# Detecting Code Plagiarism on Student Submissions

Yanyan Jiang and Chang Xu



SPAR Group, Institute of Computer Software
Department of Computer Science and Technology, Nanjing University

May 19, 2018

# Needle[1]: Detecting Code Plagiarism on Student Submissions

---

[1]Needle is a utility project in Project-N ◈, the computer systems lab series at Nanjing University: processors (NPC, NOOP), emulator (NEMU), operating system (Nanos), and compiler (NCC).

# Outline

# Students Copy Code[2]

| T1 | change comments, names, or cases | `cur->lineno` $\rightarrow$ `head->line` |
|----|----------------------------------|------------------------------------------|
| T2 | reformat or reorder code fragments | indention/code style change, etc. |
| T3 | add or delete redundant elements | add unused variables or redundant computations |
| T4 | refactor program constructs | `for` $\rightarrow$ `while`, function split/merge, function rewrites, ... |

---

[2]Neal R. Wagner. Plagiarism by student programmers. UTSA Tech Rep, 2000.

# Students Copied Lots of Code

## Without plagiarism control

A conservative post-mortem analysis (a single anonymous lab) showed that

- ~82% students plagiarized
- ~42% of the plagiarized copies are of similarity > 99%.

## With *Needle* ← this talk's topic, *penalties*, and *honor grades*

Reduced plagiarism rate

- (average) ~5% project groups plagiarized in a single lab
- ~20% project groups plagiarized in a semester

# Outline

## Problem Definition

Calculate the (minimum) editing distance between two *optimized* program binaries[3] $P_1$ and $P_2$:

1. Reordering of two functions with zero cost
2. Deletion of an instruction costing $c_d$
3. Modification of an instruction (either instruction type or operands) costing $c_m$
4. Insertion of an arbitrary instruction costing $c_i$.[4]

Programs with small $d(P_1, P_2)$ are likely to be plagiaristic copies
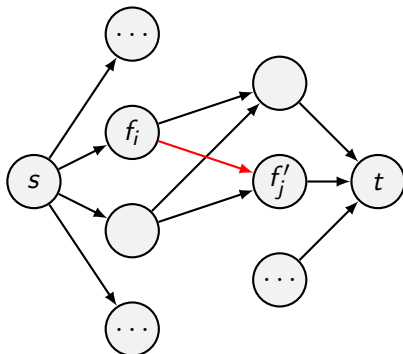
- Unfortunately, computing $d(P_1, P_2)$ is NP-Complete

---

[3]Compiler optimizations are good at "normalizing" local semantics changes.

[4]$c_d = c_m = c_i = 1$ suffices for code plagiarism detection in practice.

# The Network Flow Solution

*"A lower-bound estimation of the editing distance $d(P_1, P_2)$"*



Let each unit of flow $(l_i, r_j)$ to denote "embedding an instruction in function $f_i$ to function $f'_j$". The similarity between $P_1$ and $P_2$ is denoted by the *maximum weighted matching*.

Similarity between functions (windowed longest-common-subsequence):

$$\sigma(f_i, f_j') = \max_{k \in \{1,2,\ldots,|f_j'|\}} \text{LCS}(f_i, f_j'[k : k + \omega])$$

A weighted bipartite $G(L, R, c, w)$
- Capacity $c(\ell_i, r_j) = \sigma(f_i, f_j')$
  - cannot embed too many instructions from $f_i \rightarrow f_j'$
- Weight $w(\ell_i, r_j) = \left( 1 + e^{-\alpha \cdot \frac{\max\{\sigma(f_i, f_j'), \sigma(f_j', f_i)\}}{\min\{|f_i|, |f_j'|\}} + \beta} \right)^{-1}$
  - embedding $f_i \rightarrow f_j'$ is more profitable with a larger $\sigma(f_i, f_j')$

Program Similarity $\sigma(P_1, P_2) = \dfrac{\text{MaximumWeightFlow}(G, c, w)}{\sum_{i \in [n]} |f_i|} \in [0, 1]$

# Why It Works?

| | | |
|---|---|---|
| T1 | change comments, names, or cases | do not affect the compiled binary |
| T2 | reformat or reorder code fragments | matching is order insensitive |
| T3 | add or delete redundant elements | instructions can be aligned by LCS |
| T4 | refactor program constructs | limited changes to the optimized binary; function split/merge can be detect by matching |

# Outline

# An Anonymous Programming Assignment: Case Study

In all 79 students' submissions

- 65 of 79 (82%) submissions are plagiaristic (by a manual inspection)
- 42 of 79 (398 pairs) simiarlity > 99% ("nearly perfect embedding")
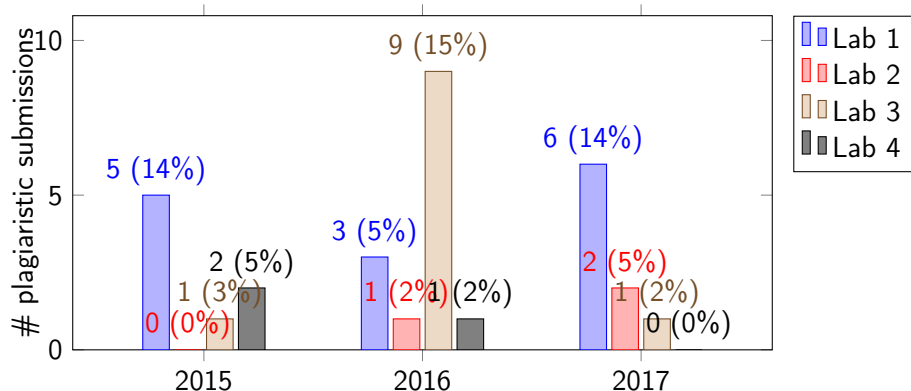
## Interview Results

The teaching assistant: "...Looking at their lab reports, I knew that many plagiarized. However, it would be impossible to handle such many cases ..."

A student: "...There are some self-motivated students who worked on the assignment. However, it's too easy to get a working copy from the Internet. Some simply copied it, and the others submitted modified versions ..."

- The "Principles and Techniques of Compilers" course
    - C-- $\rightarrow$ IR $\rightarrow$ MIPS32 assembly

- Code plagiarism control
    - let the students know (important)
    - use both tools (Needle, and a faster tool for cross-semester)
    - manual inspection and interview

# Cheater's Code Examples

```
1  cur->lineno = temp->lineno;
2  strcpy(cur->type, type);
3  cur->isLexical = 0;
4  cur->children = temp;
```

```
head->number_signal = 0;          1
head->line = temp->line;          2
strcpy(head->type,type);          3
head->child_left = temp;          4
```

```
1  $$->is_root=1;
2  $$->no_leaves=1;
3  $$->leaves[0]=(Node*)$1;
4  if(exit_error==0)
5    {print_tree($$,0);}
```

```
$$->final=0;                      1
$$->num_children=1;               2
$$->children=(Node**)malloc       3
  (sizeof(Node*)*$$->             4
     num_children);
$$->children[0]=(Node*)$1;        5
if(!wrong) printNode($$,0);       6
```

```
1  temp->line = a->line;
2  temp->lChild = a;
3  while(num > 1){
4    a->rChilds = va_arg(list,
        node*);
5    a = a->rChilds;
6    num--;
7  }
```

```
p_node->left_child = temp;        1
p_node->line = temp->line;        2
for(int i=0;i < num-1;++i){       3
  temp->right_child =             4
    va_arg(valist, struct         5
       Node*);
  temp = temp->right_child;       6
}                                 7
```

A piece of code to trick code plagiarism detectors

```
1  ...
2  draw_string("Game␣Over");
3  ...
```

```
...                          1
draw_string("G");            2
draw_string("A");            3
draw_string("M");            4
draw_string("E");            5
draw_string("␣");            6
draw_string("O");            7
draw_string("V");            8
draw_string("E");            9
draw_string("R");            10
...                          11
```

# Outline

- Plagiarism is a seductive way for students to obtain unethical grades
  - ~20% of the project groups plagiarized in the study
  - students attacked the plagiarism detection tool
  - cross-semester plagiarism (~30%) is also common

- Using plagiarism control policies *alleviate* but not *eliminate* the issue
  - there is a heavy burden on the instructors to conduct code interviews

# Thank You!