

# Interpreter-Guided JIT Compiler Test Generation

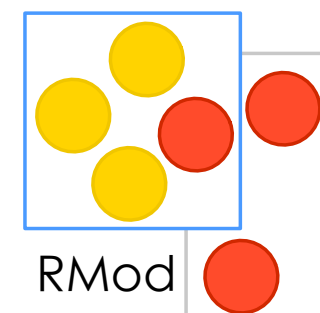
Validating the Pharo JIT compiler through  
**concolic** execution and **differential** testing

Guille Polito - Pablo Tesone - Stéphane Ducasse  
[guillermo.polito@univ-lille.fr](mailto:guillermo.polito@univ-lille.fr)  
@guillep

PLDI'22 — San Diego

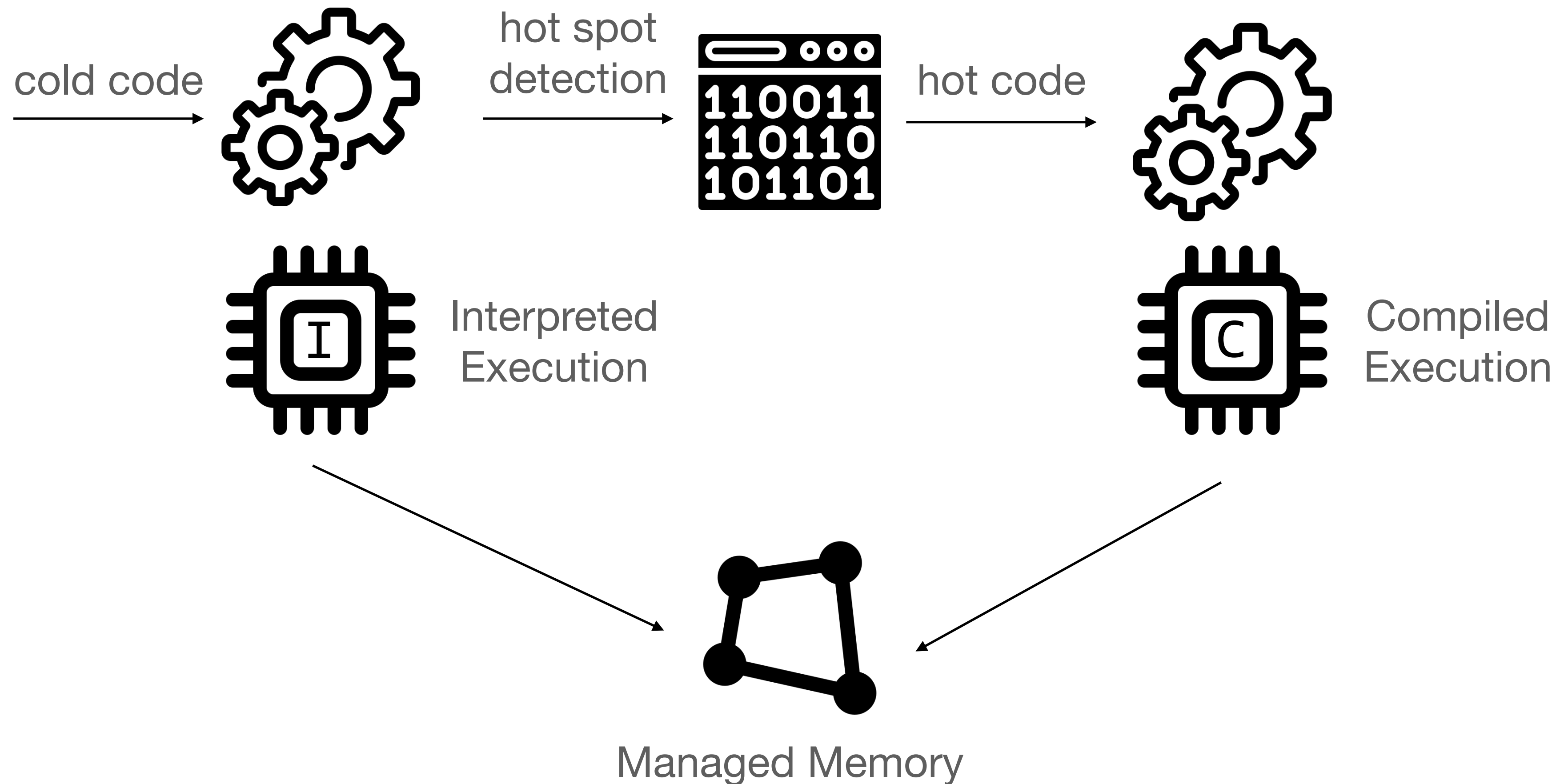


*Inria*

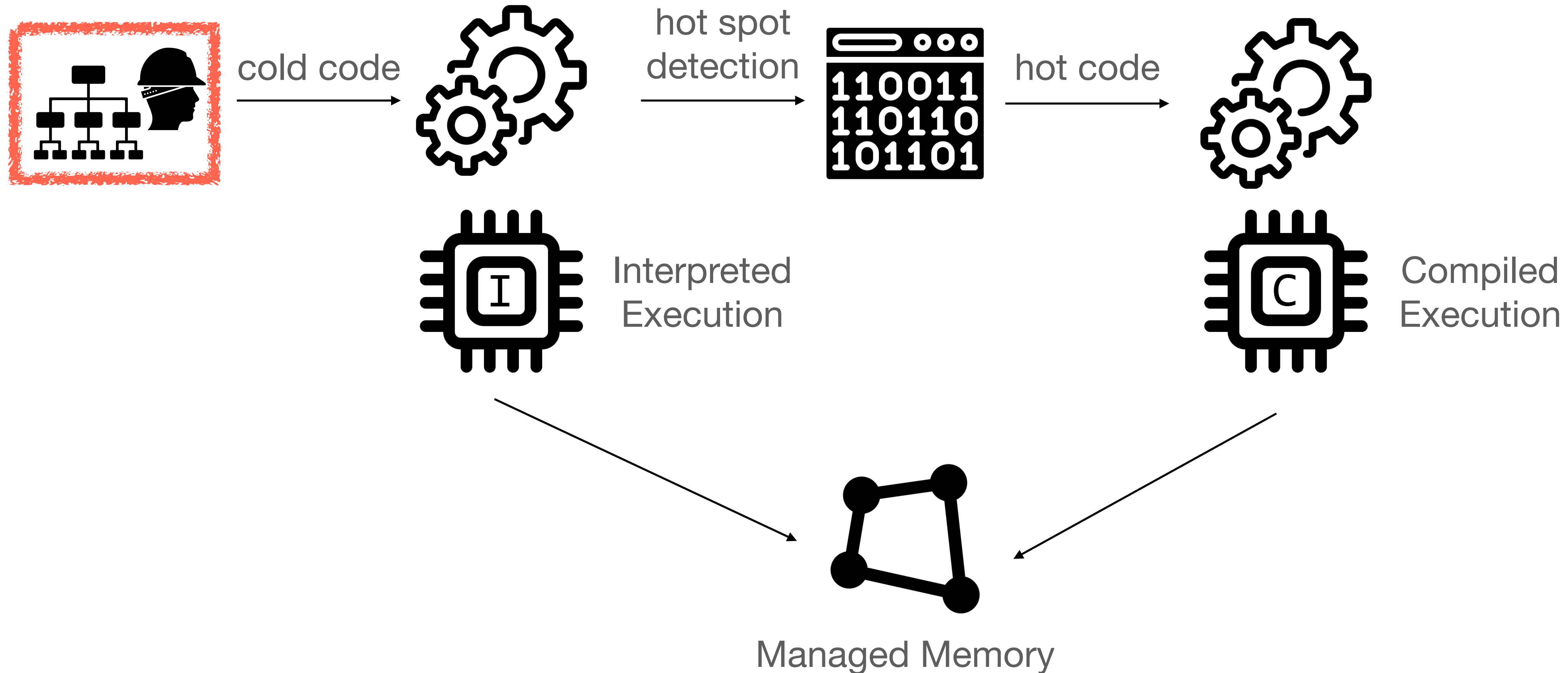


**Université  
de Lille**

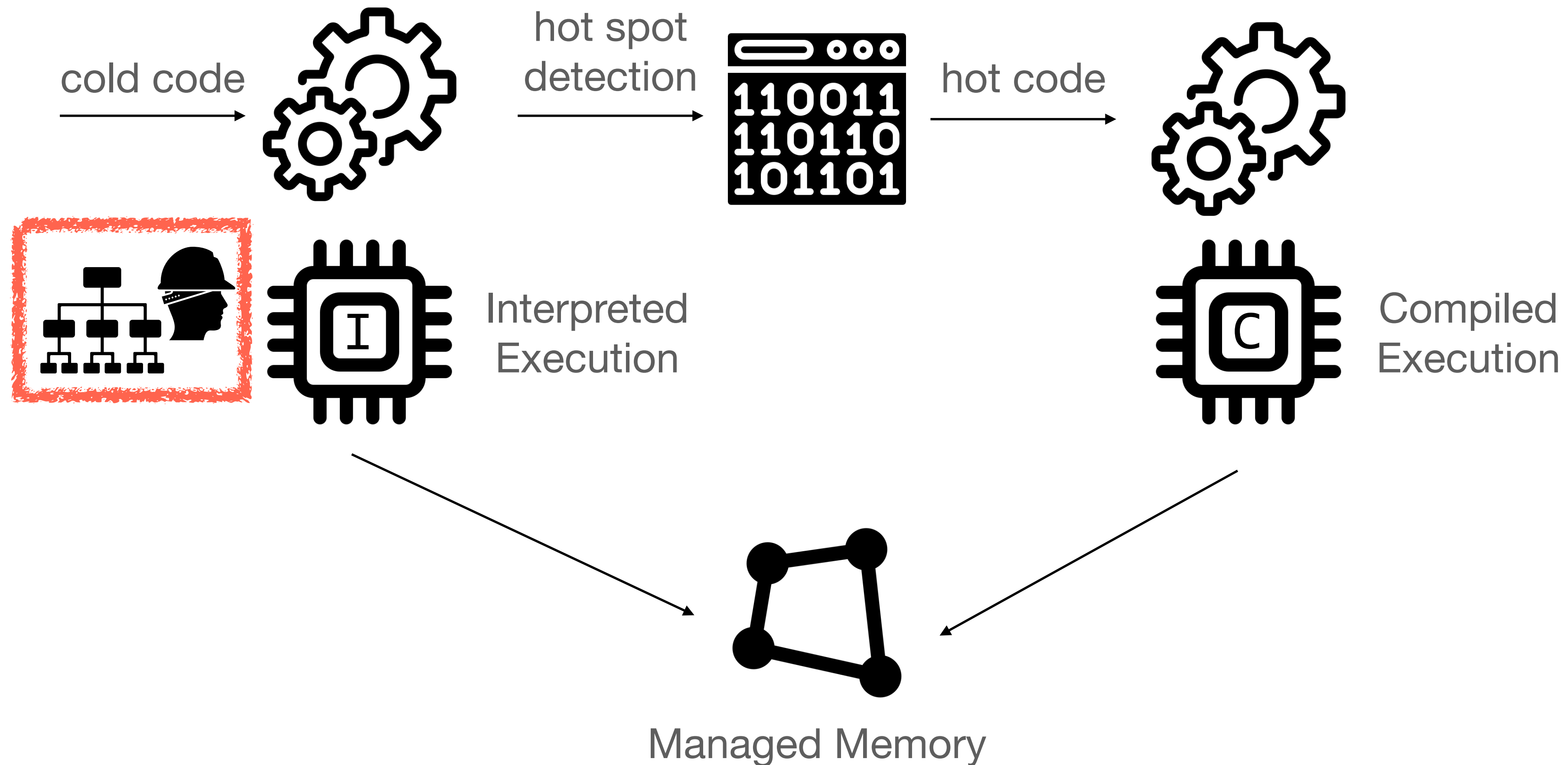
# Virtual Machine Execution Engine



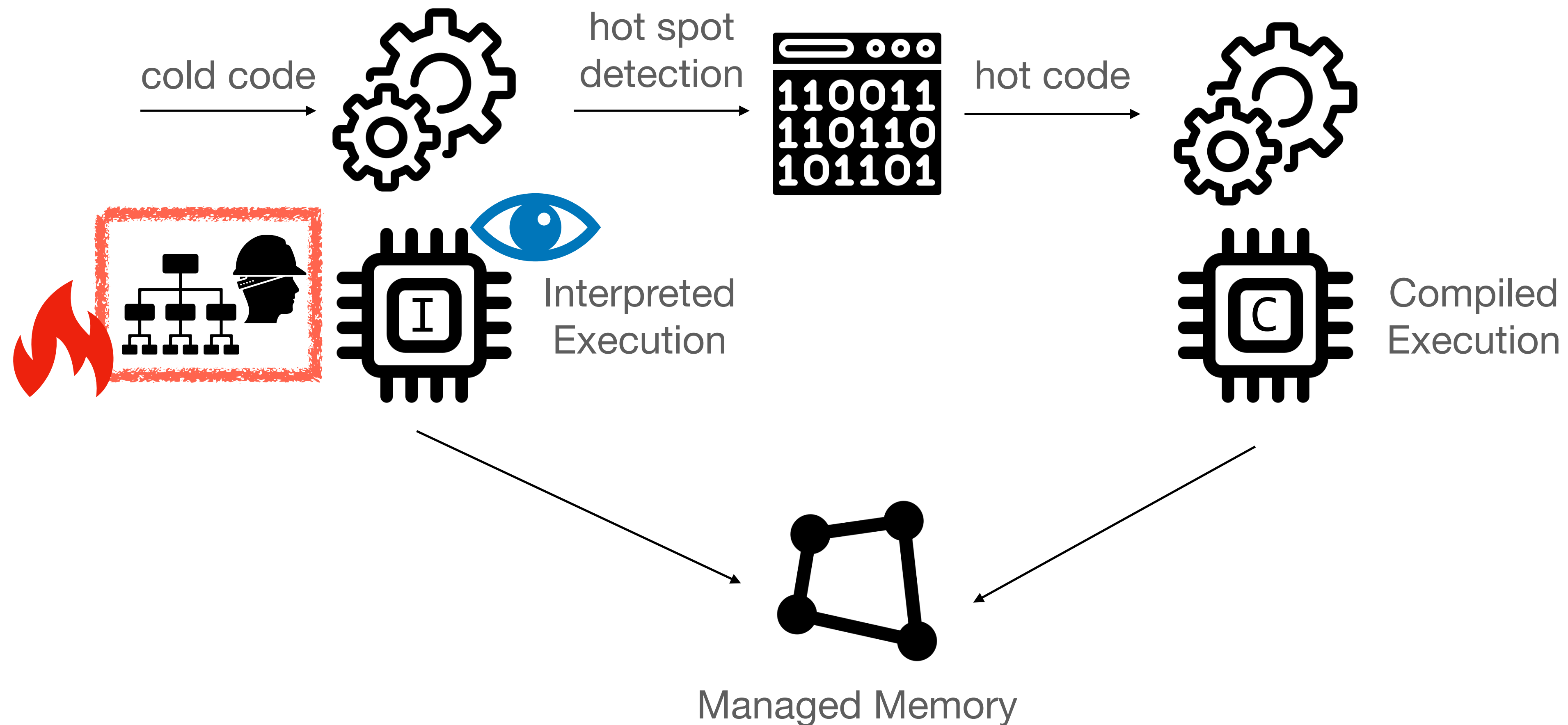
# Virtual Machine Execution Engine



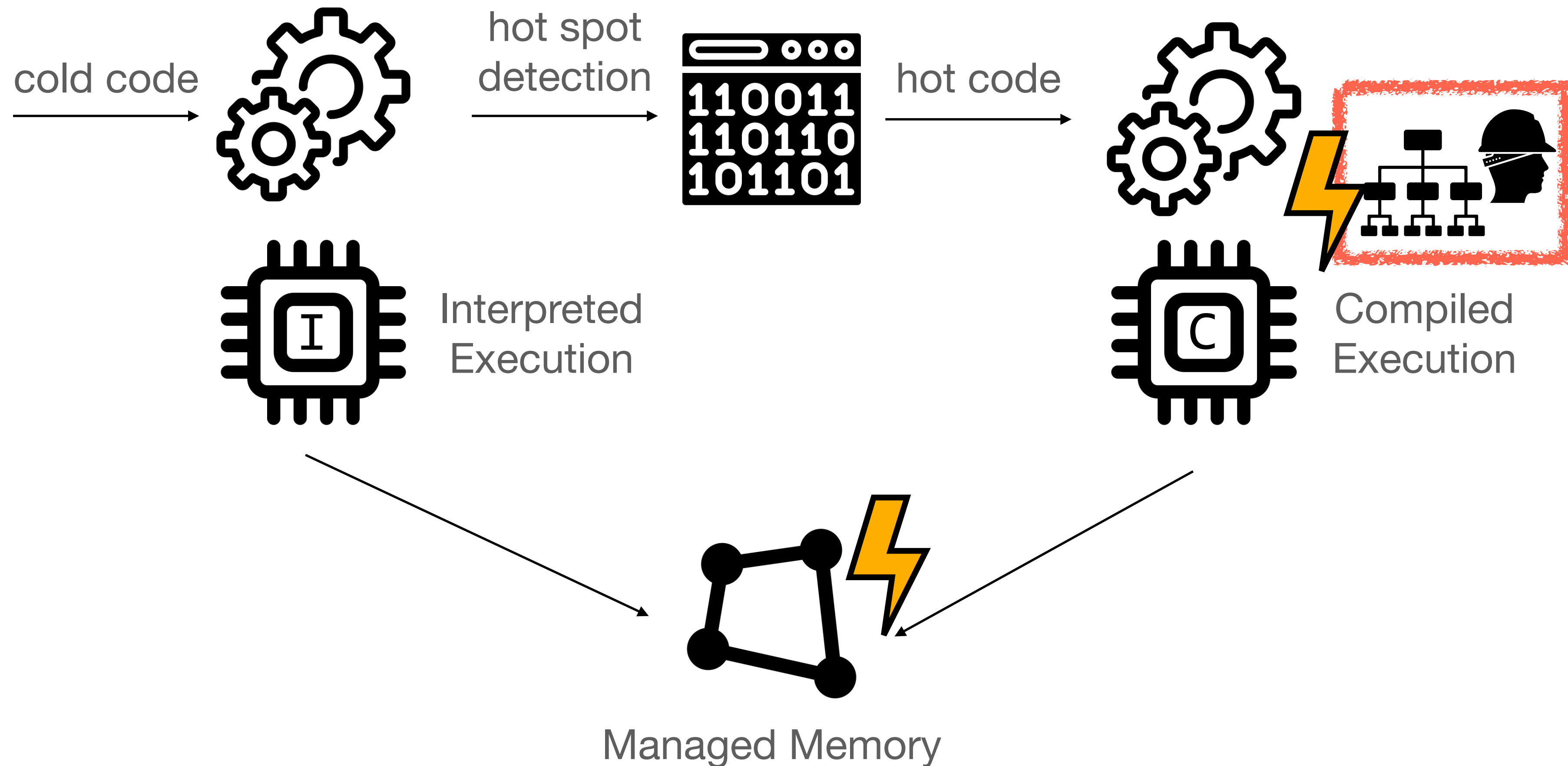
# Virtual Machine Execution Engine



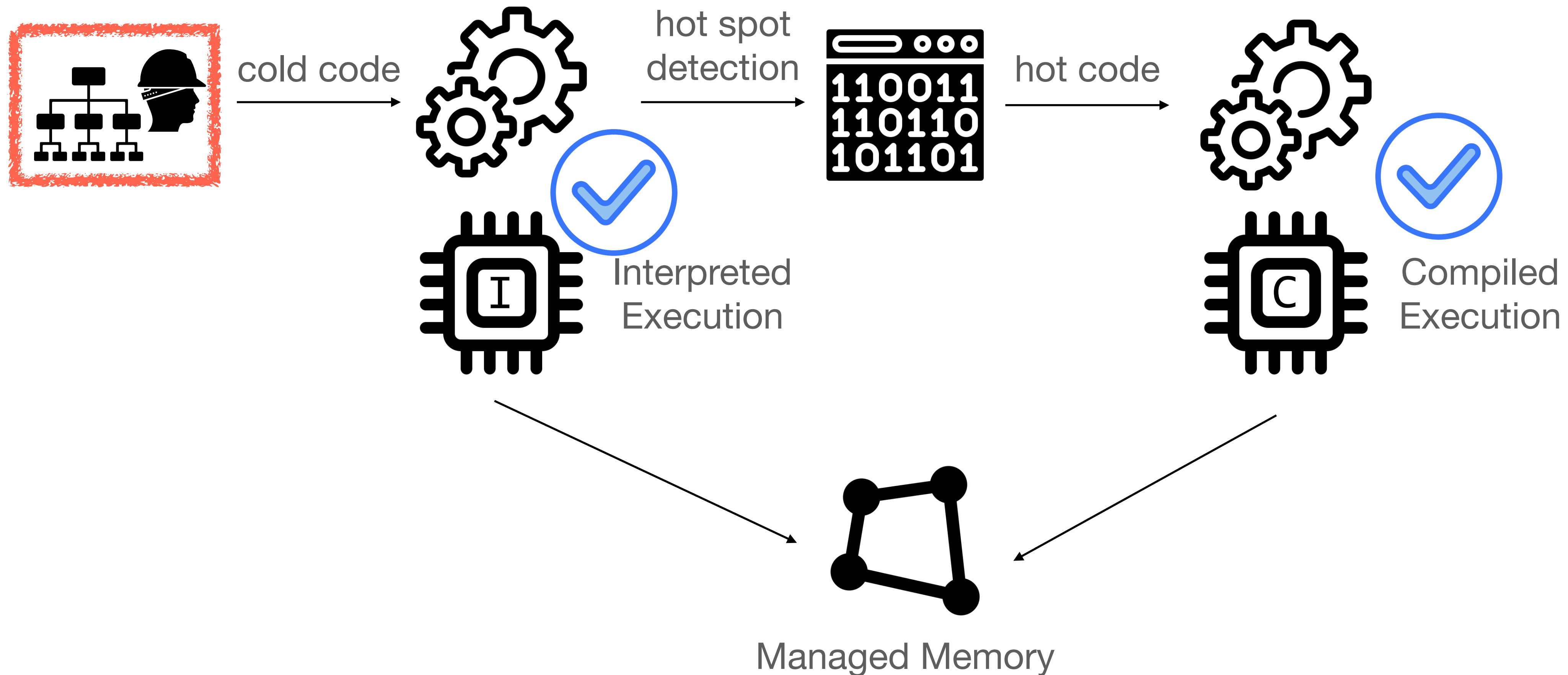
# Virtual Machine Execution Engine



# Virtual Machine Execution Engine

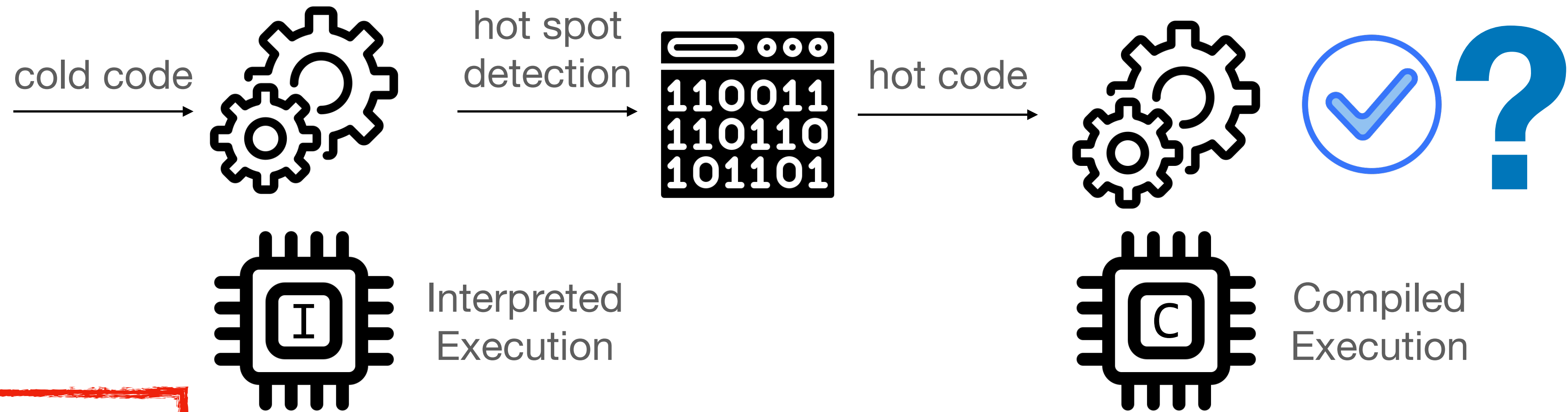


# How can we automatically test VMs?





# Challenges of VM Test Generation



## Challenge 1: Test Diversity

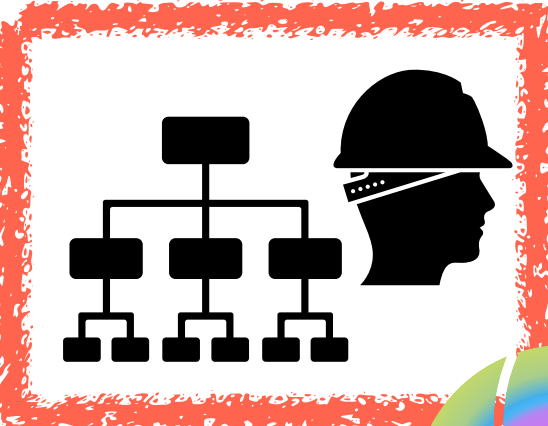
Do they cover different code *regions/branches/paths*?

## Challenge 2: Test Oracles

How do we determine what is the *expected output* of a generated test?



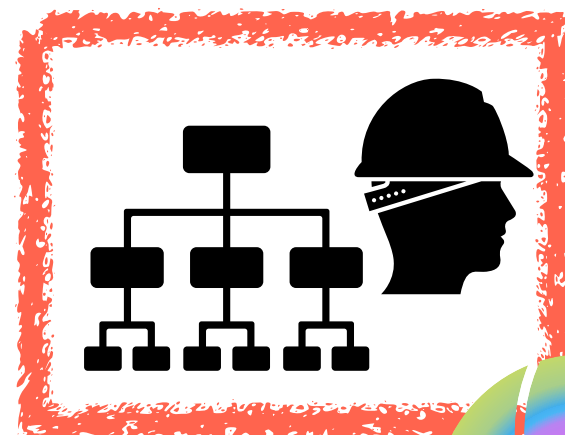
# Black Box Testing + Fuzzing



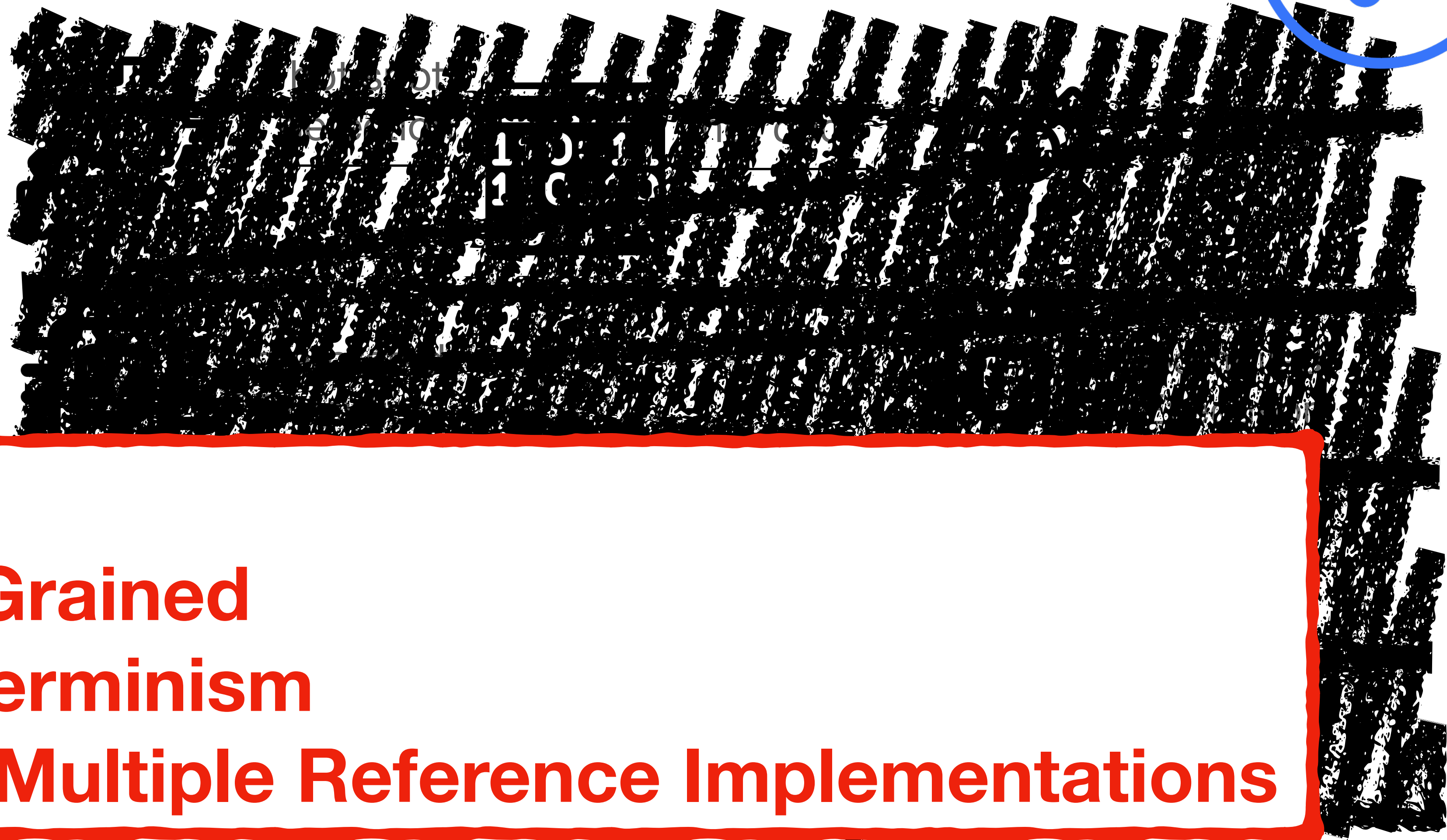
cold code



# Black Box Testing + Fuzzing

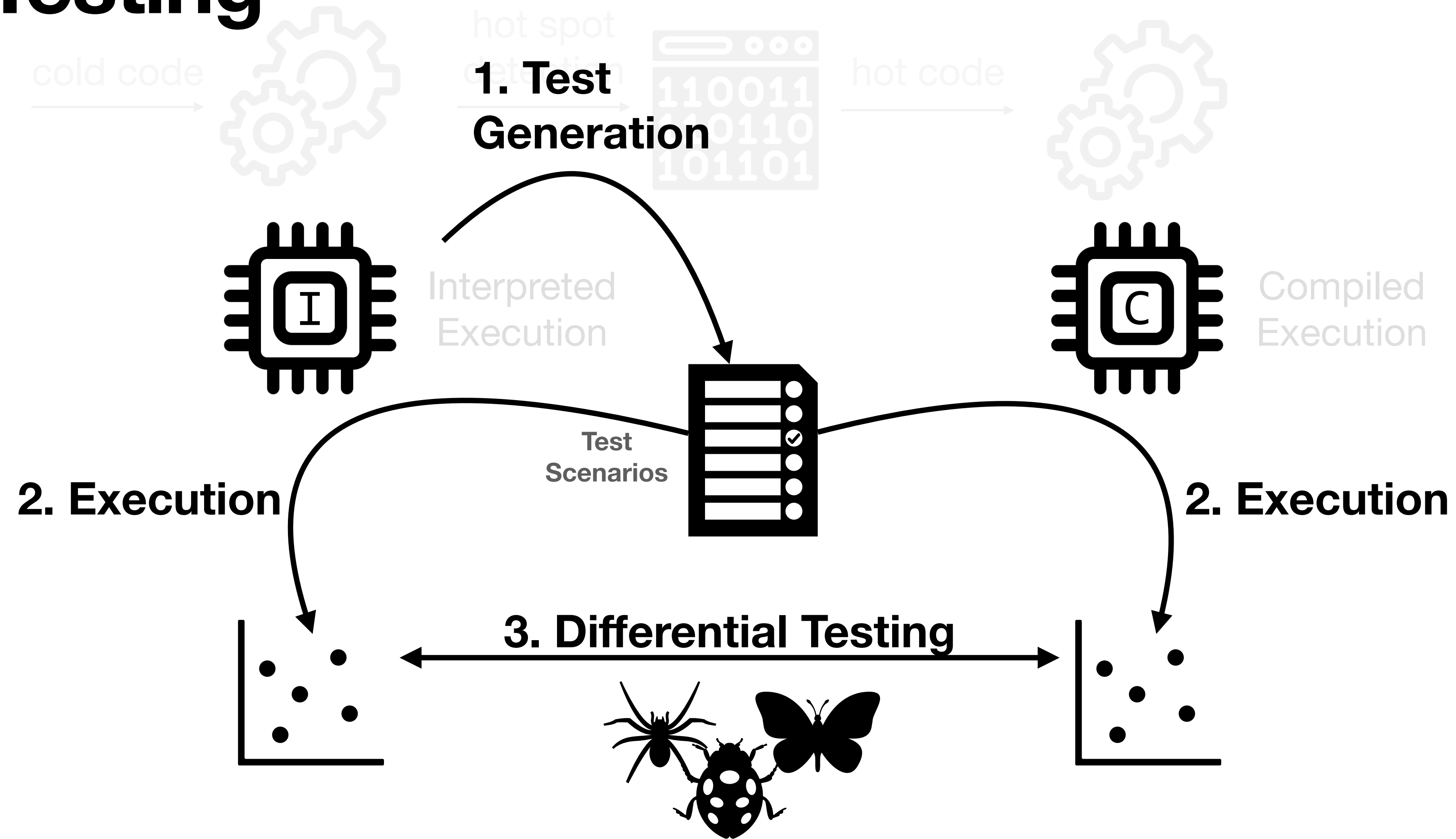


cold code



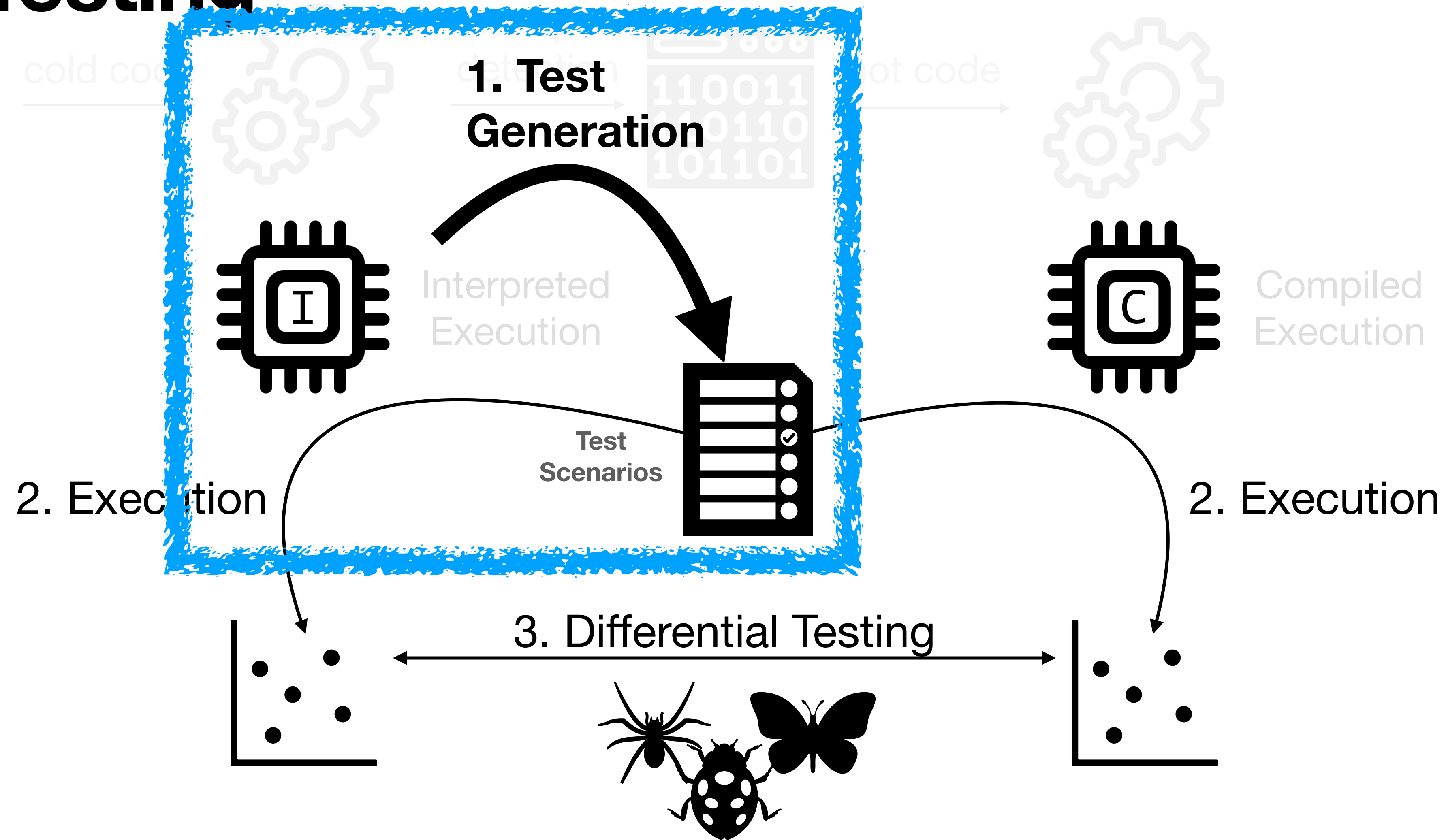
- X Slow**
- X Coarse Grained**
- X Non Determinism**
- X Require Multiple Reference Implementations**

# Interpreter-Guided Automatic JIT Compiler Unit Testing

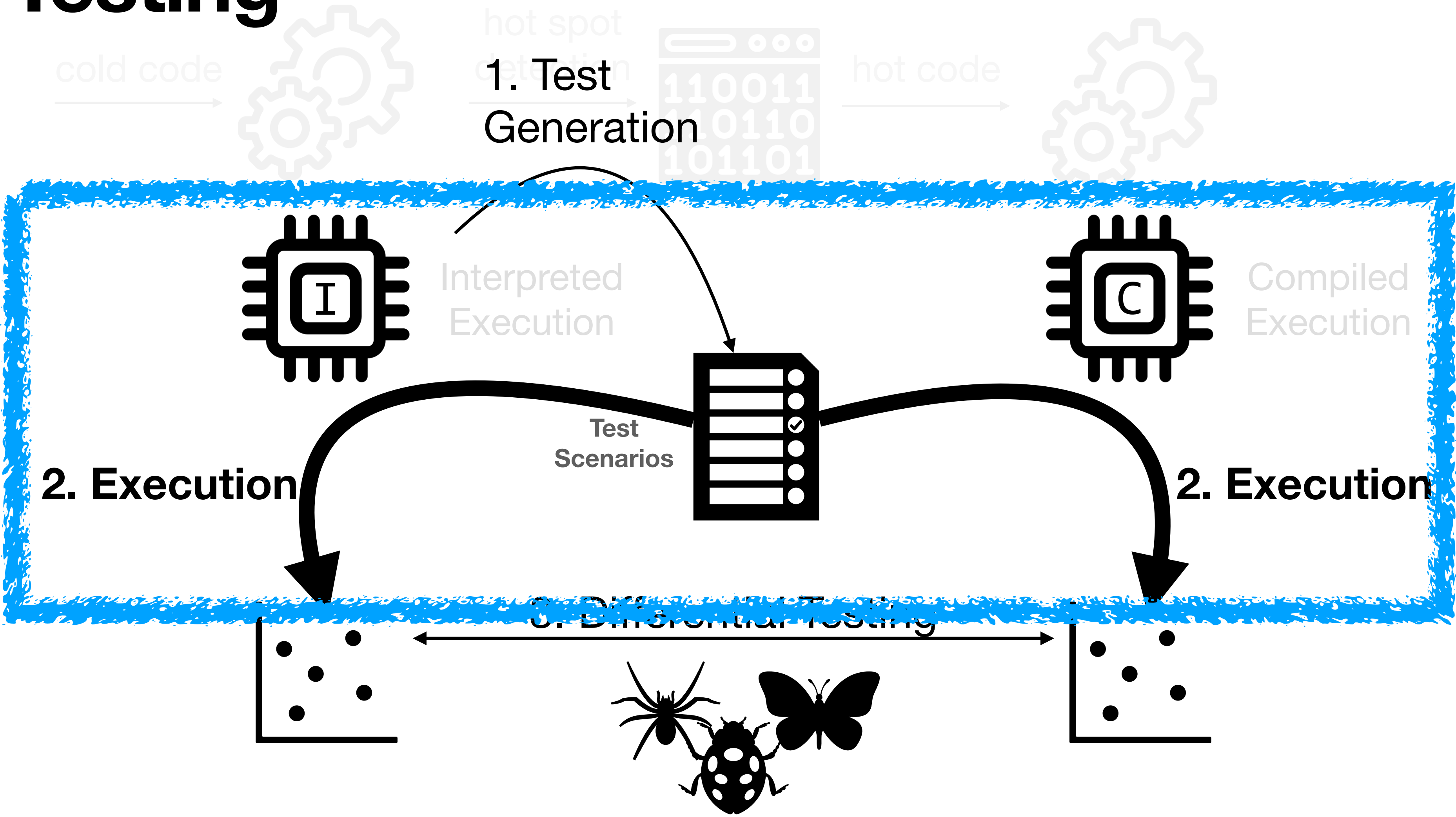




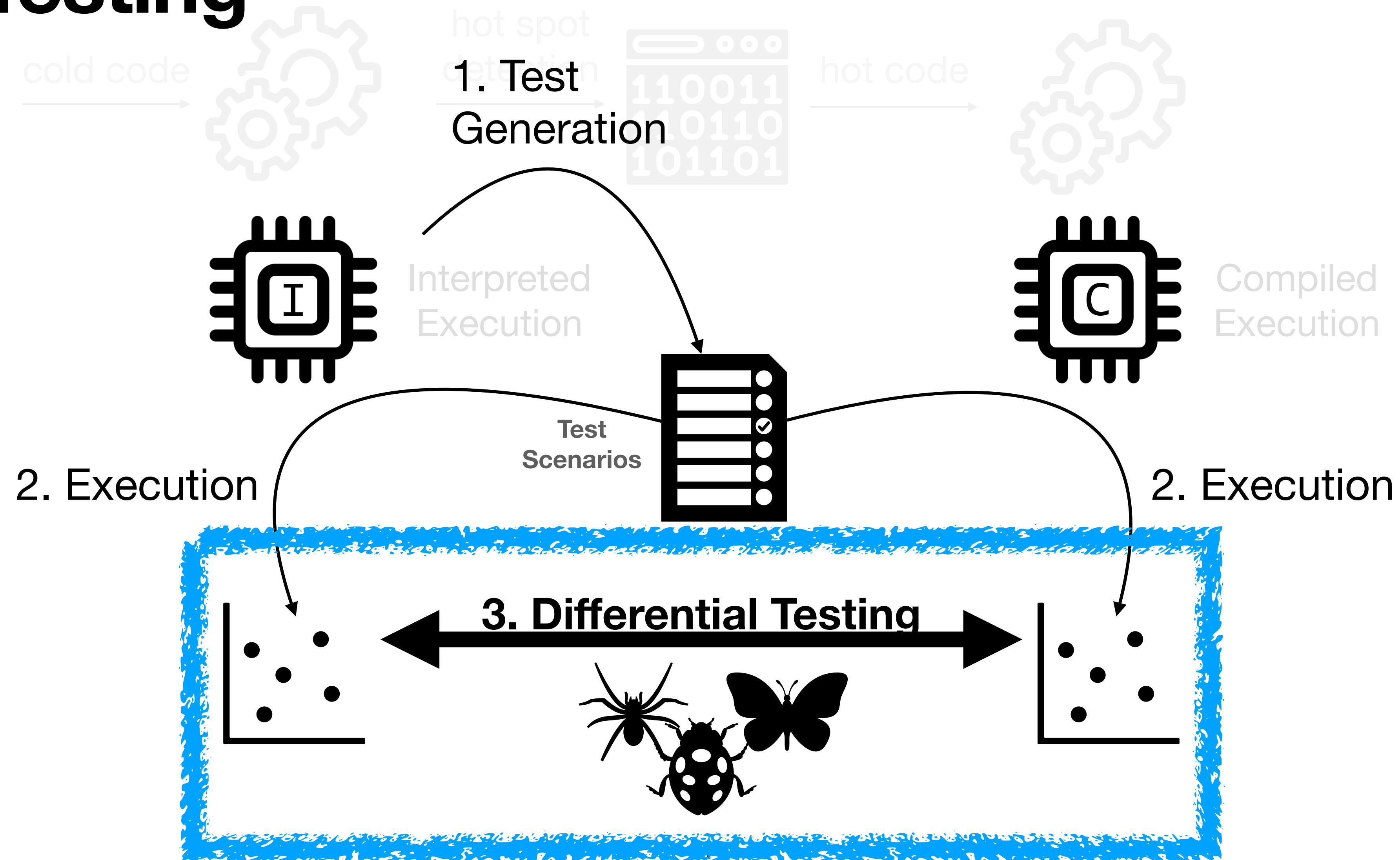
# Interpreter-Guided Automatic JIT Compiler Unit Testing



# Interpreter-Guided Automatic JIT Compiler Unit Testing



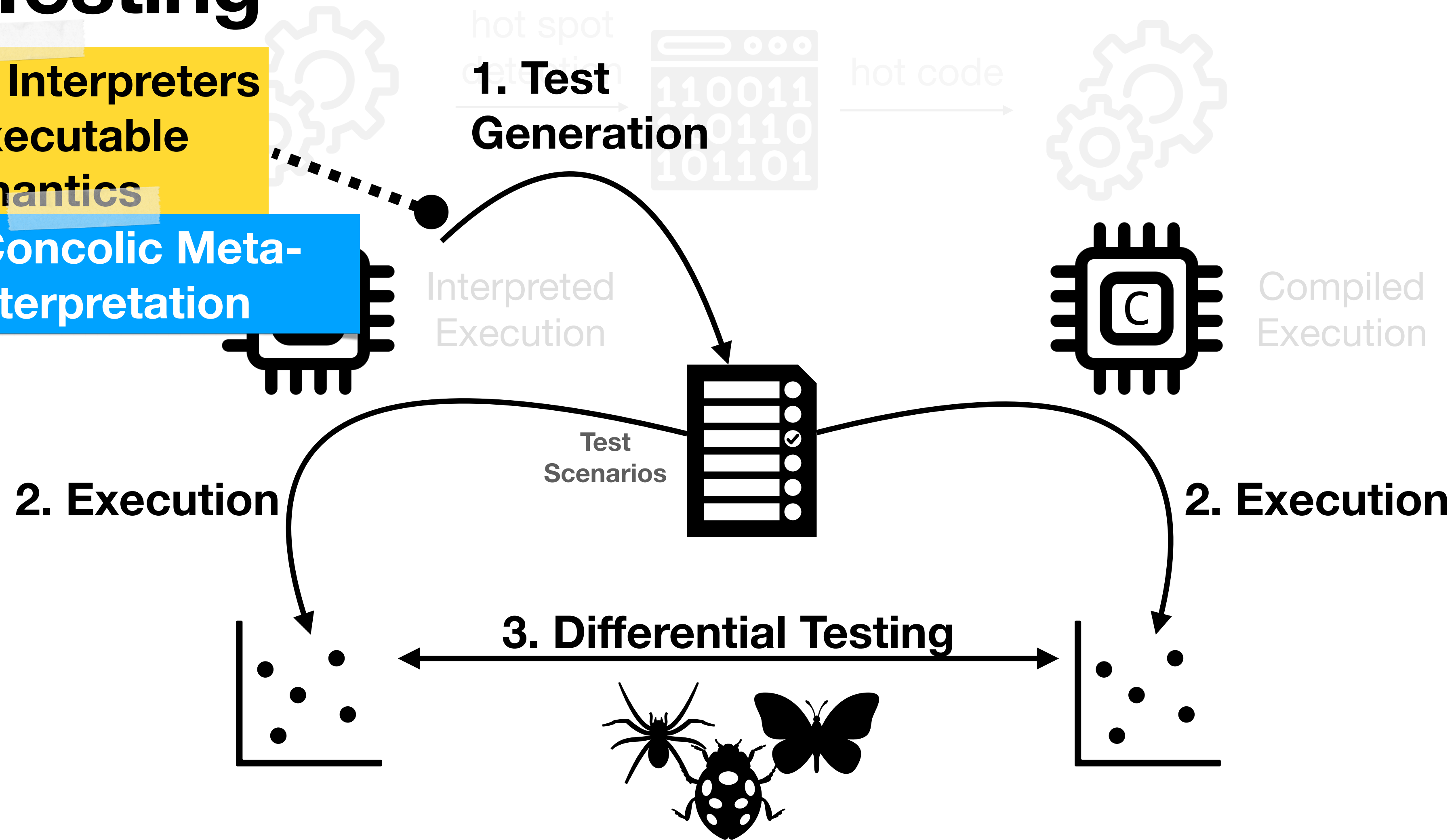
# Interpreter-Guided Automatic JIT Compiler Unit Testing



# Interpreter-Guided Automatic JIT Compiler Unit Testing

**Insight 1: Interpreters are Executable Semantics**

**=> Concolic Meta-Interpretation**





# Interpreter are Executable Semantics

## Pharo VM Example

```
1 Interpreter >> bytecodePrimAdd
```

```
2 | rcvr arg result |
```

```
3 rcvr := self internalStackValue: 1.
```

```
4 arg := self internalStackValue: 0.
```

```
5 (objectMemory areIntegers: rcvr and: arg) ifTrue: [
```

```
6 result := (objectMemory integerValueOf: rcvr) + (
```

```
7 objectMemory integerValueOf: arg).
```

```
8 "Check for overflow"
```

```
9 (objectMemory isIntegerValue: result) ifTrue: [
```

```
10 self normalSend
```

```
11 internalPop: 2
```

```
12 thenPush: (objectMemory integerObjectOf: result).
```

```
13 self fetchNextBytecode "success"]].
```

```
14 "Slow path, message send"
```

```
self normalSend
```

If both operands are integers

If their sum does not overflow

Else, slow path => message send



# Interpreter-Guided Automatic JIT Compiler Unit Testing

**Insight 1: Interpreters are Executable Semantics**

**=> Concolic Meta-Interpretation**

**1. Test Generation**

Interpreted Execution

hot code

Compiled Execution

**2. Execution**

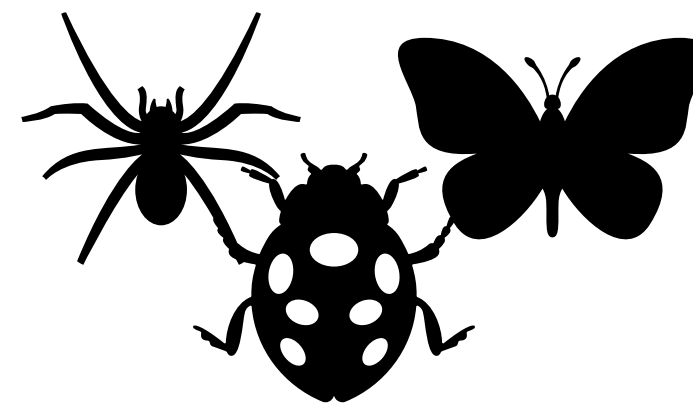
Test Scenarios

**2. Execution**

**Insight 2: Interpreters and Compiler Share Semantics**

**=> Differential Testing**

**3. Differential Testing**



# Interpreter VS Compiled Code

## Pharo VM Example

1 **Interpreter** >> bytecodePrimAdd

2 | rcvr arg result |

3 rcvr := **self** internalStackValue: 1.

4 arg := **self** internalStackValue: 0.

5 (objectMemory areIntegers: rcvr and: arg) **True**: [

6 result := (objectMemory integerValueOf: rcvr) + (

7 objectMemory integerValueOf: arg).

7 "Check for overflow"

8 (objectMemory isIntegerValue: result) **True**: [

9 **self** normalSend

10 internalPop: 2

11 thenPush: (objectMemory integerObjectOf: result).

12 **self** fetchNextBytecode: "success"]].

13 "Slow path, message send"

14 **self** normalSend

1 ... # previous bytecode **IR**

2 checkSmallInteger t0

3 jumpzero notsmi

4 checkSmallInteger t1

5 jumpzero notsmi

6 t2 := t0 + t1

7 jumplfNotOverflow continue

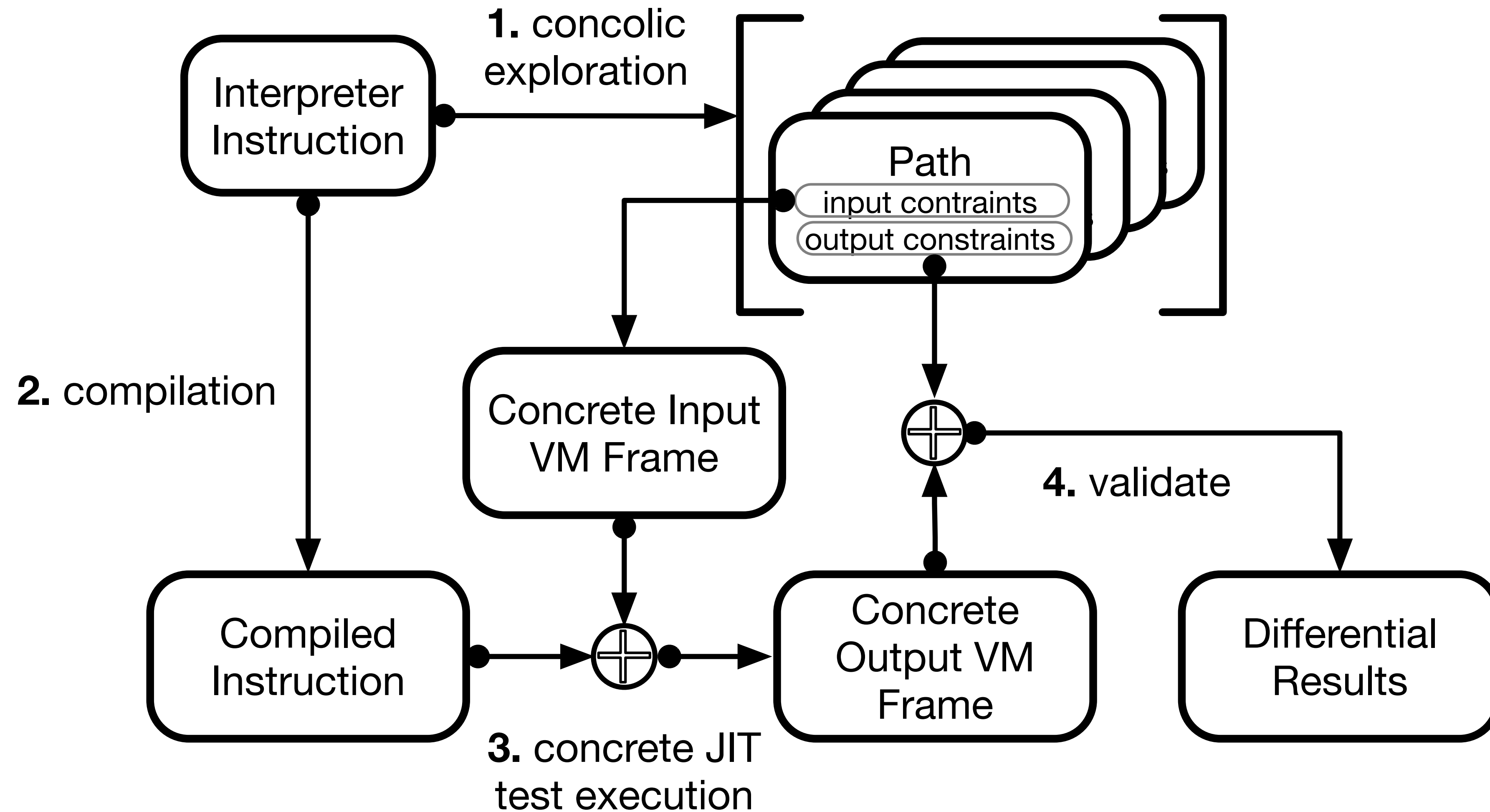
8 notsmi: #slow case first send

9 t2 := send #+ t0 t1

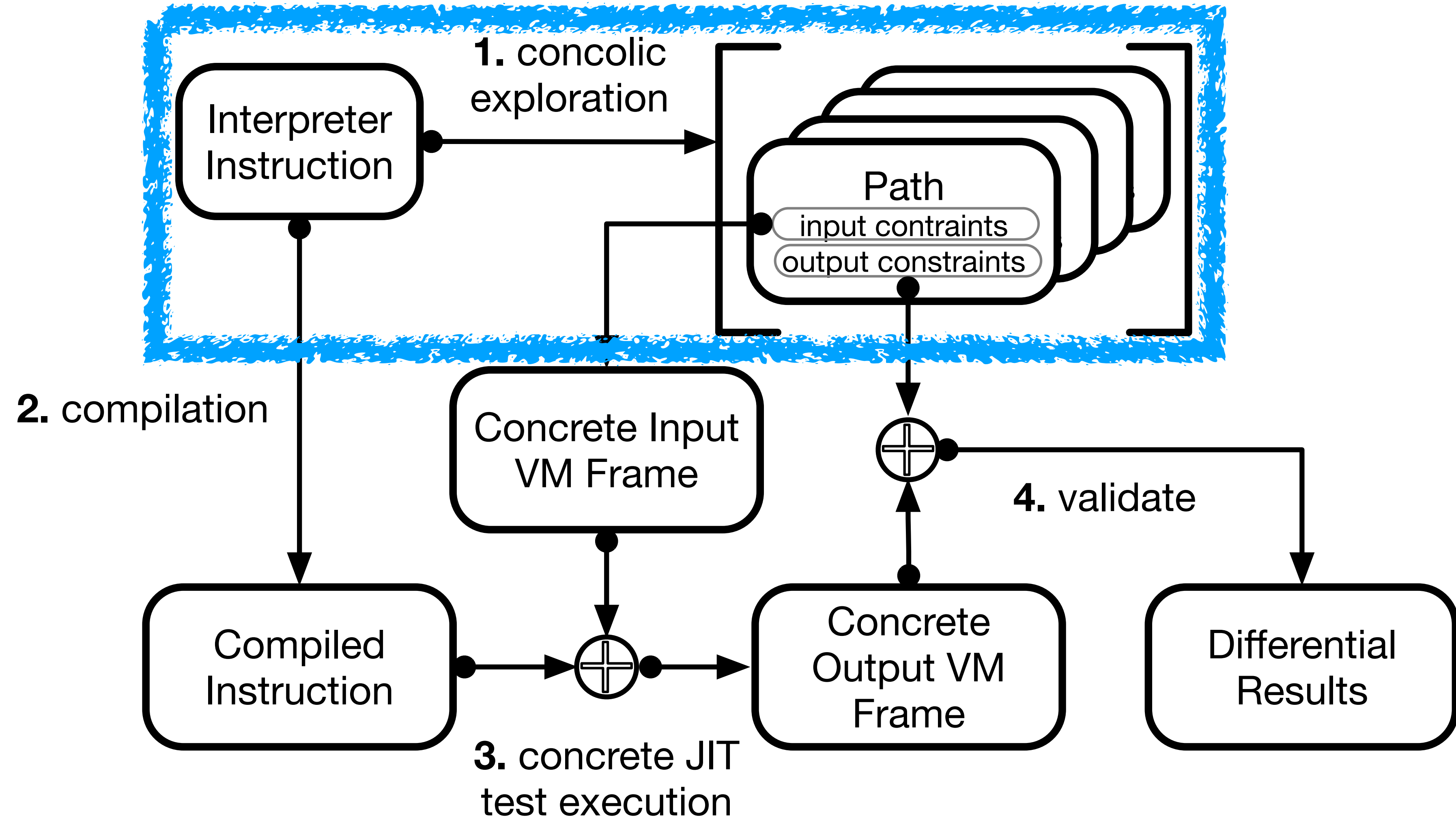
10 continue:

11 ... # following bytecode **IR**

# Implementation View

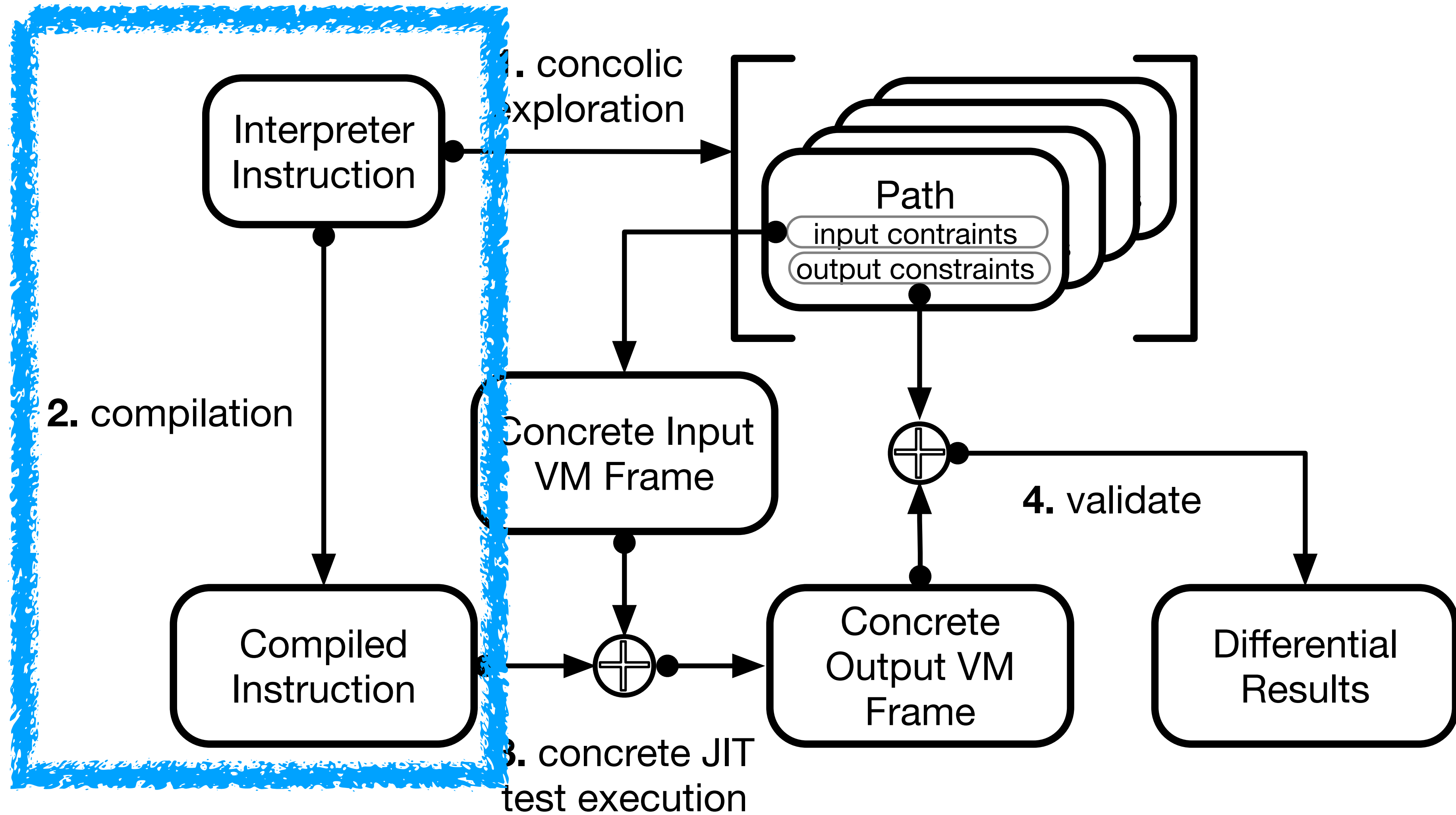


# Implementation View

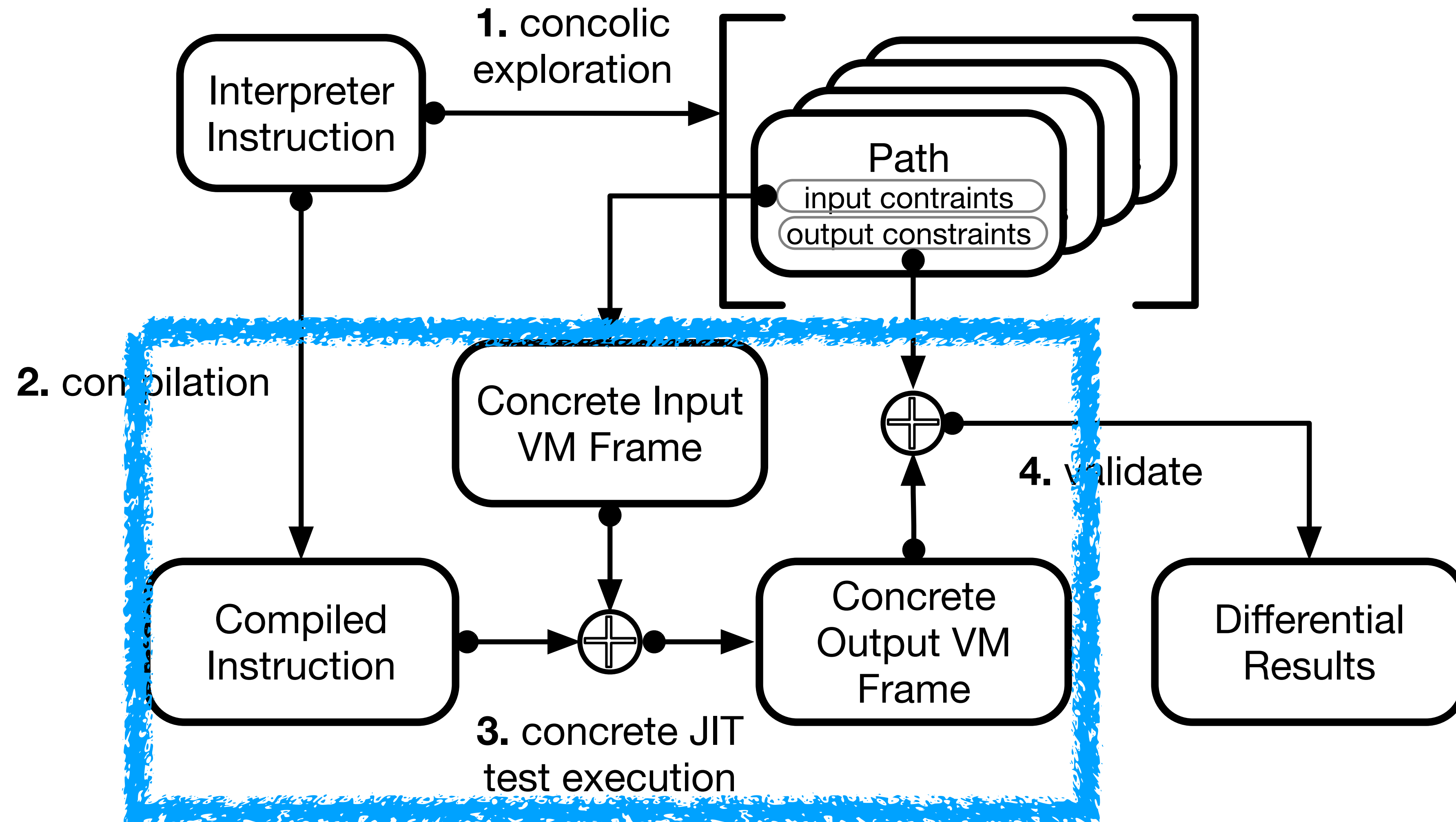




# Implementation View

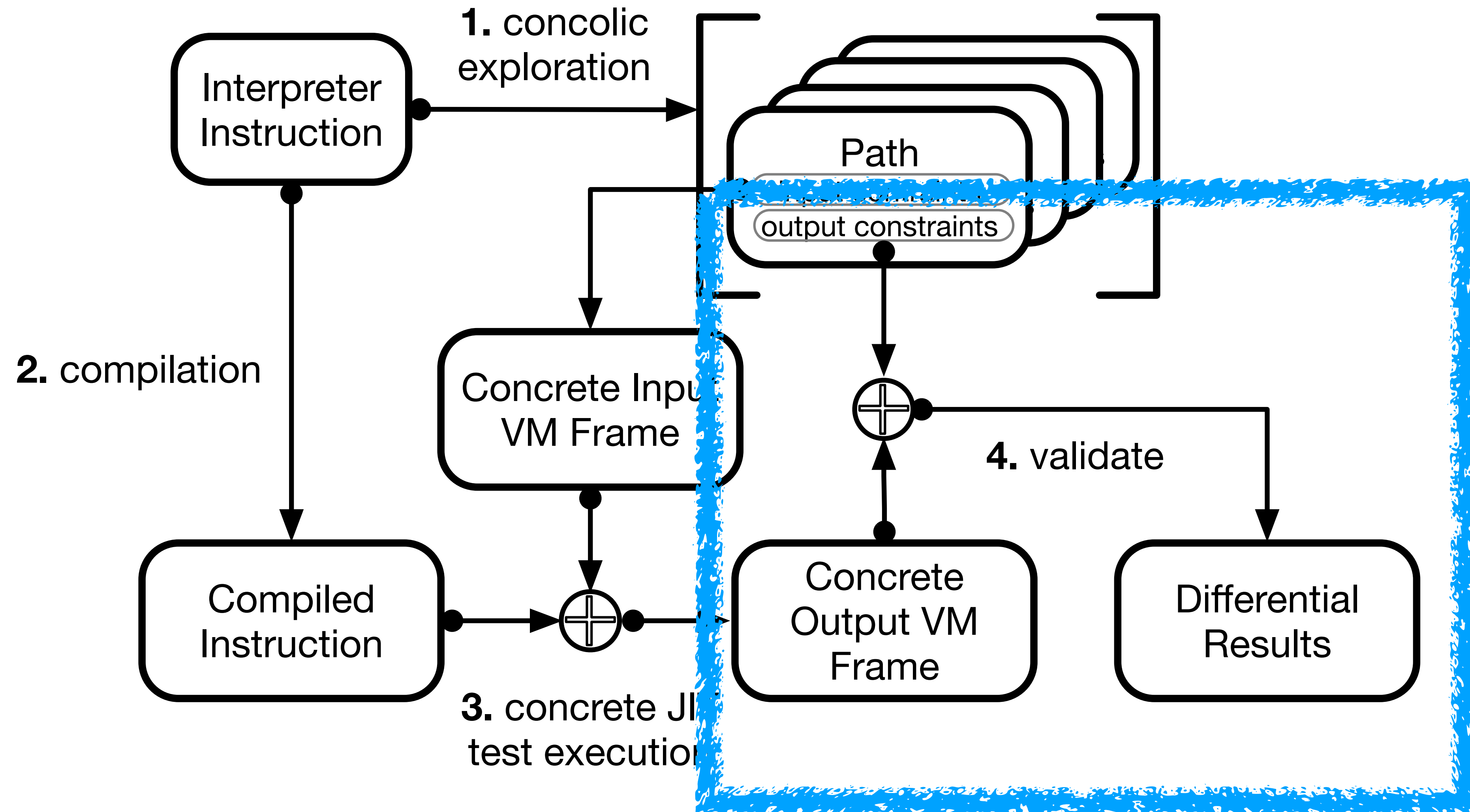


# Implementation View



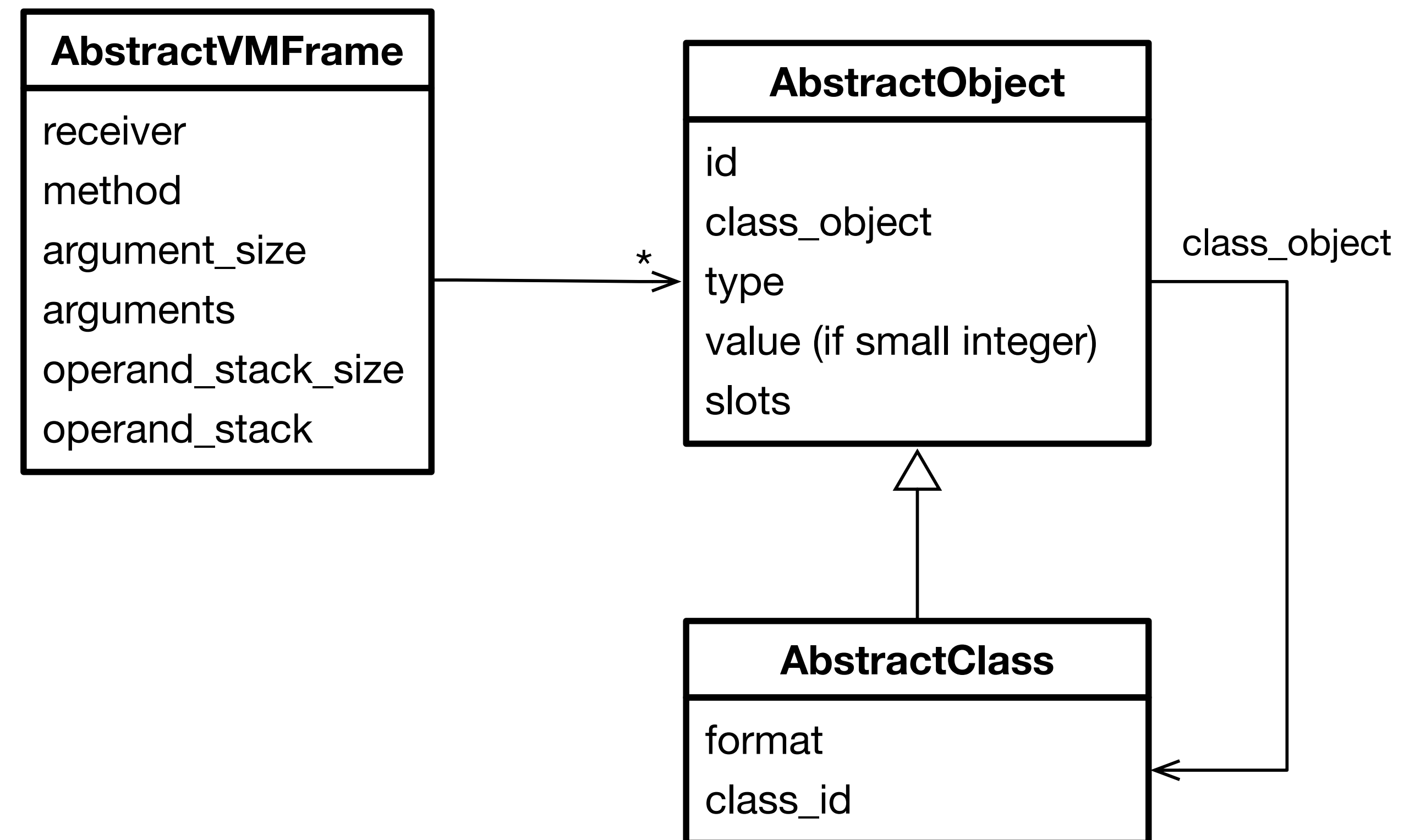


# Implementation View



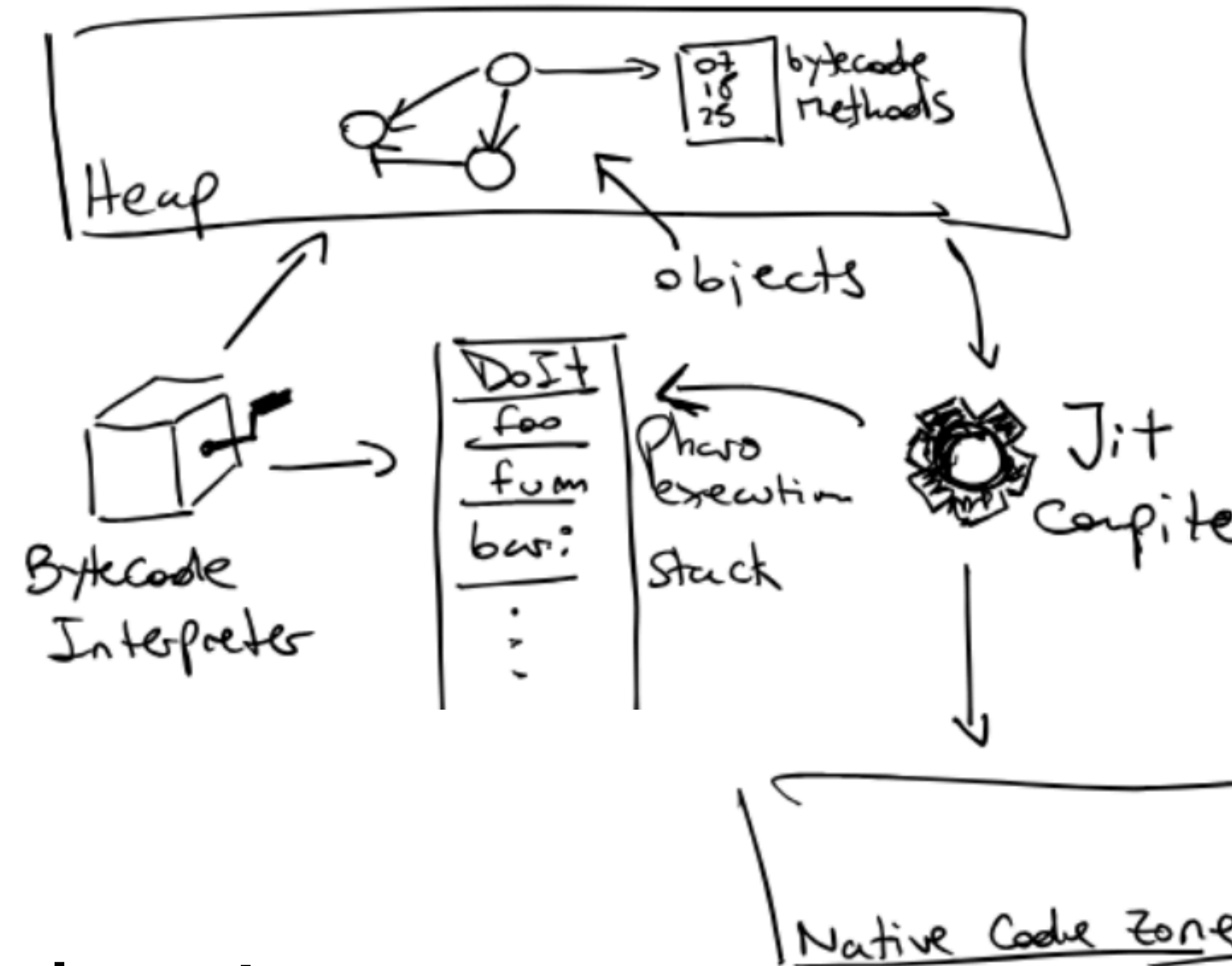
# Concolic Meta-Interpretation Model

- Models VM behaviour during concolic execution
  - Frame
  - Objects + types
  - Classes
- Then flattened into SAT solver equations



# Experimental Context: The Pharo VM

- Interpreted-compiled mixed execution
- Some numbers:
  - 255 stack based bytecodes
  - ~340 primitives/native methods
  - 146 different IR instructions
  - x86, x86-64, ARMv7, ARMv8, RISC-V
- Industrial consortium:
  - **28 International companies, 26 academic partners**



# Previous Manual Testing Effort

- No useful unit tests by ~06/2020
- Large manual testing effort during 2020 while porting to ARM64bits
  - Extended VM simulation with a (TDD compatible) unit testing infrastructure
  - **450+** written tests on the interpreter and the garbage collector\*
  - **580+** written tests on the JIT compiler\*
- Parametrisable for 32 and 64bits, **ARM32, ARM64, x86, x86-64**

\* Numbers by 05/2021



# Evaluation

- 3 bytecode compilers + 1 native method compiler
- 4928 tests generated
- **478 differences**

Compiler	# Tested Instructions	# Interpreter Paths	# Curated Paths	# Differences (%)
Native Methods (primitives)	112	2024	1520	440 (28,95%)
Simple Stack BC Compiler	175	1308	1136	18 (1,59%)
Stack-to-Register BC Compiler	175	1308	1136	10 (0,88%)
Linear-Scan Allocator BC Compiler	175	1308	1136	10 (0,88%)
<b>Total</b>	<b>637</b>	<b>5948</b>	<b>4928</b>	<b>478 (9,7%)</b>





# Analysis of Differences through Manual Inspection

- 91 causes, *6 different categories*
- Errors both in the interpreter AND the compilers
- 14 causes of ***segmentation faults!***

Family	# Cases
Missing interpreter type check	1
Missing compiled type check	13
Optimisation difference	10
Behavioral difference	5
Missing Functionality	60
Simulation Error	2



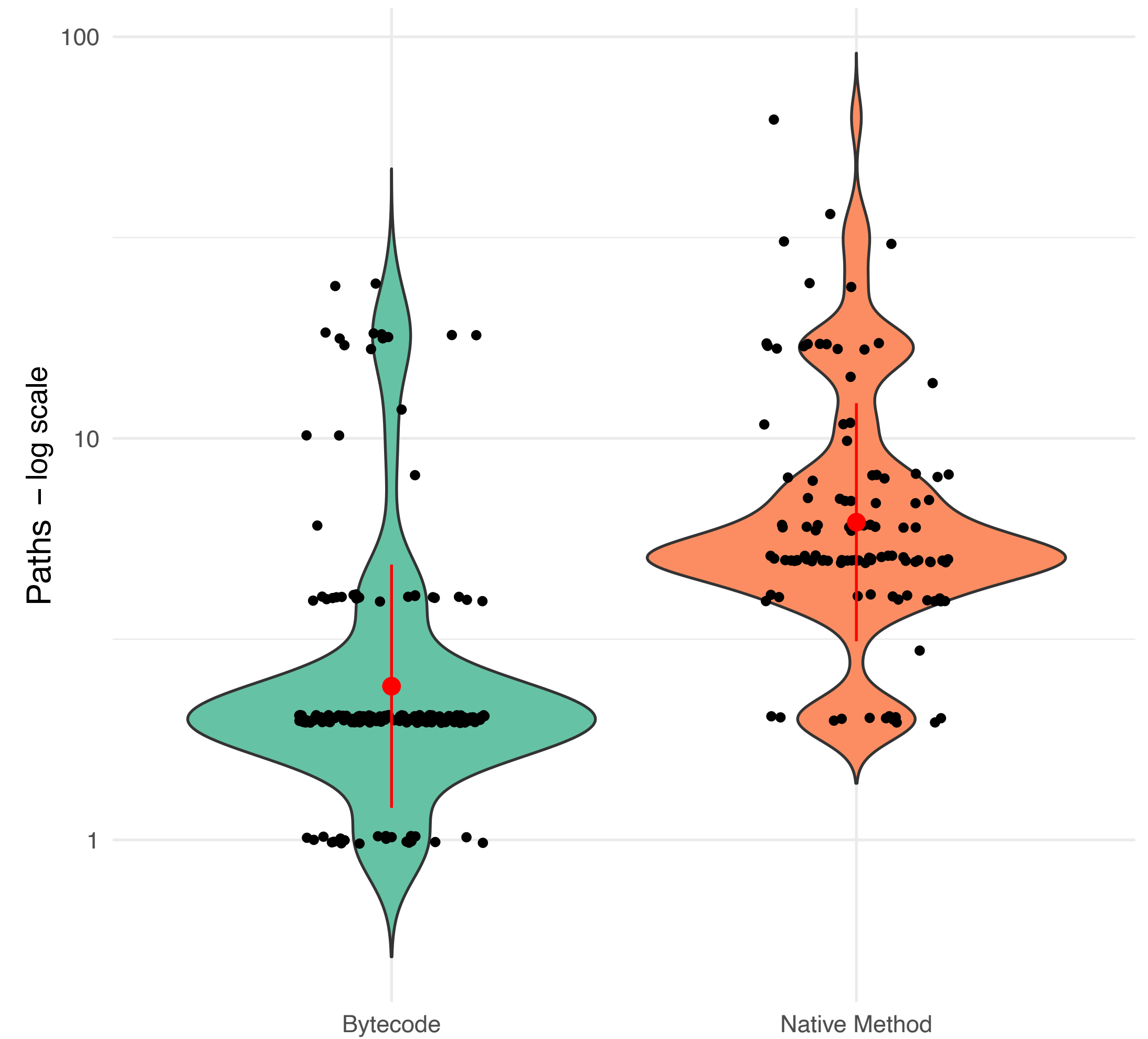
# Characterising Concolic Execution

## Paths per instruction

- Native methods present in average more paths than bytecode instructions

=> longer time to explore

=> potentially more bugs



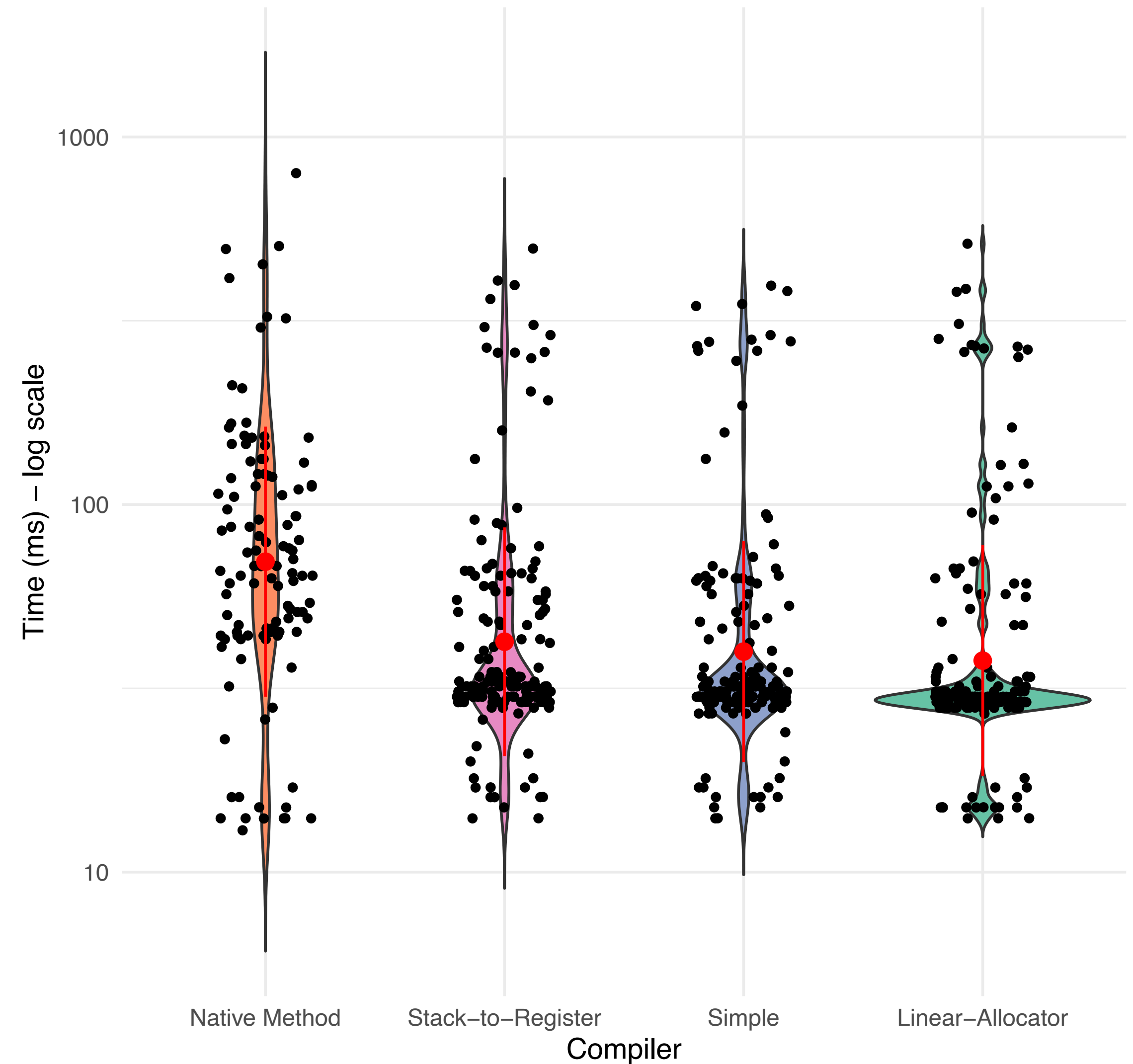
**Paths per  
Type of Instruction**





# Practical and Cheap

- Test generation ~5 minutes
- Total run time of ~10 seconds
  - Avg 30ms per instruction



# More in the article!

- Discovered Bugs
- Concolic Model
- Testing Infrastructure

## Interpreter-Guided Differential JIT Compiler Unit Testing

Guillermo Polito  
Univ. Lille, CNRS, Inria, Centrale Lille,  
UMR 9189 CRISTAL, F-59000 Lille  
France  
guillermo.polito@univ-lille.fr

Stéphane Ducasse  
Univ. Lille, Inria, CNRS, Centrale Lille,  
UMR 9189 CRISTAL  
France  
stephane.ducasse@inria.fr

Pablo Tesone  
Pharo Consortium  
Univ. Lille, Inria, CNRS, Centrale Lille,  
UMR 9189 CRISTAL  
France  
pablo.tesone@inria.fr

### Abstract

Modern language implementations using Virtual Machines feature diverse execution engines such as byte-code interpreters and machine-code dynamic translators, a.k.a. JIT compilers. Validating such engines requires not only validating each in isolation, but also that they are functionally equivalent. Tests should be duplicated for each execution engine, exercising the same execution paths on each of them.

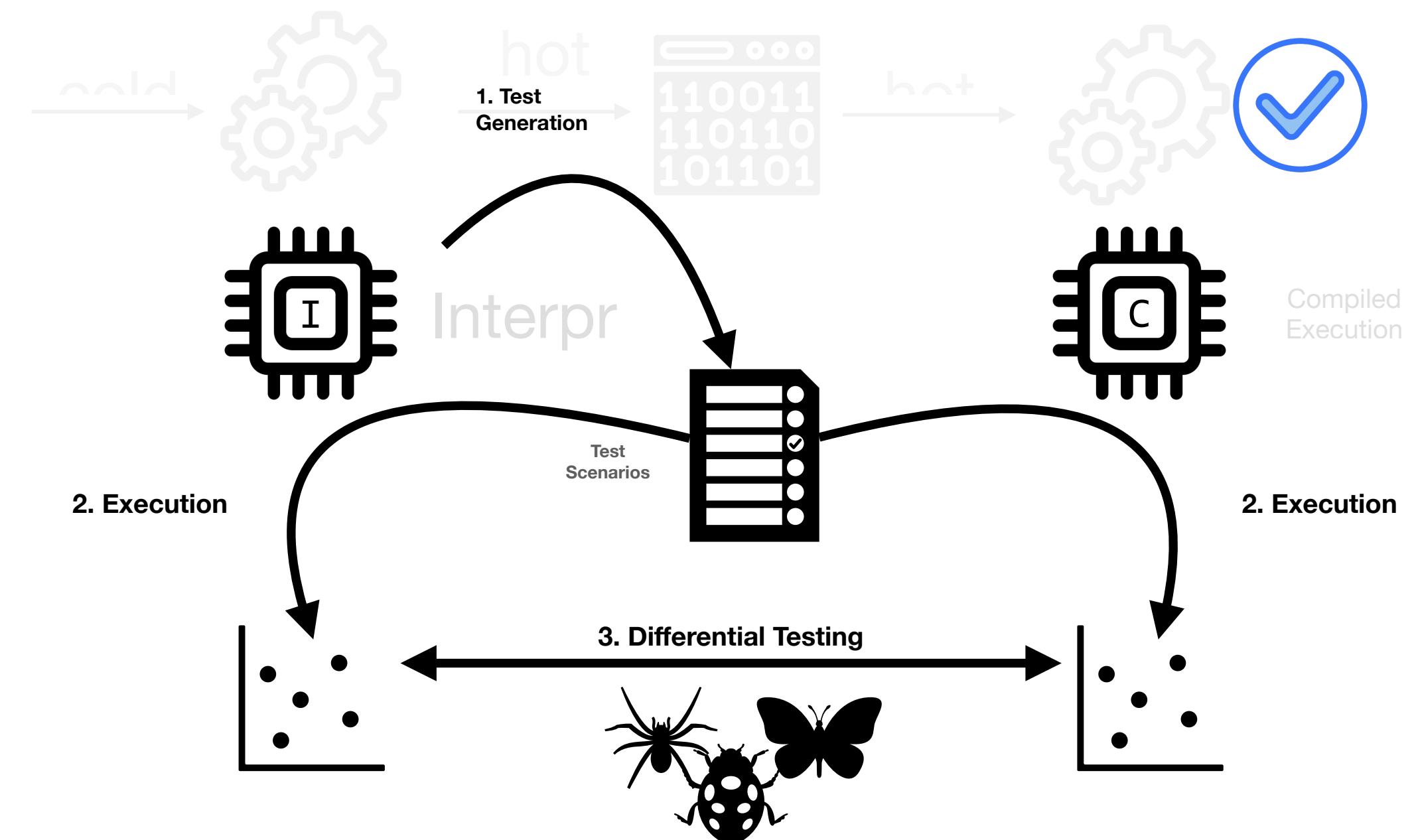
San Diego, CA, USA. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3519939.3523457>

### 1 Introduction

Modern Virtual Machines support code generation for compilation and dynamic code patching for techniques such as inline caching. They are often structured around a baseline code interpreter, a baseline JIT compiler, and a speculative

# Conclusion

- 478 differences found, 91 causes, 6 categories
- Practical:
  - 4928 tests generated in ~8 minutes
  - 4928 tests run in ~40 seconds



**Guille Polito - Pablo Tesone - Stéphane Ducasse**  
**[guillermo.polito@univ-lille.fr](mailto:guillermo.polito@univ-lille.fr)**  
**@guillep**

# Extras



# Concolic Testing through Meta-interpretation

- Idea: Guide test generation by looking at the implementation

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

Different cases  
if  $x > 100$  or  $\leq 100$ !!

Different cases  
if  $x = 1023$  or  $\neq 1023$



# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?





# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?
0	0	$x \leq 100$	





# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?
0	0	$x \leq 100$	$x > 100$



# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?
0	0	$x \leq 100$	$x > 100$
101	0		

Godefroid et al. DART: Directed Automated Random Testing. PLDI' 05



# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?
0	0	$x \leq 100$	$x > 100$
101	0	$x > 100, y \neq 1023$	



# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?
0	0	$x \leq 100$	$x > 100$
101	0	$x > 100, y \neq 1023$	$x > 100, y == 1023$



# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){  
  if (x > 100){  
    if (y == 1023){  
      segfault(!!)  
    }  
  }  
}
```

x	y	constraints	next?
0	0	$x \leq 100$	$x > 100$
101	0	$x > 100, y \neq 1023$	$x > 100, y == 1023$
101	1023		



# Concolic Testing by Example

- **Concrete** + **Symbolic** execution
- Goal: automatically discover *all* execution *paths*

```
int f(int x, int y){
  if (x > 100){
    if (y == 1023){
      segfault(!!)
    } } }
```

x	y	constraints	next?
0	0	$x \leq 100$	$x > 100$
101	0	$x > 100, y \neq 1023$	$x > 100, y == 1023$
101	1023	$x > 100, y \neq 1023$	finished!





# Example

Argument 0 (type)	Argument 1(type)	Path
0 (integer)	0 (integer)	isInteger(arg0), isInteger(arg1), isInteger(arg0+arg1)
0xFFFFFFFF (integer)	1 (integer)	isInteger(arg0), isInteger(arg1), isNotInteger(arg0+arg1)
0 (integer)	object1 (object)	isInteger(arg0), isNotInteger(arg1)
object1 (object)	0 (integer)	isNotInteger(arg0), isInteger(arg1)
object1 (object)	object2 (object)	isNotInteger(arg0), isNotInteger(arg1)

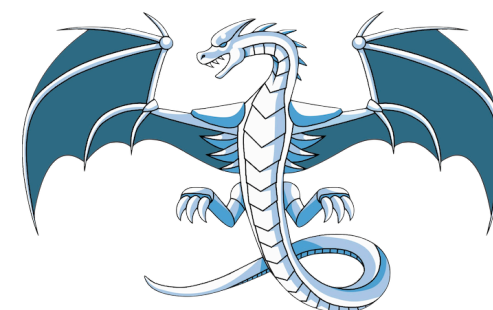
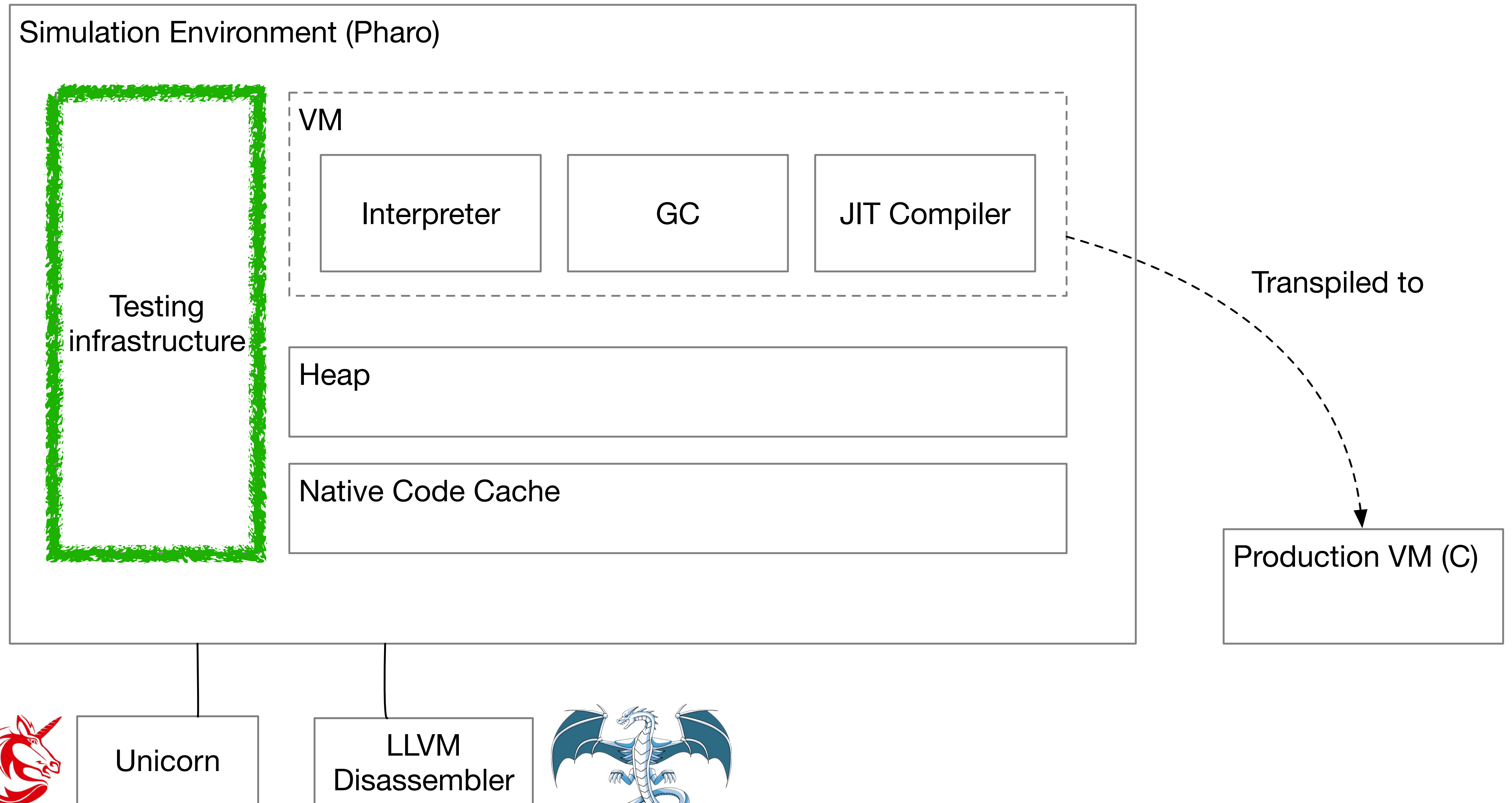
```
1 Interpreter >> bytecodePrimAdd
2 | rcvr arg result |
3 rcvr := self internalStackValue: 1.
4 arg := self internalStackValue: 0.
5 (objectMemory areIntegers: rcvr and: arg) ifTrue: [
6     result := (objectMemory integerValueOf: rcvr) + (
7         objectMemory integerValueOf: arg).
8     "Check for overflow"
9     (objectMemory isIntegerValue: result) ifTrue: [
10         self
11         internalPop: 2
12         thenPush: (objectMemory integerObjectOf: result).
13         ^ self fetchNextBytecode "success"]].
14 "Slow path, message send"
15 self normalSend
```

```
1 ... # previous bytecode IR
2     checkSmallInteger t0
3     jumpzero notsmi
4     checkSmallInteger t1
5     jumpzero notsmi
6     t2 := t0 + t1
7     jumpIfNotOverflow continue
8 notsmi: #slow case first send
9     t2 := send #+ t0 t1
10 continue:
11 ... # following bytecode IR
```

**Listing 1.** Excerpt of the byte-code interpretation implementing addition in the Pharo Virtual Machine.

**Listing 2.** Illustration of the Intermediate Representation instructions created when compiling the byte-code instruction in Listing 1.

# Simulation + Testing Environment



# Unit Testing Infrastructure Comparison

	Real Hardware Execution	Full-System Simulation	Unit-Testing
Feedback-cycle speed	Very low	Low	High
Availability	Low	High	High
Reproducibility	Low	Low	High
Precision	High	Low	Low
Debuggability	Low	High	High





