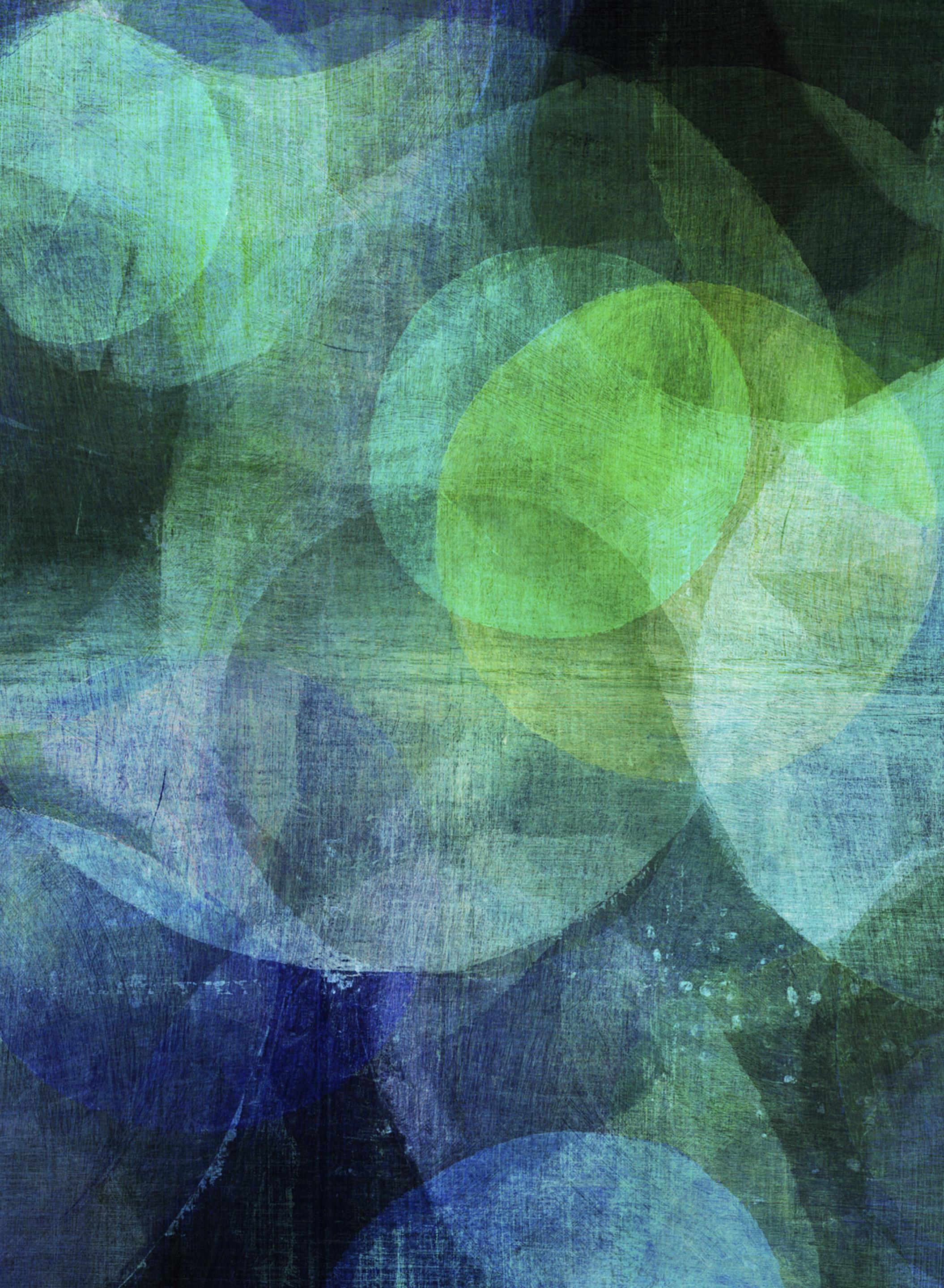


IMPROVING DYNAMIC CODE ANALYSIS BY CODE ABSTRACTION

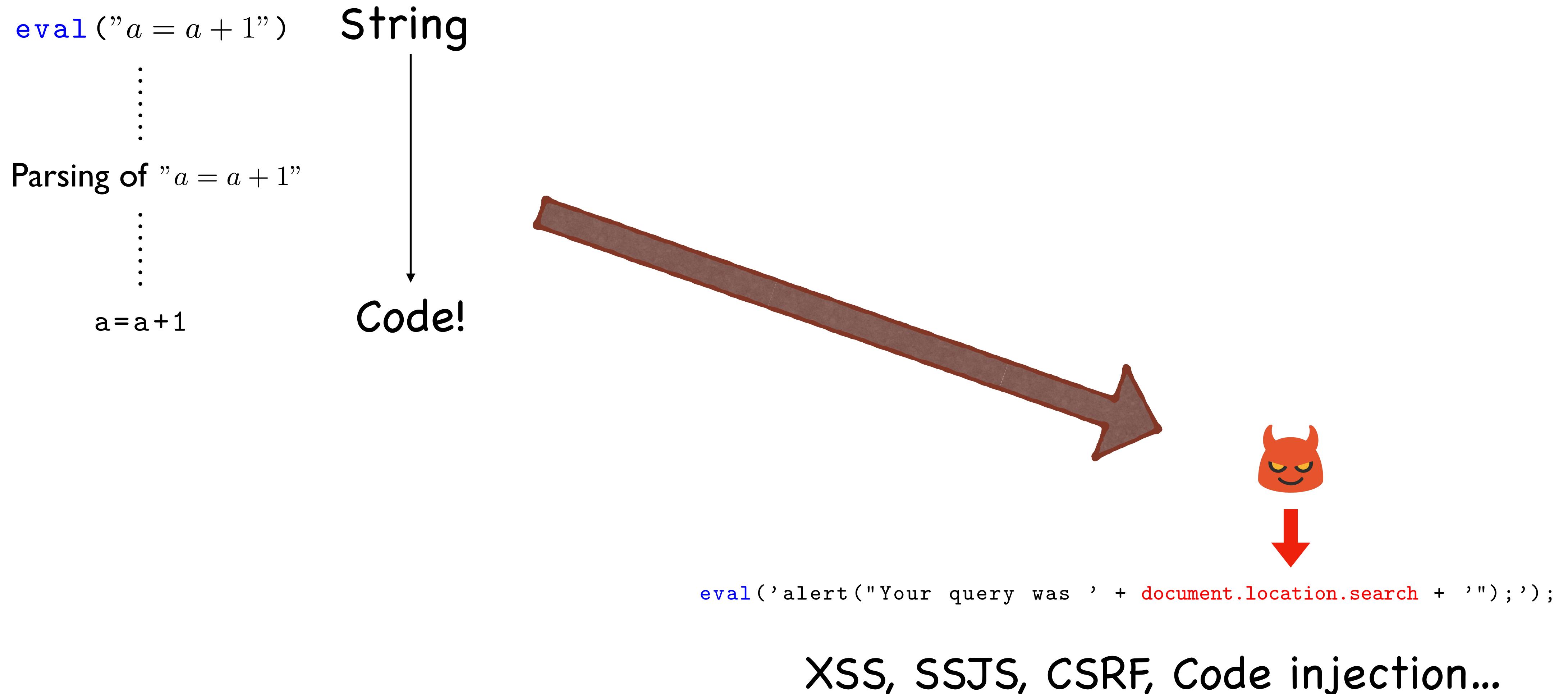
I. Mastroeni and V. Arceri
Univ. of Verona and Univ. of Venezia

MOTIVATION

Dynamic languages do have
string-to-code statements



JAVASCRIPT HAS EVAL ALLOWING EVIL THINGS



HOW MUCH IS EVAL POPULAR?

Best practice:

“eval is the most misused feature of JavaScript. Avoid it.”

Benevolent code

G. Richards et al. “The eval that men do”

Malicious code

We have analyzed a data-set of JavaScript malware samples:

57% have at least
one eval call

10% have nested
eval calls

THE STATIC ANALYSIS PROBLEM

- Static analysis assumption: The program code is static
- Eval brakes this assumption!

```
        while (x++ < 3)
            os = os + "x=x+1; y=10; x=x+1;" ;

x = 0;
while (x >= 0) {
    if (x % 2 == 0)
        x = x + 1;
    else
        x = x + 2;
}

if (x > 10)
    os = "while(x>5){x=x+1; y=x;}";
if (x == 5)
    os = "hello";
if (x == 8)
    os = "while(x;";

eval(os);
```

x = x + 1; y = 10; x = x + 1;

Derived at run-time
May contain other eval calls

STATE OF THE ART IN ANALYZING DYNAMIC CODE

- Before introducing the sound static analyzer:
 - Several static analyzers have been proposed for dynamic languages (PHP, Python, Ruby...)
 - String analyses fails to precisely analyze string manipulation programs
 - Almost none of them faced the problem of analyzing string-to-code statements
 - They ignore the eval effects (inducing unsoundness) or forbid eval usage
- A sound static analyzer for dynamic code (TOPS'21): Ingredients
 - A string abstract domain keeping enough information for extracting code
 - Design and implementation of an algorithm for extracting executable code from string abstractions
 - IDEA: Analyze code, when an eval is met, extract executable code and recursively call the static analyzer on this code!

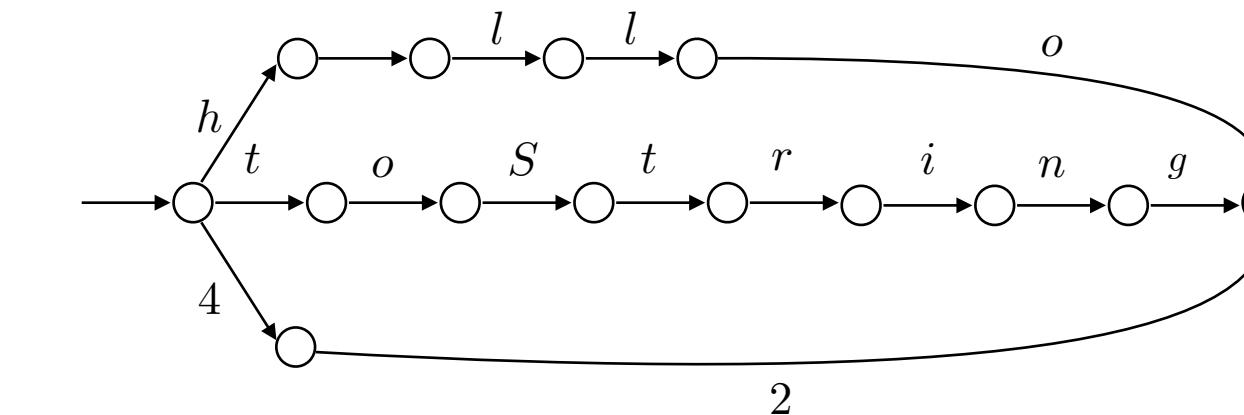
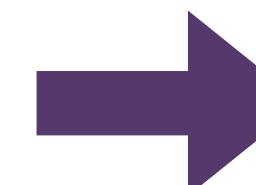
THE EXISTING SOLUTION

A sound abstract interpreter
for dynamic code



THE STRING ABSTRACTION: FINITE STATE AUTOMATA

```
{ "42",  
  "toString",  
  "hello"}
```



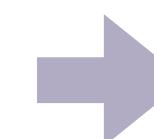
The guard value is
statically unknown



```
if (...)  
  x = "a++;"
```

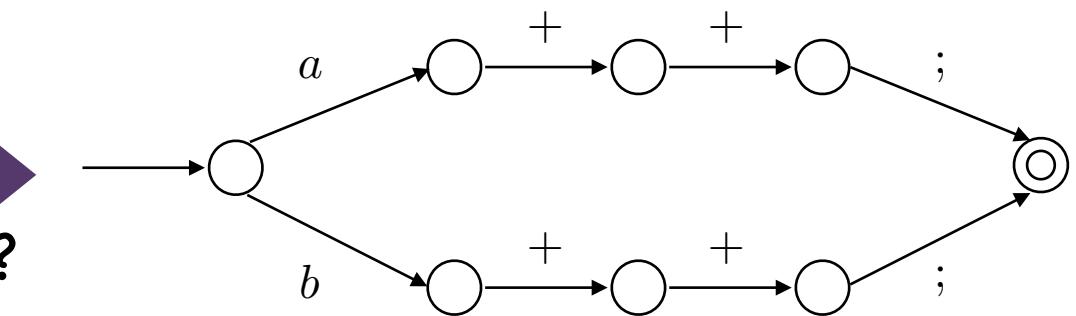
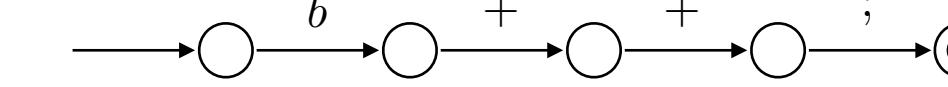
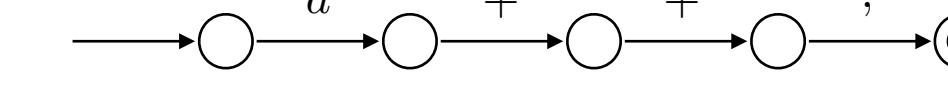


```
else  
  x = "b++;"
```



```
eval(x);
```

Which is the value of x before the eval execution?



THE STRING ABSTRACTION: FINITE STATE AUTOMATA

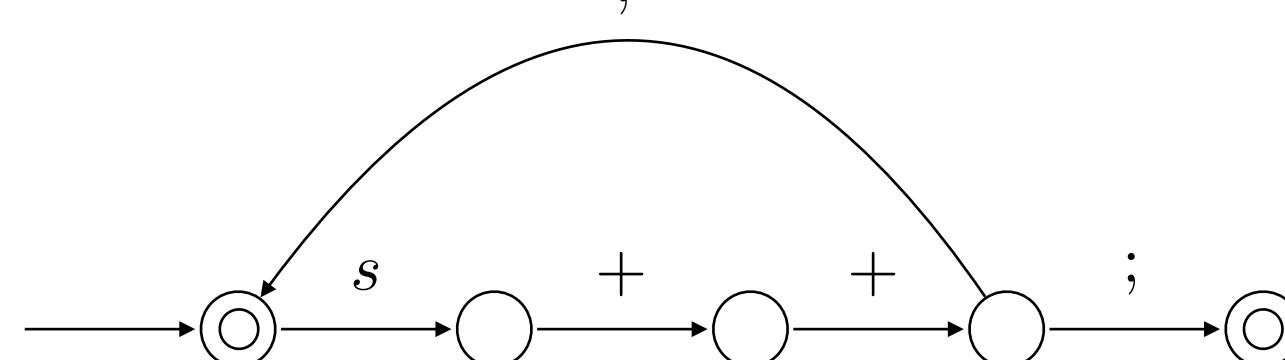
The guard value is statically unknown

→ `while (...) x = x + "s++;"` → "s++; s++; s++; s++; ..."

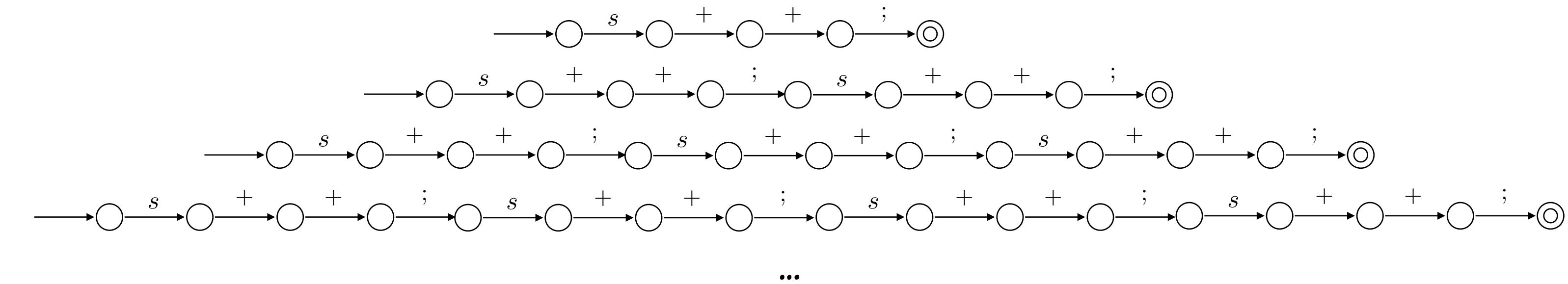
`eval(x)`

Which is the value of `x` before the `eval` execution?

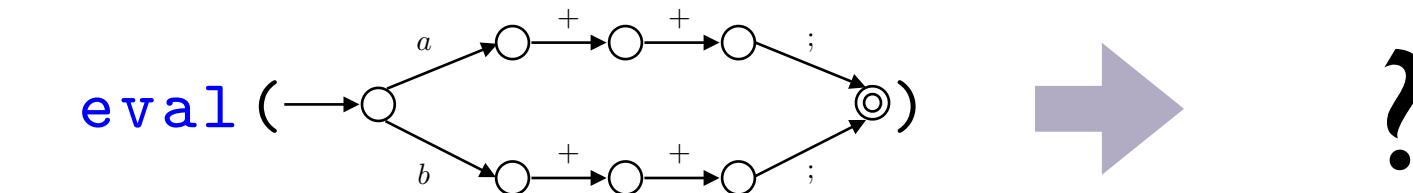
We want to ensure:
• Soundness
• Convergence



Widening



STATIC ANALYSIS OF EVAL



- How do we extract executable code from automata? 1) We abstract the argument!

```
while (x++ < 3)
    os = os + "x=x+1; y=10; x=x+1; ";

if (x > 10)
    os = "while(x>5){x=x+1; y=x;}";
```

```
if (x == 5)
    os = "hello";

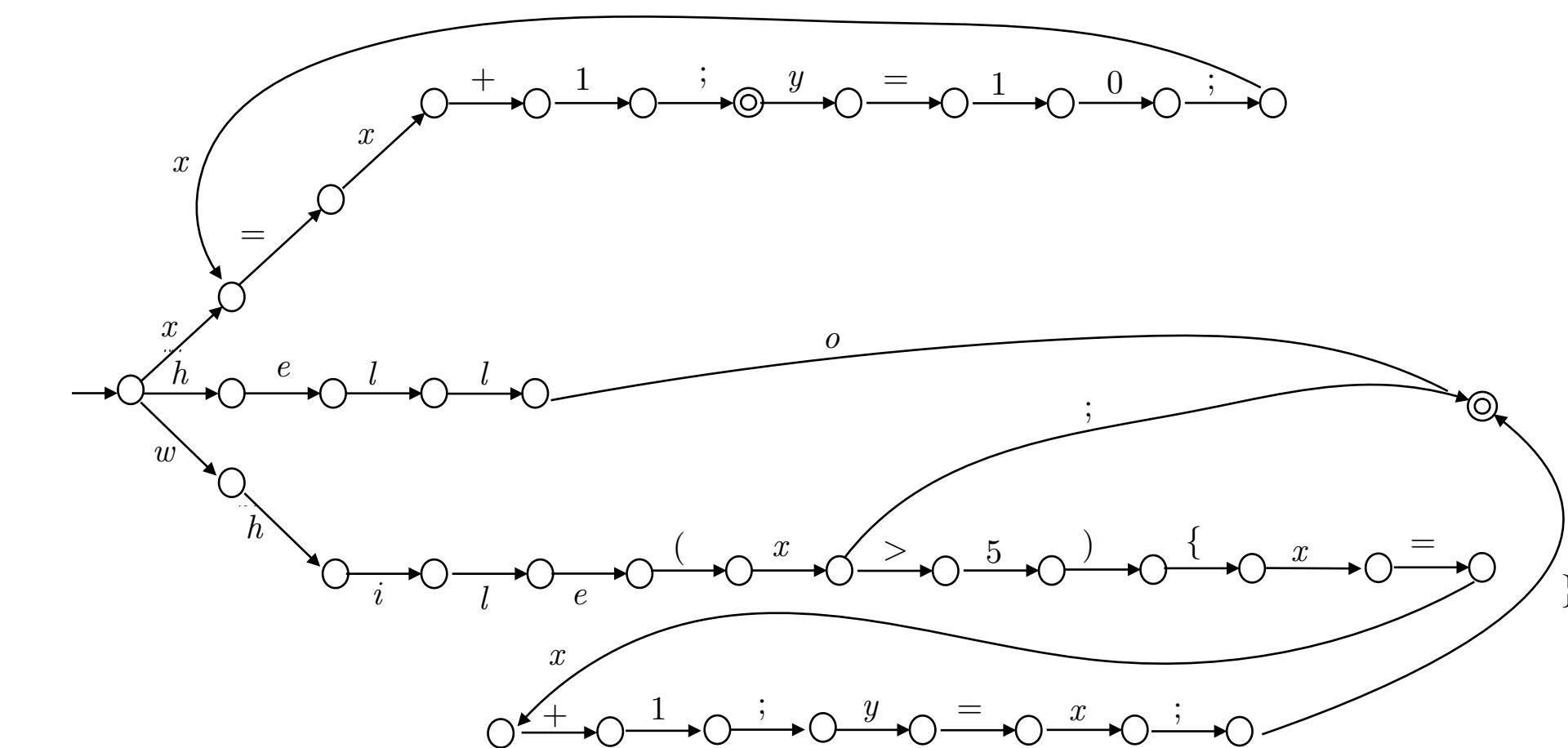
if (x == 8)
    os = "while(x;>";
```

OS

```
eval(os);
```

The value of x is statically unknown!

Not executable

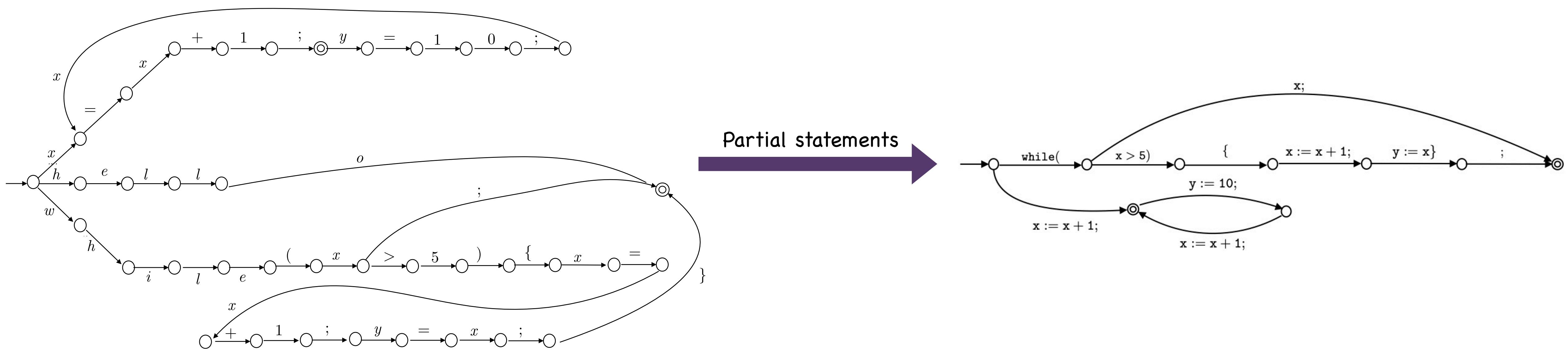


FROM STRINGS TO SEQUENCES OF PARTIAL STATEMENTS

- 2) Automaton of characters \rightarrow Automaton of partial statements

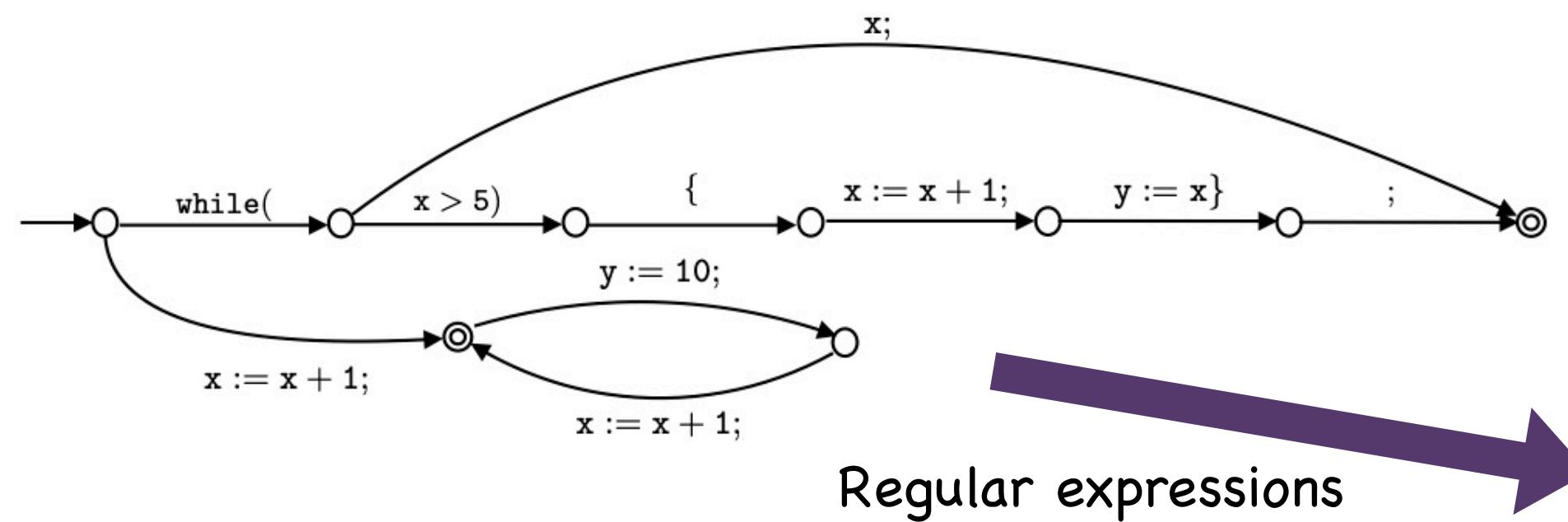
A partial statement is a prefix statement ending with a punctuation symbol or a punctuation symbol

$x=x+1;$ $\text{while}($ $x>5)$ $\{$ $\}$

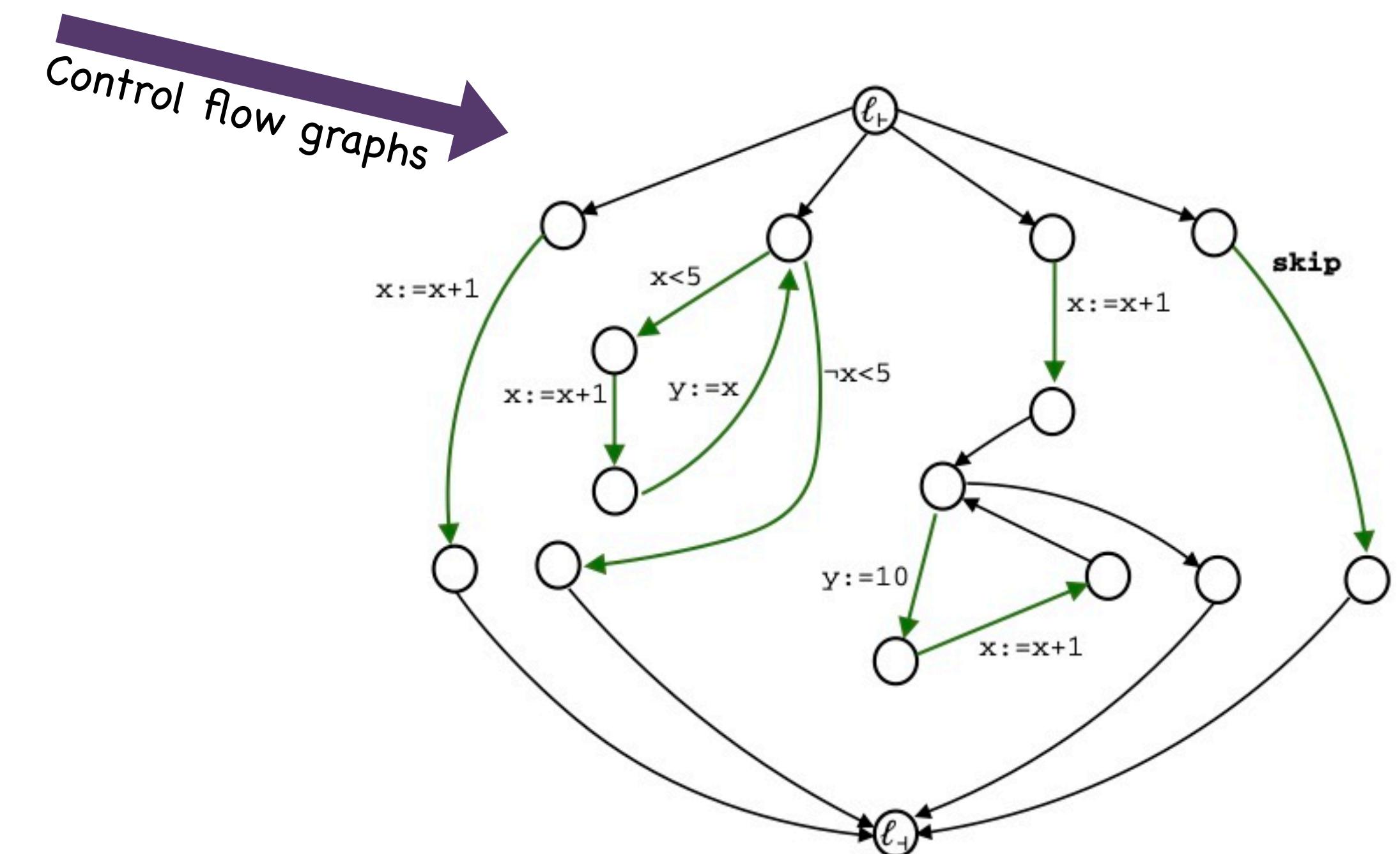


FROM PARTIAL STATEMENTS TO CODE

- 3) Automaton of partial statements \Rightarrow Regular Expressions \Rightarrow CFG

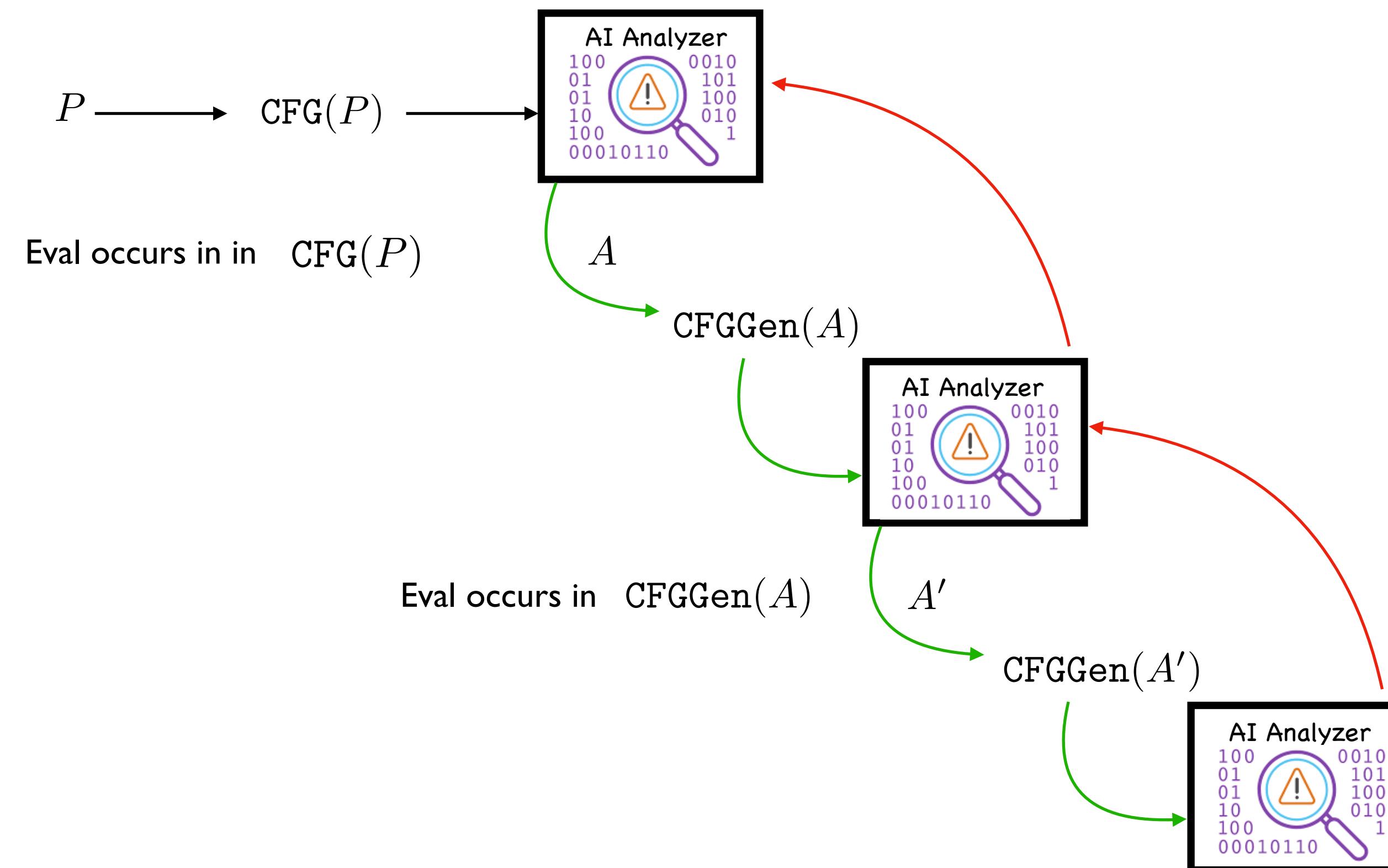


$x := x + 1; \parallel \text{while}(x > 5)\{x := x + 1; y := x\}; \parallel x := x + 1; (y := 10; x := x + 1;)^* \parallel \text{while}(x;$



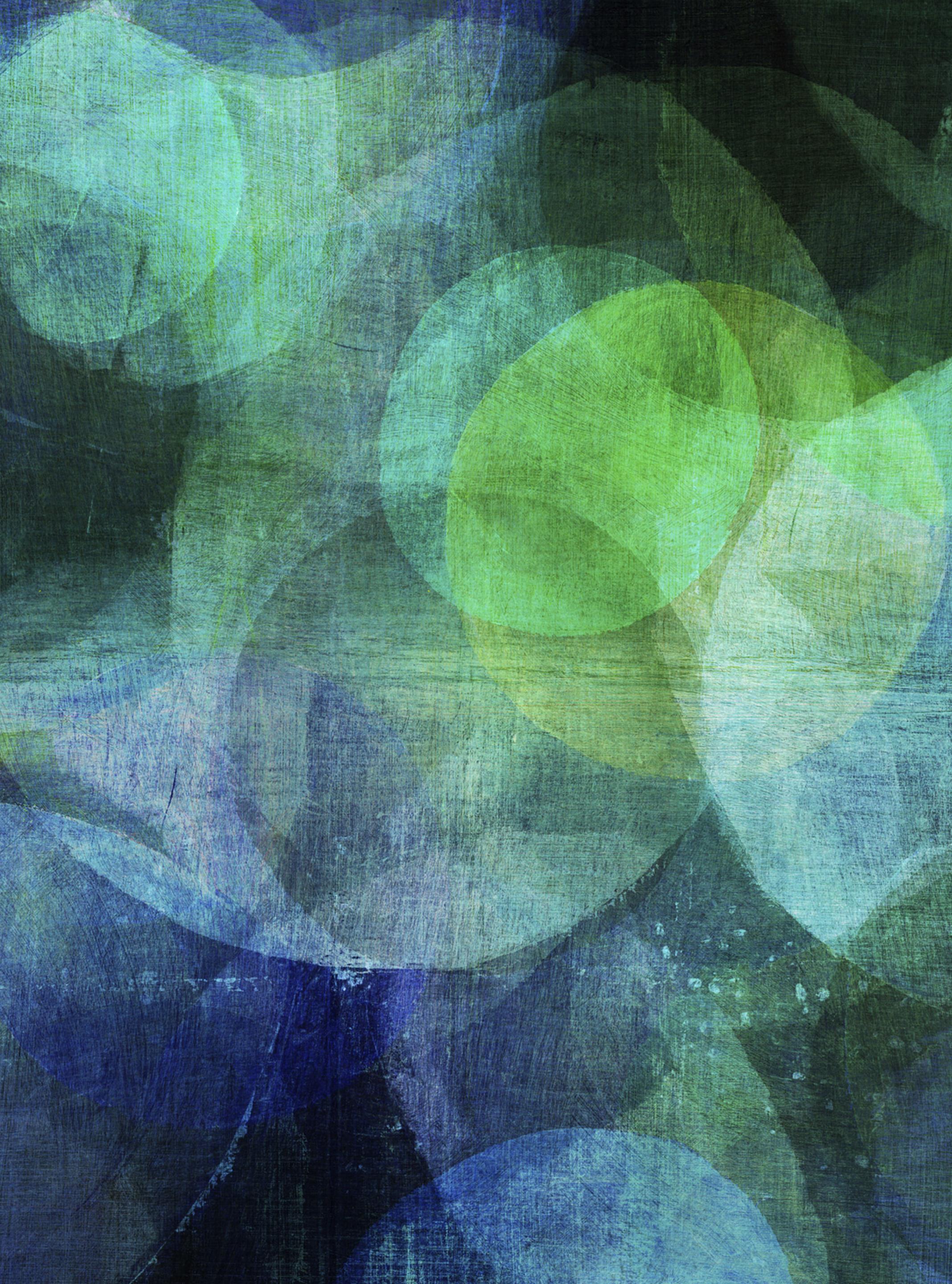
ANALYZING DYNAMICALLY GENERATED CODE

- 4) Recursively call the analyzer on the generated CFG!



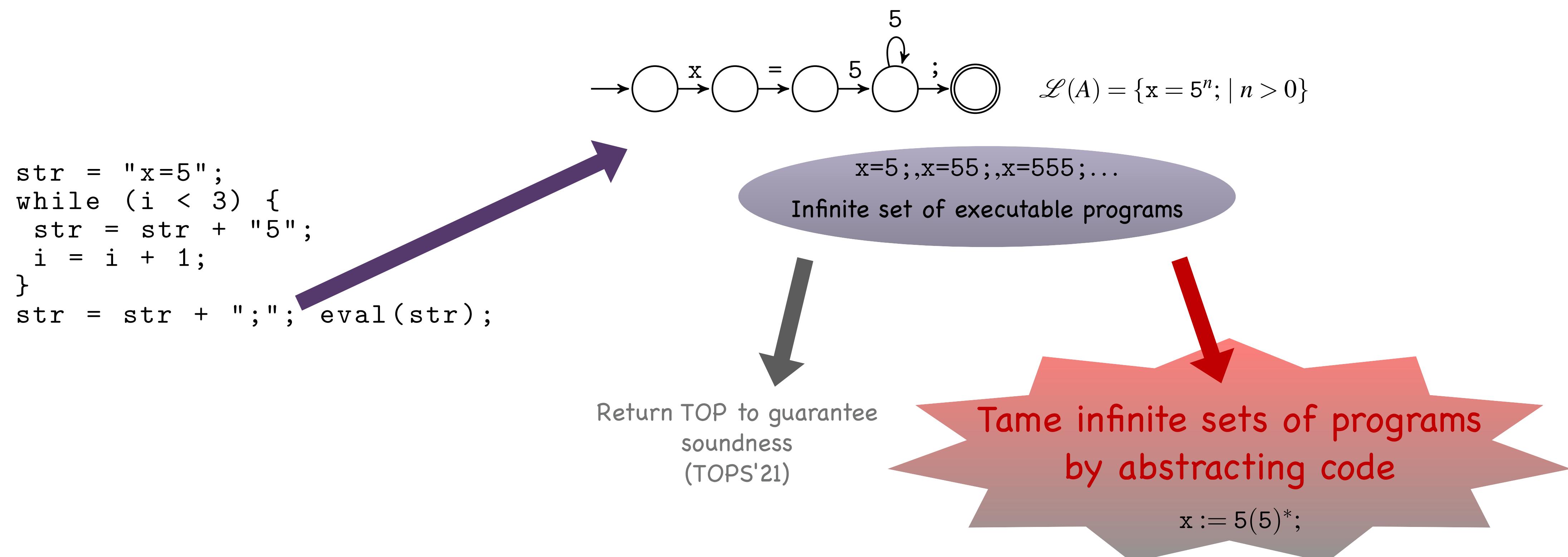
THE ANALYSIS LIMITS

Soundness vs Imprecision



LIMITATIONS... A SOLUTION

- In order to guarantee soundness we lose precision!
- When the executable code contained in the automaton is **infinite** we have a problem!





HOW CAN WE IMPROVE PRECISION?

By performing the analyzer on **abstract** code!

SEMANTIC- DRIVEN CODE ABSTRACTION

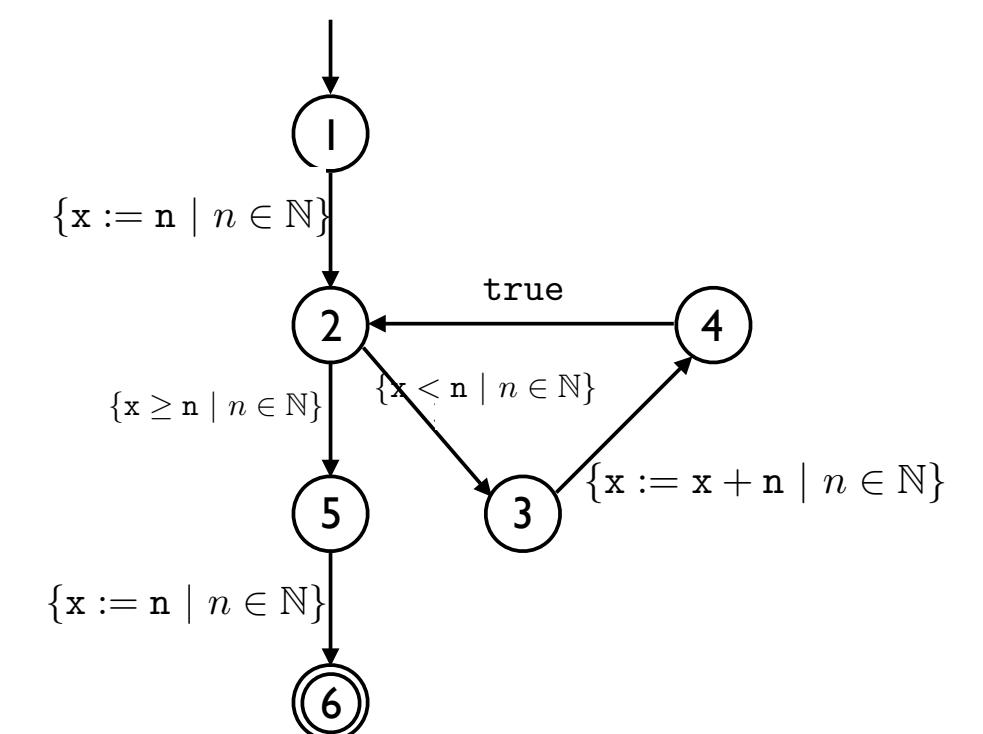
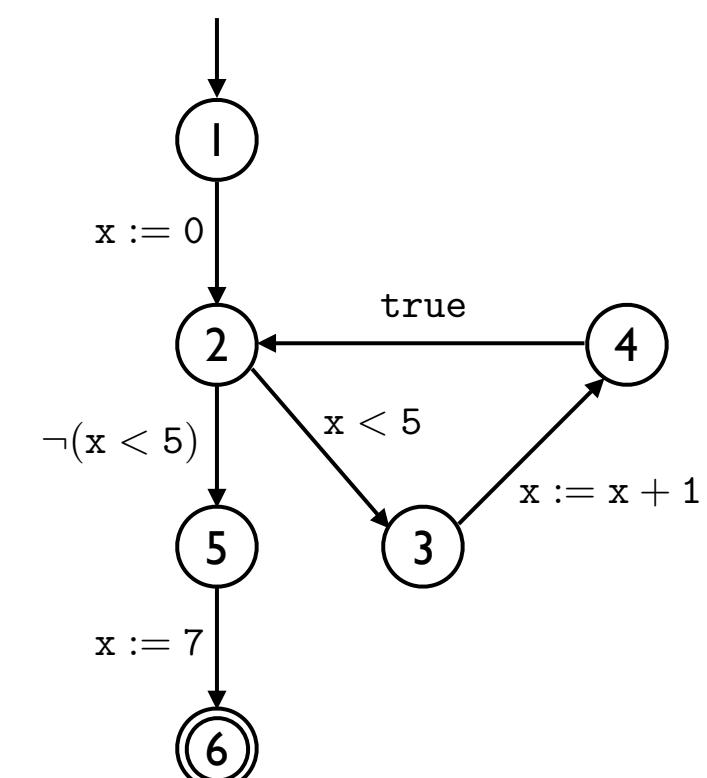
Taming "infinite" programs



CODE ABSTRACTION

- We model code as control flow graphs (CFG)
- We leave unchanged the control structure and we abstract the effects on the state: **labels as sets of statements**

```
 $\ell_1$  x := 0;  
 $\ell_2$  while (x < 5)  
  { $\ell_3$  x := x + 1  $\ell_4$ };  
 $\ell_5$  x := 7  $\ell_6$ 
```



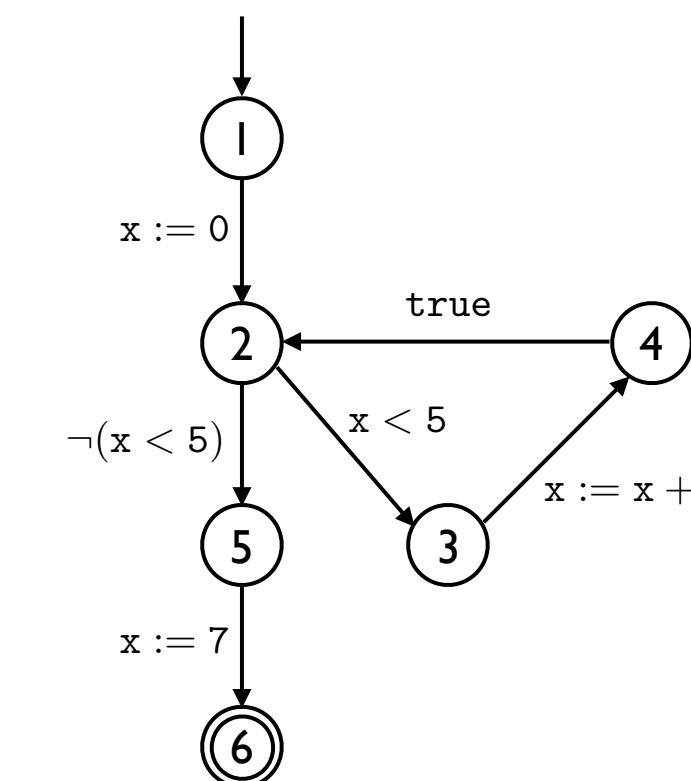
It represents infinite programs

CODE ABSTRACTION VS SEMANTIC ABSTRACTION

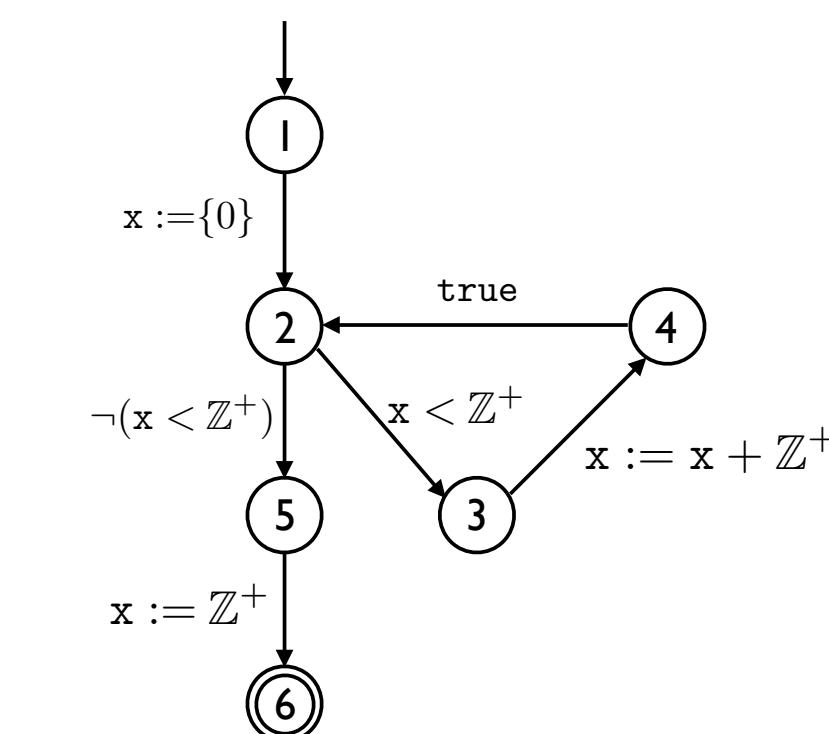
- Semantics naturally induces a code abstraction: By abstractly observing code we cannot distinguish between statements with involving the same abstract values

$\ell_1 x := 0;$
 $\ell_2 \text{while } (x < 5)$
 $\quad \{\ell_3 x := x + 1 \ell_4\};$
 $\ell_5 x := 7 \ell_6$

$\rho \in uco(\mathbb{V})$
 Semantic abstraction



Sign abstraction



$\bar{\Upsilon}[\rho](x := e) \stackrel{\text{def}}{=} x := \hat{\rho}(e) \stackrel{\text{def}}{=} \{x := e' \mid e' \in \hat{\rho}(e)\}$
 $\bar{\Upsilon}[\rho](b) \stackrel{\text{def}}{=} \hat{\rho}(b)$
 Code abstraction

Induced code abstraction $\bar{\Upsilon}[\rho]$ it can be embedded in the original analysis for ρ

$$\begin{aligned}
 \hat{\rho}(a) &: \left\{ \begin{array}{l} \hat{\rho}(a_1 \text{op } a_2) \stackrel{\text{def}}{=} \{a' \text{op } a'' \mid a' \in \hat{\rho}(a_1), a'' \in \hat{\rho}(a_2)\} \stackrel{\text{def}}{=} \hat{\rho}(a_1) \text{op } \hat{\rho}(a_2) \\ \hat{\rho}(x) \stackrel{\text{def}}{=} x, \quad \hat{\rho}(n) \stackrel{\text{def}}{=} \{m \mid m \in \rho(\{n\})\} \end{array} \right. \\
 \hat{\rho}(b) &: \left\{ \begin{array}{l} \hat{\rho}(b_1 \text{bop } b_2) \stackrel{\text{def}}{=} \hat{\rho}(b_1) \text{bop } \hat{\rho}(b_2), \quad \hat{\rho}(\neg b) \stackrel{\text{def}}{=} \neg \hat{\rho}(b) \\ \hat{\rho}(x) \stackrel{\text{def}}{=} x, \quad \hat{\rho}(\text{true}) \stackrel{\text{def}}{=} \{t \mid t \in \rho(\text{true})\}, \quad \hat{\rho}(\text{false}) \stackrel{\text{def}}{=} \{t \mid t \in \rho(\text{false})\} \end{array} \right. \\
 \hat{\rho}(s) &: \left\{ \begin{array}{l} \hat{\rho}(\text{concat}(s_1, s_2)) \stackrel{\text{def}}{=} \text{concat}(\hat{\rho}(s_1), \hat{\rho}(s_2)), \\ \hat{\rho}(\text{substr}(s, a_1, a_2)) \stackrel{\text{def}}{=} \text{substr}(\hat{\rho}(s), \hat{\rho}(a_1), \hat{\rho}(a_2)) \\ \hat{\rho}(x) \stackrel{\text{def}}{=} x, \quad \hat{\rho}(" \delta ") \stackrel{\text{def}}{=} \{\delta \mid \delta \in \rho(\sigma)\} \end{array} \right.
 \end{aligned}$$

WHICH CODE ABSTRACTIONS CAN BE ANALYZED?

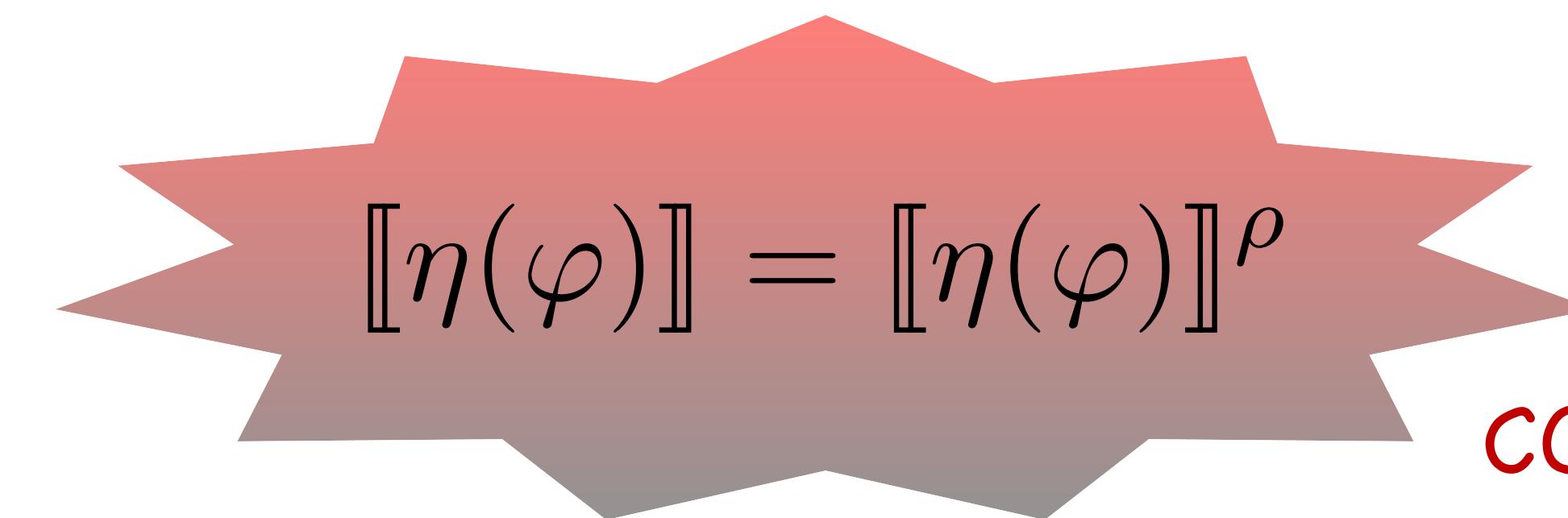
- Given a static analysis observing the semantic abstraction ρ , which code abstractions η can be still be analyzed without changing the analysis?
- Let φ denotes CFG labels

$$[\![\eta(\varphi)]\!] = [\![\eta(\varphi)]\!]^\rho$$

COMPLETENESS

WHICH CODE ABSTRACTIONS CAN BE ANALYZED?

- Given a static analysis observing the semantic abstraction ρ , which code abstractions η can be still be analyzed without changing the analysis?
- Let φ denotes CFG labels

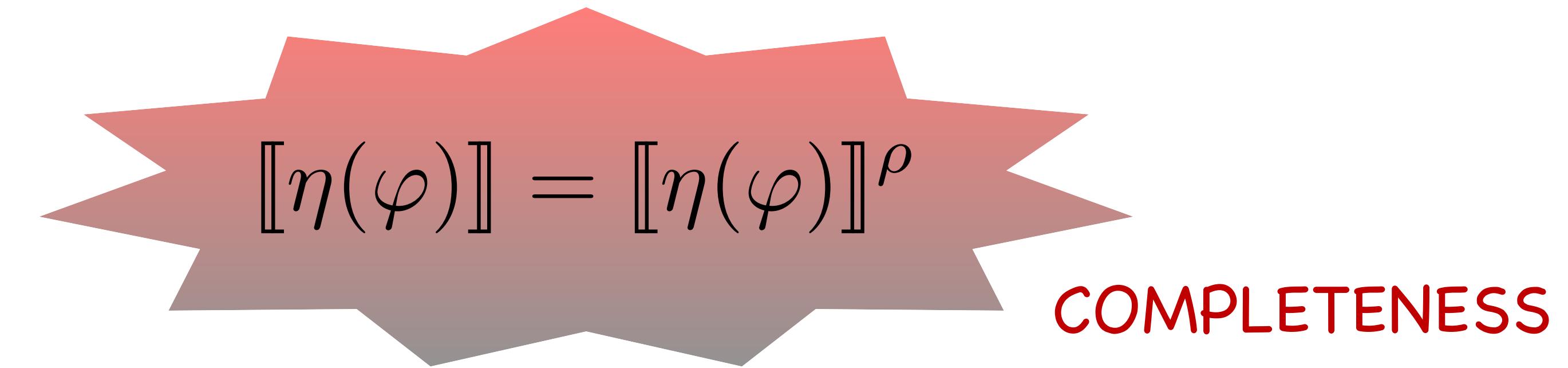
$$\llbracket \eta(\varphi) \rrbracket = \llbracket \eta(\varphi) \rrbracket^\rho$$


COMPLETENESS

- Example: $\rho=\text{Sign}$, $\eta(x:=5)=\{x:=n \mid 0 < n < 6\}$ ✗
- Example: $\rho=\text{Sign}$, $\eta(x:=5)=\{x:=n \mid n \geq 0\}$ ✓

WHICH CODE ABSTRACTIONS CAN BE ANALYZED?

- Given a static analysis observing the semantic abstraction ρ , which code abstractions η can be still be analyzed without changing the analysis?
- Let φ denotes CFG labels

$$\llbracket \eta(\varphi) \rrbracket = \llbracket \eta(\varphi) \rrbracket^\rho$$


COMPLETENESS

- The code abstraction $\bar{\eta}_\uparrow \stackrel{\text{def}}{=} \bar{\Upsilon}[\rho] \circ \eta$ always satisfies this equation, for any initial code abstraction η !
- We can make any abstract code analyzable by a given analysis!

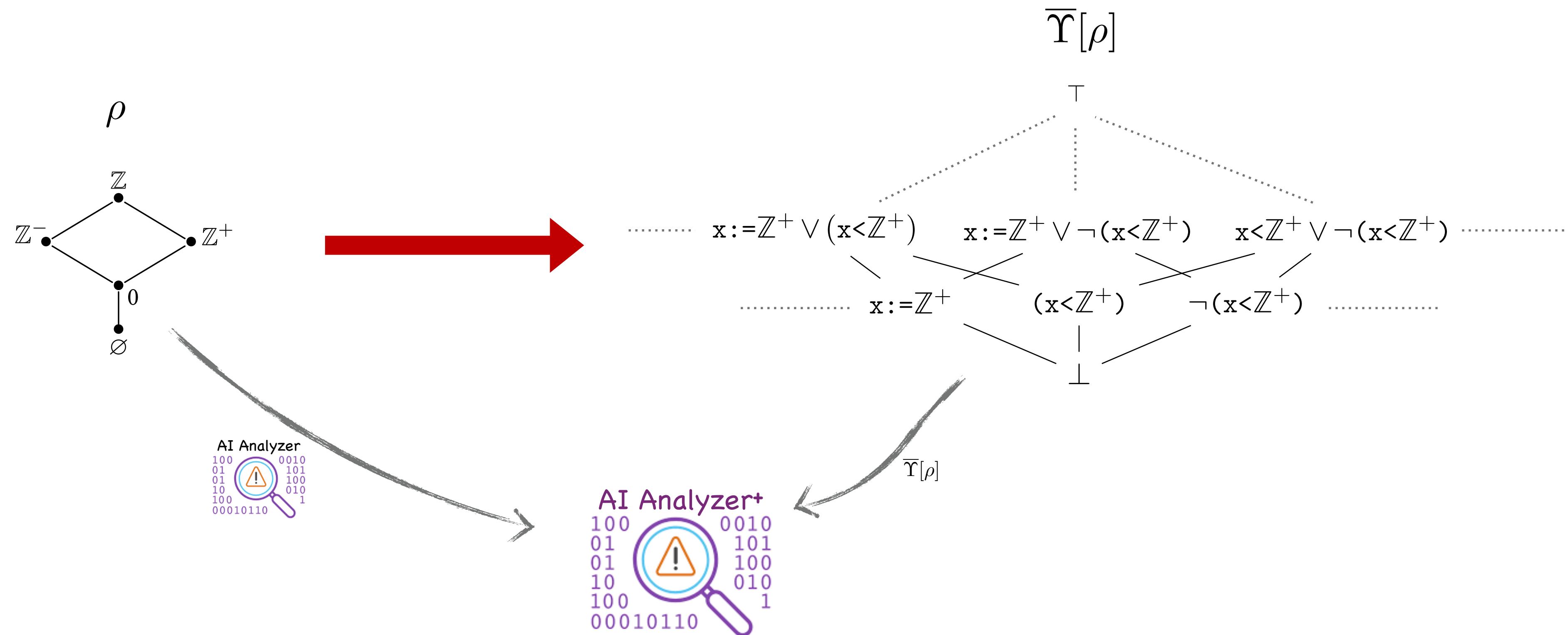
AN IMPROVED DYNAMIC CODE ANALYSIS

Analyzing abstract code



AN IMPROVED DYNAMIC CODE ANALYSIS

- Consider a static analysis observing the semantic abstraction ρ , interpreting also code abstracted by $\overline{\Upsilon}[\rho]$



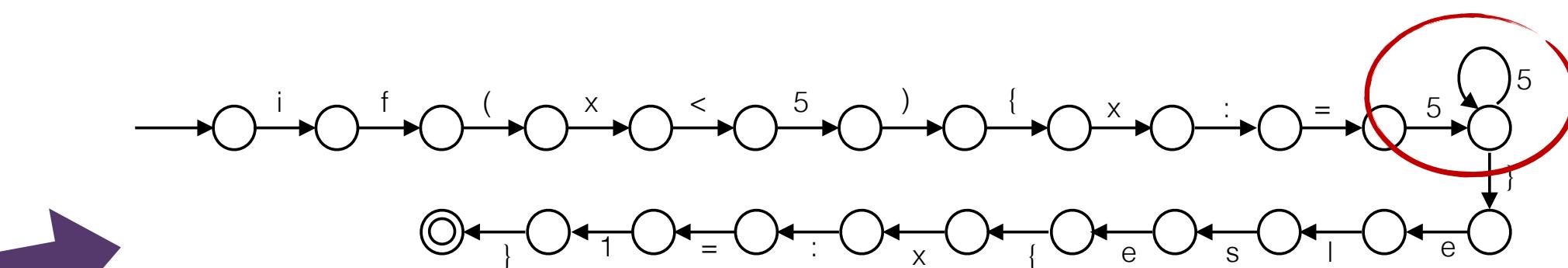
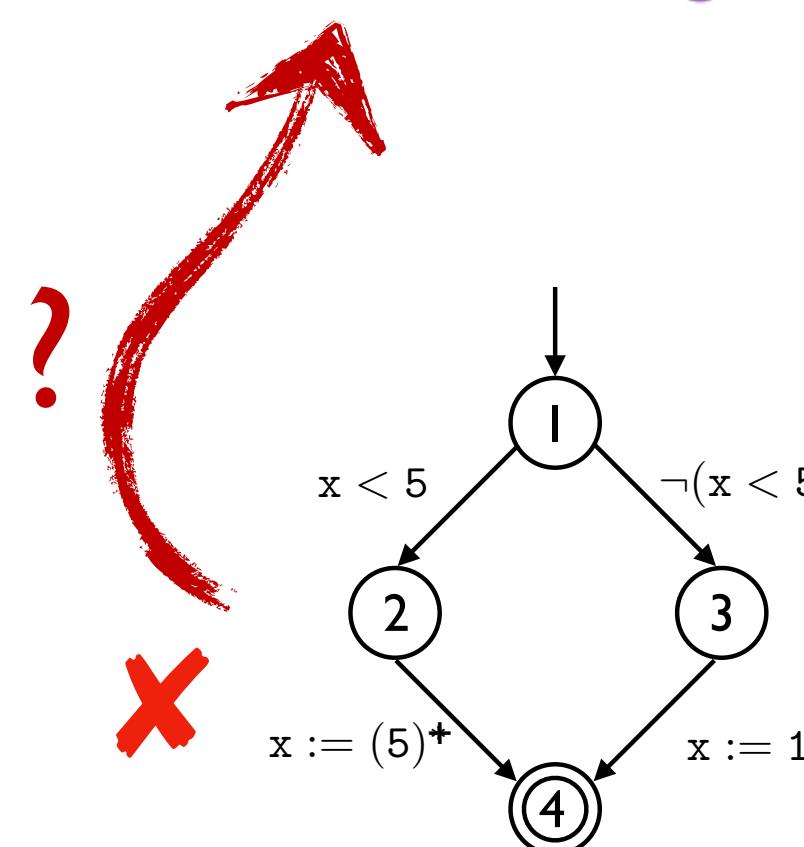
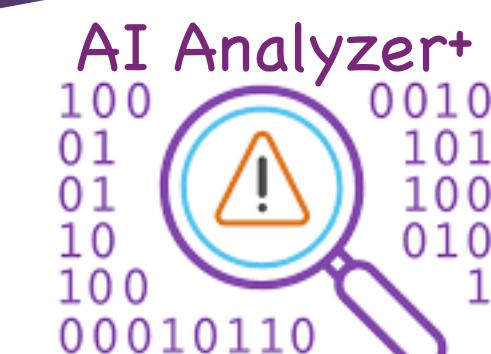
FROM AUTOMATA TO (ABSTRACT) CFG

- Analyze the program and when eval is met extract the automaton recognizing the language of its argument

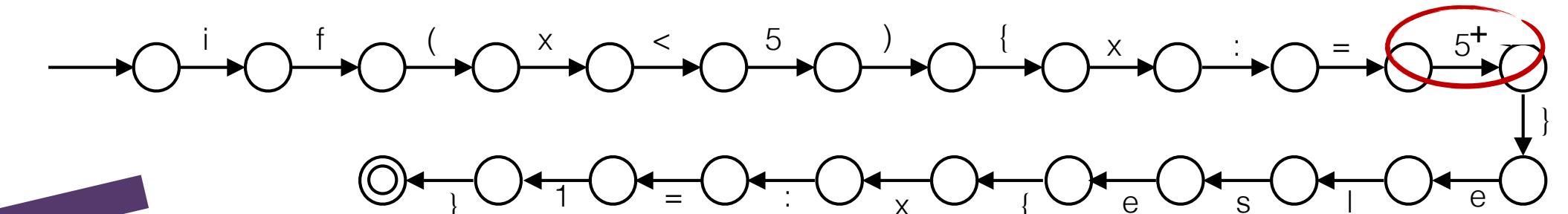
```

l1 str := "x:=5"; l2 i := 0;
l3 while (i < N) {
l4 str := str + "5";
l5 i := i+1; l6
}
l7 str := "if(x<5){" + str
    + "else{x:=1};";
l8 eval(str)l9

```



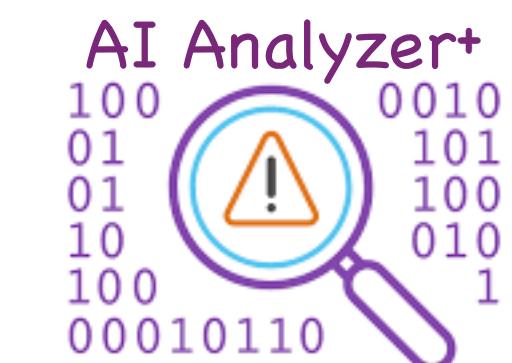
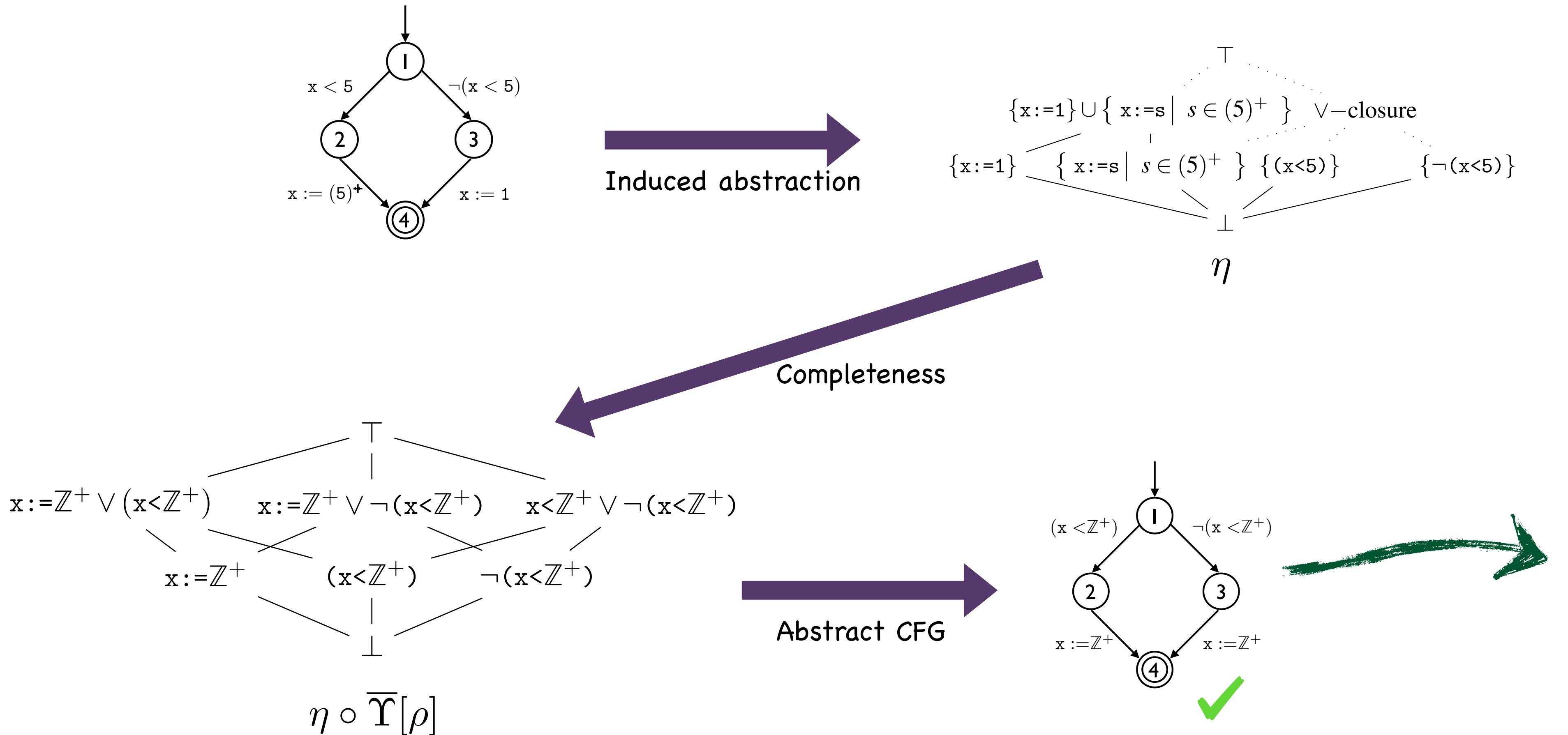
ReduceCycles



CFG generation algorithm

MAKING (ABSTRACT) CODE ANALYZABLE

- The abstract d=code generated induces a code abstraction η





CONCLUDING

Where we are and where we can go

"YOU CAN ANALYZE THE CODE YOU DON'T SEE!"

