



# EDAP15: Program Analysis

#### **Christoph Reichenbach**

### Welcome!

#### EDAP15: Program Analysis

- Instructor: Christoph Reichenbach christoph.reichenbach@cs.lth.se
- Teaching Assistant: Noric Couderc noric.couderc@cs.lth.se

#### Course Homepage: http://cs.lth.se/EDAP15

### **Course Format**

#### Completely Online

- Lectures
  - Lectures are recorded, available later
  - Later lectures may be *flipped*, will announce
  - Questions:
    - Ask in Zoom text chat
    - Online forum
- Online Quizzes
- Group Exercises
- Oral Exam

# Topics

- Concepts and techniques for understanding programs
  - Analysing program structure
  - Analysing program behaviour
- ► Language focus: Teal, a teaching language
  - Concepts generalise to other mainstream languages:
    - Imperative
    - Object-Oriented

# Goals

#### Understand:

- What is program analysis (not) good for?
- What are strenghts and limitations of given analyses?
- How do analyses influence on each other?
- ▶ How do programming language features influence analyses?
- What are some of the most important analyses?

### Be able to:

- Implement typical program analyses
- Critically assess typical program analyses

### Resources

### Course website (http://cs.lth.se/EDAP15)

- Links to everything listed here
- List of expected skills
- Slides
- Announcements
- Textbooks
- Moodle
  - Quizzes
  - Videos
  - Forum

### Course git (GitLab)

Homework assignments

# Textbooks

### Static Program Analysis

Møller & Schwartzbach

- Optional
- PDF online from authors

### Principles of Program Analysis

Nielson, Nielson & Hankin



- Optional
- 3 copies in the library
- Theory-driven

# How to Pass

### > 2020-11-04 18:00: Form Groups of 2

- Must be completed this Wednesday, 18:00
- Contact Noric if you can't find a partner

### Homework projects:

- Who: Groups
- What: Implement program analyses in Teal
  - HW1/HW6 are standalone
  - HW2–4 build on each other
- Start: Homework up every Friday
- Grading:
  - Submit solutions in course git
  - Explain solution to TA
- Deadline:
  - ► HW1-5: Thursday 20:00, 13 days after homework is up
  - **HW6**: 2021-01-12, 20:00
- Final Exam starting 2021-01-15
  - Admission: Passed all homework projects
  - Format: oral exam

# Uses of Program Analysis



# **Categories of Program Analyses**

	Manual / Interactive	Automatic
Static Analysis		
<ul> <li>Examines structure</li> <li>Sees entire program (mostly)</li> </ul>	<ul> <li>Interactive Theorem Provers</li> </ul>	<ul> <li>(Most) Type Checkers</li> <li>Static Checkers (FindBugs, SonarQube,)</li> <li>Compiler Optimisers</li> </ul>
Dynamic Analysis		
<ul> <li>Examines behaviour</li> <li>Sees interactions program ↔ world</li> </ul>	► Debuggers	<ul> <li>Unit Tests</li> <li>Benchmarks</li> <li>Profilers</li> </ul>
		Our Focus

# Summary

- Program analyses used in Software Tools:
  - ► IDEs
  - Compilers
  - Bug Checkers
  - Run-time systems

• • •

- Main Categories:
  - Static Analysis:

Examine program structure

Dynamic Analysis:

Examine program run-time behaviour

Automatic Analysis:

"Black Box": Minimal user interaction

### Manual / Interactive Analysis:

User in the loop

Advanced manual analyses exploit automatic analysis

# **Examples of Program Analysis**

Questions:

```
'Is the program well-formed?'
```

```
gcc -c program.c
javac Program.java
At least for C C++ lava
```

At least for C, C++, Java; not so easy for JavaScript!

'Does my factorial function produce the right input in the range 0–5?'

#### Java

```
@Test // Unit Test
public void testFactorial() {
    int[] expected = new int[] { 1, 1, 2, 6, 24, 120 };
    for (int i = 0; i < expected.length; i++) {
        assertEquals(expected[i], factorial(i));
    } }</pre>
```

# Let's Analyse a Program!

At least 2 of 3 resuls were wrong: "False Positives"

# A First Challenge, Continued

```
user@host$ grep 'gets(' program.c
    gets(input_buffer);
user@host$
```

- More precise: no false positives!
- Will this catch all calls to gets?

```
C: program2.c
#include <stdio.h>
void f(char* target_buffer) {
    char *(*dummy)(char*) = gets;
    dummy(target_buffer);
}
```

String search not smart enough: "False Negative"

# A First Challenge, Continued Again

```
C: program2.c
  #include <stdio.h>
  void f(char* target_buffer) {
     char *(*dummy)(char*) = gets;
     dummy(target buffer);
  }
user@host$ cc -c program.c -o program.o
user@host$ nm program.o
                # check symbol table in compiled program
0000000000000000 T f
                 U gets ← Aha!
                 U _GLOBAL_OFFSET_TABLE_
 user@host$
```

Using a more powerful analysis yielded better results

# A First Challenge, Solved?

#### C: program3.c

Dynamic library loading: gets will not show up in symbol table

Fancier program  $\implies$  harder analysis

# Analysis vs. Property-of-Interest

#### Distinguish:

- Property of interest: P
  - All lines that reference gets
  - Does the program type-check?
  - Where does the program spend most execution time?
- Analysis  $\mathcal A$  that should find P
  - Want for all programs  $\varphi$ :

$$P(\varphi) = \mathcal{A}(\varphi)$$

- Key questions:
  - What data does P depend on?
  - How can A compute P from that data?
  - ▶ How can A access that data?
    - ► Can A access all relevant data?

# And How Good Is It?

- ► As we saw, program analyses may be incorrect
- ► We often describe them with *Information Retrieval* terminology:

	$\mid \mathcal{A}(arphi)$	not $\mathcal{A}(arphi)$		
$P(\varphi)$	True Positive	False Negative		
not $P(\varphi)$	False Positive	True Negative		
▶ #Reports = #True Positives + #False Positives				
Precision = #True Positives #Reports "What fraction of our reports was in P?"				
► Recall = #True Positives #True Positives #True Positives #False Negatives "What fraction of P did we find?"				

#False Negatives is usually impossible to determine in program analysis

# Summary

- Purpose of program analysis  $\mathcal{A}$ :
  - Compute Property-of-interest P
- Program Analysis is nontrivial
  - $\blacktriangleright$  Programs can hide information that  ${\cal A}$  wants
  - $\blacktriangleright$  Analysis  ${\cal A}$  can misunderstand parts of the program

# Soundness and Completeness

Can we always build a  $\mathcal{A}$  with  $\mathcal{A}(\varphi) = P(\varphi)$ ?

Connection to Mathematical Logic:

- Assume P is boolean
- $\mathcal{A}$  is **sound** (with respect to P) iff:

 $\mathcal{A}(\varphi) \implies \mathcal{P}(\varphi)$ 

(Perfect Precision)

• A is **complete** (with respect to *P*) iff:

 $P(\varphi) \implies \mathcal{A}(\varphi)$  (Perfect Recall)

•  $\mathcal{A}(\varphi) = \mathcal{P}(\varphi)$  iff  $\mathcal{A}$  is both sound & complete

What if  $P(\varphi)$  checks whether  $\varphi$  terminates?

### The Bottom Line

# "Everything interesting about the behaviour of programs is undecidable."

- H.G. Rice [1953], paraphrased by Anders Møller

We must choose:

- Soundness
- Completeness
- Decidability
- ... pick any two.

### Soundness and Completeness: Caveat



▶ *Beware*: "sound" and "complete" be confusing:

- Example:  $P(\varphi)$  is " $\varphi$  has a bug"
- If you now want to check  $\overline{P}$ , the *negation* of *P*:
  - $\overline{P}(\varphi)$  is " $\varphi$  does not have a bug"
  - ▶  $\overline{\mathcal{A}_{\text{complete}}}$  (= run  $\mathcal{A}_{\text{complete}}$  and invert output) is sound wrt  $\overline{P}$

### Soundness and Completeness: Caveat



- ▶ *Beware*: "sound" and "complete" be confusing:
  - Example:  $P(\varphi)$  is " $\varphi$  has a bug"
  - If you now want to check  $\overline{P}$ , the *negation* of P:
    - $\overline{P}(\varphi)$  is " $\varphi$  does not have a bug"
    - $\overline{\mathcal{A}_{\text{complete}}}$  (= run  $\mathcal{A}_{\text{complete}}$  and invert output) is sound wrt  $\overline{P}$
    - $\overline{\mathcal{A}}_{\text{sound}}$  is complete wrt  $\overline{P}$

### Summary

Given property P and analysis A:
A is sound if it triggers only on P P = "program has bug": A reports only bugs
A is complete if it always triggeres on P P = "program has bug": A reports all bugs
If P is nontrivial (i.e., depend on behaviour):



# **Building a Program Analysis**



# **Building a Program Analysis**



# Gathering Our Tools

Language Tools Theories Definitions Compilers Principles THE of Program Analysis - C 🗶 jastadd **PhASAR** Analysis Frame-Astrée FindBugs works The Java® Language Specification oot Java SE 8 Edition James Gosling Bill Joy Guy Steele Gilad Bracha Alex Backley 2015-02-13 Hardware

# Language Definitions

- Pure theory
- Define structure (syntax) and meaning (semantics) of language
- Abstracts over many details
- Syntax example:

e ::= zero | one  $| \langle e \rangle + \langle e \rangle$   $| \langle e \rangle - \langle e \rangle$   $| neg \langle e \rangle$   $| (\langle e \rangle)$   $| log \langle e \rangle$ 

▶ **Property of Interest**: Does a given program  $\varphi \in e$  compute a positive number?

First, we must understand the language semantics

# Language Definitions: Semantics

### Language Definitions also specify *Semantics*:

- Static Semantics:
  - Connect parts of the program structure (variables, functions, classes, ...)
  - Enforce restrictions (e.g., via type checking)

#### Dynamic Semantics:

- Specify program run-time behaviour
- ▶ We will (mostly) treat semantics informally in this course
- ► Here: assume "obvious" semantics

# Simplifying the Lanugage

- Let's make it easier to analyse the language
- ▶ We don't need parentheses for the analysis
- ▶ *a*-*b* = *a*+neg *b* 
  - $\Rightarrow$  Abstraction (we join similar problems into one)
- log is too difficult
  - $\Rightarrow$  Restrict to sub-language (give up on some problems)

$$e ::= zero$$

$$| one$$

$$| \langle e \rangle + \langle e \rangle$$

$$| neg \langle e \rangle$$

Simplification helps us get started, but restricting to a sublanguage can quickly render an analysis impractical

# Finding a Good Theory

- ▶ Recall:  $P(\varphi)$  should detect if  $\varphi$  computes positive result
- There are many theories for program analysis
- ▶ We pick Abstract Interpretation (Patrick & Radhia Cousot):
  - Map all values to a simpler abstract domain
  - Map all operations so they respect the abstraction
- For example: classify programs into abstract domain containing:
  - ▶  $D^0$ : Computes 0
  - ► D<sup>+</sup>: Computes a positive value
  - ► D<sup>-</sup>: Computes a negative value

▶ Notation:  $\varphi \rightsquigarrow^{D} a$ , where *a* is one of  $D^{0}$ ,  $D^{+}$ ,  $D^{-}$ 

**Semantics** 

$$\begin{array}{rcl} \ominus D^{0} & = & D^{0} \\ \ominus D^{+} & = & D^{-} \\ \ominus D^{-} & = & D^{+} \\ \ominus D^{?} & = & D^{?} \end{array} \end{array}$$

$$\begin{array}{rcl} a_{1} \oplus a_{2} & = & \left\{ \begin{array}{c|c} & D^{+} & D^{0} & D^{-} \\ \hline D^{+} & D^{+} & D^{+} & D^{?} \\ D^{0} & D^{+} & D^{0} & D^{-} \\ D^{-} & D^{?} & D^{-} & D^{-} \end{array} \right\}$$

$$\begin{array}{rcl} D^{?} \oplus a = D^{?} = a \oplus D^{?} \end{array}$$

 $\operatorname{zero} \rightsquigarrow^{\scriptscriptstyle D} D^0$  one  $\rightsquigarrow^{\scriptscriptstyle D} D^+$ 

 $\frac{\text{if } \mathbf{x} \rightsquigarrow^{D} \mathbf{a} \text{ then}}{\text{neg } \mathbf{x} \rightsquigarrow^{D} \ominus \mathbf{a}} \qquad \qquad \frac{\text{if } \mathbf{x} \rightsquigarrow^{D} \mathbf{a}_{1} \text{ and } \mathbf{y} \rightsquigarrow^{D} \mathbf{a}_{2} \text{ then}}{\mathbf{x} + \mathbf{y} \rightsquigarrow^{D} \mathbf{a}_{1} \oplus \mathbf{a}_{2}}$ 

# **Correspondence: Concrete and Abstract**



#### Also:

- ► ⊖ *"is compatible with"* neg
- $\oplus$  "is compatible with" +

Abstract Interpretation explores these ideas in great detail

### **Building a Better Analysis**

**Property of Interest**: Does a given program  $\varphi \in e$  compute a positive number?

$$e ::= zero$$

$$| one$$

$$| \langle e \rangle + \langle e \rangle$$

$$| \langle e \rangle - \langle e \rangle$$

$$| neg \langle e \rangle$$

$$| (\langle e \rangle)$$

$$| log \langle e \rangle$$

• How else could we analyse this for a given program  $\varphi$ ?

- Just run  $\varphi$  and check the result
- In general: only feasible for very restricted languages
- Infeasible as soon as we add parameters or recursion

# Summary

- ▶ We can mathematically formalise syntax and semantics
- Semantics derive from syntax
- Abstract Interpretation: Theory for program analysis
  - Map program semantics into abstract domain
  - ► Map operations to *compatible* operations on abstract domain
  - Challenge: remain precise yet decidable
  - Foundation to other static analysis theories

# Some Theories: Dynamic

- Testing
- Input Generation and Fuzzing Call program with "random" inputs to see if it crashes
  - Dynamic Invariant Detection Record program behaviour, describe concisely (e.g., "always returns zero")
  - Test Generation

Automatically build test suite to capture current behaviour

#### Monitoring

Capture which parts of the code execute

#### Benchmarking

Measure resource usage (time, memory, bandwidth)

# Some Theories: Static

- Abstract Interpretation
  - Type Theory
    - 3 : Int
  - Dataflow Analysis Theory What values can my variables store?
  - ► Hoare Logic

If x > y holds before I run  $\mathbf{x} := \mathbf{x} \cdot \mathbf{1}$ , is  $x \ge y$  true?

Symbolic Execution

Use SMT solving to find properties of variables and possible execution paths

Model Checking

Given some formal model (e.g., a FSM), does it describe my code?

Statistical Analysis









36 / 48

# Static vs. Dynamic Program Analyses

	Static Analysis	Dynamic Analysis
Principle	Analyse program	Analyse program execution
	structure	
Input	Independent	Depends on input
Hardware/OS	Independent	Depends on hardware and OS
Perspective	Sees everything	Sees that which actually happens
<b>Completeness</b> (bug-finding)	Possible	Must try all possible inputs
Soundness (bug-finding)	Possible	Always, for free





# Summary

- **Preprocessor**: Transforms source code before compilation
- Static compiler: Tranlates source code into executable (machine or intermediate) code
- Interpreter: Step-by-step execution of source or intermediate code
- Dynamic (JIT) compiler: Translates code into machine-executable code
- Loader: System tool that ensures that OS starts executing another program
- Linker: System tool that connects references between programs and libraries
  - Static linker: Before running
  - Dynamic linker: While running
- Machine code: Code that is executable by a machine
- **Static Analysis**: Analyse program without executing it
- **Dynamic Analysis**: Analyse program execution

# Java lexing



# Java lexing & parsing



# Parsing in general

Translate text files into meaningful in-memory structures

- CST = Concrete Syntax Tree
  - ► Full "parse", cf. language BNF grammar
  - Not usually materialised in memory
- AST = Abstract Syntax Tree
  - Standard in-memory representation
  - Avoids syntactic sugar from CST, preserves important nonterminals as (AST nodes)
  - Converts useful tokens into attributes
- ► The AST is the most common **Intermediate Representation** (IR) of program code
  - Effective for frontend analyses
  - Other IRs focus e.g. on optimisations in the backend

#### Program analysis starts on the AST

# **In-Memory Representation**



Typical in-memory representations for this AST:

- Algebraic values (functional)
- Records (imperative)

# **In-Memory Representation**



# **Program Analysis**

We run numerous code analyses on the AST:

- Name Analysis:
  - ▶ Which name *use* binds to which *declaration*?
- Type Analysis:
  - What are the types of all expressions?
- Static Correctness Checks:
  - Are there type errors?
  - Is a variable unused?
  - Are we initialising all variables?
- Optimisations:
  - Can we speed up the program somehow?

#### Advanced static correctness checks increasingly common in compilers

# Name Analysis



### **Type Analysis**



# Summary

- Compiler represents programs in intermediate representations (IRs)
- Compiler can be separated into:
  - ► Frontend: process incoming source code, generate IR
  - Middle-end: optimise IR
  - Back-end: translate IR into executable code
- Parser matches concrete syntax tree (CST), generates abstract syntax tree (AST)
- ► Typical analyses on AST:
  - ► Name analysis: which variable use belongs to which definition?
  - Type analysis: do variable/operator/function types agree? Any implicit conversions needed?

. . .

### Outlook

#### Remember:

- Join Moodle if you have not done so yet
- ▶ Form groups by Wednesday, 18:00!
- ▶ Next Lecture (Wednesday, same channel):
  - Type-Based Analysis

http://cs.lth.se/EDAP15