# Lecture 1: Introduction to Program Analysis

17-355/17-655/17-819: Program Analysis

Rohan Padhye  and Jonathan Aldrich

February 2, 2021

* Course materials developed with Claire Le Goues

# Introductions



Prof. Rohan Padhye

Prof. Jonathan Aldrich

TA Priya Varra

**Carnegie Mellon University**
School of Computer Science

# Learning objectives

- Provide a high level definition of program analysis and give examples of why it is useful.

- Sketch the explanation for why all analyses must approximate.

- Understand the course mechanics, and be motivated to read the syllabus.

- Describe the function of an AST and outline the principles behind AST walkers for simple bug-finding analyses.

- Recognize the basic WHILE demonstration language and translate between WHILE and While3Addr.
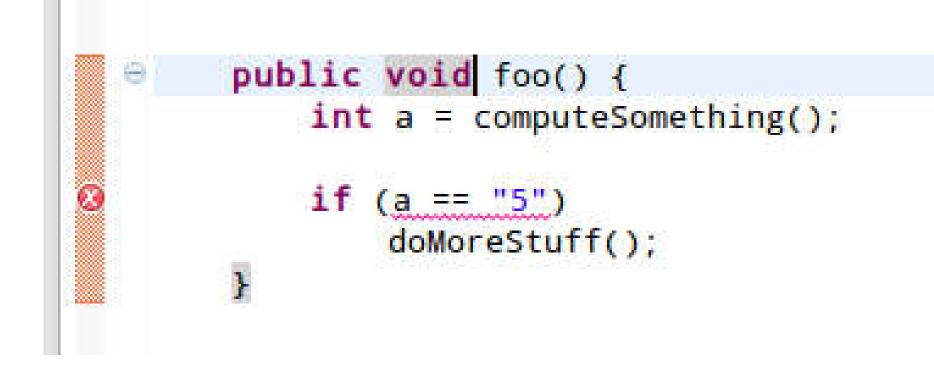
# What is this course about?

- Program analysis is the systematic examination of a program to determine its properties.

- From 30,000 feet, this requires:
  - Precise program representations
  - Tractable, systematic ways to reason over those representations.

- We will learn:
  - How to unambiguously define the meaning of a program, and a programming language.
  - How to prove theorems about the behavior of particular programs.
  - How to use, build, and extend tools that do the above, automatically.

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Why might you care?
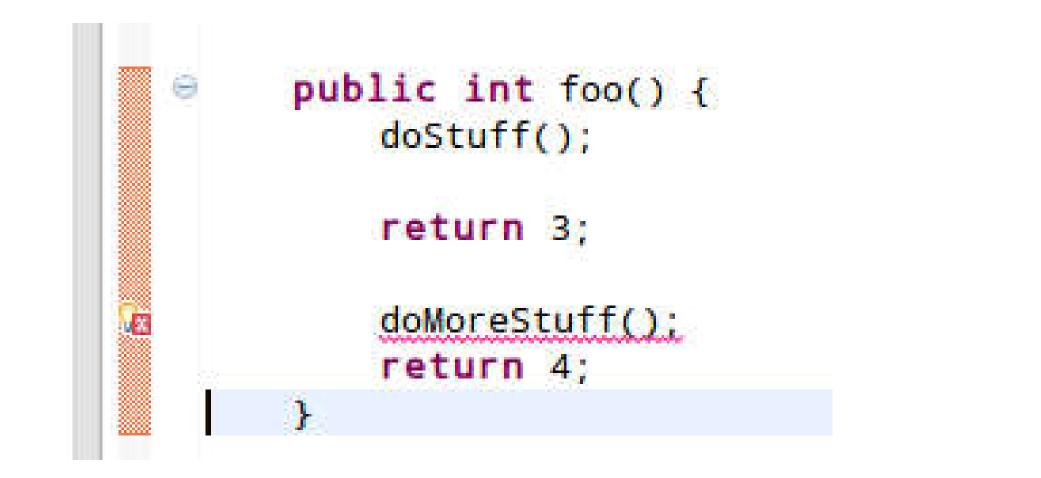
- Program analysis, and the skills that underlie it, have implications for:
  - Automatic bug finding.
  - Language design and implementation.
  - Program synthesis.
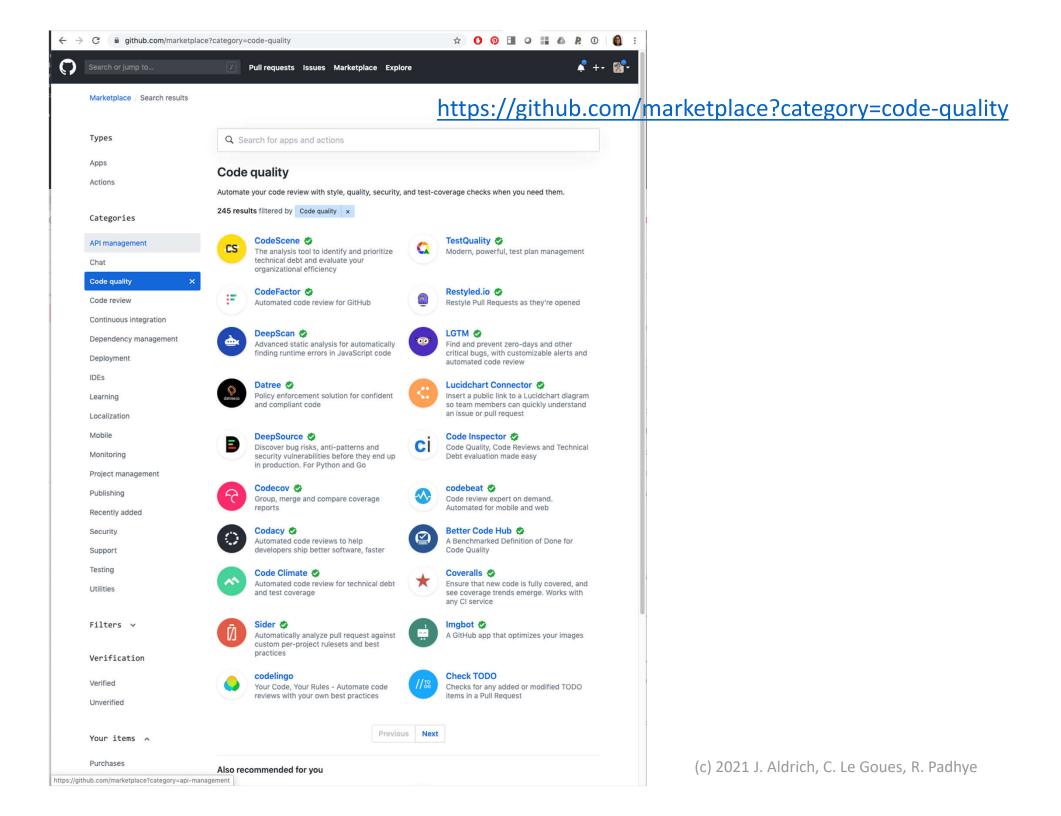  - Program transformation (refactoring, optimization, repair).

```java
public int foo() {
    doStuff();

    return 3;

    doMoreStuff();
    return 4;
}
```

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

https://github.com/marketplace?category=code-quality

(c) 2021 J. Aldrich, C. Le Goues, R. Padhye

8

```
package com.google.devtools.staticanalysis;

public class Test {
```

▾ Lint            Missing a Javadoc comment.
Java
1:02 AM, Aug 21

Please fix                                                              Not useful

```
  public boolean foo() {
    return getString() == "foo".toString();
```

▾ ErrorProne     String comparison using reference equality instead of value equality
StringEquality       (see http://code.google.com/p/error-prone/wiki/StringEquality)
1:03 AM, Aug 21

Please fix

**Suggested fix attached:** show                                        Not useful

```
  }

  public String getString() {
    return new String("foo");
  }
}
```

//depot/google3/java/com/google/devtools/staticanalysis/Test.java

```
package com.google.devtools.staticanalysis;           package com.google.devtools.staticanalysis;

                                                      import java.util.Objects;

public class Test {                                   public class Test {
  public boolean foo() {                                public boolean foo() {
    return getString() == "foo".toString();              return Objects.equals(getString(), "foo".toString());
  }                                                     }

  public String getString() {                           public String getString() {
    return new String("foo");                             return new String("foo");
  }                                                     }
}                                                     }
```

**Apply**    Cancel

# IS THERE A BUG IN THIS CODE?

```
1. /* from Linux 2.3.99 drivers/block/raid5.c */
2. static struct buffer_head *
3. get_free_buffer(struct stripe_head * sh,
4.                  int b_size) {
5.    struct buffer_head *bh;
6.    unsigned long flags;
7.    save_flags(flags);
8.    cli(); // disables interrupts
9.    if ((bh = sh->buffer_pool) == NULL)
10.       return NULL;
11.   sh->buffer_pool = bh -> b_next;
12.   bh->b_size = b_size;
13.   restore_flags(flags); // re-enables interrupts
14.   return bh;
15. }
```

ERROR: function returns with interrupts disabled!
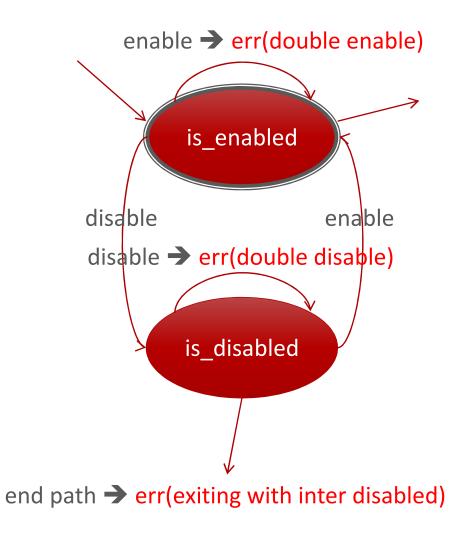
Example from Engler et al., *Checking system rules Using System-Specific, Programmer-Written Compiler Extensions*, OSDI '000

```
1. sm check_interrupts {
2. // variables; used in patterns
3. decl { unsigned } flags;
4. // patterns specify enable/disable functions
5. pat enable = { sti() ; }
6.                | { restore_flags(flags); } ;
7. pat disable = { cli() ; }
8. //states; first state is initial
9. is_enabled : disable ➜ is_disabled
10.    | enable ➜ { err("double enable"); }
11.;
12. is_disabled : enable ➜ is_enabled
13.    | disable ➜ { err("double disable"); }
14.//special pattern that matches when
15.// end of path is reached in this state
16.    | $end_of_path$ ➜
17.        { err("exiting with inter disabled!"); }
18.;
19.}
```



enable ➜ err(double enable)

is_enabled

disable          enable

disable ➜ err(double disable)

is_disabled

end path ➜ err(exiting with inter disabled)

```
1. /* from Linux 2.3.99 drivers/block/raid5.c */
2. static struct buffer_head *
3. get_free_buffer(struct stripe_head * sh,
4.                 int b_size) {
5.    struct buffer_head *bh;
6.    unsigned long flags;
7.    save_flags(flags);
8.    cli(); // disables interrupts
9.    if ((bh = sh->buffer_pool) == NULL)
10.        return NULL;
11.    sh->buffer_pool = bh -> b_next;
12.    bh->b_size = b_size;
13.    restore_flags(flags); // re-enables interrupts
14.    return bh;
15. }
```

Initial state: is_enabled

Example from Engler et al., *Checking system rules Using System-Specific, Programmer-Written Compiler Extensions*, OSDI '000

```
1. /* from Linux 2.3.99 drivers/block/raid5.c */
2. static struct buffer_head *
3. get_free_buffer(struct stripe_head * sh,
4.                 int b_size) {
5.   struct buffer_head *bh;
6.   unsigned long flags;
7.   save_flags(flags);
8.   cli(); // disables interrupts
9.   if ((bh = sh->buffer_pool) == NULL)
10.     return NULL;
11.   sh->buffer_pool = bh -> b_next;
12.   bh->b_size = b_size;
13.   restore_flags(flags); // re-enables interrupts
14.   return bh;
15.}
```

Transition to: is_disabled

Example from Engler et al., *Checking system rules Using System-Specific, Programmer-Written Compiler Extensions*, OSDI '000

```
1. /* from Linux 2.3.99 drivers/block/raid5.c */
2. static struct buffer_head *
3. get_free_buffer(struct stripe_head * sh,
4.                 int b_size) {
5.    struct buffer_head *bh;
6.    unsigned long flags;
7.    save_flags(flags);
8.    cli(); // disables interrupts
9.    if ((bh = sh->buffer_pool) == NULL)
10.       return NULL;
11.   sh->buffer_pool = bh -> b_next;
12.   bh->b_size = b_size;
13.   restore_flags(flags); // re-enables interrupts
14.   return bh;
15. }
```

Final state: is_disabled

Example from Engler et al., *Checking system rules Using System-Specific, Programmer-Written Compiler Extensions*, OSDI '000

```
1. /* from Linux 2.3.99 drivers/block/raid5.c */
2. static struct buffer_head *
3. get_free_buffer(struct stripe_head * sh,
4.                      int b_size) {
5.    struct buffer_head *bh;
6.    unsigned long flags;
7.    save_flags(flags);
8.    cli(); // disables interrupts
9.    if ((bh = sh->buffer_pool) == NULL)
10.       return NULL;
11.   sh->buffer_pool = bh -> b_next;
12.   bh->b_size = b_size;
13.   restore_flags(flags); // re-enables interrupts
14.   return bh;
15. }
```

Transition to: is_enabled

Example from Engler et al., *Checking system rules Using System-Specific, Programmer-Written Compiler Extensions*, OSDI '000

```
1. /* from Linux 2.3.99 drivers/block/raid5.c */
2. static struct buffer_head *
3. get_free_buffer(struct stripe_head * sh,
4.                 int b_size) {
5.    struct buffer_head *bh;
6.    unsigned long flags;
7.    save_flags(flags);
8.    cli(); // disables interrupts
9.    if ((bh = sh->buffer_pool) == NULL)
10.       return NULL;
11.   sh->buffer_pool = bh -> b_next
12.   bh->b_size = b_size;
13.   restore_flags(flags); // re-enables interrupts
14.   return bh;
15. }
```

Final state: is_enabled

Example from Engler et al., *Checking system rules Using System-Specific, Programmer-Written Compiler Extensions*, OSDI '000

# Behavior of interest...

- Is on uncommon execution paths.

  o Hard to exercise when testing.

- Executing (or analyzing) all paths is infeasible

- **Instead: (abstractly) check the entire possible state space of the program.**

# What is this course about?

- **Program analysis is *the systematic examination of a program to determine its properties*.**

- From 30,000 feet, this requires:
  - Precise program representations
  - Tractable, systematic ways to reason over those representations.

- We will learn:
  - How to unambiguously define the meaning of a program, and a programming language.
  - How to prove theorems about the behavior of particular programs.
  - How to use, build, and extend tools that do the above, automatically.

(c) 2021 J. Aldrich, C. Le Goues, R. Padhye

# The Bad News: Rice's Theorem

"Any nontrivial property about the language recognized by a Turing machine is undecidable."

Henry Gordon Rice, 1953

# Proof by contradiction (sketch)

Assume that you have a function that can determine if a program *p* has some nontrivial property (like `divides_by_zero`):

```
1.  int silly(program p, input i) {
2.    p(i);
3.    return 5/0;
4.  }
5.  bool halts(program p, input i) {
6.   return divides_by_zero(`silly(p,i)`);
7.  }
```

**Carnegie Mellon University**
School of Computer Science

|  | **Error exists** | **No error exists** |
| --- | --- | --- |
| **Error Reported** | True positive (correct analysis result) | False positive |
| **No Error Reported** | False negative | True negative (correct analysis result) |

Sound Analysis:

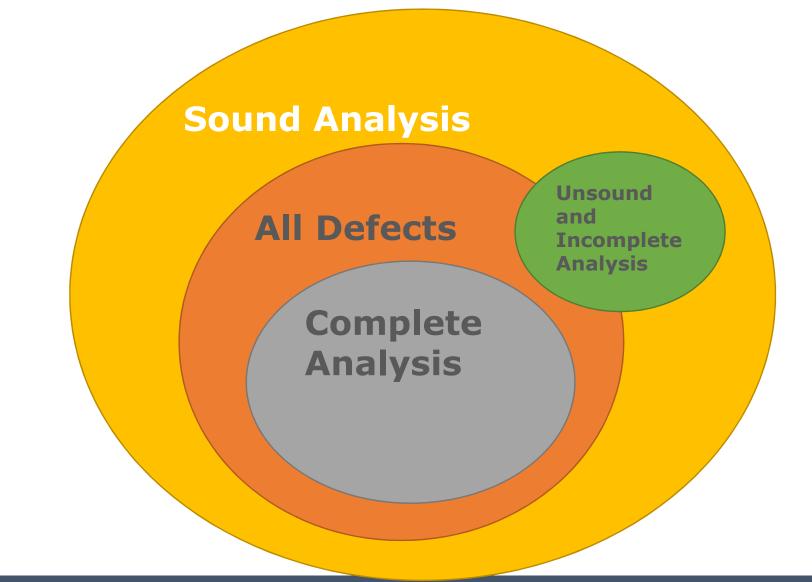      reports all defects

      -> no false negatives

      typically overapproximated

Complete Analysis:

      every reported defect is an actual defect

      -> no false positives

      typically underapproximated

# What is this course about?

- **Program analysis is** *the systematic examination of a program to determine its properties*.

- From 30,000 feet, this requires:
  - Precise program representations
  - Tractable, systematic ways to reason over those representations.

- We will learn:
  - How to unambiguously define the meaning of a program, and a programming language.
  - How to prove theorems about the behavior of particular programs.
  - How to use, build, and extend tools that do the above, automatically.

**Carnegie Mellon University**
School of Computer Science

# What is this course about?

- Program analysis is *the systematic examination of a program to determine its properties*.

- Principal techniques:
  - **Dynamic:**
    - **Testing:** Direct execution of code on test data in a controlled environment.
    - **Analysis:** Tools extracting data from test runs.
  - **Static:**
    - **Inspection:** Human evaluation of code, design documents (specs and models), modifications.
    - **Analysis:** Tools reasoning about the program without executing it.
  - …and their combination.

# Course topics

- Program representation
- Abstract interpretation: Use abstraction to reason about possible program behavior.
  - Operational semantics.
  - Dataflow Analysis
  - Termination, complexity
  - Widening, collecting
  - Interprocedural analysis
  - Datalog
  - Control flow analysis
- Hoare-style verification: Make logical arguments about program behavior.
  - Axiomatic semantics
  - Separation logic: modern bug finding.

- Symbolic execution: test all possible executions paths simultaneously.
  - Concolic execution
  - Test generation
- SAT/SMT solvers
- Program synthesis
- Dynamic analysis
- Fuzzing
- Program repair
- Model checking (briefly) : reason exhaustively about possible program states.
  - Take 15-414 if you want the full treatment!
- We will basically *not* cover types.

**Carnegie Mellon University**
School of Computer Science

# Fundamental concepts

- Abstraction.
    - Elide details of a specific implementation.
    - Capture semantically relevant details; ignore the rest.

- The importance of semantics.
    - We prove things about analyses with respect to the semantics of the underlying language.

- Program proofs as inductive invariants.

- Implementation
    - You do not understand analysis until you have written several.

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Course mechanics

# When/what.

- Lectures 2x week (T,Th – hybrid in-person + virtual).
  - Active learning exercise(s) in every class
  - Lecture notes for review

- Recitation 1x week (Fr – virtual).
  - Lab-like, very helpful for homework.
  - Be ready to work

- Homework, midterm exams, project.

- There is an optional textbook.

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# Communication

- We have a website and a Canvas site, with Piazza enabled.
  - Follow the link from the main Canvas page/syllabus to sign up for Piazza.

- Please:
  - Use Piazza to communicate with us as much as possible, unless the matter is sensitive.
  - Make your questions *public* as much as possible, since that's the literal point of Piazza.

- We have office hours! Or, by appointment.

# "How do I get an A?"

- 15% in-class participation and exercises
- 40% homework
  - Both written (proof-y) and coding (implementation-y).
  - First one (mostly coding) to be released by Friday!
- 25% midterm exam
- 20% final project
  - There will be some options here.
- No final exam; exam slot used for project presentations.
- We have late days and a late day policy; read the syllabus.

institute for SOFTWARE RESEARCH

**Carnegie Mellon University**
School of Computer Science

# CMU can be a pretty intense place.

- A 12-credit course is expected to take ~ 12 hours a week.

- We aim to provide a rigorous but tractable course.
  - More frequent assignments rather than big monoliths
  - Midterm exam to cover core material from first half of course

- Please keep us apprised of how much time the class is actually taking and whether it is interfacing badly with other courses.
  - We have no way of knowing if you have three midterms in one week.
  - Sometimes, we misjudge assignment difficulty.

- If it's 2 am and you're panicking…put the homework down, send us an email, and go to bed.

# What is this course about?

- Program analysis is *the systematic examination of a program to determine its properties*.

- From 30,000 feet, this requires:
  - Precise program representations
  - Tractable, systematic ways to reason over those representations.

- We will learn:
  - How to unambiguously define the meaning of a program, and a programming language.
  - How to prove theorems about the behavior of particular programs.
  - How to use, build, and extend tools that do the above, automatically.

# Our first representation: Abstract Syntax

- A tree representation of source code based on the language grammar.

- Concrete syntax: The rules by which programs can be expressed as strings of characters.
  - Use finite automata and context-free grammars, automatic lexer/parser generators

- Abstract syntax: a subset of the parse tree of the program.

- (The intuition is fine for this course; take compilers if you want to learn how to parse for real.)

# WHILE abstract syntax

- Categories:
  - $S \in$ **Stmt**          statements
  - $a \in$ **Aexp**          arithmetic expressions
  - $x, y \in$ **Var**          variables
  - $n \in$ **Num**          number literals
  - $P \in$ **BExp**          boolean predicates
  - $l \in$ **labels**          statement addresses (line numbers)

- Syntax:
  - ```
    S   ::= x := a | skip | S₁ ; S₂
    ```
    $S \quad ::= x := a \mid \text{skip} \mid S_1 ; S_2$

    $\quad\quad\quad \mid \text{if } P \text{ then } S_1 \text{ else } S_2 \mid \text{while } P \text{ do } S$
  - $a \quad ::= x \mid n \mid a_1 \ op_a \ a_2$
  - $op_a ::= + \mid - \mid * \mid / \mid \ldots$
  - $P \quad ::= \text{true} \mid \text{false} \mid \text{not } P \mid P_1 \ op_b \ P_2 \mid a1 \ op_r \ a2$
  - $op_b ::= \text{ and } \mid \text{or} \mid \ldots$
  - $op_r ::= < \mid \leq \mid = \mid > \mid \geq \mid \ldots$

Concrete syntax is similar, but adds things like (parentheses) for disambiguation during parsing

# Example WHILE program

```
y := x;
z := 1;
while y > 1 do
    z := z * y;
    y := y - 1
```

# Exercise: Building an AST

```
y := x;
z := 1;
while y > 1 do
   z := z * y;
   y := y - 1
```

# Practice: Building an AST for C code

```c
void copy_bytes(char dest[], char source[], int n) {
    for (int i = 0; i < n; ++i)
        dest[i] = source[i];

}
```

# Our first static analysis: AST walking

- One way to find "bugs" is to walk the AST, looking for particular patterns.
  - Walk the AST, look for nodes of a particular type
  - Check the neighborhood of the node for the pattern in question.

- Various frameworks, some more language-specific than others.
  - Tension between language agnosticism and semantic information available.
  - Consider "grep": very language agnostic, not very smart.

- One common architecture based on Visitor pattern:
  - class Visitor has a visitX method for each type of AST node X
  - Default Visitor code just descends the AST, visiting each node
  - To find a bug in AST element of type X, override visitX

- Other more recent approaches based on semantic search, declarative logic programming, or query languages.

# Example: shifting by more than 31 bits.

```
For each instruction I in the program
   if I is a shift instruction
      if (type of I's left operand is int
            && I's right operand is a constant
            && value of constant < 0 or > 31)
         warn("Shifting by less than 0 or more
               than 31 is meaningless")
```

Dashboard / Java queries                                                    ...

# Inefficient empty string test

https://help.semmle.com/wiki/display/JAVA/Inefficient+empty+string+test

Created by Documentation team, last modified on Mar 28, 2019

**Name:** Inefficient empty string test

**Description:** Checking a string for equality with an empty string is inefficient.

**ID:** java/inefficient-empty-string-test

**Kind:** problem

**Severity:** recommendation

**Precision:** high

**Query: InefficientEmptyStringTest.ql**                              › Expand source

When checking whether a string s is empty, perhaps the most obvious solution is to write something like `s.equals("")` (or `"".equals(s)`). However, this actually carries a fairly significant overhead, because `String.equals` performs a number of type tests and conversions before starting to compare the content of the strings.

## Recommendation

The preferred way of checking whether a string s is empty is to check if its length is equal to zero. Thus, the condition is `s.length() == 0`. The `length` method is implemented as a simple field access, and so should be noticeably faster than calling `equals`.

Note that in Java 6 and later, the `String` class has an `isEmpty` method that checks whether a string is empty. If the codebase does not need to support Java 5, it may be better to use that method instead.

```
1   // Inefficient version
2   class InefficientDBClient {
3       public void connect(String user, String pw) {
4           if (user.equals("") || "".equals(pw))
5               throw new RuntimeException();
6           ...
7       }
8   }
9
10  // More efficient version
11  class EfficientDBClient {
12      public void connect(String user, String pw) {
13          if (user.length() == 0 || (pw != null && pw.length() == 0))
14              throw new RuntimeException();
15          ...
16      }
17  }
```

Hint: doub

## Query: InefficientEmptyStringTest.ql

```
/**
 * @name Inefficient empty string test
 * @description Checking a string for equality with an empty string is inefficient.
 * @kind problem
 * @problem.severity recommendation
 * @precision high
 * @id java/inefficient-empty-string-test
 * @tags efficiency
 *       maintainability
 */

import java

from MethodAccess mc
where
  mc.getQualifier().getType() instanceof TypeString and
  mc.getMethod().hasName("equals") and
  (
    mc.getArgument(0).(StringLiteral).getRepresentedString() = "" or
    mc.getQualifier().(StringLiteral).getRepresentedString() = ""
  )
select mc, "Inefficient comparison to empty string, check for zero length instead."
```

44

# Practice: String concatenation in a loop

- Write pseudocode for a simple syntactic analysis that warns when string concatenation occurs in a loop
  - In Java and .NET it is more efficient to use a StringBuffer
  - Assume any appropriate AST elements

institute for
SOFTWARE
RESEARCH

**Carnegie Mellon University**
School of Computer Science

# WHILE abstract syntax

- Categories:
  - $S \in$ **Stmt**           statements
  - $a \in$ **Aexp**           arithmetic expressions
  - $x, y \in$ **Var**           variables
  - $n \in$ **Num**           number literals
  - $P \in$ **BExp**           boolean predicates
  - $l \in$ **labels**           statement addresses (line numbers)

- Syntax:
  - ```
    S    ::= x := a | skip | S₁ ; S₂
    ```
  - ```
         |   if P then S₁ else S₂ | while P do S
    ```
  - ```
    a    ::= x | n | a₁ opₐ a₂
    ```
  - ```
    opₐ ::= + | - | * | / | …
    ```
  - ```
    P    ::= true | false | not P | P₁ op_b P₂ | a1 op_r a2
    ```
  - ```
    op_b ::=  and | or | …
    ```
  - ```
    op_r ::= < | ≤ | = | > | ≥ | ...
    ```

# WHILE3ADDR:
# An Intermediate Representation

- Simpler, more uniform than WHILE syntax

- Categories:
  - $I$ ∈ **Instruction** instructions
  - $x, y$ ∈ **Var**          variables
  - $n$ ∈ **Num**          number literals

- Syntax:
  - $I$ ::= $x$ := $n$ | x := y | $x$ := $y$ $op$ $z$
          |   goto $n$ | if $x$ $op_r$ 0 goto $n$
  - $op_a$ ::= + | - | * | / | …
  - $op_r$ ::= < | ≤ | = | > | ≥ | ...
  - P ∈ **Num** → $I$

# Practice: Translating to WHILE3ADDR

- Categories:
  - $I \in$ **Instruction** instructions
  - $x, y \in$ **Var**          variables
  - $n \in$ **Num**          number literals

- Syntax:
  - $I$ ::= $x$ := $n$ | x := y | $x$ := $y$ $op$ $z$
    |  goto $n$ | if $x$ $op_r$ 0 goto $n$
  - $op_a$ ::= + | - | * | / | …
  - $op_r$ ::= < | ≤ | = | > | ≥ | ...
  - P $\in$ **Num** $\rightarrow$ $I$

# While3Addr Extensions (more later)

- Syntax:
  - $I$ ::= $x := n$ | x := y | $x := y \; op \; z$
    | goto $n$ | if x op$_r$ 0 goto $n$

    | $x := f(y)$
    | return x
    | $x := y.m(z)$
    | $x := \&p$
    | $x := {}^*p$
    | $^*p := x$
    | $x := y.f$
    | $x.f := y$

# For next time

- Get on Piazza and Canvas

- Answer the survey (location, time zone, in-person interest) we will send you!

- Read lecture notes and the course syllabus

- Homework 1 will be released later this week, and is due next Thursday.

- Discussion: what works well for remote/hybrid instruction?
  - Suggestions for Lecture?  Recitations?  Homework?
  - Feel free to forward suggestions after class too

**Carnegie Mellon University**
School of Computer Science

institute for
SOFTWARE
RESEARCH