

# Cross-Architecture Lifter Synthesis

Rijnard van Tonder and Claire Le Goues

# Cross-Architecture Lifter Synthesis

Rijnard van Tonder and Claire Le Goues

# What is a Lifter?

# What is a Lifter?

x86

ARM

...

Lower level code

# What is a Lifter?

x86

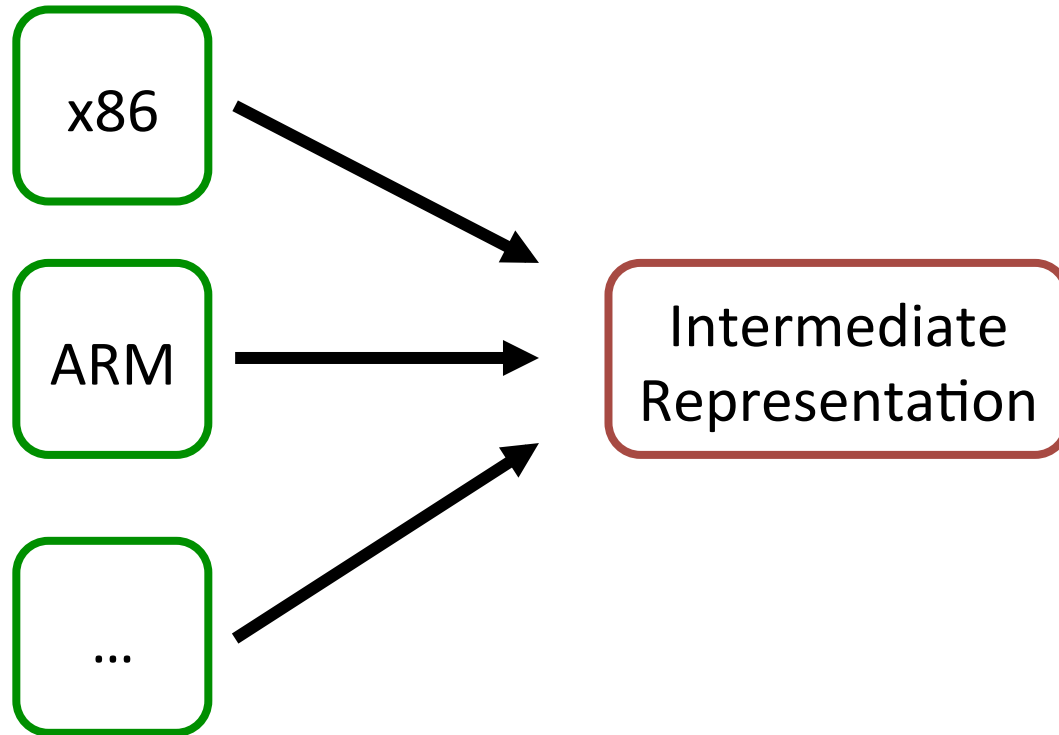
ARM

...

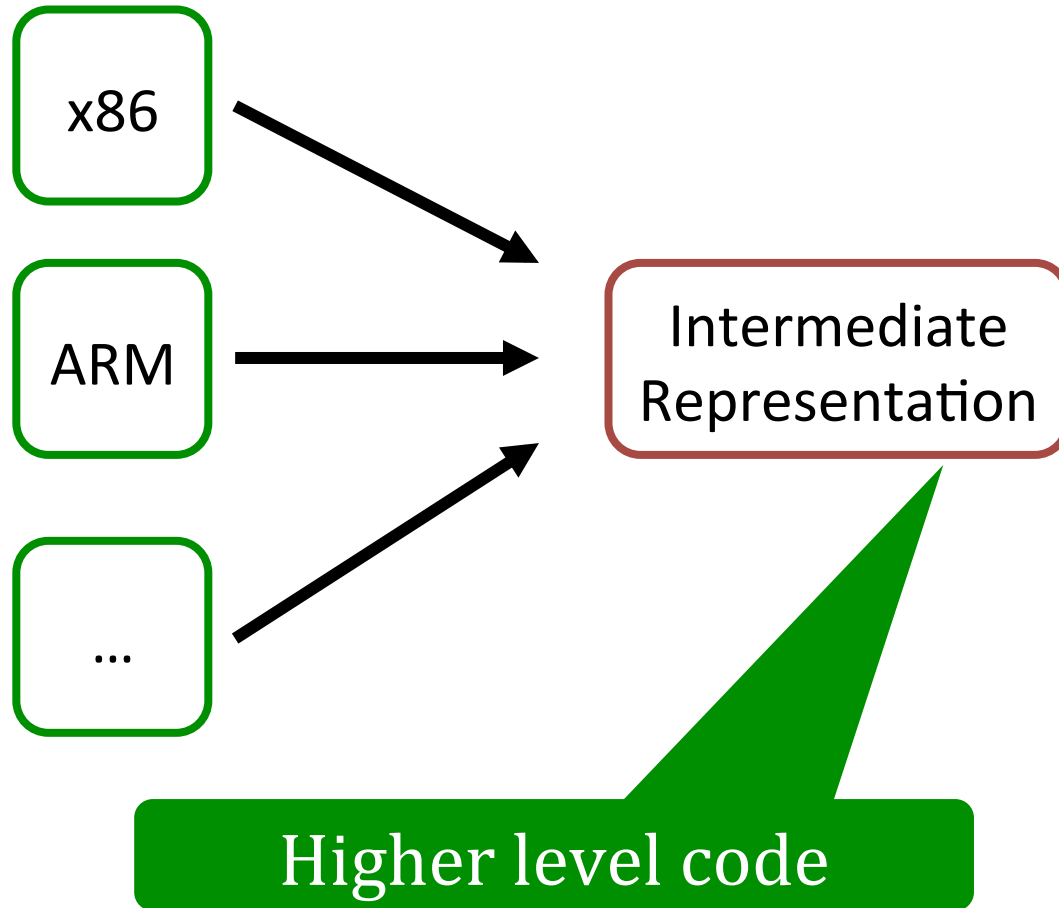
Lower level code

Architecture-specific

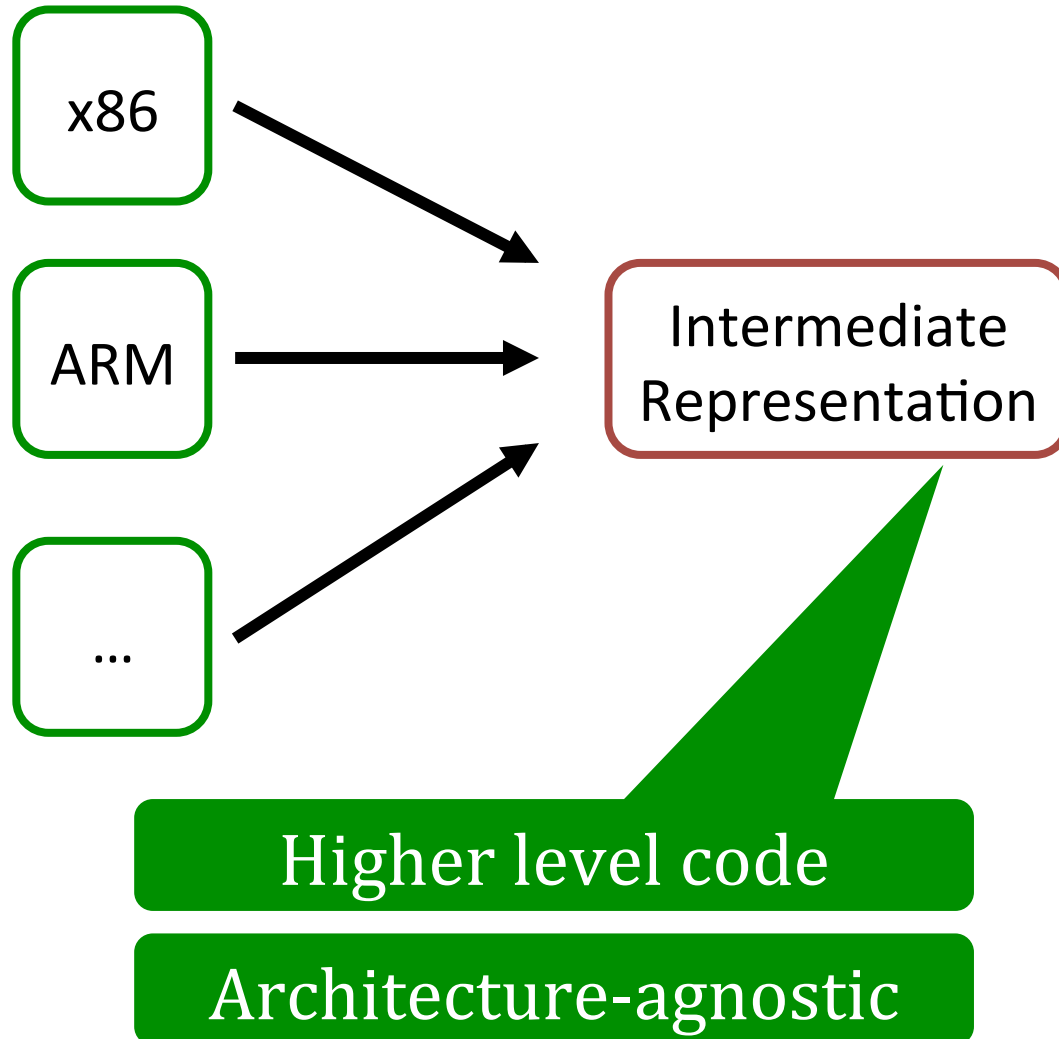
# What is a Lifter?



# What is a Lifter?



# What is a Lifter?





# Why Lifters?

# Why Lifters?

- Binary analysis

# Why Lifters?

- Binary analysis
- Simpler semantic abstraction
  - E.g., symbolic execution

# Why Lifters?

- Binary analysis
- Simpler semantic abstraction
  - E.g., symbolic execution
- Reuse analysis components
  - Control flow graph construction
  - Single Static Assignment (SSA) conversion

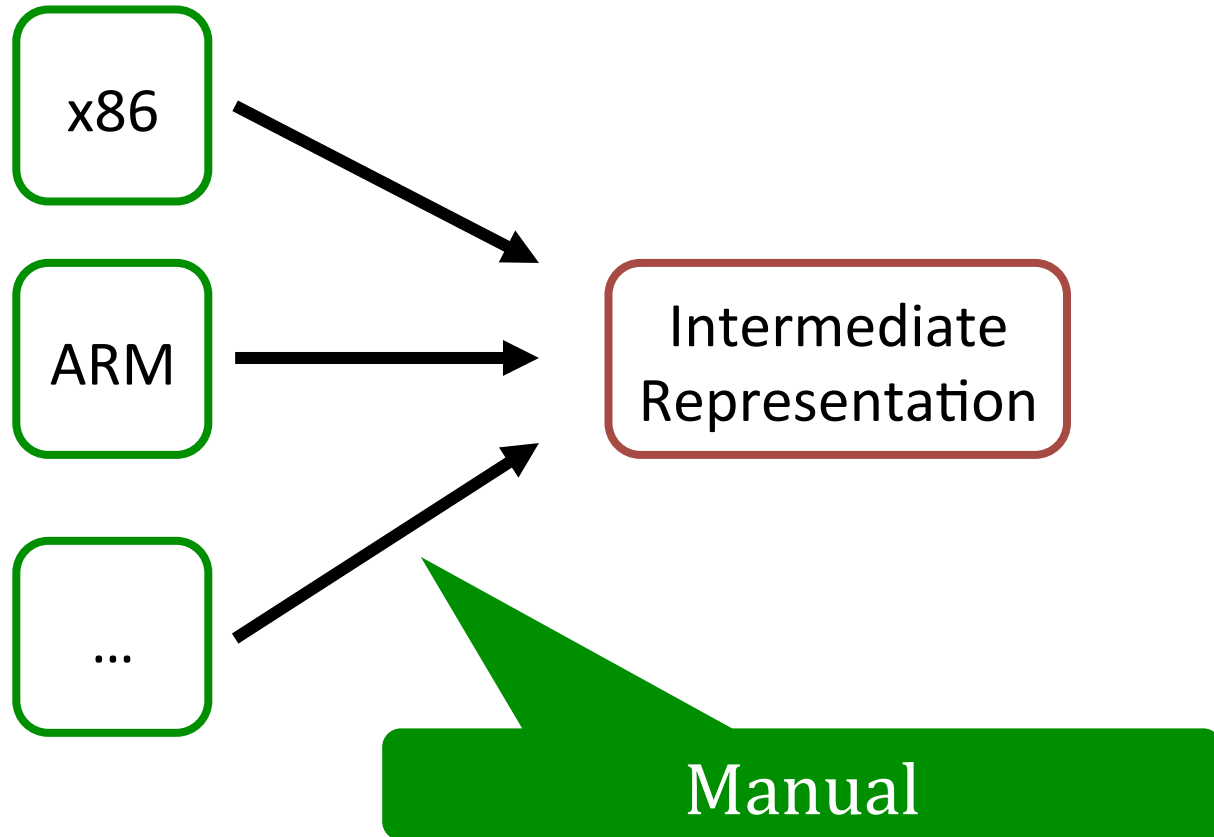
# Compilers and Decompilers

- Compilers lose information
  - Translation is not a bijection

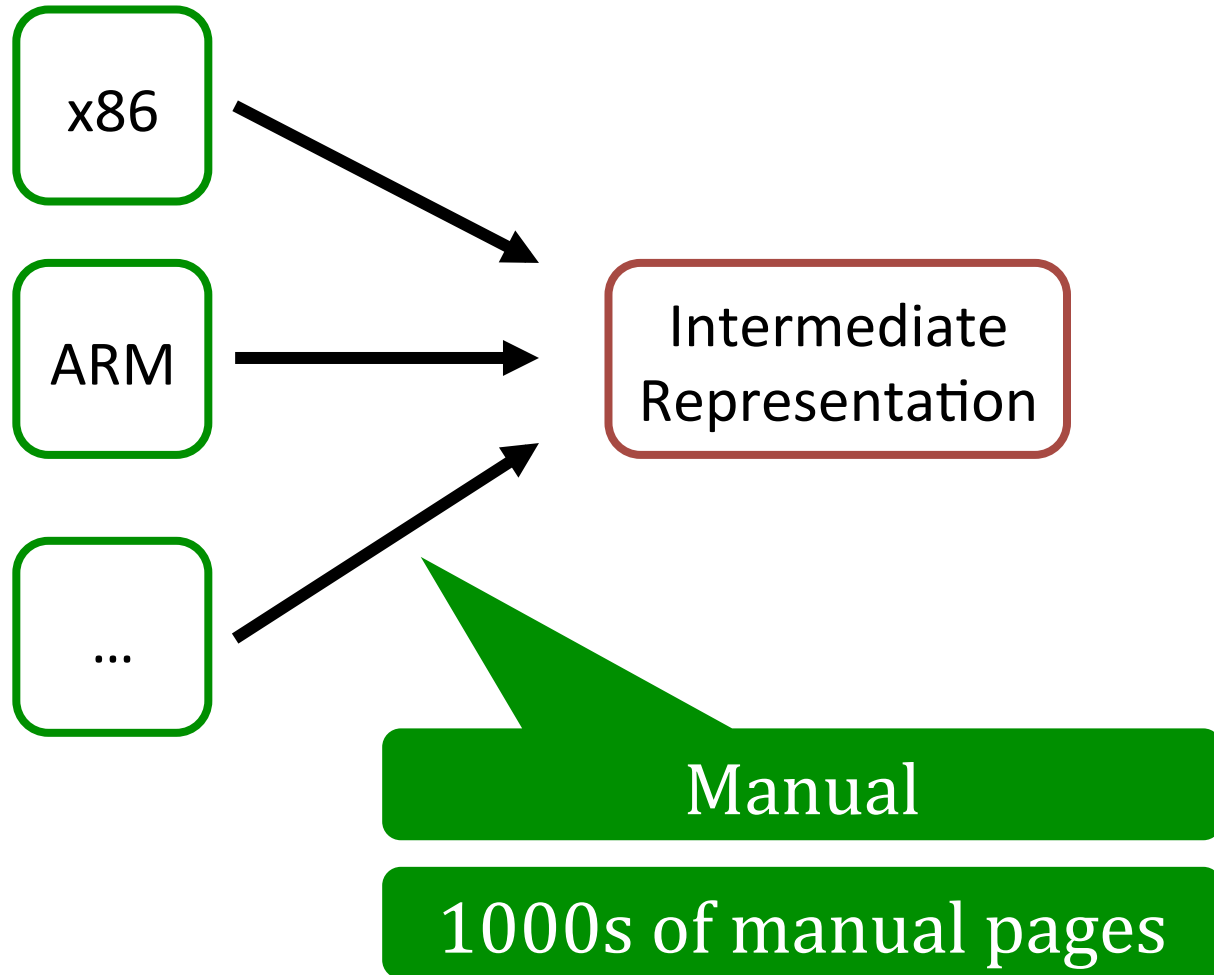
# Compilers and Decompilers

- Compilers lose information
  - Translation is not a bijection
- Decompilers recover more features
  - Often architecture-specific (e.g., ABI)

# The Problem: Translating Multiple Architectures

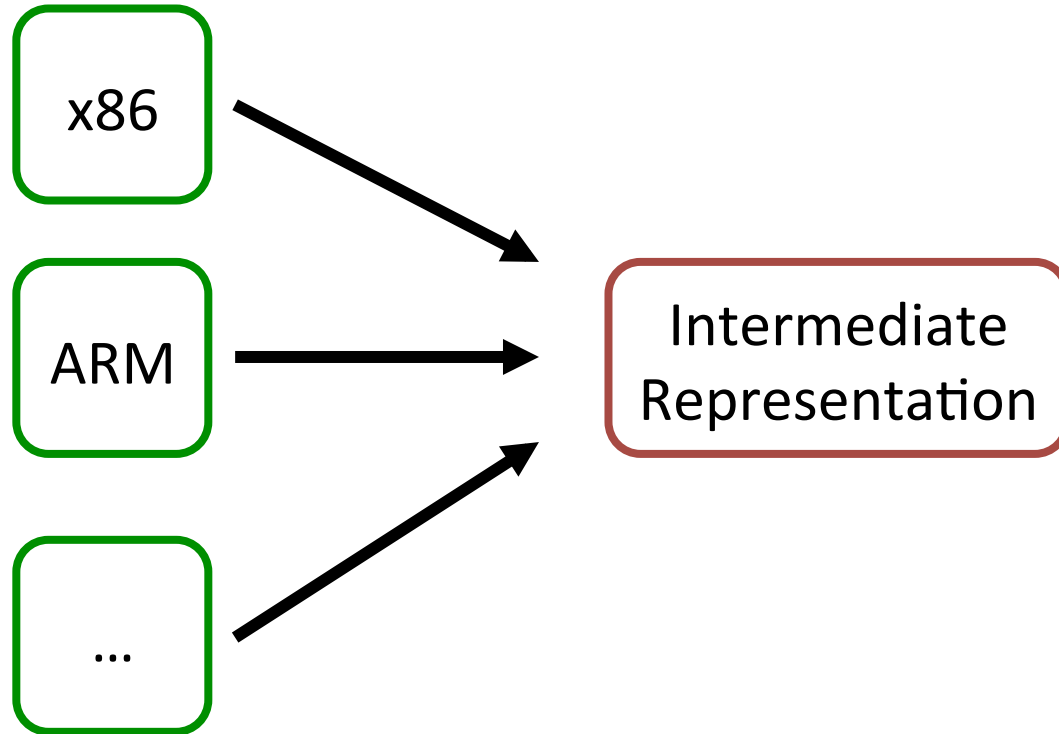


# The Problem: Translating Multiple Architectures

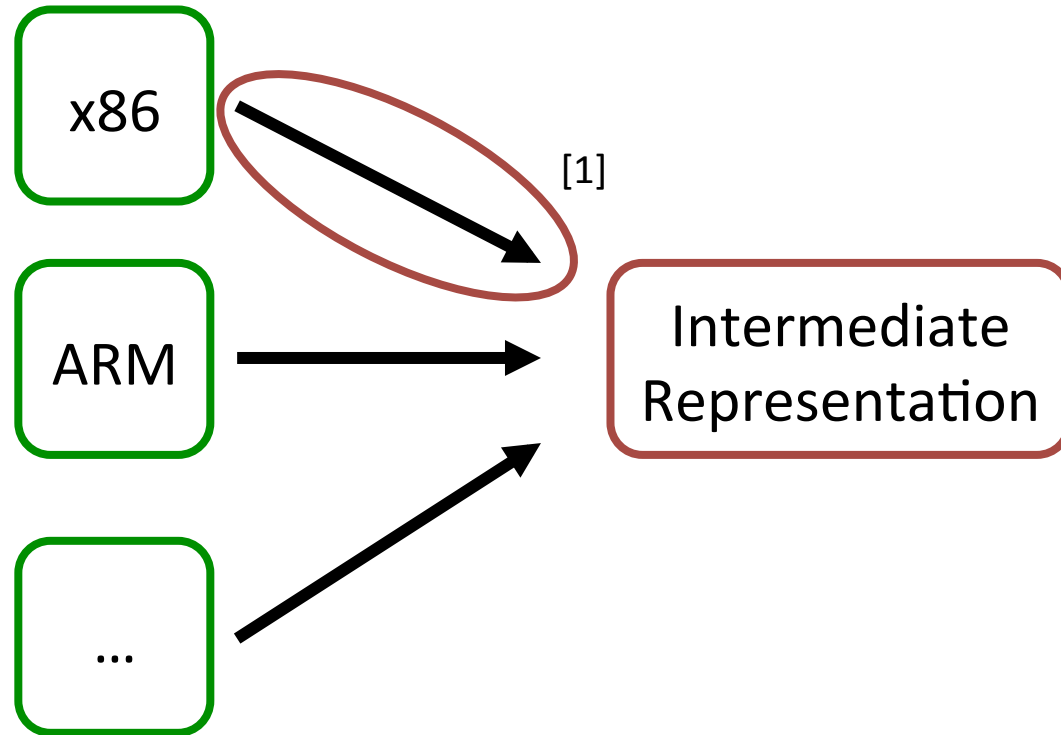




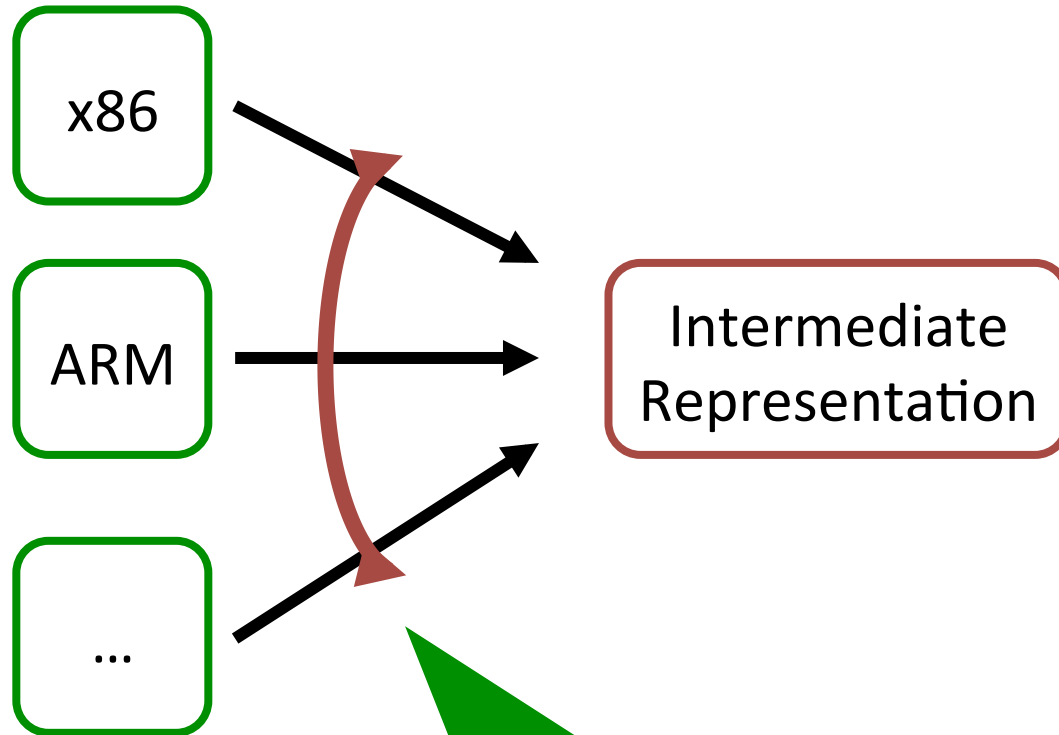
# Our Goal: Automate across Architectures



# Our Goal: Automate across Architectures

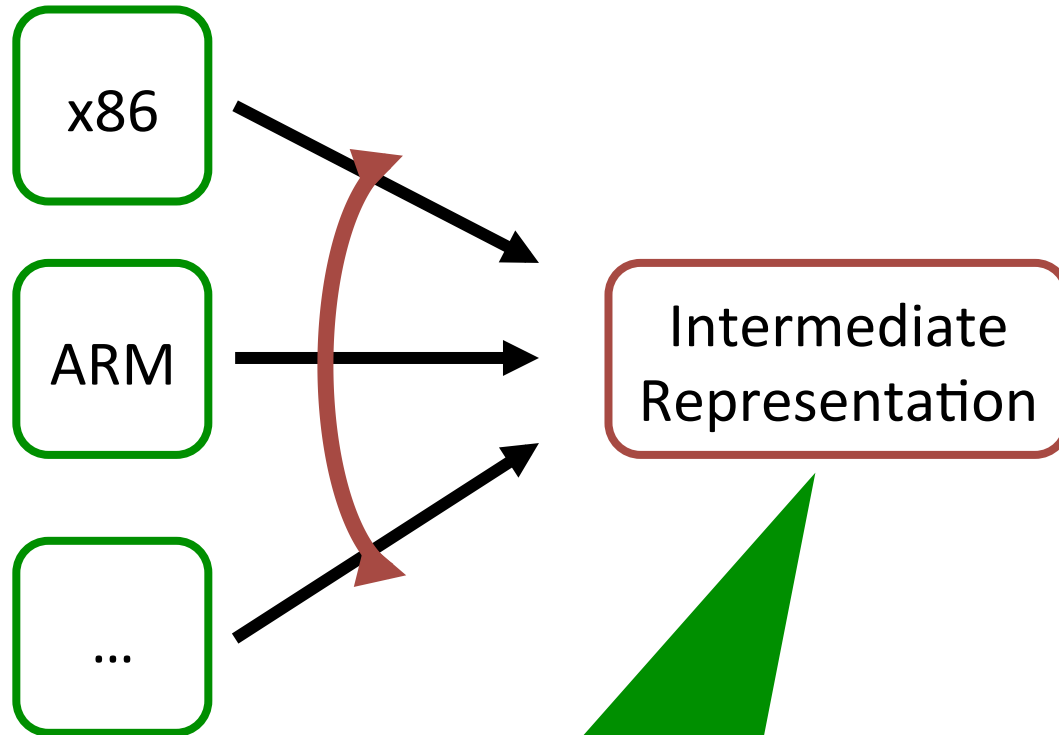


# Our Goal: Automate across Architectures



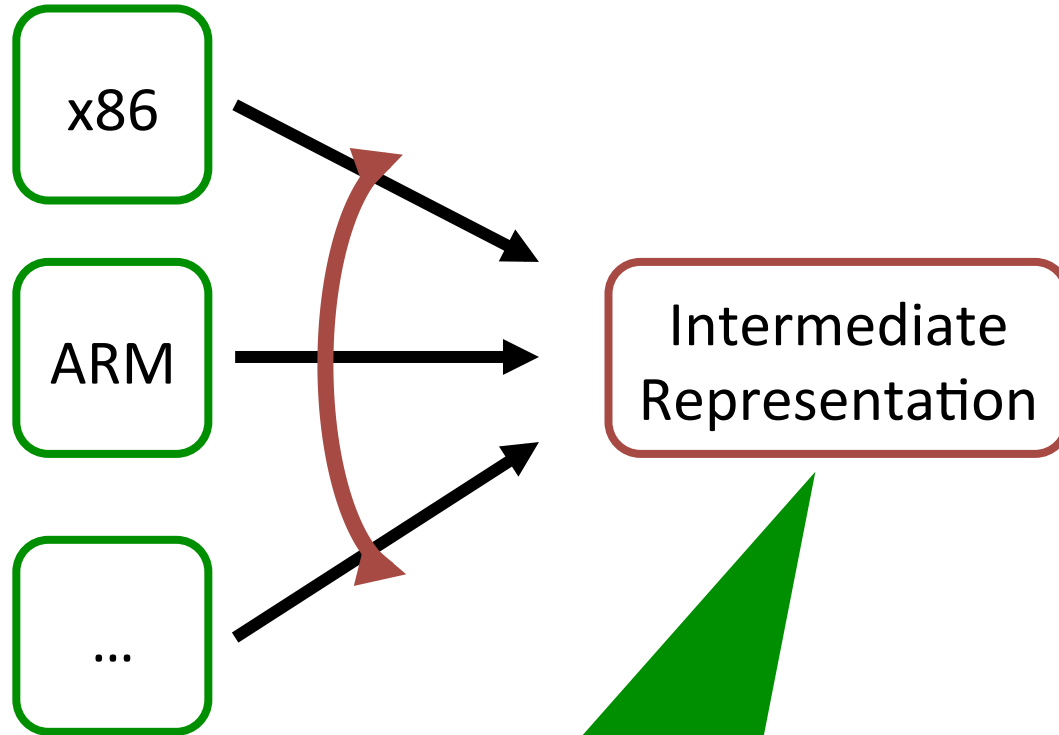
**Cross-Architecture reuse**

# Our Goal: Automate across Architectures



**Bug-finding Analyses “for free”**

# Our Goal: Automate across Architectures



Dataflow framework

VC Generation

Bug-finding Analyses “for free”

Taint analysis

# What We Do

# What We Do

- IR translation as a Syntax-Guided Synthesis problem

# What We Do

- IR translation as a Syntax-Guided Synthesis problem

$$\forall x . \varphi(x, P(x))$$



# What We Do

- IR translation as a Syntax-Guided Synthesis problem

$$\forall x . \varphi(x, P(x)) \equiv \forall x . \text{oracle}(x) = P(x)$$

# What We Do

- IR translation as a Syntax-Guided Synthesis problem

$$\forall x . \varphi(x, P(x)) \equiv \forall x . \text{oracle}(x) = P(x)$$



Correctness specification

# What We Do

- IR translation as a Syntax-Guided Synthesis problem

$$\forall x . \varphi(x, P(x)) \equiv \forall x . \text{oracle}(x) = P(x)$$



Dynamic Input/Output Pairs

# What We Do

- IR translation as a Syntax-Guided Synthesis problem

$$\forall x . \varphi(x, P(x)) \equiv \forall x . \text{oracle}(x) = P(x)$$



Dynamic Input/Output Pairs

IR Sketches

# What We Do

- IR translation as a Syntax-Guided Synthesis problem

$$\forall x . \varphi(x, P(x)) \equiv \forall x . \text{oracle}(x) = P(x)$$

Key Insight:  
Learn IR Sketches from  
existing Lifter Productions



IR Sketches

# Learning Sketch Templates: Example

# Learning Sketch Templates: Example

**ARM**

**add R3, R0**

# Learning Sketch Templates: Example

**ARM**  
**add R3, R0**

**SOURCE**

**Native Instruction**



# Learning Sketch Templates: Example

**ARM**  
**add R3, R0**

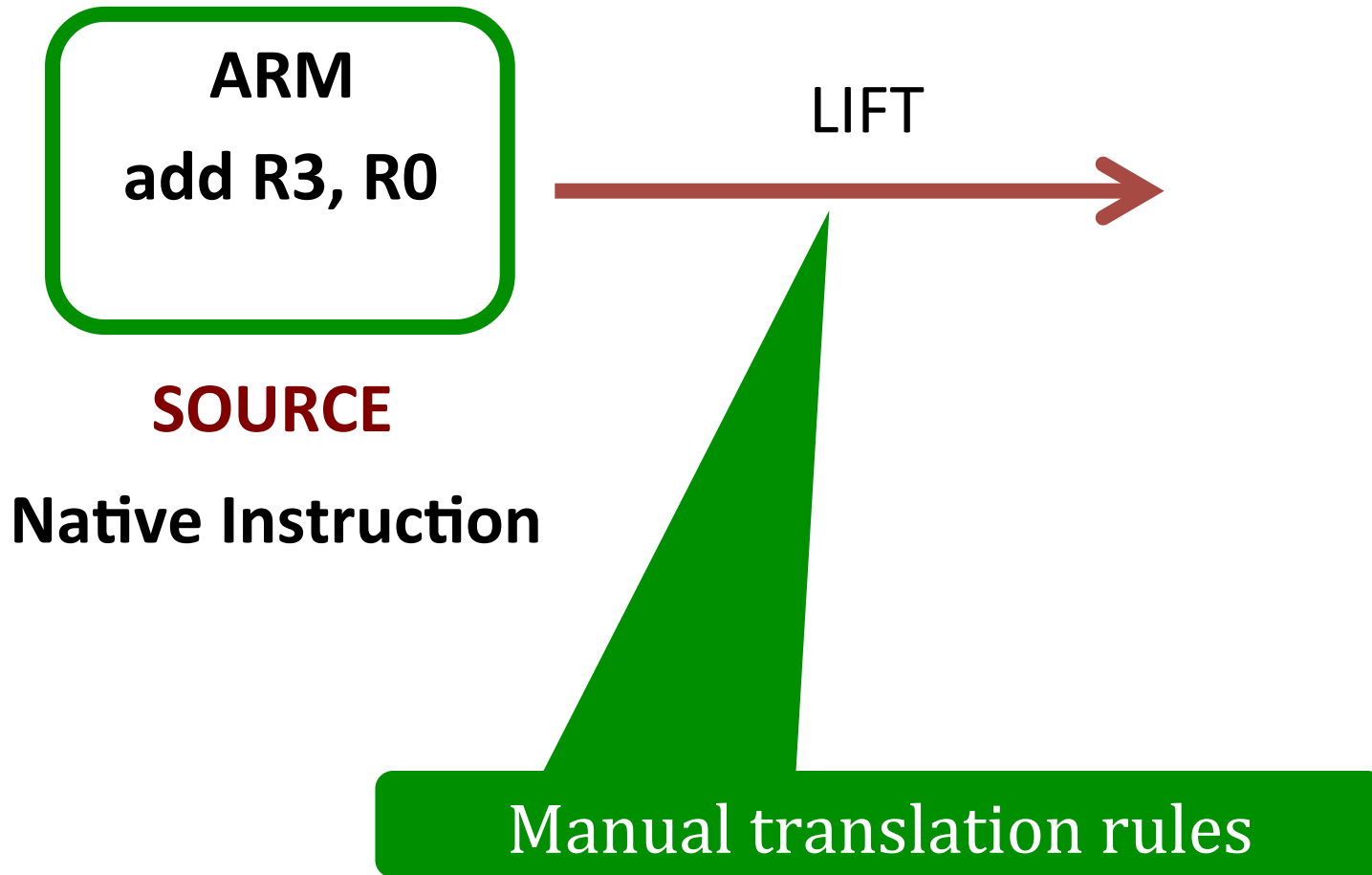
**SOURCE**

**Native Instruction**

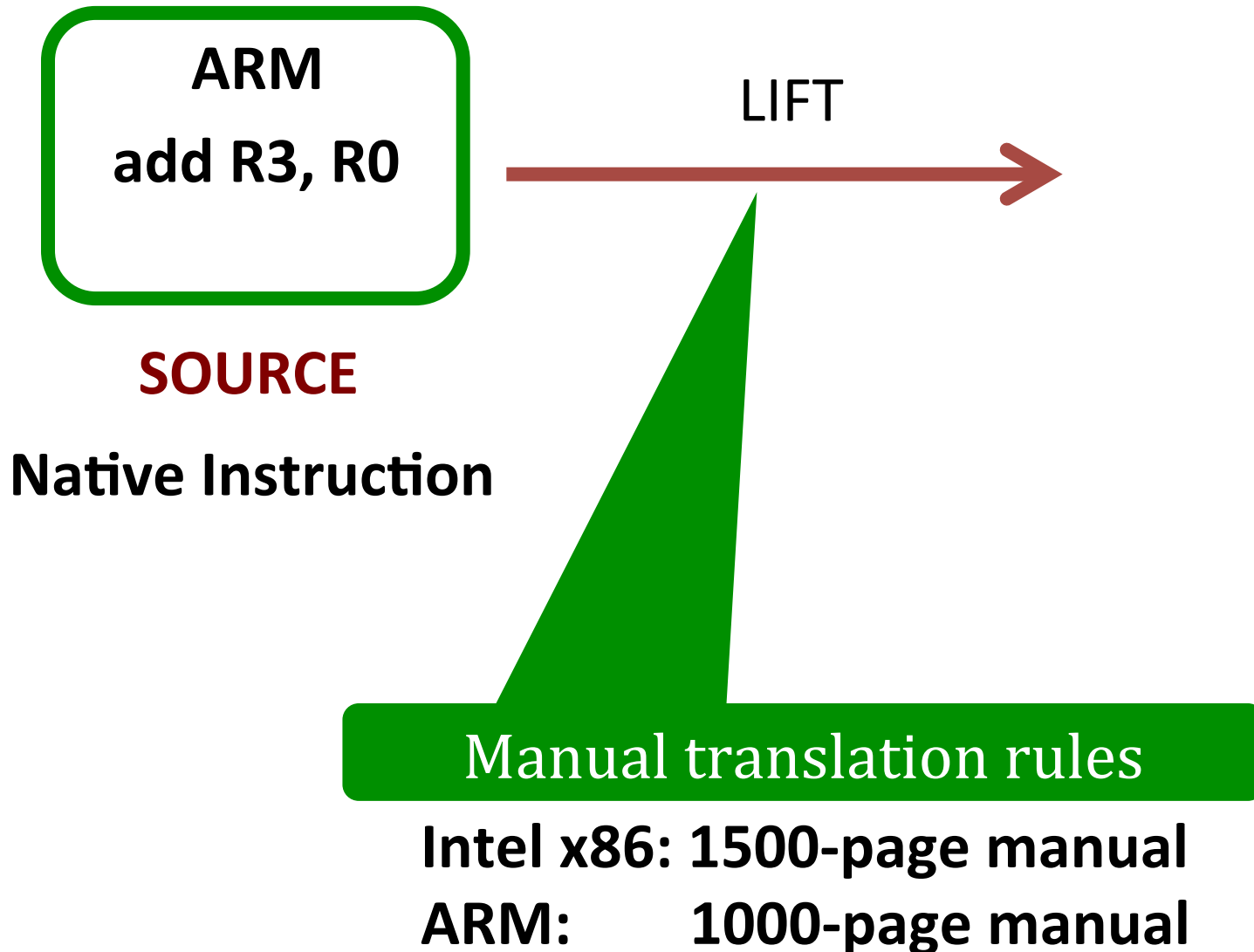


**Op codes and register names**

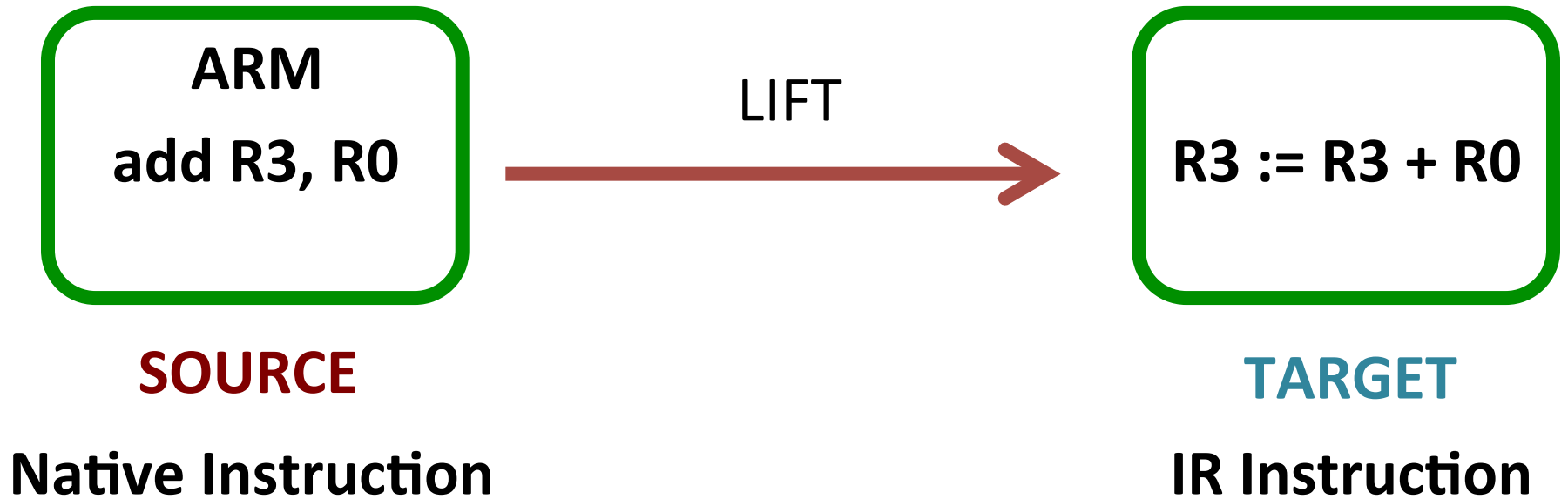
# Learning Sketch Templates: Example



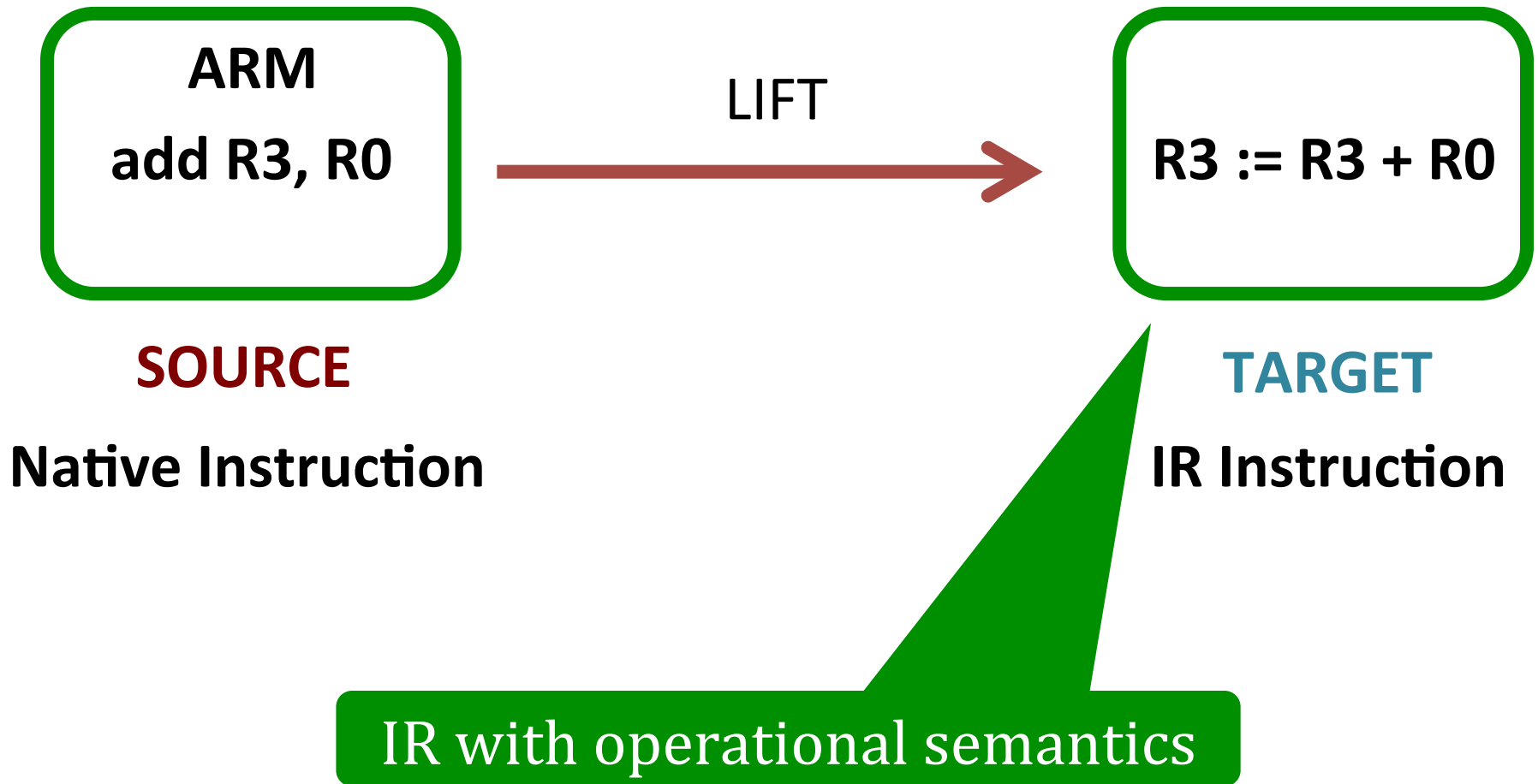
# Learning Sketch Templates: Example



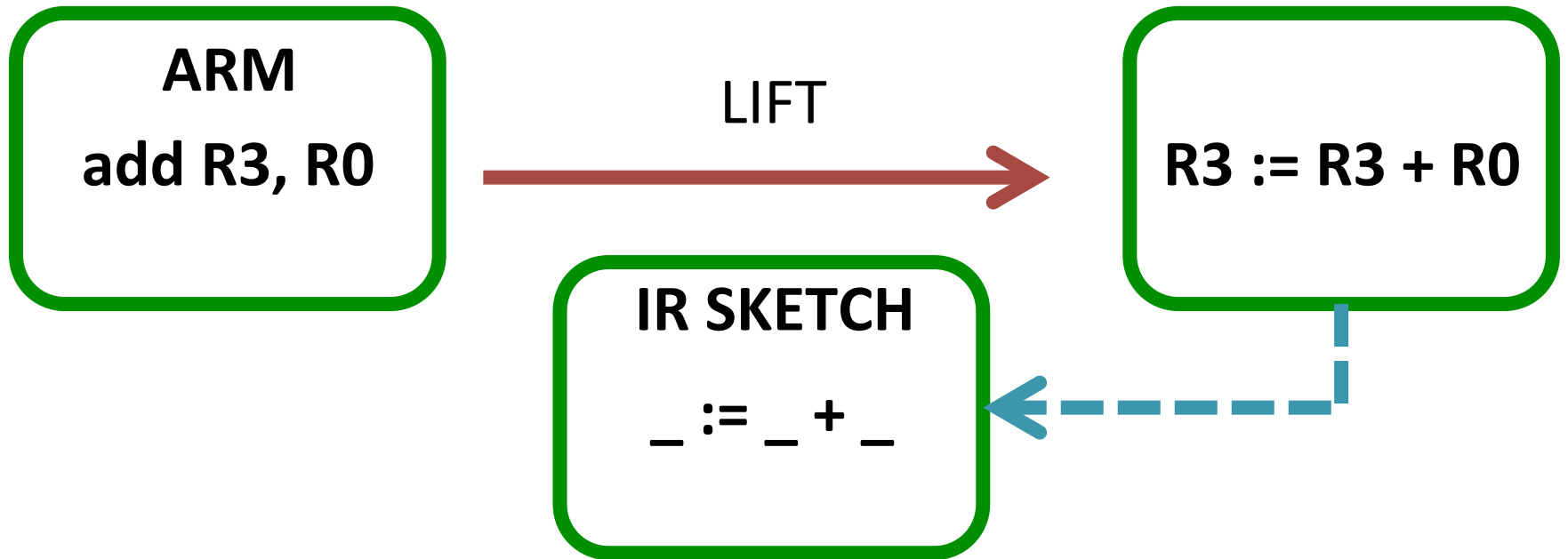
# Learning Sketch Templates: Example



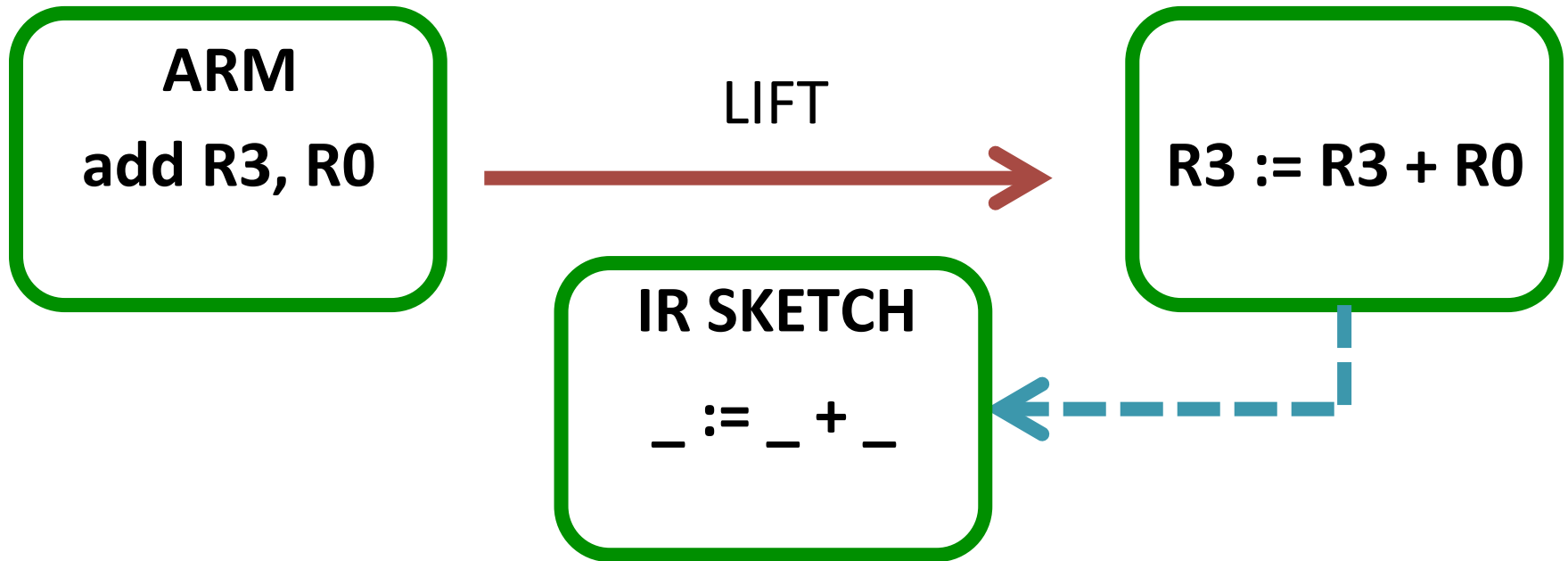
# Learning Sketch Templates: Example



# Learning Sketch Templates: Example



# Learning Sketch Templates: Example



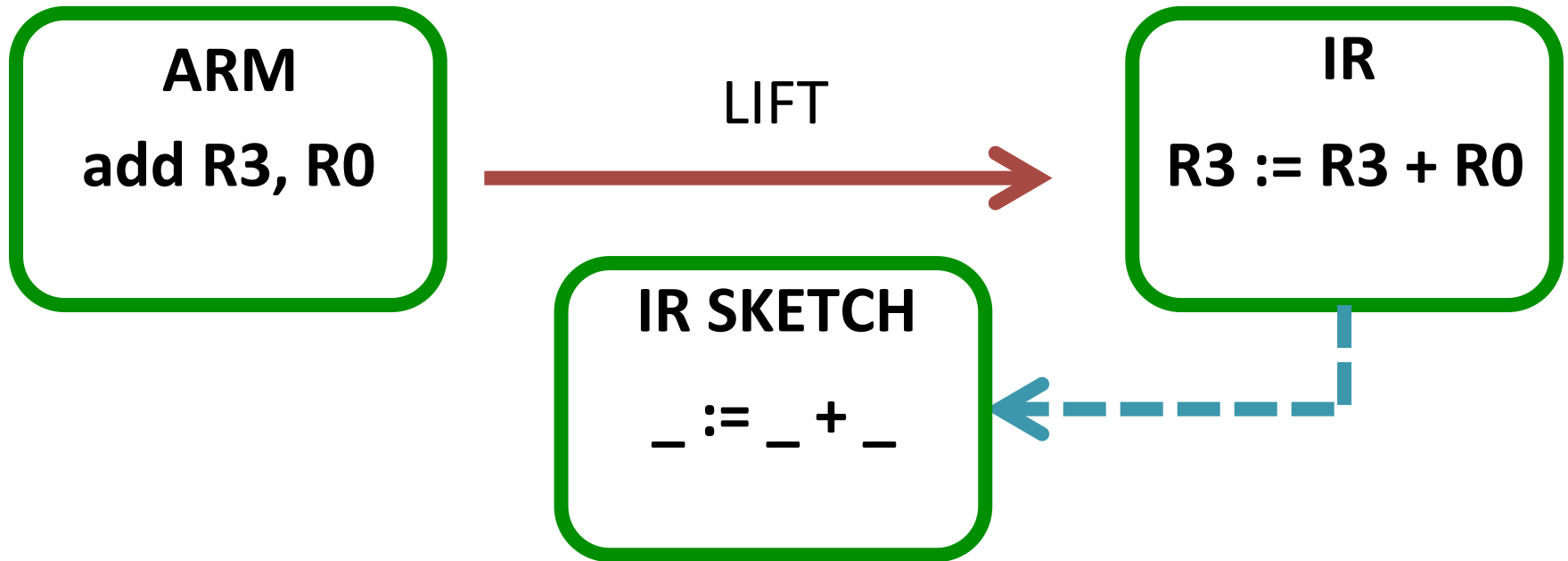
Create sketches from concrete expressions

# Learning Sketch Templates: Example

- Learn templates from lifter output for a **supported** architecture to synthesize one for an **unsupported** architecture

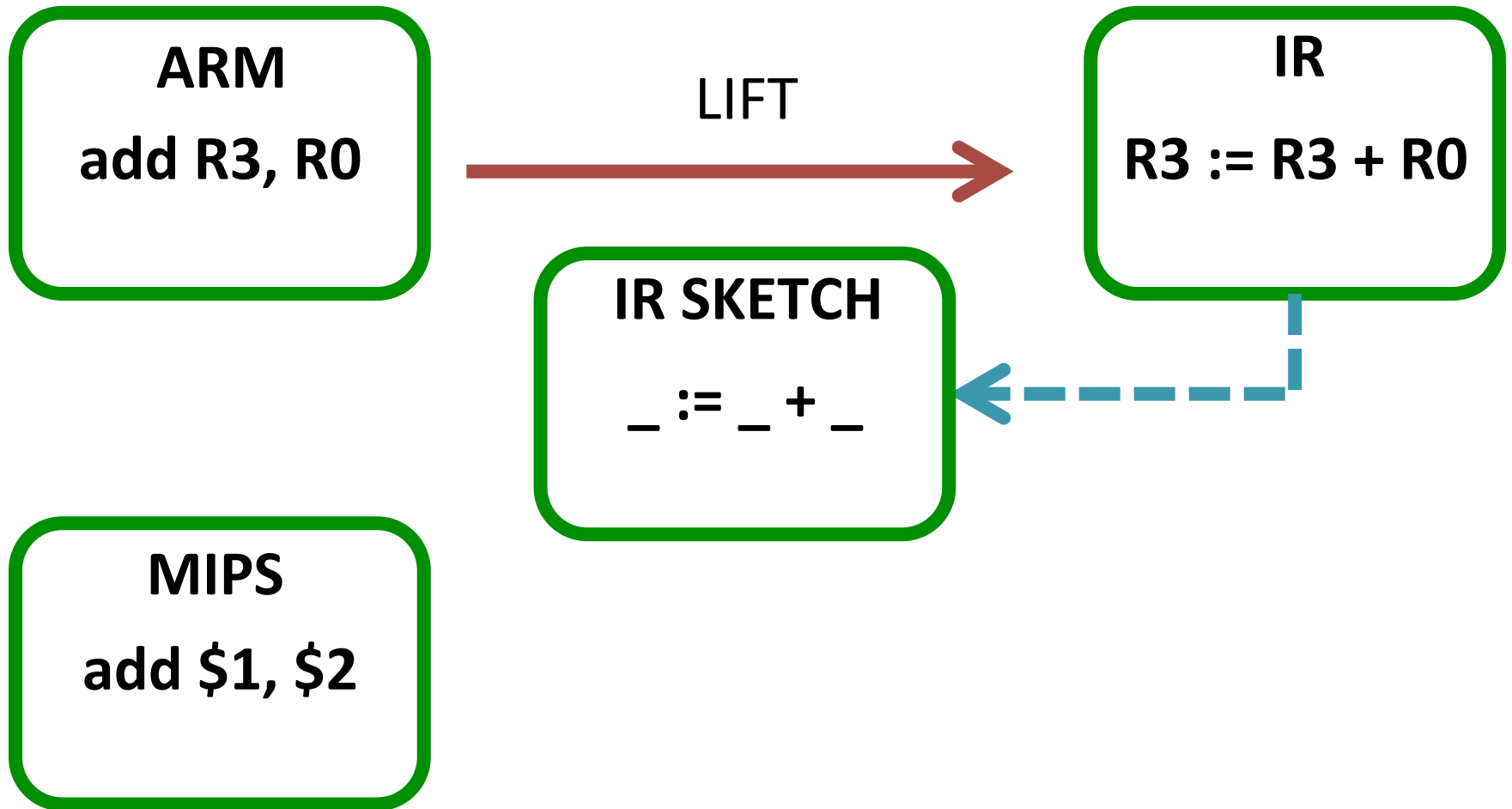


# Learning Sketch Templates: Example

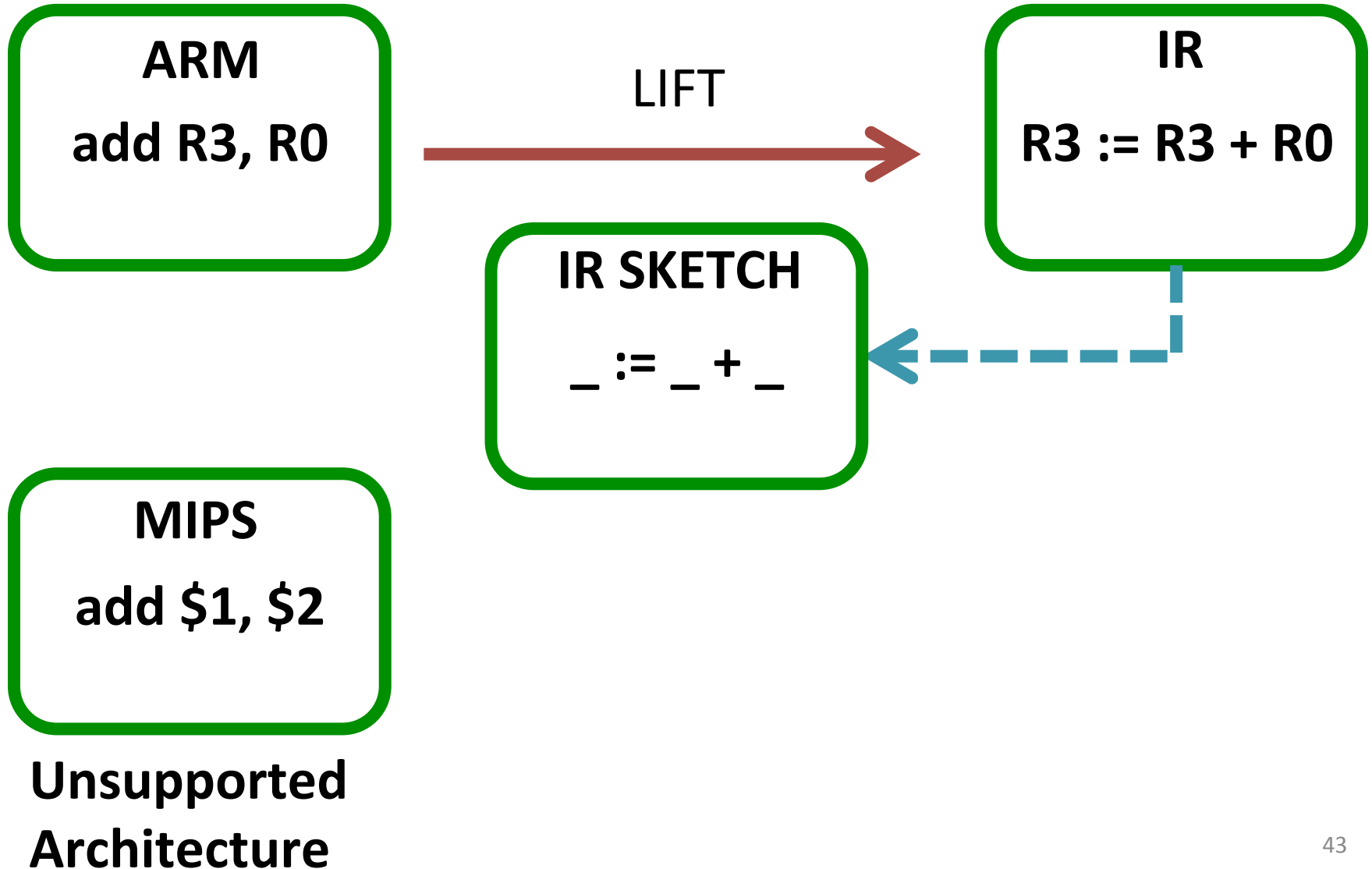


Create sketches from concrete expressions

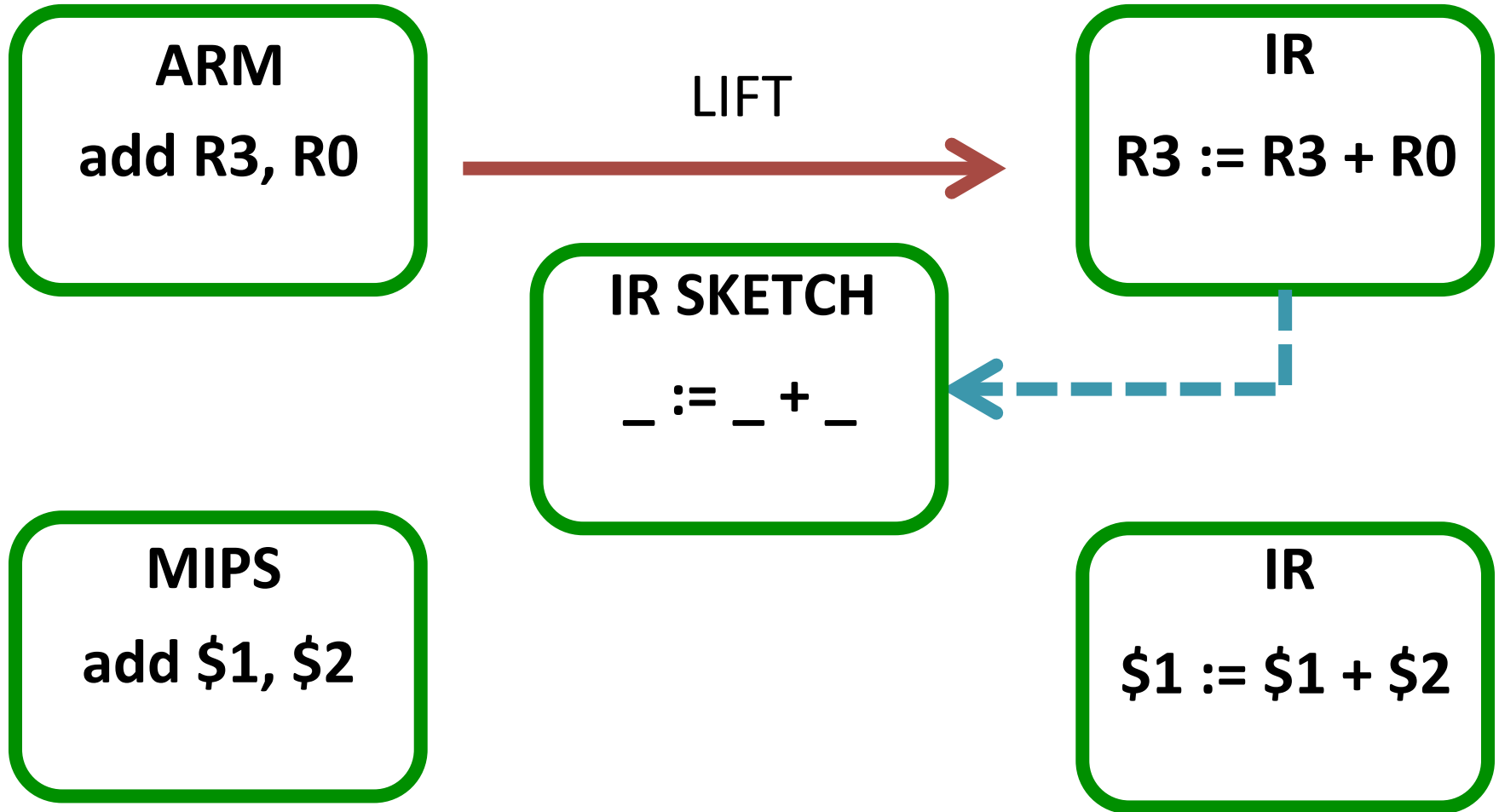
# Learning Sketch Templates: Example



# Learning Sketch Templates: Example



# Learning Sketch Templates: Example



Unsupported  
Architecture

# Learning Sketch Templates: Example

**MIPS**

**add \$1, \$2**

**Unsupported  
Architecture**

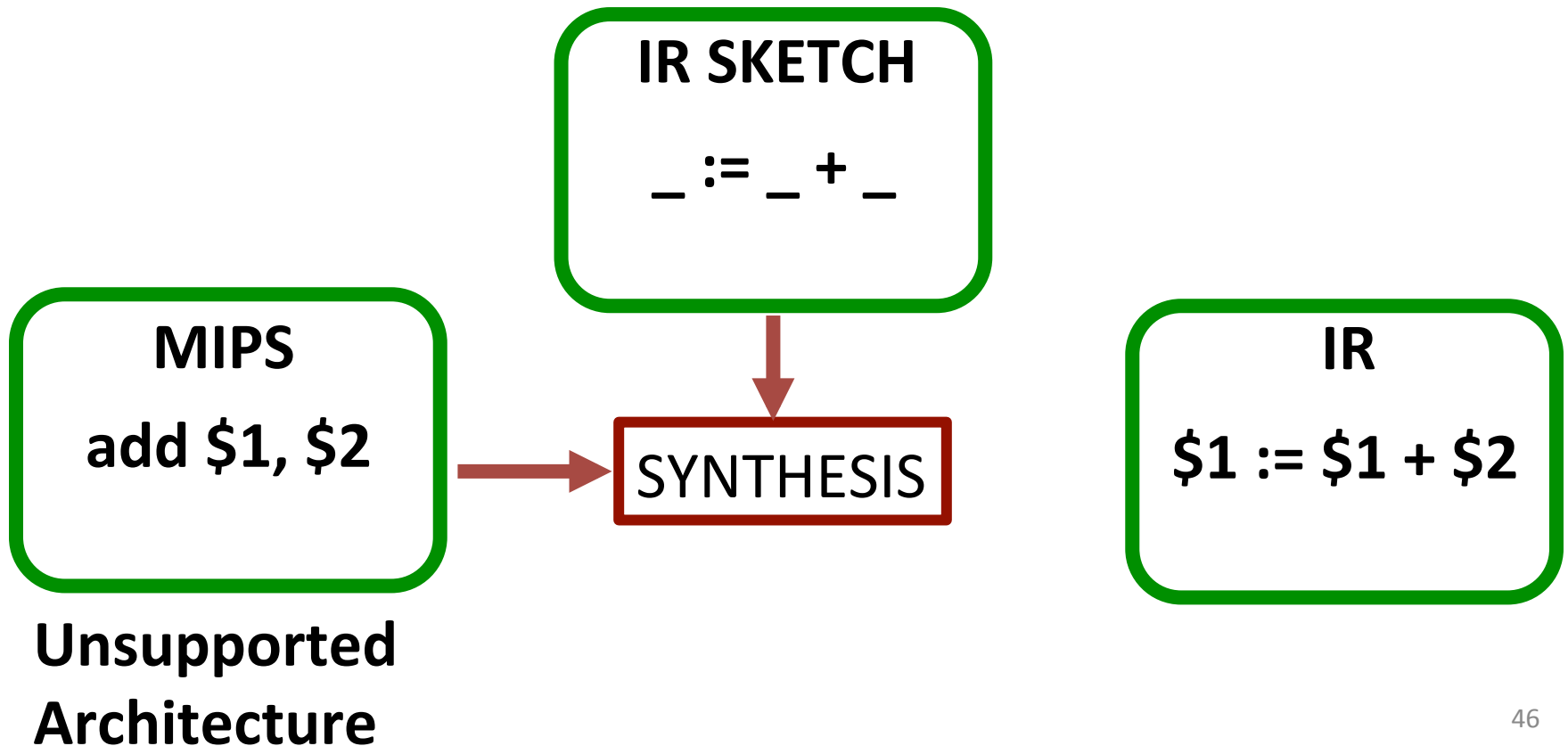
**IR SKETCH**

**\_ := \_ + \_**

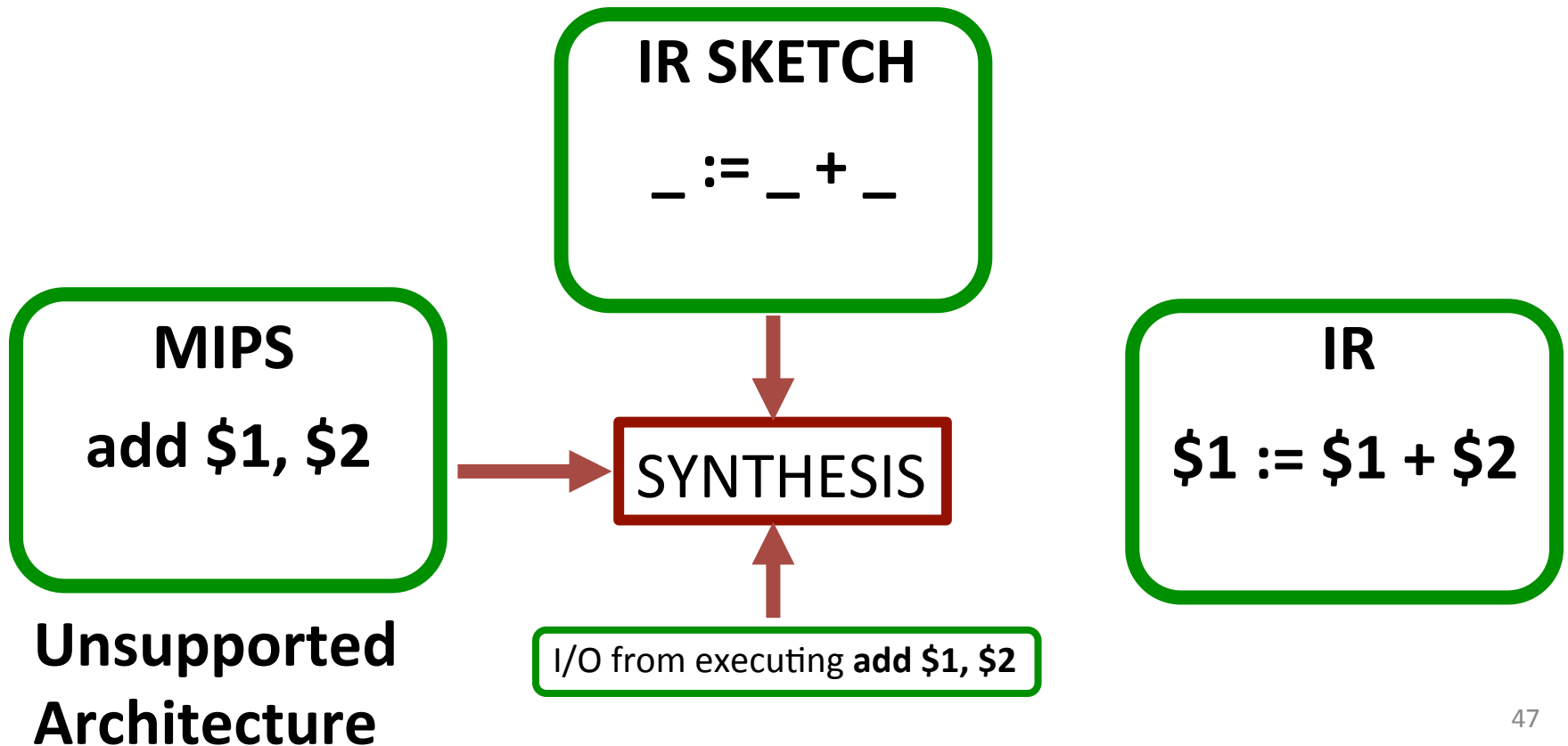
**IR**

**\$1 := \$1 + \$2**

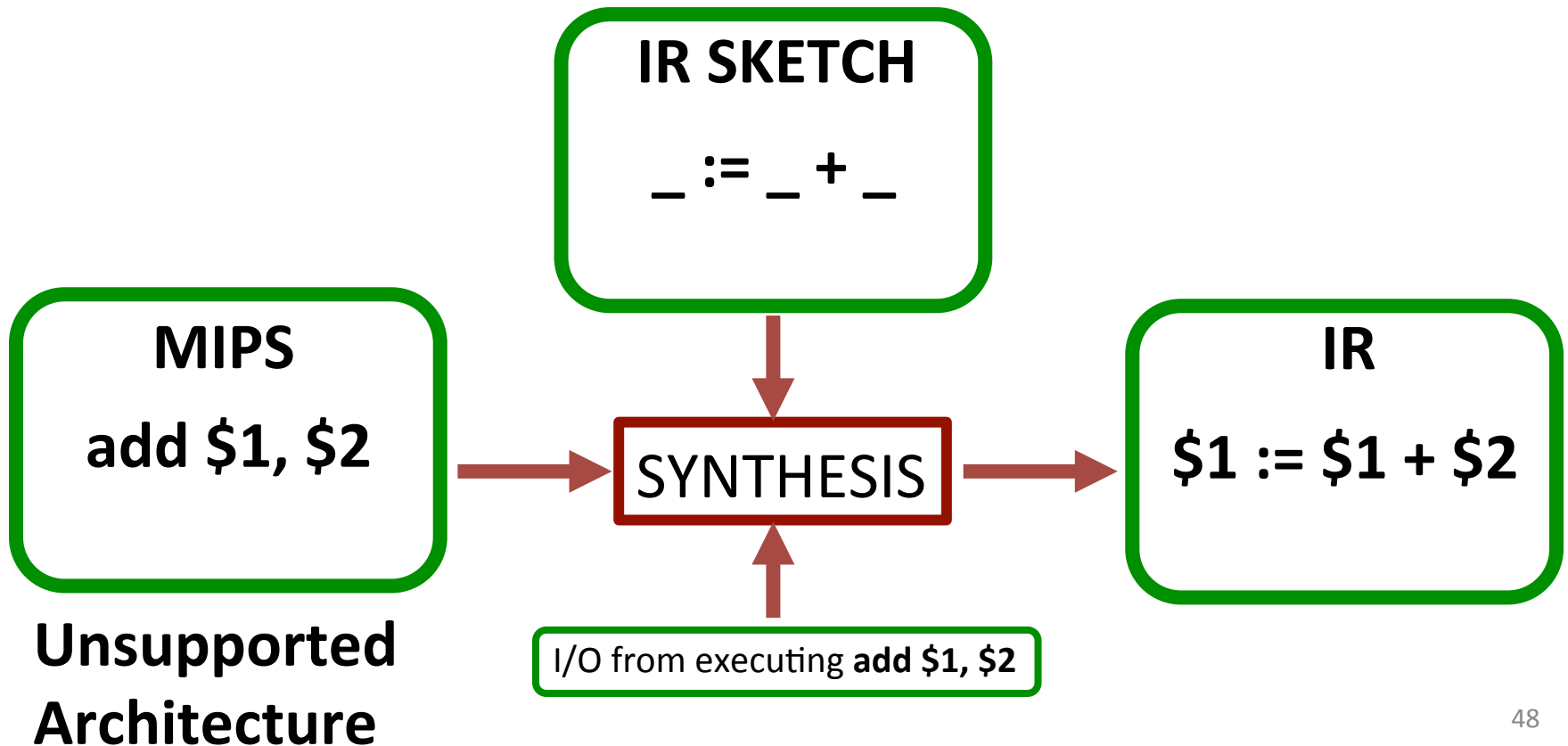
# Learning Sketch Templates: Example



# Learning Sketch Templates: Example

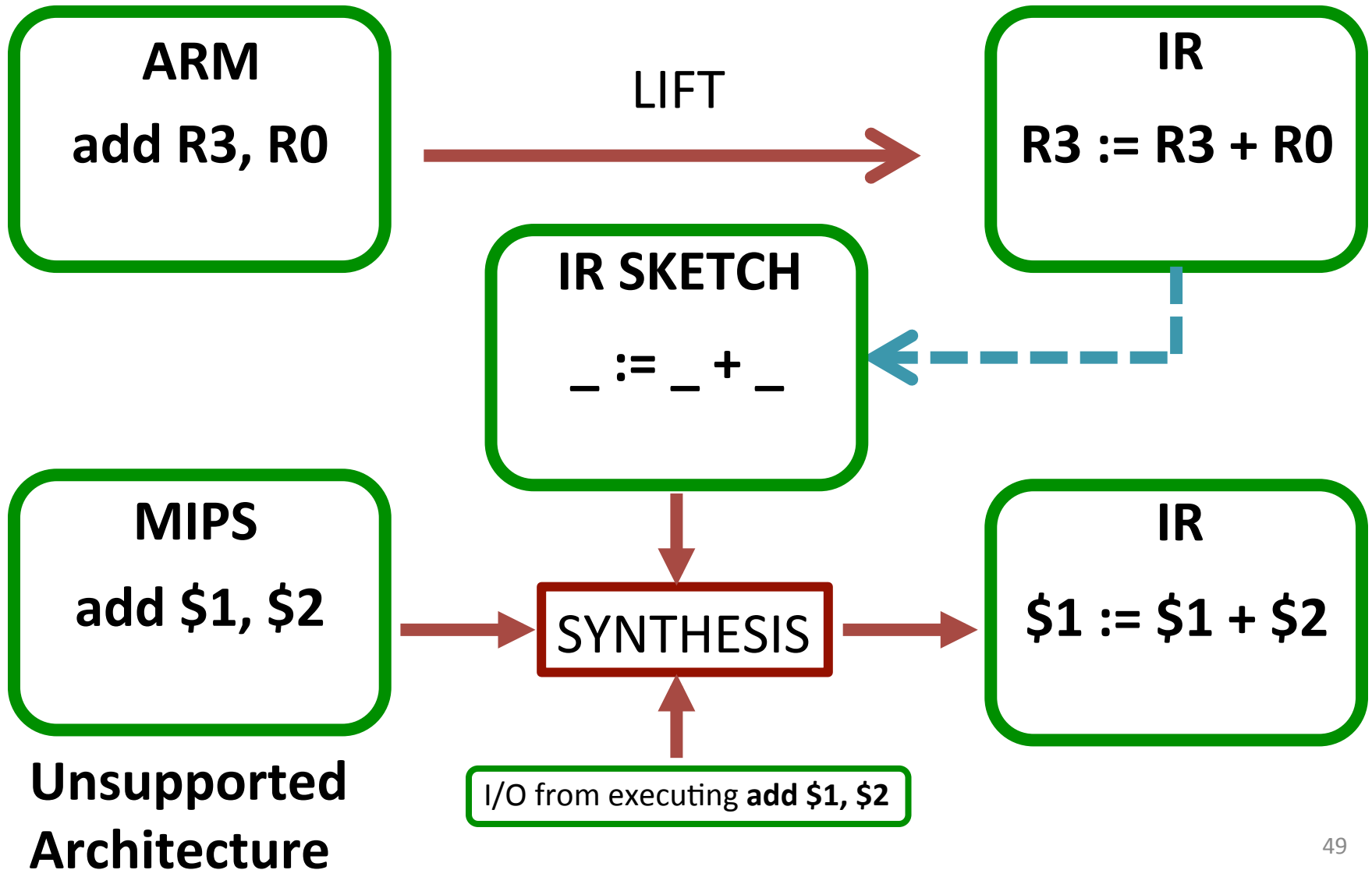


# Learning Sketch Templates: Example





# Learning Sketch Templates: Example



# Learning Sketch Templates: Example

- Learn templates from lifter output for a **supported** architecture to synthesize one for an **unsupported** architecture

# Learning Sketch Templates: Example

- Learn templates from lifter output for a **supported** architecture to synthesize one for an **unsupported** architecture
- Code is “natural”  $\longrightarrow$  shared semantic properties

# Learning Sketch Templates: Example

- Learn templates from lifter output for a **supported** architecture to synthesize one for an **unsupported** architecture
- Code is “natural”  $\longrightarrow$  shared semantic properties
- Sketch templates exploit structural qualities to guide synthesis

# Synthesis Approach

# Synthesis Approach

## Native Behavior

- I/O pairs from native execution

# Synthesis Approach

## Native Behavior

- I/O pairs from native execution
- QEMU and PIN traces

# Synthesis Approach

## Native Behavior

- I/O pairs from native execution
- QEMU and PIN traces

## IR Target

- Binary Analysis Platform (BAP) IR



# Synthesis Approach

## Native Behavior

- I/O pairs from native execution
- QEMU and PIN traces

## IR Target

- Binary Analysis Platform (BAP) IR
- IR Interpreter

# Synthesis Approach

## Native Behavior

- I/O pairs from native execution
- QEMU and PIN traces

$$\langle \sigma_{\mathbb{T}}, \mathcal{I}_{\mathbb{T}}, \emptyset \rangle \xrightarrow{\mathbb{T}} \langle \sigma'_{\mathbb{T}}, -, \mathcal{E}_{\mathbb{T}} \rangle$$

## IR Target

- Binary Analysis Platform (BAP) IR
- IR Interpreter

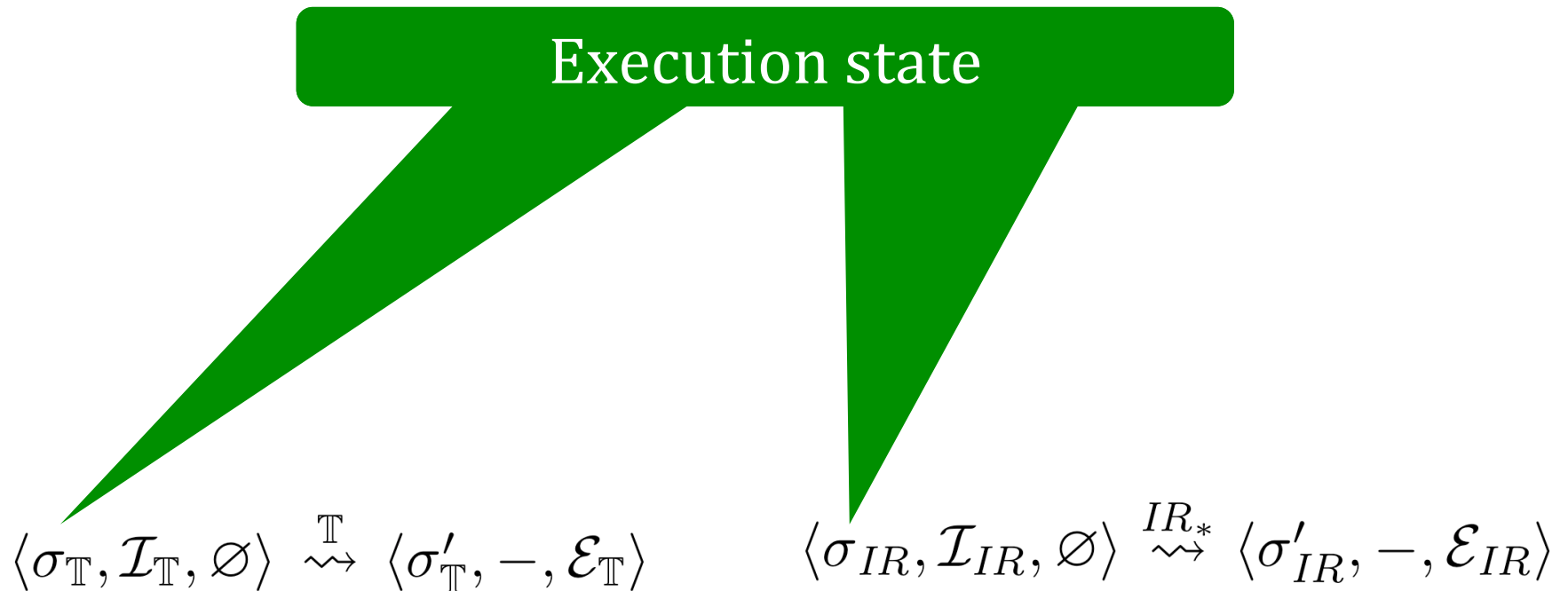
$$\langle \sigma_{IR}, \mathcal{I}_{IR}, \emptyset \rangle \xrightarrow{IR_*} \langle \sigma'_{IR}, -, \mathcal{E}_{IR} \rangle$$

# Synthesis Approach

$$\langle \sigma_{\mathbb{T}}, \mathcal{I}_{\mathbb{T}}, \emptyset \rangle \overset{\mathbb{T}}{\rightsquigarrow} \langle \sigma'_{\mathbb{T}}, -, \mathcal{E}_{\mathbb{T}} \rangle$$

$$\langle \sigma_{IR}, \mathcal{I}_{IR}, \emptyset \rangle \overset{IR_*}{\rightsquigarrow} \langle \sigma'_{IR}, -, \mathcal{E}_{IR} \rangle$$

# Synthesis Approach



# Synthesis Approach



Instruction

$$\langle \sigma_{\mathbb{T}}, \mathcal{I}_{\mathbb{T}}, \emptyset \rangle \overset{\mathbb{T}}{\rightsquigarrow} \langle \sigma'_{\mathbb{T}}, -, \mathcal{E}_{\mathbb{T}} \rangle$$


$$\langle \sigma_{IR}, \mathcal{I}_{IR}, \emptyset \rangle \overset{IR_*}{\rightsquigarrow} \langle \sigma'_{IR}, -, \mathcal{E}_{IR} \rangle$$

# Synthesis Approach

Events from execution

$$\langle \sigma_{\mathbb{T}}, \mathcal{I}_{\mathbb{T}}, \emptyset \rangle \xrightarrow{\mathbb{T}} \langle \sigma'_{\mathbb{T}}, -, \mathcal{E}_{\mathbb{T}} \rangle$$

$$\langle \sigma_{IR}, \mathcal{I}_{IR}, \emptyset \rangle \xrightarrow{IR_*} \langle \sigma'_{IR}, -, \mathcal{E}_{IR} \rangle$$

# Synthesis Approach

$$R2 := R0$$

$$[(R, \text{REG}, R0, 0x1), (W, \text{REG}, R2, 0x1)]$$

Events from execution

$$\langle \sigma_{\mathbb{T}}, \mathcal{I}_{\mathbb{T}}, \emptyset \rangle \xrightarrow{\mathbb{T}} \langle \sigma'_{\mathbb{T}}, -, \mathcal{E}_{\mathbb{T}} \rangle$$

$$\langle \sigma_{IR}, \mathcal{I}_{IR}, \emptyset \rangle \xrightarrow{IR_*} \langle \sigma'_{IR}, -, \mathcal{E}_{IR} \rangle$$

# Inferring Semantics from I/O pairs

**SOURCE**

**Native Instruction**



**add R3, R0**



# Inferring Semantics from I/O pairs

## SOURCE

### Native Instruction

**add R3, R0**

## Execution

### Input

R3 = 2

R0 = 2

### Output

R3 = 4

R0 = 2

# Inferring Semantics from I/O pairs

## SOURCE

### Native Instruction

**add R3, R0**

## Execution

### Input

R3 = 2

R0 = 2

### Output

**R3 = 4**

R0 = 2

# Synthesizing Translation

## TARGET

IR Instruction

??

## Execution

Input

R3 = 2

R0 = 2

Output

R3 = 4

R0 = 2

# Synthesizing Translation

## TARGET

IR Instruction

**R3, R0, R3**

## Execution

Input

R3 = 2

R0 = 2

Output

R3 = 4

R0 = 2

# Syntax-guided Synthesis

TARGET

IR Instruction

R3, R0, R3

\_ := \_ + \_

# Syntax-guided Synthesis

TARGET

IR Instruction

R3, R0, R3

\_ := \_ + \_

## Sketch Templates

\_ := \_ + \_

\_ := \_ - \_

\_ := \_ \* \_

\_ := \_ << \_

...



R3 := R0 + R3

R3 := R0 - R3

R3 := R0 \* R3

R3 := R0 << R3

...

# Syntax-guided Synthesis

TARGET

IR Instruction

R3, R0, R3

\_ := \_ + \_

## Sketch Templates

\_ := \_ + \_

\_ := \_ - \_

\_ := \_ \* \_

\_ := \_ << \_

...



R3 := R0 + R3

R3 := R0 - R3

R3 := R0 \* R3

R3 := R0 << R3

...

Exhaustive Enumeration

# Syntax-guided Synthesis

TARGET

IR Instruction

R3, R0, R3

\_ := \_ + \_

## Sketch Templates

\_ := \_ + \_

\_ := \_ - \_

\_ := \_ \* \_

\_ := \_ << \_

...



R3 := R0 + R3

R3 := R0 - R3

R3 := R0 \* R3

R3 := R0 << R3

...

Exhaustive Enumeration

Permute adjacent operands



# Syntax-guided Synthesis

TARGET

IR Instruction

**R3 := R0 + R3**

IR Interpreter

Input

R3 = 2

R0 = 2

Output

R3 = 4

R0 = 2

# Syntax-guided Synthesis

## Native Execution

Input

R3 = 2  
R0 = 2

Expected Output

R3 = 4  
R0 = 2



## IR Interpreter

Input

R3 = 2  
R0 = 2

Output

R3 = 4  
R0 = 2

# Validating Correctness

## TARGET

IR Instruction

**R3 := R0 \* R3**

## IR Interpreter

Input

R3 = 2

R0 = 2

Expected Output

R3 = 4

R0 = 2

# Validating Correctness

## TARGET

IR Instruction

**R3 := R0 \* R3**

## IR Interpreter

Input

R3 = 2

R0 = 2

Expected Output

R3 = 4

R0 = 2

**Oops...**

# Validating Correctness

TARGET

IR Instruction

**R3 := R0 \* R3**

**Oops...**

# Validating Correctness

## TARGET

IR Instruction

**R3 := R0 \* R3**

## Native Execution

Input

**R3 = 3**

R0 = 2

Expected Output

**R3 = 5**

R0 = 2

# Validating Correctness

## TARGET

IR Instruction

**R3 := R0 \* R3**

**R3 = 6**

## Native Execution

Input

**R3 = 3**

R0 = 2

Expected Output

**R3 = 5**

R0 = 2

# Validating Correctness

## TARGET

IR Instruction

$R3 := R0 * R3$

$R3 = 6$



Invalidate  
sketch

## Native Execution

Input

$R3 = 3$

$R0 = 2$

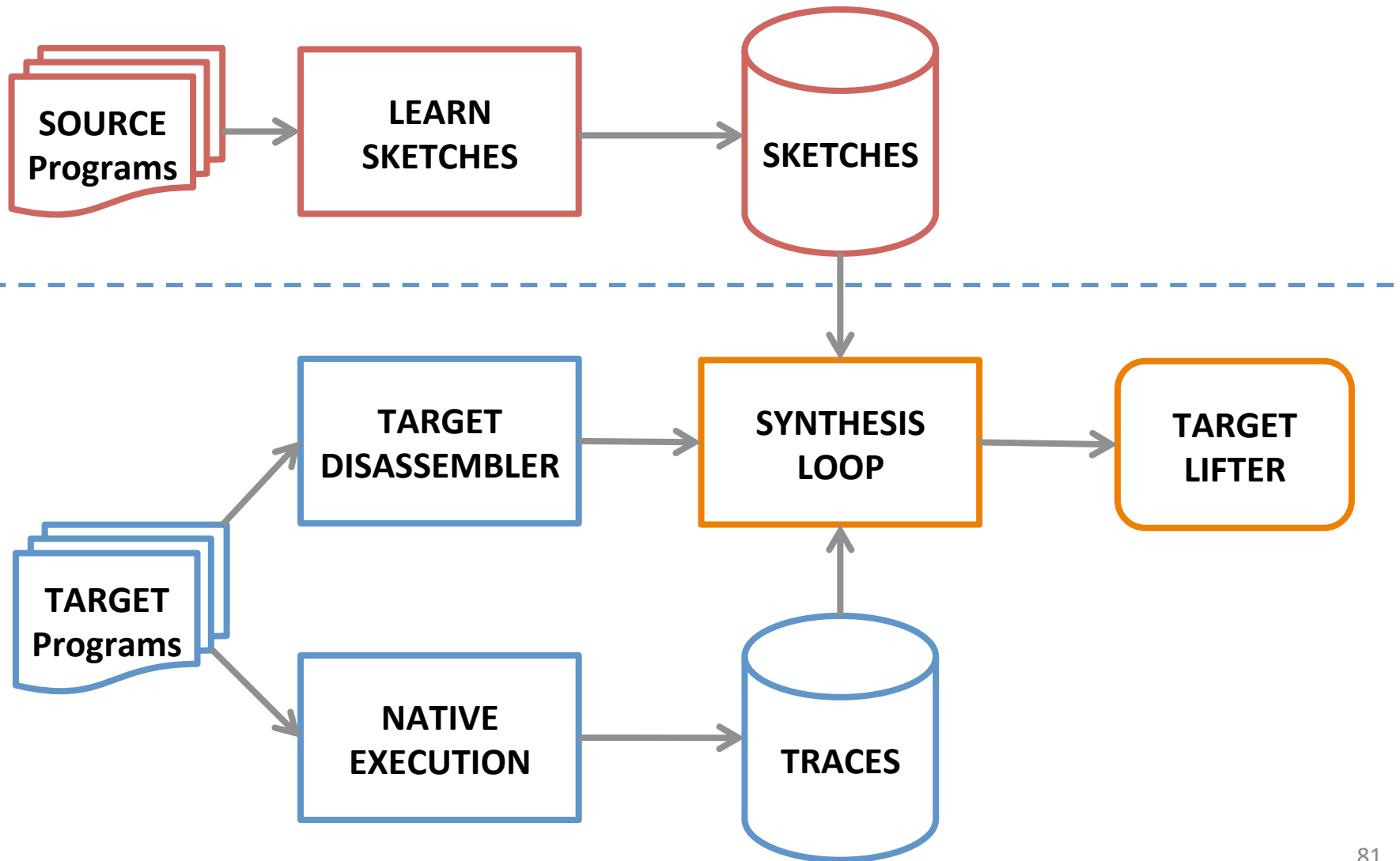
Expected Output

$R3 = 5$

$R0 = 2$



# Lifter Synthesis System



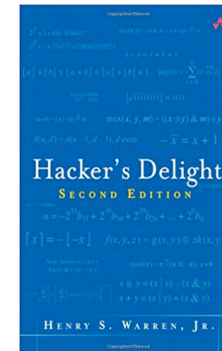
# Experimental Setup

# Experimental Setup

- Synthesize unsupported lifter target (MIPS)

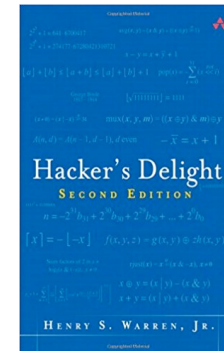
# Experimental Setup

- Synthesize unsupported lifter target (MIPS)
- Collect dynamic traces from arithmetic-heavy benchmark
  - 5 Hacker's delight programs



# Experimental Setup

- Synthesize unsupported lifter target (MIPS)
- Collect dynamic traces from arithmetic-heavy benchmark
  - 5 Hacker's delight programs
- Mine IR sketches from Coreutils, lifted from ARM and x86



# Results: Analysis Reuse

# Results: Analysis Reuse

- Taint analysis for warn-unused-result bugs

# Results: Analysis Reuse

- Taint analysis for warn-unused-result bugs
  - Forget to check memory is allocated
  - Forget to check privileges are dropped



# Results: Analysis Reuse

- Taint analysis for warn-unused-result bugs
  - Forget to check memory is allocated
  - Forget to check privileges are dropped
- Ran on 30 binaries of COTS D-Link router

# Results: Analysis Reuse

- Taint analysis for warn-unused-result bugs
  - Forget to check memory is allocated
  - Forget to check privileges are dropped
- Ran on 30 binaries of COTS D-Link router
- 29 bugs, 2 false positives

# Results: Analysis Reuse

Name	#	Functions
pppd	1	fwrite
iptables	3	fwrite
rdnssd	1	setsockopt
ntpclient	1	send
speedtest	2	system, fgets
timer	2	read, shutdown
wakeOnLanProxy	1	shutdown
wcnd	8	system

# Results: Analysis Reuse

Name	#	Functions
pppd	1	fwrite
iptables	3	fwrite
rdnssd	1	setsockopt
ntpclient	1	send
speedtest	2	system, fgets
timer	2	read, shutdown
wakeOnLanProxy	1	shutdown
wcnd	8	system

```
uint32_t data[12];
struct timeval now;

memset(data, 0, sizeof(data));
data[0] = htonl (
    ( LI << 30 ) | ( VN << 27 ) | (
    ( STRATUM << 16) | ( POLL << 8 )
);
data[1] = htonl(1<<16); /* Root Delay
data[2] = htonl(1<<16); /* Root Disper
gettimeofday(&now,NULL);
data[10] = htonl(now.tv_sec + JAN_1970);
data[11] = htonl(NTPFRAC(now.tv_usec));
send(sd,data,48,0);
```

# Results: Analysis Reuse

Name	#	Functions
pppd	1	fwrite
iptables	3	fwrite
rdnssd	1	setsockopt
ntpclient	1	send
speedtest	2	system, fgets
timer	2	read, shutdown
wakeOnLanProxy	1	shutdown
wcnd	8	system

```
uint32_t data[12];
struct timeval now;

memset(data, 0, sizeof(data));
data[0] = htonl (
    ( LI << 30 ) | ( VN << 27 ) | (
    ( STRATUM << 16) | ( POLL << 8 )
);
data[1] = htonl(1<<16); /* Root Delay
data[2] = htonl(1<<16); /* Root Disper
gettimeofday(&now, NULL);
data[10] = htonl(now.tv_sec + JAN_1970);
data[11] = htonl(NTPFRAC(now.tv_usec));
send(sd, data, 48, 0);
```

No false positives on OpenSSL

# Results: Synthesis

# Results: Synthesis

- End-to-end synthesis takes 58 seconds
  - mining sketches, processing traces, and lifter synthesis
- 29 sketches per instruction, on average
  - Synthesis converges after 2 input output pairs, on average
- Lifts 84.8% of native instructions.
  - Missed example: “load upper immediate”

# Generalizing across Architectures



# Generalizing across Architectures

- IR Sketches mined from ARM and x86 are common to both

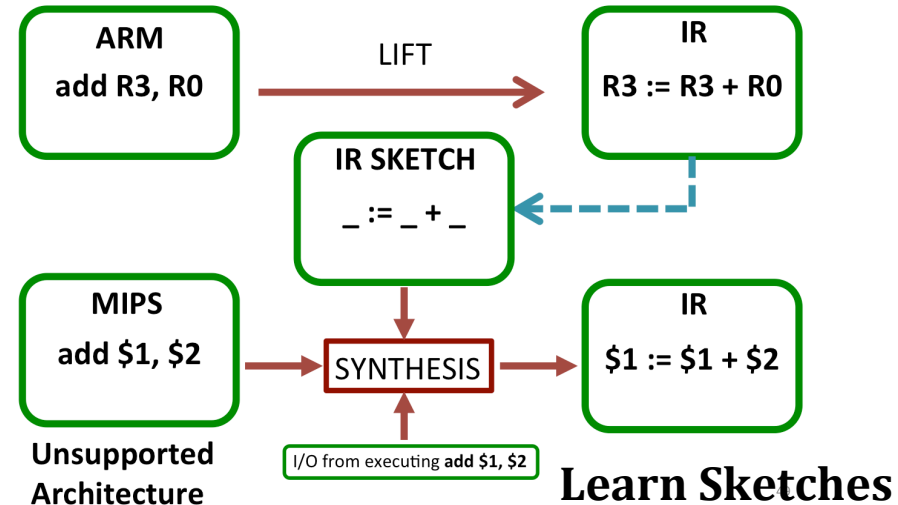
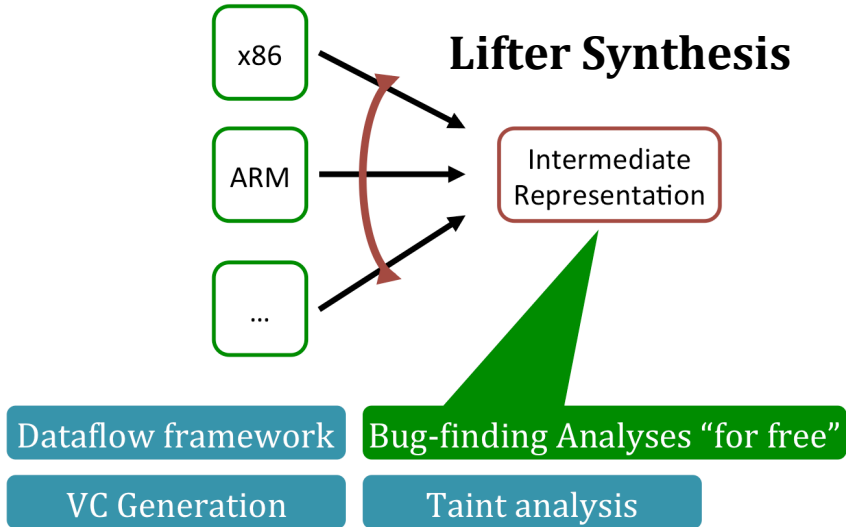
%	ARM-IR	%	X86-IR
20.9	<code>_v := _v</code>	11.9	<code>_v := _i</code>
14.5	<code>_v := _i</code>	8.6	<code>jmp _i</code>
8.9	<code>jmp _i</code>	5.1	<code>_v := mem[_v + _i]</code>
7.0	<code>_v := mem[_v + _i]</code>	4.4	<code>mem[_v + _i] := _v</code>
5.6	<code>_v := _v = _i</code>	4.2	<code>_v := _v = _i</code>
5.6	<code>_v := hi : 1[_v]</code>	4.2	<code>_v := hi : 1[_v]</code>
5.5	<code>mem[_v + _i] := _v</code>	4.0	<code>mem[_v + _i] := _v</code>
4.6	<code>_v := _v - _i</code>	3.9	<code>_v := _v - _i</code>

- Similar MIPS accuracy from either ARM or x86 sketches

# Discussion

- How good “out-of-box”?
- Improve recovery with nondeterministic sketch search and generation
- Expand to more architectures
  - SPARC, PPC, ...
- Complement automatic synthesis with (reduced) manual effort
  - Verify and fix manual translation

# Summary



## Verify with Dynamic Traces

$R2 := R0$

$[(R, \text{REG}, R0, 0x1), (W, \text{REG}, R2, 0x1)]$

Events from execution

$$\langle \sigma_T, \mathcal{I}_T, \emptyset \rangle \xrightarrow{T} \langle \sigma'_T, -, \mathcal{E}_T \rangle \quad \langle \sigma_{IR}, \mathcal{I}_{IR}, \emptyset \rangle \xrightarrow{IR_*} \langle \sigma'_{IR}, -, \mathcal{E}_{IR} \rangle$$

Name	#	Functions
pppd	1	fwrite
iptables	3	fwrite
rdnssd	1	setsockopt
ntplib	1	send
speedtest	2	system, fgets
timer	2	read, shutdown
wakeOnLanProxy	1	shutdown
wcnd	8	system

```
uint32_t data[12];
struct timeval now;

memset(data, 0, sizeof(data));
data[0] = htonl (
    ( LI << 30 ) | ( VN << 27 ) | (
        ( STRATUM << 16 ) | ( POLL << 8 )
    );
data[1] = htonl(1<<16); /* Root Delay
data[2] = htonl(1<<16); /* Root Disper:
gettimeofday(&now, NULL);
data[10] = htonl(now.tv_sec + JAN_1970);
data[11] = htonl(NTPFRAC(now.tv_usec));
send(sd, data, 48, 0);
```

## Analysis Reuse for Previously Unsupported Instruction Set