

Automated Program Repair through the Evolution of Assembly Code

Eric Schulte

University of New Mexico

08 August 2010

Introduction

We present a method of *automated program repair* through the evolution of *assembly code* compiled from *extant software*.

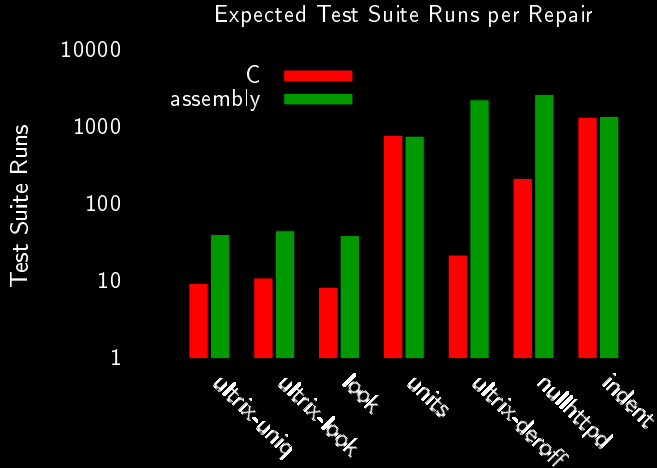
- previous work demonstrated the repair of C programs through evolution of C statements
- we extend previous work by operating at the level of assembly commands

Benefits

Benefits of the assembly code level of representation include

- General** applicable to any language which compiles to Java byte code or x86 assembly code
- Expressive** the small scale of assembly instructions is capable of expressing repairs not possible at the C statement level
- Coverage** the small alphabet of assembly commands and large number of commands in assembly programs provides access to a larger subset of the space of possible programs
- Robust** the functionality of assembly programs is more robust to mutation

Comparative Performance



Repair at the Assembly level is roughly 17% slower than at the C level.

Technical Approach

Stages

preprocessing

fault localization

mutation

evolution

Technical Approach

Stages

C Statement Tree

Buggy C Program

preprocessing

fault localization

mutation

evolution

Technical Approach

Stages

C Statement Tree

Buggy C Program

CIL

CIL intermediate

preprocessing

fault localization

mutation

evolution

Technical Approach

Stages

C Statement Tree

Buggy C Program

CIL

CIL intermediate

walk path

weighted C statement tree

preprocessing

fault localization

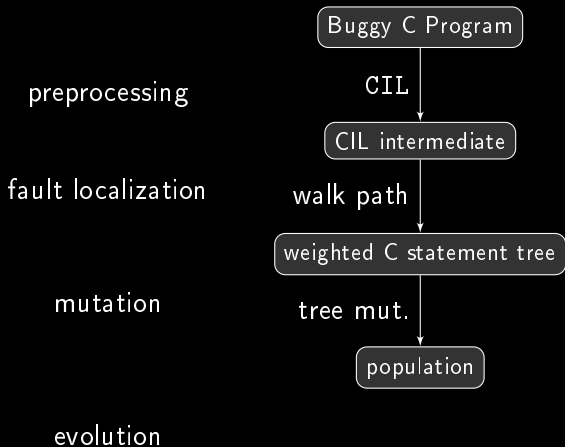
mutation

evolution

Technical Approach

Stages

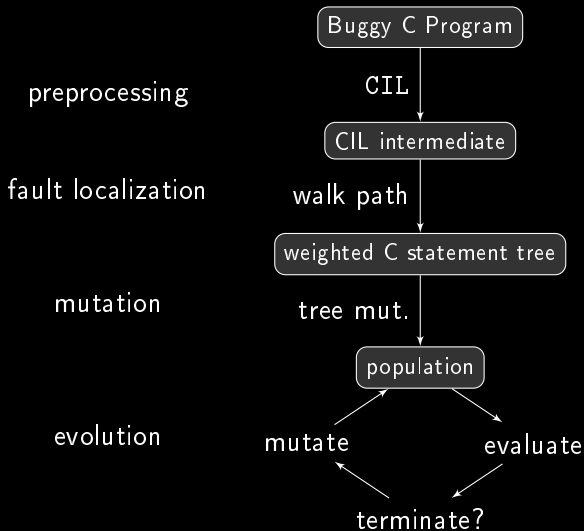
C Statement Tree



Technical Approach

Stages

C Statement Tree



Technical Approach

Stages

C Statement Tree

Assembly Linear Genome

preprocessing

Buggy C Program

Buggy C, Java, etc... Program

CIL

CIL intermediate

fault localization

walk path

weighted C statement tree

mutation

tree mut.

population

evolution

mutate

evaluate

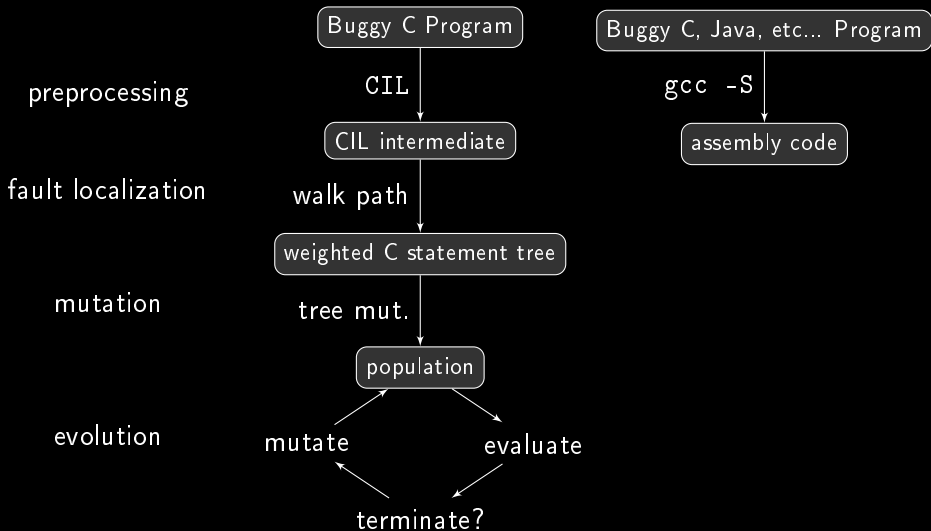
terminate?

Technical Approach

Stages

C Statement Tree

Assembly Linear Genome



Technical Approach

Stages

C Statement Tree

Assembly Linear Genome

preprocessing

Buggy C Program

Buggy C, Java, etc... Program

CIL

gcc -S

CIL intermediate

assembly code

fault localization

walk path *deterministic*

sample pc *stochastic*

weighted C statement tree

mutation

tree mut.

population

evolution

mutate

evaluate

terminate?

Technical Approach

Stages

C Statement Tree

Assembly Linear Genome

preprocessing

Buggy C Program

Buggy C, Java, etc... Program

CIL

gcc -S

CIL intermediate

assembly code

fault localization

walk path

sample pc

weighted C statement tree

weighted assembly genome

mutation

tree mut.

population

evolution

mutate

evaluate

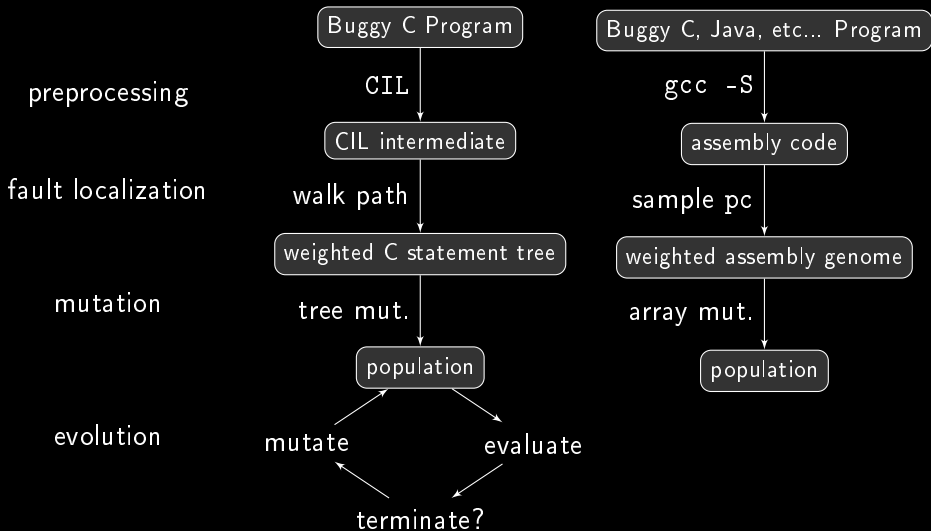
terminate?

Technical Approach

Stages

C Statement Tree

Assembly Linear Genome

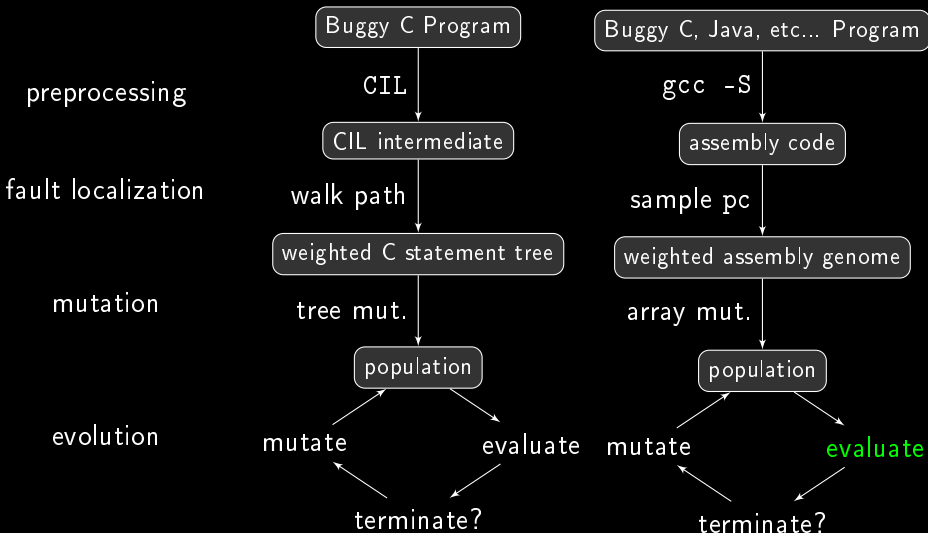


Technical Approach

Stages

C Statement Tree

Assembly Linear Genome



Representation and Genetic Operators

Linear Genome

Individuals are *linear genomes* of weighted assembly commands.

```
("main:" "pushl %ebp" "movl %esp" ... )
```

Genetic Operators

insert selects an instruction, selects a location, copies the instruction and inserts in the location

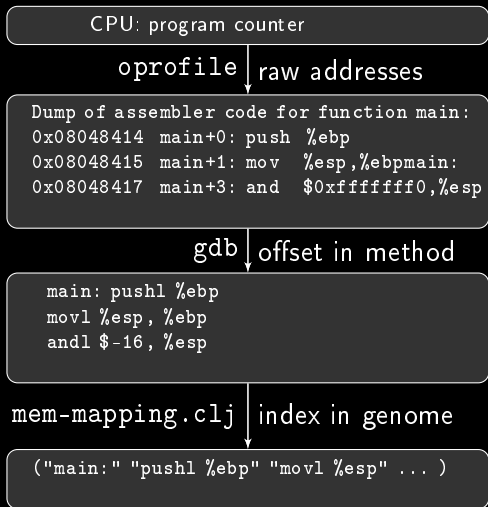
delete selects an instruction and deletes it

swap selects two instructions and swaps them

crossover selects a crossover point for each of two individuals and exchanges all instructions after that point between the individuals

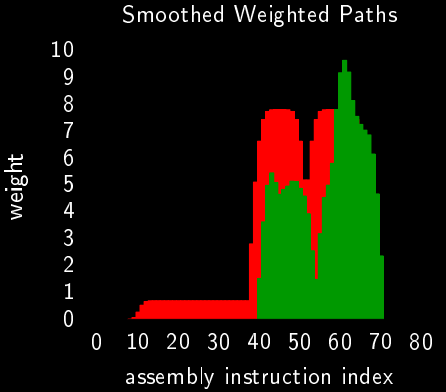
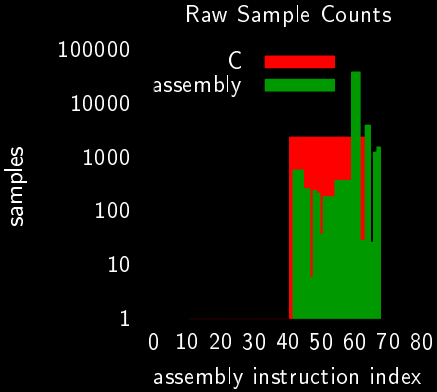
Fault Localization

Technique



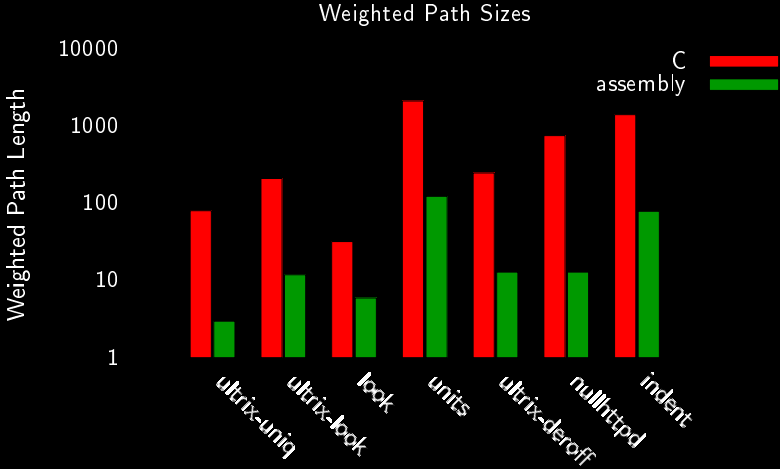
Fault Localization

Results



Fault Localization

Comparative Path Size



Generality to Multiple Languages

Input: Integer a

Input: Integer b

Output: $\gcd(a, b)$ or \perp

```
1: if  $a \equiv 0$  then  
2:   print  $a$   
3: end if  
4: while  $b \neq 0$  do  
5:   if  $a > b$  then  
6:      $a \leftarrow a - b$   
7:   else  
8:      $b \leftarrow b - a$   
9:   end if  
10: end while  
11: print  $a$ 
```

	C	Haskell	Java
Program Length	79	885	33
Total Solutions	2	15	13
Unique Solutions	2	10	1

Table: GCD Repair Results by Language

Figure: A Buggy version of Euclid's Algorithm

Mutational Robustness & Neutral Spaces

Concepts

Mutational Robustness robustness of *phenotype* to changes in *genotype*
in the case of computer programs

phenotype behavior of the program

genotype representation of the program

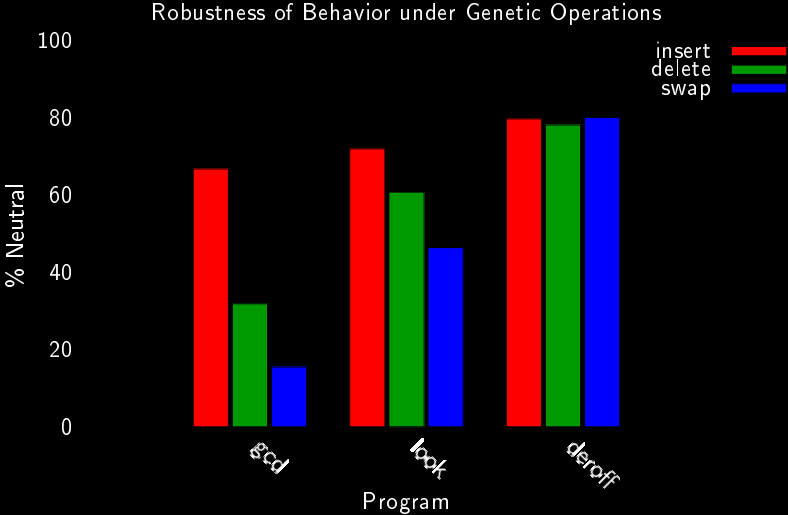
Neutral Space contiguous region of *genotype space* with constant
phenotype

Mutational Robustness & Neutral Spaces

Relevance

- Mutational Robustness directly influences evolvability
- A large neutral space allows for significant diversity in
 - representation
 - behavior
 - efficiency
 - size
 - etc. . .
- A *neutral neighborhood* defines which functionalities are reachable in a single genotypic step

Robustness under Genetic Operators



Future Work

Investigate Mutational Robustness

- differences across *languages*, *algorithms* and *representations*
- properties of the *neutral spaces* of assembly programs

Improve GP Technique

- apply ongoing work on the C statement level
- homologous crossover, steady state mutation, etc. . .

Non-Repair Evolution

- machine specific optimization
- disruption of mono-culture (N-variant)

Discussion

- Questions
- Comments
- Suggestions