplicity correction. The roles of these parameters are described in more detail by Runge et al. [11] and we do not vary them in our analysis.

3. SIMULATION DESIGN

3.1. Simulation Design: One Spatial Dimension

In order to assess the performance of PCMCI on our one-dimensional model, we fixed a grid size of N=10 and considered five variables observed at each grid cell, V,W,X,Y,Z. Only variables V and Y exhibited spatial dependence: with a left-to-right dependence at one lag and a right-to-left dependence at two lags ($V_{i-1,t-1} \rightarrow V_{i,t}$ and $V_{i+2,t-2} \rightarrow V_{i,t}$ and similarly for Y). Our simulation design is depicted in Figure 3-1.

Runge et al. [16] note that temporal autocorrelation is typically a severe difficulty for causal discovery algorithms. The PCMCI algorithm was developed specifically to abate these difficulties [16]. To assess the performance of PCMCI, we sampled autocorrelation, which we call coefficient a, from the range $\{0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9\}$, with a common autocorrelation used for all variables and grid cells. We consider many of these high degrees of autocorrelation, as autocorrelation is a notable aspect of the climate science questions motivating this study.

We sampled both within-variable spatial and between-variable dependence coefficients, which we call coefficient c, from $\{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9\}$, to assess the performance of PCMCI under a range of dependence structures. These dependence coefficients were held constant at all grid cells. Innovations $(\eta_{i,t})$ were sampled from the standard normal distribution. We generated time series with T time samples ranging from $\{50, 150, 250, 350, 475, 575, 675, 775, 900, 1000\}$.

Parameter combinations that failed to exhibit stable dynamics were excluded from our analysis. We ran 30 replicate simulation runs for each stable parameter combination. The number of possible simulation runs is 30,000, however, because most coefficient combinations were not stable, the number of runs completed was 4,500. The specific coefficients used are detailed in Section 4.2.

3.2. Simulation Design: Two Spatial Dimensions

In order to characterize the performance of PCMCI in a variety of regimes, we considered the following simulation parameters for our two-dimensional model:

- Number of Time Samples (T): $\{50, 150, 250, 350, 475, 575, 675, 775, 900, 1000\}$
- Grid Size (*N*): $\{4x4, 5x5, 6x6, 7x7, 8x8, 9x9, 10x10\}$
- Innovation Scale ($\sigma = \text{sd}(\eta_{i,t})$): $\{0.1, 0.5, 1.0, 2.0, 4.0\}$
- Neighborhood Dependence Density (NDD): $\frac{1}{9}$, $\frac{2}{9}$, $\frac{3}{9}$, $\frac{4}{9}$, $\frac{5}{9}$, $\frac{6}{9}$, $\frac{7}{9}$, $\frac{8}{9}$, $\frac{9}{9}$

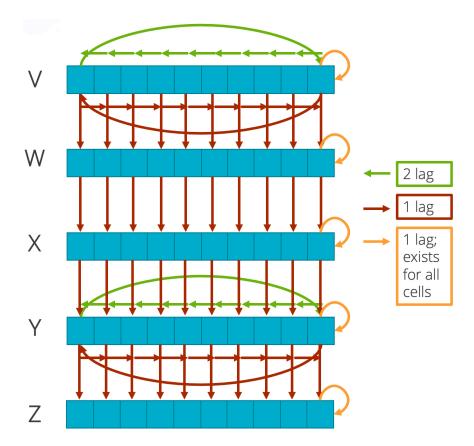


Figure 3-1. A causal graph for the one-dimensional simulation model. The five variables V,W,X,Y,Z are each observed on 10 grid cells. Each variable exhibits temporal autocorrelation (orange), while only V and Y exhibit spatial (left/right) dependencies. Cross-variable dependencies exist at every grid cell according to the causal structure $V \to W \to X \to Y \to Z$. Both the cross-variable and left-to-right dependencies occur at one lag (red), while the right-to-left dependencies occur at two lags (green). This graph has $50 = 5 \times 10$ nodes and 130 edges: the time series causal graph would have $100 = 50 \times (\text{max lag = 2})$ nodes.

Here, σ controls the scale of Gaussian innovations added to each element of X_t , and the NDD measures the number of causal parents implied by the NDM. When NDD = $\frac{1}{9}$, there is only one dependence between neighboring grid cells¹; increasing NDD adds more dependencies; NDD = $1 = \frac{9}{9}$ implies a fully connected (local) causal system. For each of these 3,150 parameter combinations, we generated 30 random stationary NDMs, yielding a total of 94,500 NDMs, from which we generated 94,500 time series.

In order to simulate these dynamics, statistically stationary NDMs are sampled using Algorithm 1. In order to avoid causal signals that are too small to be detected, we additionally only considered NDMs whose non-zero elements had magnitude at least 0.1. Because the NDMs selected were guaranteed to be stable, we encountered no numerical difficulties in our data generation process.

We generated the innovations $\eta_{i,t}$ from a suitable mean-zero normal distribution and used a Gaussian condition independence test in PCMCI. If a specific distribution for $\eta_{i,t}$ is not assumed, non-parametric independence tests can be used, though these have a higher sample complexity and require longer observational series (greater T).

¹Sometimes dependence is between a grid cell and itself, such that nodes are autocorrelated and there is no cross-dependence.

4. RESULTS

4.1. Performance Measures

To compare the PCMCI-estimated causal graphs with the underlying SCM-implied causal graphs, we report discovery performance using several measures of classification accuracy; in particular we show the F_1 -score and the Matthews Correlation Coefficient. Additional accuracy measures appear in the Appendix to this report.

The F_1 score is a popular measure of classification accuracy, which attempts to balance the precision and recall of a classifier. Specifically, the F_1 score is defined as [31]:

$$F_1 = \frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \text{Harmonic Mean}(\text{Precision}, \text{Recall})$$
(4.1)

where precision and recall are defined as

$$Precision = \frac{TP}{TP + FP} \tag{4.2}$$

$$Recall = \frac{TP}{TP + FN} \tag{4.3}$$

and TP, FP, and FN are the counts of true positives, false positives, and false negatives, respectively. F_1 ranges from 0.0 to 1.0, with 0.0 indicating perfect disagreement, that is the estimated graph is the complement of the true graph, and 1.0 indicating exact graph recovery.

We note that the F_1 score is undefined when TP = 0, as both Precision and Recall are 0, which would occur if there are no links in the true graph (i.e. all variables are independent). We note that the F_1 score can equivalently be expressed as [32]:

$$F_1 = \frac{2 \times TP}{2 \times TP + FP + FN}.\tag{4.4}$$

As such we, define F_1 to be 1.0 if FP, FN = 0 as the estimated graph is correctly fully sparse and 0.0 if FP > 0 or FN > 0.

We additionally report the Matthews Correlation Coefficient (MCC), also called the ϕ coefficient. Unlike F_1 , MCC depends on true negatives and is symmetric in the positive and negative labels: that is, if we

¹In our context, positives refer to the existence of a link while negatives refer to absence of a causal link. In other contexts, it may be more natural to refer to the absence of a causal link as a scientific finding, as the baseline assumption is that dependencies exist among all measured variables. The MCC measurement we report is invariant to this switch of labels.

compare the *complements* of the true graph and the estimated graph, representing causal independence, we get the same MCC. Chicco [32] defined MCC as follows²:

$$MCC = \frac{(TP \times TN - FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$
(4.5)

which ranges from [-1,1]. MCC = -1 implies the model is perfectly incorrect, MCC = 0 indicates a level of accuracy consistent with random guessing, and MCC = 1 indicates perfect graph recovery.

As before, we take care to define MCC for the case of sparse graphs (causal independence). MCC is undefined in any of these four cases:

- 1. if TP = 0 **AND** FP = 0
- 2. if TP = 0 **AND** FN = 0
- 3. if TN = 0 **AND** FP = 0
- 4. if TN = 0 **AND** FN = 0

We handle these cases separately, assigning values of $\{-1,0,+1\}$ as appropriate to the causal discovery problem.

- 1. If TP = 0 **AND** FP = 0, then the estimated graph is fully sparse:
 - a) if FN = 0, then the true graph is also fully sparse and we take MCC = 1;
 - b) if $FN \neq 0$ **AND** TN = 0, then the true graph is fully connected, the estimated graph missed all causal relationships, and we take MCC = -1;
 - c) if $FN \neq 0$ **AND** $TN \neq 0$, then some, but not all, of the causal independence relationships of the estimated graph are false and we take MCC = 0.
- 2. If TP = 0 **AND** FN = 0, then the true graph is fully sparse:
 - a) if FP = 0, then the estimated graph is also fully sparse and we take MCC = 1;
 - b) if $FP \neq 0$ **AND** TN = 0, then the estimated graph is fully connected, which is exactly wrong, and we take MCC = -1;
 - c) if $FP \neq 0$ **AND** $TN \neq 0$, then the estimated graph has implies some spurious causal dependencies and we take MCC = 0.
- 3. If TN = 0 **AND** FP = 0, then the true graph is fully connected:
 - a) if TP = 0, then the estimated graph is fully sparse, which is exactly wrong, and we take MCC = -1;
 - b) if $TP \neq 0$ **AND** FN = 0, then estimated graph is fully connected and we take MCC = 1;
 - c) if $TP \neq 0$ **AND** $FN \neq 0$, then the estimated graph omitted some, but not all, causal relationships and we take MCC = 0.

²Derived from an earlier definition by Matthews [33].

- 4. If TN = 0 **AND** FN = 0, then estimated graph is fully connected:
 - a) if TP = 0, then the true graph is fully sparse, so the algorithm is perfectly incorrect, and we take MCC = -1;
 - b) if $TP \neq 0$ **AND** FP = 0, then the true graph is also fully connected and we take MCC = 1.
 - c) if $TP \neq 0$ **AND** $FP \neq 0$, some, but not all, of the estimated causal dependencies are spurious and we take MCC = 0.

		c			
		0.1	0.2	0.3	0.4
	0.1	X	X	X	X
	0.2	X	X	X	
	0.3	X	X	X	
a	0.4	X	X		
	0.5	X	X		
	0.6	X			
	0.7	X			

Table 4-1. The stable autocorrelation (a) and cross-correlation (c) dependence coefficients identified for the one-dimensional model.

4.2. One-dimensional model

The results of the simulation study described in Section 3.1 are shown in Figure 4-1. For each simulation, we provided PCMCI with the correct maximum lag ($\tau_{\text{max}} = 2$) and set the threshold parameters to relatively stringent values ($\alpha_{PC} = 0.01$, $\alpha_G = 0.01$). Internal to PCMCI, we used a Gaussian partial correlation test for independence testing, as our data was generated from a linear-Gaussian VAR.

Recall from Section 3.1 that autocorrelated dependence is labeled coefficient a, and within-variable and between-variable cross-correlation dependence is labeled coefficient c. As Figure 4-1 shows, only a minority of a and c dependence coefficients were found to be stable. a was able to reach as high as 0.7, while c was only able to reach as high as 0.4. Table 4-1 shows the specific a and c combinations that were identified as stable in this model. The specific stable coefficient combinations would likely change with a different model formulation, e.g., different spatial dependence structures.

Figure 4-1 shows that PCMCI performed better with more time samples, but performance was limited by the particular a and c coefficients. The algorithm performed better where either coefficient was larger, but particularly when c was larger. For example, when c = 0.1, more time samples made little to no difference in performance beyond 250 samples.

In Figure 4-2, we show PCMCI performance as a function of autocorrelation. Figures 4-2a, 4-2b, and 4-2c depict this when c = 0.1, c = 0.2, and c = 0.3, respectively. Again we see that F_1 score increases as the a coefficient increases. Note the differently scaled Y-axes between the panels; the F_1 score reaches higher magnitudes when c is larger. This suggests that within the confines of a stable system, larger autocorrelation increase the signal-to-noise ratio, making the dynamics more easily identifiable. It does not appear that autocorrelation specifically is a detriment to structure identification.

In Figure 4-3, we show PCMCI performance as a function of cross-correlation. Figures 4-3a, 4-3b, and 4-3c depict this when a=0.1, a=0.2, and a=0.3, respectively. We more clearly see that F_1 score increases as the c coefficient increases. Note the differently scaled Y-axes between the panels; performance reaches higher magnitudes when a is larger. Like autocorrelation, larger cross-correlation increases performance, likely because of an improved signal-to-noise ratio. Larger autocorrelation and larger cross-correlation combined results in the best performance.

Finally, in Figure 4-4, we show PCMCI performance as a function of T time samples. Each data point includes all a and c values. We observe a clear pattern that PCMCI performance increases as a function of T, regardless of coefficient values.

4.3. Two-dimensional model

In this section, we present the results of the simulation study described in Sections 3.2. Recall that, for the two-dimensional simulations, we had only a single variable and that the complexity of the problem was controlled by the 3-by-3 neighborhood dynamics matrix, suitably expanded for the larger grid. For each simulation, we provided PCMCI with the correct maximum lag ($\tau_{max} = 1$) and set the threshold parameters to relatively stringent values ($\alpha_{PC} = 0.01$, $\alpha_G = 0.01$). Internal to PCMCI, we used a Gaussian partial correlation test for independence testing, as our data was generated from a linear-Gaussian VAR.

In Figure 4-5, we examine the effect of grid size (N) on both the F_1 and MCC scores, with other parameters fixed to $\sigma = 1.0$, $NDD = \frac{3}{9}$, and T = 1000. While we observe a high degree of variance in this plot, it is clear that performance degrades on larger grid sizes, though at a relatively slow rate if we recall that the problem dimensionality increases *quadratically* in N. As F_1 and MCC are highly correlated, we only depict MCC in subsequent figures. Appendix A features alternate performance measures.

Figure 4-6 depicts the effect of varying the sample length (T). We clearly observe a sub-linear growth in accuracy, as would be expected from the decreasing marginal information of additional samples.³ Figure 4-7 further depicts the effect of T for various values of grid size, N, and connectivity (NDD). Here we observe that neither grid size nor connectivity have significant impact on PCMCI performance, but that, as expected, there is a small decrease in performance as the grid size increases.

Figure 4-8 highlights the effect of graph density on PCMCI performance. From this plot, it is clear that PCMCI performance is marginally impacted by number of causal relationships increases, and increasing T removes these minimal effects. Comparing results columnwise, we again observe a relatively limited effect of grid size on our results. While Figure 4-8, clearly indicates that PCMCI is able to recover the true graph in the large sample limit, this provides limited guidance for analysts considering the use of causal discovery from data of limited sample size.

In Figure 4-9, we attempt to answer the question "how many samples will I need to expect success"? Because the threshold for "success" is problem dependent, we instead estimate the probability of MCC > m for various values of m. For moderately stringent thresholds ($m \approx 0.7$), we see that T = 500 samples appear sufficient for even large grid sizes, while even T = 1000 samples may be insufficient at highly stringent thresholds (m = 0.9). From these plots it is clear that, while average MCC performance may not vary significantly in grid size, the *dependability* of PCMCI clearly decreases rapidly in N.

Finally, Figure 4-10 investigates the effect of the innovation scale ($\sigma = \operatorname{sd}(\eta_{i,t})$) on PCMCI performance. Empirically, we observe no systematic effect of σ on performance: we hypothesize that this is because σ controls the magnitude of both the additive Gaussian innovations and the signal component $\tilde{A}\operatorname{vec}(X_t)$, leaving the effective signal-to-noise ratio of the problem unchanged. While we do not show this analytically

³Via general statistical principles, we expect MSE $\propto T^{-1/2}$, and note that MCC is a non-linear, but monotonic, function of estimation accuracy.

for the causal discovery problem, we do note that a similar phenomenon occurs in the estimation of VAR coefficients.⁴

Additional results, including analysis of the True Positive Rate (TPR), True Negative Rate (TNR), False Positive Rate (FPR), and False Negative Rate (FNR) appear in Appendix A. Those plots indicate few false positives across different simulation regimes and that decreases in MCC are primarily driven by false negatives, indicating large numbers of samples are necessary to correctly identify causal effects. While varying the PCMCI thresholding parameters α_{PC} and α_{G} may adjust the balance of false negatives and false positives, we do not explore the effect of those parameters in this work.

Additionally recalling that the variance of the OLS estimator is given by $Cov(vec(\hat{\boldsymbol{\beta}})) = (vec(\boldsymbol{X})vec(\boldsymbol{X})^T)^{-1} \otimes \sigma^2 \boldsymbol{I}$, we have $Cov(vec(\hat{\boldsymbol{\beta}})) \approx [\sigma^2(\boldsymbol{I} - \boldsymbol{A} \otimes \boldsymbol{A})^{-1}]^{-1} \otimes \sigma^2 \boldsymbol{I} = (\boldsymbol{I} - \boldsymbol{A} \otimes \boldsymbol{A})^{-1} \otimes \boldsymbol{I}$ which does not depend on σ .

⁴Briefly, let $\mathbf{X}_t = \mathbf{A}\mathbf{X}_{t-1} + \mathbf{\eta}$ for $\mathbf{\eta} \sim \mathcal{N}(0, \sigma^2 I)$. Then $\operatorname{Cov}(\mathbf{X}_t) = \operatorname{Cov}(\mathbf{A}\mathbf{X}_{t-1} + \mathbf{\eta})$ $= \mathbf{A}\operatorname{Cov}(\mathbf{X}_t)\mathbf{A}^T + \sigma^2 \mathbf{I}$ $\implies \operatorname{vec}(\operatorname{Cov}(\mathbf{X}_t)) = \sigma^2 (\mathbf{I} - \mathbf{A} \otimes \mathbf{A})^{-1}.$

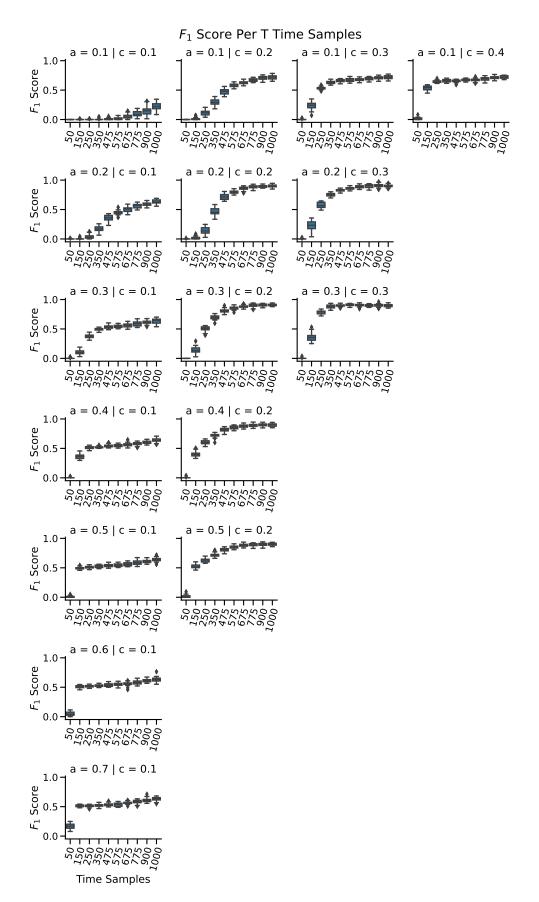


Figure 4-1. F_1 scores for the One-Dimensional model of Sections 3.1 and 3.1. Coefficients a and c represent autocorrelation and cross-correlation dependence coefficients, respectively, where cross-correlation relates to both variable-to-variable and cell-to-cell dependencies. Only stable coefficient combinations are shown. PCMCI performs better with more time samples, larger a values, and larger c values. When a or c are sufficiently large, the system becomes unstable.

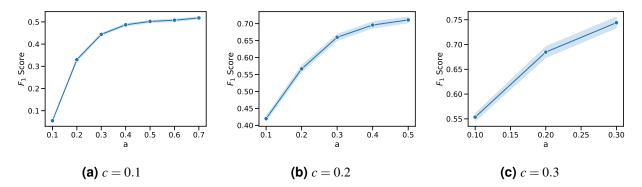


Figure 4-2. F_1 score results from the one-dimensional spatial example with varying autocorrelation, a, and constant cross-correlation, c. Each data point includes all possible T time samples. Note the different Y axes.

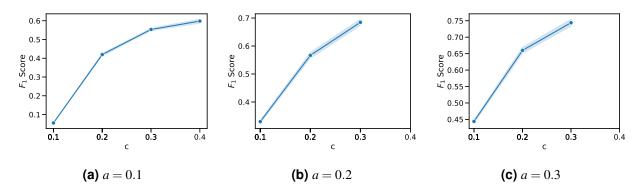


Figure 4-3. F_1 score results from the one-dimensional spatial example with varying cross-correlation, c, and constant autocorrelation, a. Each data point includes all possible T time samples. Note the different Y axes.

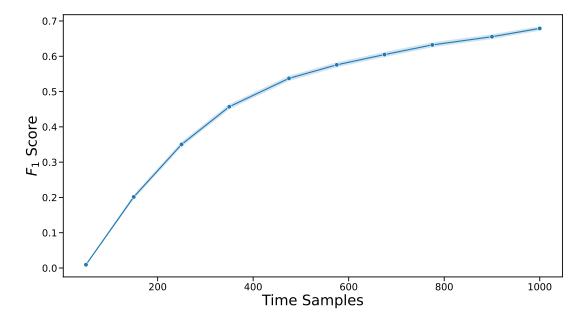


Figure 4-4. F_1 score results from the one-dimensional spatial example with varying T time samples. Each data point includes all possible a and c dependence coefficients.

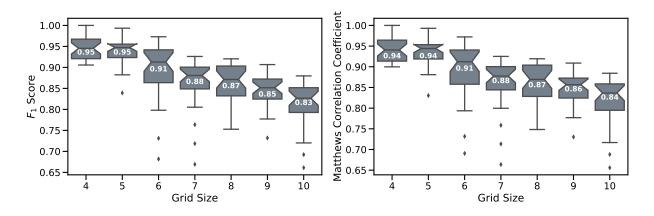


Figure 4-5. Effect of grid size (N) on PCMCI F_1 and MCC scores. Both metrics decrease relatively slowly in N. Other simulation parameters are fixed to $\sigma = 1.0$, $NDD = \frac{3}{6}$, and T = 1000.

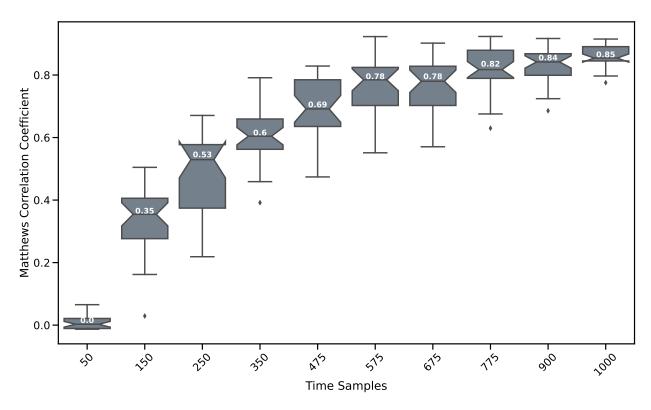


Figure 4-6. Effect of increasing sample size (T) on PCMCI performance (MCC). Performance increases sublinearly in T, with T>575 being necessary to obtain acceptable performance (MCC >0.7). Box labels report median MCC across replicates. Other simulation parameters fixed as N=10, $\sigma=1.0$, and NDD $=\frac{6}{9}$.

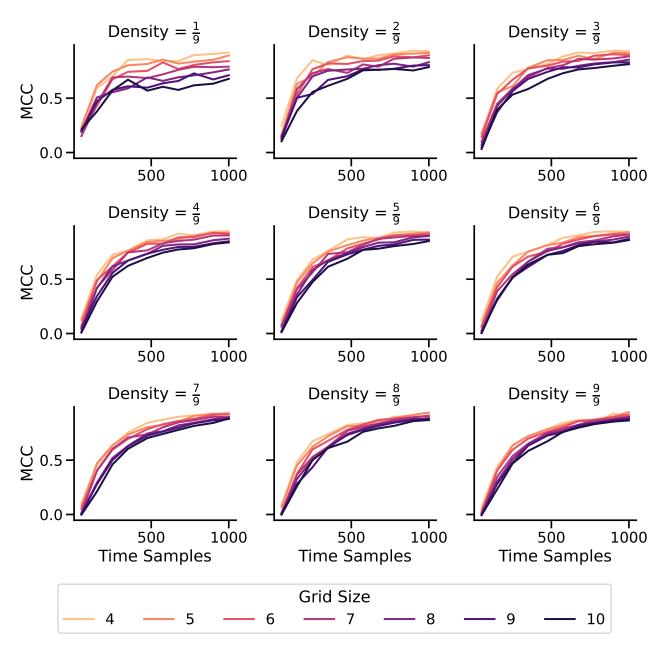


Figure 4-7. Effect of sample size, (T) grid size, (N), and neighborhood dependence density on PCMCI performance (MCC). For sufficiently large sample sizes, PCMCI is able to consistently recover the true graph structure; the effect of grid size and NDD are less pronounced than T. Values shown are mean performance over 30 replicates.

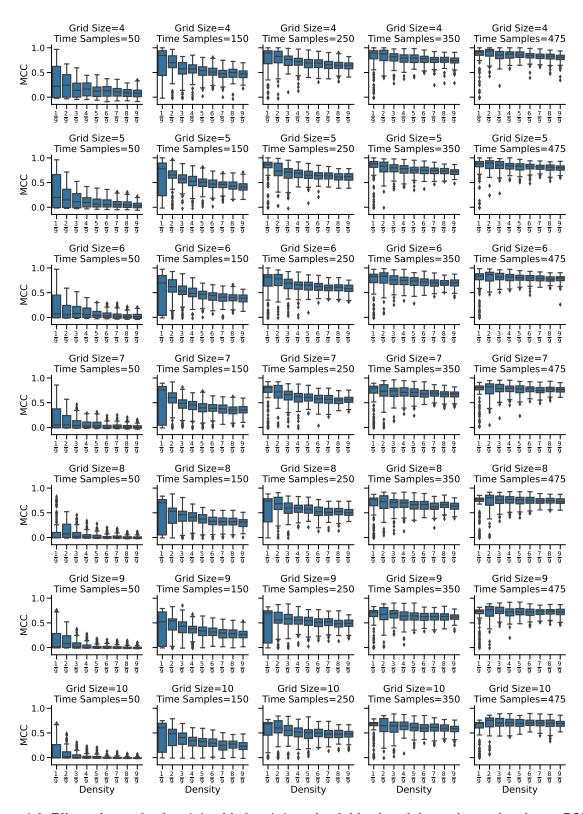


Figure 4-8. Effect of sample size, (T) grid size, (N), and neighborhood dependence density on PCMCI performance (MCC). For sufficiently large sample sizes, PCMCI is able to consistently recover the true graph structure; the effect of grid size and NDD are limited. Values shown are mean performance over 30 replicates. $\sigma=1$ for all simulations.

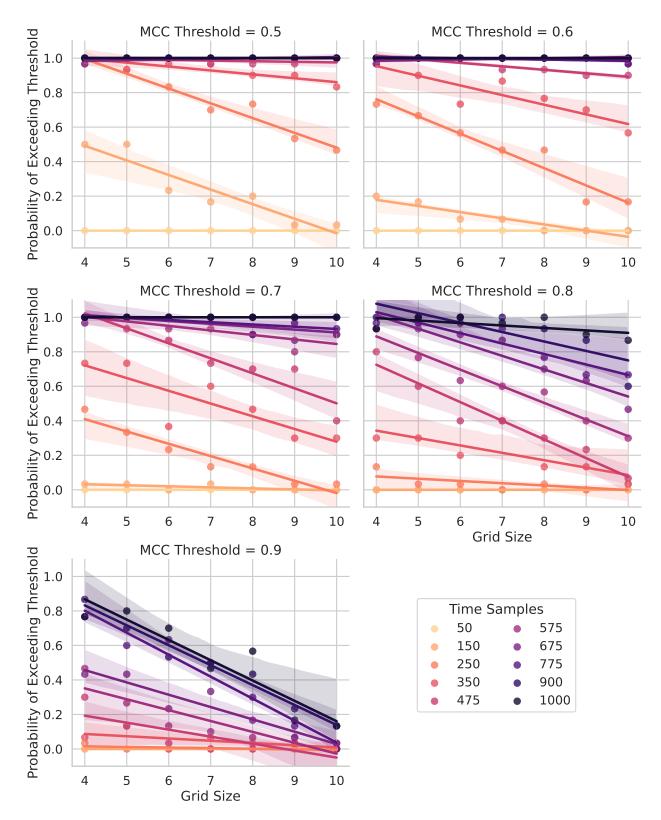


Figure 4-9. Probability of PCMCI Success as a function of grid size N and sample size T, with success defined as MCC above a user-defined threshold. Results are empirical probabilities over 30 replicates: σ and neighborhood density are fixed to 1.0 and $\frac{6}{9}$ respectively. Lines depict a simple linear model of grid size on success probability, with shaded regions depicting (non-multiplicity adjusted) confidence intervals.

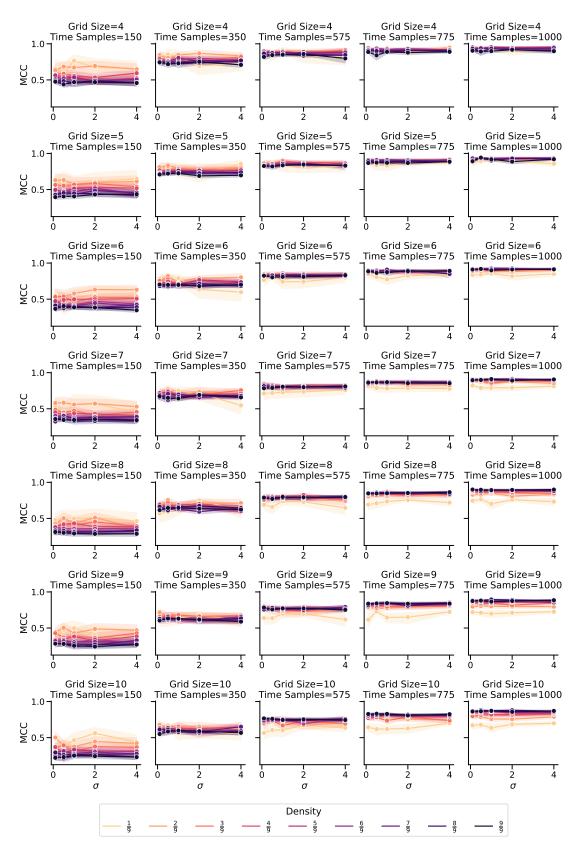


Figure 4-10. Effect of Innovation Magnitude (σ) on PCMCI performance (MCC). Changing σ appears to have no systematic effect on PCMCI performance.

5. DISCUSSION

In this work, we investigated the performance of the PCMCI causal network discovery algorithm for linear-VAR systems in one- and two-dimensional space. We varied the length of the observed time series, the size of the underlying grid, and the density of the underlying causal graph, and found significant effects of each. Our results provide a robust characterization of PCMCI performance on spatiotemporal systems and highlight several avenues of future inquiry.

Most notably, we found that $T \approx 1000$ samples were necessary to for consistent high-accuracy causal discovery across the various scenarios we considered (see Figure 4-4 and Figure 4-9). While this is consistent with the asymptotic consistency of PCMCI, these extreme sample sizes are unrealistic for the climate data analytics motivating this study. Note that we restricted our analysis to linear-Gaussian systems, which enables PCMCI to reduce the difficult problem of testing conditional independence to the relatively easier problem of estimating partial correlations. While it is possible to use PCMCI with more general conditional independence tests, these tests have a far higher sample complexity, and would require a far greater sample size to achieve consistent performance.

By contrast, the effect of the grid size was relatively minor, suggesting that performance gains may attainable through clever use of this spatial structure. Changing the number of true causal effects had notable impacts on certain performance measures, but further work is needed to determine whether this scenario is inherently more difficult for causal discovery or whether it is an artifact of the specific accuracy measures we used, *e.g.*, the number of true positives for an empty graph.

We note that in our study of the one-dimensional model, we found that PCMCI tolerated high autocorrelation well. This result is somewhat unexpected, given previous work showing that causal discovery algorithms tend to handle autocorrelation poorly. However, PCMCI was developed to be robust to autocorrelation [10]. The clearest conclusion, apart from the aforementioned benefits of more time samples, was that larger causal dependence coefficients were beneficial, regardless of whether they were autocorrelational or cross-correlational coefficients.

Finally, our study of the two-dimensional model also provided several computational advances that may be of independent interest, including characterization of the sliding dot product and VAR representations of our model, an easy-to-implement check for stability of the resulting VAR process, and an effective algorithm for sampling from the space of stable dynamics.

As shown in Figure 4-9, the probability of "successful" graph recovery is highly sensitive to both the sample size and the grid size. As the number of potential causal parents for a single grid cell increases quadratically in N, this is perhaps unavoidable. More generally, causal discovery algorithms are known to suffer from the curse of dimensionality, particularly when applied on the grid-level in spatiotemporal systems as the both the potential causal parents and the number of grid cells studied increase rapidly in the grid size [10, 15, 22, 25].

In the climate context, the underlying grids are far larger than those considered in this study, necessitating extremely large sample sizes. Unfortunately, our causal stationarity assumptions (Assumption T2 and S2) are less likely to hold over these extended time frames. To avoid this problem, some works have artificially reduced the problem dimensionality by replacing grid cells with pre-defined regions of climatological interest [15, 21, 22, 25]. They made attempts to benchmark their results with either simulated or theoretical expectations. However, their simulations were not of grid-cell-level causal dynamics, as ours are, and their studies on natural climate data could not be benchmarked rigorously. Finally, we note that these approaches are only appropriate for long-term climate analyses in which well-defined spatially-stable statistically-regular modes are the objects of study. We do not expect these approaches to perform well when studying "one-off" climate events, in which relevant regions are rarely known *a priori*, making dimension reduction a far more challenging task.

Finally, we note that our study only considered samples from the stationary distribution of a linear system driven by Gaussian innovations. As a result, our simulated data is itself Gaussian and does not reflect structures that may be found in climate data, *e.g.*, the El Niño Southern Oscillation (ENSO) or, on shorter scales, major storms. It is unclear how PCMCI would perform when applied to these stable structures, as they have complex spatiotemporal dynamics.

Causal discovery is an important aspect of modern climate research and there is a need for algorithms that can scalably and accurately determine causal structure from grid-level data. While PCMCI is quite datahungry on large grids and observational climate data are quite limited, additional insights can be gleaned from the analysis of large simulation ensembles. Currently, PC-family algorithms do not incorporate spatial structure: in future work, we hope to investigate the use of spatial structure to reduce the dimensionality of the causal discovery problem.

Causal discovery remains a challenging task, particularly in the climate domain. As simulation and observational data continues to grow in size and scope, there is a pressing need for approaches that can perform robustly at a range of time- and spatial-scales, ranging from storm tracking to diffusion of volcanic aerosols to long-term natural and anthropogenic climate changes. The benchmarking techniques and simulations of this paper give insight into the weaknesses of current approaches and suggest new avenues of causal discovery research.

BIBLIOGRAPHY

- [1] Randal S Olson, William La Cava, Patryk Orzechowski, Ryan J Urbanowicz, and Jason H Moore. PMLB: a large benchmark suite for machine learning evaluation and comparison. *BioData mining*, 10:1–13, 2017.
- [2] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255. IEEE, 2009.
- [3] Jeyan Thiyagalingam, Mallikarjun Shankar, Geoffrey Fox, and Tony Hey. Scientific machine learning benchmarks. *Nature Reviews Physics*, 4(6):413–420, 2022. doi:10.1038/s42254-022-00441-7.
- [4] Jeyan Thiyagalingam, Kuangdai Leng, Samuel Jackson, Juri Papay, Mallikarjun Shankar, Geoffrey Fox, and Tony Hey. Scimlbench: A benchmarking suite for ai for science, 2021. URL https://github.com/stfc-sciml/sciml-bench.
- [5] Osman Balci. Verification, Validation, and Certification of Modeling and Simulation Applications. *Proceedings of the 2003 Winter Simulation Conference*, 2003, 1:150–158, 2003. doi:10.1109/wsc.2003.1261418.
- [6] William L. Oberkampf and Christopher J. Roy. *Verification and Validation in Scientific Computing*. Cambridge University Press, 2010. ISBN 9780511760396. doi:10.1017/cbo9780511760396.016.
- [7] National Research Council. Assessing the Reliability of Complex Models: Mathematical and Statistical Foundations of Verification, Validation, and Uncertainty Quantification. The National Academies Press, Washington, DC, 2012. ISBN 978-0-309-25634-6. doi:10.17226/13395.
- [8] R G Sargent. Verification and validation of simulation models. *Journal of Simulation*, 7(1):12–24, 2013. ISSN 1747-7778. doi:10.1057/jos.2012.20.
- [9] Jonas Peters, Dominik Janzing, and Bernhard Schlkopf. *Elements of Causal Inference: Foundations and Learning Algorithms*. The MIT Press, 2017. ISBN 9780262037310.
- [10] J. Runge. Causal network reconstruction from time series: From theoretical assumptions to practical estimation. *Chaos: An Interdisciplinary Journal of Nonlinear Science*, 28(7):075310, 2018. ISSN 1054-1500. doi:10.1063/1.5025050.
- [11] Jakob Runge, Peer Nowack, Marlene Kretschmer, Seth Flaxman, and Dino Sejdinovic. Detecting and quantifying causal associations in large nonlinear time series datasets. *Science Advances*, 5(11): 4996—5023, 2019. doi:https://doi.org/10.1126/sciadv.aau4996. URL http://advances.sciencemag.org/.
- [12] Judea Pearl and Dana Mackenzie. *The Book of Why*. Basic Books, New York, 2018. ISBN 978-0-465-09760-9.

- [13] Michael Eichler. In AISTATS 2010: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, volume 9 of Proceedings of Machine Learning Research, pages 193–200, 2010. URL https://proceedings.mlr.press/v9/eichler10a.html.
- [14] Judea Pearl. Causal Diagrams for Empirical Research. *Biometrika*, 82(4):669–688, 1995. ISSN 0006-3444. doi:10.2307/2337329.
- [15] Jakob Runge, Vladimir Petoukhov, Jonathan F. Donges, Jaroslav Hlinka, Nikola Jajcay, Martin Vejmelka, David Hartman, Norbert Marwan, Milan Paluš, and Jürgen Kurths. Identifying causal gateways and mediators in complex spatio-temporal systems. *Nature Communications*, 6(1):8502, 2015. doi:10.1038/ncomms9502.
- [16] Jakob Runge, Sebastian Bathiany, Erik Bollt, Gustau Camps-Valls, Dim Coumou, Ethan Deyle, Clark Glymour, Marlene Kretschmer, Miguel D. Mahecha, Jordi Munoz-Mari, Egbert H. van Nes, Jonas Peters, Rick Quax, Markus Reichstein, Marten Scheffer, Bernhard Scholkopf, Peter Spirtes, George Sugihara, Jie Sun, Kun Zhang, and Jakob Zscheischler. Inferring causation from time series in Earth system sciences. *Nature Communications*, 10(1), 2019. ISSN 20411723. doi:10.1038/s41467-019-10105-3.
- [17] Peter Spirtes, Clark Glymour, and Richard Scheines. *Causation, Prediction, and Search*, volume 81 of *Lecture Notes in Statistics*. Springer, 1993. doi:10.1007/978-1-4612-2748-9.
- [18] Shohei Shimizu, Patrik O. Hoyer, Aapo Hyvarinen, and Antti Kerminen. A Linear Non-Gaussian Acyclic Model for Causal Discovery. *Journal of Machine Learning Research*, 7(72):2003–2030, 2006. URL https://www.jmlr.org/papers/volume7/shimizu06a/shimizu06a.pdf.
- [19] Jie Sun and Erik M. Bollt. Causation entropy identifies indirect influences, dominance of neighbors and anticipatory couplings. *Physica D: Nonlinear Phenomena*, 267:49–57, 2014. ISSN 0167-2789. doi:10.1016/j.physd.2013.07.001.
- [20] Jie Sun, Dane Taylor, and Erik M Bollt. Causal Network Inference by Optimal Causation Entropy. *SIAM Journal on Applied Dynamical Systems*, 14(1):73–106, 2015. doi:10.1137/140956166.
- [21] Marlene Kretschmer, Dim Coumou, Jonathan F. Donges, and Jakob Runge. Using Causal Effect Networks to Analyze Different Arctic Drivers of Midlatitude Winter Circulation. *Journal of Climate*, 29(11):4069–4081, 2016. ISSN 0894-8755. doi:10.1175/jcli-d-15-0654.1.
- [22] Peer Nowack, Jakob Runge, Veronika Eyring, and Joanna D. Haigh. Causal networks for climate model evaluation and constrained projections. *Nature Communications 2020 11:1*, 11(1):1—11, 2020. ISSN 2041-1723. doi:10.1038/s41467-020-15195-y. URL http://www.nature.com/articles/s41467-020-15195-y.
- [23] Zachary S. Kaufman, Nicole Feldl, Wilbert Weijer, and Milena Veneziani. Causal Interactions Between Southern Ocean Polynyas and High-Latitude Atmosphere-Ocean Variability. *Journal of Climate*, 33(11):4891–4905, 2020. ISSN 0894-8755. doi:10.1175/jcli-d-19-0525.1.

- [24] Christopher Krich, Jakob Runge, Diego G. Miralles, Mirco Migliavacca, Oscar Perez-Priego, Tarek El-Madany, Arnaud Carrara, and Miguel D. Mahecha. Estimating causal networks in biosphere–atmosphere interaction with the PCMCI approach. *Biogeosciences*, 17(4):1033–1061, 2020. doi:10.5194/bg-17-1033-2020.
- [25] Xavier-Andoni Tibau, Christian Reimers, Andreas Gerhardus, Joachim Denzler, Veronika Eyring, and Jakob Runge. A spatiotemporal stochastic climate model for benchmarking causal discovery methods for teleconnections. *Environmental Data Science*, 1, 2022. doi:10.1017/eds.2022.11.
- [26] Andreas Gerhardus and Jakob Runge. LPCMCI: Causal Discovery in Time Series with Latent Confounders. In *Advances in Neural Information Processing Systems*, volume 33, pages 12615–12625. Curran Associates, Inc., 2020. doi:10.5194/egusphere-egu21-8259. URL https://proceedings.neurips.cc/paper/2020/file/94e70705efae423efda1088614128d0b-Paper.pdf.
- [27] Jakob Runge. Discovering contemporaneous and lagged causal relations in autocorrelated nonlinear time series datasets. In *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence*, volume 124 of *Proceedings of Machine Learning Research*, pages 1388–1397. PMLR, 2020. URL https://proceedings.mlr.press/v124/runge20a.html.
- [28] Yi Deng and Imme Ebert-Uphoff. Weakening of atmospheric information flow in a warming climate in the Community Climate System Model. *Geophysical Research Letters*, 41(1):193–200, 2014. ISSN 1944-8007. doi:10.1002/2013gl058646.
- [29] Imme Ebert-Uphoff and Yi Deng. Causal Discovery from Spatio-Temporal Data with Applications to Climate Science. *2014 13th International Conference on Machine Learning and Applications*, pages 606–613, 2014. doi:10.1109/icmla.2014.96.
- [30] James D. Hamilton. Time Series Analysis. Princeton University Press, 1994. ISBN 9780691042893.
- [31] Nancy Chinchor. MUC-4 Evaluation Metrics. In *Proceedings of the 4th Conference on Message Understanding*, MUC4'92, page 22–29, USA, 1992. Association for Computational Linguistics. ISBN 1558602739. doi:10.3115/1072064.1072067. URL https://doi.org/10.3115/1072064.1072067.
- [32] Davide Chicco. Ten quick tips for machine learning in computational biology. *BioData Mining*, 10(1):35, 2017. ISSN 1756-0381. doi:10.1186/s13040-017-0155-3.
- [33] B.W. Matthews. Comparison of the predicted and observed secondary structure of t4 phage lysozyme. *Biochimica et Biophysica Acta (BBA) Protein Structure*, 405(2):442–451, 1975. ISSN 0005-2795. doi:https://doi.org/10.1016/0005-2795(75)90109-9. URL https://www.sciencedirect.com/science/article/pii/0005279575901099.

APPENDIX A. Additional Simulation Results: Two-Dimensional Model

In this section, we depict various performance rates of PCMCI in our two-dimensional simulation study (Section 3.2). Here we report:

$$FDR = \frac{FP}{TP + FP}$$
(False Discovery Rate, Figure A-1)
$$TPR = \frac{TP}{TP + FN} = \frac{TP}{P}$$
(True Positive Rate, Figure A-2)
$$FNR = \frac{FN}{TP + FN} = \frac{FN}{P}$$
(False Negative Rate, A-3)
$$TNR = \frac{TN}{TN + FP} = \frac{TN}{N}$$
(True Negative Rate, A-4)
$$FPR = \frac{FP}{TN + FP} = \frac{FP}{N}$$
(False Positive Rate, A-5)

where FDR is the false discovery rate; TP, FP, TN, FN are the number of true positives, false positives, true negatives, and false negatives, respectively; and P, N are the number of edges and non-edges in the true graph.

As with the one-dimensional model, PCMCI exhibits a bias towards non-discovery, with low true and false positive rates across scenarios. The FPR is almost always kept near 0, indicating that we can have a high degree of confidence in the causal effects identified by PCMCI, but that it has limited statistical power at moderate sample sizes.

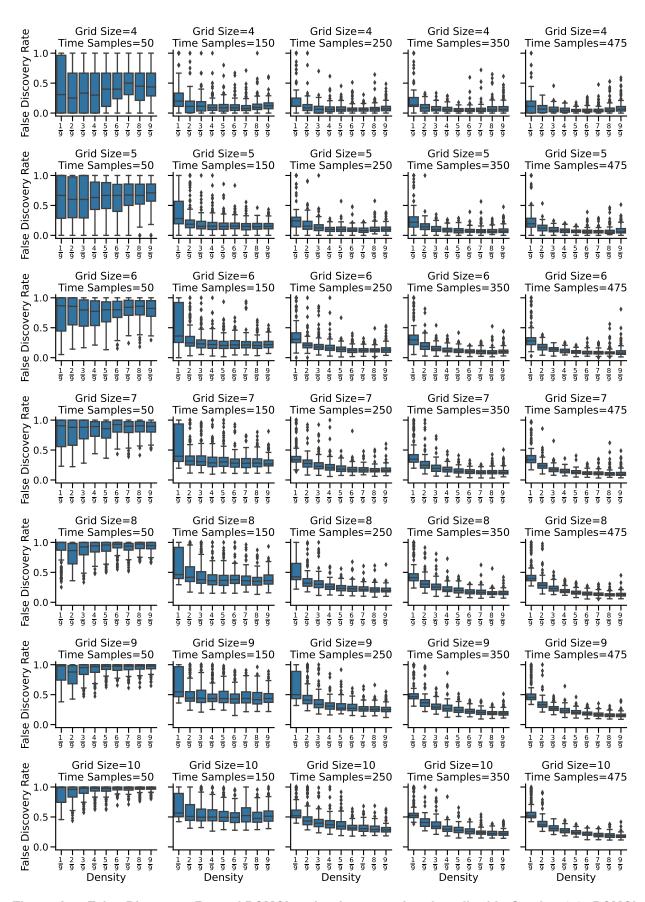


Figure A-1. False Discovery Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits low FDR for T>50. FDR decreases with the number of causal effects (density) and with increasing time samples.

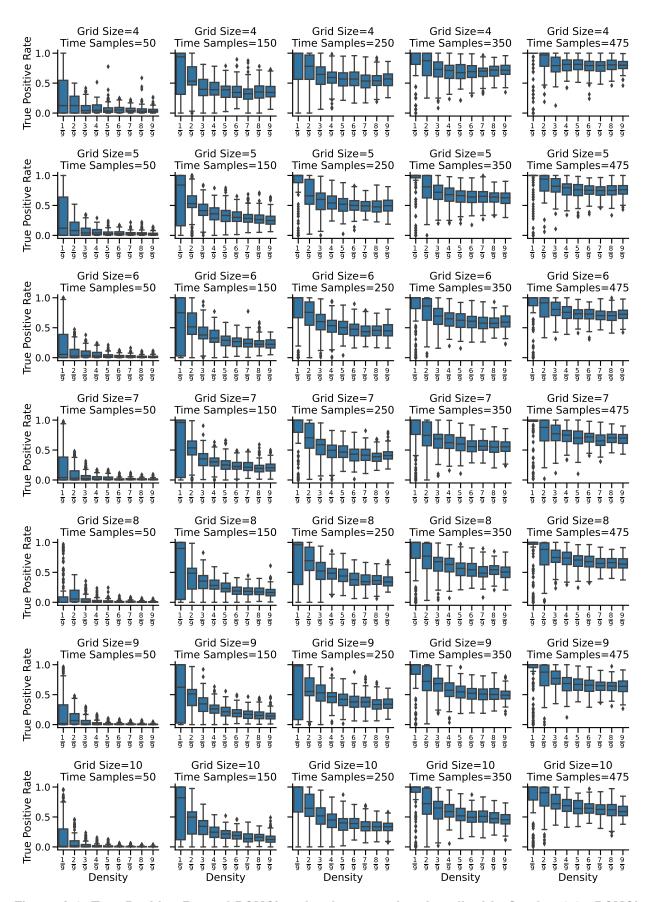


Figure A-2. True Positive Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits low true positive rates for $T_4 < 350$. TPR decreases with the number of causal effects and with increasing grid sizes.

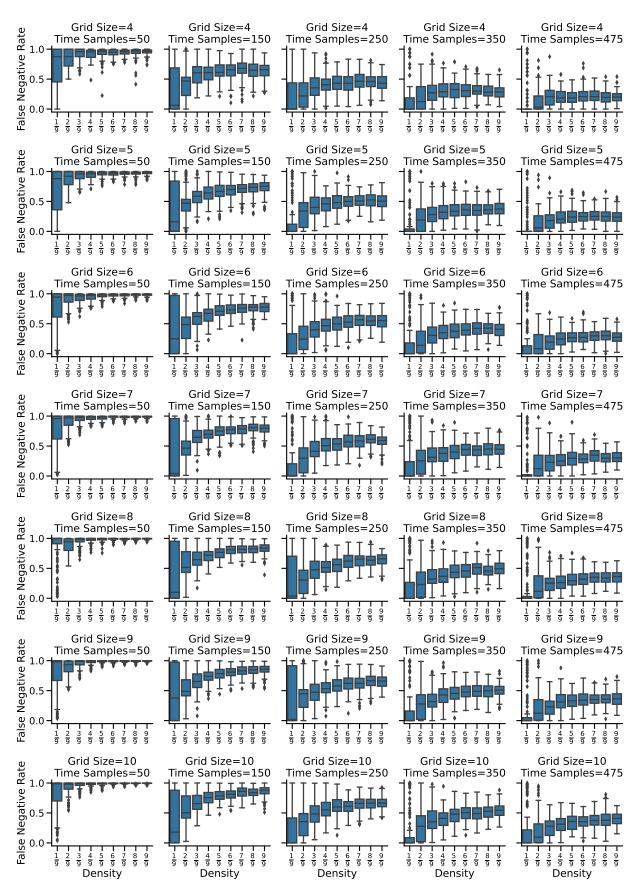


Figure A-3. False Negative Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits relatively high false negative rates in all scenarios, indicating low statistical power. FNR generally increases with the number of causal effects and with increasing grid sizes.

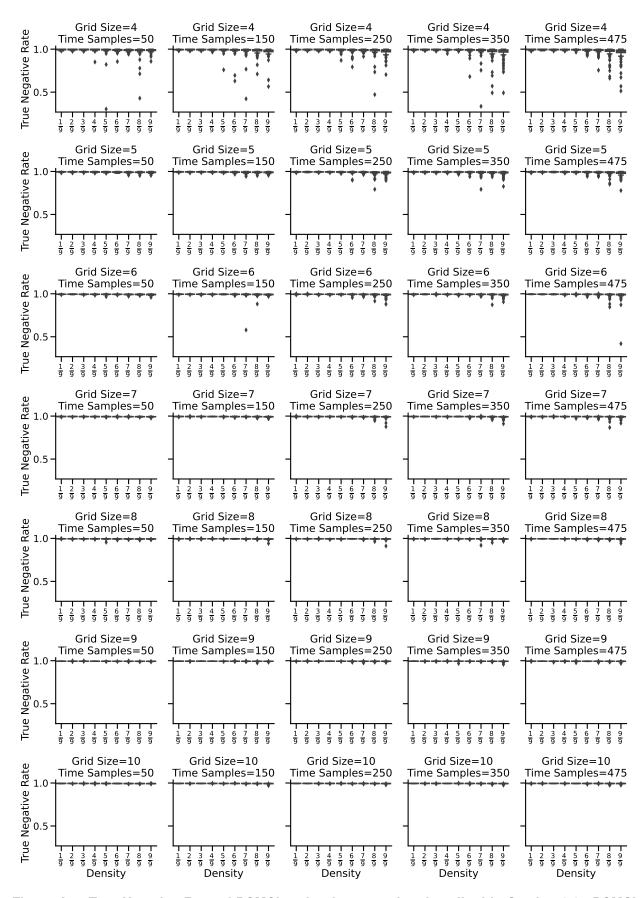


Figure A-4. True Negative Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits near perfect true negative rates in all scenarios. To the extent it varies, TNR decreases with the number of causal effects and with decreasing grid sizes.

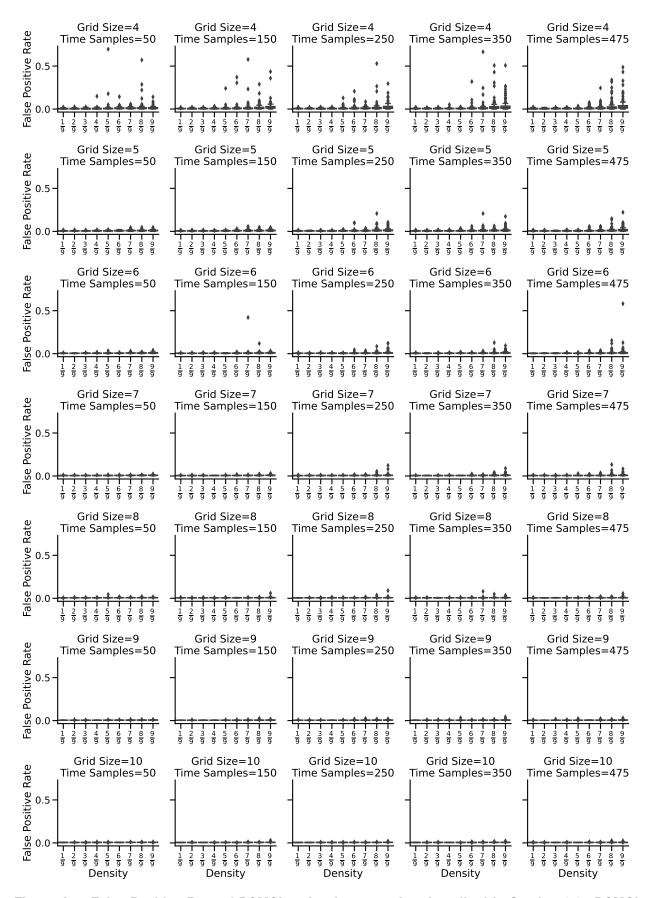


Figure A-5. False Positive Rate of PCMCI under the scenarios described in Section 3.2. PCMCI consistently exhibits near perfect false positive rates in all scenarios. To the extent it varies, FPR increases with the number of causal effects and with decreasing grid sizes.

DISTRIBUTION

Email—Internal

Name	Org.	Sandia Email Address
Ronald Oldfield	1441	raoldfi@sandia.gov
Matt Peterson	1441	mgpeter@sandia.gov
Kara Peterson	1442	kjpeter@sandia.gov
Jay Brown	5493	jbrown2@sandia.gov
Aubrey Eckert	5573	acecker@sandia.gov
Lyndsay Shand	5573	lshand@sandia.gov
Irina Tezaur	8734	ikalash@sandia.gov
Meredith G.L. Brown	8931	merbrow@sandia.gov
Diana Bull	8931	dlbull@sandia.gov
Technical Library	1911	sanddocs@sandia.gov

Hardcopy—Internal

Number of Copies	Name	Org.	Mailstop
1	L. Martin, LDRD Office	1910	0359

Hardcopy—External

Number of Copies	Name(s)	Company Name and Company Mailing Address	
1			



Sandia National Laboratories is a multimission laboratory managed and operated by National Technology & Engineering Solutions of Sandia LLC, a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

Space-Time Causal Discovery in Earth Sys-

tem Science: A Local Stencil Learning Ap-

proach

Publication Notes 7.1

Citation: Nichol, J. Jake, et al. "Space-Time Causal Discovery in Earth System

Science: A Local Stencil Learning Approach." Journal of Geophysics Research:

Machine Learning and Computation, under review.

Publication date: N/A

Conference: N/A

Formatting: The original text has been preserved as much as possible while still

adhering to the formatting requirements of this dissertation.

Data and Software Availability: The paper is currently under review and not yet

publicly available.

Funding: This work is supported by Sandia Earth Science Investment Area Labo-

ratory Directed Research and Development funding. Sandia National Laboratories

is a multimission laboratory managed and operated by National Technology and

Engineering Solutions of Sandia, LLC., a wholly owned subsidiary of Honeywell

157

International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA-0003525.

7.2 Abstract

Causal discovery tools enable scientists to infer meaningful relationships from observational data, spurring advances in fields as diverse as biology, economics, and climate science. Despite these successes, the application of causal discovery to space-time systems remains immensely challenging due to the high-dimensional nature of the data. For example, in climate sciences, modern observational temperature records over the past few decades regularly measure thousands of locations around the globe. To address these challenges, we introduce Causal Space-Time Stencil Learning (Causal Space-Time Stencil Learning (CaStLe)), a novel meta-algorithm for discovering causal structures in complex space-time systems. CaStLe leverages regularities in local space-time dependencies to learn governing global dynamics. This local perspective eliminates spurious confounding and drastically reduces sample complexity, making space-time causal discovery practical and effective. For causal discovery, CaStLe flexibly accepts any appropriately adapted time series causal discovery algorithm to recover local causal structures. These advances enable causal discovery of geophysical phenomena that were previously unapproachable, including non-periodic, transient phenomena such as volcanic eruption plumes. Regularities in local space-time dependencies are transformed into informative spatial replicates, which actually improves CaStLe's performance when applied to ever-larger spatial grids. We successfully apply CaStLe to discover the atmospheric dynamics governing the climate response to the 1991 Mount Pinatubo volcanic eruption. We provide validation experiments to demonstrate the effectiveness of CaStLe over existing causal-discovery frameworks on a range of geophysics-inspired benchmarks while identifying the method's limitations and domains where its assumptions may not hold.

Plain Language Summary

We introduce a new method for learning the dynamics of causal systems, that is, the physical rules that define a system's behavior. While this task, *causal discovery*, is not new, existing tools are ill-suited for many large geophysics datasets. Current state-of-the-art approaches use statistical techniques to search for causal relationships between all aspects of a system, examining billions of possible causal effects, or simplifying the data by focusing on the most important variables. Instead of an exhaustive search or oversimplifying the data, we incorporate basic physical principles—requiring effects to be "local" and "uniform"—to massively simplify the causal discovery problem. We demonstrate that our approach can recover known geophysical dynamics by applying it to the 1991 Mt. Pinatubo eruption, validating its ability to uncover space-time causal structure from observational data.

7.3 Introduction

Explaining the causal dynamics that govern geophysical phenomena is paramount in the Earth sciences. Climate models, for example, critically depend on understanding both local and global causal pathways to model the complex Earth system. Understanding short- and long-term consequences of the Earth system's behavior is essential for future model development, our scientific knowledge, and preparing for the future. More specifically, in atmospheric science, we know the initial state of specific wind modes, such as the quasi-biennial oscillation or the Brewer-Dobson circulation, dramatically affects the later evolution and impact of volcanic eruptions, major wildfires, or geoengineering efforts such as stratospheric aerosol injection (Hitchman et al., 1994; Jones et al., 1998; Aquila et al., 2014; Gray et al., 2018).

Traditional statistical methodologies, while providing valuable insights, often fall short of capturing the complex causal relationships inherent in geophysical systems. Causal models are hard-won and often represent the culmination of many decades of research. Causal discovery tools aim to accelerate the discovery of these relationships using statistically-rigorous techniques to separate predictable, but indirect, statistical relationships from direct causal connections. Causal discovery has been successful across the sciences, providing new understandings of climate, biological, genetic, neural, and other dynamical systems (Ebert-Uphoff and Deng, 2012; Sugihara et al., 2012; Neto et al., 2010; Zhang et al., 2011b; Kamiński et al.,

2001; Tsonis et al., 2017). However, applying existing causal methods to space and time structured data remains limited due to the complexity and scale of such systems.

This work presents a novel causal discovery methodology that overcomes these challenges to recover networks describing local causal structures from gridded data. A fundamental insight driving the present work is that in many complex systems, global phenomena—whether climate teleconnections, brain functional networks, or ecosystem dynamics—emerge from countless repeated and structured local interactions. We can better understand how complex global patterns arise by accurately capturing these foundational local structures.

Today's Earth science measurement and modeling capabilities provide a wealth of data for studying our planet's complex dynamics. However, due to the immense complexity of these dynamics, simple analyses provide only a limited understanding of the data. Causal discovery tools offer the ability to understand finer mechanistic details via causal graphs' simplicity, interpretability, and flexibility. causal discovery is a field that utilizes algorithmic causal inference to identify causal models as dependencies between fields of interest, which are often represented as a directed acyclic graph. Causal graphs let us analyze the space-time evolution of fields of interest and causal discovery can estimate them without requiring hypothesized physical models. Insights gleaned from causal discovery can further inform physical models, validate simulations against observational data, and identify future research questions.

While causal discovery show considerable promise for addressing problems in the Earth sciences, the enormous size and scope of Earth science data have limited its applications. For example, atmospheric data often contains hundreds of thousands of grid cells, each with several orders of magnitude fewer observations in time. That imbalance is one aspect of the curse of dimensionality (Bellman, 1957; Bühlmann and Geer, 2011), where high dimensionality relative to sample size challenges conventional statistical methods and renders many forms of inference, including causal discovery, unreliable without dimensionality reduction. Despite these obstacles, causal discovery has been successfully applied in Earth science (Deng and Ebert-Uphoff, 2014; Runge et al., 2015c; Capua et al., 2019, 2020; Nowack et al., 2020a; Krich et al., 2020; Galytska et al., 2022; Tibau et al., 2022; O'Kane et al., 2024; Zhao et al., 2024), primarily via dimensionality reduction techniques to reduce the number of relationships to estimate. Those contributions identified teleconnection pathways to recover large, periodic climate modes and their effects. While a dimensionality reduction approaches can be practical, the analysis of local effects has been considered challenging and generally avoided due to the curse of dimensionality (Ebert-Uphoff and Deng, 2012; Runge et al., 2015c; Nowack et al., 2020a).

In contrast to dimensionality reduction methods that marginalize large amounts of information, our work leverages the known locality in space-time systems to harness *informative spatial replicates*, i.e., repeating space-time relationships, without loss of local structural information, to identify local causal graphs. These advances

enables us to approach problem classes in space-time systems that are typically intractable with prior art—both in terms of performance and algorithmic efficiency. We highlight two features of Causal Space-Time Stencil Learning (CaStLe) that are useful contributions to causal discovery for geoscience problems: the ability to learn grid-level relationships instead of regional relationships from reduced dimensional data (e.g. principal components or modes) and the ability to handle dynamic, advective processes.

Prior causal discovery work in Earth science has primarily focused on large-scale regional phenomena, such as the El Niño Southern Oscillation. These patterns—generally consistent in their spatial distribution and periodic in nature—are well suited to global dimensionality reduction techniques, which project fields onto a small number of modes. While global teleconnections are crucial research areas, they ultimately emerge from local causal interactions. However, dimensionality reduction sacrifices critical local information, making it impossible to see how local structures give rise to global patterns. CaStLe reduces problem complexity in a fundamentally different way: By identifying and leveraging the repeating local structures, it preserves the relationships at the grid level while remaining applicable to spacetime systems that exhibit multiscale organization.

Typical dimensionality reduction approaches to causal discovery decrease the data space from many grid cells to a few regional modes and uses many observations, resulting in a *little p*, *large n* problem, where p is the number of variables and n is the number of data points. In contrast, phenomena that evolve dynami-

cally in space or occur rarely, like volcanic plumes, are harder to analyze and often have few data points. Such problems are *large p, little n*. CaStLe makes causal discovery of the space-time evolution of these phenomena tractable for the first time by leveraging the gridded sample space, avoiding the marginalization that reduces many grid cells into a single time series per regional mode, and recovering interpretable space-time causal structures.

This work's primary case study is the 1991 Mount Pinatubo eruption. It injected a plume of aerosols into the stratosphere, which then advected around the tropical zone before dispersing northward and eventually diffusing around the globe. This example demonstrates the characteristics of the unique, transient problem class, has an established research history, and exhibits dynamics verifiable with a known causal driver: stratospheric wind.

We introduce a new Earth system causal network, the *causal stencil graph*, which describes local space-time causal structures between adjacent locations, and a new estimation methodology, **Ca**usal **S**pace-Time **Stencil Learning** (CaStLe), that is capable of describing local mechanistic pathways in space and time between grid cells. Grid-level causal discovery in high dimensional space-time data has been previously considered intractable due to the curse of dimensionality (Nowack et al., 2020a; Tibau et al., 2022). Though demonstrated with climate model output, our methodology applies to any space-time system where local physical interactions drive global behavior, including fluid dynamics, biological pattern formation, or material transport processes.

CaStLe combines modern causal discovery with classical physics-based principles, namely spatial and temporal locality, to accurately perform causal discovery on large spatial domains. Our novel local-coordinate-space projection does not marginalize any data points, such that local causal information is lost, which is a common sacrifice of other space-time dimension reduction techniques such as weighted averaging or principal component analysis (PCA). This preservation of local information is crucial because global-scale phenomena in complex systems emerge from interactions at smaller scales. By mapping these foundational causal pathways, CaStLe provides insights not just into immediate local effects but also into how these effects propagate and combine to create larger-scale patterns.

With these advances, CaStLe achieves remarkable improvements over state-of-the-art space-time causal discovery approaches. CaStLe is a flexible framework that can be implemented by adapting any given time series causal discovery algorithm to the stencil approach. Our approach performs excellently in high-dimensional data regimes, making it capable of describing the local space-time evolution of transient phenomena transporting over many grid cells.

The Earth system is rich with transient phenomena examples including forest fires, monsoons, coastal erosion, salt or freshwater incursions, inter-tropical convergence zone shifts, and atmospheric rivers. Aside from elucidating underlying dynamics, CaStLe can be used to identify and characterize causal change points, such as polar vortex disruption and ocean current disruptions. Additionally, understanding these local dynamic structures can give further insights into

the construction and evolution of important macro phenomena such as the El Niño Southern Oscillation, the Quasi-Biennial Oscillation, and the Madden-Julian Oscillation. Table 1 in the Appendix summarizes the capabilities of CaStLe and their relevance to specific Earth science applications. These capabilities address analytical needs that have been challenging or infeasible with previous causal discovery approaches.

The remainder of this paper is organized as follows: Section 7.4 provides a brief background on causal discovery and its use in Earth science; Section 7.5 describes our case studies in the HSW-V and E3SMv2-SPA models and available data; Section 7.6 explains our novel CaStLe methodology; Section 7.7 demonstrates CaStLe's ability to recover known volcanic aerosol evolution in climate models of different resolution; and finally, Section 7.8 illustrates CaStLe's computational, and performance improvements over the state-of-the-art methods with synthetic data experiments.

Contributions

We introduce the CaStLe approach to causal discovery from space-time data. CaStLe allows the discovery of causal structures in high-dimensional spatial data, avoiding the need for dimension reduction techniques that dominate causal discovery of space-time data, e.g., the work by Nowack et al. (2020a). By working in the raw data space, CaStLe's causal graphs are *inherently interpretable* and do not require mapping structures from the dimension-reduced space back onto the

original data. We provide a theoretical analysis of CaStLe, showing that it has attractive computational and statistical properties and, rather remarkably, that CaStLe's accuracy actually increases on larger spatial domains. We apply CaStLe to two simulations of a major volcanic eruption and demonstrate how it can be used to better understand how stratospheric winds mediate the climate response to volcanic activity. Our first study is of a relatively simplified model to validate the methodology with proxy ground-truth. In our second study, we consider a more realistic model and find that CaStLe still provides consistent and valuable results, demonstrating its value for realistic atmospheric dynamics. Finally, extensive numerical experiments measure the advantages of CaStLe and demonstrate: i) significantly improved performance over existing causal discovery methods on a set of vector autoregressive (VAR) benchmarks; and ii) the use of CaStLe to identify the governing dynamics of Burgers' non-linear partial differential equation (PDE). While our case studies utilize climate model data, the methodology is domain-agnostic and can be applied to any high-dimensional space-time system meeting our locality and stationarity assumptions.

7.4 Background: Causal Discovery and Formal Mathematical Scope

Here, we provide a brief overview of the causal discovery field and the mathematical scope of our contributions. For a broader overview of causal discovery and its applications to Earth science, see the reviews by Glymour et al. (2019),

Runge et al. (2019b), and Runge et al. (2023), and the book by Peters et al. (2017). Additionally, we outline the mathematical constraints and assumptions that define where our methodology can be applied in the class of space-time systems.

Causal discovery is a field of causal inference that seeks to recover causal dynamics from observational data. In the parlance of causal inference, *observational data* is data that is passively observed rather than data to which treatments (e.g. manipulations) have been applied. Observational data can be natural (e.g. physical observations) or synthetic (e.g. simulations). The present work exclusively pertains to untreated data, so we will use *observational* in this way.

While correlation does not imply causation, causal discovery is built upon Reichenbach's common cause principle (Reichenbach, 1956): if two quantities are correlated then one must cause the other or there is a third causal driver of the two. causal discovery generally has two output classes: a causal graph/network (Pearl, 1995a) or a structural causal model (Pearl, 1998). We focus on causal graphs, which are networks of variables (nodes) connected by edges that denote a causal dependence. Causal graphs can be more appealing than structural equation models because they are human-interpretable and do not require prior knowledge of the underlying causal function. In the study of Earth science, causal graphs may often be preferred to visually describe space-time relationships on the globe. Our contribution produces a novel type of causal graph, the causal space-time stencil, which is detailed in Section 7.6 and an example of which is in panel 4 of Figure 7.2.

7.4.1 Related Work: Causal Structure Learning

In recent decades, causal inference has been developed into a rigorous mathematical framework (Rubin, 1974; Pearl, 2000; Pearl et al., 2016). These developments made algorithmic discovery of causal structures from observational data possible (Spirtes et al., 1993; Peters et al., 2017; Glymour et al., 2019). Causal structures can be modeled with two common forms: structural causal models (SCMs) and causal graphs. Both describe a functional relationship between a variable X_j and its causal parents, denoted $\mathcal{P}(j)$.

For example, if X_i causes X_j , then it is said X_i is a parent of X_j and $i \in \mathcal{P}(j)$. Formally, Peters et al. (2017, p.83) defines an SCM as follows:

A structural causal model (SCM) consists of a collection of d (structural) assignments

$$X_j := f_j(\boldsymbol{X}_{\mathscr{P}(j)}, \boldsymbol{\eta}_j), \qquad j = 1, \dots, d,$$

where $\mathcal{P}(j) \subseteq \{1,...,d\} \setminus \{j\}$ are called **parents of** X_j : and a joint distribution $\mathbf{P}_{\eta} = P_{\eta_1,...,\eta_d}$ over the noise variables, which we require to be jointly independent; that is \mathbf{P}_{η} is a product distribution [in our notation].

An SCM admits a unique causal graph, where $X_j \to X_i$ if $j \in \mathcal{P}(i)$ and $j \not\to X_i$ if $j \notin \mathcal{P}(i)$. While discovery of an SCM requires hypothesizing all f_j 's, discovering a causal graph can be done without knowing the exact functions. Because a causal graph does not imply a specific function between variables, each may imply

multiple SCMs. This does limit some of the inferential power of causal graphs, in exchange for more versatility.

Algorithms for discovering causal graphs have two primary classes: constraint-based and score-based algorithms. Constraint-based methods use statistical tests to compute conditional independence relationships between sets of variables. Once a set of independence relationships is established, it utilizes causal assumptions and reasoning to connect the variables with directed links. Score-based approaches are similar but use score optimization to determine causal dependence between variables. Both constraint-based and score-based algorithms produce causal graphs because they operate on graphical structures and independence relations rather than the explicit parametric relationships between variables required to specify a complete SCM.

Early causal discovery algorithms developed as two parallel traditions. The temporal Granger causality (Granger, 1969) methodology was an early innovation using time series data to determine if the past history of *X* aids the prediction of *Y* better than *Y*'s history alone. If so, then *X Granger causes Y*. Independently, the constraint-based PC algorithm (named for its authors Peter and Clark) (Glymour and Scheines, 1986) and FCI (Spirtes and Glymour, 1991) developed out of the inductive causation (Pearl and Verma, 1992) framework and the earlier SGS algorithm (Spirtes and Glymour, 1991), significantly improving the efficiency of causal discovery using statistical structures in observed data. In time, other structural algorithms developed, such as LiNGAM (Shimizu et al., 2006), utilizing asymme-

tries in non-linear and non-Gaussian data for inferences, and NOTEARS (Zheng et al., 2018), a graph score-optimization-based method. Eventually, these two traditions converged as structural methods were developed to take advantage of temporally ordered data. Key advances included: hMRF (Liu et al., 2010), which uses hidden Markov models for estimation and is grounded in Granger causal structures, PCMCI (Runge et al., 2019a) (and related PCMCI+ and LPCMCI), which improves PC to handle autocorrelated dependencies better, and DYNOTEARS (Pamfil et al., 2020), which extends the NOTEARS method to time series. More recently, a third tradition, causal representation learning, developed out of machine learning (ML) to leverage causal reasoning in ML models (Schölkopf et al., 2021). While still a developing field, it shows particular promise for estimating relationships in the presence of latent confounding.

The directed nature of time provides a powerful asymmetry to leverage, often sufficient to overcome the difficulties of autocorrelation, automatically orienting discovered relationships in time. In contrast, spatial data lacks an obvious uniform directional structure and poses challenges for causal discovery. As discussed in Section 7.3, while some approaches have incorporated domain-specific spatial constraints for point-measurement networks, none have developed a generalizable framework that leverages fundamental physical principles of locality to enable scalable causal discovery in high-dimensional gridded space-time systems.

Causal Discovery in Earth Science

We present a brief review of causal discovery for Earth science to position CaStLe within the literature. Please also see the extensive reviews by Runge et al. (2023) and Ali et al. (2024).

Ebert-Uphoff and Deng (2012) were the first to apply a causal discovery algorithm, PC-stable (Colombo and Maathuis, 2014), to the climate science domain. They were able to find a grid-cell-level causal teleconnection network in 50 year daily geopotential height data using the PC algorithm. Ebert-Uphoff and Deng (2014); Deng and Ebert-Uphoff (2014) further explored application requirements and climatological interpretations of the geopotential height analysis. In each paper, they note grid challenges related to the high expense of many grid cells, aggregation effects, and cell spacing. The first paper limits the number of grid cells to 800, while the subsequent analyses limited grid cells to 200 to minimize computational costs. While their results are compelling, they use extensive decadal data and recover patterns common to all 50 years. The fundamental difference between our work and Ebert-Uphoff and Deng's work is that they recover causal graphs from recurring atmospheric phenomena with sufficiently large datasets on relatively coarse-grained grids, whereas CaStLe is recovers networks of isolated phenomena with many more grid cells and many fewer time samples per cell.

Runge et al. (2015c) introduced an alternative approach to causal discovery of space-time Earth science data. They reduced the dimensionality with varimax-rotated principal component analysis prior to applying the causal discovery al-

gorithm, producing a graph relating discrete, potentially remote, regions. Their causal graph is most similar to a teleconnection network between large areas on the globe. Nowack et al. (2020a) utilized that framework to evaluate CMIP5 models. Particularly of note, they point out the challenges and strengths of Ebert-Uphoff and Deng (2012)'s grid-cell-level approach, "... while an analysis at the grid-cell-level is more granular which, however, carries the challenges of higher dimensionality, will have a strong redundancy among neighbouring grid cells, and grid-level metrics will require handling varying spatial resolution among data sets."

Tibau et al. (2022) built on the dimensionality reduction approach, augmenting it to output grid-cell-level networks. They specifically delineate *mode-level* (dimensionality reduction or cell aggregation) and grid-level causal discovery. Their augmentation is called Mapped-PCMCI, which first applies dimensionality reduction, then computes a mode-level causal network with PCMCI, and finally maps the grid cells within the modes to each other using the network previously constructed. Their resulting network is one consisting of edges between grid cells, but the method assumes that cells within modes are fully connected, i.e., each cell is dependent on all of its neighbors. In contrast, our work specifically seeks inter-cell spatial relationships. Finally, they also describe the failure of a traditional causal discovery approach for grid-cell-level data, "[if] we apply PCMCI directly at the grid-level, the low power of this high-dimensional and redundant estimation problem (see Section 2.2.2) leads to most links being missing."

Boussard et al. (2023) and Brouillard et al. (2024) developed the Causal Discovery with Single-parent Decoding (CDSD) algorithm within the causal representation learning framework and applied it to the climate science field. Like CaStLe, CDSD performs well in high-dimensional data settings but through a different mechanism. It performs dimensionality reduction by learning latent variables and enforcing a "single-parent" constraint where each grid cell belongs to exactly one latent factor. This naturally clusters grid cells into coherent, often contiguous regions and enables the discovery of causal relationships between these larger-scale patterns. In contrast to CaStLe's grid-level structure learning, CDSD identifies broader teleconnection pathways between regional climate modes. Thus, while CaStLe preserves the original grid structure to capture fine-grained causal dynamics, CDSD abstracts to a higher level by mapping the native grid space to an identifiable latent representation before performing causal discovery.

Several studies have addressed local-scale phenomena. Pfleiderer et al. (2020) applied causal discovery to identify precursors to seasonal hurricane frequency. They utilized the precursors to inform a predictive model. Polkova et al. (2021) identified local drivers of marine cold-air outbreaks in the Barents Sea. These demonstrate that existing causal discovery approaches can be valuable for seasonal and sub-seasonal phenomena. However, both marginalized large regions prior to analysis, reducing the space's dimensionality, and did not evaluate the space-time evolution of phenomena nor grid-level dynamics.

There are some examples of causal discovery algorithms leveraging spatial in-

formation. Zhu et al. (2016) developed pg-Causality that applies space-time pattern mining and a Gaussian Bayesian Network to seek local dependencies in the space-time propagation of air quality data. Sheth et al. (2022) developed STCD for understanding hydrological systems. They constrained the discovery of spatial structures by only allowing higher elevation nodes to be parents of lower elevation nodes because water follows the gravity gradient. While both cleverly use mined or known spatial structure to inform their causal discovery, they are both limited to use in sparse point-measured data from static base stations rather than gridded data. Further, these methods enforce constraints as filtering mechanisms, whereas CaStLe actively leverages spatial structure to enhance statistical power. Neither address the scalability challenges in high-dimensional gridded data.

Parallel Approaches in Neuroscience: Causal Discovery for High-Dimensional Spatial-Temporal Data

Other scientific domains face similar challenges with high-dimensional space-time data. Neuroscience, for example, needs to study mechanisms in brain interactions, and fMRI images may contain thousands to millions of pixels. The anatomy of the brain also exhibits locality constraints. Ramsey (2014) made computational optimizations to the Greedy Equivalence Search algorithm, including sparsity constraints and limiting the distance of potential parents, to recover graphs with millions of nodes. Saetia et al. (2021) marginalized regions of interest in the brain using spatial averaging and then applied the PCMCI algorithm to construct causal graphs. There is a common interest in recovering graphs of high-dimensional grid-

level data throughout the sciences. Developing more tools that enhance the estimation and interpretability of causal graphs in these spaces will help advance our understanding of space-time structures across the sciences.

What is clear from prior work is that grid-level analyses are challenging, both statistically and computationally, due to how many grid cell dependencies need to be estimated, the enormous number of observations needed, and the redundant information content of nearby cells. As we present in the following sections, CaStLe adds to the literature as it overcomes the statistical and computational limitations of grid-level analysis by leveraging the known physical structure of spatial information to produce interpretable graphs describing local causal structures.

7.4.2 PDE-Like Systems

We seek to perform causal discovery from space-time data governed by consistent physical laws. As detailed in Section 7.6, CaStLe operates via two phases. The first restructures the given space-time data into a lower-dimensional local neighborhood space without marginalization or loss of any data points; the second is the causal discovery step. This section details the assumptions required for efficient use of spatial replicates that enable CaStLe's first phase, scalability properties, performance in high-dimensional settings, and interpretability. We note that the assumptions necessary for the second phase will be inherited from our metaalgorithm's chosen causal discovery method. In general, they will be the causal Markov condition, faithfulness, and often causal sufficiency, which we define for-

mally in Appendix A.2.

We take PDE-like models as our starting point, and assume that all behavior in the given space are driven by a fixed set of dynamics that apply at infinitesimal time and spatial scales. Specifically, we assume that, for data observed in discrete space and time, the evolution of a single grid cell is controlled only by the values of its immediate spatial neighbors at the previous time step. Using causal discovery, we seek to determine which neighbors have a causal impact on a given grid cell and the direction of that relationship. Our analytical framework has similarities to the sparse identification framework initially developed by Brunton et al. (2016), though our approach builds upon causal discovery rather than sparse regression. Because our approach can use non-linear conditional independence tests, we can avoid the difficult dictionary construction step associated with sparse regression methods.

In contrast to causal discovery methods, other current research also focuses on approximating ordinary differential equations or PDE-like systems with operator learning approaches, such as operator neural networks (Li et al., 2020; Pathak et al., 2022; Hart et al., 2023). These Fourier Neural Operators (FNO) focus on generating accurate models of the PDE-like evolution of key variables over time and space. Their assumptions are rooted in several of the same fundamental physical principles of how PDEs propagate effects in space and time as CaStLe: locality in space and time and spatial stationarity. While CaStLe is not meant to be a predictive model, it captures important relationships between grid cells in an inter-

pretable fashion, providing insights into the underlying causal structures.

7.4.3 Causal Discovery of Physical Dynamics: Dynamical Constraints

We state here four key assumptions that capture what we describe as a PDE-like system X_t :

- **T1**) Temporal Locality: for any $\tau \neq 1$, $X_{i,t-\tau} \not\to X_{j,t}$ for any spatial coordinates (i,j)
- **T2**) Temporal Causal Stationarity: the dynamics governing the evolution of X_t do not change over time. That is, $X_{i,t-1} \to X_{j,t} \Leftrightarrow X_{i,t-1+\tau} \to X_{j,t+\tau}$ for any time offset τ .
- **S1**) Spatial Locality: if (i, j) are not neighbors (in a problem-specific sense) then $X_{i,t_1} \not\to X_{j,t_2}$ for any t_1, t_2 .
- **S2**) Spatial Causal Stationarity: the dynamics governing the evolution of X_t do not change over space. That is, $X_{i,t-1} \to X_{j,t} \Leftrightarrow X_{i+s,t-1} \to X_{j+s,t}$ for any spatial offset s.

Here, \rightarrow denotes the absence of a direct causal relationship between two variables.

Therefore, if an SCM exists for a given system, then it will have a functional shape constrained by our assumptions: $X_t = f(X_{t-1}, \eta_t)$, for some vector of noise, η_t . In the context of an SCM, the constraints are: temporal locality (T1) adds lagged relationships between parent and child variables; spatial locality (S1) restricts possible parents to those in the spatial neighborhood of each variable (grid

cell), that is, f_i is only a function of the neighborhood of i (f_i depends only on $X_{\mathscr{P}(i)}$); and temporal/spatial causal stationarity (T2 & S2) require that there be only one function, f, for all space and time in the window/region of analysis.

Building on physical principles, Assumption T1 implies that causal dependencies follow the "arrow of time" while S1 disallows "action at a distance." Assumptions T2 and S2 serve to ensure that there is a consistent causal structure to target. Assumption S1 further requires that f_i is only a function of the neighborhood of i (f_i depends only on $X_{\mathcal{P}(i)}$). We refer the reader to the book by Peters et al. (2017) for a more detailed discussion of how SCMs can be used to model physical systems.

We deliberately chose lag-1 temporal relationships in assumption T1 because they reflect fundamental physical principles: In the discretized form of PDEs, each element depends on the future state of the immediate past of its neighboring elements. The symmetry of the radius-1 neighborhood in assumption S1 and the single lag constraint in assumption T1 captures the essential causal dynamics in physical processes when temporal and spatial data resolutions are appropriately balanced.

While not descriptive of all possible systems, we assert these locality and stationarity assumptions are descriptive of any system governed or modeled after PDEs, cellular automata (Bhattacharjee et al., 2020), or Tobler's First Law of Geography (Miller, 2004; Walker, 2022). These assumptions reflect fundamental principles of locality and consistency that apply across numerous domains, from

fluid dynamics to reaction-diffusion systems. However, for these to hold in practicality, one must also assume sufficient data is available to characterize locality and dynamics are smooth and non-turbulent, relative to the analysis frame. These assumptions imply that there is an optimal balance between temporal and spatial resolution sufficient to impose space-time locality. The exact value of this scaling is problem-dependent, as more rapidly evolving systems require higher temporal resolution, and we do not explore it further here. However, we note that similar concerns are well-studied in the design of numerical differential equation solvers where spatial and temporal discretizations must be chosen suitably consistently.

Section 7.6 and A detail how these assumptions are essential for our methodology, CaStLe, and discuss their limitations. Section 7.6.6 discusses strategies for managing those limitations. While CaStLe's framework assumptions (T1, S1, T2, S2) enable efficient use of space-time samples, the algorithm adapted for CaStLe's parent-identification phase will have additional causal assumptions.

Interestingly, CaStLe's spatial locality assumption (S1) creates an environment where, when properly implemented, causal sufficiency can be satisfied by construction. When we focus on learning only the parents of the center cell while including all potential spatial neighbors in the analysis, we automatically satisfy causal sufficiency for that specific node if S1 holds. While reliant on S1 holding, this is significant because causal discovery is notoriously the most challenging causal discovery assumption to ensure in real-world settings (Spirtes et al., 1993; Raghu et al., 2018). As we discuss in Section 7.6.5, sufficiency may be

relaxed depending on which causal discovery algorithm is adapted for the parentidentification phase. However, satisfying it by construction may enable implementation choices with fewer compromises.

In the following sections, we discover grid-cell-level causal graphs under these five assumptions. Assumptions T1 and S1 allow us to significantly reduce the scope of the problem, as there are only 9 possible parents of a grid cell in 2D (8 neighbors and itself). Assumptions T2 and S2 suggest that we only need to determine a single local causal graph, because spatial stationarity allows us to extend it to the entire domain.

7.5 Data: The 1991 Mt. Pinatubo Eruption

Mount Pinatubo's eruption in 1991 was a massive, natural intervention in the climate, with effects that had a relatively high signal-to-noise ratio. The event launched 20 Tg of SO_2 gas into the atmosphere (Guo et al., 2004a,b; Kremser et al., 2016). The sulfate aerosols that resulted from these gases remained in the stratosphere for approximately two years, leading to stratospheric warming of ~ 1.5 K and surface cooling of 0.2-0.5K (Dutton and Christy, 1992; Labitzke and McCormick, 1992; Parker et al., 1996a; Soden et al., 2002). This aerosol injection has recently been the object of much study, with some authors suggesting it as a natural proxy for proposed stratospheric aerosol injection (SAI) responses to global climate change (Trenberth and Dai, 2007). Recent work continues to characterize the nature of the response to the Pinatubo eruption, with the timing and

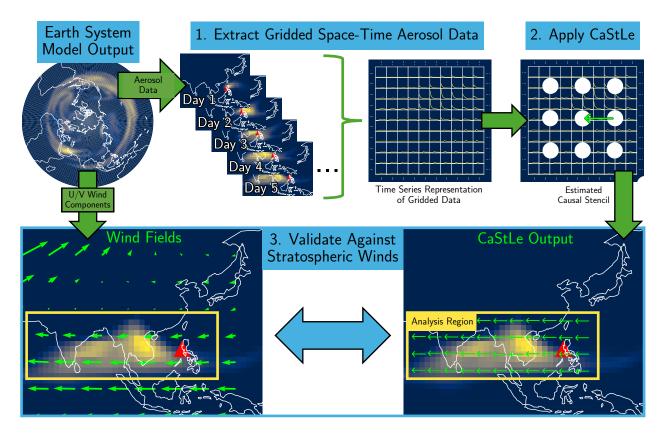


Figure 7.1: Schematic overview of the key elements of CaStLe and the process followed in its application to Mount Pinatubo's eruption of stratospheric aerosols. Beginning with Earth system model output, Step 1. is to collect stratospheric wind and aerosol data. Step 2. is to apply our novel CaStLe meta-algorithm to the aerosol data to obtain a causal graph describing the space-time evolution of the aerosols. Finally, we use the wind fields to help validate the causal graph results in Step 3.

spatial structure of the surface response being essential factors to inform policy decisions (Weylandt and Swiler, 2024).

Large volcanoes can impact climate quantities, such as surface temperatures, on timescales from months to years (Parker et al., 1996b; Robock, 2000; Timmreck, 2012; Marshall et al., 2022). However, to evaluate whether CaStLe could recover the initial advection dynamics of volcanic aerosols, we focused on the period shortly after the eruption that includes stratospheric aerosol transport. The recent paper by Marshall et al. (2022) indicates: "Although global-scale climatic impacts

following the formation of stratospheric sulfate aerosol are well understood, many aspects of the evolution of the early volcanic aerosol cloud and regional impacts are uncertain." This initial spread of aerosols in the stratosphere is a geophysical process, falling between synoptic weather patterns and longer-term impacts.

We utilized models of the event, combining stratospheric aerosol and wind data, as case study to illustrate the analysis possible with CaStLe. Figure 7.1 is a high-level illustrative schematic of the this work's key ideas: We collect gridded space-time data, e.g. aerosol optical depth (AOD) measurements, and apply it to CaStLe to learn a causal stencil graph. We then map the stencil to the original grid space. Finally, we compare the data to ground-truth. o be clear, the ground-truth in our later case studies is a proxy, referring to the models' understood underlying dynamics, not the true realization of AOD in Earth's atmosphere or a mathematical representation of the dynamics. In Section 7.7, we compare to the wind fields carrying AOD as a proxy ground-truth. In Section 7.8, we compare CaStLe results from synthetic data to mathematically-known ground-truth.

7.5.1 Held-Suarez-Williamson-Volcanic

For our first case study, we utilized the limited-variability ensemble approach of the Held-Suarez-Williamson-Volcanic (HSW-V) model (Hollowed et al., 2024). HSW-V is an atmosphere-only model built in the Department of Energy's Energy Exascale Earth System Model version 2 (E3SMv2) (Golaz et al., 2022). HSW-V does not set out to replicate the historical Mt. Pinatubo eruption or any other, but

uses the Mt. Pinatubo's eruption characteristics "to produce a plausible realization of a volcanic event, simulated with a minimal forcing set" (Hollowed et al., 2024). The model was developed specifically to facilitate basic research of attribution methodologies by providing realistic source-to-impact pathways of eruption quantities. We use this model to create a realistically complex dataset of stratospheric aerosol and wind dynamics with a clear ground-truth to demonstrate the capabilities of CaStLe and the correctness of its results.

We gathered aerosol optical depth (AOD), sulfate, and zonal (U) and meridional (V) wind fields for analysis. Only AOD is provided to CaStLe, while the sulfate, U, and V wind components are used for validating results, as detailed in Section 7.7. AOD is a derived quantity that measures the extinction of a beam of light through the atmosphere by atmospheric aerosols, i.e., it describes the amount of light occluded by atmospheric particles. One of the simplifying aspects of HSW-V is that all aerosol particles originate from SO₂ gas ejected by the volcano; this avoids confusing signals from other sources, such as smoke and dust, in the atmosphere.

The data collected from the HSW-V ensemble run are on a 2° grid with 6-hourly average observations. We selected AOD in grid cells between -20.00° to 40.00° N and -120.00° to 140.00° E, comprising 3,900 grid cells. We used the first three weeks post-eruption for our analysis.

7.5.2 Mt. Pinatubo in E3SMv2-SPA

For our second case study, we considered a simulation of the Mt. Pinatubo eruption in the fully coupled E3SMv2 model augmented with Stratospheric Prognostic Aerosol capability (E3SMv2-SPA) as detailed and validated by Brown et al. (2024). E3SMv2-SPA includes atmosphere, land, ocean, sea ice, land ice, and river components. AOD, U, and V wind fields are analogously collected from this dataset. However, in this model, aerosols are a natural feature, thus complicating the analysis of aerosol optical depth.

Data were collected on a daily temporal resolution for a 1° spatial grid. We selected grid cells between -30.00° to 60.00° N and -180.00° to 180.00° E. Analysis covered the first six months. Because this data has a coarser temporal resolution and finer spatial resolution than our study of HSW-V, we coarsened the CaStLe spatial grid to a 3° grid, resulting in 3,600 total grid cells. This helps ensure that the motion of aerosol particles between grid cells is measured within the one-day sample period.

7.6 Methodology: Causal Discovery with CaStLe

7.6.1 Notation

We first introduce notation used in the remainder of this paper. Data is observed on a spatial domain \mathcal{D} , which we typically take to be a finite subset of the real plane, \mathbb{R}^2 . The causal structure generating this data can be represented by a directed

acyclic graph $\mathscr{G}=(\mathscr{V},\mathscr{E})$, where $\mathscr{V}=\mathscr{D}$. CaStLe represents local causal structure with a *stencil*, which we identify as a graph $\tilde{\mathscr{G}}=(\tilde{\mathscr{V}},\tilde{\mathscr{E}})$ in a reduced coordinate space $(|\tilde{\mathscr{V}}=9|)$. In both the original and reduced spaces, let $\mathscr{P}(v)$ be the *potential* causal parents of v and let $\mathscr{P}(v)$ be the *actual* causal parents of v. We take \mathscr{D} to be points on a regular grid of size $N\times N$, observed over T time steps, giving data $\boldsymbol{X}\in\mathbb{R}^{N^2\times T}$. When transformed to the reduced space used by CaStLe, the resulting data matrix will be denoted $\tilde{\boldsymbol{X}}\in\mathbb{R}^{T(N-2)^2\times 9}$. Quantities estimated from data are denoted with a hat, e.g., $\hat{\mathscr{P}}(v)$. We provide additional background on the interpretation of the causal graphs $\mathscr{G},\tilde{\mathscr{G}}$ in Section 7.4.1 and formally specify the mapping between \boldsymbol{X} and $\tilde{\boldsymbol{X}}$, or equivalently, between \mathscr{V} and \mathscr{V} , in Section 7.6.3.

7.6.2 Causal Space-Time Stencil Learning

We now introduce the CaStLe paradigm for the causal discovery of local space-time dynamics. Under our assumptions, CaStLe identifies a *sketch* of the local causal dynamics, which we call a stencil. This stencil can then be used to construct the causal graph for the entire system (S2). The stencil is estimated in a reduced coordinate space, where we only examine the direct neighbors of a given grid cell (S1). We can pool information across time (T2) and space (S2) in order to estimate the stencil accurately, and the problem is tractable because we only seek causal parents which are local in time (T1). As we will see, this combination of reduced search space and pooled information provides a powerful approach to causal discovery and enables accurate causal discovery from high-dimensional

grid-cell-level data.

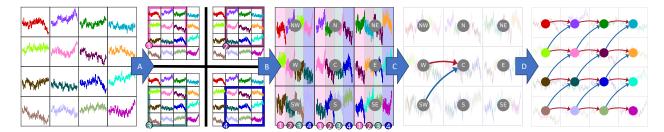


Figure 7.2: Illustration of CaStLe (Algorithm 1) as applied to space-time data on a 4×4 grid. Step A (§7.6.3): for every interior grid cell, its 3×3 (Moore) neighborhood is selected. (Note, all four 4×4 grids in the second panel are identical.) Step B (§7.6.3): Data are represented in a reduced coordinate space obtained by appending time series from each neighborhood according to its position relative to the neighborhood's center. Step C (§7.6.3): during the Parent Identification Phase (PIP), a causal discovery algorithm is used to estimate the parents of the center time series; the resulting graph forms the causal stencil. Step D (§7.6.3): the estimated stencil is expanded to its equivalent representation in the original space. Note that each time chunk (colored intervals in the center panel) in the reduced space corresponds to an interior grid cell of the original data, and that each edge in the final causal graph reflects to a stencil edge learned during the PIP. See §7.6.3 for details.

Having motivated the CaStLe approach to causal discovery from space-time data in Section 7.4.2, we now state it formally as Algorithm 1, describe its computational steps, and then analyze its statistical and computational properties.

7.6.3 The CaStLe Meta-Algorithm

Steps A-B: Projection to a Reduced Coordinate Space

CaStLe begins by transforming the given data from its original domain into a reduced coordinate space that captures the underlying causal dynamics' locality and spatial homogeneity. In this transformation, all data points are preserved, i.e., no marginalization or truncation occurs. This process is represented as Steps A and B in Figure 7.2 and Algorithm 1. In Step A, the local 3×3 (Moore) neighborhood of

Algorithm 1 CaStLe for Space-Time Data in 2D ($\mathscr{D} \subseteq \mathbb{R}^2$)

Inputs:

- Parent-Identification Phase subroutine PIP
- Gridded space-time data $\mathbf{X} \in \mathbb{R}^{T \times N^2}$
- 1. Step A: Extract 3 × 3 Moore Neighborhoods
 - For each interior point in the original space, construct local view of the data $\mathbf{X}_i = [X_{\cdot \mathscr{P}(i)}] \in \mathbb{R}^{T \times 9}$
- 2. Step B: Construct Reduced Space Data Matrix

$$\tilde{\boldsymbol{X}} = [\boldsymbol{X}_1^\top \boldsymbol{X}_2^\top \dots \boldsymbol{X}_{(N-2)^2}^\top]^\top \in \mathbb{R}^{T(N-2)^2 \times 9}$$

3. Step C: Perform Parent-Identification in Reduced Space

$$\mathsf{PIP}(\tilde{\pmb{X}}) = \tilde{\mathscr{E}} = (\hat{\mathscr{P}}(\mathtt{C}) \times \mathbb{R}^9) \subseteq \mathscr{P}(\mathtt{C}) \times \mathbb{R}^9$$

- 4. Step D: Expand Stencil Graph to Original Coordinate Space:
 - $\mathscr{E} = \emptyset \subseteq \mathscr{V}^2 \times \mathbb{R}$
 - For each $(p, w) \in \hat{\mathscr{E}}$:

$$\mathscr{E} = \mathscr{E} \cup \{ (p(v), v, w) \text{ for } v \in \mathscr{V} \}$$

Outputs:

- Graph Stencil, $\tilde{\mathscr{E}}$
- • Estimated Causal Graph, $\mathscr{G} = (\mathscr{V},\mathscr{E})$

each interior cell is selected, and each cell is labeled by its location relative to the center cell (S, NW, E, etc.). This process creates $(N-2)^2$ sub-views in $\mathbf{X}_i \in \mathbb{R}^{T \times 9}$.

In Step B, these views are concatenated along the time dimension to create a reduced coordinate space data matrix $\tilde{X} \in \mathbb{R}^{T(N-2)^2 \times 9}$. Note, when concatenating the subviews, data are aligned by their coordinates relative to the neighborhood center so that, e.g., data from all NW cells are aligned upon concatenation, even though they originally come from different spatial locations. Although this transformation results in specific time series segments appearing in multiple reduced space cells, these repetitions do not eventually create spurious dependencies in the causal stencil, as CaStLe only seeks lag-1 dependencies. The repeated segments are well-separated in the temporal dimension, and no chunks appear in different cells in the same interval.

We depict this process on a 4×4 grid in the first half of Figure 7.2. In Step A, the four interior cells are sequentially highlighted, and their local neighborhoods are extracted, which are depicted in boxes colored according to the center used. In Step B, the local data views are concatenated to form one set of time series, with each temporal *chunk* reflecting the color of the center cell of the underlying data view.

Step C: Parent-Identification Phase

CaStLe next examines the reduced coordinate space data representation, \tilde{X} , to identify the stencil of the local causal dynamics. This is done by applying an

augmentation of an arbitrary time series causal discovery algorithm to identify the parents of the center cell, C. We emphasize that we only seek the parents of C, not a full causal structure, in this step and refer to it as the *Parent Identification Phase* (PIP). Under assumption S1 (locality), all parents of C are present at this step, satisfying causal sufficiency, ensuring more accurate estimation of the causal stencil. By contrast, while the data of the parents for the exterior cells, e.g. W, is included in the reduced data space matrix, \tilde{X} , it spreads across multiple columns, and accurate parent identification is not possible. The output of this process is a set of (up to) 9 weighted edges, corresponding to the parents of C (the eight neighboring cells and C itself).

We depict the PIP in Step C of Figure 7.2, where two parents of C are identified: W, which has a positive dependence on C, and SW, exhibiting negative dependence. Note that while the PIPs we implemented in testing—see Section 7.8.1—had no trouble with the *seams* connecting each time *chunk* in the reduced space, we propose an improved testing implementation in E to alleviate potential statistical testing issues.

Step D: Graph Reconstruction in the Original Space

Finally, CaStLe uses the stencil constructed in Step C to reconstruct the causal graph in the original data space, in a process that essentially reverses Steps A and B. Specifically, for each edge identified in $\tilde{\mathcal{E}}$, corresponding edges are added to grid cell in the original domain. We depict this in the final step of Figure 7.2 where

the stencil is repeated throughout the entire 4×4 space, copying the two parents of C identified in Step C, to create a causal graph in the original space. Note also that we use the stencil to identify parents for both interior and boundary cells, omitting edges that go "off-grid" when applying the stencil to boundary cells.

7.6.4 Theoretical Properties

CaStLe has many advantages over classical causal discovery algorithms in gridded space-time settings. By reducing the causal discovery problem to identifying the causal parents of the center cell (C) in the reduced space, CaStLe achieves significant improvements in both the computation necessary to infer the causal graph and the statistical quality of that graph. As previewed in Section 7.4.2, the PIP's focus on identifying only the parents of the center cell creates an important connection to the causal discovery assumption of causal sufficiency. Because we include all spatial neighbors (as defined by our locality assumption S1) in the conditioning set, all potential parents of the center cell are present in the analysis. If our spatial locality assumption holds, causal sufficiency is automatically satisfied within each local stencil analysis. This represents a key advantage of the CaStLe framework - while the Markov condition and faithfulness remain necessary assumptions for the PIP algorithm, our implementation leverages spatial structure to ensure causal sufficiency by construction.

Below, we briefly outline the theoretical implications and their contributions to CaStLe's remarkable performance and algorithmic improvements. Their deriva-

tion, a deeper analysis, and a discussion on graph estimation asymptotic consistency are provided in B. We discuss CaStLe's asymptotic consistency in C, which shows that CaStLe converges on the correct causal stencil as grid size increases, given a PIP consistent in increasing time samples. These properties illustrate the mathematical justification for CaStLe's empirical correctness and improvement over the state of the art shown in the following sections.

CaStLe yields significant improvements to both time complexity, a measure of an algorithm's computation time as it scales with input size (e.g., number of time steps, graph nodes), and statistical complexity, a measure of estimation performance given larger sample sizes. Following the complexity analysis of Kalisch and Bühlmann (2007), we show that traditional causal discovery approaches are bounded by $\mathcal{O}(np^32^p) = \mathcal{O}(T(N^2)^32^{N^2}) = \mathcal{O}(TN^62^{N^2})$, for T time samples and $N \times N = N^2$ grid cells. Since CaStLe computes on the smaller reduced coordinate space, and only seeks causal parents of one node, rather than parents of all nodes, several terms become constants, resulting in $\mathcal{O}(np^32^p) = \mathcal{O}(T(N-2)^2 \times 10^{-5})$ $9^3 \times 2^9$) = $\mathcal{O}(TN^2)$. CaStLe's computational complexity is $\mathcal{O}(TN^2)$, a major improvement over existing approaches. For more details on this derivation, see Appendix B.1. By leveraging locality and spatial replicates, CaStLe identifies causal structure for the entire graph ($\mathcal{O}(N^4)$ possible edges) in N^2 time. Kalisch and Bühlmann (2007, Appendix B) show that the probability of the PC algorithm incorrectly estimating the true graph is bounded by $\approx \mathcal{O}(N^{2N^2})$, whereas we find that CaStLe's error probability scales as $\approx \mathcal{O}\left(\frac{N^2T}{e^{N^2T}}\right)$. From this, as the grid size grows

larger, we see PC is less likely to estimate the correct causal graph, while CaStLe is more likely to estimate the correct graph. Furthermore, both of these effects are exponential, implying significant performance differences even on moderately sized graphs; this change from a regime of exponential decay to super-exponential growth in graph recovery performance makes local causal graph recovery feasible, finally enabling the tools of causal discovery to scalably explore grid-level Earth science dynamics in commonly high-dimensional settings.

7.6.5 Methodological Limitations

CaStLe's assumptions may pose challenges in some domains of interest, and violations of these assumptions can affect the CaStLe output. For example, large-scale homogeneity can be difficult to achieve in geosciences, which is the primary rationale for the spatial-blocking strategy that we implement for our application in Section 7.7. Locality assumptions (T1 & S1) create a framework where the causal Markov condition can be effectively applied to local structures, while causal stationarity assumptions (T2 & S2) create consistency in these structures across space and time. However, the PIP algorithm we use within CaStLe additionally requires standard causal discovery assumptions, particularly the causal Markov condition and faithfulness, which is a separate non-trivial assumption. We list causal sufficiency as an assumption, however, if the others hold then it follows that all of the causal parents of the stencil's center are in its immediate neighborhood, so sufficiency is satisfied by construction. Alternatively, causal sufficiency may be relaxed

if the chosen PIP is an algorithm that does not rely on sufficiency, such as the FCI algorithm (Glymour et al., 2019). As such, violations of CaStLe's assumptions relate directly to violations of the causal Markov condition, faithfulness, and causal sufficiency. Both Spirtes et al. (1993, p. 29) and Runge (2018a) discuss assumption violations in causal discovery and some examples of how they manifest in resulting graphs. We have included a more detailed discussion on each assumption and their limitations in A.

7.6.6 Strategies for Addressing Limitations

To address the limitations of CaStLe's assumptions, several practical strategies can be employed. One effective approach is the use of spatial blocking to create subdivisions where dynamics are more uniform, thus mitigating the violation of spatial causal stationarity (S2). The selection and size of these blocks are highly domain-dependent and can be guided by subject matter expertise. An automated approach may be sufficient for certain dynamics, such as stratospheric dynamics, but more manual approaches may be necessary for surface-level dynamics where blocks are chosen based on topological assumptions. In specific areas of interest, blocks can be manually created to avoid topological boundaries such as coastlines, rivers, and mountain ranges, ensuring that the assumptions of spatial homogeneity are better satisfied.

Additionally, strategies such as variograms can be used to test for spatial statistical stationarity, providing heuristics for effective blocking. In future work, an

iterative block size estimation approach could be considered. Varying the block size serves as a form of *stability check*, a technique widely applied in ML to ensure robustness of discoveries to parameter choices and modeling assumptions (Allen et al., 2023). However, it is important to note that there may not always be a single optimal block size due to the complex nature of spatial dynamics. Instead, there may be a range of valuable block sizes depending on the needs for analysis and the limitations of the setting. Because CaStLe is data efficient, it may be better to tend towards smaller blocks, which are more likely to be homogeneous, but possibly at the cost of some interpretability.

Deep learning and space-time feature engineering approaches may be fruitful directions for future research on automated block-identification. Methods such as δ -MAPS (Fountalis et al., 2018), feature extraction with convolutional neural networks (Nukavarapu et al., 2023), and spatiotemporal cluster analysis (Davis et al., 2025) are strong starting points. These computational approaches could automate the identification of optimal spatial blocks, reducing reliance on manual delineation and subject matter expertise while preserving the statistical properties necessary for valid causal discovery with CaStLe.

By employing these strategies and acknowledging their limitations, the robustness and applicability of CaStLe in various domains can be significantly enhanced, allowing for more accurate causal discovery in complex space-time systems. In general, more data at higher spatial and temporal resolutions will make satisfying the assumptions easier. The appeal of CaStLe is when one is interested in smallscale local dynamics, it is preferable to analyze raw gridded data directly, because marginalization can introduce statistical artifacts.

I provides an empirical investigation of how violations of each assumption affect CaStLe's performance when applied to our E3SMv2-SPA case study. Our analysis reveals that CaStLe is surprisingly robust to moderate assumption violations. While violations of spatial and temporal causal stationarity (particularly with overly large blocks or extended time intervals) introduce more false positives and reduce interpretability, CaStLe often still identifies key true causal pathways. This robustness to moderate assumption violations further expands the practical utility of CaStLe in realistic Earth science applications where perfect adherence to assumptions is rarely possible.

7.7 Results: Discovering Atmospheric Dynamics in Global Climate Models

As described in Section 7.5, we applied CaStLe to output of the Held-Suarez-Williamson-Volcanic atmosphere model, tuned to accurately reproduce the observed Pinatubo response (Hollowed et al., 2024), and the E3SMv2-SPA model including the eruption. In this section, we describe how we applied CaStLe to these case studies and present the results.

7.7.1 Validation with HSW-V

We first note important implementation considerations, particularly how CaStLe's assumptions are satisfied. In general, if assumptions T1, T2, S1, and S2 are uncertain, either because of data availability or dynamical instability, then assumptions can be verified using subject matter expertise. In this study of Mt. Pinatubo, we describe how we carefully managed each assumption prior to applying CaStLe.

In order to be sure CaStLe's assumptions of temporal locality, temporal causal stationarity, and spatial locality (T1, T2, and S1) held in the dataset's 2° grid resolution (corresponding to approximately 214 km at 15 degrees N), we used atmospheric wind speeds at the time of the eruption, which were recorded at 25 m/s on average at 30 hPa; cf. Figure 1 in Thomas et al. (2009). That speed translates to a theoretical maximal aerosol travel distance of 540 km over a 6-hour period, meaning aerosols should move fast enough to traverse one 2° grid cell per time step.

Spatial causal stationarity, assumption S2, is indeed violated considering the globe holistically. We resolved this challenge by using a spatial blocking strategy to create subdivisions in which dynamics were more uniform, and applied CaStLe within each separately. As noted in Section 7.6.6, the selection of blocks and their size is a potential challenge and is highly domain-dependent. We conducted a sensitivity analysis of block sizes, which is presented in H, and determined that dynamics were consistent in various of block sizes. We chose a middle size, $20^{\circ} \times$

20°, for this analysis to balance more nuanced outputs (smaller sizes) with less risk of false positives (larger sizes). This case study was selected for its relatively simple advective dynamics to clearly validate CaStLe and demonstrate its results in an atmospheric setting. We observe that stratospheric winds vary smoothly and slowly, without hard boundaries, which enables us to use a regular grid of blocks. Other settings, such as surface level analyses, the blocking strategy will certainly require special treatment to avoid analysis across hard dynamical boundaries, such as coastlines and mountain ranges. In H, we also demonstrate that blocking alone is not sufficient for non-CaStLed approaches to succeed.

We chose CaStLe's PIP to be the PC-Stable-Single algorithm because in our validation experiments in Section 7.8.1, we found it to be the marginally more effective PIP. However, those experiments showed any tested PIP algorithm is effective. PC-Stable-Single is the PC-Stable causal discovery algorithm (Colombo and Maathuis, 2014) adapted to find the causal parents of only one node; its pseudocode is provided in L. Specific CaStLe parameterizations are given in G. In J, we present similar results using DYNOTEARS for CaStLe's PIP.

Our proxy ground-truth in this case study was stratospheric winds that cause suspended aerosols to advect through space. We display dominant wind fields throughout the space to validate the resulting graphs. Our dataset included wind components in 72 pressure levels in the HSW-V dataset, so we display column-averages of the levels at the levels where volcanic sulfate was most prevalent. Specifically, we chose pressure levels containing more than 5.00 µg of sulphate

Kg air, which were between \sim 6-114.00 hPa. With this, we effectively captured the stratosphere and 56% of all sulfate aerosols in all atmosphere levels. By comparing winds in at the stratospheric levels where most of the sulfur was present, we can directly compare CaStLe's discovery of AOD's space-time evolution to wind data in the same locations.

Comparing the wind and recovered stencils in Figure 7.3, it is clear to see that CaStLe is able to accurately reconstruct the prevailing stratospheric winds using only AOD observations. As these wind fields are the key drivers of aerosol dispersal, it is clear that CaStLe can accurately capture the dynamics dictating the spatial pattern of the Pinatubo response. The CaStLe stencils best capture the underlying wind fields when AOD levels are high. When there are few particles in a region, it is challenging to determine wind by solely observing dispersal patterns. We also observe a zonal (East-West) pattern driving the aerosol dispersion, with Pinatubo aerosols transported nearly fully around the equator within 3 weeks, while meridional (North-South) dispersion taking much longer. This alignment between CaStLe-derived causal structures and observed wind patterns demonstrates the method's effectiveness in reconstructing the physical mechanisms driving aerosol transport, particularly in regions with sufficient particle density to enable clear detection of dispersal trajectories.

Comparative Analysis of CaStLe Versus Traditional Approaches on HSW-V

The current state-of-the-art causal discovery methods cannot tractably approach this study of Mt. Pinatubo's aerosol short-term evolution. As described in Section 7.3, dimensionality reduction techniques commonly used to make them tractable are suitable for spatially static, periodic space-time patterns. However, they are not good solutions for studying a dynamic, transient pattern because modes derived from those techniques are space-timely invariant. Moreover, they are meant to capture large-scale teleconnections, rather than local dynamics that eventually give rise to global phenomena such as teleconnections. For a detailed demonstration of why dimensionality reduction approaches, such as PCA and PCA-varimax, are insufficient for capturing local causal structures in space-time systems like volcanic eruption plumes, see F.

Traditional approaches attempted without dimensionality reduction suffer from the *curse of dimensionality* when applied to short-term global-scale phenomena because there are more grid cells than temporal observations. They also struggle to identify local connections in the massive search space they seek, where every grid cell may be dependent on any other grid cell; i.e., they are not constrained by local causal structure. Finally, their efficiency scales poorly as the grid size gets larger, requiring a lot of time to execute on relatively small grids. We present specifics below and discuss time complexity in depth in Section 7.6.4 and Appendix B.1.

Here, we demonstrate the disparity in performance between traditional approaches and CaStLe for our HSW-V case study using the PC algorithm. The reasons for

the disparity are explored in Sections 7.3 and 7.4. Because PC did not terminate within 48 hours on the full spatial region studied in Section 7.7.1, we restricted the analysis space the area between 20.00° to 50.00° N and 55° W to 120° E in the first 8.5 days after the eruption. On the 2° grid, the given space is equivalent to a 35×35 grid, or 1,225 grid cells. Since temporal observations were 6-hourly, there were 34 time series samples per grid cell.

Figure 7.4 shows the results of the PC causal algorithm and CaStLe-PC-Stable applied to a large section of grid cells for the HSW-V problem. Figure 7.4a illustrates that PC is incapable of reconstructing a graph with any meaningful physical interpretation. There are some local dynamics found, but they are dominated by the many links across disparate locations. PC was implemented here with the partial correlation conditional independence test, a test alpha-value of 0.00001, and a p-value threshold of 0.05 to remove links below that threshold in the final graph. P-values were corrected using the Benjamini-Hochberg procedure prior to final thresholding.

In Figure 7.4b, CaStLe was applied to 10°-by-10° blocks, rather than the 20°-by-20° blocks in Figure 7.3. The smaller block size enables more link density and nuanced results, with the possibility of more mistakes. In this illustration, we chose to display the stencils mapped back to the original space for each block to compare to PC more fairly and demonstrate how much more sparse CaStLe's results are. We found that CaStLe was again able to recover the westward aerosol transport from Mt. Pinatubo. Because HSW-V only models aerosols from the volcano, there is

little to no aerosol signal outside the plume, and results in these areas will be less reliable.

Additionally, the run-time of the PC algorithm is demonstrably poorer than CaStLe. The PC algorithm experiment in Figure 7.4a PC took 65 minutes to execute for a 35×35 grid size. In contrast, the CaStLe experiment in Figure 7.4b completed all blocks serially in 0.46 seconds on the same data. Further, for each of the panels in Figure 7.3, CaStLe computed the 39 stencils for the 3,900 grid cells in a total of 10 seconds. These empirical data points are explained by CaStLe's improved theoretical properties, as detailed in Section 7.6.4 and B.

7.7.2 Extending to More Complexity: E3SMv2-SPA Modeled Aerosols

Given the intended simplicity of the HSW-V model, we also evaluated a simulation of the Mt. Pinatubo eruption in E3SMv2-SPA. More complex graphs arise with a more complex model, providing an opportunity for more nuanced analysis and discovery, but with a higher chance of false positives and false negatives. E3SMv2-SPA is a fully coupled model, so AOD results from many sources including the volcanic eruption and Saharan dust. As such, we expect results to be somewhat noisier, however, as we demonstrate below, CaStLe is still able to identify important features of transport. Because of this additional complexity, we focus on CaStLe as an exploratory tool and leave additional analysis to future work. However, even with the added complexity, CaStLe can obtain compelling results consistent with dominant stratospheric winds as well as the dynamics dis-

covered in our study of HSW-V.

We used 15° spatial blocks so that CaStLe operates on a 5×5 grid space per block. This size strikes a balance in the trade-off that a smaller block-grid enables more nuance in the final output, and larger block-grids take advantage of more spatial replicates to multiply sample size. We chose to study the eruption in two distinct 20-day intervals spanning a six month period to understand the changing evolution of the plume.

Similarly to HSW-V, we utilize the U and V wind fields to visually validate the CaStLe results. In this case, we did not average over multiple altitudes, instead opting to simply use the 50 hPa wind fields; this altitude was shown in Brown et al. (2024, Figure S6) to contain significant levels of the sulfate aerosols.

Figure 7.5 depicts the results of our experiment on E3SM. Again, we applied CaStLe-PC-Stable to construct causal stencils for each given spatial block. We selected two intervals of interest from our results to show here. Day 15 is June 15, 1991, the day of the eruption, so the top row of Figure 7.5 is the first 20 days after the eruption. The bottom row was selected to illustrate later dynamics when aerosols have circumnavigated the tropical zone and more northward advection is present. Days 175-195 are November 22 to December 12, 1991, a little over six months after the eruption.

In the more challenging setting of the fully-coupled E3SMv2-SPA model, our results in the first weeks are still generally consistent with those in HSW-V presented in Section 7.7.1, showing that CaStLe is largely robust to greater com-

plexity. We note that visually identifying the sulfate aerosol plume is much more difficult in this case as the background AOD is quite strong. A solution may be to apply CaStLe to AOD anomalies (computed by subtracting grid cell long-term AOD means from the signal in the analysis period), thus potentially removing background variability from the analysis. However, our goal in this work is to present CaStLe as applied to raw data to illustrate what it can and cannot accomplish in complex, heterogeneous settings.

Regardless, we observe that tropical westward advection is present throughout both studied time periods, but more complexity is present in other regions, in part due to the background AOD. Six months later, the aerosols and winds are in a different regime. We observe northward and southward causal structures in the northern latitudes matching dominant wind fields in the area, with CaStLe stencils still consistent in the tropics. Additionally, CaStLe recovers dynamics moving aerosols northwards above central Asia and southwards through western North America. Causal structures are recovered more often and more accurately where stronger winds coincide with more aerosol presence, building a map of significant aerosol movement. A more complex model and smaller block sizes illustrate more nuanced dynamics, and there is more to learn from these; however, we leave deeper atmospheric dynamics analysis to future work.

7.8 Validation and Benchmarking

In this section, we demonstrate the effectiveness of the CaStLe approach to spacetime causal discovery, highlighting its ability to identify structure in low-signal and data-sparse regimes. We first demonstrate the benefits the CaStLe approach can provide to *any* causal discovery algorithm using a synthetic linear-Gaussian dynamics benchmark; we then apply CaStLe to an important non-linear PDE problem, showing that we can determine the underlying advective forcing.

7.8.1 Evaluating CaStLe: A Comparative Analysis

We demonstrate the effectiveness of CaStLe using a set of local interaction models (LIMs), building upon the comparison framework introduced by Nichol et al. (2023). In summary, we defined a stencil for each experiment that dictates how each grid cell depends on its nine neighbors (including itself). A LIM is a special case of an SCM, which simulates the evolution of a gridded space by computing the current state of each grid cell based on a predefined function of the historical states of its neighbors. In the linear case, this is most simply accomplished with vector autoregression (VAR) models, where the coefficient is sparse, only containing nonzero entries where a desired dependence exists between neighbors. The function is defined by a linear function of coefficients in the given stencil. Our results appear in Figure 7.6, which shows that CaStLe provides significant improvements in graph recovery regardless of the causal discovery algorithm used in

the parent identification phase.

Data: Benchmark Construction

In order to compare different causal discovery algorithms with a common set of benchmarks, we begin by generating coefficient matrices parameterizing spatially homogeneous and statistically stationary VAR(1)s that satisfy our key assumptions S1 and S2. We generate coefficient matrices for these VARs, \tilde{M} , using the following sampling scheme:

- 1. Generate a random 3×3 *local dynamics matrix*, M, with d non-zero elements, one of which is the central element (autocorrelation). Each of the d non-zero elements, $\{a_i\}_{i=1}^d$, have a random value $1.0 \ge \text{coefficient}_i \ge s_*$.
- 2. Expand M to \tilde{M} on a grid of size $N \times N$ (cf. Step D of Algorithm 1 or Figure 2-2 of Nichol et al. (2023))
- 3. If $|\lambda_{\max}(\tilde{M})| \geq 1$, scale \tilde{M} by $|\lambda_{\max}(\tilde{M})|$.
- 4. If $m < s_* \ \forall m \in \tilde{M}$, reject, else accept.

where $|\lambda_{\max}(\tilde{M})|$ is the maximum absolute eigenvalue of \tilde{M} , which when above 1.0 indicates the system is numerically unstable (Strang, 2016, p.307). We note that this process is essentially an accept-reject scheme used to sample from the set of statistically stationary & spatially homogeneous VARs on a 2D grid with minimum signal strengths $s_* \geq 0.1$ and fixed sparsity levels in the range $d \in \{1, 2, \dots, 9\}$.

After each \tilde{M} is generated, we create a single realization, using standard Gaussian noise applied independently, cell-wise at each time step.

Method Comparison: Highlighting CaStLe's Strengths

On each realization, we apply one of three causal discovery algorithms, in both CaStLed and non-CaStLed form: i) the PC algorithm of Spirtes and Glymour (1991) as adapted to time series by Runge et al. (2019a, Algorithm S1 with q=1); ii) PCMCI, an autocorrelated time series extension of PC developed by Runge et al. (2019a); and iii) the DYNOTEARS approach of Pamfil et al. (2020), itself a time series adaption of the NOTEARS approach of Zheng et al. (2018). We additionally compare each of these against a simple sparse VAR approach, where we estimate VAR coefficients directly using ordinary least squares (OLS) and truncate coefficients with magnitude less than s_* ; this approach is not necessarily causal, but it is the exact model of our data generating process and provides a useful point of comparison.

We compare the estimated graph structure with the true graph derived from the sparsity pattern of \tilde{M} and report the average Matthews' Correlation Coefficient (MCC) (Matthews, 1975) and F_1 score over 30 replicates. We used an adapted MCC formula derived by Nichol et al. (2023), which accounts for edge cases in which the denominator would be zero, but is otherwise defined as:

$$MCC = \frac{(TP \times TN - FP \times FN)}{\sqrt{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}}$$
(7.1)

where TP, FP, TN, and FN are true positive count, false positive count, true negative count, and false negative count, respectively. Here, a positive is a graph edge that exists, and a negative is a graph edge that does not exist. The MCC graph similarity measure is sometimes preferable to the more common F_{β} Score (β is chosen such that recall is considered β times as important as precision), which is dependent on the ratio of positive to negative test cases; we treat link positives equally to link negatives, hence our preference for MCC. Figure 7.6 includes the F_1 score due it its common use in causal discovery, but results are similar.

In Figure 7.6, we depict CaStLe performance results on a 2D VAR with ground-truth link density $d=\frac{4}{9}$. We show two extremes of sample size: a low-sample regime of T=10, which is barely enough to identify the local dynamics of 9 cells, and a high-sample regime of T=150. Our results are quite striking: in the low-sample regime, the CaStLed versions of each algorithm can accurately infer graph structure, with near-perfect performance on grids of size 10×10 . By contrast, the performance of the non-CaStLed versions is essentially no better than random guessing, with only the sparse VAR able to exhibit any skill, and then only on small grids. In the high-sample regime, the CaStLed variants perform well on all grid sizes, with CaStLe-PC consistently achieving perfect recovery; the non-CaStLed variants perform better, as expected, but their performance still decays quickly as the spatial grid grows.

While the stronger performance of the CaStLed variants is noteworthy, the exhibited trends are even more important and highlight the true strength of the CaS-

tLe approach: CaStLed approaches *improve* on larger grids while traditional approaches suffer. While Figure 7.6 shows results for the fixed link density $d = \frac{4}{9}$, we present results for all other link densities in K.

Having established CaStLe's strong performance on linear dynamics, we also validated its effectiveness on non-linear systems that more closely resemble realistic physical processes in Earth science. Specifically, we applied CaStLe to the advection-diffusion dynamics of Burgers' equation, a fundamental non-linear PDE that models a combination of advective and diffusive processes. Unlike our VAR benchmarks, which are discrete linear models with random initializations, Burgers' equation presents continuous non-linear dynamics that allow us to evaluate CaStLe's ability to recover spatial propagation patterns under controlled conditions. Our analysis demonstrates that CaStLe successfully identifies the underlying advection angle across a range of diffusion conditions, further supporting its applicability to complex space-time systems. This non-linear validation's complete methodology and results are presented in D.

7.9 Discussion

We have introduced CaStLe, a novel causal discovery meta-algorithm tailored for analyzing grid-level space-time data sets arising in Earth science. CaStLe can be directly applied to grid-level data and does not require pre-processing and spatial dimension reduction, allowing it to capture dynamics in the natural domain of the data rather than a derived (PCA-type) space. This distinction is crucial be-

cause global-scale phenomena across many complex systems—whether climate teleconnections, ecological patterns, or fluid dynamics—emerge from networks of local causal interactions that are often lost in dimensionality reduction approaches. While demonstrated with Earth science case studies, CaStLe is fundamentally domain-agnostic, applicable to any space-time system governed by local physical interactions, from fluid dynamics and heat transfer to biological pattern formation.

CaStLe can overcome the limitations of existing causal discovery approaches in Earth science's space-time data, filling a significant gap. By leveraging realistic assumptions of locality and homogeneity, CaStLe creates "spatial replicates" to substitute large observational domains for lengthy time series. This process transforms the spatial causal discovery problem from the high-dimensional (many variables, few observations) to the low-dimensional (few variables, many observations) regime, allowing accurate and efficient discovery of underlying causal dynamics. A key aspect of CaStLe is the causal *stencil* graph, a simplified representation of the local dynamics driving larger global behaviors. This notion of a stencil is particularly well-suited for systems able to be modeled by PDEs, as PDE-type dynamics inherently enforce both locality and homogeneity, as well as the sufficiency assumptions necessary for causal discovery to be *truly causal*.

We used these insights to identify the space-time evolution of volcanic aerosols that erupted from Mount Pinatubo in the HSW-V and E3SMv2-SPA models. We found that CaStLe found the expected path of advection in both models and more nuanced dynamics, including northward and southward dispersion, in E3SMv2-

SPA. We showed that CaStLe outperforms its peers in the causal discovery of synthetic benchmarks generated by vector autoregressive structural causal models. Additionally, as detailed in D, we found that CaStLe could accurately identify the advection angle in our Burgers' equation benchmark, demonstrating that it can filter out the "noise" of diffusion.

Our brief theoretical analysis of CaStLe in Section 7.6.4 and in B, demonstrates two regimes of consistent estimation for CaStLe, i.e., CaStLe recovers the true causal dynamics: long time series $(T \to \infty)$ or large grid sizes $(N \to \infty)$. This starkly contrasts existing approaches, whose performance rapidly deteriorates as $N \to \infty$. Several other important theoretical questions remain open, including the optimal relationship between sampling rates and grid resolution, behavior under mild violation of the key assumptions, and the correct target of inference for systems without clear advective dynamics (e.g., the chemical evolution of atmospheric aerosols).

We have focused on space-time data observed on regular 2D grids, but we believe that this assumption can be relaxed to adapt CaStLe for a broader range of observational structures. CaStLe can also be adapted to multivariate space-time data (more than one observation at each point) by including more comeasured variables in CaStLe's transformation of the region to the reduced coordinate space, enabling causal discovery of the space-time interactions of multiple species on the grid-level, which is a particularly exciting avenue of future research and application to Earth system dynamics. Developing data-driven methods for evaluating

block sizes based on output robustness will enable more automatic application of CaStLe, requiring less subject matter expertise. Finally, causal representation learning is a nascent field combining the estimation power of machine learning with the strength of inference of causal discovery. Applying these techniques in CaStLe's parent-identification phase or for discovering spatial embeddings for regional block analysis is an exciting potential direction for future work.

Because our assumptions are readily satisfied by many physical systems, CaStLe can be applied quite broadly in the physical sciences. It may find value in any space-time system in which quantities at every point in space impact their adjacent spatial neighbors. In the Earth system, it may be of particular interest for studying forest fires, ocean dynamics, salt/fresh water incursions, and coastal erosion, for example. For atmospheric rivers, CaStLe could identify pathways of moisture transport and evolution; for wildfire spread, it could reveal causal relationships between local weather conditions and fire behavior; for drought propagation, it could track how soil moisture deficits spread across regions. CaStLe's preservation of local causal structures while efficiently handling high-dimensional data offers advantages over approaches requiring dimension reduction. For datasets where the temporal sampling is too coarse relative to the spatial resolution, extending to a radius-2 neighborhood might be appropriate while still maintaining our core assumption of locality. This extension would preserve the fundamental CaStLe methodology—only the dimensionality of the reduced coordinate space would increase. Additionally, CaStLe provides a promising framework for Earth system model evaluation (Nowack et al., 2020a; Nichol et al., 2021), potentially identifying where models produce correct outcomes through incorrect causal mechanisms.

While climate science typically studies large, long-term phenomena, the community increasingly recognizes the importance of understanding multi-scale interactions (Diffenbaugh et al., 2005; Palu, 2019; Agarwal et al., 2019; Zhang et al., 2022). Teleconnections present an exciting challenge for future applications of CaStLe. These statistical dependencies between distant regions appear to violate locality but physically result from countless local interactions that are often unobserved or unmodeled. A two-stage methodology could be effective for tackling this challenge. First, apply CaStLe to discover local causal stencils, and then apply a complementary causal discovery technique to connect the discovered local processes across scales. This approach could bridge the gap between local and global causal discovery in climate science.

Complex space-time systems present apex challenges for causal discovery, combining chaotic dynamics, high dimensionality, noisy observational records, and complex underlying physical processes. CaStLe represents the first successful application of causal graph discovery to learn grid-cell-level causal structures in Earth systems. By preserving local causal structures while efficiently handling high-dimensional data, CaStLe presents a path toward connecting micro-scale interactions with macro-scale phenomena, potentially offering new insights into how global patterns emerge from local causal mechanisms. There are rich future research directions, including multivariate analysis and automated block size selec-

tion. The feasible discovery of local causal stencils presents an exciting new frontier for causal discovery of space-time data, particularly in the Earth sciences.

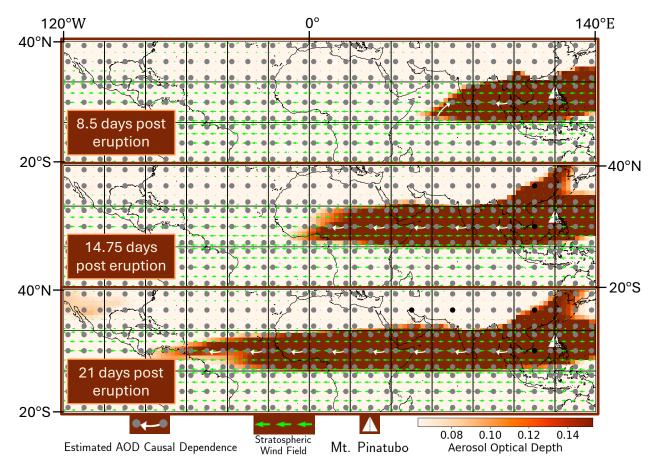


Figure 7.3: Application of CaStLe-PC-Stable to HSW-V simulation of the 1991 Mt. Pinatubo eruption. The stencils estimated by CaStLe (white) capture the underlying high-altitude wind fields (green) using only satellite-measured AOD, with near perfect accuracy in high aerosol regions (red-orange). Autodependencies are shown with black nodes where grid cells cause themselves, and gray nodes where there is no autodependence. All links represent a six hour time lag, the time resolution of the HSW-V dataset. On longer horizons (bottom row), CaStLe is able to recover equatorial wind currents as far away as South America, half-way around the world from Mt. Pinatubo (white triangle). CaStLe accurately identifies the prevailing westerly atmospheric winds because it was able to identify the space-time dependence between neighboring grid cells. Additional details are given in Section 7.7.

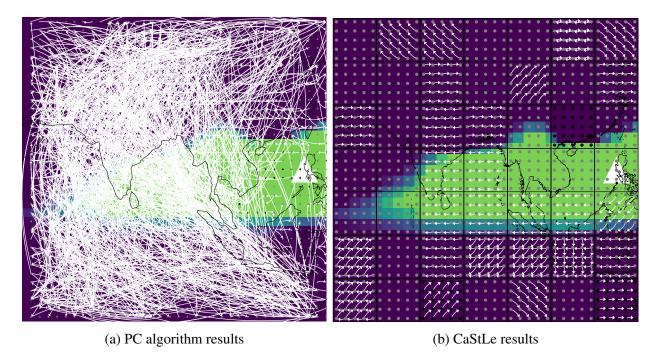


Figure 7.4: Causal maps inferred from the PC algorithm applied naively to all grid cells and CaStLe's equivalent results immediately to the west of Mt. Pinatubo; a 35×35 grid between -20.00° to 50.00° N and 55.00° to 125.00° E in a 8.5 day span after the eruption. All links represent a six hour time lag, the time resolution of the HSW-V dataset. As expected, PC struggled with the high dimensionality and the discovered dependencies do not conform to the ground-truth understanding that aerosols advected towards the west. It also fails to identify local dynamics, instead drawing most connections over great distances. The PC analysis was computed in 729 minutes on 1,600 grid cells, while the CaStLe analysis was computed in 0.46 seconds.

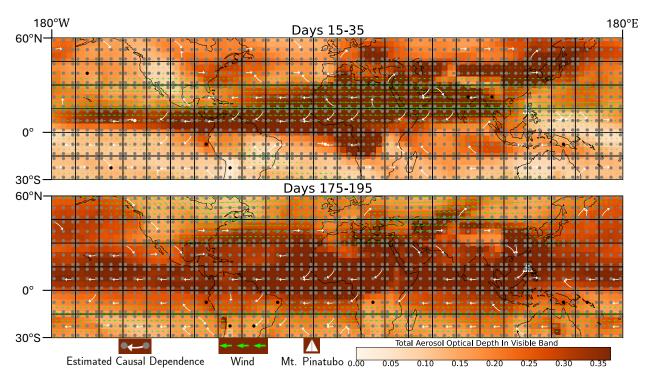


Figure 7.5: Application of CaStLe-PC-Stable to E3SMv2-SPA simulation of the 1991 Mt. Pinatubo eruption. The stencils estimated by CaStLe (white) capture the underlying high-altitude wind fields (green) using only total aerosol optical depth (AOD). Autodependencies are shown with black nodes where grid cells cause themselves, and gray nodes where there is no autodependence. All links represent a one day time lag, the time resolution of the E3SMv2-SPA dataset. The heatmap depicts AOD from any source at 50 hPa. The top panel depicts learning from the first 20 days after eruption, which began on day 15. The bottom panel depicts learning approx 6 months after the eruption over a 20-day time period. In the more challenging setting of the fully-coupled E3SMv2-SPA model, our results in the first weeks are still generally consistent with those in HSW-V presented in Section 7.7.1, showing that CaStLe is largely robust to greater complexity. In the bottom panel, the aerosols and winds are in a different regime. CaStLe stencils are still consistent in the tropics and now begin to recover dynamics pushing aerosols northwards above central Asia and southwards through western North America. A more complex model and smaller block sizes illustrate more nuanced dynamics, and there is more to learn from these, however, we leave deeper atmospheric dynamics analysis to future work.

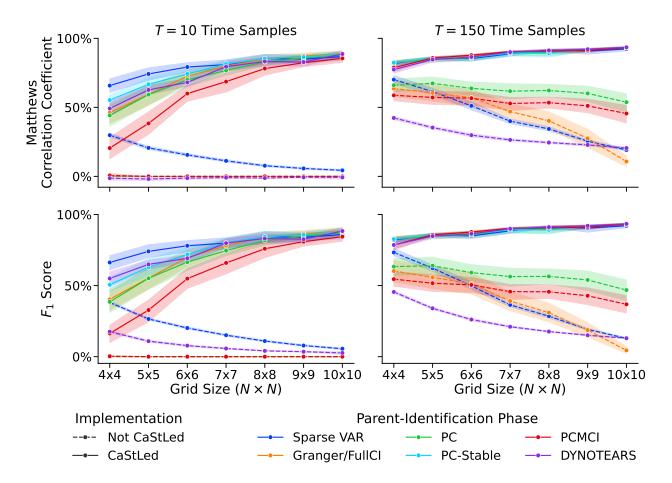


Figure 7.6: Comparison of CaStLed and non-CaStLed causal discovery approaches on linear-Gaussian dynamics, including Granger causality or FullCI (orange), PC (green), PCMCI (red), and DYNOTEARS (purple), as well as a statistical model of the data generating process (blue) presented with both MCC and F_1 metrics. In the low-sample size regime (T=10, left) CaStLed approaches can accurately recover the underlying causal graph, with performance increasing on larger grid sizes (solid lines); by contrast, non-CaStLed approaches are unable to perform better than mere chance (dashed lines). Even a model based on the underlying data generating process (Sparse VAR, blue) is significantly outperformed by its CaStLed counterpart. In the high-sample size regime (T=150, right), non-CaStLe approaches have improved performance but still compare unfavorably with their CaStLed counterparts.

Appendices

Appendices

Table 1: Capabilities of CaStLe for Earth science applications. This table summarizes the key methodological advantages of CaStLe and their relevance to specific Earth science phenomena, highlighting applications where grid-level causal discovery enables analyses that were previously infeasible with prior causal discovery approaches.

Capability	Description	Relevant Applications
Local mechanism discovery	from local causal interactions. Previous approaches use dimensionality reduction,	Volcanic plume transport (Sjolte et al., 2021), wildfire propagation & plume transport (Baranowski et al., 2021), atmospheric rivers (Payne et al., 2020; Baño-Medina et al., 2025; Higgins et al., 2025)
Transient, non-periodic phenomena	CaStLe effectively identifies grid-level causal pathways.	Volcanic eruptions, heat waves (Keellings and Moradkhani, 2020), wildfires (Driscoll et al., 2024)
High- dimensional data settings	replicates to make high-	Gridded Earth science data from: regional climate modeling, satellite observation analysis, climate reanalysis products (Ali et al., 2024, Table 3)
•	models and observations at the grid level, poten- tially	

A Understanding Assumptions

In this section, we outline the key assumptions underpinning the CaStLe framework and their relationship to causal discovery assumptions.

A.1 CaStLe Assumptions

CaStLe operates via two complementary sets of assumptions:

- 1. CaStLe Framework Assumptions (T1, S1, T2, S2): These enable efficient use of spatiotemporal data by leveraging locality and stationarity to transform a high-dimensional problem into a tractable one.
- 2. **Causal Discovery Assumptions**: The causal discovery algorithm used within CaStLe's Parent Identification Phase requires its own set of assumptions typically the Causal Markov Condition, Faithfulness, and Causal Sufficiency.

While these assumption sets are conceptually distinct and serve different purposes, they work together to enable scalable causal discovery in high-dimensional space-time systems.

In review, our framework introduces four key assumptions to capture a "PDE-like" system X_t , creating an environment where local space-time dynamics can be efficiently learned:

T1) Temporal Locality: restricts causal influence the most recent past state, one time lag, aligning with how PDEs are discretized.

- T2) Temporal Causal Stationarity: ensures consistent causal structure over time.
- **S1**) Spatial Locality: limits causal influence to immediate spatial neighbors.
- **S2**) Spatial Causal Stationarity: ensures consistent causal structure across space.

These assumptions enable CaStLe to leverage "spatial replicates"—treating each local neighborhood as providing information about the same underlying causal process. This transforms what would be a high-dimensional, data-sparse problem (many variables, few observations) into a data-rich problem (few variables, many observations).

A.2 Causal Discovery Assumptions

Separately, the causal discovery algorithm used within CaStLe's PIP require its own assumptions. The three foundational assumptions of causal discovery are detailed in Runge (2018a) and in Spirtes et al. (1993, Ch. 3):

• Causal Markov condition: for $\mathbf{X} = \{X_1, X_2, \dots, X_n\}$, each variable X_i is conditionally independent of its non-effects given its direct causes $\mathcal{P}(X_i)$:

$$X_i \perp \!\!\!\perp \!\!\!\perp X \setminus \mathscr{P}(X_i) \mid \mathscr{P}(X_i)$$

- A variable is conditionally independent of its non-effects given its direct causes.
- Faithfulness: if X_i and X_j are statistically dependent, then there exists a direct

causal link or a common cause:

$$X_i \not\perp \!\!\! \perp X_j \implies \exists$$
 a direct cause or common cause of X_i and X_j

Conversely, if X_i and X_j are conditionally independent given their parents $\mathscr{P}(X_i)$ and $\mathscr{P}(X_j)$:

$$X_i \perp \!\!\! \perp X_j \mid \mathscr{P}(X_i), \mathscr{P}(X_j) \Longrightarrow \text{ no direct causal link between } X_i \text{ and } X_j$$

- All conditional independencies in the data arise from the causal structure (no accidental cancellations).
- Causal sufficiency: all common causes of observed variables are also observed.

A.3 Relationship Between Assumption Sets

While CaStLe assumptions (T1-S2) and causal discovery assumptions serve different purposes, there are important interactions between them:

- CaStLe assumptions create an environment where causal discovery becomes tractable in some high-dimensional gridded settings.
- CaStLe assumptions do not guarantee causal discovery assumptions will be satisfied.
- For example, even in perfectly stationary systems (T2, S2 satisfied), faithfulness can be violated through counteracting mechanisms, as demonstrated in

Runge (2018a).

• Similarly, the Causal Markov Condition is a property of the joint distribution that cannot be derived from locality assumptions.

Instead of replacing causal discovery assumptions, CaStLe's assumptions create a context where causal discovery methods can be applied efficiently to high-dimensional space-time data.

CaStLe's Implementation and Causal Sufficiency

One meaningful connection exists between CaStLe's implementation and causal discovery assumptions: When CaStLe focuses on identifying only the parents of the center cell while including all potential spatial neighbors (per assumption S1), causal sufficiency is automatically satisfied for that specific node by construction - assuming S1 holds true.

This is a significant benefit, as causal sufficiency is typically the most difficult assumption to guarantee in practice (Spirtes et al., 1993; Raghu et al., 2018). While CaStLe cannot guarantee faithfulness or the Markov condition holds, its design cleverly leverages spatial structure to address causal sufficiency within each local analysis.

A.4 Potential Violations and Their Manifestations

Violations of CaStLe's assumptions can occur in various ways, leading to different manifestations in the causal discovery process. Violations of CaStLe's assump-

tions can affect results in different ways:

- 1. Violations of Temporal/Spatial Locality (T1, S1): If causal effects extend beyond immediate neighbors, CaStLe will miss these connections, creating false negatives.
- 2. Violations of Stationarity (T2, S2): If dynamics change across space or time, CaStLe's stencil will represent only an average pattern, potentially creating both false positives and negatives.
- 3. Even with CaStLe assumptions holding, traditional faithfulness violations can occur through cancellation effects or deterministic relationships.

Below, we provide examples of how these assumptions can be violated and their potential impacts, drawing on the discussion by Runge (2018a).

Temporal and Spatial Locality (T1, S1)

- *General Violation*: These assumptions can be violated by any process that introduces dependencies beyond immediate temporal or spatial neighbors.
- Example Time Aggregation: Time aggregation can violate temporal locality by introducing dependencies across multiple time steps. Runge (2018a) discusses how time aggregation can cause such violations (Section IV.B, Example 4). Figure 5 in Runge (2018a) illustrates the impact of time aggregation on causal inference.

• Example - Spatial Aggregation: Similarly, spatial aggregation can violate spatial locality by introducing dependencies across non-neighboring spatial units.

Temporal and Spatial Causal Stationarity (T2, S2)

- *General Violation*: These assumptions can be violated by any process that introduces changes in the causal relationships over time or space.
- Example Counteracting Mechanisms: Counteracting mechanisms or heterogeneous processes can violate these stationarity assumptions. If the data contains opposing generating processes (e.g., different hemispheres in climate data), the faithfulness assumption may be violated. This results in unstable and inconsistent causal relationships. Runge (2018a) discusses such violations in Section IV.C, Example 5, and provides an illustration in Figure 6.

Understanding potential violations and their manifestations is crucial for applying our framework effectively in realistic scenarios. Section 7.6.6 outlines practical strategies to mitigate these violations.

B Statistical and Time Complexity

In this section, we elaborate on Section 7.6.4 and provide a more detailed discussion of the time-complexity (Appendix B.1) and statistical (Appendix B.2) properties of CaStLe. Additionally, we provide analyses giving conditions under which

CaStLe is (asymptotically) guaranteed to recover the true causal graph, independent of the specific PIP used.

B.1 Time Complexity

Steps A, B, and D of CaStLe consist primarily of copying and rearranging of data, so we focus our analysis on the complexity of Step C, which dominates the runtime of CaStLe. Because CaStLe can use a variety of PIPs within Step C, we begin with a general analysis of the worst-case time complexity of causal discovery algorithms. Throughout, recall that a runtime complexity $\mathcal{O}(f(n))$ implies there exists a fixed constant $C \ge 0$ such that that the algorithm terminates in at most Cf(n) steps for any input of size n.

Kalisch and Bühlmann (2007) and Runge (2018a) discuss the time complexity of causal discovery, particularly the PC algorithm. Much of constraint-based causal discovery is descendant of PC, and it represents a valuable baseline for comparing the computational complexity of CaStLe and prior work. Causal discovery is largely bounded by how long it requires to determine independence between nodes (bounded by samples and size of conditioning sets of nodes) and how many times it needs to do so (generally bounded by the number of nodes). Runge (2018a) cite the time complexity of a single conditional independence test using ordinary least squares (linear partial correlation), while Kalisch and Bühlmann (2007) explore bounds on the number of tests in PC. Our analysis is consistent with theirs, which we derive from first principles.

Consider causal discovery in p-dimensions (p measured variables) with n samples; suppose further that it is known, a priori, that any node in the causal graph has at most degree q: that is, no element has more than q causal parents. An exhaustive search for the causal parents of a single node will require evaluating $\sum_{i=0}^{q} \binom{p}{i} = \mathcal{O}(2^p)$ possible sets of parents; repeating this process for all p nodes evaluation of up to $\mathcal{O}(p2^p)$ possible causal graphs. If we construct graphs using statistical tests for linear partial (conditional) correlation, each test can be performed in $\mathcal{O}(np\min\{n,p\}) = \mathcal{O}(np^2)$ time (the time required to fit an OLS regression to n observations and p variables using a direct method such as an SVD or QR factorization), yielding an overall runtime of

$$\mathscr{O}(np^2 * p2^p) = \mathscr{O}(np^3 2^p).$$

This analysis is quite loose, and as Runge (2018a) notes, the complexity of a *sin-gle* linear conditional independence test can be reduced to $\mathcal{O}(np^2q^2)$ when efficient algorithms are used. Far stronger guarantees can be provided for specific causal discovery algorithms that more efficiently search the space of possible graphs. Regardless, even this rough analysis will be sufficient to demonstrate the algorithmic improvements attained by CaStLe.

We now consider the specific context of causal discovery from gridded time series data. Here, we have n = T total observations and have $p = N^2$ features of our data. Direct application of causal discovery to this data gives a worst-case

complexity of

$$\mathscr{O}(np^32^p) = \mathscr{O}(T(N^2)^32^{N^2}) = \mathscr{O}(TN^62^{N^2}),$$

so the complexity of standard causal discovery methods grows *super-exponentially* with the size of the grid. For the purposes of direct comparison to CaStLe, where $p = N^2$, we assume PC's $\tau_{max} = 1$. By contrast, the reduced space where CaStLe's PIP operates has $T(N-2)^2$ observations and only p = 9 features, yielding a *polynomial* worst-case runtime of

$$\mathscr{O}(np^32^p) = \mathscr{O}(T(N-2)^2 * 9^3 * 2^9) = \mathscr{O}(TN^2).$$

Even for grids of relatively modest size, this improvement can be significant: consider a small 30×30 grid; at 1° resolution, this covers approximately 1.5% of the globe. Unstructured causal discovery methods need to consider approximately 30^6*2^{30} possible graphs, while CaStLe needs to evaluate only $9^3*2^9 = 373,248$ graphs, representing an improvement of approximately 2×10^{12} -fold. Specific PIPs may provide less dramatic improvements, but it is clear that CaStLe can be expected to be millions-if not billions-of times more efficient than existing approaches.

Note that in our application scenarios, CaStLe is always applied to a square $N \times N$ grid. However, more generally we can consider p grid cells. Traditional causal discovery will be bounded by

$$\mathcal{O}(Tp^32^p),$$

while CaStLe will be bounded by

$$\mathcal{O}(Tp)$$
.

Thus, if grid cells scale linearly, CaStLe scales linearly in both samples and grid cells.

B.2 Statistical Consistency

Statistically, we see that CaStLe can achieve significantly improved estimation performance compared to a full graph inference approach. Rather than give a general analysis, we rely on the prior work of Kalisch and Bühlmann (2007) to compare CaStLe-PC with the standard PC algorithm. Using the same definitions of n, p, q as in our previous analysis, Kalisch and Bühlmann (2007, Appendix B) show that the probability of the PC algorithm incorrectly estimating the causal graph incorrectly is bounded above by

$$P[\hat{\mathscr{G}} \neq \mathscr{G}] = \mathscr{O}\left(p^{q+2}(n-q)e^{-c(n-q)}\right).$$

In our setting, this gives an error probability of

$$\mathscr{O}\left(p^{q+2}(n-q)e^{-c(n-q)}\right) = \mathscr{O}\left((N^2)^{N^2+2}(T-N^2)e^{-c(T-N^2)}\right) = \mathscr{O}\left(N^{2N^2}e^{cN^2}*Te^{-cT}\right)$$

for PC applied in the original data space. It is clear that this quantity grows rapidly in N, consistent with the intuition that causal discovery algorithms struggle when

applied to larger spatial domains. By contrast, this analysis implies that the error probability of CaStLe-PC scales as

$$\mathscr{O}\left(p^{q+2}(n-q)e^{-c(n-q)}\right) = \mathscr{O}\left(9^{9+2}(T(N-2)^2 - 9)e^{-c(T(N-2)^2 - 9)}\right) = \mathscr{O}\left(\frac{TN^2}{e^{TN^2}}\right)$$

Quite surprisingly, this *decreases* with the graph size (*N*), implying that CaStLe actually achieves *better performance* when applied to larger spatial domains. We demonstrate the remarkable practical effect of this scaling in Section 7.8.1. Similar improvements can be shown for any base causal discovery algorithm (and associated PIP) for which precise estimates of statistical convergence rates are available.

C Asymptotic Consistency

We examine the asymptotic consistency of CaStLe, with a particular focus on the Parent Identification Phase (PIP). Asymptotic consistency is a fundamental property that ensures the accuracy of causal graph estimates as the number of observations increases. We begin by establishing the technical assumptions necessary for our analysis, specifically those related to the p-values generated by the PIP for edge existence. These assumptions are critical for maintaining control over both false positive and false negative rates, thereby ensuring the reliability of our causal inferences. The central theorem we present demonstrates that, under these conditions, CaStLe achieves asymptotic consistency as the number of nodes approaches infinity. In the case of Bayesian score optimization causal discovery,

such as DYNOTEARS, Bayesian posterior probabilities can be used in lieu of p-values with suitable minor modifications to the combination procedure. The proof is structured into three parts, addressing the independence of observations, the application of Fisher's method for combining p-values, and the implications of using overlapping regions. Through this analysis, we aim to reinforce the validity of our algorithm and its effectiveness in uncovering causal relationships in gridded space-time data structures.

Technical Assumption (P1):

- The Parent Identification Phase, $PIP(\cdot)$, produces p-values for edge existence, which satisfy the following:
 - For every non-edge (i,j) $(j \notin \mathscr{P}(i))$, $\mathbb{P}(p_{\mathrm{PIP}}^{(i,j)} \leq u) = u$ for all $u \in [0,1]$; that is $p_{\mathrm{PIP}}^{(i,j)} \sim \mathscr{U}([0,1])$ is uniformly distributed.
 - For every edge (i, j) $(j \notin \mathcal{P}(i))$ and every $T > T_0$, there exists $\pi_{(i, j)}^T(u) > 0$ such that $\mathbb{P}(p_{\text{PIP}}^{(i, j)} \leq u) \leq \max\{0, u \pi_{(i, j)}^T(u)\} < u$ for all $u \in [0, 1]$.

Taken together, these require that the $PIP(\cdot)$ control the false positive rate at the nominal significance level used and that the false negative rate is less than the false positive rate.

Here, T_0 is a minor technical assumption to allow the PIP to have non-trivial accuracy: we use it to exclude trivial cases like T=1, in which no time series causal discovery mechanism can be accurate.

Additionally, note that we typically assume that the $PIP(\cdot)$ is asymptotically

consistent, so that $\pi_{(i,j)}^T(u)$ is bounded above 0 for all u as $T \to \infty$. This can be used to prove T-asymptotic consistency of CaStLe, but in this section we aim only to prove N-asymptotic consistency.

Theorem: Suppose \mathscr{D} is an $\mathbb{R}^{T \times N \times N}$ realization of a data-generating process satisfying T1-S2. Suppose also that $\mathsf{PIP}(\cdot)$ is a parent-identification-phase satisfying P1. Then, there exists a T_0 such that for any $T \geq T_0$, CaStLe is asymptotically consistent as $N \to \infty$; that is, the causal graph estimated by CaStLe converges to the true causal graph generating \mathscr{D} with probability 1.

Proof. This proof proceeds in three parts:

- First, we argue that, for large *N*, well-separated (non-overlapping) spatial regions can be considered IID realizations.
- Next, we argue that the application of Fisher's method leads to asymptotic consistency of CaStLe.
- Finally, we argue that "infill" of the overlapping regions does not invalidate the asymptotic consistency.

At a high level, we argue that, because it is T-asymptotically consistent, there exists some T_0 where the PIP has non-trivial power. We then apply standard statistical methods for combining several weak p-values to obtain a global strong p-value. The technical bookkeeping of our argument serves primarily to deal with the fact that we use overlapping spatial regions and cannot assume independence of

the individual *p*-values; we overcome this by selecting regions that are sufficiently spatially separated to be statistically independent on the time scale considered.

Without loss of generality, we focus on asymptotically consistent estimation of a single edge, say (East, Center). Extension to all 9 stencil edges follows immediately by a standard union bound argument.

Part I: For analytical simplicity, we divide the spatial region into square regions of size $(5+2T) \times (5+2T)$. On a grid of size $N \times N$, there are $B_{N,T} = \lfloor N/(5+2T) \rfloor$ such regions. We apply the PIP(·) to the center 3×3 region of each region separately, obtaining $B_{N,T}$ p-values for the existence of the edge. Because these central regions are separated by (at least) 2T + 2 grid cells and causal effects exist at a distance of at most 2T under our data generating model, these p-values can be treated as statistically independent. (This is essentially the same argument used by Goerg and Shalizi (2013), though their application is quite different.)

Part II: Given $B_{N,T}$ independent p-values, we then apply Fisher's method for combining p-values. Specifically, given a set of p-values for edge non-existence, Fisher's method controls the *familywise error-rate*, rejecting the global null (no edges anywhere). By our assumption of spatial homogeneity, if an edge exists in at least one region, it must exist everywhere, so Fisher's method precisely tests for edge existence in the stencil.

Recall that Fisher's method constructs a test statistic $T = -2\sum_{b=1}^{B} \log p_b$ and tests it against a null χ_B^2 distribution. We consider two cases:

1. If the edge does not exist, each p-value is $\mathscr{U}([0,1])$ by construction and the

test statistic T follows its null distribution. So long as the global significance level used for Fisher's test α_{Fisher} is converging to 0 as $N \to \infty$, we have asymptotic consistency for edge absence.

2. If the edge does exist, each p-value is less than α with probability $(1+c)\alpha$ for some c strictly positive. We then have that T has a non-central χ^2 distribution, which is asymptotically distinguishable from a (central) χ^2 at all significance levels as $N \propto B \to \infty$.

Taken together, these guarantee the the output of Fisher's method is asymptotically consistent for both edge presence and edge absence.

Part III: In practice, we apply CaStLe not to disjoint regions but to overlapping regions. As discussed elsewhere, the region-discretization strategy and the use of Fisher's method are such that this does not cause "cross-contamination" or invalid tests of edge existence. We note here that this strategy also does not invalidate asymptotic consistency of CaStLe. Specifically, we note that, with overlapping regions, the *p*-values used in Fisher's method may no longer be assumed independent.

In this case, however, this is not an issue as they exhibit positive dependence (as they are taken from overlapping data). As such, the true degrees of freedom of T under the null are less than the nominal degrees of freedom; this leads Fisher's method to be (if anything) overly conservative in finite samples. Hence, for the case of edge absence, the nominal significance level is understated and we retain consistency as long as we take $\alpha_{\text{Fisher}} \xrightarrow{N \to \infty} 0$; for the case of edge presence, it

suffices to note that the true sampling distribution is still asymptotically distinguishable from the null (since each individual p-value is powerful), so we retain consistency.

We note that Fisher's method may not be the optimal method for combining p-values. In particular, Holm's method allows for arbitrary dependence of the p-values, likely yielding better performance at finite N, but we do not pursue this approach here as the implementation and theoretical analysis are somewhat more difficult. As with Fisher's method, Holm's method controls the error rate of the global null which, under our assumptions of causal stationarity, is precisely the correct null for accurate stencil estimation.

Additionally, we note that the p-values produced by the PIP under the null do not need to precisely satisfy a uniform distribution; conservative p-values decrease the value of Fisher's statistic T, thereby lowering the rate of false positives.

Remark: If $PIP(\cdot)$ is strongly asymptotically consistent as $T \to \infty$, it must satisfy assumption P1.

Proof. We argue by contradiction. Suppose that $PIP(\cdot)$ were not asymptotically consistent and that the false positive rates and false negative rates of the PIP were equal (or worse, the false negative rate was greater than the false positive rate). Specifically, assume that there exists a true edge (i, j) and some $\pi_- > 0$ such that $\mathbb{P}(p_{PIP}^{(i,j)} \leq u) > \pi_- + u$ for all T and all u. For the PIP to guarantee no false positives, we must take $\alpha \to 0$ as $T \to \infty$. But this would imply that there remains an asymptotic π_- probability of a false negative $(\mathbb{P}(p_{PIP}^{(i,j)} \leq \alpha) > \alpha + \pi_i \geq \pi_- > 0)$,

contradicting our assumption of asymptotic consistency.

D Application to Non-Linear Dynamics: Continuous Systems via Burgers' Equation

This appendix extends our validation of CaStLe to non-linear dynamical systems through application to Burgers' equation, demonstrating the method's effectiveness beyond the linear systems discussed in the main text.

Having established the strong performance of CaStLe on discrete models of linear dynamics, we turn to a far more challenging domain: continuous models with non-linear PDEs. Specifically, motivated by our interest in turbulent atmospheric dynamics, we consider Burgers' equation, a PDE used to model a combination of advective (directed flow) and diffusive processes (Burgers, 1948). While initially developed to model fluid flows, Burgers' equation has been successfully applied to a variety of fields, such as turbulence, non-linear wave propagation, traffic flow, cosmology, gas dynamics, and more (Bonkile et al., 2018). In the following experiments, we again implemented CaStLe's PIP with the PC-Stable-Single algorithm.

We note that the interaction of PDE dynamics with causal language is rather subtle: while PDEs are imbued with a "forward" direction in time, the actual numerical methods used to solve them include "forward" and "backward" steps in the underlying integrators as well as sophisticated interpolation schemes. Our focus here is not on finding a causal model for the PDE solution per se, but on identifying

the structure of the underlying advection. This choice is motivated in part by the results of Rubenstein et al. (2018), who explored the related problem of identifying causal models from deterministic ordinary differential equations (ODEs). As they note, there is not generally a single causal graph corresponding to an ODE, with different models being appropriate at equilibrium or under various interventions. Given the additional complexity of PDEs, we believe that identifying the underlying advection angle provides the most meaningful causal representation of Burgers-type dynamics, particularly as it relates to our volcanic eruption aerosol case study.

D.1 Burgers' Equation: Model and Parameters

In two dimensions, Burgers' equation can be written as:

$$\frac{\partial u}{\partial t} + u \left(\alpha \frac{\partial u}{\partial x} + \beta \frac{\partial u}{\partial y} \right) = c \left(\frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} \right) + f$$
(2)
Advective Dynamics
Diffusive Dynamics

where α, β are the advection coefficients in the x, y directions, capturing directed flow dynamics; c is the diffusion coefficient; and f is a forcing term representing additional mass being injected into the system. In order to create a closed system with no exogenous forcings, we take f = 0 uniformly throughout this section.

The left panel of Figure D1 shows three different solutions to Burgers' equation at different advection angles (θ), advection strength ($M = \sqrt{\alpha^2 + \beta^2}$), and diffusivities (c), each with the same initial conditions. Examining the time evolution of

these solutions (left to right), we see that the high-advection low-diffusion systems (top) exhibit a clear direction of flow, while it is far more difficult to find direction in low-advection high-diffusion systems (bottom). We take inferring the angle of advection as our principal task: given an observed solution u to Equation (2), can we determine the angle of the underlying advective dynamics?

D.2 Advection Angle Estimation

Given a CaStLe-estimated stencil, we infer the angle of underlying advection in the following manner: i) identify each potential parent edge of C with a vector, taking the angle of the underlying edge in the reduced space as direction and the (signed) strength of the underlying relationship as magnitude; ii) sum these vectors to obtain an aggregate estimate of the advective dynamics; iii) take the angle of the vector sum as an estimate of the underlying advection angle. In pseudo-code, we can write this as

$$\hat{\theta} = \mathtt{atan2}\left(\sum_{l \in \mathscr{P}(\mathtt{C})} e_l \sin \theta_l, \sum_{l \in \mathscr{P}(\mathtt{C})} e_l \cos \theta_l\right).$$

Here atan2 is the signed arctangent function, $\mathcal{P}(C) = \{NW, N, \dots, W\}$ represents all potential parents the center cell, e_l represents the strength of that edge (0 for non-present edges), and θ_l represents the angle of that edge (135°, 90°, ..., 180°). This process allows us to estimate all angles instead of just the eight angles present in the stencil structure.

D.3 Experimental Setup

In order to assess the effectiveness of CaStLe-PC in a variety of regimes, we generate (approximate) solutions to Equation (2) with 500 angles sampled uniformly from $[0^{\circ}, 360^{\circ})$, advection magnitudes varying from 1 to 10 and diffusion coefficients from 0.05 to 0.5. The diffusion-free ("noiseless") case of c=0 is numerically unstable. To compute the simulated Burgers' dynamics, we use MATLAB's default PDE solver (pdesolve) on a circular mesh of radius 3 and 100 time steps equally spaced between t=0 and t=1. Then we interpolated the finite-element solution onto a grid of size 25×25 , covering the square $[-1,1]^2$, yielding spatial points that are approximately 0.1 units apart. We restrict our solution to avoid any boundary conditions. Finally, we apply CaStLe-PC and the aforementioned advection angle estimation method, and compare the estimated angle to the true angle. We demonstrate three realizations of this process in the left-hand panel of Figure D1.

Angle Estimation Results

Our results appear in the right panel of Figure D1, where we plot the difference in the true and estimated angle, taking care to account for the "wrapping" behavior of angle-valued data. We see that stronger advection (higher SNR) consistently leads to improved estimation (downward trend within each group), with estimated angles consistently within 10° for advection magnitude 5 or greater. Comparing across different levels of the diffusion coefficient c, we note that higher c increases

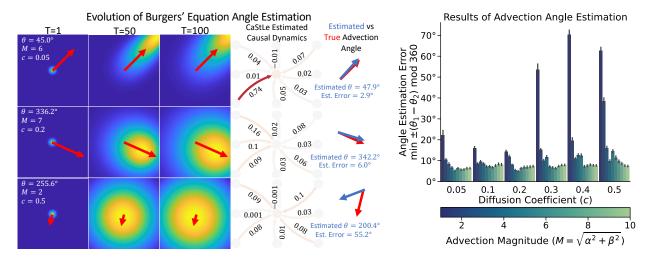


Figure D1: Application of CaStLe-PC to advection estimation from non-linear PDE dynamics. In the left panel, the first three columns depict realizations of Burgers' equation under different advection-to-diffusion regimes; the fourth column depicts the causal stencil identified by CaStLe-PC; and the final column compares the estimated advection angle with the true advection angle. The right panel depicts the accuracy of CaStLe-PC under various signal-to-noise conditions. Each combination of advection and diffusion rates were tested with 500 angles sampled uniformly from $[0^{\circ},360^{\circ})$. In low-diffusion (high SNR) scenarios, CaStLe-PC can identify the underlying advection clearly (top row of left panel and yellow-green columns in right panel). By contrast, in low-advection (low SNR) scenarios, CaStLe-PC struggles to accurately identify the underlying advective dynamics (bottom row of left panel and blue bars in right panel). Even in highly diffusive scenarios, CaStLe-PC is able to accurately estimate the underlying advection when it is sufficiently strong (around $M/c \geq 20$) as shown in the middle row of the left panel. Additional details are given in D.

the angle estimation error, as we would expect in the higher-noise regimes. For low advection magnitude and $c \ge 0.3$, we see an average error approaching the "pure guessing" value of 90° . Even at high diffusion levels (c = 0.5), moderate advection magnitudes of 5-6 are sufficient to ensure accurate estimation. From these, we see that CaStLe-PC is able to consistently recover advection structure across a wide range of SNR regimes. As demonstrated in F, traditional dimension reduction approaches such as PCA and PCA-varimax, when combined with standard causal discovery methods, fail to accurately capture the advection dynamics in Burgers' equation, particularly in identifying the correct advection angle. This highlights

CaStLe's unique ability to preserve and extract meaningful causal structures from nonlinear PDE systems that would otherwise be lost through dimensionality reduction.

The takeaway from these results is that CaStLe can not only generalize to continuous, non-linear models of advection and diffusion, but it can successfully infer the direction of causality in any advective-diffusive system, given that the diffusion is not so large as to dominate advection. Further, each simulation has only one signal surrounded by large areas without data or causal information. Despite this sparsity and the presence of regions where diffusive information flow might suggest incorrect advection angles, CaStLe successfully identifies the correct advection angle when analyzing the full space. CaStLe is asked to learn from the full space, but successfully hones in on the correct advection angle. With these results, we believe CaStLe can be applied to a broad range of space-time systems with advective-diffusive properties to better understand their dynamics.

E Proposed Modification of Statistical Methods for CaStLed Data

While essentially any consistent PIP may be used in Step C, we anticipate that most PIPs will be derived from already existing causal discovery algorithms. Often, these algorithms are statistical in nature and it may be inappropriate to apply them directly to \tilde{X} due to the *seams* connecting each time *chunk*. For a statistical method, which computes a *p*-value for each potential edge (smaller *p*-values leading to

present edges), we suggest the following chunk testing modification:

- 1. For each chunk $b \in \{1, ..., (N-1)^2\}$, let p_b be the p-value resulting from the PIP applied to that chunk.
- 2. Compute $T = -2\sum_{b} \ln p_b$
- 3. Let $p_{\text{agg}} = 1 \chi_{2(N-1)^2}^2(T)$ where $\chi_k^2(x)$ is the cumulative distribution function (CDF) of a χ^2 random variable with k degrees of freedom evaluate at x.
- 4. If $p_{\text{agg}} < p_*$, identify a parent.

This method adapts Fisher's classical method for combining independent p-values to our setting. In practice, however, we have found that for sufficiently large T, this *chunking* is unnecessary as the proportion of *seams* in \tilde{X} goes to zero, and the PIP identifies the correct causal structure despite the small fraction of points of misspecification (1/T).

F Limitations of Dimensionality Reduction for Space-Time Causal Discovery

We demonstrate the limitations of dimensionality reduction approaches such as PCA and PCA-varimax when applied to space-time causal discovery of advective-diffusive processes. Causal discovery methods in Earth science often employ these techniques to reduce the high dimensionality of gridded data before applying causal discovery algorithms. While effective for identifying large-scale telecon-

nections, we show that these approaches fail to capture the local causal structures that are essential for understanding space-time dynamics at the grid-cell level.

To illustrate these limitations, we apply PCA and PCA-varimax dimension reduction followed by PCMCI causal discovery—the procedure described by Runge et al. (2015c), Nowack et al. (2020a), and Tibau et al. (2022) and employed in subsequent work—to each of our case studies: Burgers' equation, HSW-V, and E3SMv2-SPA. Our analysis reveals that while dimensionality reduction techniques can identify dominant modes of variability, they struggle to preserve the spatial relationships between neighboring grid cells, thus obscuring the local causal pathways that CaStLe is specifically designed to recover.

For the PCMCI step, we explored multiple lag values in our experiments and found that the results were consistently unable to capture the directional advection structure regardless of lag parameter choice. This suggests that the limitation is a fundamental constraint of the dimensionality reduction approach. In the results below, we show the simplest case with a maximum lag of 1.

Figure F1 shows the PCA analysis of Burgers' equation, where four EOFs capture approximately 91% of variance but the resulting PCMCI causal graph fails to recover the directional advection process, demonstrating PCA's inability to preserve local causal structures. Figure F2 shows similar limitations with PCA-Varimax applied to the same Burgers' equation data, where despite the rotation enhancing spatial localization of patterns, the causal graph still cannot represent the known directional advection dynamics. Figure F3 illustrates PCA applied to

the HSW-V volcanic aerosol dataset, where four EOFs explain 85% of variance but produce a causal graph that misrepresents the known transport mechanisms. Figure F4 demonstrates that even with varimax rotation, which provides more spatially distinct patterns in the HSW-V dataset, the resulting causal graph cannot capture the directional flow of volcanic aerosols. The EOFs were reordered according to the identified centroids' longitude to improve interpretability. Figure F5 shows the application of PCA to the E3SMv2-SPA climate model data, where nine EOFs account for 87% of variance, yet the PCMCI causal graph fails to detect the underlying atmospheric circulation patterns. Figure F6 reveals that PCA-Varimax rotation of the E3SMv2-SPA data, with EOFs similarly reordered by longitudinal position for interpretability, still fails to recover the known directional transport processes, further confirming the limitations of dimensionality reduction for space-time causal discovery.

PCA Analysis of Burgers' Equation Solution EOF 1 (44.2% of Variance) **EOF 1 Time Series ≻** 50 0 20 40 80 40 . 20 60 80 100 PDE Value EOF 2 (27.2% of Variance) **EOF 2 Time Series** EOF 2 Score **≻** 50 0 40 60 80 20 20 40 60 80 100 PDE Value EOF 3 (13.4% of Variance) **EOF 3 Time Series** EOF 3 Score ≻ 50 0 -40 20 80 0,000 0.025 0.050 20 40 60 80 100 PDE Value EOF 4 (6.7% of Variance) **EOF 4 Time Series** EOF 4 Score ≻ 50 2 0 -60 20 40 80 Χ ,003 60 80 100 Ó 20 40 Time Step PDE Value PCMCI Causal Graph **Explained Variance** Individual Cumulative 80 % of variance 60 20

Figure F1: PCA study of Burgers' equation solution ($\theta = 45^{\circ}$, M = 6, c = 0.05). Four empirical orthogonal functions (EOFs) capture $\approx 91\%$ of 46 griance, with spatial patterns (left) and temporal evolution (right). The bottom panels show explained variance distribution and PCMCI causal graph, which fails to accurately represent the known directional advection process in the underlying PDE highlighting limitations of this approach for local causal structures in space-time systems

EOF

-0.4 0.0 0.4 auto-MCI

PCA-Varimax Analysis of Burgers' Equation Solution EOF 1 (26.6% of Variance) **EOF 1 Time Series** ≻ 50 0 20 40 80 . 20 40 60 80 100 PDE Value EOF 2 (32.1% of Variance) **EOF 2 Time Series** EOF 2 Score 5.0 2.5 ≻ 50 0 -80 20 40 60 0.0 20 40 60 80 100 PDE Value EOF 3 (13.9% of Variance) **EOF 3 Time Series** EOF 3 Score ≻ 50 0 -20 40 80 0030 20 40 60 80 100 PDE Value EOF 4 (19.0% of Variance) **EOF 4 Time Series** ≻ 50 0 -80 20 40 60 003 ò 20 80 100 40 60 000 Time Step PDE Value PCMCI Causal Graph **Explained Variance** Individual -- Cumulative 80 % of variance 60 20 -0.4 0.0 0.4 auto-MCI EOF

Figure F2: PCA-Varimax study of Burgers' equation solution ($\theta = 45^{\circ}$, M = 6, c = 0.05). Four empirical orthogonal functions (EOFs) capture 2491% of variance, with spatial patterns (left) and temporal evolution (right). The bottom panels show explained variance distribution and PCMCI causal graph, which fails to accurately represent the known directional advection process in the underlying PDE, highlighting limitations of this approach for local causal structures in space-time

PCA Analysis of HSW-V

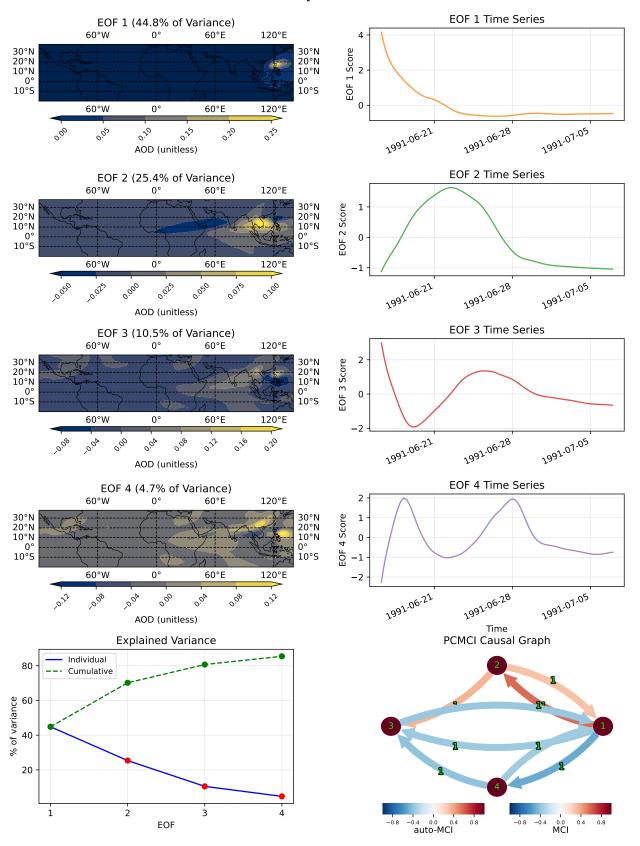


Figure F3: PCA study of the HSW-V dataset, in the time interval 21 days post-eruption. Four empirical orthogonal functions (EOFs) capture 2485% of variance, with spatial patterns (left) and temporal evolution (right). The bottom panels show explained variance distribution and PCMCI causal graph, which fails to accurately represent the known directional advection process in the underlying system, highlighting limitations of this approach for local causal structures in space-

PCA-Varimax Analysis of HSW-V

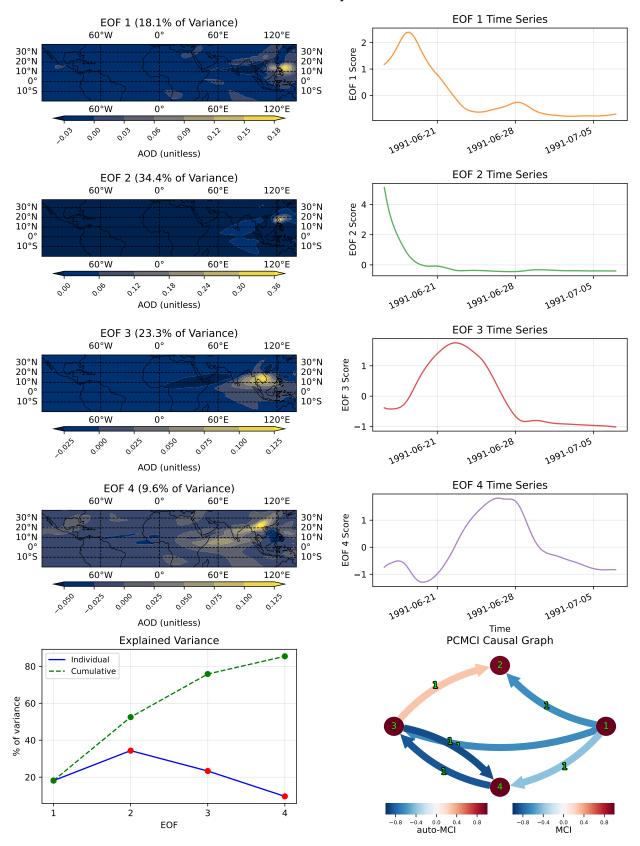


Figure F4: PCA-Varimax study of the HSW-V dataset, in the time interval 21 days post-eruption. Four empirical orthogonal functions (EOFs) capture ≈85% of variance, with spatial patterns (left) and temporal evolution (right). Since varimax rotation does not preserve the explained variance ordering, we reordered EOFs according to the identified centroid's longitude. The bottom panels show explained variance distribution and PCMCI causal graph, which fails to accurately represent

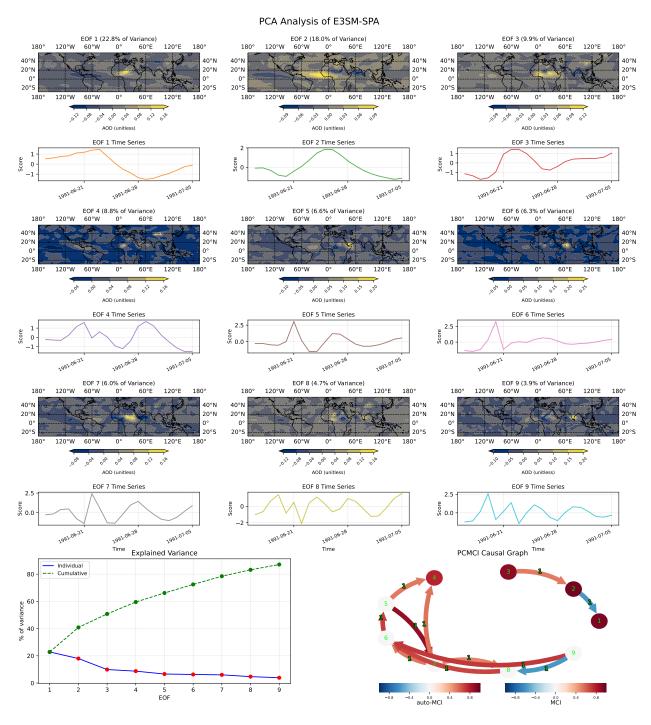


Figure F5: PCA study of the E3SMv2-SPA dataset, in the time interval of days 15-35. Nine empirical orthogonal functions (EOFs) capture $\approx 87\%$ of variance, with spatial patterns (left) and temporal evolution (right). The bottom panels show explained variance distribution and PCMCI causal graph, which fails to accurately represent the known directional advection process in the underlying system, highlighting limitations of this approach for local causal structures in spacetime systems.

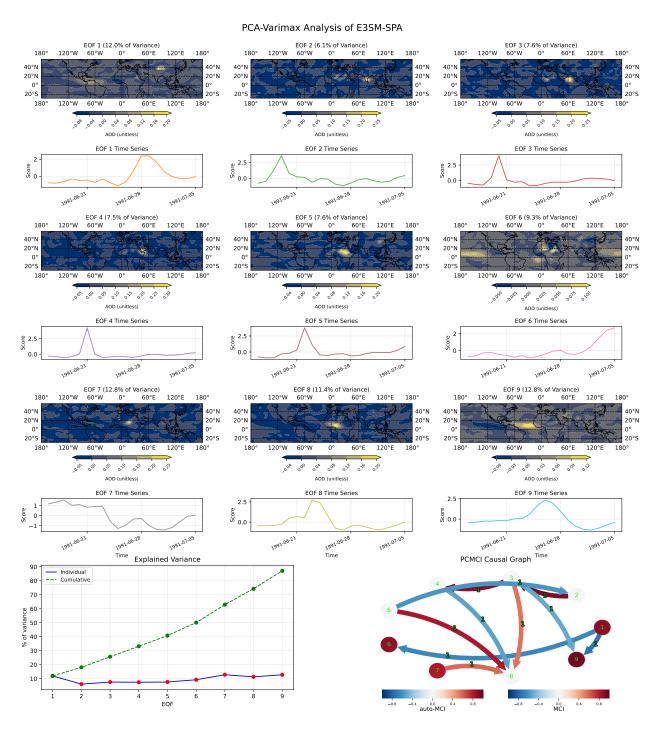


Figure F6: PCA-Varimax study of the E3SMv2-SPA dataset, in the time interval of days 15-35. Nine empirical orthogonal functions (EOFs) capture $\approx 87\%$ of variance, with spatial patterns (left) and temporal evolution (right). Since varimax rotation does not preserve the explained variance ordering, we reordered EOFs according to the identified centroid's longitude. The bottom panels show explained variance distribution and PCMCI causal graph, which fails to accurately represent the known directional advection process in the underlying system, highlighting limitations of this approach for local causal structures in space-time systems.

G Additional experimental details for Section 7.7

CaStLe inherits several of the runtime parameters of the underlying PIP used. In Section 7.7, we set these values at relatively stringent threshold to highlight the most robust and important dynamics and to yield a highly interpretable graph; additional weaker dynamics can be recovered by relaxing these choices at the (potential) cost of additional false positive edges and less interpretability. Data-driven optimization of these parameters is difficult, though the validation strategies suggested by Allen et al. (2023) may be useful here. Specifically, we set a p-value threshold of 1×10^{-5} and removed estimated partial correlations of magnitude less than 0.35; we note here that, due to the adaptive search heuristics used by the PIP, the p-value threshold applied here is not a proper measure of statistical significance, but only a *heuristic* measure of estimated strength. We note that our resulting interpretations are generally quite robust to specific choices of these values.

H Analysis of Spatial Blocking

Here, we briefly investigate two impacts of spatial blocking, of the kind used in Section 7.7. Spatial blocking is a process in which regions of the global space are separated into blocks where CaStLe is applied individually and independently. This can be done for the sake of interpretability and to help ensure the spatial causal structure is uniform and homogeneous in the blocked space, satisfying Assumption

S2.

First, we consider the impact of block size on the HSW-V case study. In our demonstration in Section 7.7.1, we approached block size heuristically, and we chose a relatively large block size to demonstrate correctness saliently. We found that results are generally robust to larger and smaller block sizes in the HSW-V case. In Figure H1, we show that the recovered dynamics in each stencil are generally the same over space for each block size. We see that larger block sizes are easier to interpret at a glance, while smaller sizes describe more nuance. We also found that results were generally robust to block size in the E3SMv2-SPA case.

Second, we consider the impact of a blocking strategy for causal discovery generally by comparing results of the PC algorithm to one block in E3SMv2-SPA to CaStLe-PC's results from the same data. Our comparison of CaStLe and the PC algorithm in Figure 4 make it clear that CaStLe captures the spatial evolution of Mt. Pinatubo's plume much more effectively and about 80,000 times faster. However, one may be concerned that sparsity and correctness could be achieved with blocking alone. In Figure H2a, PC struggles to estimate an interpretable and physically meaningful graph of the dependence structure in this area because of the signal redundancy between nonadjacent grid cells and that there are only 20 observations per grid cell and 25 grid cells. Figure H2b illustrates much better performance from CaStLe, in which CaStLe learns a stencil from the region and projects it back into the original grid space.

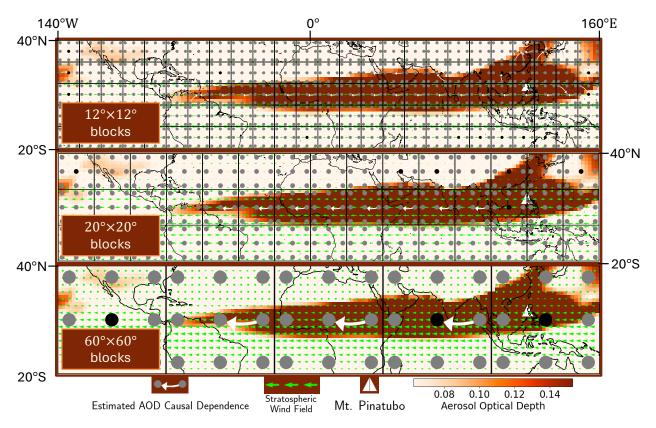
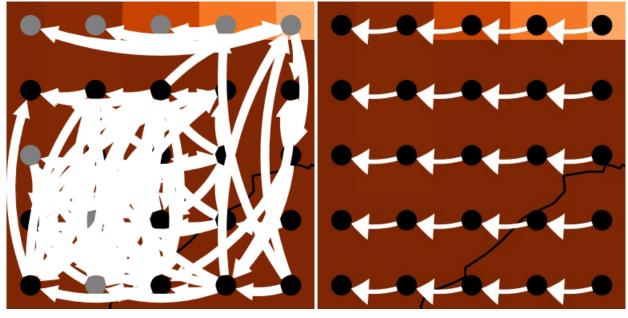


Figure H1: Results of CaStLe applied to HSW-V 21 days after the Mt. Pinatubo eruption with three different block sizes, $12^{\circ} \times 12^{\circ}$, $20^{\circ} \times 20^{\circ}$, and $60^{\circ} \times 60^{\circ}$. We find that results are generally consistent over the same area for each block size, with smaller block sizes allowing for additional nuance in some areas. Note that the $20^{\circ} \times 20^{\circ}$ block panel is similar to the results shown in Figure 3, but more longitudes were added to get a space factorable by more integers, such as 12, 20, and 60.



(a) PC algorithm results

(b) CaStLe results

Figure H2: The PC algorithm and CaStLe applied to E3SMv2-SPA in the $15^{\circ} \times 15^{\circ}$ block between 15.00° to 30.00° N and 75° to 90° E. from the day of the eruption to 20 days later. PC struggles to estimate an interpretable and physically meaningful graph of the dependence structure in this area. In contrast, CaStLe is able to identify an interpretable dependence structure that represents the local dynamics within the space.

I Analysis of Assumption Violation Examples

Here, we evaluate the impacts of potential violations of CaStLe's assumptions in our study of E3SMv2-SPA from Section 7.7.2.

I.1 Time Resolution is Too Coarse (Assumption T1)

The dataset's time resolution can determine if the temporal locality assumption (T1) holds. If the time resolution is too coarse, the temporal causal structures may be marginalized out or unmeasured. Dependencies between neighboring grid cells may not be manifested in the sparse time sampling. Here, we explore how our study of E3SMv2-SPA from Section 7.7.2 changes after coarsening the temporal resolution.

We coarsened the time resolution by two, from a daily to a two-daily resolution.

Figure I1 demonstrates that CaStLe finds much fewer links when the time resolution is too coarse. However, the links that are detected are mostly consistent with known advective processes.

I.2 Time Interval is Too Long (Assumption T2)

When the time interval is too long, there may be too many causal structures in the data. This violates temporal causal stationarity (T2). Here, we investigate such a scenario.

We first computed causal stencils for an extended period, between day 15, the day of the eruption, to day 65. This is 30 days longer than our initial analysis from

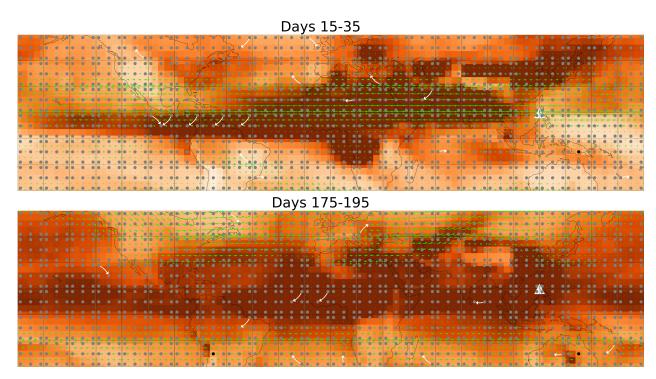


Figure I1: Results of using a coarsened temporal resolution (two-daily) in the E3SMv2-SPA study. CaStLe finds many fewer links in this setting. It is clear that when time is too coarse, causal structures fail to be detected. However, the remaining links that are found are largely true positives, suggesting that CaStLe is relatively robust to coarser time sampling.

the start of the eruption.

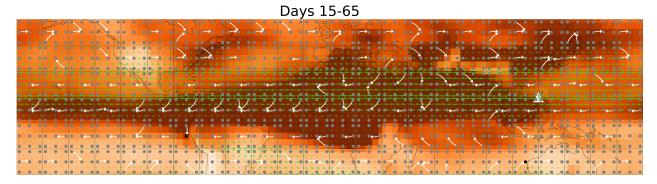


Figure I2: Results of applying CaStLe to a longer time interval from day 15 to 65. CaStLe identifies more links, indicating it is learning too many causal structures in the data, but still finds many of the true positives we found in our initial study. This indicates that many of the blocks in this interval have temporal causal stationarity, leading CaStLe to perform adequately.

We then computed causal stencils for the entire period between day 15 to day 215, roughly six months later.

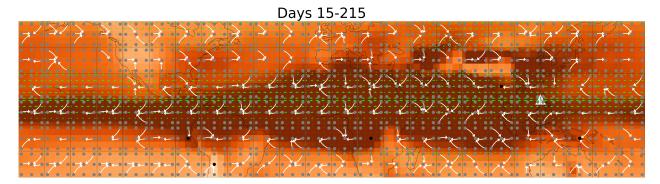


Figure I3: Results of applying CaStLe to a time interval that is too long and contains too many causal structures, day 15 to 200. We see that CaStLe identifies many links in each block. Comparing them to the winds is ineffective because the wind arrows are averages over the whole period rather than reflections of how they change in time, which CaStLe is learning from. With such a density of links, it is further challenging to know which are correct and which are spurious.

Figure I2 shows that when the time interval is longer, CaStLe identifies more links, indicating it is learning too many causal structures in the data, but still finds many of the true positives we found in our initial study. Figure I3 demonstrates the challenges of applying CaStLe to a time interval that contains too many difference

causal structures. CaStLe identifies many links, creating uninterpretable stencils. The winds are a poor comparison because each arrow is a temporal average for that location, which is not representative over the entire interval. CaStLe may be capturing many spurious links or capturing all of the many fluctuating dynamics over the interval. Resulting is are uninterpretable stencils with unknown true and false positives. However, there are some blocks in the equatorial regions with sparse stencils. That indicates that dynamics were relatively stationary over the period.

I.3 Grid Resolution is Too Coarse (Assumption S1)

An appropriate grid resolution is important for satisfying the spatial locality assumption (S1). If the grid is too coarse then the underlying spatial structure may be marginalized out or unmeasured. If it is too small, causal relationships may appear outside the stencil neighborhood, requiring a radius-2 neighborhood implementation. Here, we investigate a grid resolution that is too coarse.

We coarsened the grid to 9° , rather than the 3° used in Section 7.7.2. Given that, to maintain 5×5 grid cells per block, each block is again $45^{\circ} \times 45^{\circ}$.

In Figure I4, we see that CaStLe performs very well overall. There are few false positives and it clearly captures the overall advection dynamics of the system.

We also coarsened the grid to 18° , resulting in $90^{\circ} \times 90^{\circ}$ blocks. In Figure I5, CaStLe performs well in the early time interval, clearly identifying the east-to-west advection pattern. However, in the later time interval, it finds no spatial structures

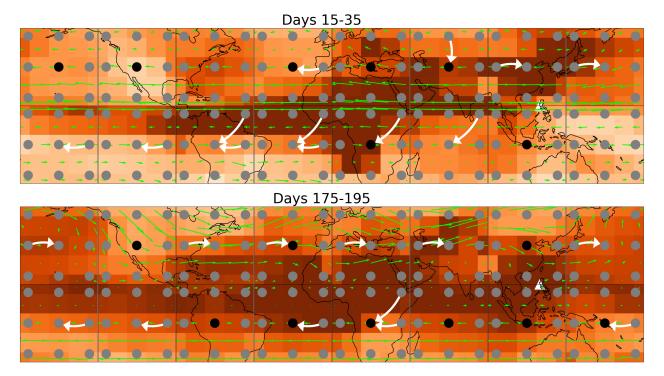


Figure I4: Results of using a coarse grid (9°) in the E3SMv2-SPA study. We find that CaStLe performs very well overall. There are few false positives and it clearly captures the overall advection dynamics of the system.

apart from autodependencies in each block. This is likely because the east-to-west advection is weaker in this period and the grid is too coarse to capture the narrower bands of northward advection that dominates the interval.

We find that CaStLe is very robust to this assumption violation. It captures all of the most dominant advection patterns, while struggling to find smaller, weaker ones.

I.4 Block Sizes are Too Large (Assumption S2)

In H, we found that CaStLe's output was robust to very large and very small block sizes. Spatial blocks are intended to isolate regions such that only one underlying spatial causal structure exists in the block. If the blocks are too large, then As-

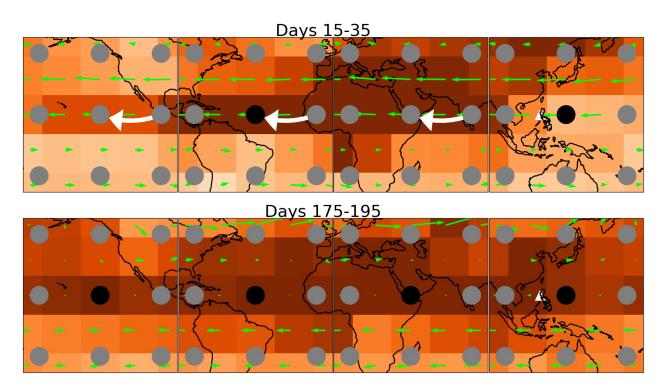


Figure I5: Results of using a coarse grid (18°) in the E3SMv2-SPA study. CaStLe performs well in the early time interval, clearly identifying the east-to-west advection pattern. However, in the later time interval, it finds no spatial structures apart from autodependencies in each block. This is likely because the east-to-west advection is weaker in this period and the grid is too coarse to capture the narrower bands of northward advection that dominates the interval.

sumption S2 may be violated.

In Figure I6, we used block sizes equal to $45^{\circ} \times 45^{\circ}$. Here, each block has 15×15 grid cells. This is in contrast to the 5×5 grid cell, $15^{\circ} \times 15^{\circ}$ blocks used in Section 7.7.2.

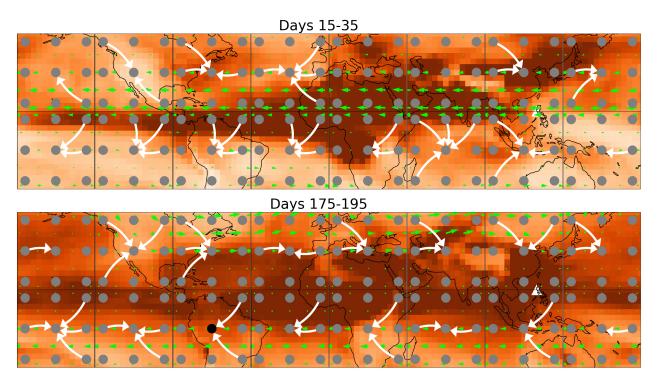


Figure I6: Results of using block sizes too large in the E3SMv2-SPA study. We see that many true positives are found, but many false positives as well. CaStLe seems to identify multiple contradictory causal structures within many cells, which may lead to more spurious links discovered. Even where links appear correct, they are largely uninterpretable in the presence of contradictions.

We find that while true positives remain, several false positives appear. Some positives may be the result of identifying multiple causal structures correctly within the space, while others may be confused results found because of the high density of links. In further testing with intermediate block sizes, we found that CaStLe is moderately robust to this assumption violation. As block sizes approach a more appropriate size, false positives diminish and true positives remain.

J Additional GCM Results

Figure J1 depicts results of implementing CaStLe with the Bayesian score optimization causal discovery algorithm, DYNOTEARS. We also presented results of DYNOTEARS applied to our VAR benchmark in Section 7.8.1. Here, we show that CaStLe-DYNOTEARS is able to recover comparable results to the CaStLe-PC-stable results shown in Section 7.7.1.

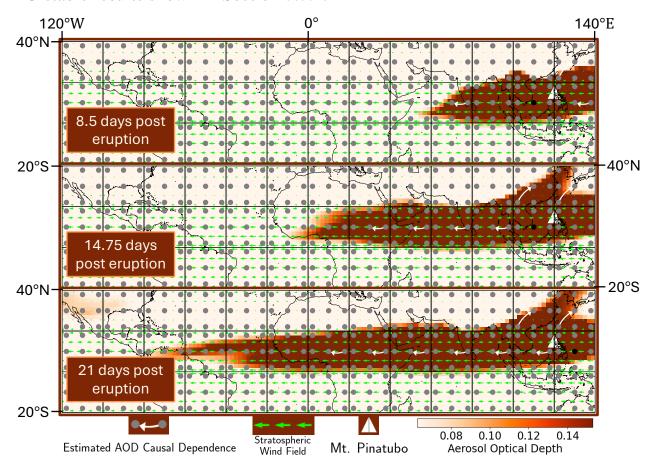


Figure J1: Application of CaStLe-DYNOTEARS to HSW-V simulation of the 1991 Mt. Pinatubo eruption. The stencils estimated by CaStLe (white) capture the underlying high-altitude wind fields (green) using only satellite-measured AOD, with near perfect accuracy in high aerosol regions (redorange). On longer horizons (bottom row), CaStLe is able to recover equatorial wind currents as far away as South America, half-way around the world from Mt. Pinatubo (white triangle). CaStLe accurately identifies the prevailing westerly atmospheric winds because it was able to identify the space-time dependence between neighboring grid cells.

K Additional VAR Results

In Section 7.8.1, we demonstrated the strong performance of CaStLe on VARgenerated space-time data with fixed sparsity level d=4; in particular, CaStLed variants uniformly improve over the performance of equivalent unstructured causal discovery algorithms. We repeat this analysis for a variety of sparsity levels in Figures K1 and K2 for the MCC and F_1 score similarity metrics, respectively. As in Figure 7.6, the CaStLed variants continue to significantly outperform across all sparsity levels, d; furthermore, as noted above, we observe that CaStLe can correctly estimate the underlying grid even on as few as T=10 time samples when a sufficiently large grid is observed; non-CaStLe methods struggle on larger grid sizes, consistent with our analyses in the previous section. A time limit of 48 hours of wall-clock time was applied for each individual graph estimation: performance properties of methods that did not terminate during this window are not shown (e.g., DYNOTEARS) with d=6; T=10; N=10).

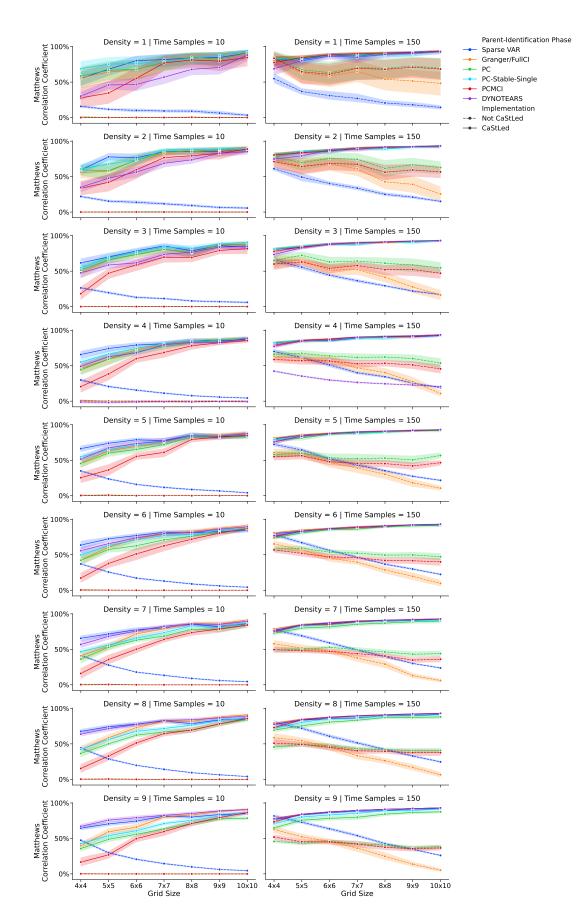


Figure K1: Matthews correlation coefficient (MCC) comparison between CaStLed and non-CaStLed causal discovery approaches on 2D VAR dynamics for each sparsity level, including Granger causality (orange), PC (green), PC-Stable-Single (cyan), PCMCI (red), DYNOTEARS (purple), and a statistical model of the data generating process (blue). See Section 7.8.1 for experi-

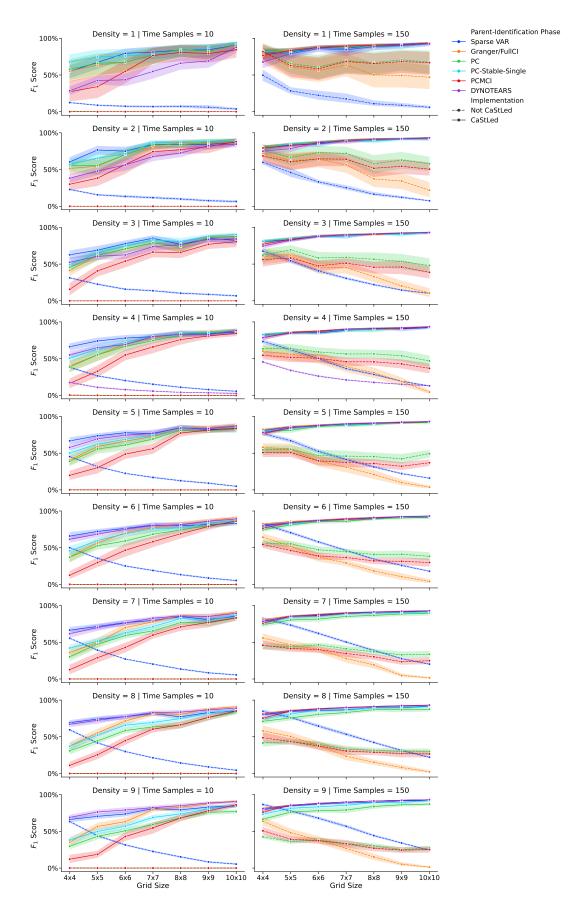


Figure K2: F_1 score comparison between CaStLed and non-CaStLed causal discovery approaches on 2D VAR dynamics for each sparsity level, meluding Granger causality (orange), PC (green), PC-Stable-Single (cyan), PCMCI (red), DYNOTEARS (purple), and a statistical model of the data generating process (blue). See Section 7.8.1 for experimental details.

L PC-Stable-Single

For the convenience of the reader, we include pseudo-code for the PC-Stable-Single algorithm of Runge et al. (2019a), itself an adaptation of the PC-Stable algorithm of Colombo and Maathuis (2014). We use this as the PIP used for the CaStLe-based analyses shown in Sections 7.7.1, 7.7.2, and D. As our experiments in the proceeding section show, PC-Stable-Single exhibits small, but consistent improvements over alternative PIP choices.

Open Research Section

The data generated and used for our HSW-V, VAR, and PDE experiments in Sections 7.7.1, 7.8.1, and D are available on Zenodo via https://doi.org/10.5281/zenodo.12701546 with GNU Lesser General Public License v3.0 or later (Nichol, 2024). The data used for the E3SMv2-SPA experiments in Section 7.7.2 can be found in Brown et al. (2024). The code for generating data, running experiments, and generating figures can be found here https://github.com/jjakenichol/CaStLe.

Acknowledgments

We thank Kara Peterson, the Deputy Principal Investigator of the CLDERA (CLimate impact: Determining Etiology thRough pAthways) project at Sandia National Laboratories (SNL), for helping to make this work possible. We also thank Joey

Algorithm 2 PC-stable-single

Precondition: Time series dataset $\mathbf{X} = \{X^1, X^2, ..., X^N\}$, selected variable X^j , maximum time lag τ_{max} (default $\tau_{max} = 1$), significance threshold α_{PC} , maximum condition dimension p_{max} (default $p_{max} = N_{\tau_{max}}$), maximum number of combinations q_{max} (default $q_{max} = 1$), conditional independence test function *I*. 1: **function** $CI(X, Y, \mathbf{Z})$ Test $X \perp \!\!\!\perp Y | \mathbf{Z}$ using test statistic measure I 2: **return** p-value, test statistic value I 3: 4: Initialize set of parents $\widehat{\mathscr{P}}(X_t^j) = \{X_{t-\tau}^i : i \in \{1,...,N\}, \tau \in \{1,...,\tau_{max}\}\}$ 5: Initialize dictionary of test statistic values $I^{min}(X_{t-\tau}^i \to X_t^i) = \infty \ \forall X_{t-\tau}^i \in \widehat{\mathscr{P}}(X_t^j)$ 6: **for** $p = 0, ..., p_{max}$ **do**

```
if |\mathscr{P}(X_t^J)| - 1 < p then
 7:
 8:
                Break for-loop
                                                                                                    for all X_{t-\tau}^i in \widehat{\mathscr{P}}(X_t^j) do
 9:
10:
                for all lexicographically chosen subsets \mathscr{S} \subseteq \widehat{\mathscr{P}}(X_t^i) \setminus \{X_{t-\tau}^i\}, with |\mathscr{S}| = p do
11:
                     q = q + 1
12:
                     if q >= q_{max} then
13:
                           Break from inner for-loop
14:
                     Run CI test to obtain (p-value, I) \leftarrow CI(X_{t-\tau}^i, X_t^i, \mathscr{S})
15:
                     if |I| < I^{min}(X_{t-\tau}^i \to X_t^i) then
                                                                                  ⊳ Store min. I of parent among all tests
16:
                          I^{min}(X_{t-\tau}^i \to X_t^i) = I
17:
                     if p-value > \alpha_{PC} then
                                                                       \triangleright Removed only after all X_{t-\tau}^i have been tested
18:
                          Mark X_{t-\tau}^i for removal from \widehat{\mathscr{P}}(X_t^i)
19:
                          Break from inner loop
20:
          Remove non-significant parents from \widehat{\mathscr{P}}(X_t^i)
21:
          Sort parents in \widehat{\mathscr{P}}(X_t^i) by I^{min}(X_{t-\tau}^i \to X_t^i) from largest to smallest
22:
```

23: **return** $\widehat{\mathscr{P}}(X_t^i)$

Hart at SNL for helping with 2D Burgers' equation modeling and Tom Ehrmann at SNL for his help in understanding the Atmospheric dynamics we sought to capture. Finally, we thank everyone on CLDERA's simulation team, especially Benj Wagman, Hunter Brown, and Joe Hollowed, for developing the E3SMv2-SPA and HSW-V models, preparing the data, and providing their expertise.

This work was supported by the Laboratory Directed Research and Development program at Sandia National Laboratories, a multi-mission laboratory managed and operated by National Technology & Engineering Solutions of Sandia, LLC (NTESS), a wholly owned subsidiary of Honeywell International Inc., for the U.S. Department of Energy's National Nuclear Security Administration (DOE/NNSA) under contract DE-NA0003525. This written work is authored by employees of NTESS. The employees, not NTESS, own the right, title, and interest in and to the written work and is responsible for its contents. Any subjective views or opinions that might be expressed in the written work do not necessarily represent the views of the U.S. Government. The publisher acknowledges that the U.S. Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this written work or allow others to do so, for U.S. Government purposes. The DOE will provide public access to results of federally sponsored research in accordance with the DOE Public Access Plan.

This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

8 M-CaStLe: Uncovering Local Causal Structures in Multivariate Space-Time Gridded Data

8.1 Publication Notes

Citation: Nichol, J. Jake, et al. "Space-Time Causal Discovery in Earth System

Science: A Local Stencil Learning Approach." Unsubmitted.

Publication date: N/A

Conference: N/A

Formatting: The original text has been preserved as much as possible while still adhering to the formatting requirements of this dissertation.

Data and Software Availability: The paper is currently being prepared for submission and is not yet publicly available.

8.2 Abstract

Causal discovery tools propose to solve one of science's most important and challenging problems, the identification of underlying structure from observed phenomena. Many systems prohibit the feasible or ethical application of more robust methods, such as randomized control trials. In particular, space-time systems, such as the Earth system or ecological systems, are attractive for causal discovery because they could suffer costly alterations if they are manipulated haphazardly. However, space-time systems are challenging to evaluate because their discretized representation as gridded space-time data is often very high-dimensional—possessing many more grid cells than temporal observations. The CaStLe meta-algorithm introduced by Nichol et al. (2024) proposed to solve that problem in scenarios satisfying their assumptions. However, it is limited to univariate analysis, identifying the space-time structure underlying a single quantity.

In this work, we present Multivariate Causal Space-Time Stencil Learning (M-CaStLe), a multivariate extension to CaStLe. We adapt the two phases of CaStLe to first collect the multiple variables in the repeating local neighborhood information in space-time gridded data, and second evaluate the causal parents of variables in the local neighborhood structure. M-CaStLe produces a multivariate causal stencil graph, which extends the CaStLe stencil to represent each variable at each location of the Moore neighborhood. We've added a decomposition method for interpreting the multivariate stencil in terms of just spatial dynamics or just inter-variable dynamics with the spatial graph and reaction graph, respectively. To evaluate M-CaStLe, we developed a multivariate space-time vector autoregression model (VAR) benchmark methodology. The multivariate space-time VARs provide data generation and ground-truth causal stencils for direct evaluation of M-CaStLe.

Our experiments demonstrate that M-CaStLe achieves high precision across

varying numbers of variables and grid sizes, indicating reliable identification of true positive links. However, recall decreases with an increasing number of variables, suggesting more complex systems have more challenging signals to identify. Further analysis shows that recall improves with stronger signal strengths, even in systems with up to 200 variables, indicating good performance in very high variable regimes. Comparisons with the PC algorithm reveal that M-CaStLe-PC consistently outperforms PC in high-dimensional settings, highlighting M-CaStLe's robustness in complex multivariate systems.

8.3 Introduction

Causal discovery is a set of causal inference tools for estimating the underlying structure in observed phenomena. While optimal causal estimation requires randomization, in many settings it is infeasible or unethical to apply (Runge et al., 2019b; Glymour et al., 2019). Thus, causal discovery for space-time systems is critical for scientific inquiry of complex emergent phenomena in physical systems because they often present challenges for randomization. For example, we have one Earth and randomly intervening in its systems is both prohibitively expensive and unethical due to unknown downstream effects. Likewise, neuroscience and ecology are prohibitive to random intervention.

Since the advent of Granger causality (Granger, 1969), the Rubin causal model (Rubin, 2019), causal graphs (Pearl et al., 2016), and the PC algorithm (Spirtes et al., 1993) (named for its authors, Peter and Clark), causal inference and causal

discovery of observed data have developed into a rigorous mathematical framework. Today, causal discovery has become a rich literature with many algorithms and applications throughout the sciences (Glymour et al., 2019; Runge et al., 2023), including the health, Earth, and social sciences (Ebert-Uphoff and Deng, 2012; Cooper et al., 2015; Runge et al., 2019b; Nowack et al., 2020a; Feder et al., 2022; Zanga et al., 2022; Sadeghi et al., 2023). Finally, causal representation learning is an exciting nascent field is developing that merges the flexibility and predictive power of machine learning with causal discovery techniques (Schölkopf et al., 2021).

This work presents a causal discovery approach for space-time systems with gridded data. Unlike space-time systems with point data, such as city-level data, gridded datasets generally enable the analysis of continuous effects over space, since they are regular and complete throughout the grid. However, such systems come with dimensionality challenges. Frequently, the number of grid cells scales faster than the number of temporal samples per grid cell (Runge et al., 2019b). Further challenging their analysis, such systems usually have multiple interacting variables per grid cell that are of scientific interest.

For example, in the Earth system, several interacting quantities may be measured over tens of thousands of grid cells, with hundreds of observations per variable in each grid cell. Atmospheric data often contains hundreds of thousands of grid cells, each with several orders of magnitude fewer observations in time. That imbalance is one aspect of the *curse of dimensionality* (Bellman, 1957; Bühlmann

and Geer, 2011), where high dimensionality relative to sample size challenges conventional statistical methods and renders many forms of inference, including causal discovery, unreliable without dimensionality reduction.

Dimensionality reduction, such as principal component analysis (PCA) (Greenacre et al., 2022; Weylandt and Swiler, 2024), marginalizes large regions of grid cells into several one-dimensional time series. Each time series is then used for individual variables in the chosen causal discovery algorithm (Runge et al., 2015c). This procedure is effective for identifying large-scale patterns such as climate teleconnections (Tibau et al., 2022), but eliminates local grid-level interactions by construction. While large-scale patterns are important aspects of study in complex systems, the nature of their emergence is also important to understand. Local interactions determine the location and magnitude of larger patterns and other midscale phenomena, such as weather and seasonal patterns in atmospheric sciences.

Nichol et al. (2024) developed Causal Space-Time Stencil Learning (CaStLe), which is capable of grid-level causal discovery of high-dimensional space-time data. CaStLe can efficiently identify local causal relationships of a given quantity in space-time systems where traditional approaches fail. However, many scientific questions in complex space-time systems require analysis of multiple quantities per grid cell, such as temperature and soil moisture in Earth system monitoring of drought conditions (Sun et al., 2021) or infection dynamics in epidemiological modeling using infection severity, duration of infection, and population age (Ganesan and Subramani, 2021; Paul et al., 2021).

In this work, we propose an extension to the CaStLe meta-algorithm enabling multivariate space-time causal discovery of grid-level data. We show that Multivariate Causal Space-Time Stencil Learning (M-CaStLe) can effectively capture the causal relationships in multivariate space-time systems. Our results demonstrate that M-CaStLe is capable of accurately estimating local multivariate space-time structures from gridded data, outperforming the PC algorithm, especially in high-dimensional settings. This suggests that M-CaStLe is a robust tool for causal discovery in complex multivariate systems, providing valuable insights into the underlying dynamics of such systems.

8.3.1 Background and Motivation

CaStLe is a meta-algorithm for causal discovery in high-dimensional space-time systems. By leveraging local causal regularities, CaStLe transforms the causal discovery problem from a high-dimensional space with many variables and limited observations to a low-dimensional embedding with fewer variables and more abundant observations. This transformation enhances the efficiency and accuracy of causal discovery, facilitating the identification of causal relationships in their natural context. The present work extends of CaStLe, aiming to broaden its applicability to multivariate space-time dynamics, making it a versatile tool for analyzing various space-time systems in the physical sciences.

8.3.2 Foundations of the CaStLe Framework

In many natural and engineered systems, complex global behaviors emerge from simple local interactions that follow consistent physical dynamics. Nichol et al. (2024) called such systems *partial differential equation (PDE)-like* because they exhibit consistent dynamics defined by interactions between adjacent points in space, with smooth transitions between dynamical boundaries and equilibria. These are characterized by a set of fundamental assumptions that constrain their dynamics:

- **T1**) Temporal Locality: for any $\tau \neq 1$, $X_{i,t-\tau} \not\to X_{j,t}$ for any spatial coordinates (i,j)
- **T2)** Temporal Causal Stationarity: the dynamics governing the evolution of X_t do not change over time. That is, $X_{i,t-1} \to X_{j,t} \Leftrightarrow X_{i,t-1+\tau} \to X_{j,t+\tau}$ for any time offset τ .
- **S1**) Spatial Locality: if (i, j) are not neighbors (in a problem-specific sense) then $X_{i,t_1} \not\to X_{j,t_2}$ for any t_1, t_2 .
- **S2**) Spatial Causal Stationarity: the dynamics governing the evolution of X_t do not change over space. That is, $X_{i,t-1} \to X_{j,t} \Leftrightarrow X_{i+s,t-1} \to X_{j+s,t}$ for any spatial offset s.

Here, $\not\rightarrow$ denotes the absence of a direct causal relationship between two variables. Nichol et al. (2024, Appendix A) describes these assumptions in detail, including ways they may be violated and their Appendix I demonstrates some examples of their violations and CaStLe's robustness to them. To apply causal discovery, the causal assumptions *the causal Markov condition* and *faithfulness* (Spirtes et al., 1993) must be additionally assumed. Because of the locality assumptions, the commonly required *causal sufficiency* assumption may be relaxed (Nichol et al., 2024).

Such systems exhibit both temporal and spatial locality. Temporal locality (T1) dictates that state transitions depend only on the immediate past, preventing "backward causation" and respecting the arrow of time. Spatial locality (S1) ensures that interactions occur only between proximate elements, eliminating action at a distance.

The governing dynamics in these systems demonstrate invariance across both time and space. Temporal causal stationarity (T2) means the rules of evolution remain constant throughout the analysis period—the same causes produce the same effects regardless of when they occur. Spatial causal stationarity (S2) implies that these rules apply uniformly across the domain—the physical location of an element does not alter how it responds to its neighbors. While many macro-scale spaces contain multiple sets of equilibrium dynamics, there are typically micro-scale regions containing stationary spatial causality.

These systems can be represented through structural causal models (SCM) of the form:

$$X_{i,t} = f_i(X_{\mathcal{N}(i),t-1}, \eta_{i,t})$$
(8.1)

Where $X_{\mathcal{N}(i),t-1}$ represents the states of elements in the neighborhood of i at the previous time step, and $\eta_{i,t}$ captures stochastic innovations. Under spatial causal stationarity, the functional form f_i is identical for all i, reducing to a single function f that applies throughout the domain. In short, this space-time SCM implies grid cells exhibit Granger-causal dynamics, which imply that each grid cell's temporal information content encodes the past-history of itself and its immediate neighbors.

This framework encompasses numerous well-studied systems including those governed by partial differential equations, cellular automata, and various lattice models in statistical physics. The approach provides a powerful foundation for both forward simulation and inverse problems—identifying the underlying causal structure from observed spatiotemporal data.

CaStLe not only seeks to identify local causal dynamics but also to do so for high-dimensional systems. In some cases, it may be enough to apply causal discovery independently to small groups of local grid cells; however, in many systems of study, more grid cells are present than observations within each. To accomplish discovery in this regime, we need to efficiently use all the dynamical information in a system.

CaStLe leverages the inherent locality and stationarity to collect time series representing the space-time replicates in such systems. Every grid cell's time series encodes the causal influence of its neighbors, and they can be used as informative

replicates of the system's local dynamics. CaStLe processes a set of grid cells, collecting each one's data on its local dependence, then learns the causal structure of the grid cells and their neighborhoods.

CaStLe's first phase is to form the Locally Encoded Neighborhood Structure (LENS), an embedding representing the Moore neighborhood–a 3×3 matrix of a grid cell and its eight immediate neighbors. The LENS contains concatenated time series from each grid cell's Moore neighborhood so that the local dynamics from each neighbor is repeated. The embedding is a 3×3 matrix, with each entry representing the North West, North, North East, West, center, East, South West, South, and South East grid positions of the Moore neighborhood. Each entry of the embedding contains long concatenated time series collected from throughout the original grid space. Each time series is of length $T\times(N-2)^2$, for the grid's dimension N and T time samples per grid cell. The embedding does not marginalize any data, so no information loss occurs, as would happen during other dimensionality reduction techniques. Figure 8.1 is a conceptual diagram depicting using the local Moore neighborhood to construct the LENS.

Once the embedding is constructed, CaStLe's second phase, the Parent-Identification Phase (PIP) applies an adapted causal discovery algorithm to the embedding. Any time series causal discovery algorithm may be adapted by requiring it to treat the embedding's center grid cell as special: it may be the only child in the resulting causal graph; parents are unrestricted. This adaptation has multiple effects: it creates a graph of the generalized ancestry for each grid cell, eliminates would-be

Locally Encoded Neighborhood Structure

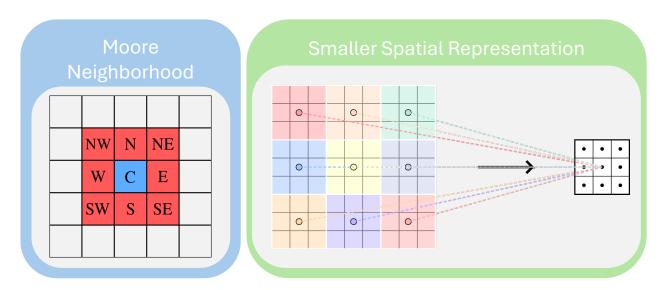


Figure 8.1: A conceptual diagram of the LENS that CaStLe constructs for learning underlying local causal dynamics in gridded data. This encoding transforms the original grid space into a local neighborhood structure without marginalization, preserving all of the local relationships in the gridded time series data.

unobserved confounding between the embedding's outer grid cells and their neighbors beyond the embedding, and increases computational and statistical efficiency, which is detailed below. The result of the PIP on the embedding is the *causal sten-cil graph*, a representation of the local causal dynamics between all grid cells in the system.

8.3.3 Theoretical Properties and Empirical Validation of CaStLe

Nichol et al. (2024) showed that CaStLe exhibits significant performance and efficiency improvements for grid-level causal discovery. It successfully reconstructed known volcanic aerosol dynamics, driven by stratospheric winds, in the weeks after the Mount Pinatubo eruption of 1991. We demonstrated its general performance

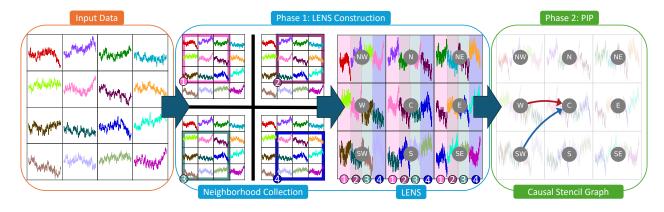


Figure 8.2: A demonstration of the full CaStLe process to produce a causal stencil graph on an example input 4×4 gridded space-time system. In the LENS phase, neighborhood information is collected from each of the interior grid cells, which are then concatenated to form the LENS. Finally, the PIP phase applies an adapted time series causal discovery algorithm to learn the space-time parents of the center node. The learned stencil depicts the underlying space-time structure of each grid cell in the original data.

on advective-diffusive dynamics with a Burgers' equation simulation study. We compared it to existing causal discovery algorithms with ground-truth defined by space-time vector autoregression model (VAR) models.

Because CaStLe constructs the LENS, a lower spatial-dimension embedding, and the PIP limits potential causal children to only the center node, the number of variables and possible links are both fixed to nine. That property enables much more efficient causal discovery. Computational complexity is a measurement of the asymptotic bounding on how many computational resources are required for increasingly large input sizes. The PC algorithm has a computational complexity bounded by $\mathcal{O}(Tp^32^p)$, when applied to an p grid cells, with T time samples per cell. We showed that CaStLe is bounded by $\mathcal{O}(Tp)$.

CaStLe also exhibits improved sample complexity, which measures the asymptotic bounds on how many samples are required to ensure correct graph estima-

tion. The probability of the PC algorithm incorrectly estimating the true graph is bounded by $\approx \mathcal{O}(p^p)$. In contrast, we find that CaStLe's error probability scales as $\approx \mathcal{O}\left(\frac{pT}{e^{pT}}\right)$. From this, as the grid size grows larger, we find that PC is less likely to estimate the correct causal graph, while CaStLe is more likely to estimate the correct graph.

Nichol et al. (2024) also demonstrated several empirical results of CaStLe with benchmarks and realistic climate model output studies. It was shown that CaStLe can robustly capture the transport patterns of volcanic aerosols emitted by the 1991 Mount Pinatubo eruption. It outperformed the PC algorithm in terms of accuracy and execution time, largely because PC naively sought causal relationships between all grid cells without the benefits of the LENS. CaStLe was also robust to moderate assumption violations. The VAR benchmark study compared CaStLe to popular time series causal discovery methods, including the PC algorithm (Spirtes and Glymour, 1991), PCMCI (Runge et al., 2019a), and DYNOTEARS (Pamfil et al., 2020). They found that CaStLe variants performed well, with better results on larger grids, while non-CaStLe algorithms struggled and performed more poorly on larger grids. The Burgers' equation study evaluated CaStLe's performance in different advection speed and diffusivity regimes via advection-diffusion partial differential equation (PDE) model output. CaStLe performed well except in settings where diffusion dominated, making advection signals unrecoverable.

8.3.4 Research Gap and Motivation for Multivariate Extension

Nichol et al. (2024) showed that CaStLe can reconstruct the local space-time causal structure between grid cells of one quantity, e.g., atmospheric aerosols. While helpful in understanding the underlying dynamics of a species transporting or propagating in a complex environment, it leaves learning impacts of that transport to later inference and analysis. Such a manual or post-hoc multivariate inference becomes complex as the number of variables increases.

For example, in Nichol et al. (2024), CaStLe identified the space-time evolution of volcanic aerosols in the stratosphere from the Mt. Pinatubo eruption. Given the rich literature of that eruption, we know that the volcano's SO₂ output increased stratospheric temperatures and decreased tropospheric temperatures for two-to-five years (Dutton and Christy, 1992; Labitzke and McCormick, 1992; Parker et al., 1996a; Soden et al., 2002). The eruption's SO₂ did not directly impact temperature, the plume of gas underwent chemical and physical evolutions, forming H₂SO₄ and advecting and diffusing around the globe. However, univariate CaStLe needs to analyze each chemical species separately and cannot determine interactions between species.

To estimate the space-time dynamics of each variable separately and then infer variable interactions afterward potentially introduces errors and does not have the benefit of joint estimation, which is available in time series causal discovery, such as PCMCI (Runge et al., 2019a). Furthermore, learning space-time causal structures from each variable independently may miss cross-variable confounding, leading to space-time estimation errors and incorrect inference of the underlying physical process.

Joint estimation of space-time dynamics and variable interactions can enable more complex analyses. For example, SO₂ follows a chemical and physical causal pathway to mediate temperature. SO₂ reacts with water molecules to become H₂SO₄. Finally, H₂SO₄ interacts with incoming solar radiation, which impacts temperatures. Understanding the local space-time dynamics of these aerosol species as they transport around the globe may help explain local temperature impacts. Domains outside of atmospheric chemistry and Earth systems science where estimating grid-level multivariate interactions in space-time systems (MacEachren et al., 1999; Haas, 2002) would be valuable are computational fluid dynamics (Wimer et al., 2023), spatiotemporal pharmacokinetics (Guarin et al., 2021; Klingelhuber et al., 2024), and computational chemistry (Higham, 2008; Owen et al., 2024).

Multivariate interactions are challenging to estimate at the grid-level because the high-dimensionality of datasets observed from space-time systems becomes more challenging with more variables because each variable entails p more grid cells to estimate for the same T observations per grid cell per variable. CaStLe solves the high-dimensional challenge in many univariate space-time systems. Extending its capabilities to discover variable interactions simultaneously with space-time dynamics for each variable enables robust discovery of how they interact in space and time. Doing so while maintaining the interpretability of the graphs at

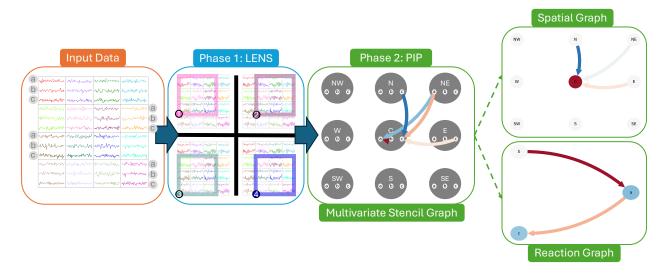


Figure 8.3: A schematic diagram of the input, computational phases, and output of M-CaStLe. Similar to CaStLe's procedure (c.f. Figure 8.2), the first phase collects local neighborhood information into the LENS, which now collects information for each variable's time series in each grid cell. The second phase applies the PIP to every variable at every position in the LENS to determine which variables cause the center variables from each location in the LENS. Finally, the resulting multivariate stencil graph can be decomposed into the *spatial graph* and *reaction graph* for improved interpretability and potential analysis.

scale is also challenging. Multivariate stencil graphs need to contain many more nodes for each variable and still describe local structures.

8.3.5 Contributions

M-CaStLe solves these challenges by adapting both phases of the original CaStLe meta-algorithm. The first phase, which restructures the given gridded data into the LENS, is adapted to restructure multivariate data to preserve space-time and intervariable relationships. The univariate PIP sought causal relationships terminating in only one node (the center). The multivariate PIP is adapted to find parents and children of multiple sets of nodes for each variable.

These advances enable the simultaneous estimation of space-time dynamics and

inter-variable interactions, providing a more comprehensive understanding of complex systems. This capability is particularly valuable in fields such as atmospheric science, where understanding the interplay between different chemical species and their impact on climate is crucial. We validated M-CaStLe through extensive experiments on synthetic benchmarks. Our results demonstrate that M-CaStLe outperforms existing methods in terms of accuracy and computational efficiency, particularly in high-dimensional settings. The empirical validation shows that M-CaStLe can robustly capture the causal structure of multivariate space-time systems, making it a powerful tool for scientific discovery and analysis.

8.3.6 Paper Organization

The remainder of this paper is organized as follows: in Section 8.4 we introduce M-CaStLe, our multivariate extension to CaStLe; Section 8.5 discusses our benchmark's experimental setup with VARs; Section 8.6 presents a rigorous analysis of the multivariate results of M-CaStLe benchmarked on multivariate models of space-time dynamics; and finally we discuss the presented work and future directions in Section 8.7.

8.4 Methods

Multivariate CaStLe (M-CaStLe) extends CaStLe's capabilities to discover local space-time causal structures in multivariate data. M-CaStLe produces the multivariate causal space-time stencil graph, which describes how a *set of variables*

interact within their Moore neighborhood over time. The multivariate stencil is often challenging to interpret immediately. To improve interpretability, we present the multivariate stencil the *reaction graph* and the *spatial graph*, which decompose the multivariate stencil output by M-CaStLe into a graph of inter-variable relationships (without a spatial aspect) and a graph of spatial relationships (without variable relationships).

M-CaStLe adapts both phases of the CaStLe meta-algorithm to enable construction of a LENS containing multiple variables and successful causal discovery of space-time and inter-variable dependencies within the LENS. Input data consists of V variables measured on an $N \times N$ grid over T time steps, yielding a tensor $\mathbf{X} \in \mathbb{R}^{N \times N \times V \times T}$. Figure 8.3 depicts each step of M-CaStLe. In this example, we illustrate a simple 4×4 original grid space, G, which has V = 3 locally interacting variables, a, b, and c, with T = 500 time samples.

8.4.1 Phase 1: The Locally Encoded Neighborhood Structure (LENS)

Phase 1 collects neighborhoods in the same fashion as the univariate CaStLe, but it now collects multiple time series per spatial location in the Moore neighborhood for each variable. The univariate LENS is a 3×3 matrix where each element contains one time series of length $T\times (N-2)^2$. Since M-CaStLe has V variables, the multivariate LENS is a 3×3 matrix where each element contains V time series of length $T\times (N-2)^2$. In short, it is a tensor in $\mathbb{R}^{3\times3\times V\times L}$, where $L=T\times (N-2)^2$ is the length of each concatenated time series. In Figure 8.3, Phase 1 depicts the pro-

cess the LENS construction follows to collect time series from each Moore neighborhood as the window slides across *G*. It collects all three variables from each grid cell within the neighborhood window and concatenates them to the LENS, according to their position relative to the center of the neighborhood window and the respective variable in each position. Like the univariate LENS, there is no marginalization or loss of data, and its structure allows it to be fully invertible. We do not have a reason to invert the procedure in this analysis, but it illustrates that no information loss occurs.

8.4.2 Phase 2: The Parent-Identification Phase (PIP)

In univariate CaStLe, the PIP adapts a given time series causal discovery algorithm, such as DYNOTEARS (Pamfil et al., 2020), to seek the parents of only the center node in the LENS. To adapt this approach to M-CaStLe, we do the same for each variable in the center node. Rather than allowing one child in the discovery process, we now allow V children. This has the effect of every variable in every position in the LENS having a potential causal effect on every variable in the center position. Resulting is a multivariate stencil, such as the one depicted in the third panel of Figure 8.3. This example illustrates a stencil of three variables with dependencies between each over space and time.

8.4.3 Interpretability: Decomposing the Multivariate Stencil

While the stencil in Figure 8.3 may be interpretable after careful viewing, multivariate stencils of more variables or with more dependencies can be challenging to parse visually. For that reason, we have developed a decomposition scheme to analyze the variable interactions and the spatial structure of all variables separately. The far right of Figure 8.3 illustrates the spatial graph and reaction graph corresponding to the stencil to their left.

Computing the stencil decomposition is straightforward and similar for both the spatial and reaction graphs. To compute the spatial graph, the stencil links are aggregated along the variable dimension, and the location from which they originate is preserved. For example, in Figure 8.3, two links are coming from the NE position to the center, a negative dependence (light blue) via $a \rightarrow a$ and a positive dependence (orange) via $a \rightarrow c$, and both of those are aggregated to find one weakly negative link NE \rightarrow C in the spatial graph. Note that there is a $b_{center} \rightarrow a_{center}$ link in the stencil and that it is represented as an autodependence link in the spatial graph, illustrated by the center node's coloring. The node and link colors directly associate with continuous link dependence strength that is output by M-CaStLe, but we omit that detail for the example in Figure 8.3.

The reaction graph is computed by aggregating stencil links along the spatial dimension while preserving the variable dimension. For example, in Figure 8.3, there are two links $c \to c$ in both the N and E locations, strongly negative (blue)

from the N and weakly positive (red) from the E. Those are aggregated to form the light-blue c node in the reaction graph. Resulting is a graph of variables that represents the aggregate strengths of dependencies from any direction.

To aggregate the stencil link coefficients, we use Fisher's z-transformation. It stabilizes the variance of the correlation coefficients, making them more suitable for averaging. The process involves converting each coefficient into a z-score, computing the arithmetic mean of the z-scores, and then converting the average z-score back to a correlation coefficient using the inverse Fisher's z-transformation. This method ensures that the combined value accurately reflects the underlying dependencies between variables.

8.5 Benchmarking M-CaStLe with VARs

We developed random and stable multivariate space-time systems with two spatial dimensions using mathematically defined ground-truth causal stencil graphs to evaluate the performance of M-CaStLe with a variety of system parameters.

8.5.1 Background: Univariate Space-Time VARs

Our methodology for generating data builds upon the work used by Nichol et al. (2024), which is fully detailed by Nichol et al. (2023). They developed a procedure for generating benchmark datasets of stable 2D space-time systems through the systematic construction of coefficient matrices parameterizing VARs of order 1 (VAR(1)s). Causal graphs have a direct mapping from VARs (Peters et al., 2017;

Runge et al., 2019a), which enables precise benchmark comparisons between VAR modeled data and causal discovery estimated graphs.

A system on an $N \times M$ grid with T time samples, $X \in \mathbb{R}^{N \times M \times T}$ with elements $X_{i,j,t}$, can be modeled by a VAR(1) with

$$\boldsymbol{X}_{t} = \boldsymbol{A}\boldsymbol{X}_{t-1} + \boldsymbol{\eta}_{t}, \tag{8.2}$$

where A is the coefficient matrix encoding linear dependencies between all variables in the system and η represents independent *innovations* on X for each variable at each time step. In this case, innovations are modeled with a unit normal distribution.

The space-time VAR methodology initializes a 3×3 matrix defining local grid-level dynamics between neighbors, called the neighborhood dependence matrix (NDM). Random NDMs of predetermined sparsity, d, are generated to describe how every grid cell in the space is dependent on the grid cells in its Moore neighborhood. To simulate an entire grid, the NDM can be structurally mapped to an A matrix for the entire grid. For an $N\times M$ grid space, $A\in\mathbb{R}^{NM\times NM}$. Finally, most 2D VARs are not numerically stable. To ensure stability, $\rho(A)<1.0$, where $\rho(A)$ is the spectral radius of A (Strang, 2016, p.307). Through the NDM definition, VARs can simulate locality in physical systems.

8.5.2 Multivariate Space-Time VARs

To adapt the space-time VAR procedure for multivariate systems, we grow the NDM in a new variable dimension, which gets mapped to a larger, flat, *A* matrix. The multivariate NDM describes interactions between multiple variables at the local grid-level, enabling VAR modeling of multivariate space-time dynamics.

For a system of V variables, the multivariate dynamics are represented by set of $V \times V \times 3 \times 3$ matrices. Each 3×3 matrix corresponds to the space-time dependence structure of a particular pair of parent and child variables. Like the univariate NDM, each entry in each 3×3 matrix is a coefficient value representing the influence of the entry's spatial location in the Moore neighborhood on the center location.

The NDM is mapped to an \mathbf{A} matrix, which represents the interactions of every grid cell-variable on every other grid cell-variable. For a grid of size $N \times M$ spatial dimensions and V variables, the matrix $\mathbf{A} \in \mathbb{R}^{NMV \times NMV}$. With the computed \mathbf{A} matrix, we again enforce stability by ensuring $\rho(\mathbf{A}) < 1.0$, where $\rho(\mathbf{A})$ is the spectral radius of \mathbf{A} .

With a stable **A** matrix, experimental data can be generated for any number of grid cells, time samples, local dependencies, and variables. Although **A** is larger, the VARs still have the form of Equation 8.2. Since most **A** matrices will be unstable, our implementation uses an accept-reject scheme similar to the univariate approach of Nichol et al. (2024) to generate stable **A** matrices:

- 1. Generate a random set of 3×3 *local dynamics matrices*, $\{C_{ij}\}$, for each pair of child and parent variables, resulting in $V \times V$ matrices. Each C_{ij} has d non-zero elements, including the central element (autocorrelation), where $1 \le d \le 9$. Each of the d non-zero elements, $\{a_i\}_{i=1}^d$, have a random value $1.0 \ge$ coefficient $i \ge s_*$.
- 2. Expand $\{C_{ij}\}$ to form the matrix \mathbf{A} for a grid of size $N \times M$, resulting in $\mathbf{A} \in \mathbb{R}^{NMV \times NMV}$.
- 3. If $|\lambda_{\max}(\mathbf{A})| \geq 1$, scale \mathbf{A} by $|\lambda_{\max}(\mathbf{A})|$.
- 4. If $c < s_* \ \forall c \in \mathbf{A}$, reject, else accept.

where $|\lambda_{\max}(A)|$ is the maximum absolute eigenvalue of A. This is used to sample from the set of statistically stationary & spatially homogeneous VARs on a 2D grid with minimum signal strengths $s_* \geq 0.1$ and fixed sparsity levels in the range $d \in \{1, 2, ..., 9\}$.

8.6 Results

We present empirical results of M-CaStLe's performance on our VAR benchmarks varying: the number of variables, grid sizes, the number of graph dependencies (graph edges), and the magnitude of coefficients. These demonstrate that M-CaStLe is suitable for estimation of local multivariate space-time structures from gridded data. Additionally, we compare M-CaStLe's performance to the popular PC algorithm for causal discovery.

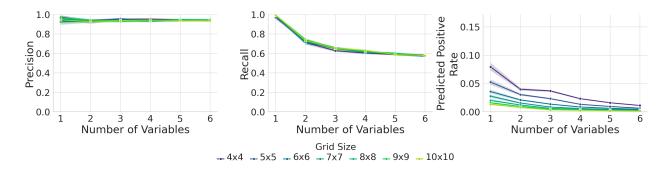


Figure 8.4: Showing precision and recall alongside predicted positive rate, a measure of how often a positive is predicted among all other predictions. As variables increase, the predicted positive rate decreases, which diminishes recall.

8.6.1 Metrics

Since VARs map directly to ground-truth causal graphs, we measured M-CaStLe's performance using binary classification measures. Let G = (V, E) be the ground-truth graph where V is the set of nodes and $E \subseteq V \times V$ is the set of edges. For any node pair $(i, j) \in V \times V$, a positive instance is defined as $(i, j) \in E$ and a negative instance as $(i, j) \notin E$. This enables our usage of precision, recall, and F_1 score, defined as follows:

$$Precision = \frac{TP}{TP + FP}$$
 (8.3)

$$Recall = \frac{TP}{TP + FN}$$
 (8.4)

$$F_1 = \frac{2 \cdot Precision \cdot Recall}{Precision + Recall}$$
(8.5)

where TP, FP, TN, and FN denote true positives, false positives, true negatives, and false negatives, respectively. Put simply, precision is the proportion of correctly detected positives to all detected positives, with a range of [0,1], where 1 is a perfect precision; recall is the proportion of correctly detected positives to how

many positives should have been detected, with a range of [0,1], where 1 is a perfect recall; and the F_1 score is the harmonic mean of precision and recall, with a range of [0,1], where 1 indicates perfect graph estimation.

8.6.2 Data Generation

We used the following data generation parameter ranges with 30 replicates each:

- Time samples T = 1000
- $N \times N$ grid sizes where $N \in [4, 5, 6, 7, 8, 9, 10]$
- Number of variables $V \in [1, 2, 3, 4, 5, 6]$
- Density $d \in (0, \dots 0.5]$
- Coefficients $c \in [0.1, 1.0]$

where density is relative to the stencil graph density: $d = \frac{L}{(3\times3\times V^2)}$ with L links, such that $d \le 1$. Since a V = 1 system can have up to L = 9, the most allowable here are L = 4. A V = 6 system may have $L \in [1, \dots 162]$. However, not all densities produced 30 stable systems after 48 hours of the accept-reject scheme described in Section 8.5.2. It is clear that there are zero systems in the limit of increasing density with a given minimum coefficient size. Appendix A details which of the above combinations successfully produced 30 systems for analysis. In total, 56,283 experiments were generated, with more experiments for systems of more variables.

8.6.3 Multivariate Performance

Figure 8.4 illustrates precision, recall, and positive prediction rate (PPR) in our experiments as the number of variables increases, with individual lines for each of the grid sizes. All available densities are marginalized in each line, with 95% confidence intervals. We found that precision is very high in all cases, regardless of the number of variables or grid size, with an average value of ≈ 0.94 . This indicates that when M-CaStLe identifies a positive link, it is likely to be a true positive. We found that recall is very high for V=1 and decreases as the number of variables increase, with the mean value ≈ 0.62 . This indicates that M-CaStLe may be relatively conservative, identifying a little more than half of the true links in the systems with more variables. However, it may also indicate limitations of the synthetic data model.

To shed some light on the recall results, we considered PPR. PPR is the fraction of all possible connections that were predicted as positive, regardless of correctness, given by

$$PPR = \frac{TP + FP}{TP + FP + TN + FN},$$
(8.6)

with a range of [0,1], where 1 indicates all possible edges were estimated (a fully dense graph). No particular PPR value necessarily indicates good performance, because it is a measure of the estimated graph's density.

In Figure 8.4, we see that recall and PPR are both decreasing as graph size increases (larger grid size and more variables). This possibly indicates that as the

graphs are getting larger, signals are more challenging to detect. We investigated the data generation model's apparent limitations in Appendix A. Figure A7 demonstrates that fewer stable systems could be generated for larger graphs, relative to their potential. A8 demonstrates that as the number of links increases among all systems, the maximum and minimum coefficients in each system quickly decrease. This indicates that the systems may be more challenging to correctly estimate, suggesting that M-CaStLe's recall may be more reflective of the data generating model than being a conservative estimator.

8.6.4 Comparison to the PC Algorithm

Nichol et al. (2024) compared CaStLe to several prior causal discovery methods and found CaStLe outperformed the others, particularly has the data dimensionality increased. In the multivariate regime, the data's dimensionality is multiplied by the number of variables. Multivariate systems should be far more challenging for causal discovery without dimensionality reduction. Here, we compare M-CaStLe to the PC algorithm, which is still in popular (Glymour et al., 2019) use and is the predecessor to most constraint-based causal discovery algorithms.

Figure 8.5 shows the F_1 score of M-CaStLe-PC and PC with increasing links, with the number of variables held constant to V = 4. The remaining V variables are given in Appendix B.1. We see that M-CaStLe-PC's F_1 score is consistently much higher than PC's. PC struggles with the very high dimensionality of the system since it is naive to the spatial and variable structure. Given that F_1 score is

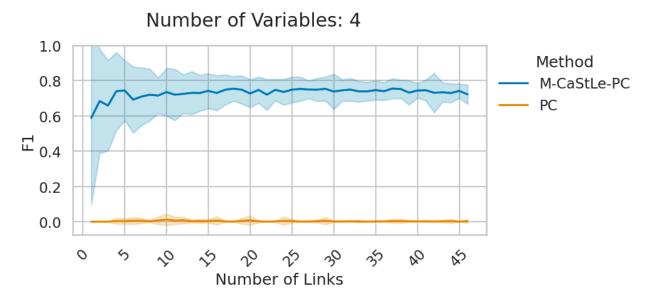


Figure 8.5: A comparison between M-CaStLe-PC and PC considering the F_1 score for V=4 as the number of links increases on a 4×4 grid. M-CaStLe-PC outperforms PC in every case because PC struggles with the very high dimensionality of the systems since it is naive to the spatial and variable structures.

the harmonic mean of precision and recall, we can see that M-CaStLe's aggregate performance is between the very high precision and relatively low recall described above.

8.6.5 Exploring Recall

To better evaluate the reason for M-CaStLe's relatively low recall, we tested it on a separate set of benchmark systems. In these, we constructed simple systems with many more variables and a range of coefficient magnitudes. The systems model a chain of dependence between each variable where there is one link per variable. The link is assigned a random parent location in the Moore neighborhood, and points to the center of the next variable. With this, we model different spatial relationships between variables but only one between variables. We explored $V \in$

 $\{10, 50, 100, 200\}$ and a set of 20 coefficients $\{c_i\}_{i=0}^{19}$ logarithmically spaced from 0.01 to 2.0, where $c_i = 0.01 \times 10^{2i/19}$. Every link had the same coefficient for each realization. Each realization had T = 1000 time samples and we restricted the grid size to 4×4 , which is the most challenging for M-CaStLe because there are fewer spatial replicates to leverage.

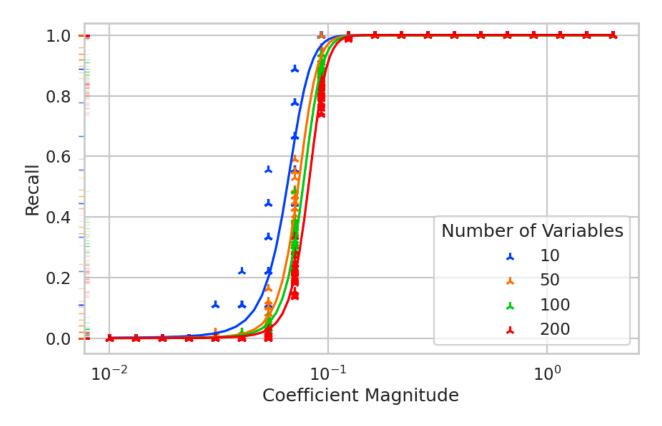


Figure 8.6: In simple chains of multivariate stencils, even with an extremely large number of variables, recall can be captured perfectly if the signal strength is large enough.

Figure 8.6 illustrates that recall increases proportionately with coefficient magnitude for all numbers of variables. Recall is 0 when coefficients are too small and 1 when they are large enough. There is an inflection interval in the coefficient magnitudes in which recall increases sharply. The three-parameter sigmoid functions fit to each set set of Vs shows that recall is ordered by V. That means that,

while high recall is achievable for up to 200 variables, systems with more variables are marginally more challenging to estimate, which conforms to our expectations. These results show that high recall is possible in high variable regimes if signals are strong enough.

8.7 Discussion

We have proposed M-CaStLe, a multivariate extension to space-time grid-level causal discovery with CaStLe (Nichol et al., 2024). M-CaStLe adapts both the Locally Encoded Neighborhood Structure construction and Parent-Identification Phase to learn inter-variable relationships in gridded space-time data. To represent these complex relationships, M-CaStLe produces a multivariate causal stencil graph that depicts which variable at each location in a Moore neighborhood causes each variable. To aid interpretation of the multivariate stencil, we introduced a decomposition method to extract spatial relationships and inter-variable relationships separately with the *spatial graph* and *reaction graph*.

Like CaStLe, M-CaStLe overcomes the limitations of high-dimensional gridded space-time systems, where there are more grid cells to estimate that time series samples in each. The inclusion of multiple variables exacerbates the highdimensional challenge, but M-CaStLe includes variable structures in the spatial replicates it leverages to form the LENS. The LENS collects repeating multivariate spatial structures to form a 3×3 spatial data representation of the underlying dynamics. With that, the PIP recovers the multivariate stencil describing the underlying causal relationships that define the system's grid-level behavior.

We developed a multivariate gridded space-time benchmark framework, building upon the work by Nichol et al. (2023). The benchmark defines mathematical structures (VAR models) representing the space-time relationships between grid cells with multiple variables per grid cell. The structures directly translate to causal graphs for ground-truth evaluation.

M-CaStLe performed well in the benchmark experiments. Its precision and recall were near 1 in systems with multiple variables when signal strengths were large enough. We applied the time series adapted PC causal discovery algorithm to the same benchmarks. We found that M-CaStLe had much better performance on multivariate systems than the PC causal discovery algorithm.

Recall suffered in highly complex systems cases because more complex systems exhibited smaller signal strengths per interaction. This supports our hypothesis that larger and more complex systems with many interacting components have fewer stable parameterizations. That is additionally supported by recent work investigating the *piranha problem* (Tosh et al., 2025), which describes the inevitable consequence that large complex systems will converge to weaker signals to maintain stability.

While we have demonstrated that M-CaStLe can identify multivariate spacetime dynamics, more work is needed to understand its application in real-world settings. Nichol et al. (2023) demonstrated CaStLe on the advective, transient dynamics of the Mount Pinatubo eruption's volcanic plume. A natural next step would characterize the atmospheric chemistry of the $SO_2 \rightarrow H_2SO_4$ pathways and how it mediates solar radiation and surface temperatures. One challenge described by Nichol et al. (2023) was having sufficient spatial and temporal data resolutions to capture the effects of interest on the grid-level. Earth system models can output data at sufficiently high resolutions, as they must compute in them to model realistic physics (Golaz et al., 2022), but input/output speeds and storage limitations may sometimes be bottlenecks. Nonetheless, as technologies improve, more expressive datasets will be available and more meaningful analysis methods will be critical for their evaluation. Further, satellite imagery now produces very high spatial resolutions, but, depending on the quantities and regions of interest, may have lower temporal sampling rates. However, as more satellites are deployed and technologies continue to improve, they will provide a greater wealth of data. Other application domains, such as computational chemistry, fluid dynamics, and spatiotemporal pharmacokinetics can modeled or observed at sufficiently high resolutions given their smaller scale in comparison to the Earth system.

While some dataset limitations still exist, Nichol et al. (2023) proposed other future research directions that may yield value in spite of those limitations. In particular, where spatial resolution is insufficiently matched temporal resolution, extending CaStLe and M-CaStLe to collect and evaluate larger neighborhoods, such as a radius-2 Moore neighborhood, could enable finding causal relationships that skip over immediately adjacent grid cells.

In this work, we have introduced M-CaStLe, a multivariate extension to the

grid-level space-time causal discovery meta-algorithm, CaStLe. M-CaStLe addresses the significant challenge of estimating causal relationships in high-dimensional space-time systems with multiple interacting variables, which traditional approaches struggle to handle effectively. By enabling the simultaneous estimation of spacetime dynamics and inter-variable interactions, M-CaStLe can enable advances in our understanding of complex systems, particularly in fields such as atmospheric science, computational fluid dynamics, computational chemistry, spatiotemporal pharmacokinetics, and epidemiological modeling. Our benchmark experiments demonstrate that M-CaStLe outperforms existing methods in accuracy, making it a robust and valuable tool for scientific discovery and analysis. Univariate CaStLe made a significant step in the analysis of high-dimensional grid-level dynamics and M-CaStLe makes multivariate space-time analysis possible. As a powerful tool for uncovering intricate causal relationships, M-CaStLe paves the way for more informed decision-making and deeper insights into the underlying mechanisms of complex phenomena.

Appendices

A Completed Data Generation Parameters

As noted in Section 8.6, not all parameter combinations generated stable systems. Here, we present the parameter ranges that did successfully generate 30 replicates to produce out results. We additionally evaluate the range of coefficient sizes generated, demonstrating the difficulty of creating complex systems with strong signals and many interdependencies.

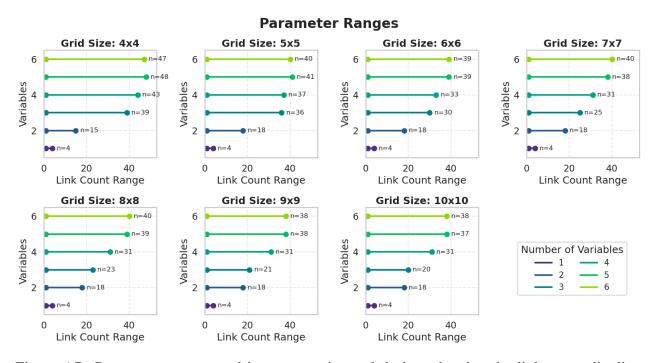


Figure A7: Parameter ranges used in our experimental design, showing the link count distribution for each grid size and variable count combination. Each horizontal line represents the span of network links tested, with each parameter combination having at least 30 replicate experiments (no values shown). Our experiments covered grid sizes from 4×4 to 10×10 and 1-6 variables per grid. All experiments used 1000 time samples and coefficient values between 0.1 and 1.0. The network density, d, defined as the ratio of actual links, d, to maximum possible links $d = \frac{L}{(3\times3\times V^2)}$, where $d \in (0, \dots 0.5]$. Not all density values produced 30 stable systems within our computational constraints, particularly at higher densities. This visualization shows which parameter combinations successfully generated sufficient replicates for statistical analysis.

Parameter ranges used in our experimental design, showing the link count dis-

tribution for each grid size and variable count combination. Each horizontal line represents the span of network links tested, with each parameter combination having at least 30 replicate experiments (n values shown). Our experiments covered grid sizes from 4×4 to 10×10 and 1-6 variables per grid. All experiments used 1000 time samples and coefficient values between 0.1 and 1.0. The network density, defined as the ratio of actual links (L) to maximum possible links in a 3×3 stencil graph (d = L/($3\times3\times$ V²)), ranged from near zero to 0.5. Not all theoretical density values produced 30 stable systems within our computational constraints, particularly at higher densities. This visualization shows which parameter combinations successfully generated sufficient replicates for statistical analysis.

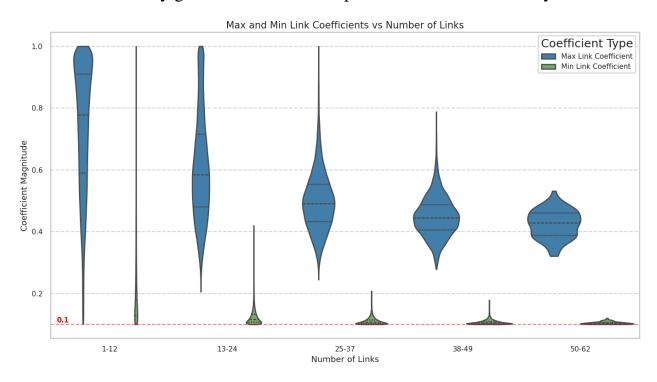


Figure A8: The relationship between link coefficients and the number of links present. As the number of links increases, maximum (blue) and minimum (green) link coefficients show a clear decreasing trend, with their distribution becoming narrower and centered around lower values. This reveals that networks with more links have weaker signals, suggesting that highly interconnected systems cannot be stable with large dependencies.

B Additional VAR Results

In this appendix, we present additional results related to the performance of our proposed method, M-CaStLe, with VAR benchmarks. We delve into various metrics that evaluate the effectiveness of M-CaStLe.

B.1 PC Comparison Results

We examined the impact of the number of variables on key performance indicators such as F_1 score, precision, and recall. We provide a comparison between M-CaStLe-PC and the time series PC algorithm. This analysis, illustrated in Figure B9, emphasizes how M-CaStLe-PC consistently outperforms PC across various scenarios, particularly as the number of links increases in a 4×4 grid. The results underscore the challenges faced by PC in high-dimensional environments, where its naive approach to spatial and variable structures limits its effectiveness.

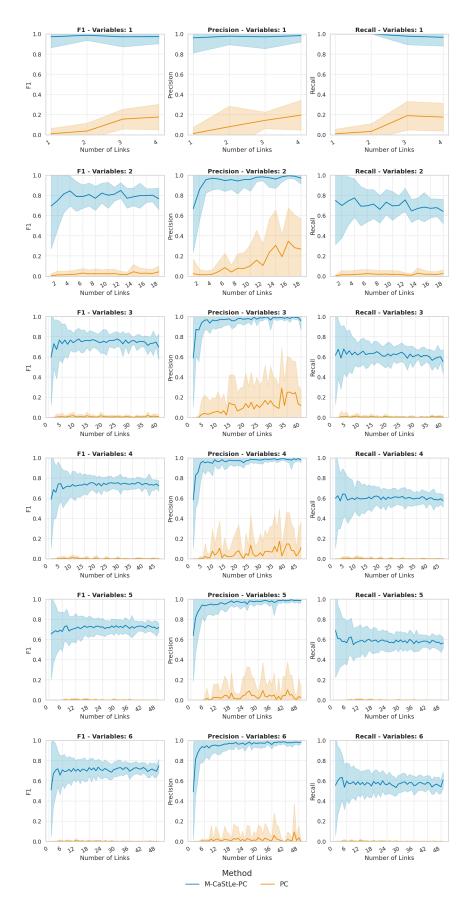


Figure B9: Comparisons between M-CaStLe-PC and PC considering the F_1 score, precision, and recall for all V as the number of links increases on a 4×4 grid. M-CaStLe-PC outperforms PC in every case because PC struggles with the very high dimensionality of the systems since it is naive to the spatial and variable structures.

9 Conclusion

This dissertation's research rests upon the shoulders of 100 years of data-driven knowledge discovery. It does so by advancing our understanding of what contemporary methods are capable of for complex systems and filling a critical research gap in the discovery of underlying local dynamics. The Causal Space-Time Stencil Learning (CaStLe) meta-algorithm developed here enables scalable causal discovery of grid-level dynamics in multiple variables for high-dimensional data—an important and elusive advancement in causal discovery research, particularly for the Earth sciences. These contributions equip scientists to approach more nuanced problems to explain the complex systems that rule our environment.

This chapter first summarizes Parts I and II of the work detailed in this dissertation and then explores exciting future avenues of research. Part I explored the foundational work I completed in exploring the capabilities of machine learning feature importance and state-of-the-art causal discovery for structure learning in the Earth sciences. Part II described my contributions to grid-level causal discovery with CaStLe and M-CaStLe.

9.1 Part I: Synthesis of Foundations Work

9.1.1 Machine Learning Feature Importance for Climate Models

Chapter 3 sought to learn if machine learning (ML) feature importance can be used to identify differences between climate model ensemble members' output data and observed data from satellite reanalysis products. In particular, I wanted to understand if I could predict and explain the Arctic's minimum yearly sea ice extent. Sea ice extent measures the square area of sea covered by ice, an important factor in Arctic life and trade vessel navigation. I trained ML models on 10 Arctic features that predict yearly sea ice extent minimums. Comparing ML model outputs between Arctic datasets gave us an understanding of their differences.

My methodology used separate random forest regression (RFR) (Breiman, 2001) models to learn from an observational dataset and five Energy Exascale Earth System Model (E3SM) (E3SM Project, 2018) simulation ensemble members. Random forests are ML models formed from aggregated decision trees. As RFR models train, they simultaneously build Gini importance values as part of the tree structures. It determines which features provide the most predictive power and encodes them in its Gini importance values. Thus, Gini importance describes how important each training feature is for the model's predictive power.

With the six trained models, I compared the calculated feature importance values to understand differences in the datasets. The baseline was data collected from satellite reanalysis products, which are observational datasets that use sophisticated

models to interpolate missing data where clouds obstructed satellites. With that, I could compare its feature importance values with those of RFR models trained on the climate model simulation runs. I found important similarities between the datasets, suggesting that the models captured some fundamental dynamics in the Arctic climate. The E3SM model runs were the most similar to each other and had some noticeable differences with the observational dataset. While both datasets identified the same six important features, the E3SM datasets consistently overweighted these features, with both ranking and magnitude discrepancies.

This work contributes to the broader climate analysis toolset by demonstrating how explainable machine learning can be used to learn about complex datasets. The work shows how physics-based models and ML can be used in tandem to learn more about critical systems in the Earth's climate. ML analyses like this can enhance climate model evaluation to improve existing model development and tuning practices. While more complex ML models proliferate, this work illustrates one important reason to maintain interpretability and explainability. Rather than simply demonstrating that discrepancies exist, analyses like this can help pinpoint potential sources of the discrepancies and lead climate model developers to the right place for refinement.

However, ML feature importance metrics are limited (Mandler and Weigand, 2024). The models themselves are subject to critical failures, such as various biases (Mehrabi et al., 2021), Simpson's paradox (Selvitella, 2017), and the Clever Hans effect (Lapuschkin et al., 2019), which can harm prediction performance or even

make them appear to have high predictive skill, whereas it performs poorly outside the given training and testing datasets (Lee and Chen, 2025). Feature importance itself can be misleading and failure-prone due to issues such as multicollinearity between features (Cammarota and Pinto, 2021). Even when everything works as intended, it is important to know that ML feature importance is not a causal description of the data's underlying generating process. However, it is rather a description of the trained model itself. (Parr and Wilson, 2021; Parr et al., 2024)

This research has been significantly extended and advanced with follow-on work by Brown et al. (2025), where several coauthors from our original study developed a novel pathway detection methodology. They went beyond a comparative analysis to create networks of connected features based on random forest feature importance to relate climate quantities. Their work builds on our initial claim that ML feature importance can be used to obtain insights into systems' underlying structure. The progression from feature importance comparisons to network construction demonstrates the continued impact of our initial insights.

The work in this chapter has become a part of a broader literature on machine learning for the Earth sciences (Labe and Barnes, 2022; Konya and Nematzadeh, 2024; Lao et al., 2024). In the subsequent chapters, I investigated causal inference frameworks to understand underlying dynamics better.

9.1.2 Causal Discovery for Climate Model Evaluation

Chapters 4 and 5 are complementary works where I explored applying a state-of-the-art causal discovery algorithm for the Arctic sea ice system. Chapter 4 discusses the research framework and methodology, and Chapter 5 discusses the implementation and results. This work extends the RFR feature importance approach in Chapter 3 to a causal discovery framework.

The PCMCI (Runge et al., 2019a) time series causal discovery algorithm was applied to Arctic climate features that may explain sea ice extent. PCMCI produces causal directed acyclic graphs (DAGs) that if its assumptions are satisfied, represent the estimated causal relationships between given features. While the RFR feature importance values describe *what* features are important, causal discovery can help answer *why* relationships exist between features. Comparing causal graphs estimated from different datasets and data sources enables a more mechanistic comparison. It can answer whether two data sources are *structurally* similar.

I used the F₁ score, the harmonic mean of precision and recall, as a similarity metric for comparing estimated causal graphs. I found that all data sources (observed and E3SM simulated) had similar graphs. However, the E3SM graphs were more dense, implying that more features were interconnected. This seems to corroborate our RFR feature importance finding that E3SM feature importances were over-weighting features, but more rigorous analysis is needed to confirm that connection more generally.

While the F₁ score is a good starting point for graph similarity, it cannot pinpoint where the graph differences are. I proposed further work to develop more node-level comparison metrics to better understand structural similarities and differences. I additionally recommended more subregional analyses—both the RFR and causal analyses evaluated quantities spanning the entire Arctic, and more meaningful insights may be gleaned from relating smaller regions within the Arctic.

Developing a better understanding of the smaller-scale processes that accumulate to produce emergent phenomena in the Earth system was the impetus for the work in Part II.

9.2 Part II: Discovery of Local Dynamics

9.2.1 Grid-Level Benchmarking of PCMCI

Chapter 6 developed grid-level space-time benchmarks for causal discovery methods and evaluated the PCMCI (Runge et al., 2019a) time series causal discovery algorithm. PCMCI was developed for highly autocorrelated time series data. It has been applied extensively in the Earth sciences (Runge et al., 2019c). However, its application methodology has been limited to regional analyses in which Earth science time series are obtained from dimensionality reduction methods such as weighted averages, principal component analysis (PCA), and related methods (Runge et al., 2015c; Tibau et al., 2022).

Our grid-level benchmark began with a 1D spatial grid and was extended to a 2D grid, for which each grid cell contained a time series with defined dependencies on its immediate neighbors. Both were structured as vector autoregression models (VARs), which enables a mathematically defined model that generates data and maps directly to a ground-truth causal graph. Using graph similarity metrics to compare PCMCI's estimated causal graph with each dataset's underlying VAR, I found PCMCI struggled to estimate the graphs well, except when it had unreal-istically high amounts of time samples per grid cell. In short, I determined that significant algorithmic advances would be needed to apply causal discovery like PCMCI at the grid-level.

The work presented computational advances as well. While using VARs for systems modeling and causal discovery benchmarking is not new (Runge et al., 2019d), my innovation was using them to model stable space-time dynamical systems with locally dependent grid cells. I produced gridded space-time data using a sliding dot product with a local neighborhood dependence matrix (NDM). In that way, they are similar to how cellular automata are defined, in which a single grid-level rule determines complex global behavior.

9.2.2 CaStLe: Grid-Level Causal Discovery

In Chapter 7, I introduced CaStLe, a grid-level causal discovery meta-algorithm. CaStLe addresses the fundamental challenge of causal discovery that I identified previously: many space-time gridded datasets are high-dimensional in practice.

High sample complexity reduces the power of causal discovery's statistical estimators. CaStLe remedies this by two central premises: underlying dynamics act locally, each grid cell influences only its neighbors, and neighboring grid cells generally exhibit similar dynamics. Through these, CaStLe leverages locality and stationarity to collect *informative spatial replicates* for local causal structures, which boosts efficiency and efficacy of the causal discovery task.

CaStLe produces a novel causal graph type, the *causal stencil graph*, which is a spatially structured graph representing a Moore neighborhood of nodes, which represent grid cells. The Moore neighborhood is a grid cell and its eight immediate neighbors. The stencil graph describes which neighbors are causal parents of the center node, enabling full representation of local causal structure.

CaStLe has two phases to estimate local grid-level structures. The first phase reorganizes the data into a smaller spatial representation, which I name the Locally Encoded Neighborhood Structure (LENS) in later work, which forms a 3×3 spatial embedding without marginalizing any data points. This embedding captures local causal structures by representing the Moore neighborhood allowing the detection of dependencies from all adjacent directions. The LENS phase multiplies the number of available samples through its collection of spatial replicates. Mathematically, this phase maps $\mathbb{R}^{M \times N \times T} \to \mathbb{R}^{3 \times 3 \times L}$ on an $M \times N$ grid over T time steps; L = T(M-1)(N-1) concatenated time series points.

The second phase is the Parent-Identification Phase (PIP), which applies an adapted time series causal discovery algorithm to target identification of the causal

parents of the LENS's center cell. Through that, the it can be determined which spatial neighbors influence the center cell. Any time series causal discovery algorithm may be implemented in this phase, given that it can be adapted. While not exhaustive of all existing algorithms, the adaptation has been trivial in our experience. Finally, once the PIP is applied to the LENS, the causal stencil graph is estimated.

I demonstrated the efficacy of CaStLe on three benchmark problems: atmospheric aerosol advection, the VAR benchmark presented in Chapter 6, and Burgers' equation, a partial differential equation (PDE) model of advection and diffusion. First, CaStLe correctly reconstructed the stratospheric aerosol advection dynamics from the 1991 Mount Pinatubo eruption with data from two climate models. VAR benchmarks enabled a careful parameter study of many different gridded systems with exact ground truth. It also contained a comparison of CaStLed methods and alternative causal discovery approaches, in which CaStLe outperformed all others. Finally, the study of Burgers' equation demonstrated that CaStLe can generally filter out diffusion "noise" to recover the primary transport mechanism. It shows that CaStLe can be applicable in many advection-transport systems, which are common in the Earth system.

Theoretical analysis showed marked improvement over the stat-of-the-art. Algorithms based on the PC algorithm will be bounded by a computational complexity of $\mathcal{O}(Tp^32^p)$, whereas CaStLe is bounded by $\mathcal{O}(Tp)$, for T time samples per p grid cells. Our analysis of CaStLe's sample complexity shows that its accuracy

does improve as grid sizes increase. This is in contrast to traditional approaches, which struggle more as grid sizes get larger.

CaStLe is generally applicable to physics-governed space-time systems that satisfy the locality and stationarity assumptions. These include many processes in Earth science, fluid dynamics, and other fields where effects propagate locally through space and exhibit consistent or smoothly changing dynamics across regions. It can be applied in extremely data-poor settings, where only short time intervals are observed. It is especially valuable in settings in which grid-level dynamics define the phenomena under study and marginalization would destroy that information. These include advective, transient, and non-periodic phenomena such as volcanic eruptions, wildfires, and traveling weather fronts. CaStLe is a flexible meta-algorithm, enabling implementation with today's best causal discovery algorithms and those of the future, including causal representation learning. It is highly extensible, being adaptable to multiple variables, more than two spatial dimensions, longer time lags, and larger local neighborhoods.

CaStLe provides another path for physical model evaluation by elucidating where and why behavior does not match intended dynamics. For the first time, grid-level processes are recoverable with causal discovery, which opens the door to future multi-scale analyses to determine how local structures give rise to emergent global patterns. However, this initial version of CaStLe is univariate—it can only estimate space-time dynamics of one quantity, such as aerosols. It would be significantly more valuable estimating the space-time dynamics of multiple vari-

ables and their interactions. This is precisely what Chapter 8 addresses.

9.2.3 M-CaStLe: Multivariate Grid-Level Causal Discovery

I followed the development of CaStLe with an extension enabling multivariate analyses simultaneously with space-time structure discovery. Chapter 8 details the methodological innovations making that possible. We adapted both phases of CaStLe, developed a method for interpretability, and benchmarked M-CaStLe.

CaStLe's LENS was adapted to include multiple variables per time series. The mapping from the given gridded space to the multivariate LENS is represented by the transformation $\mathbb{R}^{N\times N\times V\times T}\to\mathbb{R}^{3\times 3\times V\times L}$, where L=T(M-1)(N-1) denoted the length of each concatenated time series. With this, multiple variables' spacetime structures are captured. CaStLe's PIP was adapted by allowing each of the variables in the center grid cell of the LENS to be children and no other grid cells. That allows for an adapted time series causal discovery algorithm to estimate the multivariate space-time dynamics underlying the given data.

M-CaStLe was validated using the spatial VAR benchmark detailed previously with a multivariate extension. I found that M-CaStLe significantly outperforms the PC algorithm for grid-level multivariate causal discovery. It had remarkably high precision, and its recall was mediated by the size of coefficients in each VAR. Systems with more dependencies require smaller coefficients in order to be stable, but the signals become more challenging to detect amid the noise.

M-CaStLe is the first causal discovery approach to enable grid-level causal dis-

covery of multiple variables. This new capability can facilitate new research directions in physical systems such as the Earth sciences, computational chemistry, ecology, fluid dynamics, and pharmacokinetics. It presents many opportunities for interdisciplinary collaborations to analyze systems in a new way.

9.3 Connections and Research Frontiers

The research detailed in this dissertation traces a methodological journey from correlative machine learning approaches to mechanistic causal discovery frameworks for complex physical systems, with an emphasis on Earth science. The work spans multiple scales, progressing from regional analyses to tackling high-dimensional grid-level dynamics. The primary contribution, CaStLe, accomplished grid-level discovery for the first time by leveraging locality and stationarity principles simplifying the causal discovery task without sacrificing spatial information through dimensionality reduction. Instead, CaStLe maintains critical spatial structure by collecting informative spatial replicates. The resulting causal stencil graph describes local causal structures between grid cells in a highly interpretable format. M-CaStLe enables a more comprehensive system understanding by extending capabilities to multiple variables. This work provides scientists with new tools to discover how local dynamics give rise to emergent global phenomena by bridging statistical learning with physical interpretation. The following explores these connections and highlights promising research frontiers that build upon these methodological foundations.

Bibliography

- R. McGill, J. W. Tukey, and W. A. Larsen, "Variations of box plots," *The American Statistician*, vol. 32, no. 1, pp. 12–16, 1978.
- R. E. Bellman, *Dynamic Programming*. Princeton, NJ: Princeton University Press, 1957, introduced the term "curse of dimensionality" in the preface: "All this may be subsumed under the heading 'the curse of dimensionality.' Since this is a curse which has hung over the head of the physicist and astronomer for many a year...".
- P. Bühlmann and S. v. d. Geer, "Statistics for High-Dimensional Data, Methods, Theory and Applications," *Springer Series in Statistics*, pp. 99–182, 2011.
- J. Pearl and D. Mackenzie, *The Book of Why*. New York: Basic Books, 2018.
- J. Peters, D. Janzing, and B. Schlkopf, *Elements of Causal Inference: Foundations and Learning Algorithms*. Cambridge, Massachusetts: The MIT Press, 2017.
- J. J. Nichol, M. G. Peterson, K. J. Peterson, G. M. Fricke, and M. E. Moses, "Machine learning feature analysis illuminates disparity between E3SM climate models and observed climate change," *Journal of Computational and Applied Mathematics*, vol. 395, p. 113451, 10 2021.
- M. A. Hernán and J. M. Robins, Causal Inference: What if. Crc Press, 2020.

- D. B. Rubin, "Estimating causal effects of treatments in randomized and nonrandomized studies," *Journal of Educational Psychology*, vol. 66, no. 5, pp. 688–701, 1974.
- J. Robins, "A new approach to causal inference in mortality studies with a sustained exposure period—application to control of the healthy worker survivor effect," *Mathematical Modelling*, vol. 7, no. 9, pp. 1393–1512, 1986. [Online]. Available: https://www.sciencedirect.com/science/article/pii/0270025586900886
- H. F. Dorn, "Philosophy of Inferences from Retrospective Studies," *American Journal of Public Health and the Nations Health*, vol. 43, no. 6_Pt_1, pp. 677–683, 1953.
- A. R. Feinstein, "Clinical biostatistics," *Clinical Pharmacology & Therapeutics*, vol. 11, no. 2, pp. 282–292, 1970.
- A. P. Dawid, "Causal Inference Without Counterfactuals," *Journal of the American Statistical Association*, vol. 95, no. 450, p. 407, 2000.
- S. Wright, "Correlation and Causation," *Journal of Agricultural Research*, vol. 20, 1921. [Online]. Available: https://books.google.co.ao/books/about/Journal_of_Agricultural_Research. html?hl=pt-PT&id=lNNdIV_qpwIC&utm_source=gb-gplus-shareJournal

- P. Spirtes, C. Glymour, and R. Scheines, "Causation, Prediction, and Search," *Lecture Notes in Statistics*, 1993.
- J. Pearl, "Causal Diagrams for Empirical Research," *Biometrika*, vol. 82, no. 4, p. 669, 1995.
- L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001. [Online]. Available: http://dx.doi.org/10.1023/A%3A1010933404324
- S. Nembrini, I. R. König, and M. N. Wright, "The revival of the Gini importance?" *Bioinformatics*, vol. 34, no. 21, pp. 3711–3718, 2018.
- S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Advances in Neural Information Processing Systems 31*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 4765–4774. [Online]. Available: https://dl.acm.org/doi/10.5555/3295222.3295230
- B. Krishnapuram, M. Shah, A. Smola, C. Aggarwal, D. Shen, R. Rastogi, M. T. Ribeiro, S. Singh, and C. Guestrin, ""Why Should I Trust You?"," in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, ser. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, 2016, pp. 1135–1144.
- A. Shrikumar, P. Greenside, and A. Kundaje, "Learning Important Features

through Propagating Activation Differences," in *Proceedings of the 34th International Conference on Machine Learning - Volume 70*, ser. ICML'17. JMLR.org, 2017, p. 3145–3153.

- J. Pearl, "The Do-Calculus Revisited," arXiv, 2012.
- ——, "FROM BAYESIAN NETWORKS TO CAUSAL NETWORKS," *Mathematical Models for Handling Partial Knowledge in Artificial Intelligence*, pp. 157–182, 1995.
- P. Spirtes and C. Glymour, "An Algorithm for Fast Recovery of Sparse Causal Graphs," *Social Science Computer Review*, vol. 9, no. 1, pp. 62–72, 1991.
- J. Runge, P. Nowack, M. Kretschmer, S. Flaxman, and D. Sejdinovic, "Detecting and quantifying causal associations in large nonlinear time series datasets," *Science Advances*, vol. 5, no. 11, pp. 4996—5023, 2019. [Online]. Available: http://advances.sciencemag.org/
- J. Runge, "Causal network reconstruction from time series: From theoretical assumptions to practical estimation," *Chaos: An Interdisciplinary Journal of Non-linear Science*, vol. 28, no. 7, p. 075310, 2018.
- J. Arnhold, P. Grassberger, K. Lehnertz, and C. Elger, "A robust method for detecting interdependences: application to intracranially recorded EEG," *Physica D: Nonlinear Phenomena*, vol. 134, no. 4, pp. 419–430, 1999.
- G. Sugihara, R. May, H. Ye, C.-h. Hsieh, E. Deyle, M. Fogarty, and S. Munch,

- "Detecting Causality in Complex Ecosystems," *Science*, vol. 338, no. 6106, pp. 496–500, 2012. [Online]. Available: https://www.science.org/doi/abs/10.1126/science.1227079
- P. Spirtes and K. Zhang, "Causal discovery and inference: concepts and recent methodological advances," *Applied Informatics*, vol. 3, no. 1, p. 3, 2016.
- J. Pearl, "Graphs, Causality, and Structural Equation Models," *Sociological Methods & Research*, vol. 27, no. 2, pp. 226–284, 1998. [Online]. Available: https://doi.org/10.1177/0049124198027002004
- J. Runge, "Discovering contemporaneous and lagged causal relations in autocorrelated nonlinear time series datasets," in *Proceedings of the 36th Conference on Uncertainty in Artificial Intelligence*, ser. Proceedings of Machine Learning Research, vol. 124. PMLR, 2020, pp. 1388–1397. [Online]. Available: https://proceedings.mlr.press/v124/runge20a.html
- C. Glymour, K. Zhang, and P. Spirtes, "Review of Causal Discovery Methods Based on Graphical Models," *Frontiers in Genetics*, vol. 10, p. 524, 2019.
- J. Runge, S. Bathiany, E. Bollt, G. Camps-Valls, D. Coumou, E. Deyle, C. Glymour, M. Kretschmer, M. D. Mahecha, J. Munoz-Mari, E. H. v. Nes, J. Peters, R. Quax, M. Reichstein, M. Scheffer, B. Scholkopf, P. Spirtes, G. Sugihara, J. Sun, K. Zhang, and J. Zscheischler, "Inferring causation from time series in Earth system sciences," *Nature Communications*, vol. 10, no. 2553, 2019.

- S. Guo, G. J. S. Bluth, W. I. Rose, I. M. Watson, and A. J. Prata, "Re-evaluation of SO2 release of the 15 June 1991 Pinatubo eruption using ultraviolet and infrared satellite sensors," *Geochemistry, Geophysics, Geosystems*, vol. 5, no. 4, 2004.
- L. Gimeno, R. Nieto, M. Vázquez, and D. A. Lavers, "Atmospheric rivers: a minireview," *Frontiers in Earth Science*, vol. 2, p. 2, 2014.
- G. Wang, D. Zhong, T. Li, Y. Zhang, C. Meng, M. Zhang, X. Song, J. Wei, and Y. Huang, "Study on sky rivers: Concept, theory, and implications," *Journal of Hydro-environment Research*, vol. 21, pp. 109–117, 2018.
- I. Ebert-Uphoff and Y. Deng, "Causal Discovery from Spatio-Temporal Data with Applications to Climate Science," 2014 13th International Conference on Machine Learning and Applications, pp. 606–613, 2014.
- R. Fisher, *Statistical Methods for Research Workers*, ser. Biological monographs and manuals. Oliver and Boyd, 1925. [Online]. Available: https://books.google.com/books?id=I0NBAAAAIAAJ
- N. S. Hall, "R. A. Fisher and his advocacy of randomization," *Journal of the History of Biology*, vol. 40, no. 2, pp. 295–325, 2007.
- J. Splawa-Neyman, D. M. Dabrowska, and T. P. Speed, "On the Application of Probability Theory to Agricultural Experiments. Essay on Principles. Section 9," *Roczniki Nauk Rolniczych Tom X [in Polish]*; translated in Statistical

- *Science*, vol. 5, no. 4, pp. 465–472, 1923. [Online]. Available: http://www.jstor.org/stable/2245382
- D. B. Rubin, "Causal Inference Using Potential Outcomes: Design, Modeling, Decisions," *Journal of the American Statistical Association*, vol. 100, no. 469, pp. 322–331, 2005.
- P. W. Holland, "Statistics and Causal Inference," *Journal of the American Statistical Association*, vol. 81, no. 396, p. 945, 1986.
- J. Pearl, *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2000. [Online]. Available: https://books.google.com/books?id=wnGU_TsW3BQC
- ——, "Causal inference in the health sciences: a conceptual introduction," *Health services and outcomes research methodology*, vol. 2, no. 3, pp. 189–220, 2001.
- P. Tarka, "An overview of structural equation modeling: its beginnings, historical development, usefulness and controversies in the social sciences," *Quality & Quantity*, vol. 52, no. 1, pp. 313–354, 2018.
- N. Wiener, "The theory of prediction," Modern mathematics for engineers, 1956.
- C. W. J. Granger, "Investigating Causal Relations by Econometric Models and Cross-spectral Methods," *Econometrica*, vol. 37, no. 3, p. 424, 1969, granger Causality seminal paper.

- H. Reichenbach, *The Direction of Time*, ser. Dover books on physics, M. Reichenbach, Ed. Dover Publications, 1956.
- S. Shimizu, P. O. Hoyer, A. Hyvarinen, and A. Kerminen, "A Linear Non-Gaussian Acyclic Model for Causal Discovery," *Journal of Machine Learning Research*, vol. 7, no. 72, p. 2003-2030, 2006. [Online]. Available: https://www.jmlr.org/papers/volume7/shimizu06a/shimizu06a.pdf
- J. D. Ramsey, "A Scalable Conditional Independence Test for Nonlinear, Non-Gaussian Data," *arXiv*, vol. abs/1401.5031, 2014. [Online]. Available: http://arxiv.org/abs/1401.5031
- K. Zhang, J. Peters, D. Janzing, and B. Schölkopf, "Kernel-Based Conditional Independence Test and Application in Causal Discovery," in *Proceedings of the Twenty-Seventh Conference on Uncertainty in Artificial Intelligence*, ser. UAI'11. Arlington, Virginia, USA: AUAI Press, 2011, p. 804–813.
- J. Runge, "Conditional independence testing based on a nearest-neighbor estimator of conditional mutual information," vol. 84, pp. 938–947, 09–11 Apr 2018. [Online]. Available: https://proceedings.mlr.press/v84/runge18a.html
- R. Sen, A. T. Suresh, K. Shanmugam, A. G. Dimakis, and S. Shakkettai, "Model-Powered Conditional Independence Test," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS'17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 2955–2965.

- J. Houghton, Y. Ding, D. Griggs, M. Noguer, P. v. d. Linden, X. Dai, M. Maskell, and C. Johnson, "Climate Change 2001: The Scientific Basis," *Contribution of Working Group I to the Third Assessment Report of the Intergovernmental Panel on Climate Change (IPCC)*, vol. 881., p. 881, 2001.
- K. Hasselmann, "Multi-pattern fingerprint method for detection and attribution of climate change," *Climate Dynamics*, vol. 13, no. 9, pp. 601–611, 1997.
- I. Ebert-Uphoff and Y. Deng, "A new type of climate network based on probabilistic graphical models: Results of boreal winter versus summer," *Geophysical Research Letters*, vol. 39, no. 19, 2012.
- Z. S. Kaufman, N. Feldl, W. Weijer, and M. Veneziani, "Causal Interactions Between Southern Ocean Polynyas and High-Latitude Atmosphere-Ocean Variability," *Journal of Climate*, vol. 33, no. 11, pp. 4891–4905, 2020.
- M. Kretschmer, D. Coumou, J. F. Donges, and J. Runge, "Using Causal Effect Networks to Analyze Different Arctic Drivers of Midlatitude Winter Circulation," *Journal of Climate*, vol. 29, no. 11, pp. 4069–4081, 2016.
- P. Nowack, J. Runge, V. Eyring, and J. D. Haigh, "Causal networks for climate model evaluation and constrained projections," *Nature Communications* 2020 11:1, vol. 11, no. 1, pp. 1—11, 2020. [Online]. Available: http://www.nature.com/articles/s41467-020-15195-y
- J. Runge, V. Petoukhov, and J. Kurths, "Quantifying the Strength and Delay of Cli-

matic Interactions: The Ambiguities of Cross Correlation and a Novel Measure Based on Graphical Models," *Journal of Climate*, vol. 27, no. 2, pp. 720–739, 2014.

- J. Runge, R. V. Donner, and J. Kurths, "Optimal model-free prediction from multivariate time series," *Physical Review E*, vol. 91, no. 5, p. 052909, 2015.
- M. Vejmelka, L. Pokorná, J. Hlinka, D. Hartman, N. Jajcay, and M. Paluš, "Non-random correlation structures and dimensionality reduction in multivariate climate data," *Climate Dynamics*, vol. 44, no. 9-10, pp. 2663–2682, 2014.
- V. Eyring, P. M. Cox, G. M. Flato, P. J. Gleckler, G. Abramowitz, P. Caldwell, W. D. Collins, B. K. Gier, A. D. Hall, F. M. Hoffman, G. C. Hurtt, A. Jahn, C. D. Jones, S. A. Klein, J. P. Krasting, L. Kwiatkowski, R. Lorenz, E. Maloney, G. A. Meehl, A. G. Pendergrass, R. Pincus, A. C. Ruane, J. L. Russell, B. M. Sanderson, B. D. Santer, S. C. Sherwood, I. R. Simpson, R. J. Stouffer, and M. S. Williamson, "Taking climate model evaluation to the next level," *Nature Climate Change*, vol. 9, no. 2, pp. 102–110, 2019, "Other promising diagnostic developments on the horizon that should be further advanced include studies that assess responses to perturbations rather than mean climate95, and the application of innovative data science methods in Earth system science96 such as neural networks97, machine learning-based anomaly detection techniques98, graphical models and causal discovery99.".
- A. Hannachi, I. T. Jolliffe, and D. B. Stephenson, "Empirical orthogonal functions

- and related techniques in atmospheric science: A review," *International Journal of Climatology*, vol. 27, no. 9, pp. 1119–1152, 2007.
- A. Gerhardus and J. Runge, "LPCMCI: Causal Discovery in Time Series with Latent Confounders," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 12615–12625. [Online]. Available: https://proceedings.neurips.cc/paper/2020/file/94e70705efae423efda1088614128d0b-Paper.pdf
- A. Tsonis and P. Roebber, "The architecture of the climate network," *Physica A: Statistical Mechanics and its Applications*, vol. 333, pp. 497–504, 2004.
- A. A. Tsonis, K. L. Swanson, and P. J. Roebber, "What Do Networks Have to Do with Climate?" *Bulletin of the American Meteorological Society*, vol. 87, no. 5, pp. 585–595, 2006.
- X.-A. Tibau, C. Reimers, A. Gerhardus, J. Denzler, V. Eyring, and J. Runge, "A spatiotemporal stochastic climate model for benchmarking causal discovery methods for teleconnections," *Environmental Data Science*, vol. 1, p. e12, 2022.
- B. Schölkopf, F. Locatello, S. Bauer, N. R. Ke, N. Kalchbrenner, A. Goyal, and Y. Bengio, "Toward Causal Representation Learning," *Proceedings of the IEEE*, vol. 109, no. 5, pp. 612–634, 2021.
- J. Boussard, C. Nagda, J. Kaltenborn, C. E. E. Lange, P. Brouillard, Y. Gurwicz,

- P. Nowack, and D. Rolnick, "Towards Causal Representations of Climate Model Data," *arXiv*, 2023.
- P. Brouillard, S. Lachapelle, J. Kaltenborn, Y. Gurwicz, D. Sridhar, A. Drouin, P. Nowack, J. Runge, and D. Rolnick, "Causal Representation Learning in Temporal Data via Single-Parent Decoding," *arXiv*, 2024.
- P. Pfleiderer, C.-F. Schleussner, T. Geiger, and M. Kretschmer, "Robust predictors for seasonal Atlantic hurricane activity identified with causal effect networks," *Weather and Climate Dynamics*, vol. 1, no. 2, pp. 313–324, 2020.
- I. Polkova, H. Afargan-Gerstman, D. I. V. Domeisen, M. P. King, P. Ruggieri, P. Athanasiadis, M. Dobrynin, O. Aarnes, M. Kretschmer, and J. Baehr, "Predictors and prediction skill for marine cold-air outbreaks over the Barents Sea," *Quarterly Journal of the Royal Meteorological Society*, vol. 147, no. 738, pp. 2638–2656, 2021.
- J. Y. Zhu, C. Zhang, H. Zhang, S. Zhi, V. O. Li, J. Han, and Y. Zheng, "pg-Causality: Identifying Spatiotemporal Causal Pathways for Air Pollutants with Urban Big Data," *IEEE Transactions on Big Data*, vol. 4, no. 4, pp. 571–585, 2016.
- P. Sheth, R. Shah, J. Sabo, K. S. Candan, and H. Liu, "STCD: A Spatio-Temporal Causal Discovery Framework for Hydrological Systems," *2022 IEEE International Conference on Big Data* (*Big Data*), vol. 00, pp. 5578–5583, 2022.

- S. Saetia, N. Yoshimura, and Y. Koike, "Constructing Brain Connectivity Model Using Causal Network Reconstruction Approach," *Frontiers in Neuroinformatics*, vol. 15, p. 619557, 2021.
- E3SM Project, "Energy Exascale Earth System Model (E3SM)," [Computer Software] https://dx.doi.org/10.11578/E3SM/dc.20180418.36, Apr. 2018. [Online]. Available: https://dx.doi.org/10.11578/E3SM/dc.20180418.36
- J. Stroeve and D. Notz, "Changing state of Arctic sea ice across all seasons," sep 2018.
- "Arctic report card 2019," National Oceanic and Atmospheric Administration, Washington, DC, Tech. Rep., December 2019, released at the American Geophysical Union Fall Meeting in San Francisco, California, December 10, 2019. [Online]. Available: https://www.arctic.noaa.gov/Report-Card
- L. C. Smith and S. R. Stephenson, "New trans-Arctic shipping routes navigable by midcentury," *PNAS*, vol. 110, no. 13, pp. 4871–4872, 2013.
- H. Goosse, J. E. Kay, K. C. Armour, A. Bodas-Salcedo, H. Chepfer, D. Docquier *et al.*, "Quantifying climate feedbacks in polar regions," *Nature Communications*, vol. 9, no. 1919, 2018.
- F. Sevellec, A. V. Fedorov, and W. Liu, "Arctic sea-ice decline weakens the atlantic meridional overturning circulation," *Nature Climate Change*, vol. 7, pp. 604–610, 2017.

- J. Cohen, K. Pfeiffer, and J. A. Francis, "Warm Arctic episodes linked with increased frequency of extreme winter weather in the United States," *Nature Communications*, vol. 9, no. 869, 2018.
- I. Cvijanovic, B. D. Santer, C. Bonfils, D. D. Lucas, J. C. H. Chiang, and S. Zimmerman, "Future loss of Arctic sea-ice cover could drive a substantial decrease in california's rainfall," *Nature Communications*, vol. 8, no. 1947, 2017.
- E. Rosenblum and I. Eisenman, "Sea ice trends in climate models only accurate in runs with biased global warming," *Journal of Climate*, vol. 30, no. 16, pp. 6265–6278, aug 2017.
- A. G. Meehl, C. Covey, T. Delworth, M. Latif, B. Mcavaney, J. F. B. Mitchell *et al.*, "THE WCRP CMIP3 Multimodel Dataset: A New Era in Climate Change Research," *American Meteorological Society*, no. September, 2007.
- J. Stroeve, M. M. Holland, W. Meier, T. Scambos, and M. Serreze, "Arctic sea ice decline: Faster than forecast," *Geophysical Research Letters*, vol. 34, no. 9, 2007.
- M. G. A. Taylor Karl E., Stouffer Ronald J., "An Overview of CMIP5 and the Experiment Design," *American Meteorological Society*, vol. 3, no. april, pp. 485–498, 2012.
- V. Eyring, S. Bony, G. A. Meehl, C. A. Senior, B. Stevens, R. J. Stouffer *et al.*, "Overview of the Coupled Model Intercomparison Project

- Phase 6 (CMIP6) experimental design and organization," *Geoscientific Model Development*, vol. 9, no. 5, pp. 1937–1958, 2016. [Online]. Available: https://www.geosci-model-dev.net/9/1937/2016/
- J. C. Stroeve, V. Kattsov, A. Barrett, M. Serreze, T. Pavlova, M. Holland *et al.*, "Trends in Arctic sea ice extent from CMIP5, CMIP3 and observations," *Geophysical Research Letters*, vol. 39, no. 16, pp. 1–7, 2012.
- M. Ionita, K. Grosfeld, P. Scholz, R. Treffeisen, and G. Lohmann, "September Arctic Sea Ice minimum prediction a new skillful statistical approach," *Earth System Dynamics Discussions*, pp. 1–23, sep 2018. [Online]. Available: https://www.earth-syst-dynam-discuss.net/esd-2018-61/
- T. G. Reid and P. M. Tarantino, "Arctic sea ice extent forecasting using support vector regression," in *Proceedings 2014 13th International Conference on Machine Learning and Applications, ICMLA 2014.* Institute of Electrical and Electronics Engineers Inc., feb 2014, pp. 1–6.
- G. Peng, W. N. Meier, D. J. Scott, M. H. Savoie, and N. Snow, "A long-term and reproducible passive microwave sea ice concentration data record for climate studies and monitoring," *Earth System Science Data*, pp. 311–318, 2013.
- A. Schweiger, R. Lindsay, J. Zhang, M. Steele, H. Stern, and R. Kwok, "Uncertainty in modeled Arctic sea ice volume," *Journal of Geophysical Research: Oceans*, vol. 116, no. 9, pp. 1–21, 2011.

- NOAA, OAR, and ESRL-PSD, "Ncep-doe reanalysis 2," 2019, nCEP_Reanalysis 2 data provided by the NOAA/OAR/ESRL PSD, Boulder, Colorado, USA. [Online]. Available: https://www.esrl.noaa.gov/psd/
- —, "Noaa extended reconstructed sea surface temperature," 2019, NOAA_ERSST_V4 data provided by the NOAA/OAR/ESRL PSD, Boulder, Colorado, USA. [Online]. Available: https://www.esrl.noaa.gov/psd/
- J.-C. Golaz, P. M. Caldwell, L. P. Van Roekel, M. R. Petersen, Q. Tang, J. D. Wolfe, G. Abeshu, V. Anantharaj, X. S. Asay-Davis, D. C. Bader, S. A. Baldwin, G. Bisht, P. A. Bogenschutz, M. Branstetter, M. A. Brunke, S. R. Brus, S. M. Burrows, P. J. Cameron-Smith, A. S. Donahue, M. Deakin, R. C. Easter, K. J. Evans, Y. Feng, M. Flanner, J. G. Foucar, J. G. Fyke, B. M. Griffin, C. Hannay, B. E. Harrop, M. J. Hoffman, E. C. Hunke, R. L. Jacob, D. W. Jacobsen, N. Jeffery, P. W. Jones, N. D. Keen, S. A. Klein, V. E. Larson, L. R. Leung, H.-Y. Li, W. Lin, W. H. Lipscomb, P.-L. Ma, S. Mahajan, M. E. Maltrud, A. Mametjanov, J. L. McClean, R. B. McCoy, R. B. Neale, S. F. Price, Y. Qian, P. J. Rasch, J. E. J. Reeves Eyre, W. J. Riley, T. D. Ringler, A. F. Roberts, E. L. Roesler, A. G. Salinger, Z. Shaheen, X. Shi, B. Singh, J. Tang, M. A. Taylor, P. E. Thornton, A. K. Turner, M. Veneziani, H. Wan, H. Wang, S. Wang, D. N. Williams, P. J. Wolfram, P. H. Worley, S. Xie, Y. Yang, J.-H. Yoon, M. D. Zelinka, C. S. Zender, X. Zeng, C. Zhang, K. Zhang, Y. Zhang, X. Zheng, T. Zhou, and Q. Zhu, "The DOE E3SM Coupled Model Version 1: Overview and Evaluation at Standard

- Resolution," *Journal of Advances in Modeling Earth Systems*, vol. 11, no. 7, pp. 2089–2129, 2019.
- J. E. Kay, C. Deser, A. Phillips, A. Mai, C. Hannay, G. Strand *et al.*, "The community earth system model (cesm) large ensemble project: A community resource for studying climate change in the presence of internal climate variability," *Bulletin of the American Meteorological Society*, vol. 96, no. 8, pp. 1333–1349, 2015. [Online]. Available: https://doi.org/10.1175/BAMS-D-13-00255.1
- B. J. Zib, X. Dong, B. Xi, and A. Kennedy, "Evaluation and intercomparison of cloud fraction and radiative fluxes in recent reanalyses over the arctic using BSRN surface observations," *Journal of Climate*, vol. 25, no. 7, pp. 2291–2305, 2012.
- L. Breiman, "Bagging predictors," *Machine learning*, vol. 24, no. 2, pp. 123–140, 1996.
- R. E. Banfield, L. O. Hall, K. W. Bowyer, and W. P. Kegelmeyer, "A Comparison of Decision Tree Ensemble Creation Techniques," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 29, no. 1, pp. 173–180, Jan. 2007.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel *et al.*, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.

- L. Breiman, "Arcing classifiers," *Ann. Statist.*, vol. 26, no. 3, pp. 801–849, 06 1998. [Online]. Available: https://doi.org/10.1214/aos/1024691079
- S. Van Der Walt, S. C. Colbert, and G. Varoquaux, "The numpy array: a structure for efficient numerical computation," *Computing in Science & Engineering*, vol. 13, no. 2, p. 22, 2011.
- M. Waskom and the seaborn development team, "mwaskom/seaborn," Sep. 2020. [Online]. Available: https://doi.org/10.5281/zenodo.592845
- L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and Regression Trees*. Wadsworth International Group, Belmont, California, 1984.
- T. Hengl, M. Nussbaum, M. N. Wright, G. Heuvelink, and B. Gräler, "Random forest as a generic framework for predictive modeling of spatial and spatio-temporal variables," *PeerJ*, vol. 6, 2018. [Online]. Available: https://doi.org/10.7717/peerj.5518
- L. Breiman, "Rejoinder: Arcing classifiers," *The Annals of Statistics*, vol. 26, no. 3, pp. 841–849, 1998. [Online]. Available: http://www.jstor.org/stable/120059
- R. E. Banfield, L. O. Hall, K. W. Bowyer, and W. P. Kegelmeyer, "A new ensemble diversity measure applied to thinning ensembles," in *Multiple Classifier Systems*,
 T. Windeatt and F. Roli, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2003, pp. 306–316.
- U. S. Bhatt, P. Bieniek, C. Bitz, E. Blanchard-Wrigglesworth, H. Eicken,

- H. Goessling *et al.*, "2019 sea ice outlook full post-season report," February 2020, editors: Turner-Bogren, B. and H. V. Wiggins. [Online]. Available: https://www.arcus.org/sipn/sea-ice-outlook/2019/post-season
- A. C. Ordonez, C. M. Bitz, and E. Blanchard-Wrigglesworth, "Processes controlling Arctic and Antarctic sea ice predictability in the Community Earth System Model," *Journal of Climate*, vol. 31, pp. 9771–9786, 2018.
- E. Blanchard-Wrigglesworth, K. C. Armour, C. M. Bitz, and E. Deweaver, "Persistence and inherent predictability of arctic sea ice in a GCM ensemble and observations," *Journal of Climate*, vol. 24, no. 1, pp. 231–250, jan 2011.
- J. Runge, S. Bathiany, E. Bollt, G. Camps-Valls, D. Coumou, E. Deyle, C. Glymour, M. Kretschmer, M. D. Mahecha, J. Munoz-Mari, E. H. v. Nes, J. Peters, R. Quax, M. Reichstein, M. Scheffer, B. Scholkopf, P. Spirtes, G. Sugihara, J. Sun, K. Zhang, and J. Zscheischler, "Inferring causation from time series in Earth system sciences," *Nature Communications*, vol. 10, no. 1, 2019.
- P. Nowack, J. Runge, V. Eyring, and J. D. Haigh, "Causal networks for climate model evaluation and constrained projections," *Nature Communications* 2020 11:1, vol. 11, no. 1, pp. 1—11, 2020. [Online]. Available: http://www.nature.com/articles/s41467-020-15195-y
- J. Pearl, "Causal inference in statistics: An overview," *Statistics Surveys*, vol. 3, no. September, pp. 96—146, 2009.

- P. Spirtes and K. Zhang, "Causal discovery and inference: concepts and recent methodological advances," *Applied Informatics*, vol. 3, no. 1, p. 3, 2016.
- A. C. I. Assessment, *Impacts of a Warming Arctic: Arctic Climate Impact Assessment*. Cambridge University Press, 2004, aCIA Overview report.
- J. Runge, V. Petoukhov, J. F. Donges, J. Hlinka, N. Jajcay, M. Vejmelka, D. Hartman, N. Marwan, M. Paluš, and J. Kurths, "Identifying causal gateways and mediators in complex spatio-temporal systems," *Nature Communications*, vol. 6, no. 1, p. 8502, 2015.
- J. Pearl and D. Mackenzie, *The Book of Why*. New York: Basic Books, 2018.
- J. Runge, P. Nowack, M. Kretschmer, S. Flaxman, and D. Sejdinovic, "Detecting and quantifying causal associations in large nonlinear time series datasets," *Science Advances*, vol. 5, no. 11, p. eaau4996, Nov 2019. [Online]. Available: https://www.science.org/doi/10.1126/sciadv.aau4996
- J. Runge, "Quantifying information transfer and mediation along causal pathways in complex systems," *Physical Review E*, vol. 92, no. 6, p. 062829, 2015.
- M. H. Hitchman, M. McKay, and C. R. Trepte, "A climatology of stratospheric aerosol," *Journal of Geophysical Research: Atmospheres*, vol. 99, no. D10, pp. 20689–20700, 1994. [Online]. Available: https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/94JD01525
- D. B. A. Jones, H. R. Schneider, and M. B. McElroy, "Effects of the quasi-biennial

- oscillation on the zonally averaged transport of tracers," *Journal of Geophysical Research: Atmospheres*, vol. 103, no. D10, pp. 11235–11249, 1998.
- V. Aquila, C. I. Garfinkel, P. Newman, L. Oman, and D. Waugh, "Modifications of the quasi-biennial oscillation by a geoengineering perturbation of the stratospheric aerosol layer," *Geophysical Research Letters*, vol. 41, no. 5, pp. 1738–1744, 2014.
- L. J. Gray, J. A. Anstey, Y. Kawatani, H. Lu, S. Osprey, and V. Schenzinger, "Surface impacts of the Quasi Biennial Oscillation," *Atmospheric Chemistry and Physics*, vol. 18, no. 11, pp. 8227–8247, 2018.
- E. C. Neto, M. P. Keller, A. D. Attie, and B. S. Yandell, "Causal graphical models in systems genetics: A unified framework for joint inference of causal network and genetic architecture for correlated phenotypes," *The Annals of Applied Statistics*, vol. 4, no. 1, pp. 320–339, 2010.
- X. Zhang, X.-M. Zhao, K. He, L. Lu, Y. Cao, J. Liu, J.-K. Hao, Z.-P. Liu, and L. Chen, "Inferring gene regulatory networks from gene expression data by path consistency algorithm based on conditional mutual information," *Bioinformatics*, vol. 28, no. 1, pp. 98–104, 2011.
- M. Kamiński, M. Ding, W. A. Truccolo, and S. L. Bressler, "Evaluating causal relations in neural systems: Granger causality, directed transfer function and statistical assessment of significance," *Biological Cybernetics*, vol. 85, no. 2, pp. 145–157, 2001.

- A. A. Tsonis, E. R. Deyle, H. Ye, and G. Sugihara, "Convergent Cross Mapping: Theory and an Example," *Advances in Nonlinear Geosciences*, pp. 587–600, 2017.
- Y. Deng and I. Ebert-Uphoff, "Weakening of atmospheric information flow in a warming climate in the Community Climate System Model," *Geophysical Research Letters*, vol. 41, no. 1, pp. 193–200, 2014.
- J. Runge, V. Petoukhov, J. F. Donges, J. Hlinka, N. Jajcay, M. Vejmelka, D. Hartman, N. Marwan, M. Paluš, and J. Kurths, "Identifying causal gateways and mediators in complex spatio-temporal systems," *Nature Communications*, vol. 6, no. 1, p. 8502, 2015.
- G. D. Capua, M. Kretschmer, R. V. Donner, B. v. d. Hurk, R. Vellore, R. Krishnan, and D. Coumou, "Tropical and mid-latitude teleconnections interacting with the Indian summer monsoon rainfall: a theory-guided causal effect network approach," *Earth System Dynamics*, vol. 11, no. 1, pp. 17–34, 2019.
- G. D. Capua, J. Runge, R. V. Donner, B. v. d. Hurk, A. G. Turner, R. Vellore, R. Krishnan, and D. Coumou, "Dominant patterns of interaction between the tropics and mid-latitudes in boreal summer: causal relationships and the role of timescales," *Weather and Climate Dynamics*, vol. 1, no. 2, pp. 519–539, 2020.
- C. Krich, J. Runge, D. G. Miralles, M. Migliavacca, O. Perez-Priego, T. El-Madany, A. Carrara, and M. D. Mahecha, "Estimating causal networks in

- biosphere–atmosphere interaction with the PCMCI approach," *Biogeosciences*, vol. 17, no. 4, pp. 1033–1061, 2020.
- E. Galytska, K. Weigel, D. Handorf, R. Jaiser, R. H. Köhler, J. Runge, and V. Eyring, "Causal model evaluation of Arctic-midlatitude teleconnections in CMIP6," *Journal of Geophysical Research: Atmospheres*, vol. 128, no. 17, 2022.
- T. J. O'Kane, D. Harries, and M. A. Collier, "Bayesian Structure Learning for Climate Model Evaluation," *Journal of Advances in Modeling Earth Systems*, vol. 16, no. 5, 2024.
- H. Zhao, V. Kitsios, T. J. O'Kane, and E. V. Bonilla, "Bayesian Factorised Granger-Causal Graphs For Multivariate Time-series Data," *arXiv*, 2024.
- J. Runge, A. Gerhardus, G. Varando, V. Eyring, and G. Camps-Valls, "Causal inference for time series," *Nature Reviews Earth & Environment*, vol. 4, no. 7, pp. 487–505, 2023.
- J. Pearl, M. Glymour, and N. Jewell, *Causal Inference in Statistics: A Primer*. Wiley, 2016. [Online]. Available: https://books.google.com/books?id=L3G-CgAAQBAJ
- C. Glymour and R. Scheines, "Causal modeling with the TETRAD program," *Synthese*, vol. 68, no. 1, pp. 37–63, 1986.
- J. Pearl and T. S. Verma, "A statistical semantics for causation," *Statistics and Computing*, vol. 2, no. 2, pp. 91–95, 1992.

- X. Zheng, B. Aragam, P. Ravikumar, and E. P. Xing, "DAGs with NO TEARS: Continuous Optimization for Structure Learning," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, ser.
 NIPS'18. Red Hook, NY, USA: Curran Associates Inc., 2018, p. 9492–9503.
- Y. Liu, A. Niculescu-Mizil, A. Lozano, and Y. Lu, "Learning Temporal Causal Graphs for Relational Time-Series Analysis," in *Proceedings of the 27th International Conference on International Conference on Machine Learning*, ser. ICML'10. Madison, WI, USA: Omnipress, 2010, p. 687–694.
- R. Pamfil, N. Sriwattanaworachai, S. Desai, P. Pilgerstorfer, K. Georgatzis,
 P. Beaumont, and B. Aragam, "DYNOTEARS: Structure Learning from Time-Series Data," *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, vol. 108, pp. 1595–1605, 2020. [Online].
 Available: https://proceedings.mlr.press/v108/pamfil20a.html
- S. Ali, U. Hasan, X. Li, O. Faruque, A. Sampath, Y. Huang, M. O. Gani, and J. Wang, "Causality for Earth Science A Review on Time-series and Spatiotemporal Causality Methods," *arXiv*, 2024.
- D. Colombo and M. H. Maathuis, "Order-independent constraint-based causal structure learning," *Journal of Machine Learning Research*, vol. 15, no. 1, pp. 3741–3782, 2014.
- S. L. Brunton, J. L. Proctor, and J. N. Kutz, "Discovering governing equations from data by sparse identification of nonlinear dynamical systems," *Proceedings*

- of the National Academy of Sciences of the United States of America, vol. 113, no. 15, pp. 3932–3937, 2016.
- Z. Li, N. Kovachki, K. Azizzadenesheli, B. Liu, K. Bhattacharya, A. Stuart, and A. Anandkumar, "Fourier Neural Operator for Parametric Partial Differential Equations," *arXiv*, 2020.
- J. Pathak, S. Subramanian, P. Harrington, S. Raja, A. Chattopadhyay, M. Mardani, T. Kurth, D. Hall, Z. Li, K. Azizzadenesheli, P. Hassanzadeh, K. Kashinath, and A. Anandkumar, "FourCastNet: A Global Data-driven High-resolution Weather Model using Adaptive Fourier Neural Operators," *arXiv*, 2022.
- J. Hart, M. Gulian, I. Manickam, and L. P. Swiler, "Solving High-Dimensional Inverse Problems with Auxiliary Uncertainty via Operator Learning with Limited Data," *Journal of Machine Learning for Modeling and Computing*, vol. 4, no. 2, pp. 105–133, 2023.
- K. Bhattacharjee, N. Naskar, S. Roy, and S. Das, "A survey of cellular automata: types, dynamics, non-uniformity and applications," *Natural Computing*, vol. 19, no. 2, pp. 433–461, 2020.
- H. J. Miller, "Tobler's First Law and Spatial Analysis," *Annals of the Association of American Geographers*, vol. 94, no. 2, pp. 284–289, 2004.
- R. T. Walker, "GEOGRAPHY, VON THÜNEN, AND TOBLER'S FIRST LAW:

- TRACING THE EVOLUTION OF A CONCEPT," *Geographical Review*, vol. 112, no. 4, pp. 591–607, 2022.
- V. K. Raghu, J. D. Ramsey, A. Morris, D. V. Manatakis, P. Sprites, P. K. Chrysanthis, C. Glymour, and P. V. Benos, "Comparison of strategies for scalable causal discovery of latent variable models from mixed data," *International Journal of Data Science and Analytics*, vol. 6, no. 1, pp. 33–45, 2018.
- S. Guo, W. I. Rose, G. J. S. Bluth, and I. M. Watson, "Particles in the great Pinatubo volcanic cloud of June 1991: The role of ice," *Geochemistry, Geophysics, Geosystems*, vol. 5, no. 5, 2004.
- S. Kremser, L. W. Thomason, M. v. Hobe, M. Hermann, T. Deshler, C. Timmreck, M. Toohey, A. Stenke, J. P. Schwarz, R. Weigel, S. Fueglistaler, F. J. Prata, J. Vernier, H. Schlager, J. E. Barnes, J. Antuña-Marrero, D. Fairlie, M. Palm, E. Mahieu, J. Notholt, M. Rex, C. Bingen, F. Vanhellemont, A. Bourassa, J. M. C. Plane, D. Klocke, S. A. Carn, L. Clarisse, T. Trickl, R. Neely, A. D. James, L. Rieger, J. C. Wilson, and B. Meland, "Stratospheric aerosol—Observations, processes, and impact on climate," *Reviews of Geophysics*, vol. 54, no. 2, pp. 278–335, 2016.
- E. G. Dutton and J. R. Christy, "Solar radiative forcing at selected locations and evidence for global lower tropospheric cooling following the eruptions of El Chichón and Pinatubo," *Geophysical Research Letters*, vol. 19, no. 23, pp. 2313–2316, 1992.

- K. Labitzke and M. P. McCormick, "Stratospheric temperature increases due to Pinatubo aerosols," *Geophysical Research Letters*, vol. 19, no. 2, pp. 207–210, 1992.
- D. E. Parker, H. Wilson, P. D. Jones, J. R. Christy, and C. K. FOLLAND, "The impact of Mount Pinatubo on world-wide temperatures," *International Journal of Climatology*, vol. 16, no. 5, pp. 487–497, 1996.
- B. J. Soden, R. T. Wetherald, G. L. Stenchikov, and A. Robock, "Global Cooling After the Eruption of Mount Pinatubo: A Test of Climate Feedback by Water Vapor," *Science*, vol. 296, no. 5568, pp. 727–730, 2002.
- K. E. Trenberth and A. Dai, "Effects of Mount Pinatubo volcanic eruption on the hydrological cycle as an analog of geoengineering," *Geophysical Research Letters*, vol. 34, no. 15, 2007.
- M. Weylandt and L. P. Swiler, "Beyond PCA: Additional Dimension Reduction Techniques to Consider in the Development of Climate Fingerprints," *Journal of Climate*, vol. 37, no. 5, pp. 1723–1735, 2024.
- D. E. Parker, H. Wilson, P. D. Jones, J. R. Christy, and C. K. Folland, "The impact of mount pinatubo on world-wide temperatures," *International Journal of Climatology*, vol. 16, no. 5, pp. 487–497, 1996. [Online]. Available: https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/%28SICI%291097-0088% 28199605%2916%3A5%3C487%3A%3AAID-JOC39%3E3.0.CO%3B2-J

- A. Robock, "Volcanic eruptions and climate," *Reviews of Geophysics*, vol. 38, no. 2, pp. 191–219, 2000. [Online]. Available: https://agupubs.onlinelibrary. wiley.com/doi/abs/10.1029/1998RG000054
- C. Timmreck, "Modeling the climatic effects of large explosive volcanic eruptions," *Wiley Interdisciplinary Reviews: Climate Change*, vol. 3, no. 6, pp. 545–564, 2012.
- L. R. Marshall, E. C. Maters, A. Schmidt, C. Timmreck, A. Robock, and M. Toohey, "Volcanic effects on climate: recent advances and future avenues," *Bulletin of Volcanology*, vol. 84, no. 5, p. 54, 2022.
- J. P. Hollowed, C. Jablonowski, H. Y. Brown, B. R. Hillman, D. L. Bull, and J. L. Hart, "Localized injections of interactive volcanic aerosols and their climate impacts in a simple general circulation model," *EGUsphere*, vol. 2024, pp. 1–38, 2024.
- J. Golaz, L. P. V. Roekel, X. Zheng, A. F. Roberts, J. D. Wolfe, W. Lin, A. M. Bradley, Q. Tang, M. E. Maltrud, R. M. Forsyth, C. Zhang, T. Zhou, K. Zhang, C. S. Zender, M. Wu, H. Wang, A. K. Turner, B. Singh, J. H. Richter, Y. Qin, M. R. Petersen, A. Mametjanov, P. Ma, V. E. Larson, J. Krishna, N. D. Keen, N. Jeffery, E. C. Hunke, W. M. Hannah, O. Guba, B. M. Griffin, Y. Feng, D. Engwirda, A. V. D. Vittorio, C. Dang, L. M. Conlon, C. Chen, M. A. Brunke, G. Bisht, J. J. Benedict, X. S. Asay-Davis, Y. Zhang, M. Zhang, X. Zeng, S. Xie, P. J. Wolfram, T. Vo, M. Veneziani, T. K. Tesfa, S. Sreepathi, A. G. Salinger,

- J. E. J. R. Eyre, M. J. Prather, S. Mahajan, Q. Li, P. W. Jones, R. L. Jacob, G. W. Huebler, X. Huang, B. R. Hillman, B. E. Harrop, J. G. Foucar, Y. Fang, D. S. Comeau, P. M. Caldwell, T. Bartoletti, K. Balaguru, M. A. Taylor, R. B. McCoy, L. R. Leung, and D. C. Bader, "The DOE E3SM Model Version 2: Overview of the Physical Model and Initial Model Evaluation," *Journal of Advances in Modeling Earth Systems*, vol. 14, no. 12, 2022.
- H. Y. Brown, B. Wagman, D. Bull, K. Peterson, B. Hillman, X. Liu, Z. Ke, and L. Lin, "Validating a microphysical prognostic stratospheric aerosol implementation in E3SMv2 using observations after the Mount Pinatubo eruption," *Geoscientific Model Development*, vol. 17, no. 13, pp. 5087–5121, 2024.
- M. Kalisch and P. Bühlmann, "Estimating high-dimensional directed acyclic graphs with the PC-algorithm," *Journal of Machine Learning Research*, vol. 8, pp. 613–636, 2007. [Online]. Available: https://www.jmlr.org/papers/v8/kalisch07a.html
- G. I. Allen, L. Gan, and L. Zheng, "Interpretable Machine Learning for Discovery: Statistical Challenges and Opportunities," *Annual Review of Statistics and Its Application*, vol. 11, no. 1, pp. 97–121, 2023.
- I. Fountalis, C. Dovrolis, A. Bracco, B. Dilkina, and S. Keilholz, "δ-MAPS: from spatio-temporal data to a weighted and lagged network between functional domains," *Applied Network Science*, vol. 3, no. 1, p. 21, 2018.
- N. Nukavarapu, J.-A. Yang, and M. M. Jankowska, "Unsupervised Deep Learning

- Approach to Analyze Spatio-Temporal Change in Satellite Imagery," *IGARSS* 2023 2023 IEEE International Geoscience and Remote Sensing Symposium, vol. 00, pp. 2496–2499, 2023.
- W. L. Davis, M. L. Carlson, I. K. Tezaur, D. L. Bull, K. J. Peterson, and L. P. Swiler, "Spatio-temporal multivariate cluster evolution analysis for detecting and tracking climate impacts," *Journal of Computational and Applied Mathematics*, vol. 465, p. 116583, 2025.
- M. A. Thomas, M. A. Giorgetta, C. Timmreck, H.-F. Graf, and G. Stenchikov, "Simulation of the climate impact of Mt. Pinatubo eruption using ECHAM5 Part 2: Sensitivity to the phase of the QBO and ENSO," *Atmospheric Chemistry and Physics*, vol. 9, no. 9, pp. 3001–3009, 2009.
- J. J. Nichol, M. Weylandt, M. Smith, and L. Swiler, "Benchmarking the PCMCI Causal Discovery Algorithm for Spatiotemporal Systems," Sandia National Laboratories, Tech. Rep., 2023. [Online]. Available: https://www.osti.gov/biblio/1991387
- G. Strang, *Introduction to Linear Algebra*, 5th ed. Wellesley, MA: Wellesley-Cambridge Press, 2016.
- B. Matthews, "Comparison of the predicted and observed secondary structure of T4 phage lysozyme," *Biochimica et Biophysica Acta (BBA) Protein Structure*, vol. 405, no. 2, pp. 442–451, 1975. [Online]. Available: https://www.sciencedirect.com/science/article/pii/0005279575901099

- N. S. Diffenbaugh, J. S. Pal, R. J. Trapp, and F. Giorgi, "Fine-scale processes regulate the response of extreme events to global climate change," *Proceedings of the National Academy of Sciences*, vol. 102, no. 44, pp. 15774–15778, 2005.
- M. Palu, "Coupling in complex systems as information transfer across time scales," *Philosophical Transactions of the Royal Society A*, vol. 377, no. 2160, p. 20190094, 2019.
- A. Agarwal, L. Caesar, N. Marwan, R. Maheswaran, B. Merz, and J. Kurths, "Network-based identification and characterization of teleconnections on different scales," *Scientific Reports*, vol. 9, no. 1, p. 8808, 2019.
- Z. Zhang, G. Li, Y. Cai, X. Cheng, Y. Sun, J. Zhao, P. Shu, L. Ma, and Z. An, "Millennial-Scale Monsoon Variability Modulated by Low-Latitude Insolation During the Last Glaciation," *Geophysical Research Letters*, vol. 49, no. 1, 2022.
- J. Sjolte, F. Adolphi, H. Guðlaugsdóttir, and R. Muscheler, "Major Differences in Regional Climate Impact Between High- and Low-Latitude Volcanic Eruptions," *Geophysical Research Letters*, vol. 48, no. 8, 2021.
- K. Baranowski, C. Faust, P. Eby, and N. Bharti, "Quantifying the impacts of Australian bushfires on native forests and gray-headed flying foxes," *Global Ecology and Conservation*, vol. 27, p. e01566, 2021.
- A. E. Payne, M.-E. Demory, L. R. Leung, A. M. Ramos, C. A. Shields, J. J. Rutz, N. Siler, G. Villarini, A. Hall, and F. M. Ralph, "Responses and impacts of atmo-

- spheric rivers to climate change," *Nature Reviews Earth & Environment*, vol. 1, no. 3, pp. 143–157, 2020.
- J. Baño-Medina, A. Sengupta, J. D. Doyle, C. A. Reynolds, D. Watson-Parris, and L. D. Monache, "Are AI weather models learning atmospheric physics? A sensitivity analysis of cyclone Xynthia," *npj Climate and Atmospheric Science*, vol. 8, no. 1, p. 92, 2025.
- T. B. Higgins, A. C. Subramanian, P. A. G. Watson, and S. Sparrow, "Changes to Atmospheric River Related Extremes Over the United States West Coast Under Anthropogenic Warming," *Geophysical Research Letters*, vol. 52, no. 5, 2025.
- D. Keellings and H. Moradkhani, "Spatiotemporal Evolution of Heat Wave Severity and Coverage Across the United States," *Geophysical Research Letters*, vol. 47, no. 9, 2020.
- D. A. Driscoll, K. J. Macdonald, R. K. Gibson, T. S. Doherty, D. G. Nimmo, R. H. Nolan, E. G. Ritchie, G. J. Williamson, G. W. Heard, E. M. Tasker, R. Bilney, N. Porch, R. A. Collett, R. A. Crates, A. C. Hewitt, E. Pendall, M. M. Boer, J. Gates, R. L. Boulton, C. M. Mclean, H. Groffen, A. C. Maisey, C. T. Beranek, S. A. Ryan, A. Callen, A. J. Hamer, A. Stauber, G. J. Daly, J. Gould, K. L. Klop-Toker, M. J. Mahony, O. W. Kelly, S. L. Wallace, S. E. Stock, C. J. Weston, L. Volkova, D. Black, H. Gibb, J. J. Grubb, M. A. McGeoch, N. P. Murphy, J. S. Lee, C. R. Dickman, V. J. Neldner, M. R. Ngugi, V. Miritis, F. Köhler, M. Perri, A. J. Denham, B. D. E. Mackenzie, C. A. M. Reid, J. T. Rayment, A. Arriaga-

Jiménez, M. W. Hewins, A. Hicks, B. A. Melbourne, K. F. Davies, M. E. Bitters, G. D. Linley, A. C. Greenville, J. K. Webb, B. Roberts, M. Letnic, O. F. Price, Z. C. Walker, B. R. Murray, E. M. Verhoeven, A. M. Thomsen, D. Keith, J. S. Lemmon, M. K. J. Ooi, V. L. Allen, O. T. Decker, P. T. Green, A. Moussalli, J. K. Foon, D. B. Bryant, K. L. Walker, M. J. Bruce, G. Madani, J. L. Tscharke, B. Wagner, C. R. Nitschke, C. R. Gosper, C. J. Yates, R. Dillon, S. Barrett, E. E. Spencer, G. M. Wardle, T. M. Newsome, S. A. Pulsford, A. Singh, A. Roff, K. J. Marsh, K. Mcdonald, L. G. Howell, M. R. Lane, R. H. Cristescu, R. R. Witt, E. J. Cook, F. Grant, B. S. Law, J. Seddon, K. K. Berris, R. M. Shofner, M. Barth, T. Welz, A. Foster, D. Hancock, M. Beitzel, L. X. L. Tan, N. A. Waddell, P. M. Fallow, L. Schweickle, T. D. L. Breton, C. Dunne, M. Green, A.-M. Gilpin, J. M. Cook, S. A. Power, K. Hogendoorn, R. Brawata, C. J. Jolly, M. Tozer, N. Reiter, and R. D. Phillips, "Biodiversity impacts of the 2019–2020 Australian megafires," *Nature*, vol. 635, no. 8040, pp. 898–905, 2024.

- G. Goerg and C. Shalizi, "Mixed licors: A nonparametric algorithm for predictive state reconstruction," in *Proceedings of the Sixteenth International Conference on Artificial Intelligence and Statistics*, ser. Proceedings of Machine Learning Research, C. M. Carvalho and P. Ravikumar, Eds., vol. 31. Scottsdale, Arizona, USA: PMLR, 29 Apr–01 May 2013, pp. 289–297. [Online]. Available: https://proceedings.mlr.press/v31/goerg13a.html
- J. Burgers, "A Mathematical Model Illustrating the Theory of Turbulence," Ad-

- vances in Applied Mechanics, vol. 1, pp. 171–199, 1948.
- M. P. Bonkile, A. Awasthi, C. Lakshmi, V. Mukundan, and V. S. Aswin, "A systematic literature review of Burgers' equation with recent advances," *Pramana*, vol. 90, no. 6, p. 69, 2018.
- P. K. Rubenstein, S. Bongers, B. Schölkopf, and J. M. Mooij, "From Deterministic ODEs to Dynamic Structural Causal Models," in *UAI'18: Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, ser. UAI'18. AUAI Press, 2018. [Online]. Available: http://auai.org/uai2018/proceedings/papers/43.pdf
- J. Nichol, "CaStLe Data Release for JGR MLC 2024," Jul. 2024. [Online]. Available: https://doi.org/10.5281/zenodo.12701546
- J. J. Nichol, W. Michael, F. G. Matthew, M. E. Moses, D. Bull, and L. P. Swiler, "Space-time causal discovery in climate science: A local stencil learning approach," *ESS Open Archive* 172253117.78663487, 2024.
- D. B. Rubin, "Essential concepts of causal inference: a remarkable history and an intriguing future," *Biostatistics & Epidemiology*, vol. 3, no. 1, pp. 140–155, 2019.
- I. Ebert-Uphoff and Y. Deng, "Causal Discovery for Climate Research Using Graphical Models," *Journal of Climate*, vol. 25, no. 17, pp. 5648–5665, 2012.
- G. F. Cooper, I. Bahar, M. J. Becich, P. V. Benos, J. Berg, J. U. Espino, C. Gly-

- mour, R. C. Jacobson, M. Kienholz, A. V. Lee, X. Lu, R. Scheines, and team, and the Center for Causal Discovery, "The center for causal discovery of biomedical knowledge from big data," *Journal of the American Medical Informatics Association*, vol. 22, no. 6, pp. 1132–1136, 2015.
- A. Feder, K. A. Keith, E. Manzoor, R. Pryzant, D. Sridhar, Z. Wood-Doughty, J. Eisenstein, J. Grimmer, R. Reichart, M. E. Roberts, B. M. Stewart, V. Veitch, and D. Yang, "Causal Inference in Natural Language Processing: Estimation, Prediction, Interpretation and Beyond," *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 1138–1158, 2022.
- A. Zanga, E. Ozkirimli, and F. Stella, "A Survey on Causal Discovery: Theory and Practice," *International Journal of Approximate Reasoning*, vol. 151, pp. 101–129, 2022.
- A. Sadeghi, A. Gopal, and M. Fesanghary, "Causal Discovery in Financial Markets: A Framework for Nonstationary Time-Series Data," *arXiv*, 2023.
- M. Greenacre, P. J. F. Groenen, T. Hastie, A. I. D'Enza, A. Markos, and E. Tuzhilina, "Principal component analysis," *Nature Reviews Methods Primers*, vol. 2, no. 1, p. 100, 2022.
- H. Sun, H. Liu, Y. Ma, and Q. Xia, "Optical Remote Sensing Indexes of Soil Moisture: Evaluation and Improvement Based on Aircraft Experiment Observations," *Remote Sensing*, vol. 13, no. 22, p. 4638, 2021.

- S. Ganesan and D. Subramani, "Spatio-temporal predictive modeling framework for infectious disease spread," *Scientific Reports*, vol. 11, no. 1, p. 6741, 2021.
- S. K. Paul, S. Jana, and P. Bhaumik, "A Multivariate Spatiotemporal Model of COVID-19 Epidemic Using Ensemble of ConvLSTM Networks," *Journal of The Institution of Engineers (India): Series B*, vol. 102, no. 6, pp. 1137–1142, 2021.
- A. M. MacEachren, M. Wachowicz, R. Edsall, D. Haug, and R. Masters, "Constructing knowledge from multivariate spatiotemporal data: integrating geographical visualization with knowledge discovery in database methods," *International Journal of Geographical Information Science*, vol. 13, no. 4, pp. 311–334, 1999.
- T. C. Haas, "New systems for modeling, estimating, and predicting a multivariate spatio–temporal process," *Environmetrics*, vol. 13, no. 4, pp. 311–332, 2002.
- N. T. Wimer, L. Esclapez, N. Brunhart-Lupo, M. H. d. Frahan, M. Rahimi, M. Hassanaly, J. Rood, S. Yellapantula, H. Sitaraman, B. Perry, M. Martin, O. Doronina, S. N. Appukuttan, M. Rieth, and M. Day, "Visualizations of a methane/diesel RCCI engine using PeleC and PeleLMeX," *Physical Review Fluids*, vol. 8, no. 11, p. 110511, 2023.
- M. Guarin, R. Faelens, A. Giusti, N. D. Croze, M. Léonard, D. Cabooter, P. Annaert, P. d. Witte, and A. Ny, "Spatiotemporal imaging and pharmacokinetics of fluorescent compounds in zebrafish eleuthero-embryos after different routes of administration," *Scientific Reports*, vol. 11, no. 1, p. 12229, 2021.

- F. Klingelhuber, S. Frendo-Cumbo, M. Omar-Hmeadi, L. Massier, P. Kakimoto, A. J. Taylor, M. Couchet, S. Ribicic, M. Wabitsch, A. C. Messias, A. Iuso, T. D. Müller, M. Rydén, N. Mejhert, and N. Krahmer, "A spatiotemporal proteomic map of human adipogenesis," *Nature Metabolism*, vol. 6, no. 5, pp. 861–879, 2024.
- D. J. Higham, "Modeling and Simulating Chemical Reactions," *SIAM Review*, vol. 50, no. 2, pp. 347–368, 2008.
- L. D. Owen, W. Ge, M. Rieth, M. Arienti, L. Esclapez, B. S. Soriano, M. E. Mueller, M. Day, R. Sankaran, and J. H. Chen, "PeleMP: The Multiphysics Solver for the Combustion Pele Adaptive Mesh Refinement Code Suite," *Journal of Fluids Engineering*, vol. 146, no. 4, 2024.
- C. Tosh, P. Greengard, B. Goodrich, A. Gelman, A. Vehtari, and D. Hsu, "The Piranha Problem: Large Effects Swimming in a Small Pond," *Notices of the American Mathematical Society*, vol. 72, no. 01, p. 1, 2025.
- H. Mandler and B. Weigand, "A review and benchmark of feature importance methods for neural networks," *ACM Computing Surveys*, vol. 56, no. 12, pp. 1–30, 2024.
- N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, "A Survey on Bias and Fairness in Machine Learning," *ACM Computing Surveys (CSUR)*, vol. 54, no. 6, pp. 1–35, 2021.

- A. Selvitella, "The ubiquity of the Simpson's Paradox," *Journal of Statistical Distributions and Applications*, vol. 4, no. 1, p. 2, 2017.
- S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, and K.-R. Müller, "Unmasking Clever Hans predictors and assessing what machines really learn," *Nature Communications*, vol. 10, no. 1, p. 1096, 2019.
- H. Lee and S. Chen, "Systematic Bias of Machine Learning Regression Models and Correction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. PP, no. 99, pp. 1–11, 2025.
- C. Cammarota and A. Pinto, "Variable selection and importance in presence of high collinearity: an application to the prediction of lean body mass from multi-frequency bioelectrical impedance," *Journal of Applied Statistics*, vol. 48, no. 9, pp. 1644–1658, 2021.
- T. Parr and J. D. Wilson, "Partial dependence through stratification," *Machine Learning with Applications*, vol. 6, p. 100146, 2021.
- T. Parr, J. Hamrick, and J. D. Wilson, "Nonparametric feature impact and importance," *Information Sciences*, vol. 653, p. 119563, 2024.
- M. G. Brown, M. G. Peterson, I. K. Tezaur, K. J. .Peterson, and D. L. Bull, "Random forest regression feature importance for climate impact pathway detection," *Journal of Computational and Applied Mathematics*, vol. 464, p. 116479, 2025.
- Z. M. Labe and E. A. Barnes, "Comparison of Climate Model Large Ensembles

With Observations in the Arctic Using Simple Neural Networks," *Earth and Space Science*, vol. 9, no. 7, 2022.

- A. Konya and P. Nematzadeh, "Recent applications of AI to environmental disciplines: A review," *Science of The Total Environment*, vol. 906, p. 167705, 2024.
- J. Lao, X. Wang, B. Shi, B. Wang, and Z. Jiao, "Advances in the application of artificial intelligence in environmental science," *Chinese Journal of Nature*, vol. 46, no. 4, pp. 271–280, 2024.