

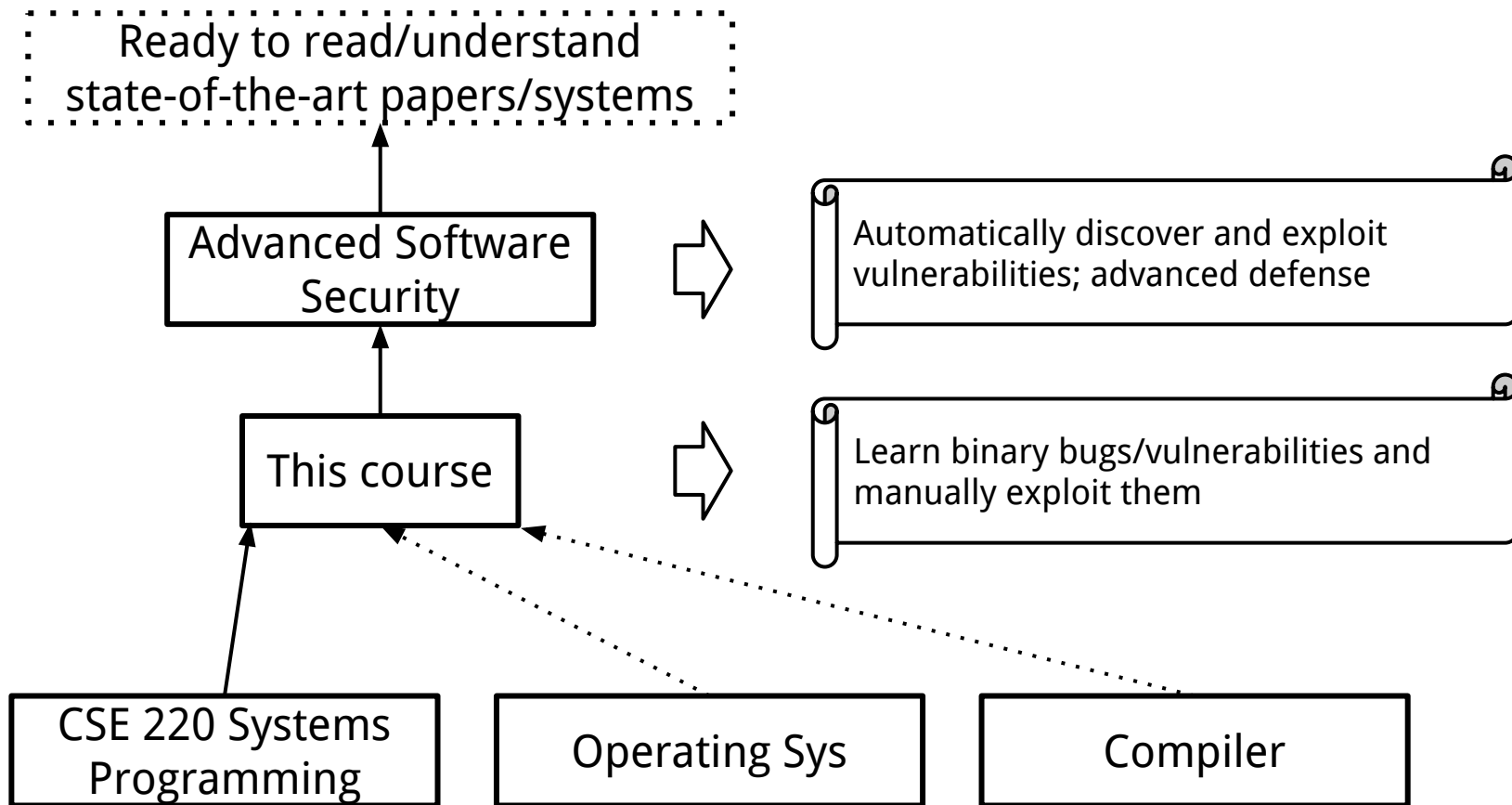
CSE 410/510 Special Topics: Software Security

Instructor: Dr. Ziming Zhao

Location: Obrian 109

Time: Monday, Wednesday 5:00PM-6:20PM

If you want to be a system/software security guy ...



From 410/510 to security research

- Other background knowledge
 - Static analysis
 - Fuzzing
 - Dynamic taint analysis
 - Symbolic execution

Static Analysis

LLVM



What is LLVM?

An open source framework for building tools

- Tools are created by linking together various libraries provided by the LLVM project and your own

An extensible, strongly typed intermediate representation, i.e. LLVM IR

- <https://llvm.org/docs/LangRef.html>

An industrial strength C/C++ optimizing compiler

- Which you might know as clang/clang++ but these are really just drivers that invoke different parts (libraries) of LLVM

LLVM: A Compilation Framework for Lifelong Program Analysis & Transformation

Chris Lattner Vikram Adve
University of Illinois at Urbana-Champaign
{lattner,vadve}@cs.uiuc.edu
<http://llvm.cs.uiuc.edu/>

ABSTRACT

This paper describes LLVM (Low Level Virtual Machine), a compiler framework designed to support *transparent, lifelong program analysis and transformation* for arbitrary programs, by providing high-level information to compiler transformations at compile-time, link-time, run-time, and in idle time between runs. LLVM defines a common, low-level code representation in Static Single Assignment (SSA) form, with several novel features: a simple, *language-independent* type-system that exposes the primitives commonly used to implement high-level language features; an instruction for typed address arithmetic; and a simple mechanism that can be used to implement the exception handling features of high-level languages (and `setjmp/longjmp` in C) uniformly

mizations performed at link-time (to preserve the benefits of separate compilation), machine-dependent optimizations at install time on each system, dynamic optimization at run-time, and profile-guided optimization between runs (“idle time”) using profile information collected from the end-user.

Program optimization is not the only use for lifelong analysis and transformation. Other applications of static analysis are fundamentally interprocedural, and are therefore most convenient to perform at link-time (examples include static debugging, static leak detection [24], and memory management transformations [30]). Sophisticated analyses and transformations are being developed to enforce program safety, but must be done at software installation time or load-time [19]. Allowing lifelong reoptimization of the program gives architects the power to evolve processors and

LLVM

LLVM is written in C++; uses STL; vector, set and map

LLVM sources are hosted on GitHub

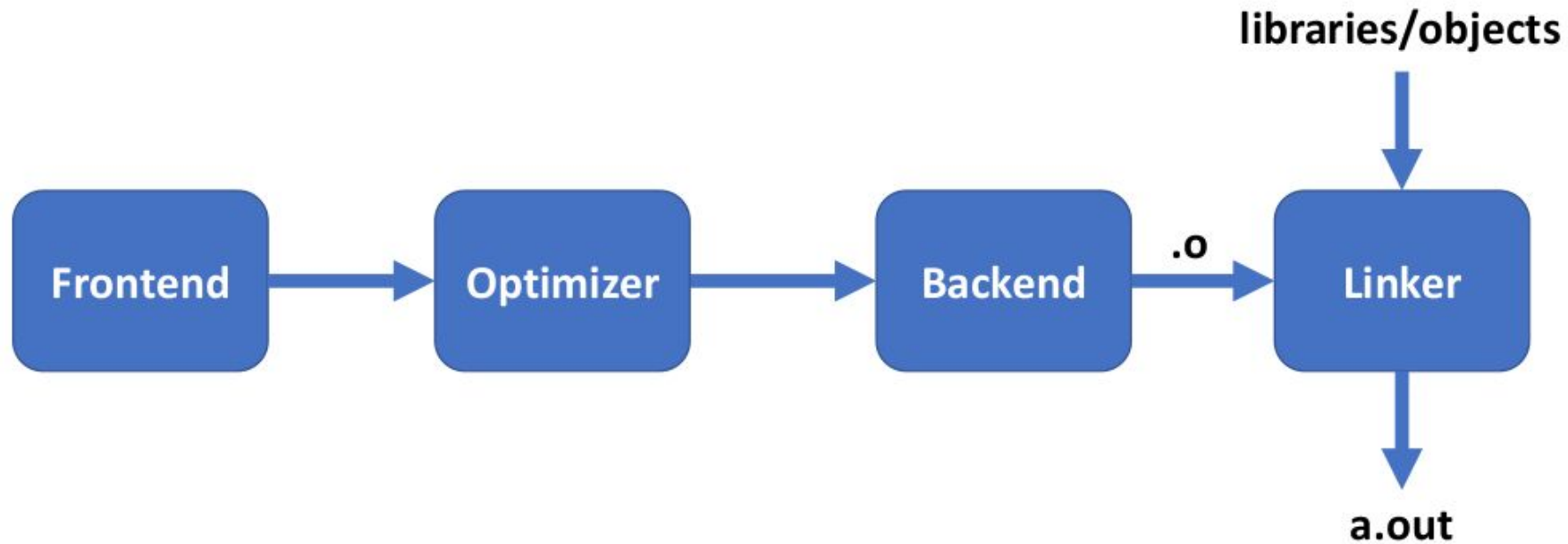
<https://github.com/llvm/llvm-project>

LLVM is split into multiple Git repositories

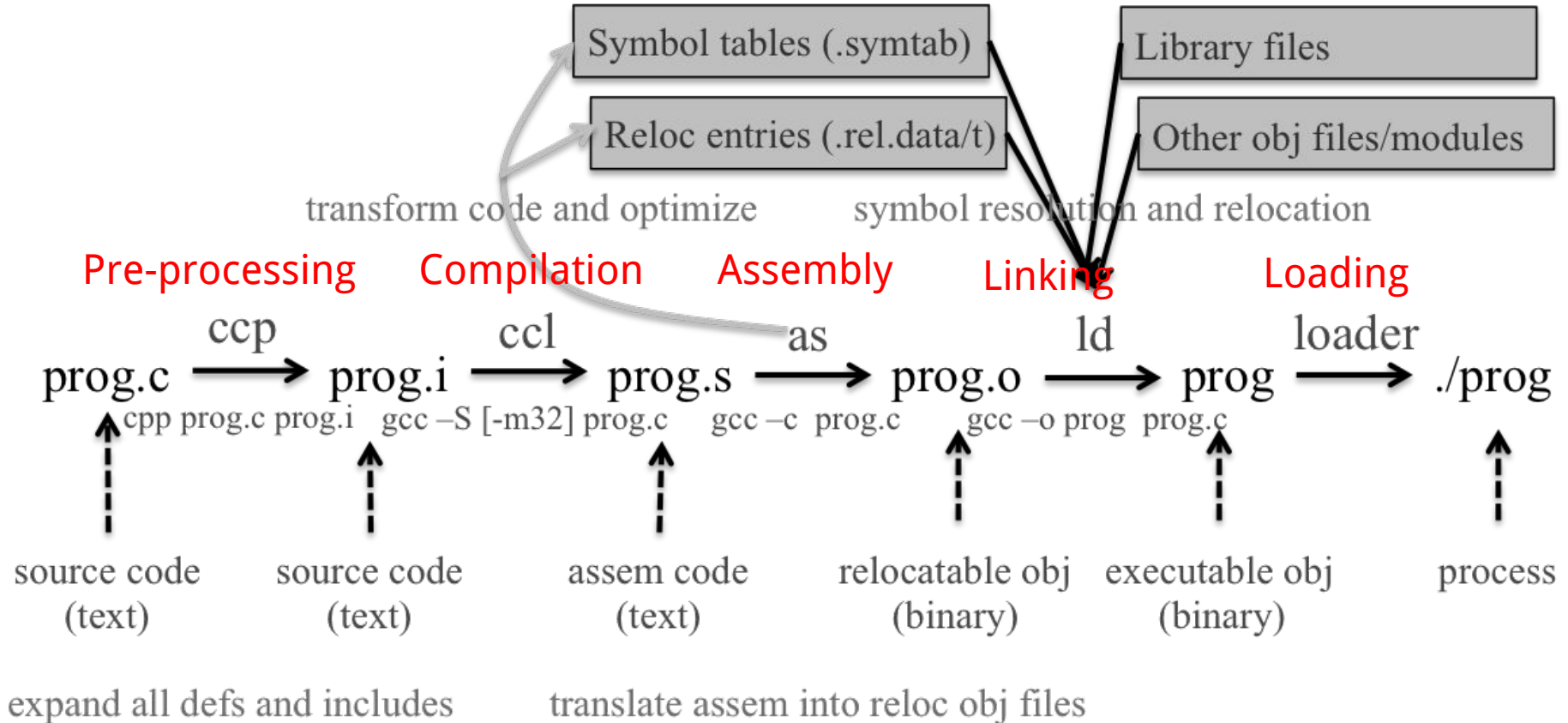
- For this class you will need the clang and llvm git repos

<https://llvm.org/>

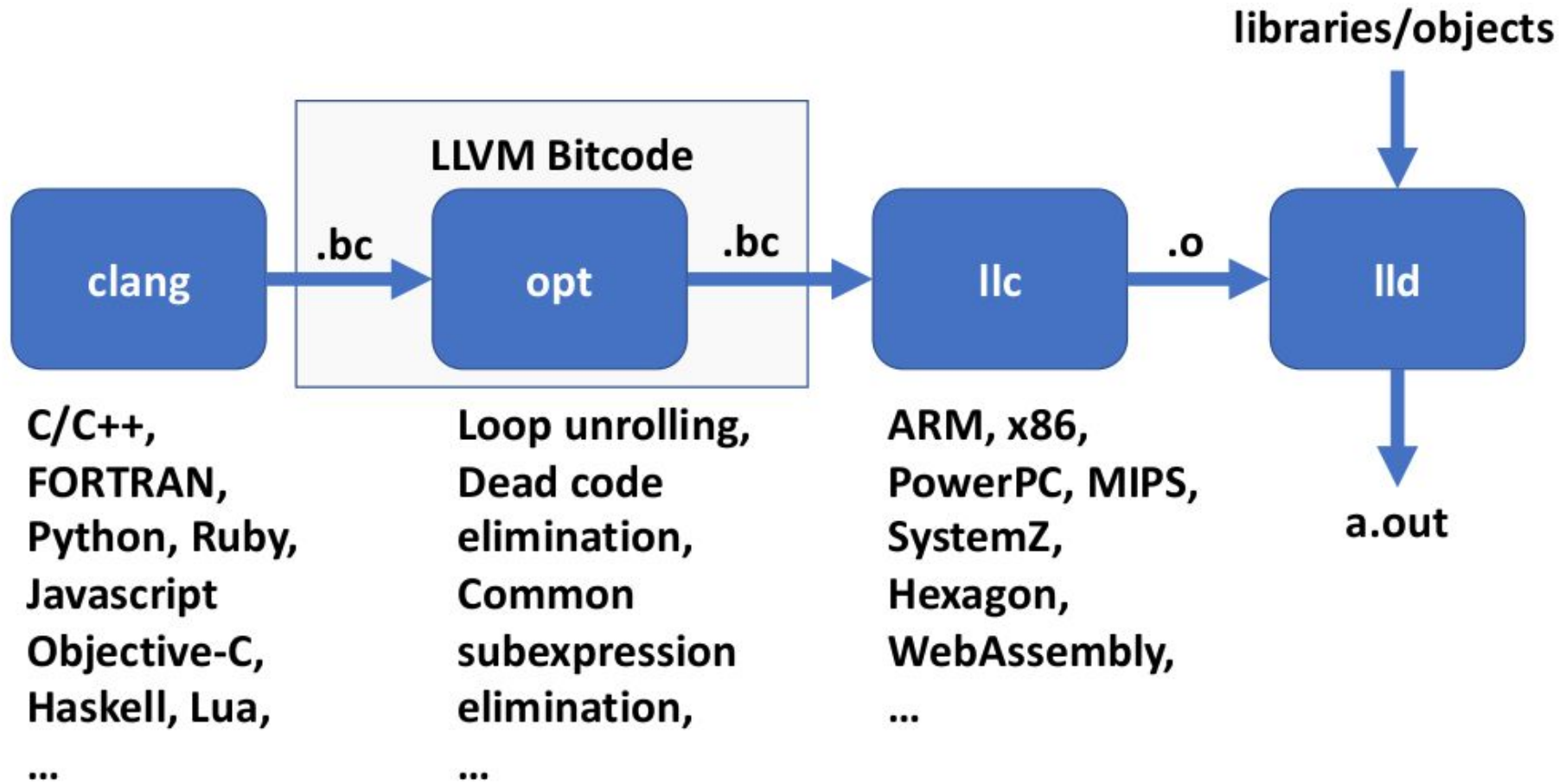
Typical Compiler Flow



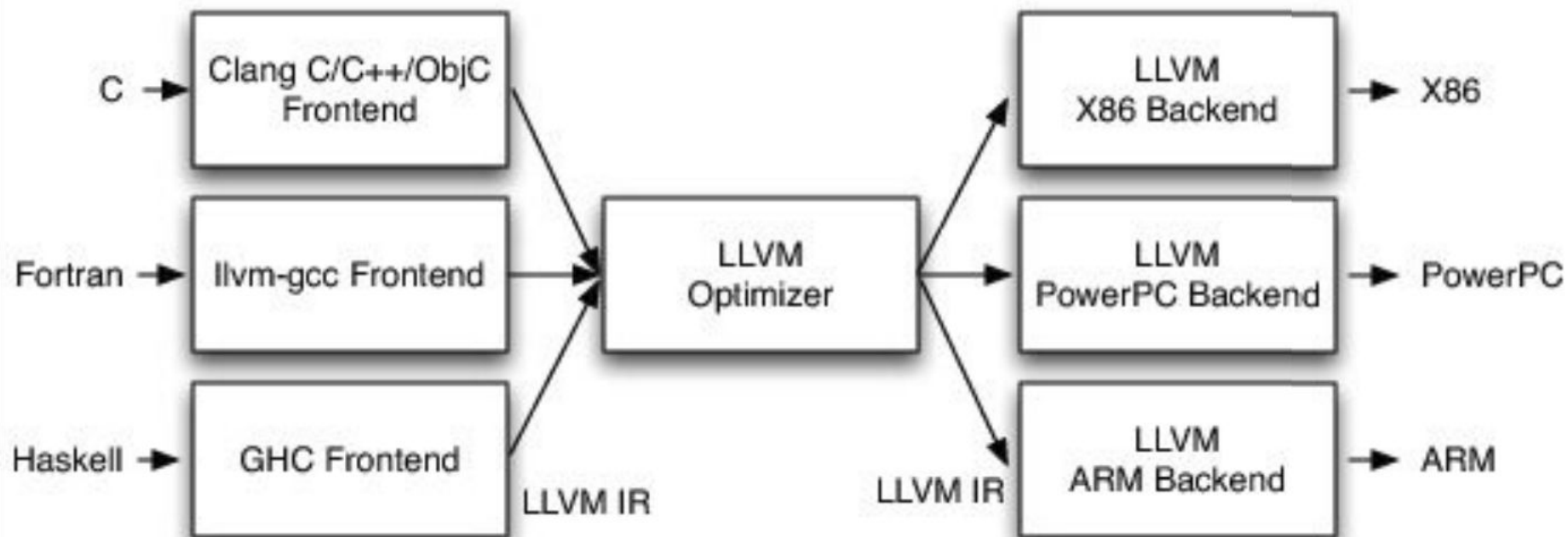
From a C program to a process



LLVM Flow



LLVM Flow



Clang/Clang++

Clang is a frontend for several C-family languages

- C and C++ being the most widely known
 - Supports C++11/14/17/20
- (Objective C/C++, OpenCL, CUDA< and RenderScript are the other C-style languages actively developed)

LLVM IR / LLVM Instruction Set

The LLVM Intermediate Representation

Some characteristics of LLVM IR

- RISC-like instruction set (3 addresses; human readable, assembly like)
- Strongly typed
- Explicit control flow
- Uses a virtual register set with infinite temporaries (%)
- In Static Single Assignment form
- Abstracts machine details such as calling conventions and stack references

LLVM IR reference is online

- <https://llvm.org/docs/LangRef.html>

The LLVM Intermediate Representation

LLVM IR is actually defined in three isomorphic forms

- the textual format above
- an in-memory data structure inspected and modified by optimizations themselves
- an efficient and dense on-disk binary "bitcode" format (.bc)

The LLVM Project also provides tools to convert the on-disk format from text to binary

- `llvm-as` assembles the textual .ll file into a .bc file containing the bitcode
- `llvm-dis` turns a .bc file into a .ll file.

Static Single Assignment (SSA) form

In compiler design, static single assignment form (often abbreviated as SSA form or simply SSA) is a property of an intermediate representation (IR), which requires that each variable be assigned exactly once, and every variable be defined before it is used.

SSA was proposed by Barry K. Rosen, Mark N. Wegman, and F. Kenneth Zadeck in POPL 1988

Static Single Assignment (SSA) form

```
y := 1  
y := 2  
x := y
```

Not SSA

```
y1 := 1  
y2 := 2  
x1 := y2
```

SSA

Different Types of Passes in LLVM

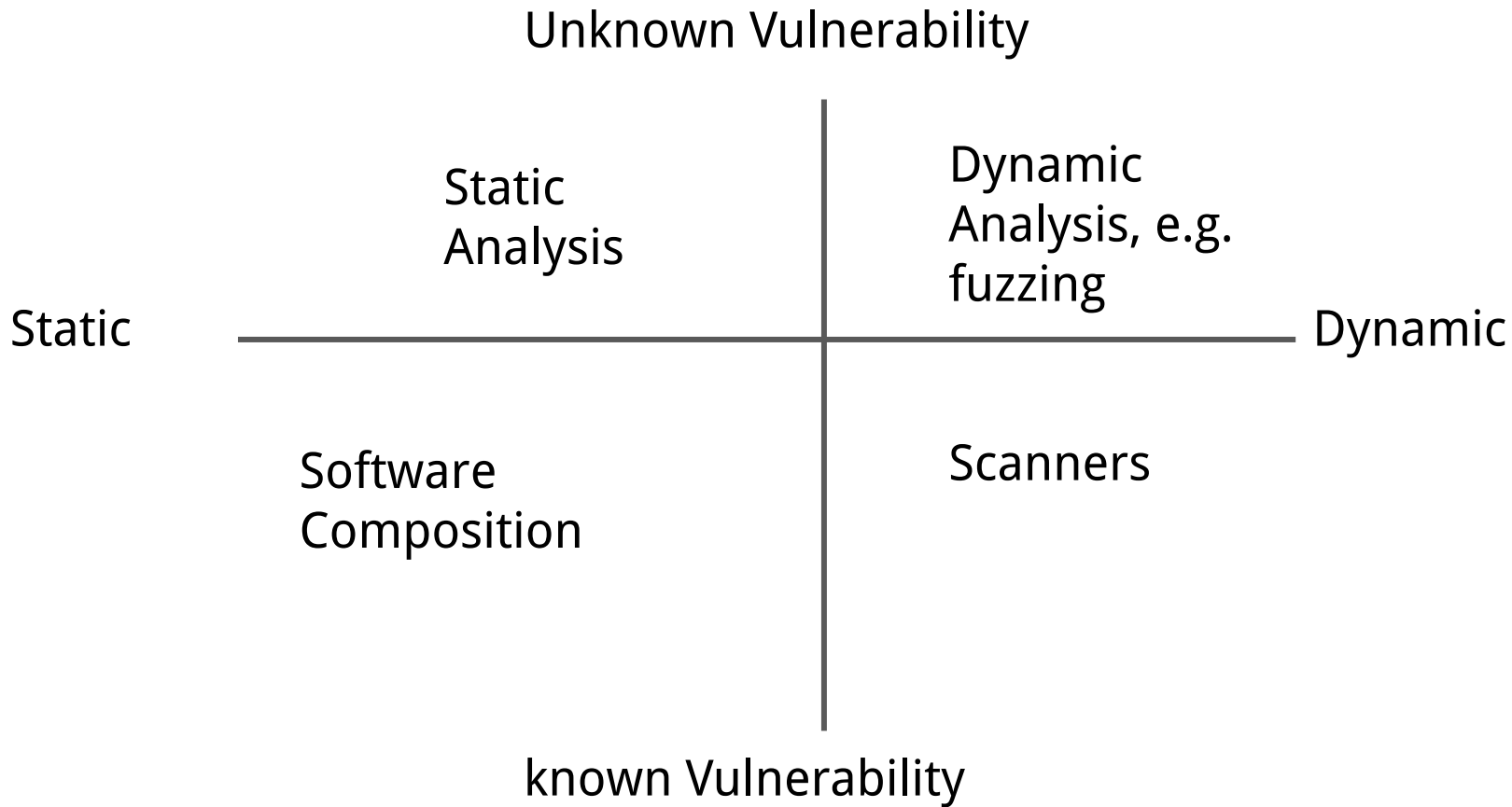
- Levels of Granularity
 - Module Pass - Can think of this as a single source file
 - Call Graph Pass - Traverses a program bottom-up
 - Function Pass - Runs over individual functions
 - Basic Block Pass - Runs over individual basic blocks within a function
 - (Immutable Pass, Region Pass, MachineFunctionPass - Less important for today)
- Analysis Passes versus Transform pass
 - Analysis Pass - Computes information that other passes can use for debugging
 - Transform Pass - Mutates the program. ■ i.e. A side effect occurs, which could invalidate other passes!

LLVM Program Structure

- Module contains Functions/GlobalVariables
 - Module is unit of compilation/analysis/optimization
- Function contains BasicBlocks/Arguments
 - Functions roughly correspond to functions in C
- BasicBlock contains list of instructions
 - Each block ends in a control flow instruction
- Instruction is opcode + vector of operands
 - All operands have types
 - Instruction result is typed

How to (automatically) Find Bugs?

- Manual code inspection - not automatically
 - peer-reviewing code before releasing it
 - pen testing
- Static program analyzers
 - automatically inspect code and flag unexpected code patterns
- Fuzzing
 - Fuzzing repeatedly **executes** an application with all kinds of input variants (dynamic analysis)
 - most effective when applied to standalone applications
 - fuzzing does not generate false alarms

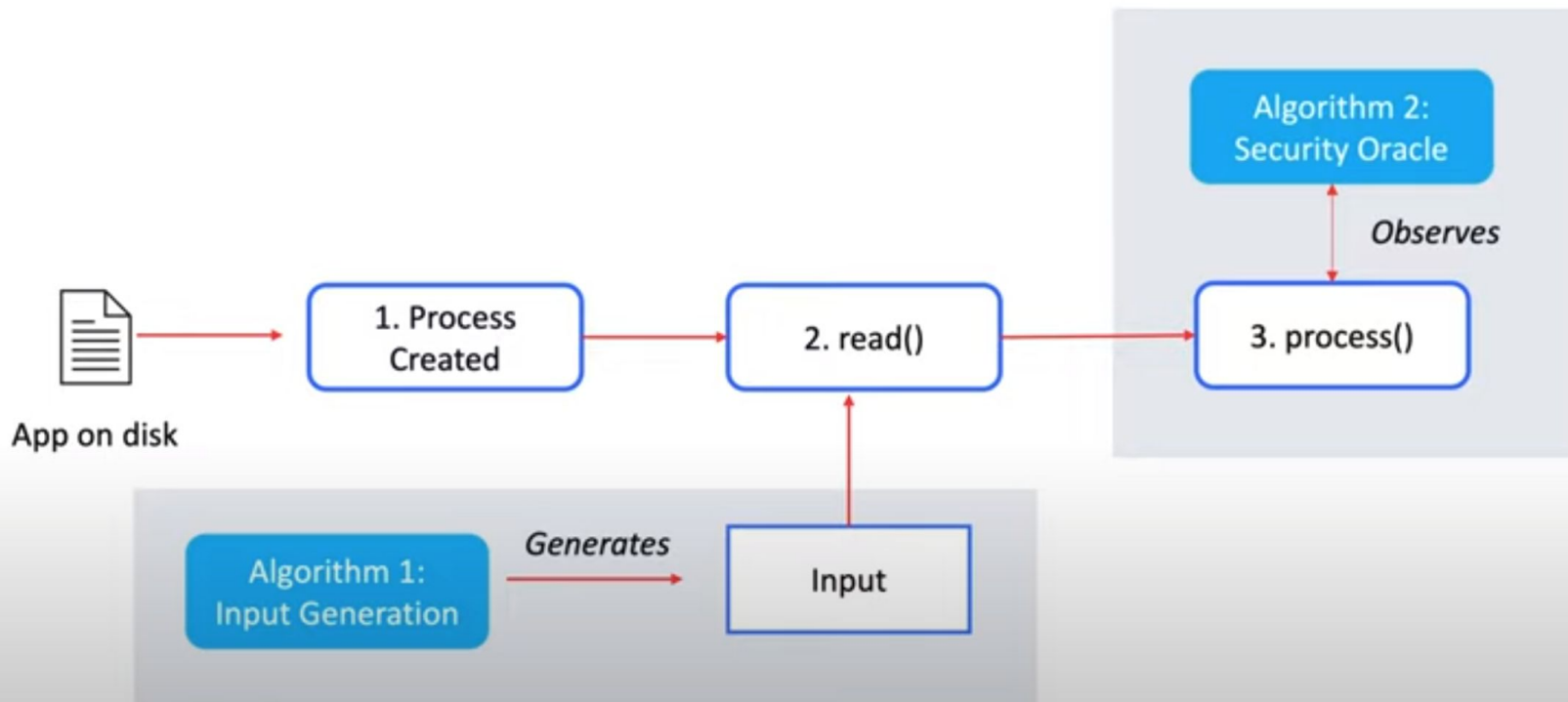


What can be fuzzed? Everything!

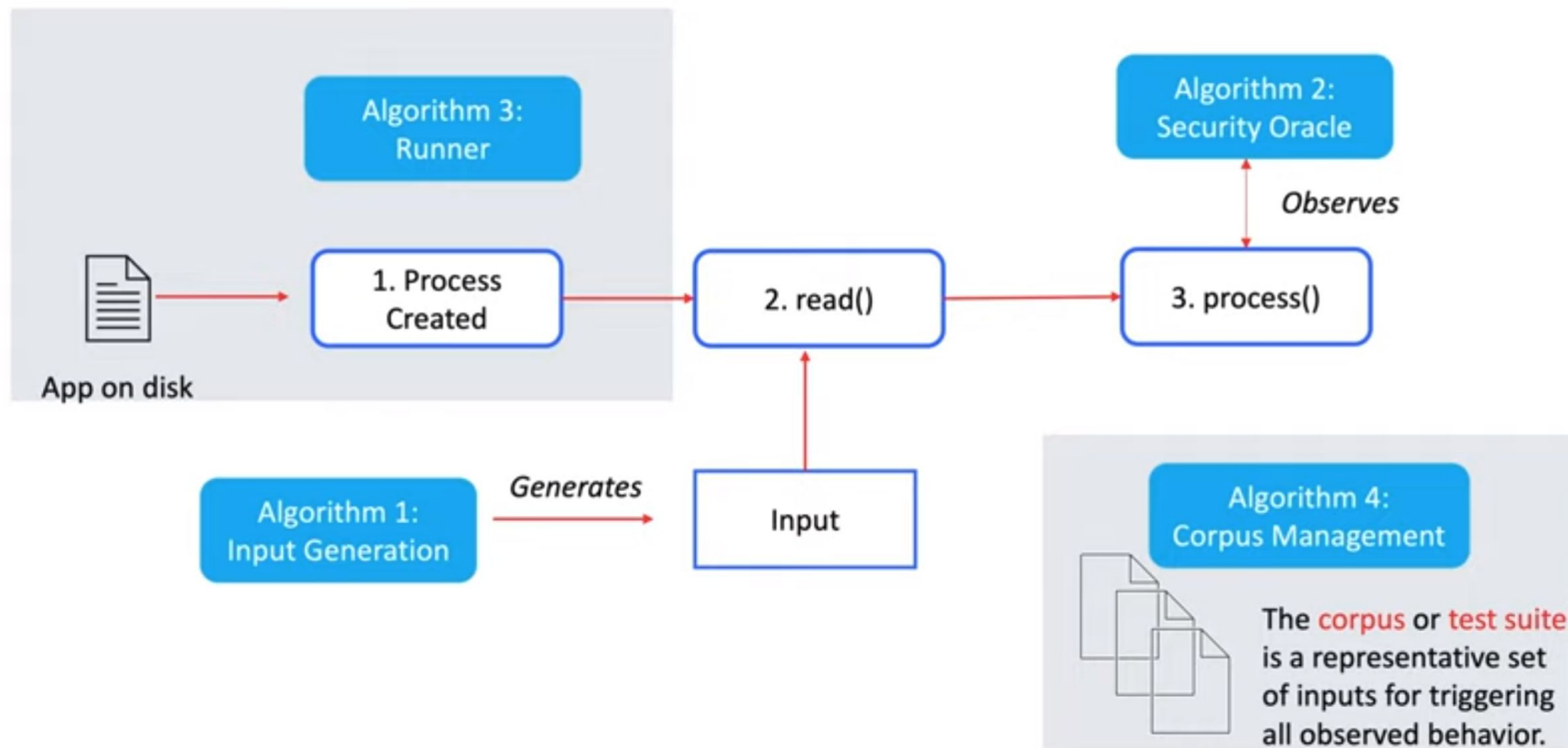
Fuzzing all kinds of software applications as long as they take some kind of input

- Documents
- Images
- Sensor reading, e.g. sound, temperature, etc.
- Videos
- Network packets
- Web pages

Fuzzing has at least two algorithms



With additional algorithms as you get more advanced



Types of Fuzzing

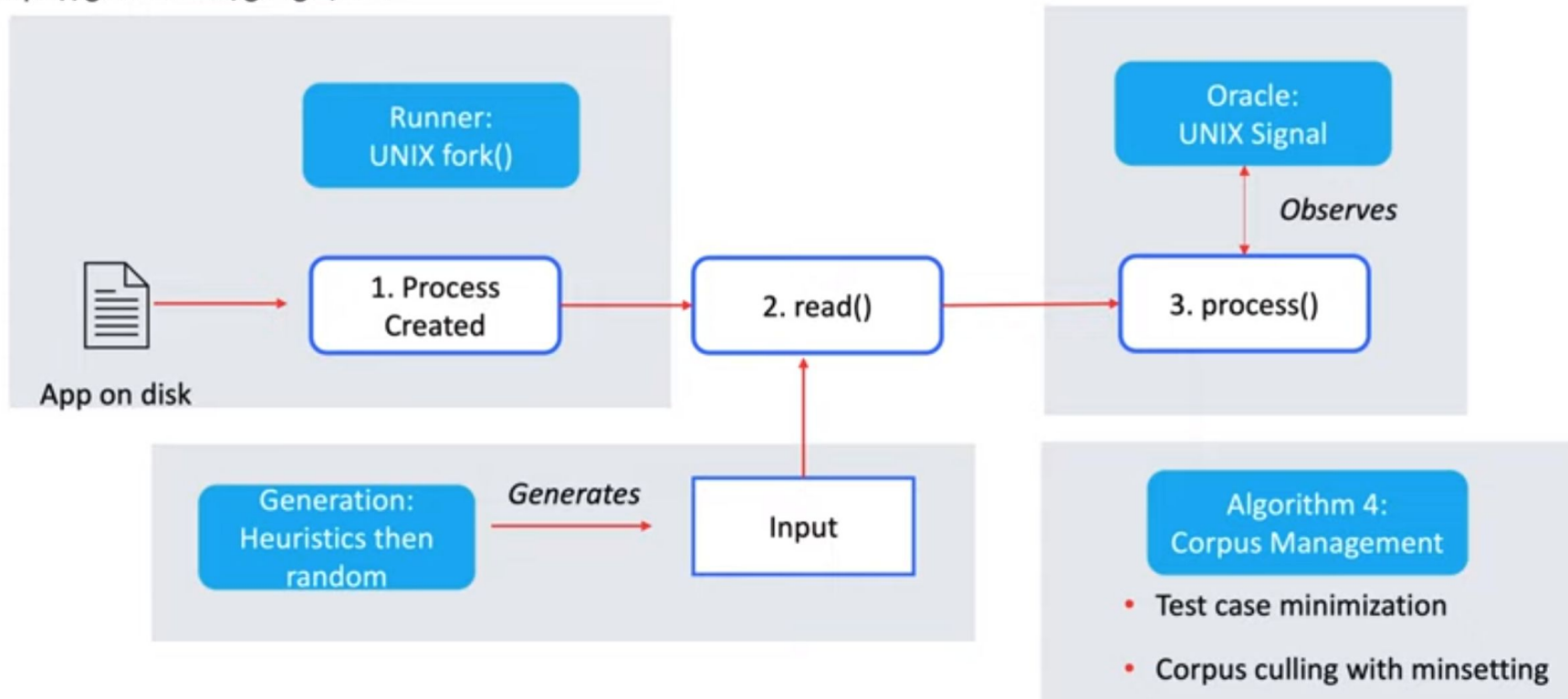
- Blackbox Fuzzing
 - randomly mutates well-formed application inputs, and then tests the application with these modified inputs
- Grammar-Based Fuzzing
 - generates many new inputs satisfying the constraints encoded by a grammar
- Whitebox Fuzzing
 - symbolically executing the program under test dynamically, gathering constraints on inputs from conditional branches encountered along the execution

Types of Fuzzing

- Generation 1: Random input or mutation (1950 - early 2000)
- Generation 2: Protocol/grammar/model fuzzing
- Generation 3: Coverage-based fuzzing
 - AFL, libfuzzer
- Generation 4: Symbolic execution

AFL: American Fuzzy Lop [2013]

<https://github.com/google/AFL>

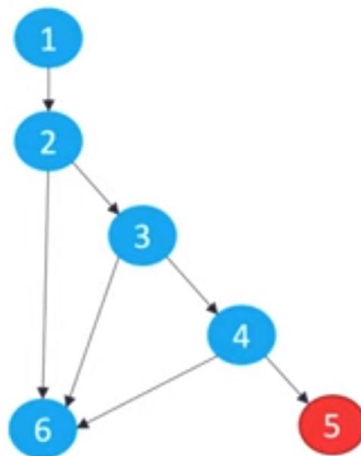


<https://www.youtube.com/watch?v=MYxfDhNa2-U&feature=youtu.be>

Coverage-based fuzzing



```
1. example(char *input){  
2.   if (input[0] == 'b')  
3.     if(input[1] == 'u')  
4.       if(input[2] == 'g')  
5.         crash();  
6.   return;  
7. }
```



A control flow graph $G = (V, E)$ has a vertex for every node statement ($v = \text{set of statements } s_i$), and an edge $E = (s_1, s_2)$ if there is a possible control transfer between statement s_1 and s_2 .

Dynamic Taint Analysis

- Dynamic analysis: monitor code as it executes
- It tracks information flow through a program at runtime between sources and sinks
- It runs a program and observes which computations are affected by predefined taint sources such as user input

Example Use Case

Unknown Vulnerability Detection.

Dynamic taint analysis can look for misuses of user input during an execution. For example, dynamic taint analysis can be used to prevent code injection attacks by monitoring whether user input is executed

<i>program</i>	$::=$	<i>stmt</i> *
<i>stmt s</i>	$::=$	<i>var</i> := <i>exp</i> store(<i>exp</i> , <i>exp</i>) goto <i>exp</i> assert <i>exp</i> if <i>exp</i> then goto <i>exp</i> else goto <i>exp</i>
<i>exp e</i>	$::=$	load(<i>exp</i>) <i>exp</i> \diamond_b <i>exp</i> \diamond_u <i>exp</i> <i>var</i> get_input(<i>src</i>) <i>v</i>
\diamond_b	$::=$	typical binary operators
\diamond_u	$::=$	typical unary operators
<i>value v</i>	$::=$	32-bit unsigned integer

Table I: A simple intermediate language (SIMPIL).

Operational Semantics

computation

$\langle \text{current state} \rangle, \text{stmt} \rightsquigarrow \langle \text{end state} \rangle, \text{stmt}'$

Context	Meaning
Σ	Maps a statement number to a statement
μ	Maps a memory address to the current value at that address
Δ	Maps a variable name to its value
pc	The program counter
ι	The next instruction

Figure 2: The meta-syntactic variables used in the execution context.

We denote by $\mu, \Delta \vdash e \Downarrow v$ evaluating an expression e to a value v in the current state given by μ and Δ . The

$$\begin{array}{c}
\frac{v \text{ is input from } src}{\mu, \Delta \vdash \text{get_input}(src) \Downarrow v} \text{ INPUT} \quad \frac{\mu, \Delta \vdash e \Downarrow v_1 \quad v = \mu[v_1]}{\mu, \Delta \vdash \text{load } e \Downarrow v} \text{ LOAD} \quad \frac{}{\mu, \Delta \vdash \text{var} \Downarrow \Delta[\text{var}]} \text{ VAR} \\
\\
\frac{\mu, \Delta \vdash e \Downarrow v \quad v' = \Diamond_u v}{\mu, \Delta \vdash \Diamond_u e \Downarrow v'} \text{ UNOP} \quad \frac{\mu, \Delta \vdash e_1 \Downarrow v_1 \quad \mu, \Delta \vdash e_2 \Downarrow v_2 \quad v' = v_1 \Diamond_b v_2}{\mu, \Delta \vdash e_1 \Diamond_b e_2 \Downarrow v'} \text{ BINOP} \quad \frac{}{\mu, \Delta \vdash v \Downarrow v} \text{ CONST} \\
\\
\frac{\mu, \Delta \vdash e \Downarrow v \quad \Delta' = \Delta[\text{var} \leftarrow v] \quad \iota = \Sigma[pc + 1]}{\Sigma, \mu, \Delta, pc, \text{var} := e \rightsquigarrow \Sigma, \mu, \Delta', pc + 1, \iota} \text{ ASSIGN} \quad \frac{\mu, \Delta \vdash e \Downarrow v_1 \quad \iota = \Sigma[v_1]}{\Sigma, \mu, \Delta, pc, \text{goto } e \rightsquigarrow \Sigma, \mu, \Delta, v_1, \iota} \text{ GOTO} \\
\\
\frac{\mu, \Delta \vdash e \Downarrow 1 \quad \Delta \vdash e_1 \Downarrow v_1 \quad \iota = \Sigma[v_1]}{\Sigma, \mu, \Delta, pc, \text{if } e \text{ then goto } e_1 \text{ else goto } e_2 \rightsquigarrow \Sigma, \mu, \Delta, v_1, \iota} \text{ TCOND} \\
\\
\frac{\mu, \Delta \vdash e \Downarrow 0 \quad \Delta \vdash e_2 \Downarrow v_2 \quad \iota = \Sigma[v_2]}{\Sigma, \mu, \Delta, pc, \text{if } e \text{ then goto } e_1 \text{ else goto } e_2 \rightsquigarrow \Sigma, \mu, \Delta, v_2, \iota} \text{ FCOND} \\
\\
\frac{\mu, \Delta \vdash e_1 \Downarrow v_1 \quad \mu, \Delta \vdash e_2 \Downarrow v_2 \quad \iota = \Sigma[pc + 1] \quad \mu' = \mu[v_1 \leftarrow v_2]}{\Sigma, \mu, \Delta, pc, \text{store}(e_1, e_2) \rightsquigarrow \Sigma, \mu', \Delta, pc + 1, \iota} \text{ STORE} \\
\\
\frac{\mu, \Delta \vdash e \Downarrow 1 \quad \iota = \Sigma[pc + 1]}{\Sigma, \mu, \Delta, pc, \text{assert}(e) \rightsquigarrow \Sigma, \mu, \Delta, pc + 1, \iota} \text{ ASSERT}
\end{array}$$

Figure 1: Operational semantics of SIMPIL.

Example 1. Consider evaluating the following program:

```
1  x := 2 * get_input(·)
```

The evaluation for this program is shown in Figure 3 for the input of 20. Notice that since the ASSIGN rule requires the expression e in $\text{var} := e$ to be evaluated, we had to recurse to other rules (BINOP, INPUT, CONST) to evaluate the expression $2 * \text{get_input}(\cdot)$ to the value 40.

Example 1. Consider evaluating the following program:

```
1  x := 2 * get_input(·)
```

The evaluation for this program is shown in Figure 3 for the input of 20. Notice that since the ASSIGN rule requires the expression e in $\text{var} := e$ to be evaluated, we had to recurse to other rules (BINOP, INPUT, CONST) to evaluate the expression $2 * \text{get_input}(\cdot)$ to the value 40.

$$\begin{array}{c}
 \frac{}{\mu, \Delta \vdash 2 \Downarrow 2} \text{CONST} \quad \frac{20 \text{ is input}}{\mu, \Delta \vdash \text{get_input}(\cdot) \Downarrow 20} \text{INPUT} \quad v' = 2 * 20 \\
 \hline
 \mu, \Delta \vdash 2 * \text{get_input}(\cdot) \Downarrow 40 \quad \text{BINOP} \quad \Delta' = \Delta[x \leftarrow 40] \quad \iota = \Sigma[pc + 1] \\
 \hline
 \Sigma, \mu, \Delta, pc, x := 2 * \text{get_input}(\cdot) \rightsquigarrow \Sigma, \mu, \Delta', pc + 1, \iota \quad \text{ASSIGN}
 \end{array}$$

Figure 3: Evaluation of the program in Listing 1.

Dynamic Taint Analysis

- Any program value whose computation depends on data derived from a taint source is considered tainted (denoted T)
- A taint policy P determines exactly how taint flows as a program executes, what sorts of operations introduce new taint, and what checks are performed on tainted values.

<i>program</i>	$::=$	<i>stmt</i> *
<i>stmt s</i>	$::=$	<i>var</i> := <i>exp</i> store(<i>exp</i> , <i>exp</i>) goto <i>exp</i> assert <i>exp</i> if <i>exp</i> then goto <i>exp</i> else goto <i>exp</i>
<i>exp e</i>	$::=$	load(<i>exp</i>) <i>exp</i> \diamond_b <i>exp</i> \diamond_u <i>exp</i> <i>var</i> get_input(<i>src</i>) <i>v</i>
\diamond_b	$::=$	typical binary operators
\diamond_u	$::=$	typical unary operators
<i>value v</i>	$::=$	32-bit unsigned integer

Table I: A simple intermediate language (SIMPIL).

<i>taint t</i>	$::=$	T F
<i>value</i>	$::=$	$\langle v, t \rangle$
τ_Δ	$::=$	Maps variables to taint status
τ_μ	$::=$	Maps addresses to taint status

Table II: Additional changes to SIMPIL to enable dynamic taint analysis.

Dynamic Taint Policies

A taint policy specifies three properties:

- (Taint Introduction) how new taint is introduced to a program
- (Taint Propagation) how taint propagates as instructions execute
- (Taint Checking) how taint is checked during execution.

$$\begin{array}{c}
\frac{}{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash \text{var} \Downarrow \langle \Delta[\text{var}], \tau_\Delta[\text{var}] \rangle} \text{T-VAR} \qquad \frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle v, t \rangle}{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash \text{load } e \Downarrow \langle \mu[v], P_{\text{mem}}(t, \tau_\mu[v]) \rangle} \text{T-LOAD} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle v, t \rangle}{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash \Diamond_u e \Downarrow \langle \Diamond_u v, P_{\text{unop}}(t) \rangle} \text{T-UNOP} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_1 \Downarrow \langle v_1, t_1 \rangle \quad \tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_2 \Downarrow \langle v_2, t_2 \rangle \quad P_{\text{bincheck}}(t_1, t_2, v_1, v_2, \Diamond_b) = \mathbf{T}}{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_1 \Diamond_b e_2 \Downarrow \langle v_1 \Diamond_b v_2, P_{\text{binop}}(t_1, t_2) \rangle} \text{T-BINOP} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle v, t \rangle \quad \Delta' = \Delta[\text{var} \leftarrow v] \quad \tau'_\Delta = \tau_\Delta[\text{var} \leftarrow P_{\text{assign}}(t)] \quad \iota = \Sigma[\text{pc} + 1]}{\tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, \text{pc}, \text{var} := e \rightsquigarrow \tau_\mu, \tau'_\Delta, \Sigma, \mu, \Delta', \text{pc} + 1, \iota} \text{T-ASSIGN} \\
\\
\frac{\iota = \Sigma[\text{pc} + 1] \quad P_{\text{memcheck}}(t_1, t_2) = \mathbf{T} \quad \tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_1 \Downarrow \langle v_1, t_1 \rangle \quad \tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_2 \Downarrow \langle v_2, t_2 \rangle \quad \mu' = \mu[v_1 \leftarrow v_2] \quad \tau'_\mu = \tau_\mu[v_1 \leftarrow P_{\text{mem}}(t_1, t_2)]}{\tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, \text{pc}, \text{store}(e_1, e_2) \rightsquigarrow \tau'_\mu, \tau_\Delta, \Sigma, \mu', \Delta, \text{pc} + 1, \iota} \text{T-STORE} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle 1, t \rangle \quad \iota = \Sigma[\text{pc} + 1]}{\tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, \text{pc}, \text{assert}(e) \rightsquigarrow \tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, \text{pc} + 1, \iota} \text{T-ASSERT} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle 1, t_1 \rangle \quad \tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_1 \Downarrow \langle v_1, t_2 \rangle \quad P_{\text{condcheck}}(t_1, t_2) = \mathbf{T} \quad \iota = \Sigma[v_1]}{\tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, \text{pc}, \text{if } e \text{ then goto } e_1 \text{ else goto } e_2 \rightsquigarrow \tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, v_1, \iota} \text{T-TCOND} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle 0, t_1 \rangle \quad \tau_\mu, \tau_\Delta, \mu, \Delta \vdash e_2 \Downarrow \langle v_2, t_2 \rangle \quad P_{\text{condcheck}}(t_1, t_2) = \mathbf{T} \quad \iota = \Sigma[v_2]}{\tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, \text{pc}, \text{if } e \text{ then goto } e_1 \text{ else goto } e_2 \rightsquigarrow \tau_\mu, \tau_\Delta, \Sigma, \mu, \Delta, v_2, \iota} \text{T-FCOND} \\
\\
\frac{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle v_1, t \rangle \quad P_{\text{gotocheck}}(t) = \mathbf{T} \quad \iota = \Sigma[v_1]}{\tau_\mu, \tau_\Delta, \mu, \Delta \vdash e \Downarrow \langle v_1, t \rangle} \text{T-GOTO}
\end{array}$$

Dynamic Taint Analysis

- TaintDroid: An Information-Flow Tracking System for Realtime Privacy Monitoring on Smartphones

OSDI 2010

ACM Transactions on Computer Systems 2014



Allow **App Name** to
take pictures and
record video?

2 of 2

DENY

ALLOW



Allow **App Name** to
access this device's
location?

2 of 2

DENY

ALLOW



Allow **App Name** to
make and manage
phone calls?

2 of 2

DENY

ALLOW



Allow **App Name** to
access photos, media,
and files on your
device?

2 of 2

DENY

ALLOW



Allow **App Name** to
access this device's
location?

DENY

ALLOW

TaintDroid

- An extension to the Android mobile-phone platform that tracks the flow of privacy-sensitive data through third-party applications
- TaintDroid assumes that downloaded, third-party applications are not trusted, and monitors—in real-time—how these applications access and manipulate users' personal data
- detect when sensitive data leaves the system via untrusted applications

Use Dynamic Taint Analysis

- Sensitive information is first identified at a taint source, where a taint marking indicating the information type is assigned
- Dynamic taint analysis tracks how labeled data impacts other data in a way that might leak the original sensitive information
- This tracking is often performed at the instruction level
- the impacted data is identified before it leaves the system at a taint sink (usually the network interface)

Multilevel Taint Analysis

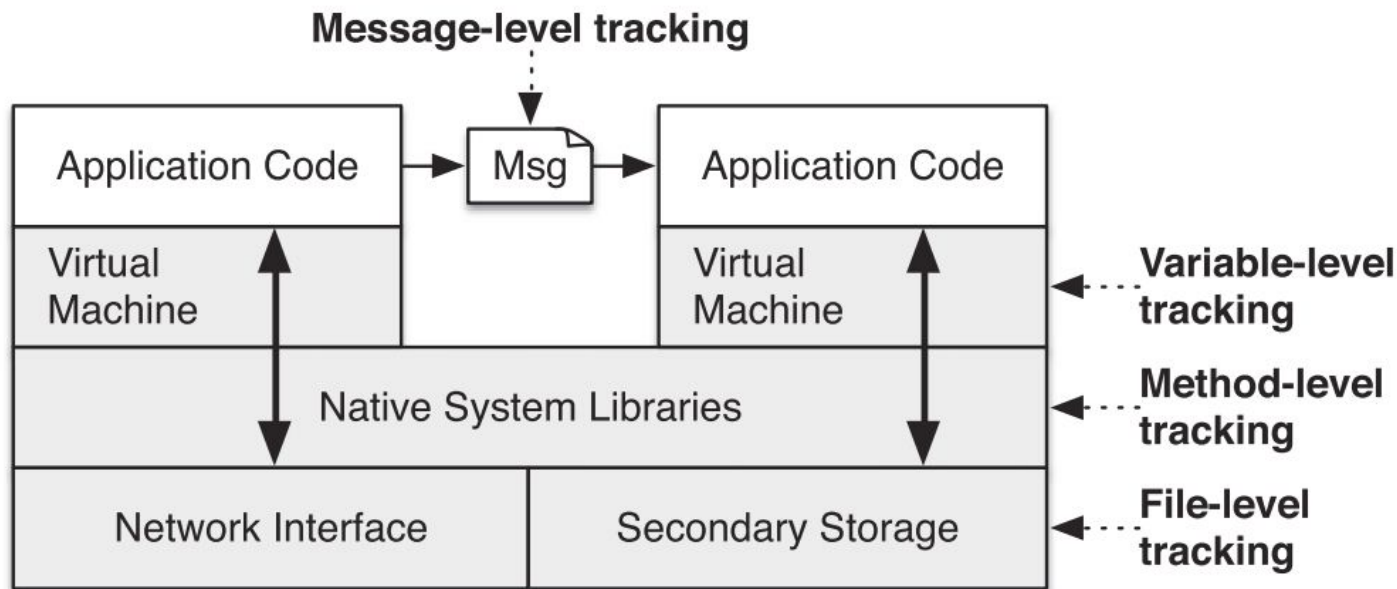


Fig. 1. Multilevel approach for performance-efficient taint tracking within a common smartphone architecture.

Android Background

Android is a Linux-based, open-source mobile-phone platform

Applications are written in Java and compiled to a custom bytecode format known as Dalvik EXecutable (DEX)

Each application executes within its own Dalvik VM interpreter instance. Each instance executes as a unique UNIX user identity to isolate applications within the Linux platform

Applications communicate via the Binder IPC subsystem

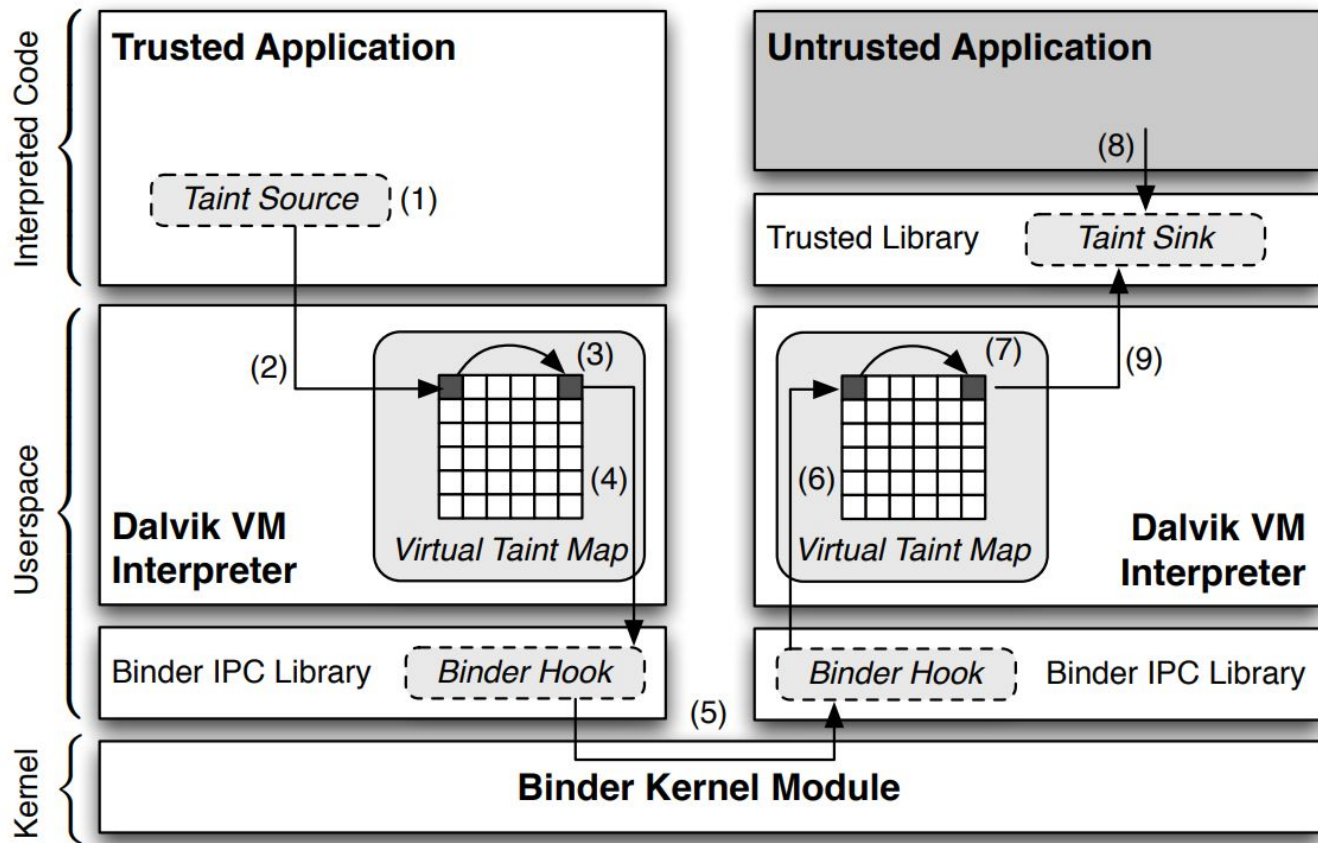


Figure 2: TaintDroid architecture within Android.

Challenges

taint tag storage

interpreted code taint propagation

native code taint propagation

IPC taint propagation, and

secondary storage taint propagation

Interpreted Code Taint Propagation

Table I. DEX Taint Propagation Logic

Op Format	Op Semantics	Taint Propagation	Description
<i>const-op</i> $v_A C$	$v_A \leftarrow C$	$\tau(v_A) \leftarrow \emptyset$	Clear v_A taint
<i>move-op</i> $v_A v_B$	$v_A \leftarrow v_B$	$\tau(v_A) \leftarrow \tau(v_B)$	Set v_A taint to v_B taint
<i>move-op-R</i> v_A	$v_A \leftarrow R$	$\tau(v_A) \leftarrow \tau(R)$	Set v_A taint to return taint
<i>return-op</i> v_A	$R \leftarrow v_A$	$\tau(R) \leftarrow \tau(v_A)$	Set return taint (\emptyset if void)
<i>move-op-E</i> v_A	$v_A \leftarrow E$	$\tau(v_A) \leftarrow \tau(E)$	Set v_A taint to exception taint
<i>throw-op</i> v_A	$E \leftarrow v_A$	$\tau(E) \leftarrow \tau(v_A)$	Set exception taint
<i>unary-op</i> $v_A v_B$	$v_A \leftarrow \otimes v_B$	$\tau(v_A) \leftarrow \tau(v_B)$	Set v_A taint to v_B taint
<i>binary-op</i> $v_A v_B v_C$	$v_A \leftarrow v_B \otimes v_C$	$\tau(v_A) \leftarrow \tau(v_B) \cup \tau(v_C)$	Set v_A taint to v_B taint \cup v_C taint
<i>binary-op</i> $v_A v_B$	$v_A \leftarrow v_A \otimes v_B$	$\tau(v_A) \leftarrow \tau(v_A) \cup \tau(v_B)$	Update v_A taint with v_B taint
<i>binary-op</i> $v_A v_B C$	$v_A \leftarrow v_B \otimes C$	$\tau(v_A) \leftarrow \tau(v_B)$	Set v_A taint to v_B taint
<i>aput-op</i> $v_A v_B v_C$	$v_B[v_C] \leftarrow v_A$	$\tau(v_B[\cdot]) \leftarrow \tau(v_B[\cdot]) \cup \tau(v_A)$	Update array v_B taint with v_A taint
<i>aget-op</i> $v_A v_B v_C$	$v_A \leftarrow v_B[v_C]$	$\tau(v_A) \leftarrow \tau(v_B[\cdot]) \cup \tau(v_C)$	Set v_A taint to array and index taint
<i>sput-op</i> $v_A f_B$	$f_B \leftarrow v_A$	$\tau(f_B) \leftarrow \tau(v_A)$	Set field f_B taint to v_A taint
<i>sget-op</i> $v_A f_B$	$v_A \leftarrow f_B$	$\tau(v_A) \leftarrow \tau(f_B)$	Set v_A taint to field f_B taint
<i>iput-op</i> $v_A v_B f_C$	$v_B(f_C) \leftarrow v_A$	$\tau(v_B(f_C)) \leftarrow \tau(v_A)$	Set field f_C taint to v_A taint
<i>iget-op</i> $v_A v_B f_C$	$v_A \leftarrow v_B(f_C)$	$\tau(v_A) \leftarrow \tau(v_B(f_C)) \cup \tau(v_B)$	Set v_A taint to f_C and obj. ref. taint

Register variables and class fields are referenced by v_X and f_X , respectively. R and E are the return and exception variables maintained within the interpreter. A , B , and C are bytecode constants.

Table 2: Applications grouped by the requested permissions (L: location, C: camera, A: audio, P: phone state). Android Market categories are indicated in parenthesis, showing the diversity of the studied applications.

Applications*	#	Permissions [†]			
		L	C	A	P
The Weather Channel (News & Weather); Cestos, Solitaire (Game); Movies (Entertainment); Babble (Social); Manga Browser (Comics)	6	x			
Bump, Wertago (Social); Antivirus (Communication); ABC — Animals, Traffic Jam, Hearts, Blackjack, (Games); Horoscope (Lifestyle); Yellow Pages (Reference); 3001 Wisdom Quotes Lite, Dastelefonbuch, Astrid (Productivity), BBC News Live Stream (News & Weather); Ring-tones (Entertainment)	14	x			x
Layar (Lifestyle); Knocking (Social); Coupons (Shopping); Trapster (Travel); Spongebob Slide (Game); ProBasketBall (Sports)	6	x	x		x
MySpace (Social); Barcode Scanner, ixMAT (Shopping)	3		x		
Evernote (Productivity)	1	x	x	x	

* Listed names correspond to the name displayed on the phone and not necessarily the name listed in the Android Market.

† All listed applications also require access to the Internet.

Findings

Table 3: Potential privacy violations by 20 of the studied applications. Note that three applications had multiple violations, one of which had a violation in all three categories.

Observed Behavior (# of apps)	Details
Phone Information to Content Servers (2)	2 apps sent out the phone number, IMSI, and ICC-ID along with the geo-coordinates to the app's content server.
Device ID to Content Servers (7)*	2 Social, 1 Shopping, 1 Reference and three other apps transmitted the IMEI number to the app's content server.
Location to Advertisement Servers (15)	5 apps sent geo-coordinates to ad.qwapi.com, 5 apps to admob.com, 2 apps to ads.mobclix.com (1 sent location both to admob.com and ads.mobclix.com) and 4 apps sent location [†] to data.flurry.com.

* TaintDroid flagged nine applications in this category, but only seven transmitted the raw IMEI without mentioning such practice in the EULA.

[†]To the best of our knowledge, the binary messages contained tainted location data (see the discussion below).

Symbolic Execution

- Builds predicates that characterize
 - Conditions for executing paths
 - Effects of the execution on program state
- Bridges program behavior to logic
- Finds important applications in
 - program analysis
 - test data generation
 - formal verification (proofs) of program correctness

Symbolic state

Values are concrete but **symbol** and **expressions over symbols**
Executing statements computes new **expressions**

Example 6. *Consider the following program:*

```
1  x := 2*get_input(.)  
2  if x-5 == 14 then goto 3 else goto 4  
3  // catastrophic failure  
4  // normal behavior
```

Only one input will trigger the failure.

$\frac{v \text{ is a fresh symbol}}{\mu, \Delta \vdash \text{get_input}(\cdot) \Downarrow v} \text{ S-INPUT}$
$\frac{\mu, \Delta \vdash e \Downarrow e' \quad \Pi' = \Pi \wedge e' \quad \iota = \Sigma[pc + 1]}{\Pi, \Sigma, \mu, \Delta, pc, \text{assert}(e) \rightsquigarrow \Pi', \Sigma, \mu, \Delta, pc + 1, \iota} \text{ S-ASSERT}$
$\frac{\mu, \Delta \vdash e \Downarrow e' \quad \Delta \vdash e_1 \Downarrow v_1 \quad \Pi' = \Pi \wedge (e' = 1) \quad \iota = \Sigma[v_1]}{\Pi, \Sigma, \mu, \Delta, pc, \text{if } e \text{ then goto } e_1 \text{ else goto } e_2 \rightsquigarrow \Pi', \Sigma, \mu, \Delta, v_1, \iota} \text{ S-TCOND}$
$\frac{\mu, \Delta, \vdash e \Downarrow e' \quad \Delta \vdash e_2 \Downarrow v_2 \quad \Pi' = \Pi \wedge (e' = 0) \quad \iota = \Sigma[v_2]}{\Pi, \Sigma, \mu, \Delta, pc, \text{if } e \text{ then goto } e_1 \text{ else goto } e_2 \rightsquigarrow \Pi', \Sigma, \mu, \Delta, v_2, \iota} \text{ S-FCOND}$

Figure 6: Operational semantics of the language for forward symbolic execution.

Statement	Δ	Π	Rule	pc
start	$\{\}$	$true$		1
$x := 2 * \text{get_input}(\cdot)$	$\{x \rightarrow 2 * s\}$	$true$	S-ASSIGN	2
if $x-5 == 14$ goto 3 else goto 4	$\{x \rightarrow 2 * s\}$	$[(2 * s) - 5 == 14]$	S-TCOND	3
if $x-5 == 14$ goto 3 else goto 4	$\{x \rightarrow 2 * s\}$	$\neg[(2 * s) - 5 == 14]$	S-FCOND	4

Table VII: Simulation of forward symbolic execution.

Challenges

- **Symbolic Memory.** What should we do when the analysis uses the μ context — whose index must be a non-negative integer — with a symbolic index?
- **System Calls.** How should our analysis deal with external interfaces such as system calls?
- **Path Selection.** Each conditional represents a branch in the program execution space. How should we decide which branches to take?

Challenges

- **Symbolic Memory.** What should we do when the analysis uses the μ context — whose index must be a non-negative integer — with a symbolic index?
- **System Calls.** How should our analysis deal with external interfaces such as system calls?
- **Path Selection.** Each conditional represents a branch in the program execution space. How should we decide which branches to take?

Challenges

- **Symbolic Memory.** What should we do when the analysis uses the μ context — whose index must be a non-negative integer — with a symbolic index?
- **System Calls.** How should our analysis deal with external interfaces such as system calls?
- **Path Selection.** Each conditional represents a branch in the program execution space. How should we decide which branches to take?