# HBase Dependency Extraction
## Team5Star

● ● ●

Daniel McVicar - 213027479
Rafay Sheikh - 213033451
Daniel Fortunato - 216796443
Adham El Shafie - 212951018
Yahya Ismail- 213235403

# Overview

# Introduction

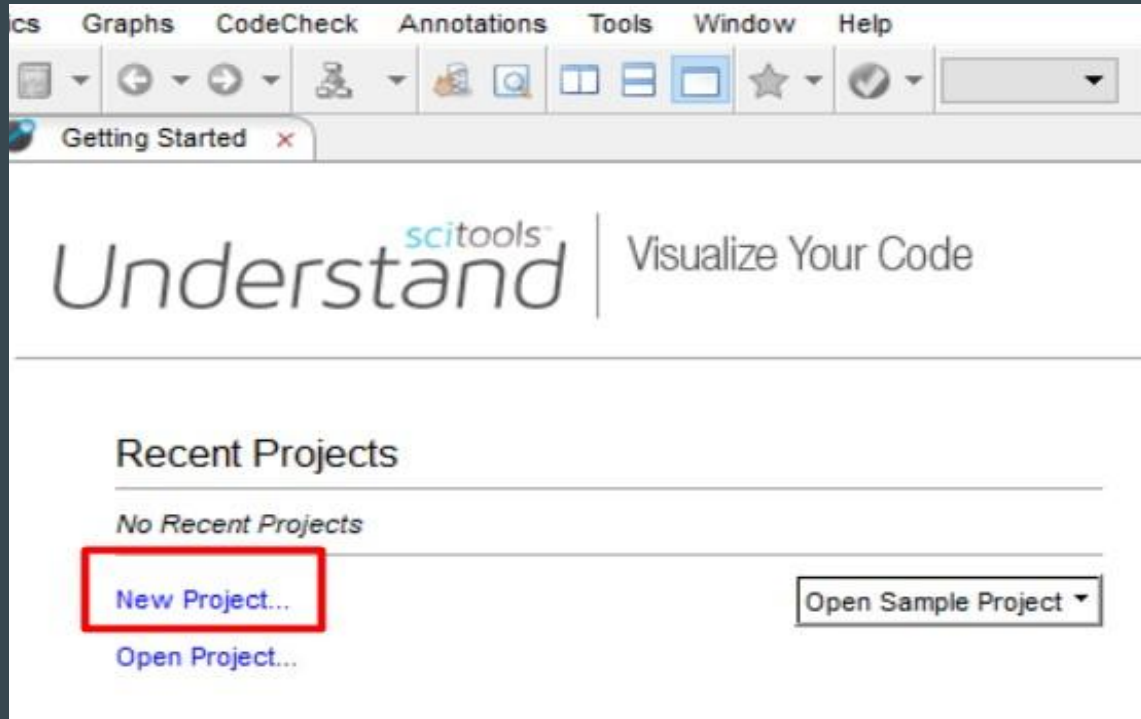Extraction of file level dependencies from source code (directory)

- File dependencies help show how a software works/functions
  - Helps engineers better understand the software to maintain, update etc
- Dependency extraction is a step in the architecture recovery pipeline
- Can be used to understand code with poor documentation
- Analyze code architecture for optimization purposes and code metrics
- 3 approaches for dependency extraction that we used:
  - Understand
  - srcML
  - DependencyExtractor

# What is Understand

- IDE to visualize and analyze static code
- Developed to aid software developers to understand and document source code
- Also helps new developers to comprehend legacy code passed down to them
- Supports a variety of source code languages

# Extracting Dependencies



Once understand is installed, opening it up results in the figure shown

# Extracting Dependencies



Rename the current project
to preferred name



Supported languages for
source code

# Extracting Dependencies



**Source Files**

Add directories that contain the source files you want analyzed. You may choose whether or not to automatically include subdirectories. Files that match languages you selected are added to the project.

You do not need to add files included from other libraries, since those can be added later when you identify include directories.

0 Files   Add a Directory   Add a File ▼

Can add directories or files related to source code



| Directory | C:\Users\xtrem\Desktop\school\4314 | ∨ | ... |
| Configured Filters | Java (*.java) | ∨ | ... |
| Additional Filters | | ∨ | |
| Exclude | .* | ∨ | |

☑ Include subdirectories

☑ Watch this directory

The directory will be watched for any on disk changes. If a file is deleted from disk, it is removed from the project. If new sources are added to this directory or subdirectory of sources is added, it is automically added to the project. Those files and directories must match above filters and settings.
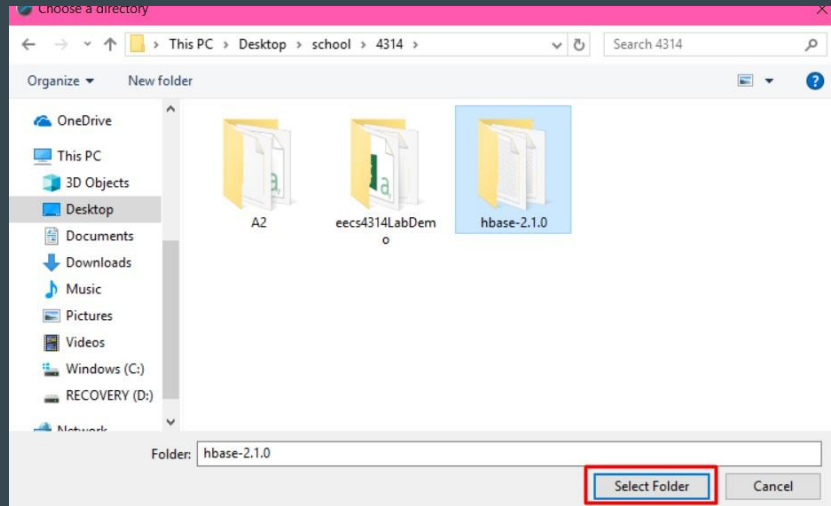
OK   Cancel

"Include subdirectories" and "watch this directory" are checked

# Extracting Dependencies



Navigate to the required source code root directory

# Extracting Dependencies



Confirm root directory



HBase is coded in java
therefore it is in the
"Configured Filters"

# Extracting Dependencies



Subdirectories of HBase
root directory shown

Can configure more settings
if need be

# Extracting Dependencies



Understand begins analyzing the
source code for errors and warnings

# Extracting Dependencies



Highlight project directory and export csv file of dependencies in HBase

Analyzed project shown in "Architecture Browser"

# Extracting Dependencies



Export a .csv (comma separated values) of file dependencies

For comparison, needed short name (file name), from file and to file only

# Extracting Dependencies





Sample of extracted .csv file shown in Excel

Settings for extracted dependencies are done and confirmed

# Representing Dependencies



Understand has a 'graph representation of dependencies' option

Generated graph dependency of HBase between first-level subdirectories

# Advantages and Disadvantages of Scitools: Understand

| Advantages | Disadvantages |
|---|---|
| Easy visualization of source code (through treemaps and dependency graphs) | The languages of the source code need to be known |
| Produces warnings and errors where there might be vulnerabilities in source code | Only has a 'set of rules' to detect these vulnerabilities but there may be more |
| Supports a variety of languages for the legacy/source code | A lot to take in the first time through |

# What is srcML

- A tool for analysing and exploring source code

- Takes source code and transforms it to XML format

- Produced a 150MB text file when run on HBase



```xml
 1  <?xml version="1.0" encoding="UTF-8" standalone="yes"?>
 2  <unit
 3      xmlns="http://www.srcML.org/srcML/src" revision="0.9.5" language="Java" filename="TableName.java">
 4      <comment type="block" format="javadoc">/**
 5   * Licensed to the Apache Software Foundation (ASF) under one
 6   * or more contributor license agreements.  See the NOTICE file
 7   * distributed with this work for additional information
 8   * regarding copyright ownership.  The ASF licenses this file
 9   * to you under the Apache License, Version 2.0 (the
10   * "License"); you may not use this file except in compliance
11   * with the License.  You may obtain a copy of the License at
12   *
13   *     http://www.apache.org/licenses/LICENSE-2.0
14   *
15   * Unless required by applicable law or agreed to in writing, software
16   * distributed under the License is distributed on an "AS IS" BASIS,
17   * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
18   * See the License for the specific language governing permissions and
19   * limitations under the License.
20   */</comment>
21      <package>package
22          <name>
23              <name>org</name>
24              <operator>.</operator>
25              <name>apache</name>
26              <operator>.</operator>
27              <name>hadoop</name>
28              <operator>.</operator>
29              <name>hbase</name>
30          </name>;
31      </package>
32      <import>import
33          <name>
34              <name>java</name>
35              <operator>.</operator>
36              <name>nio</name>
37              <operator>.</operator>
38              <name>ByteBuffer</name>
39          </name>;
40      </import>
41      <import>import
42          <name>
43              <name>java</name>
```

# Extracting Dependencies

```python
import lxml
from lxml import etree

#Takes an array of elements and checks to see if any of
#those elements suggest their parent is part of hbase.
def is_in_hbase(child_array):
    for child in child_array:
        if child.text:
            if "apache" in child.text or "hadoop" in child.text or "hbase" in child.text:
                return True
    return False

tree = etree.parse("hbase_srcml_output.xml")
root = tree.getroot()

for unit in root:
    if "unit" in unit.tag:
        file_name = unit.get("filename")
        file_name = file_name.split("\\")[-1]

        for child in unit:
            if "import" in child.tag:
                for subchild in child:
                    if is_in_hbase(subchild[:]):
                        dependency_name = subchild[-1].text + ".java"
                        print("IMPORT," + file_name + "," + dependency_name)
```
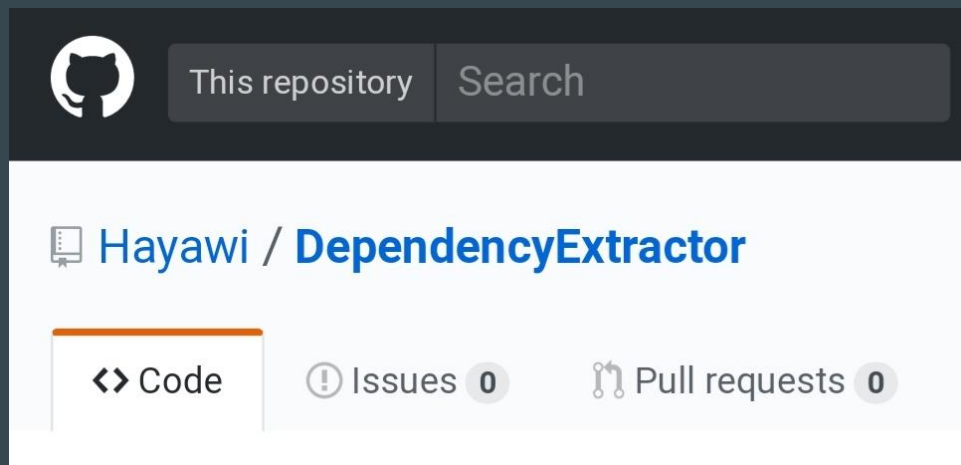
```
IMPORT,TestInstancePending.java,HBaseClassTestRule.java
IMPORT,TestInstancePending.java,SmallTests.java
IMPORT,TestInstancePending.java,ZKTests.java
IMPORT,TestHQuorumPeer.java,Configuration.java
IMPORT,TestHQuorumPeer.java,FileSystem.java
IMPORT,TestHQuorumPeer.java,Path.java
IMPORT,TestHQuorumPeer.java,HBaseClassTestRule.java
IMPORT,TestHQuorumPeer.java,HBaseConfiguration.java
IMPORT,TestHQuorumPeer.java,HBaseZKTestingUtility.java
IMPORT,TestHQuorumPeer.java,HConstants.java
IMPORT,TestHQuorumPeer.java,MediumTests.java
IMPORT,TestHQuorumPeer.java,ZKTests.java
IMPORT,HBaseZKTestingUtility.java,Configuration.java
IMPORT,HBaseZKTestingUtility.java,Path.java
IMPORT,HBaseZKTestingUtility.java,MiniZooKeeperCluster.java
IMPORT,HBaseZKTestingUtility.java,ZKWatcher.java
IMPORT,HBaseZKTestingUtility.java,InterfaceAudience.java
IMPORT,TestReadOnlyZKClient.java,Configuration.java
IMPORT,TestReadOnlyZKClient.java,HBaseClassTestRule.java
IMPORT,TestReadOnlyZKClient.java,HBaseZKTestingUtility.java
IMPORT,TestReadOnlyZKClient.java,HConstants.java
IMPORT,TestReadOnlyZKClient.java,ExplainingPredicate.java
IMPORT,TestReadOnlyZKClient.java,MediumTests.java
IMPORT,TestReadOnlyZKClient.java,ZKTests.java
IMPORT,TestReadOnlyZKClient.java,AsyncCallback.java
IMPORT,TestReadOnlyZKClient.java,CreateMode.java
IMPORT,TestReadOnlyZKClient.java,KeeperException.java
IMPORT,TestReadOnlyZKClient.java,Code.java
IMPORT,TestReadOnlyZKClient.java,ZooDefs.java
IMPORT,TestReadOnlyZKClient.java,ZooKeeper.java
IMPORT,TestRecoverableZooKeeper.java,Configuration.java
IMPORT,TestRecoverableZooKeeper.java,Abortable.java
IMPORT,TestRecoverableZooKeeper.java,HBaseClassTestRule.java
IMPORT,TestRecoverableZooKeeper.java,HBaseZKTestingUtility.java
IMPORT,TestRecoverableZooKeeper.java,HConstants.java
IMPORT,TestRecoverableZooKeeper.java,MediumTests.java
IMPORT,TestRecoverableZooKeeper.java,ZKTests.java
IMPORT,TestRecoverableZooKeeper.java,Bytes.java
IMPORT,TestRecoverableZooKeeper.java,CreateMode.java
IMPORT,TestRecoverableZooKeeper.java,KeeperException.java
IMPORT,TestRecoverableZooKeeper.java,Watcher.java
IMPORT,TestRecoverableZooKeeper.java,Ids.java
IMPORT,TestRecoverableZooKeeper.java,ZooKeeper.java
IMPORT,TestRecoverableZooKeeper.java,Stat.java
IMPORT,TestZKLeaderManager.java,Configuration.java
IMPORT,TestZKLeaderManager.java,Abortable.java
IMPORT,TestZKLeaderManager.java,HBaseClassTestRule.java
```

# Advantages and Disadvantages of srcML

| Advantages | Disadvantages |
|---|---|
| High level of detail for each class provides more information about interactions between classes | Large initial XML document is not human readable |
| XPath provides a robust search language for XML documents | Searching for detailed information requires extra programing to use XPath |
| Works with source code in any language | |

# DependencyExtractor (Using Import Statements)

- Java based application that takes in a root directory and scrubs the import statements
- Powerful and limitless, i.e. can be made to do anything
- Extremely time consuming, especially as features expand
- Reinventing the wheel

This repository | Search

Hayawi / **DependencyExtractor**

<> Code    ⚠ Issues 0    Pull requests 0

# DependencyExtractor Usecase



DependencyExtractor Use Case Sequence Diagram

# Comparison Process - Overview

Quantitative Analysis

- Gather data from all dependency tools used

- Format data with excel

- Once all data is in the same format we can begin comparing

- Use script to compare the three sets of data gathered from the tools

Qualitative Analysis
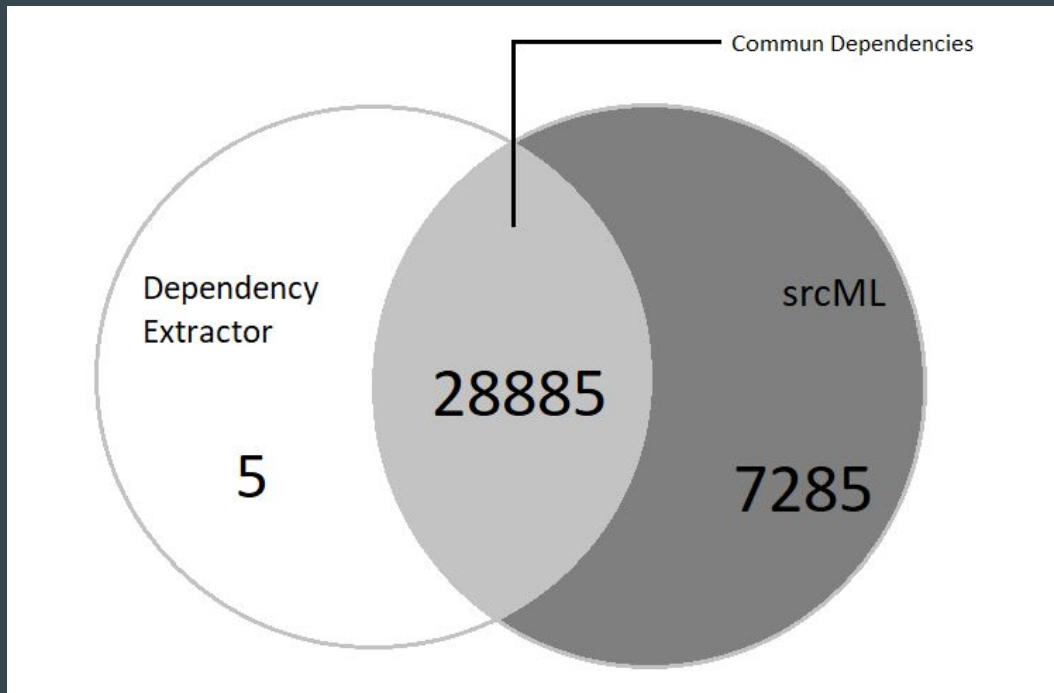
- Use sample size calculator

- Confidence level of: 95%, Confidence interval: 5, Population: 36175, indicating a required sample size of 380 random dependencies

- Use random number generator in excel to shuffle entries

- Take the first 380 samples

- Compare and contrast DependencyExtractor and srcML

- Use Understand as a control group to source conclusions on dependency relevance

# Quantitative Analysis - The script

```python
data_a = set(open("A_dependency_extractor.txt","r").readlines())
data_b = set(open("B_srcml_dependencies.txt","r").readlines())
data_c = set(open("C_understand_dependencies.csv","r").readlines())

total_nr_dep_a = len(data_a)
total_nr_dep_b = len(data_b)
total_nr_dep_c = len(data_c)

print("The total number of dependencies for data_A is: " + str(total_nr_dep_a))
print("The total number of dependencies for data_B is: " + str(total_nr_dep_b))
print("The total number of dependencies for data_C is: " + str(total_nr_dep_c))

same_count = 0
missing_count = 0
for line in data_a:
    if line in data_b:
        same_count += 1
    else:
        missing_count += 1

print("Data set A and B share " + str(same_count) + " dependencies. Data set A has " + str(missing_count) + " dependencies not mentioned in Data set B.")

same_count = 1
missing_count = 0
for line in data_b:
    if line in data_a:
        same_count += 1
    else:
        missing_count += 1

print("Data set B and A share " + str(same_count) + " dependencies. Data set B has " + str(missing_count) + " dependencies not mentioned in Data set A.")
```

23

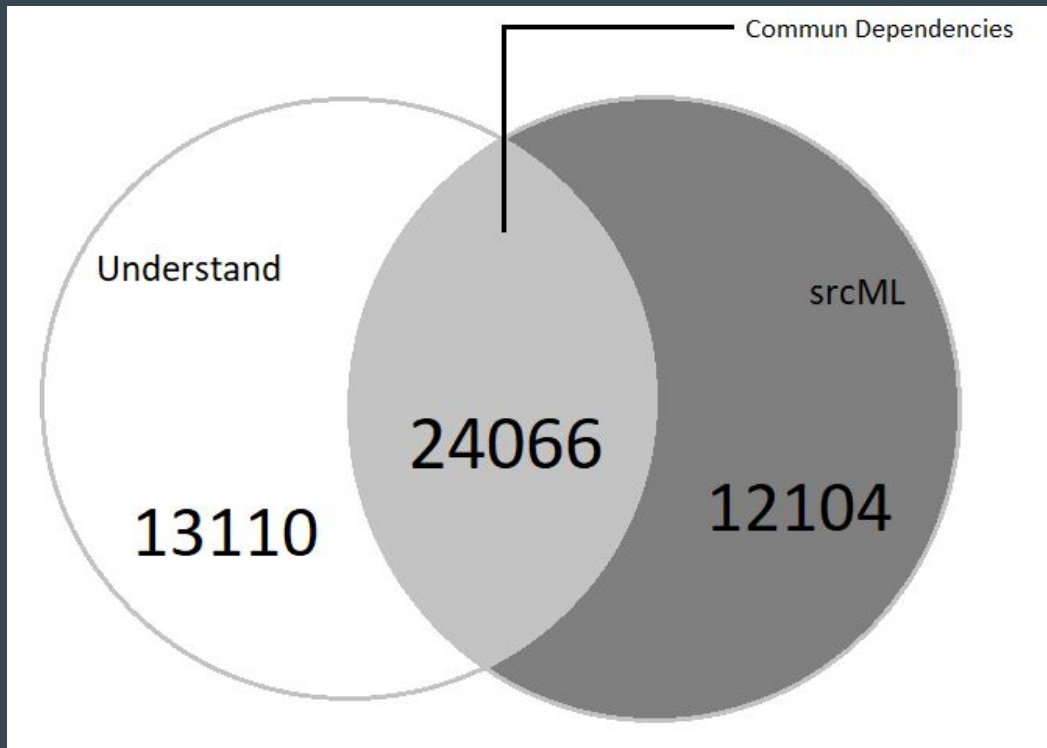# Quantitative Analysis - DependencyExtractor vs srcML



DependencyExtractor:

28890 total dependencies

srcML:

36170 total dependencies

# Quantitative Analysis - Understand vs srcML
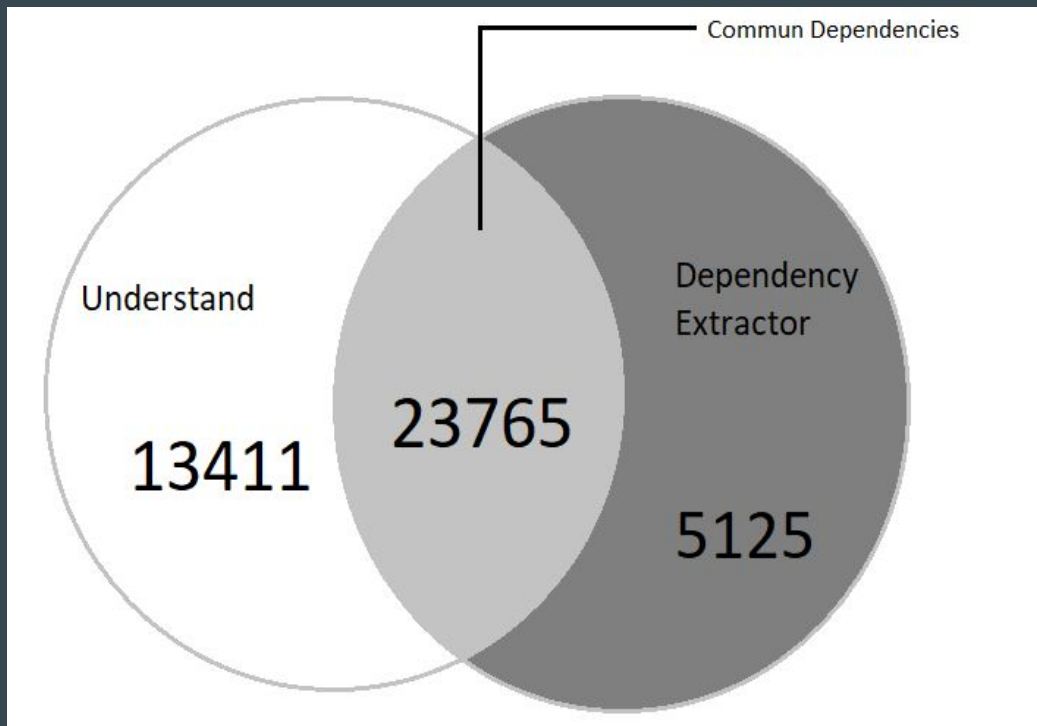


Understand:

37176 total dependencies

srcML:

36170 total dependencies

# Quantitative Analysis - Understand vs DependencyExtractor



Understand:

37176 total dependencies

DependencyExtractor:

28890 total dependencies

# Qualitative Analysis - DependencyExtractor vs srcML

1. Overlap: ~79%
2. DependencyExtractor: ~0.000138%
3. srcML: ~20%

- Large overlap between DependencyExtractor and srcML
- Very few dependencies in DependencyExtractor were unique, while 20% of srcML were unique.

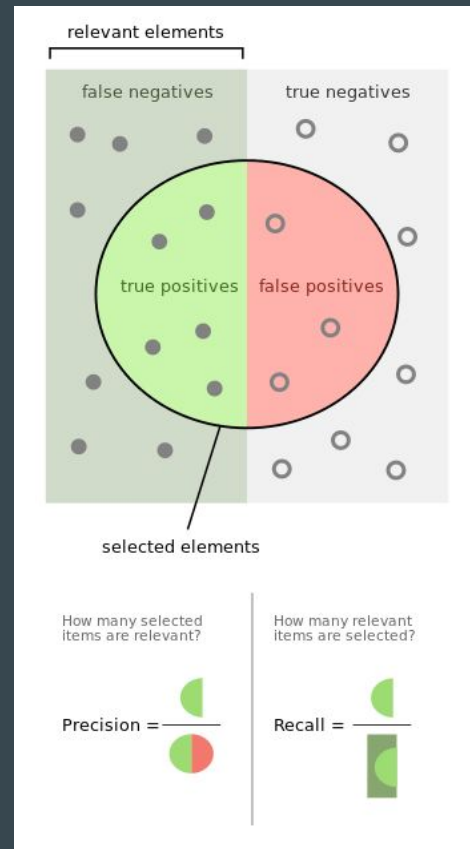Why? Is srcML returning false positives?

**Determine Sample Size**

| | |
|---|---|
| Confidence Level: | ⦿95% ○99% |
| Confidence Interval: | 5 |
| Population: | 36175 |
| Calculate | Clear |
| Sample size needed: | 380 |

# Qualitative analysis - Precision and Recall

- Precision (Positive Predictive Value) is the fraction of retrieved instances that are relevant among the retrieved instances.

- Recall (Sensitivity) is the fraction of relevant instances that have been retrieved over the total amount of relevant instances.



relevant elements

false negatives

true negatives

true positives

false positives

selected elements

How many selected items are relevant?

How many relevant items are selected?

Precision =

Recall =

# Qualitative Analysis - srcML Precision and Recall

Using Understand as our control group we will calculate precision and recall.

1. Precision: ~66.5%
2. Recall: ~64.7%

- Mediocre Precision and Recall
- srcML seems to be inaccurately returning a large set of incorrect dependencies

Is DependencyExtractor any better?

# Qualitative Analysis - DE Precision and Recall

Using Understand as our control group we will calculate precision and recall.

1. Precision: ~82.3%
2. Recall: ~63.9%

- Impressive Precision with mediocre Recall
- It seems DependencyExtractor is doing a better job at extracting accurate dependencies, which explains the odd statistics from before

srcML is retrieving a lot of irrelevant dependencies with it's shotgun approach, leading to an abysmal precision.

# Risks and Limitations

- The Dependency Extractor is only as robust as the time spent extracting dependencies (very time consuming if we want to do a perfect job)
- The srcML generates of huge XML file which is not easily human readable and requires a second parsing step to extract desired information.
- Understand needs a properly configured and ordered source code directo so that the representation is accurate
- Since the list of dependencies is so big , some might have been lost when using the extraction softwares
- The sample taken for comparison may not be representative of the entire system

# Conclusion (+Lessons Learned)

- Large scale and complex systems:
  - All 3 methods gave us different number of dependencies
  - Unlikely that different extractions methods will give similar results
  - Errors/mistakes possible
- Understand
  - Most comprehensive and inclusive
- srcMl
  - Limited number of dependencies (similar to our own program)
  - Requires parsing and need to format data
- DependencyExtractor
  - Many dependencies limited by inheritance, functions calls etc.
  - Software could have multiple coding languages so would need to account for that