

Inverse Preprocessor

Yiming Wu

October 10, 2014

1 Introduction

Program repair has long been an essential problem with heated discussion, since most software systems inevitably have bugs that need to be fixed. Developers for large software projects must confirm, triage, and localize defects before fixing these bugs and validating the fixes. To reduce human manual labor, numbers of automatic patch generation techniques or tools have been proposed.

Among these automatic tools, C oriented fixing tools count most of them because of the prevalence of C in open-source or enterprise software projects. Weimer et al. proposed a population-based technique and developed **GenProg**, leveraging genetic programming[1]. Nguyen et al. presented **SemFix** based on symbolic execution, constraint solving and program synthesis[2]. Kaleeswaran et al. exploited statistical correlation analysis to identify perspective bugs[3]. Qi et al. provided a prototype tool **RsRepair** using random search[4].

However, all C oriented automatic repair tools face a common problem: preprocessing. Although preprocessing is not explicitly stated in their papers, automatic tools like **GenProg** abstract the C program into an AST which might include preprocessing. It is impossible to know the whole picture of a C program without preprocessing to get rid of predefined macro and conditional instructions. Developers for automatic repair tools have to write a qualified preprocessor before they can truly start core analysis of the program. This will cost meaningless labor and increase chances of bugs within bug-repair tools themselves.

Furthermore, patches generated by automatic repair tools may confuse their users because these patches are targeted for preprocessed code. (Need an example here where patches loses macro name).

In this paper, we present a novel bidirectional preprocessor **BXP** that can tackle preprocessing problem for automatic repair tools. Our bidirectional preprocessor performs like a standard preprocessor for automatic repair tools to get preprocessed code. It can also generate patches for *unpreprocessed* code according to the patches for *preprocessed* code from automatic program repair tools. With the help of **BXP**, we believe the develop of automatic C program repair tools can be greatly simplified. In addition, the patches generated by repair tools are now more agreeable and graceful.

(This paragraph is left for evaluation)

(This paragraph is left for contribution list)

2 MEDG: Macro Expansion Dependency Graph

To retrieve macro calls in original code, Inverse Preprocessor needs the information of all the macro expansions in the code. During preprocessing, our preprocessor PP records those information with the help of MEDG: Macro Expansion Dependency Graph.

To grab the main idea of MEDG, let us start with a quick example. In the example code, there are four macro definitions and one macro call. Preprocessor PP takes in original example code and produces preprocessed code as well as MEDG.

```
#define g(a,b) a+b
#define h(a) a*a
#define NUM 3
#define f(a,b) g(a,b)+h(a)+NUM

...
f(a,b)
...
```

Figure 1: A Quick Example Code

Like standard preprocessor, our Preprocessor PP iterates through the whole context over and over until there is no more macro expansion. Each time Preprocessor PP expands a macro call, it will record dependency connection edge between macro-call token and expanding tokens. Each connection edge is encoded with three domain: macro name, parameter list and expansion position in original code. These information help Inverse Preprocessor getting the original code back quickly and elegantly.

The figure 3 shows the MEDG of the quick example, it goes in two iterations. In the first iteration,

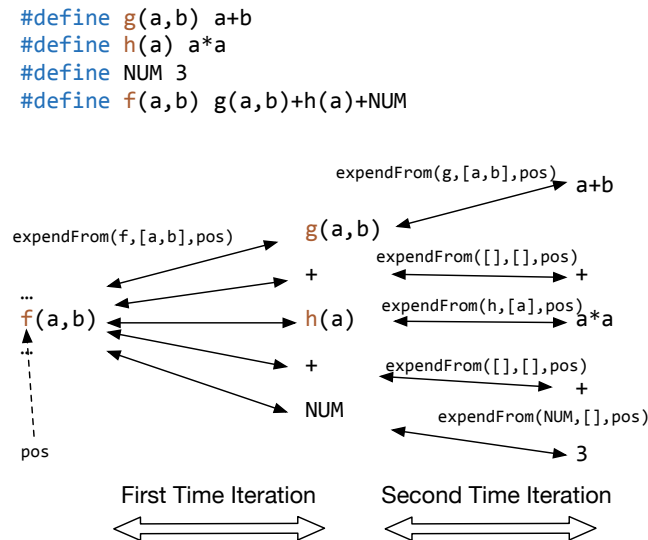


Figure 2: MEDG of Quick Example.

Here is our pseudo code of Preprocessor PP.

Algorithm 1 Preprocessor

Input: Original Code**Output:** Preprocessed Code with MEDG1: **function** ITERATEONCE(*ctx*)2: **end function**

3 Deal with Changes on Code with MEDG

In this section, we will show how to deal with changes on preprocessed code and transform the changes back after we record macro call information with MEDG. We will first give out an example, then discuss the pseudo code for dealing with changes on code.

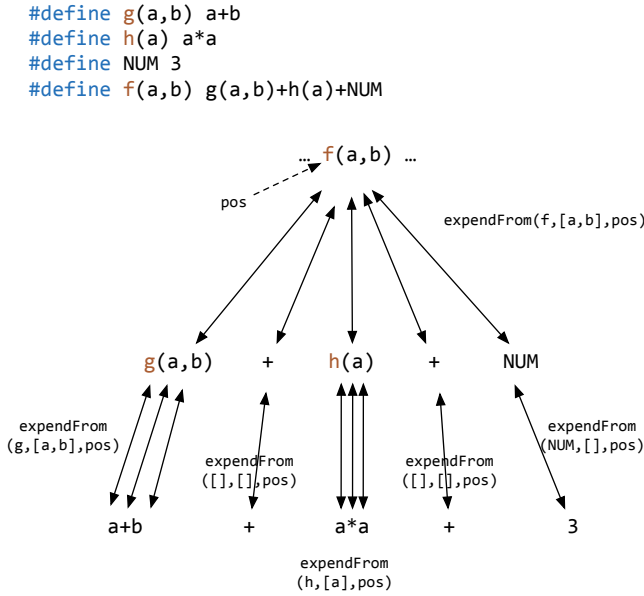


Figure 3: MEDG without Changes.

```

#define g(a,b) a+b
#define h(a) a*a
#define NUM 3
#define f(a,b) g(a,b)+h(a)+NUM

```

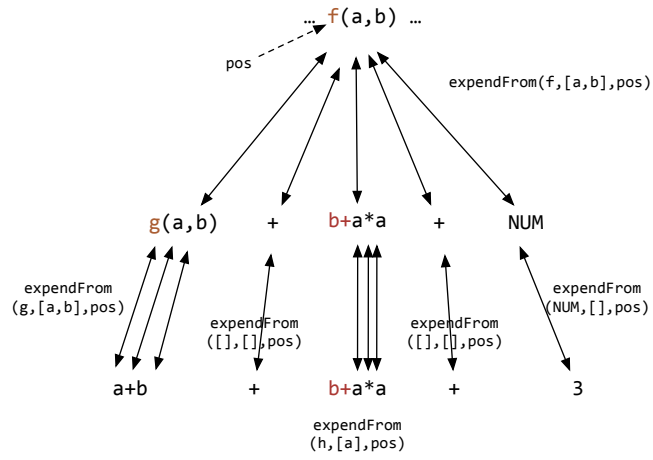


Figure 4: Adding Tokens to the Preprocessed Code.

Algorithm 2 a

Input: a**Output:** a

```
1: function MERGESORT(Array, left, right)
2:   result  $\leftarrow$  0
3:   if left < right then
4:     middle  $\leftarrow$  (left + right)/2
5:     result  $\leftarrow$  result + MERGESORT(Array, left, middle)
6:     result  $\leftarrow$  result + MERGESORT(Array, middle, right)
7:     result  $\leftarrow$  result + MERGER(Array, left, middle, right)
8:   end if
9:   return result
10: end function
11:
12: function MERGER(Array, left, middle, right)
13:   i  $\leftarrow$  left
14:   j  $\leftarrow$  middle
15:   k  $\leftarrow$  0
16:   result  $\leftarrow$  0
17:   while i < middle and j < right do
18:     if Array[i] < Array[j] then
19:       B[k ++]  $\leftarrow$  Array[i ++]
20:     else
21:       B[k ++]  $\leftarrow$  Array[j ++]
22:       result  $\leftarrow$  result + (middle - i)
23:     end if
24:   end while
25:   while i < middle do
26:     B[k ++]  $\leftarrow$  Array[i ++]
27:   end while
28:   while j < right do
29:     B[k ++]  $\leftarrow$  Array[j ++]
30:   end while
31:   for i = 0  $\rightarrow$  k - 1 do
32:     Array[left + i]  $\leftarrow$  B[i]
33:   end for
34:   return result
35: end function
```
