

Automatically Finding Patches Using Genetic Programming

Authors: Westley Weimer, ThanhVu Nguyen, Claire Le Goues,
Stephanie Forrest

Presented by: David DeBonis, Qi Lu, Shuang Yang

Department of Computer Science

University of New Mexico

February 11, 2013

Outline

- 1 Motivation
- 2 Program Representation
- 3 Genetic Operators
 - Mutation
 - Crossover
- 4 Conclusion
- 5 Discussion

Error Correction in Source Code

- Based on positive / negative test cases
- Program isolation
- Repair by insert / delete / swap
- Repeat until a variant passes all tests
- Minimize difference

Example of Errant Code

Listing 1: Euclid's greatest common divisor algorithm

```
/* requires: a >= 0, b >= 0 */
void gcd(int a, int b) {
    if (a == 0) {
        printf("%d", b);
    }
    while (b != 0)
        if (a > b)
            a = a - b;
        else
            b = b - a;
    printf("%d", a);
    exit(0);
}
```

Variants Creation and Representation

- Restricted to code substitutions from other parts
- Mutation / Crossover constrained to area relevant to error
- Abstract Syntax Tree (AST)
- Weighted program path
- Fitness = $\sum_{i=1}^N w_i$ passed test cases

Program Representation

- An *abstract syntax tree (AST)*. It includes all of the statements in the program.
- A *weighted path* through the program, *(statement, weight)*. The weighted path is a list of pairs, each pair contains a statement and a weight based on that statement's occurrences in various testcases.

The algorithm of mutation

Input: Program P to be mutated.

Input: Path $Path_P$ of interest.

Output: Mutated program variant.

```
1: for all  $\langle stmt_i, prob_i \rangle \in Path_P$  do
2:   if  $\text{rand}(0, 1) \leq prob_i \wedge \text{rand}(0, 1) \leq W_{mut}$  then
3:     let  $op = \text{choose}(\{\text{insert}, \text{swap}, \text{delete}\})$ 
4:     if  $op = \text{swap}$  then
5:       let  $stmt_j = \text{choose}(P)$ 
6:        $Path_P[i] \leftarrow \langle stmt_j, prob_i \rangle$ 
7:     else if  $op = \text{insert}$  then
8:       let  $stmt_j = \text{choose}(P)$ 
9:        $Path_P[i] \leftarrow \langle \{stmt_i; stmt_j\}, prob_i \rangle$ 
10:    else if  $op = \text{delete}$  then
11:       $Path_P[i] \leftarrow \langle \{\}, prob_i \rangle$ 
12:    end if
13:  end if
14: end for
15: return  $\langle P, Path_P, \text{fitness}(P) \rangle$ 
```

The example for Mutation

```
1 while (b != 0)
2     if (a > b)
3         a = a - b;
4     else
5         b = b + a;
6 return a;
```


The example for Mutation

```

1 while (b != 0)
2   if (a > b)
3     a = a - b;
4   else
5     b = b + a;
6 return a;

```

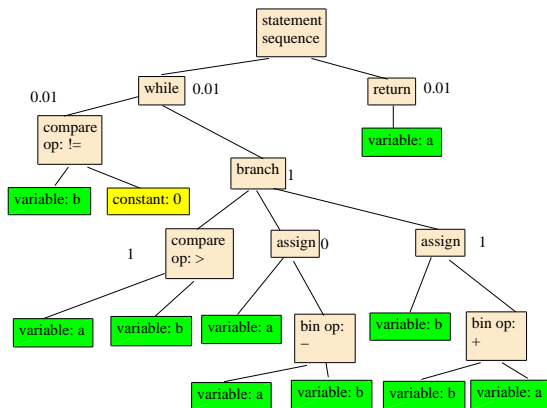
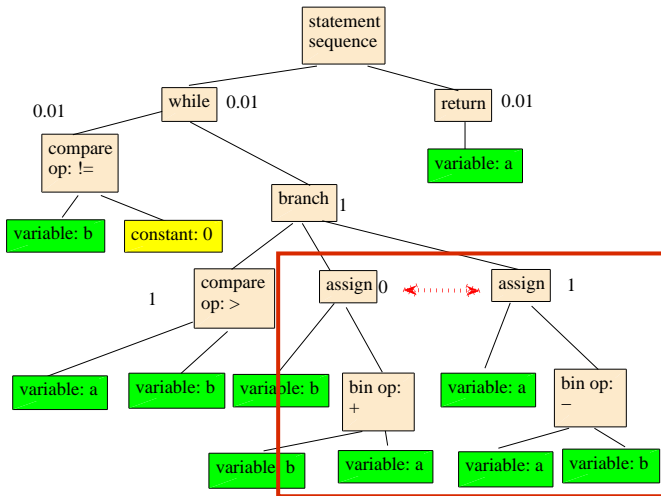


Figure: The AST of the program

Mutation: Swap

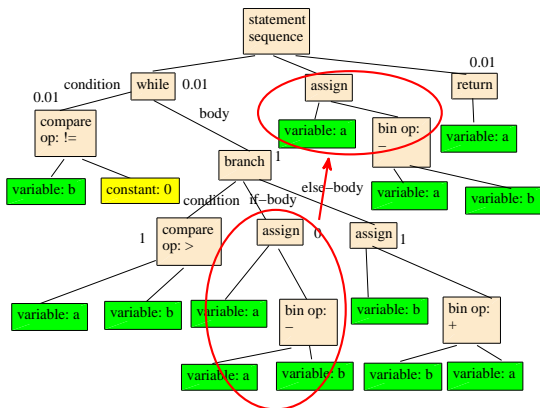


Mutation: Swap

```
1 while (b != 0)
2   if (a > b)
3     a = a - b;
4   else
5     b = b + a;
6 return a;
```

```
1 while (b != 0)
2   if (a > b)
3     b = b + a;
4   else
5     a = a - b;
6 return a;
```

Mutation: Insertion

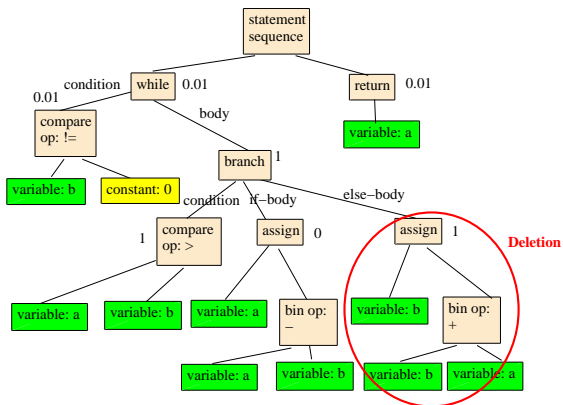


Mutation: Insertion

```
1 while (b != 0)
2   if (a > b)
3     a = a - b;
4   else
5     b = b + a;
6 return a;
```

```
1 while (b != 0)
2   if (a > b)
3     a = a - b;
4   else
5     b = b + a;
6 a = a - b;
7 return a;
```

Mutation: Deletion

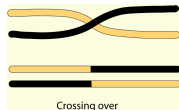


Mutation: Deletion

```
1 while (b != 0)
2   if (a > b)
3     a = a - b;
4   else
5     b = b + a;
6 return a;
```

```
1 while (b != 0)
2   if (a > b)
3     a = a - b;
4 return a;
```

Crossover



- Choose a cutoff point for each program
- Combine the "first part" of one program with the "second part" of another and vice versa
- Input: $[P_1, P_2, P_3, P_4]$ and $[Q_1, Q_2, Q_3, Q_4]$ with cutoff 2
Child: $[P_1, P_2, Q_3, Q_4]$ and $[Q_1, Q_2, P_3, P_4]$

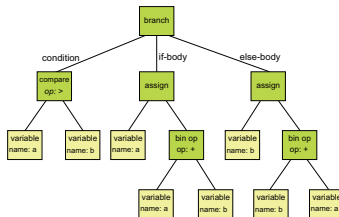
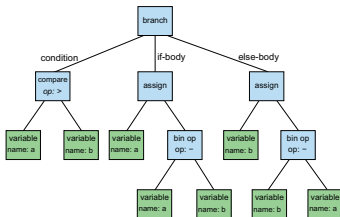
Crossover

```
1  if (a > b)
2    a = a - b;
3  else
4    b = b - a;
```

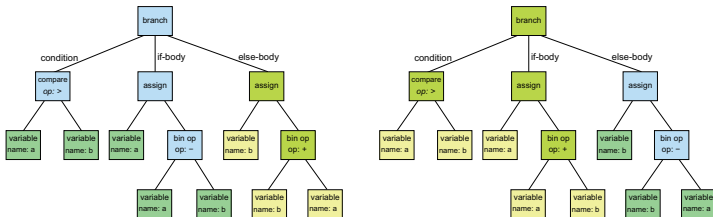
```
1  if (a > b)
2    a = a + b;
3  else
4    b = b + a;
```

- cutoff point at 4

Crossover



Crossover



Crossover

```
1  if (a > b)
2    a = a - b;
3  else
4    b = b + a;
```

```
1  if (a > b)
2    a = a + b;
3  else
4    b = b - a;
```

Crossover

Input: Parent programs P and Q .

Input: Paths $Path_P$ and $Path_Q$.

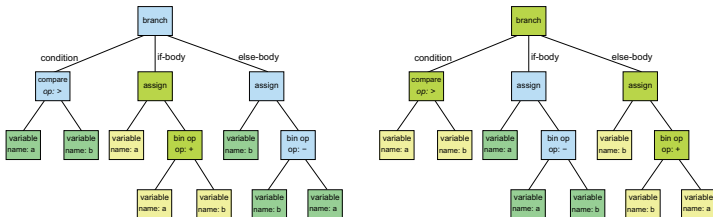
Output: Two new child program variants C and D .

```

1:  $cutoff \leftarrow \text{choose}(|Path_P|)$ 
2:  $C, Path_C \leftarrow \text{copy}(P, Path_P)$ 
3:  $D, Path_D \leftarrow \text{copy}(Q, Path_Q)$ 
4: for  $i = 1$  to  $|Path_P|$  do
5:   if  $i > cutoff$  then
6:     let  $\langle stmt_p, prob \rangle = Path_P[i]$ 
7:     let  $\langle stmt_q, prob \rangle = Path_Q[i]$ 
8:     if  $\text{rand}(0, 1) \leq prob$  then
9:        $Path_C[i] \leftarrow Path_Q[i]$ 
10:       $Path_D[i] \leftarrow Path_P[i]$ 
11:    end if
12:  end if
13: end for
14: return  $\langle C, Path_C, \text{fitness}(C) \rangle, \langle D, Path_D, \text{fitness}(D) \rangle$ 

```

Crossover



Conclusion

- Demonstrated GP approach
 - Based on AST representation
 - Mutate / Crossover algorithm
 - Minimal code changes
- Experiments over 10 test cases
 - Sizeable code bases
 - Orthogonal errors present
 - Reasonable solution found (50%)
- Efficacy to code maintenance
 - Experiments using Amazon Cloud
 - Equated to under \$8 per bug on average
 - Who wants to manually debug anyway!

Discussion

- Limitations and Assumptions
 - Only as good as your test cases
 - Not necessarily memory / computing time optimized
 - Expects redundant code segments
- How well does it scale?
 - by number of bugs...
 - by number of revisions...
 - by number of LOC...
- Will this technique work on another GP?