# MySQL Dependency Extraction

**Tabs vs. Spaces**

Varsha Ragavendran          Davood Anbarnam          Glib Sitiugin

Nisha Sharma          Ayman Abualsunun          Kevin Arindaeng

Anton Sitkovets

# Overview

- Introduction
- Compare dependency extraction techniques  for the following:
  - Understand
  - srcML
  - Our Program
- Quantitative & Qualitative Analysis for the following:
  - Understand vs. srcML
  - srcML vs. Our Program
  - Understand vs. Our Program
- Potential Risks and Limitations
- Lessons Learned
- Conclusion

# Introduction

- Dependency extraction techniques help to:
  - Analyze large open source code architecture to optimize software design
  - Easily identify  dependencies within the code
  - Understand complex code bases with poor documentation
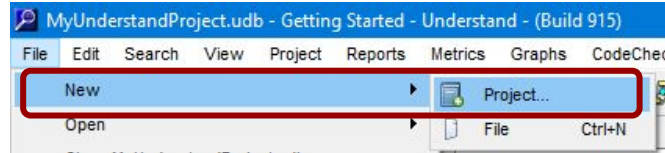
# Dependency extraction using Understand

Understand is a static analysis tool focused on:

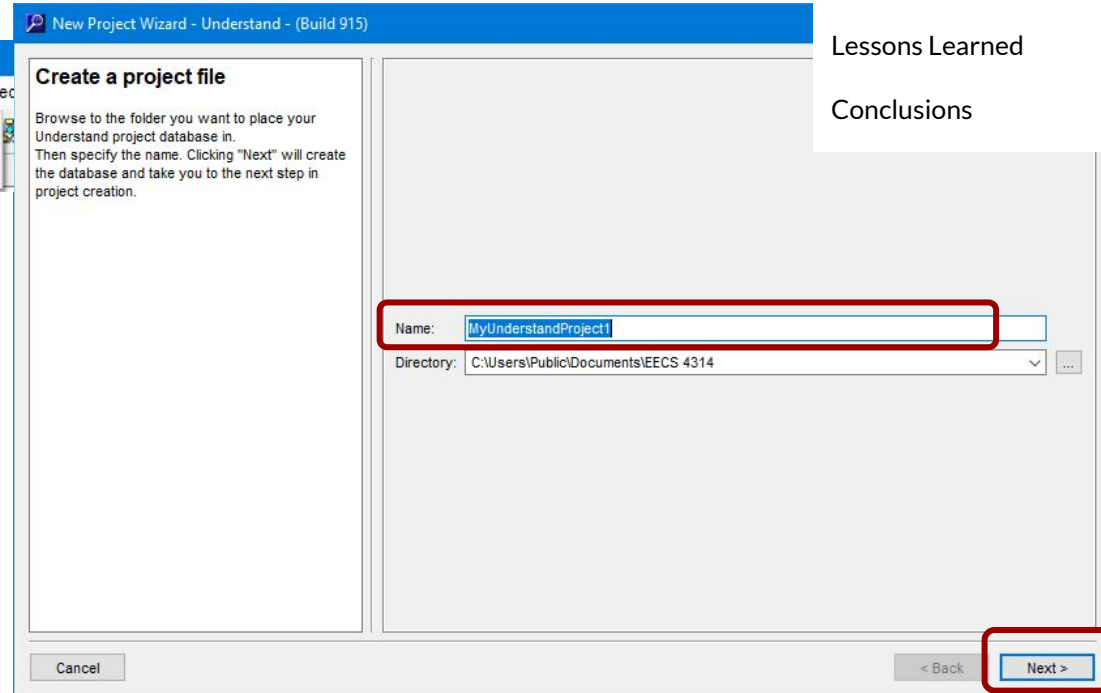- source code comprehension, metrics, and standard testing.

A file can depend on another through:

- Import

- Inheritance

- Implementation

- Method calls and object initializations

- @ Java annotations

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Dependency extraction using Understand



Create a new project
And Name it

# Dependency extraction using Understand

Select the Source Code Languages

# Dependency extraction using Understand

Import Files

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Dependency extraction using Understand



Add Directory

# Dependency extraction using Understand



9

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Dependency extraction using Understand



Export

# Dependency extraction using srcML

```
You're back!@Falalfel MINGW64 /d/Uni Stuff/EECS 4314/Assignment3
$ srcml mysql-server-mysql-8.0.2 -o mysql.xml
```

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Dependency extraction using srcML

.

```
<cpp:include>#<cpp:directive>include</cpp:directive> <cpp:file>&lt;stdio.h&gt;</cpp:file></cpp:include>

<function><type><name>int</name></type> <name>main</name><parameter_list>(<parameter><decl><type><name>i
<block>{
  <decl_stmt><decl><type><name>int</name></type> <name>i</name></decl>;</decl_stmt>
  <for>for <control>(<init><expr><name>i</name><operator>=</operator> <literal type="number">1</literal>
  <block>{
    <expr_stmt><expr><call><name>fprintf</name><argument_list>(<argument><expr><name>stdout</name></expr
    <if>if <condition>(<expr><name>i</name> <operator>&lt;</operator> <name>argc</name> <operator>-</ope
      <block type="pseudo"><expr_stmt><expr><call><name>fprintf</name><argument_list>(<argument><expr><n
  }</block></for>
  <expr_stmt><expr><call><name>fprintf</name><argument_list>(<argument><expr><name>stdout</name></expr><
  <return>return <expr><literal type="number">0</literal></expr>;</return>
}</block></function>
```
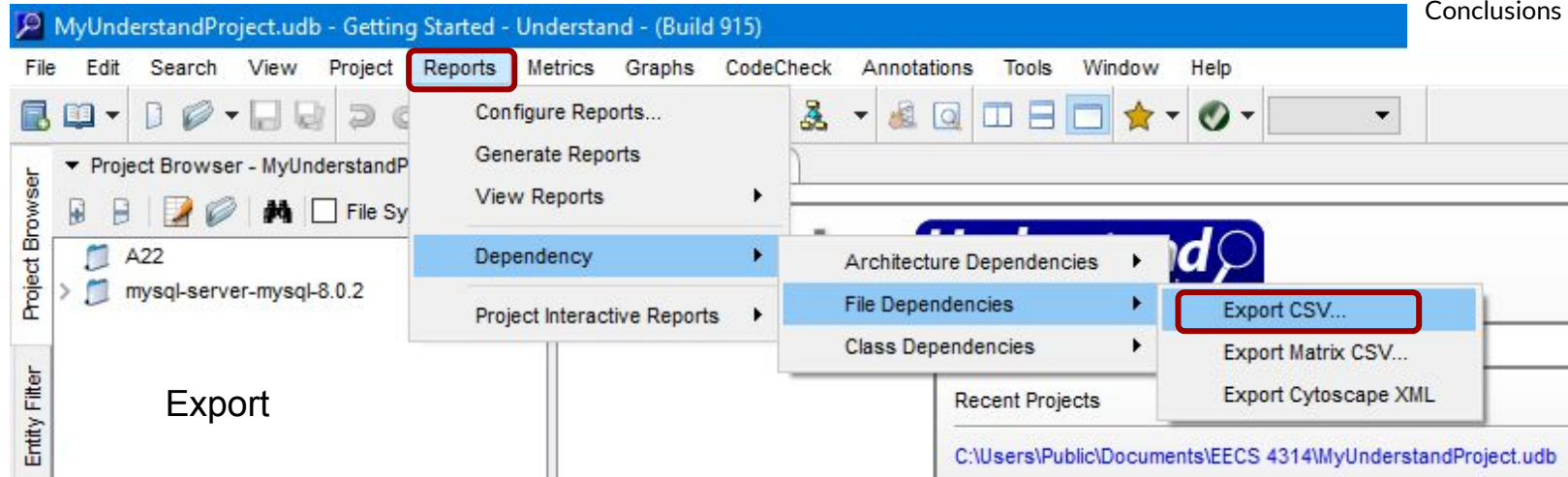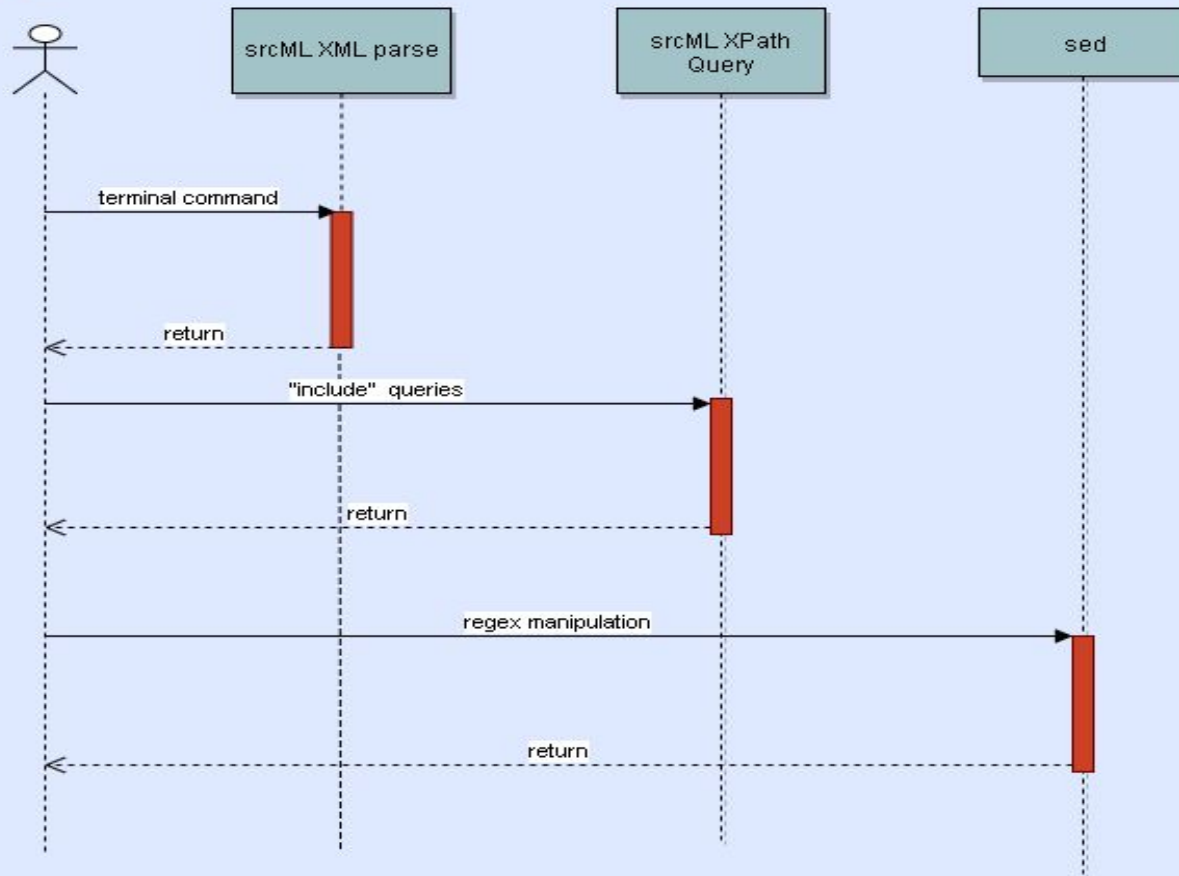
12

# Dependency Extraction using Program

Wrote a Python script that walks through the MySQL source code folder.

Checks every ".cc", ".c" and ".h" file. Parses each line of the file and writes the dependency between the current file and the files being included using "#include" to a raw.ta file.

https://github.com/azkevin/EECS4314/blob/master/A3/a3data/include.py

https://github.com/azkevin/EECS4314/blob/master/A3/a3data/test.raw.ta

```python
import os
import sys

file_ta = open(sys.argv[2], "w")

for root, dirs, files in os.walk(sys.argv[1], topdown=False):
    for name in files:
        if name[-3:] == ".cc" or name[-2:] == ".h" or name[-2:] == ".c":
            lines = open(os.path.join(root, name), "r")
            for line in lines:
                if line[:8] == "#include":
                    left_dependency = os.path.join(root, name).replace('\\', '/')[os.path.join(root, name).find("mysql-server-mysql-8.0.2"):]
                    string = "{} -> {}\n".format(left_dependency, line[10:-2])
                    if string.rfind('"') != -1:
                        string = string[0:string.rfind('"')] + "\n"
                    if (string.rfind('>') != -1) & (string.count('>') > 1):
                        string = string[0:string.rfind('>')] + "\n"
                    if string[string.find('>')+2:-1].find('/') > -1:
                        string = "{} -> {}\n".format(left_dependency, os.path.basename(os.path.normpath(string[string.find('>')+2:-1])))
                    file_ta.write(string)
file_ta.close()
```
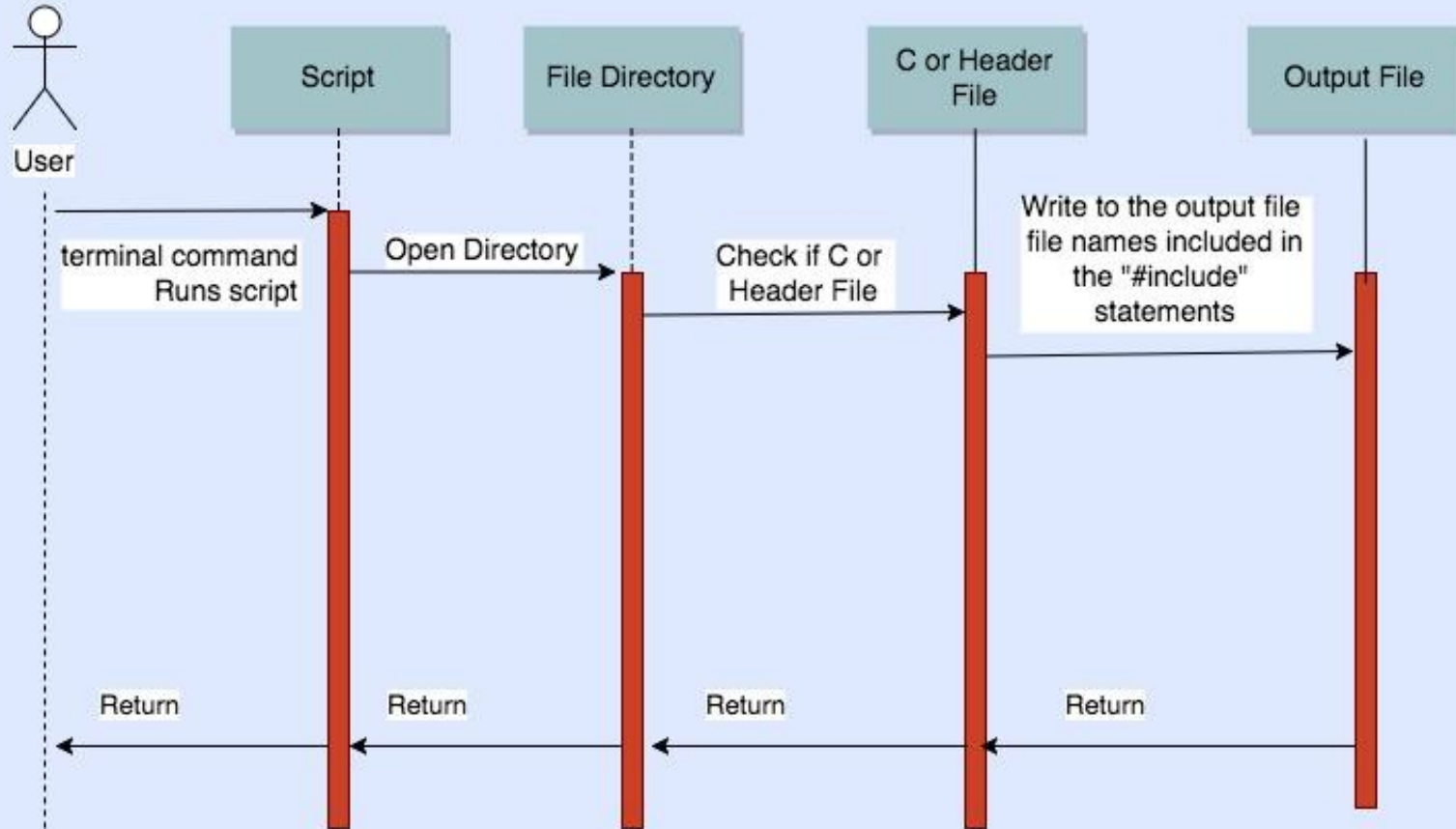
14

**Find Dependencies using Include Directives in Source Code**

15

# Finding Common Dependencies

```
public static void main(String[] args) throws FileNotFoundException, IOException
{
    String file1 = args[0];
    String file2 = args[1];
    String output_file = args[2];

    Scanner file1_scanner = new Scanner(new File(file1));
    PrintWriter writer = new PrintWriter(output_file, "UTF-8");

    while(file1_scanner.hasNextLine())
    {
        String line_f1 = file1_scanner.nextLine();
        if(!line_f1.isEmpty())
        {
            Scanner file2_scanner = new Scanner(new File(file2));
            boolean found = false;
            while(file2_scanner.hasNextLine() && !found)
            {
                String line_f2 = file2_scanner.nextLine();
                //System.out.println("F1 : " + line_f1 + "  F2 : " + line_f2);
                if(line_f1.equals(line_f2))
                {
                    writer.println(line_f1);
                    found = true;
                }
            }
            file2_scanner.close();
        }
    }
}
```
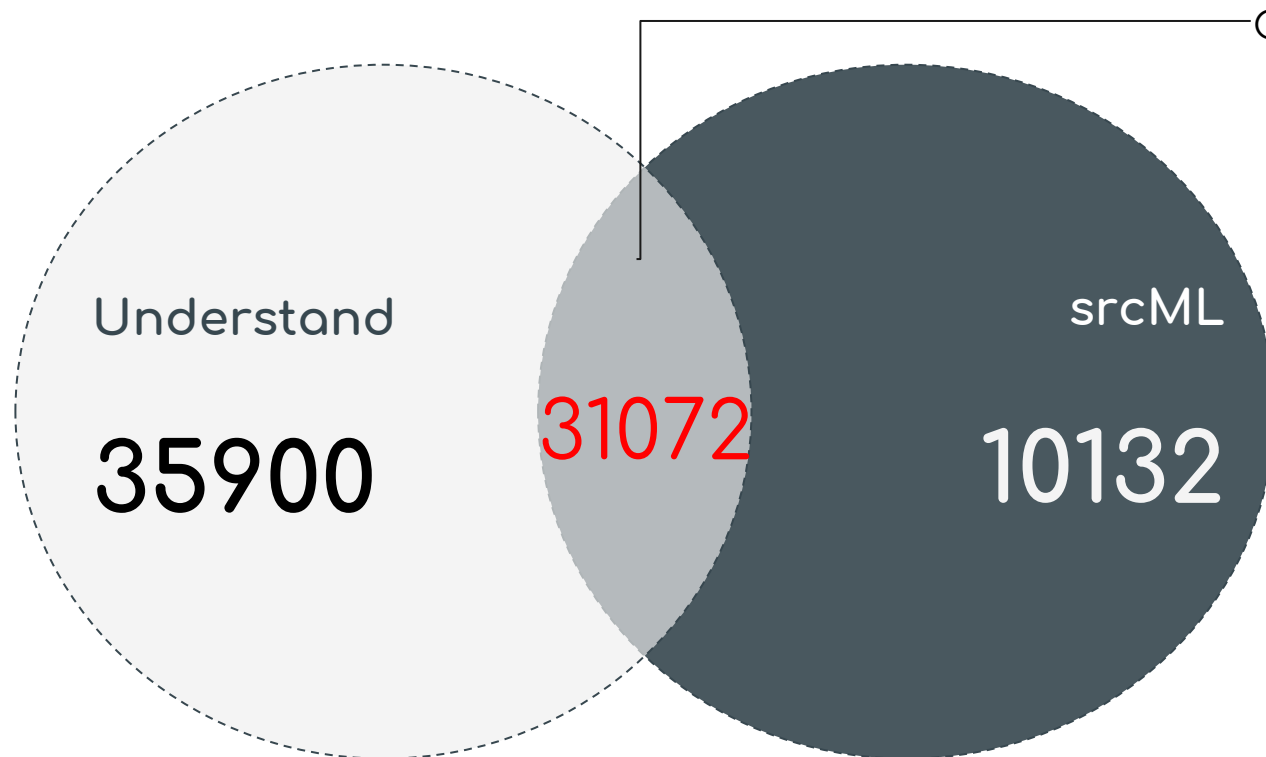
$O(n^2)$

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Quantitative analysis: Understand vs srcML



Common

Understand

srcML

31072

35900

10132

**Understand:**
**66972** total
dependencies

**srcML:**
**41204** total
dependencies

17

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Qualitative analysis: Understand vs srcML

Used sampling calculator with the following data:

1) Understand: Confidence level of 95%, Confidence interval of +/- 6.92%, total population: 76,474. Sampling size: 200

Using stratified sampling method:

a) Overlap: 31,072/76,474 * 200 ~ 81 cases.
b) Understand: 35,900/76,474 * 200 ~ 93 cases.
c) srcML: 10,132/76,474* 200 ~ 26 cases.
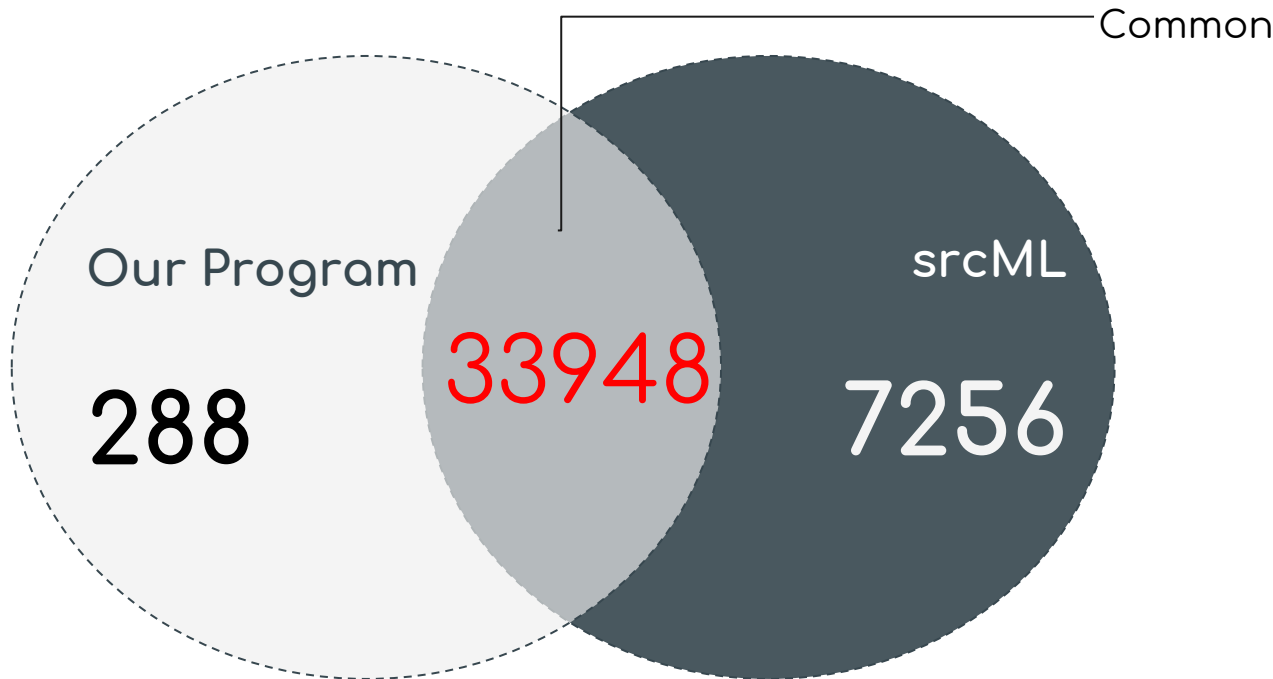
# Noticeable differences: Understand vs srcML

The srcML tool only looks for the "Include" in the .c files, sometimes dependencies can be found in a different way such as inheritance.

Understand scans the whole source code directory, including all types of files such as .txt or .yy.

19

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Quantitative Analysis: Our Program vs. srcML

https://github.com/azkevin/EECS4314/blob/master/A3/a3data/srcML_Include_Common



Common

Our Program

srcML

33948

288

7256

**Our Program:**
**34236** total
dependencies

**srcML:**
**41204** total
dependencies

20

# Qualitative Analysis: Our Program vs. srcML



**Determine Sample Size**

Confidence Level: ● 95% ○ 99%

Confidence Interval: 5

Population: 34236

[Calculate] [Clear]

Sample size needed: 380

# Qualitative Analysis: Our Program vs. srcML

Using stratified sampling method:

Overlap: (33948/41492) * 380 = ~311 cases.     ~81.8%
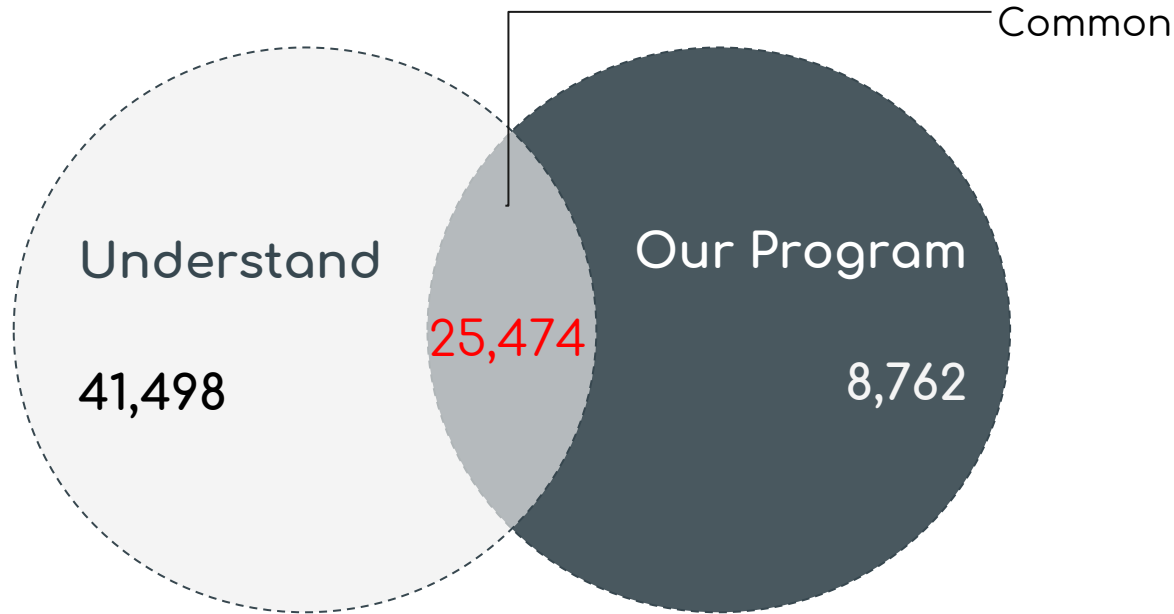
Our Program: (288/41492) * 380 = ~3 cases.     ~0.79%

srcML: (7256/41492)* 380 = ~66  cases.     ~17.4%

22

# Noticeable Differences: Our Program vs. srcML

# Quantitative Analysis: Understand vs. Our Program

Common

Understand

25,474

Our Program

41,498

8,762

## Understand:
**66972** total dependencies

## Our Program:
**34236** total dependencies

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Qualitative Analysis: Understand vs. Our Program

Used sampling calculator with the following data:

1) Understand: Confidence level of 95%, Confidence interval of +/- 6.92%, total population: 76,474. Sampling size: 200

Using stratified sampling method:

a) Overlap: 25,474/75,729 * 200 ~ 67 cases.
b) Understand: 41498/75,729 * 200 ~ 110 cases.
c) Our Program: 8762/75,729* 200 ~ 23 cases.

Dependency Extraction

Quantitative &
Qualitative Analysis

Potential Risks &
Limitations

Lessons Learned

Conclusions

# Noticeable Differences: Understand vs. Our Program

- Our Program does not detect .java, .cpp and .hpp files for dependency derivation
- This explains why Our program detects much less number of dependencies than Understand
- Our Program only looks for the "Include" in the .c files, sometimes dependencies can be found in a different way such as inheritance

# Potential Risks and Limitations

1) srcML puts an extra XML tag before each keyword statement, we have to parse that with a java code to get the "include" statements.

2) srcML skips a significant number of dependencies since it only parse extensions that it supports.

3) Can't get access to the source code of understand, comparing the tool as a black box is limited.

4) Can't analyze the whole output, pick a sample.

5) Hidden dependencies that get generated after the build.

6) The Program we developed for method 3 skips certain extensions as well.

7) The program requires Python interpreter to run which is a limitations.

# Lessons learned

# Conclusion

- Many ways to extract dependencies
- Trade-offs associated with each technique
- Do the dependencies tell the real story?
  - Compare
  - Contrast
  - Conclude

# Question Period