

# Symbolic Execution

Suman Jana

---

**Acknowledgement: Baishakhi Ray (Uva), Omar Chowdhury (Purdue), Saswat Anand (GA Tech), Rupak Majumdar (UCLA), Koushik Sen (UCB)**

A solid orange horizontal bar spanning the width of the slide, located at the bottom.

# What is the goal?

```
static OSStatus
SSLVerifySignedServerKeyExchange(SSLContext *ctx, bool isRsa, SSLBuffer signedParams,
                                uint8_t *signature, UInt16 signatureLen)
{
    OSStatus      err;
    ...

    if ((err = SSLHashSHA1.update(&hashCtx, &serverRandom)) != 0)
        goto fail;
    if ((err = SSLHashSHA1.update(&hashCtx, &signedParams)) != 0)
        goto fail;
    if ((err = SSLHashSHA1.final(&hashCtx, &hashOut)) != 0)
        goto fail;

    // code omitted for brevity...

    err = sslRawVerify(ctx,
                      ctx->peerPubKey,
                      dataToSign,           /* plaintext */
                      dataToSignLen,        /* plaintext length */
                      signature,
                      signatureLen);

    if(err) {
        sslErrorLog("SSLDecodeSignedServerKeyExchange: sslRawVerify "
                   "returned %d\n", (int)err);
        goto fail;
    }

fail:
    SSLFreeBuffer(&signedHashes);
    SSLFreeBuffer(&hashCtx);
    return err;
}
```

Oops...

Never gets called (but needed to be)...

Despite the name, always returns "it's OK!!!"

# Testing

- Testing approaches are in general manual
  - Time consuming process
  - Error-prone
  - Incomplete
  - Depends on the quality of the test cases or inputs
  - Provides little in terms of coverage
-

Can

Yes, we can.

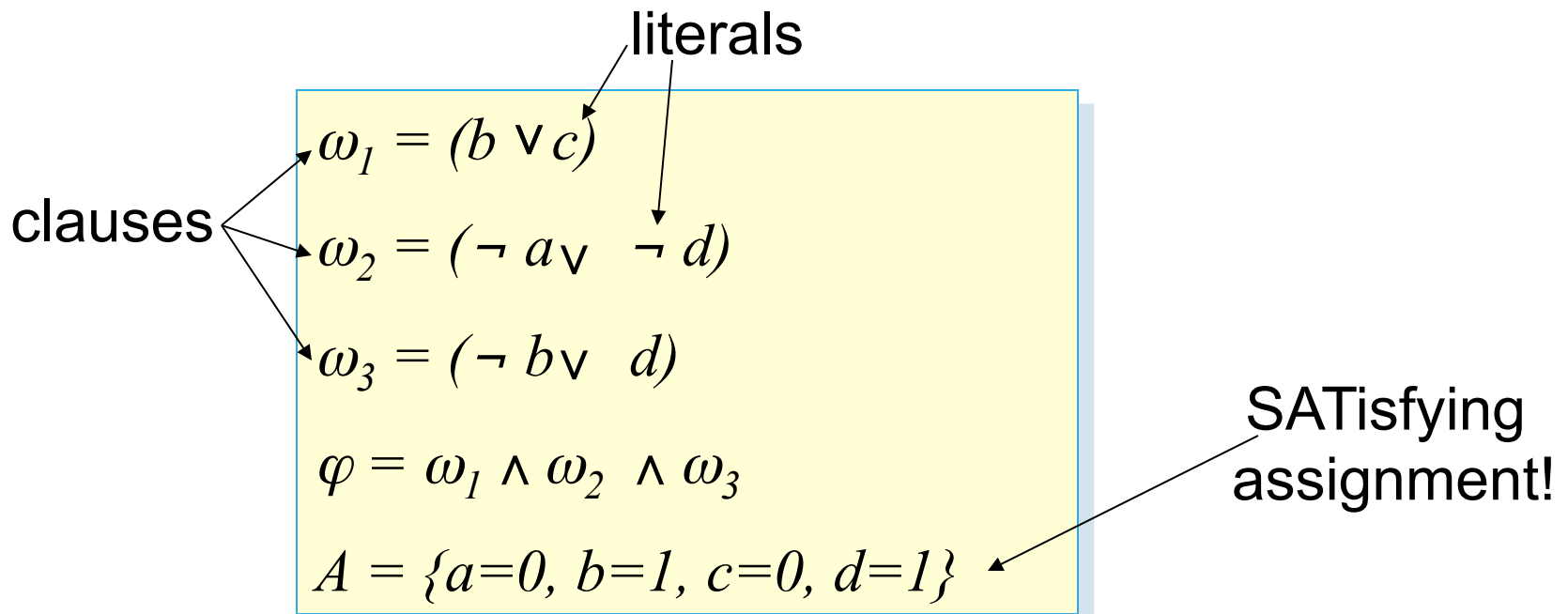
forms of

can we some

it automatic?

# Background: SAT

Given a propositional formula in CNF, find if there exists an assignment to Boolean variables that makes the formula true:



# Background: SMT (Satisfiability Modulo Theory)

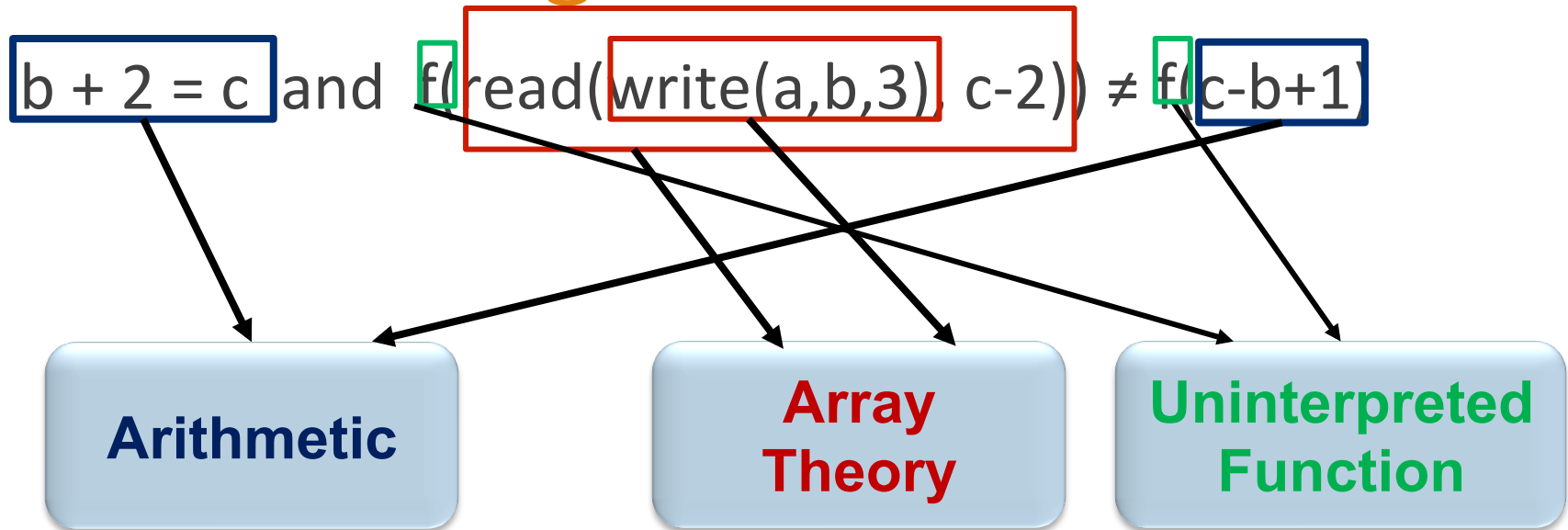
- An SMT instance is a generalization of a [Boolean SAT](#) instance
- Various sets of variables are replaced by [predicates](#) from a variety of underlying theories.

**Input:** a **first-order** formula  $\varphi$  over background theory (Arithmetic, Arrays, Bit-vectors, Algebraic Datatypes)

**Output:** is  $\varphi$  satisfiable?

- does  $\varphi$  have a model?
  - Is there a refutation of  $\varphi$  = proof of  $\neg\varphi$ ?
-

# Background: SMT



# Example SMT Solving

$b + 2 = c$  and  $f(\text{read}(\text{write}(a, b, 3), c - 2)) \neq f(c - b + 1)$

[Substituting  $c$  by  $b + 2$ ]

$b + 2 = c$  and  $f(\text{read}(\text{write}(a, b, 3), b + 2 - 2)) \neq f(b + 2 - b + 1)$

[Arithmetic simplification]

$b + 2 = c$  and  $f(\text{read}(\text{write}(a, b, 3), b)) \neq f(3)$

[Applying array theory axiom]

forall  $a, i, v$ :  $\text{read}(\text{write}(a, i, v), i) = v$

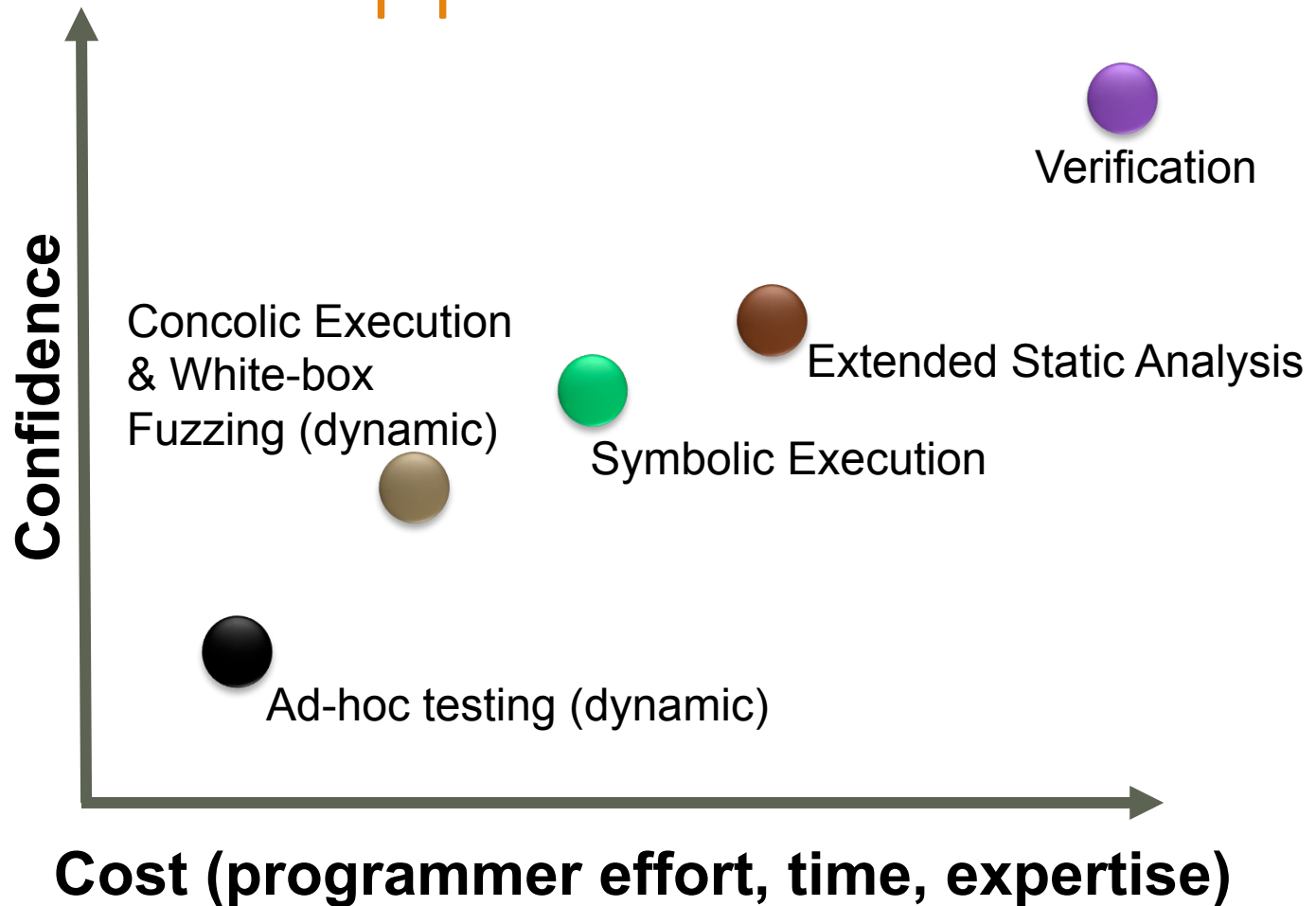
$b + 2 = c$  and  $f(3) \neq f(3)$  [**NOT SATISFIABLE**]

$\text{read} : \text{array} \times \text{index} \rightarrow \text{element}$

$\text{write} : \text{array} \times \text{index} \times \text{element} \rightarrow \text{array}$



# Program Validation Approaches



# Automatic Test Generation

## Symbolic & Concolic Execution

How do we automatically generate test inputs that induce the program to go in different paths?

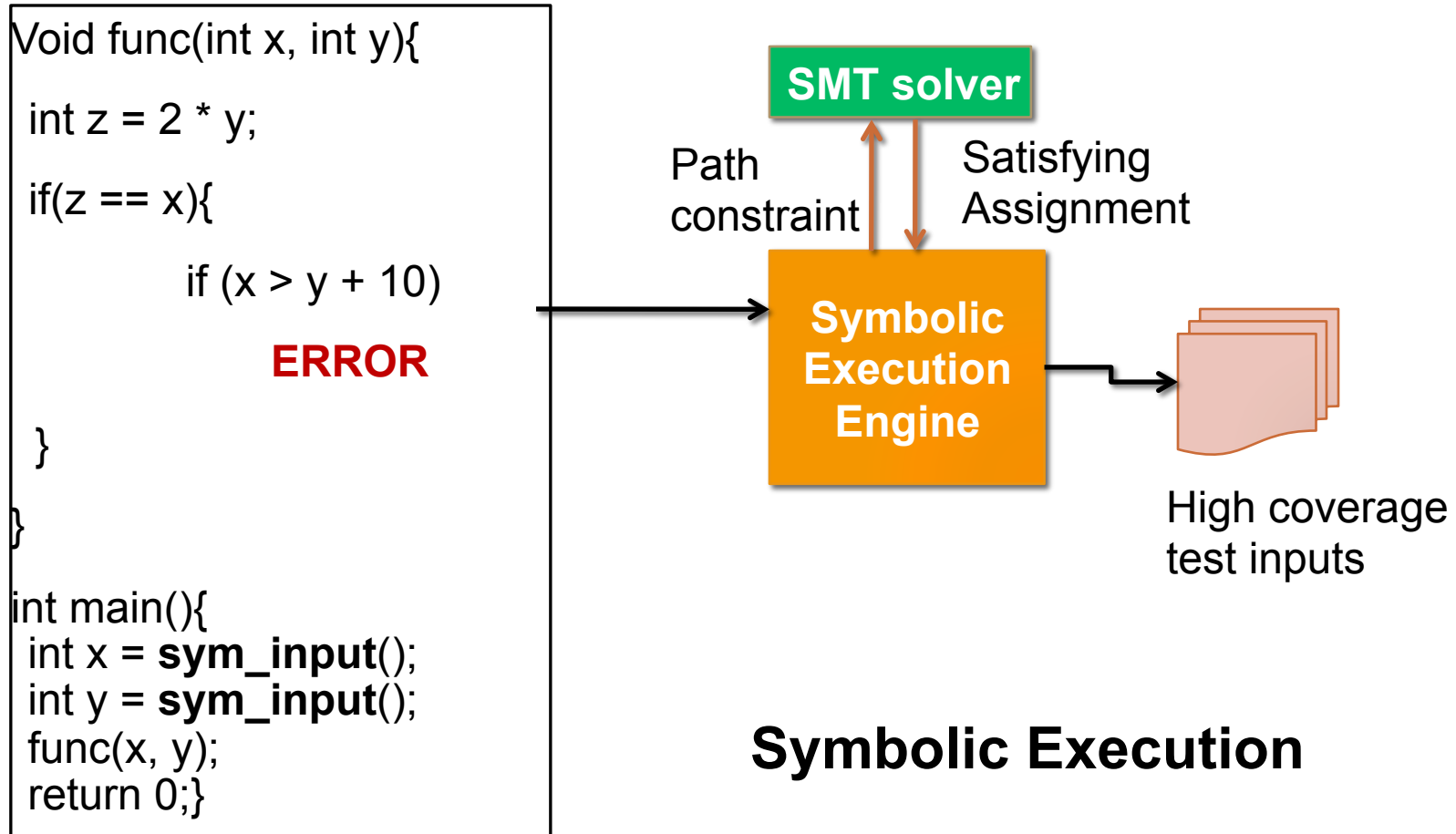
### Intuition:

- Divide the whole possible input space of the program into equivalent classes of input.
  - For each equivalence class, all inputs in that equivalence class will induce the same program path.
  - Test one input from each equivalence class.
-

# Logistics Update

- PA1 is due \*before class\* on Wednesday (8<sup>th</sup> Feb).
  - Answer the last two questions (part 2 and 3) in PA1 in as much detail as possible
  - Your piazza answers will be counted as class participation, so be a good citizen and try to help others!
  - Project timeframe:
    - List of group members due - Feb 13th before class
    - Project proposal (1 page) due - Feb 22nd before class
    - Midterm project status report (1 page) due - Mar 27th before class
    - Final report ( $\geq 6$  pages) due - May 3<sup>rd</sup>
  - There will be 4 PAs (no quizzes). Two before midterm and two after midterm.
  - The midterm will be open book/note.
-

# Symbolic Execution



# Symbolic Execution

Execute the program with symbolic valued inputs (**Goal: good path coverage**)

Represents *equivalence class of inputs* with first order logic formulas (**path constraints**)

One path constraint abstractly represents all inputs that induces the program execution to go down a specific path

Solve the path constraint to obtain one representative input that exercises the program to go down that specific path

**Symbolic execution implementations:** KLEE, Java Pathfinder, etc.

---

# More details on Symbolic Execution

Instead of concrete state, the program maintains **symbolic states**, each of which maps variables to symbolic values

**Path condition** is a quantifier-free formula over the symbolic inputs that encodes all branch decisions taken so far

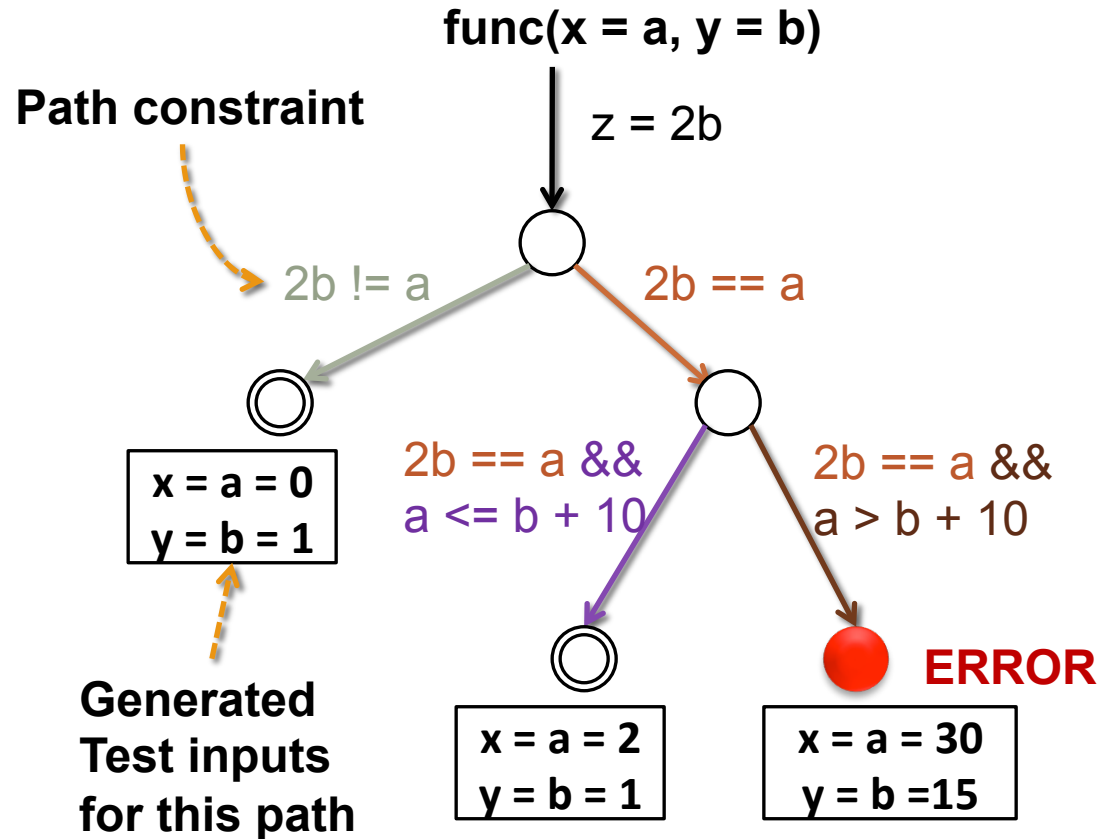
All paths in the program form its **execution tree**, in which some paths are feasible and some are infeasible

---

# Symbolic Execution

```
Void func(int x, int y){  
    int z = 2 * y;  
    if(z == x){  
        if (x > y + 10)  
            ERROR  
    }  
}  
  
int main(){  
    int x = sym_input();  
    int y = sym_input();  
    func(x, y);  
    return 0;  
}
```

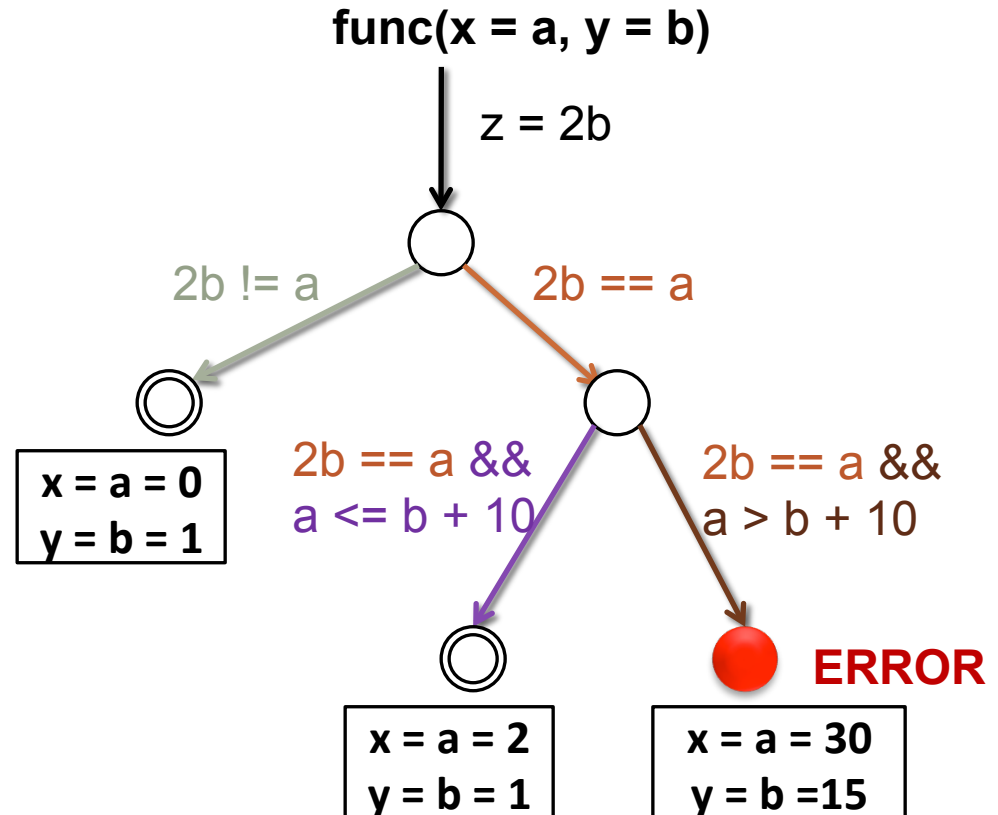
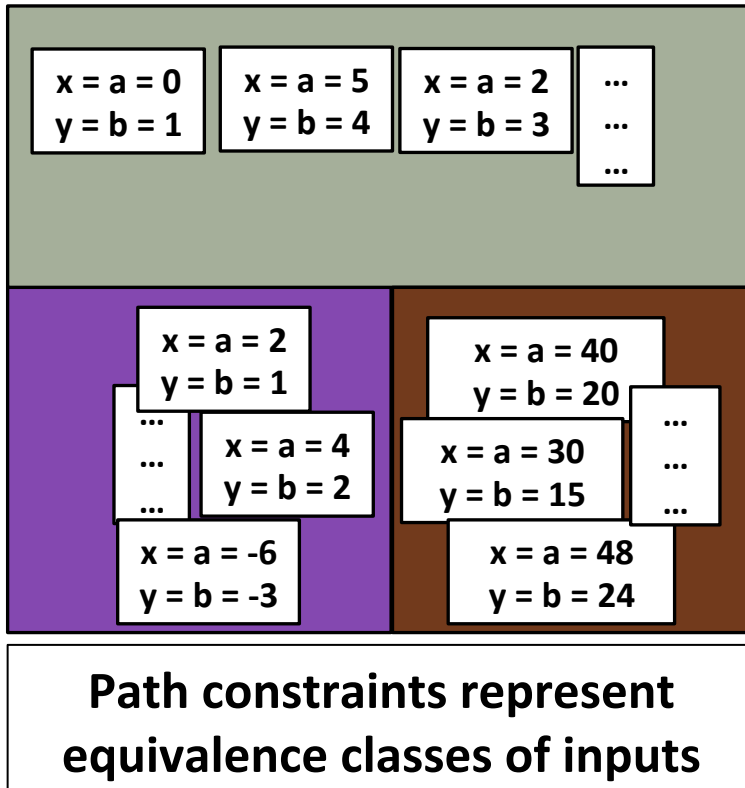
**How does symbolic execution work?**



**Note: Require inputs to be marked as symbol**

# Symbolic Execution

How does symbolic execution work?

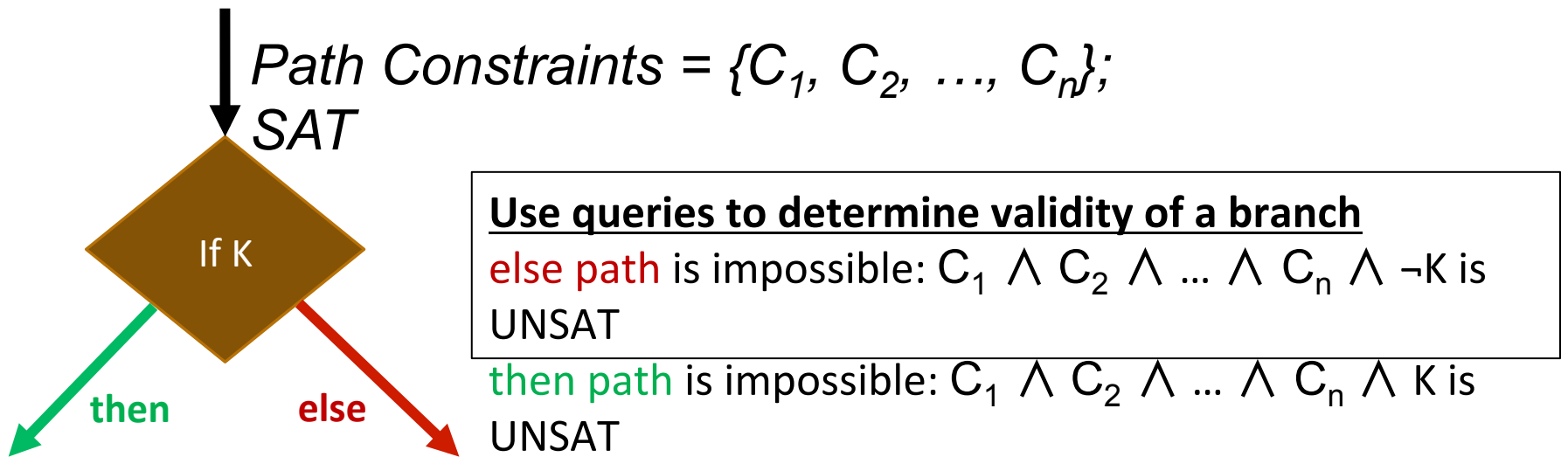




# SMT Queries

Counterexample queries (generate a test case)

Branch queries (whether a branch is valid)



# Optimizing SMT Queries

## Expression rewriting

- Simple arithmetic simplifications ( $x * 0 = 0$ )
- Strength reduction ( $x * 2^n = x \ll n$ )
- Linear simplification ( $2 * x - x = x$ )

## Constraint set simplification

- $x < 10 \ \&\& \ x = 5 \ \rightarrow \ x = 5$

## Implied Value Concretization

- $x + 1 = 10 \ \rightarrow \ x = 9$

## Constraint Independence

- $i < j \ \&\& \ j < 20 \ \&\& \ k > 0 \ \&\& \ i = 20 \ \rightarrow \ i < j \ \&\& \ i < 20 \ \&\& \ i = 20$

# Optimizing SMT Queries (contd.)

## Counter-example Cache

- $i < 10 \ \&\& \ i = 10$  (no solution)
- $i < 10 \ \&\& \ j = 8$  (satisfiable, with variable assignments  $i \rightarrow 5, j \rightarrow 8$ )

## Superset of unsatisfiable constraints

- $\{i < 10, i = 10, j = 12\}$  (unsatisfiable)

## Subset of satisfiable constraints

- $i \rightarrow 5, j \rightarrow 8$ , satisfies  $i < 10$

## Superset of satisfiable constraints

- Same variable assignments might work

# How does Symbolic Execution Find bugs?

It is possible to extend symbolic execution to help find bugs.

**How:** Dedicated checkers

- **Divide by zero example** ---  $y = x / z$  where  $x$  and  $z$  are variables and assume current PC is  $f$
- Even though we only fork in the branch leading to the division operator
- One branch in which  $z \neq 0$
- We will get constraints:

$z \neq 0$

- $f \wedge z \neq 0$  will give us concrete input values that will not cause an error.

Write a dedicated checker for each kind of bug (e.g., buffer overflow, integer overflow, integer underflow)

# Classic Symbolic Execution ---

## Practical Issues

**Loops and recursions** --- infinite execution tree

**Path explosion** --- exponentially many paths

**Heap modeling** --- symbolic data structures and pointers

**SMT solver limitations** --- dealing with complex path constraints

**Environment modeling** --- dealing with native/system/library calls/file operations/network events

**Coverage Problem** --- may not reach deep into the execution tree, specially when encountering loops.

# Solution: Concolic Execution

**Concolic** = **Concrete** + **Symbolic**

Combining Classical Testing with Automatic  
Program Analysis

Also called **dynamic symbolic execution**

The intention is to visit deep into the program execution tree

Program is simultaneously executed with concrete and symbolic inputs


Start off the execution with a random input

Specially useful in cases of remote procedure call

**Concolic execution implementations:** SAGE (Microsoft), CREST

# Concolic Execution Steps

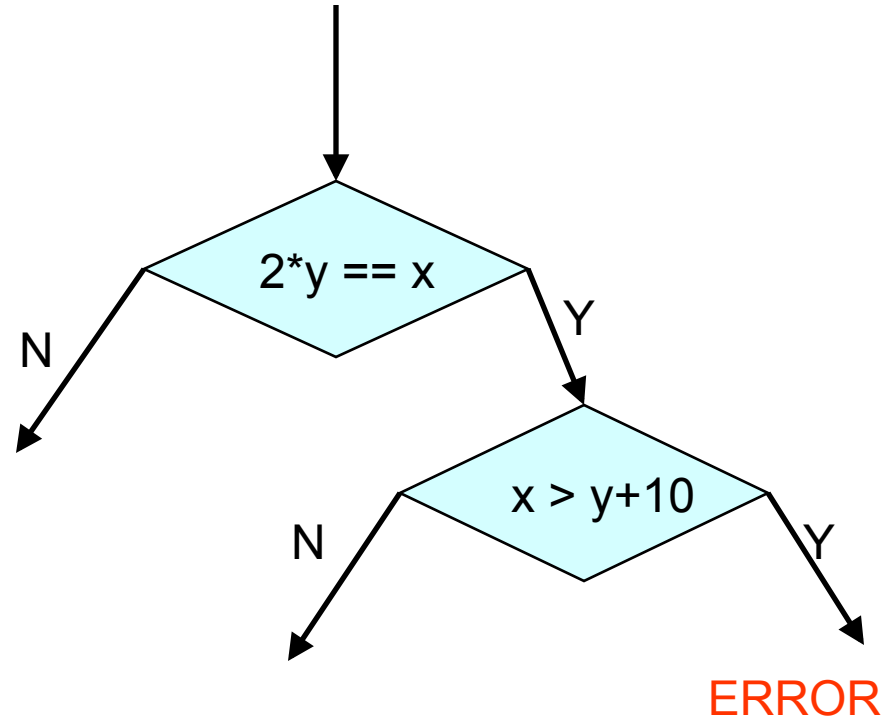
- Generate a random seed input to start execution
- Concretely execute the program with the random seed input and collect the path constraint
- Example: **a && b && c**
- In the next iteration, negate the last conjunct to obtain the constraint **a && b && !c**
- Solve it to get input to the path which matches all the branch decisions except the last one



Why not from the first?

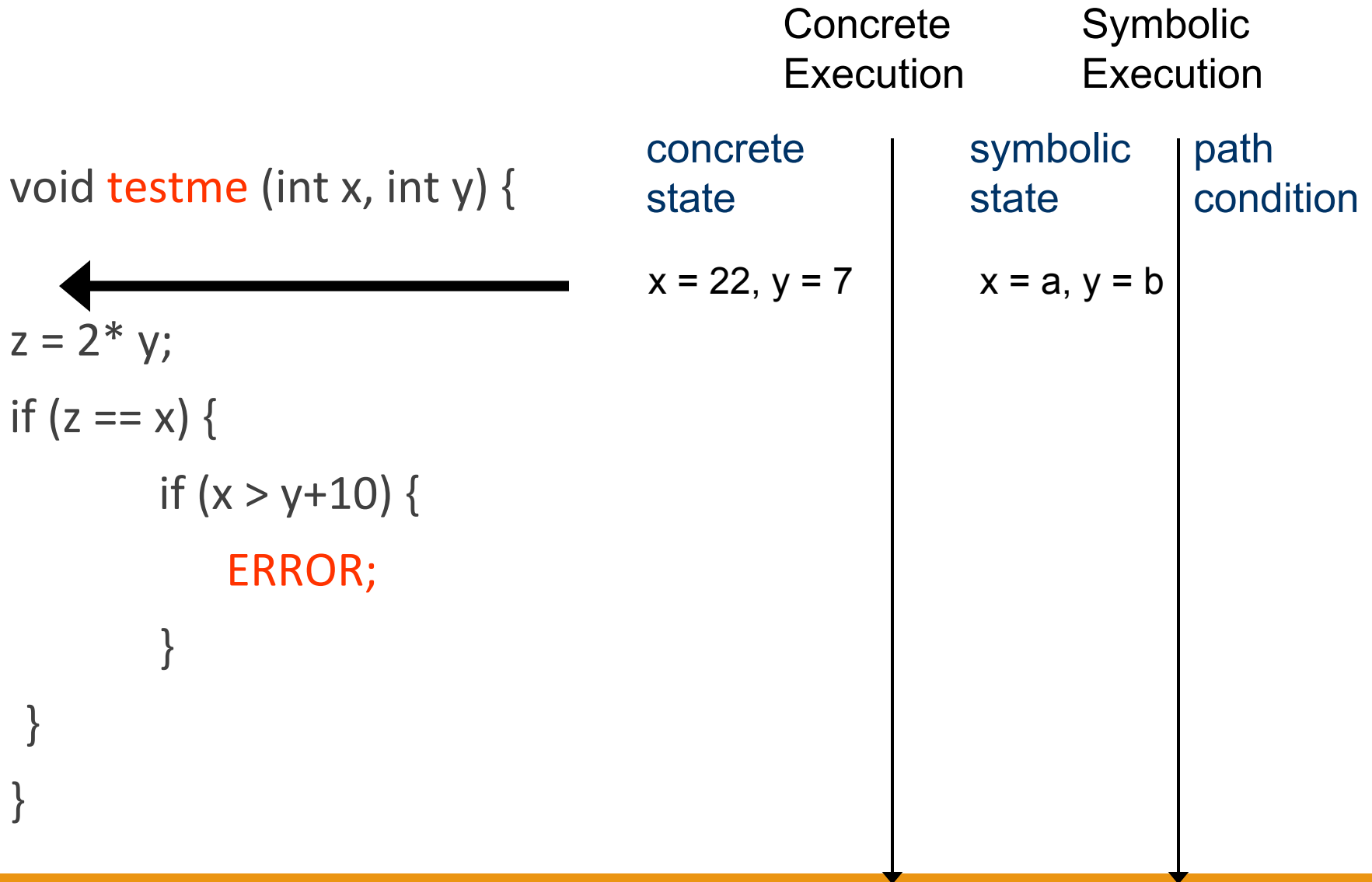
# Example

```
void testme (int x, int y)
{
  z = 2*y;
  if (z == x) {
    if (x > y+10) {
      ERROR;
    }
  }
}
```

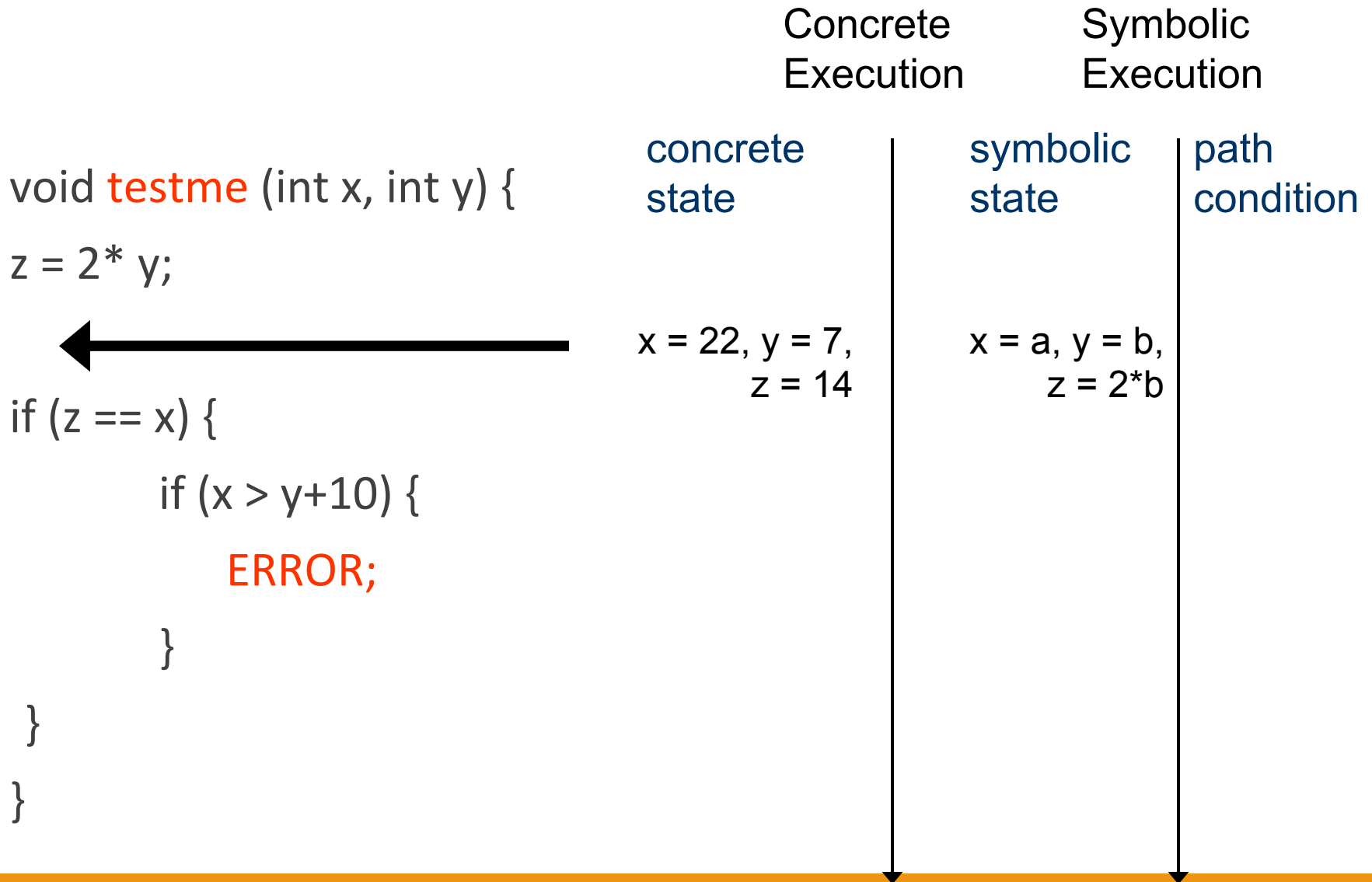




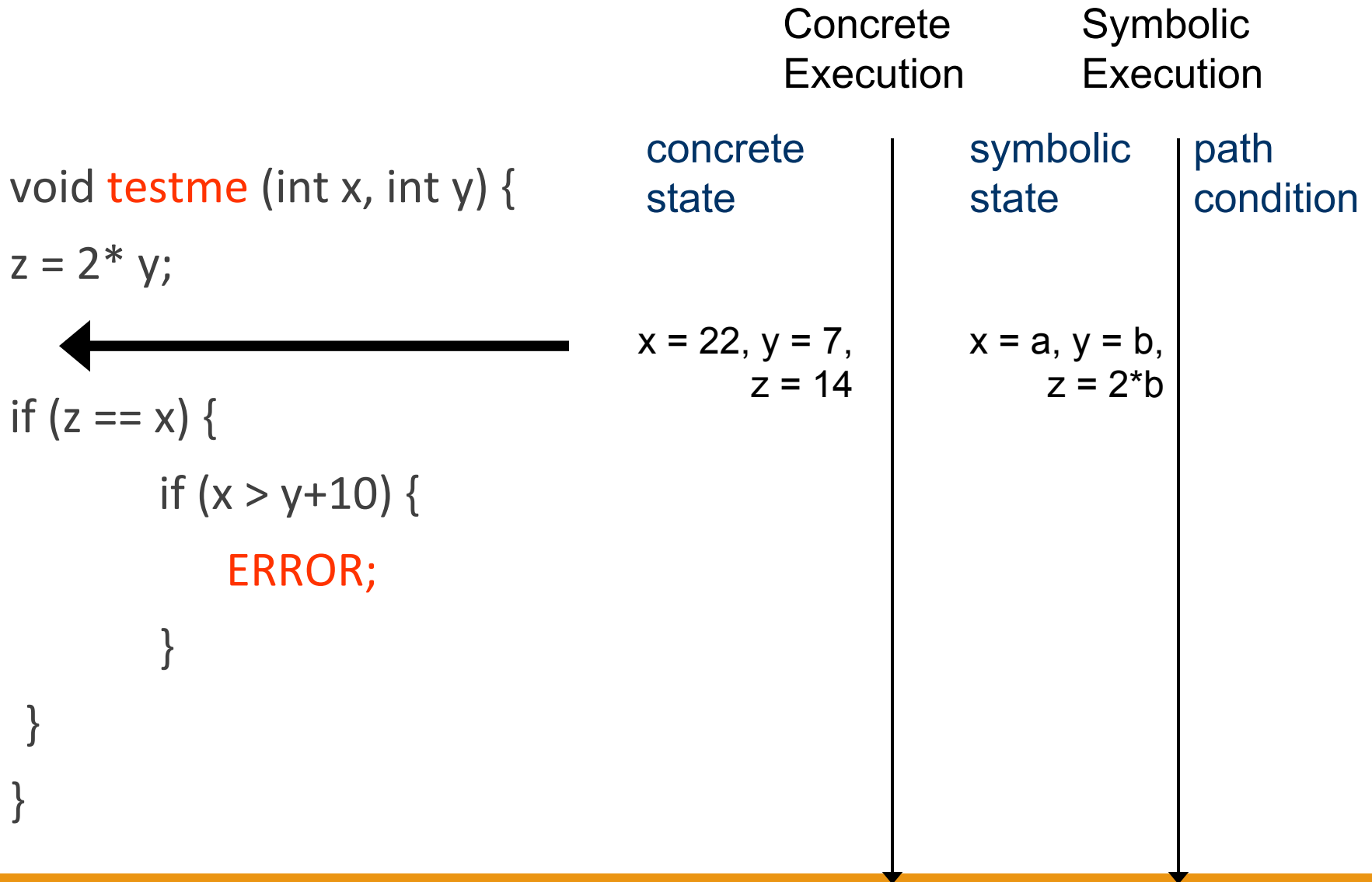
# Concolic execution example



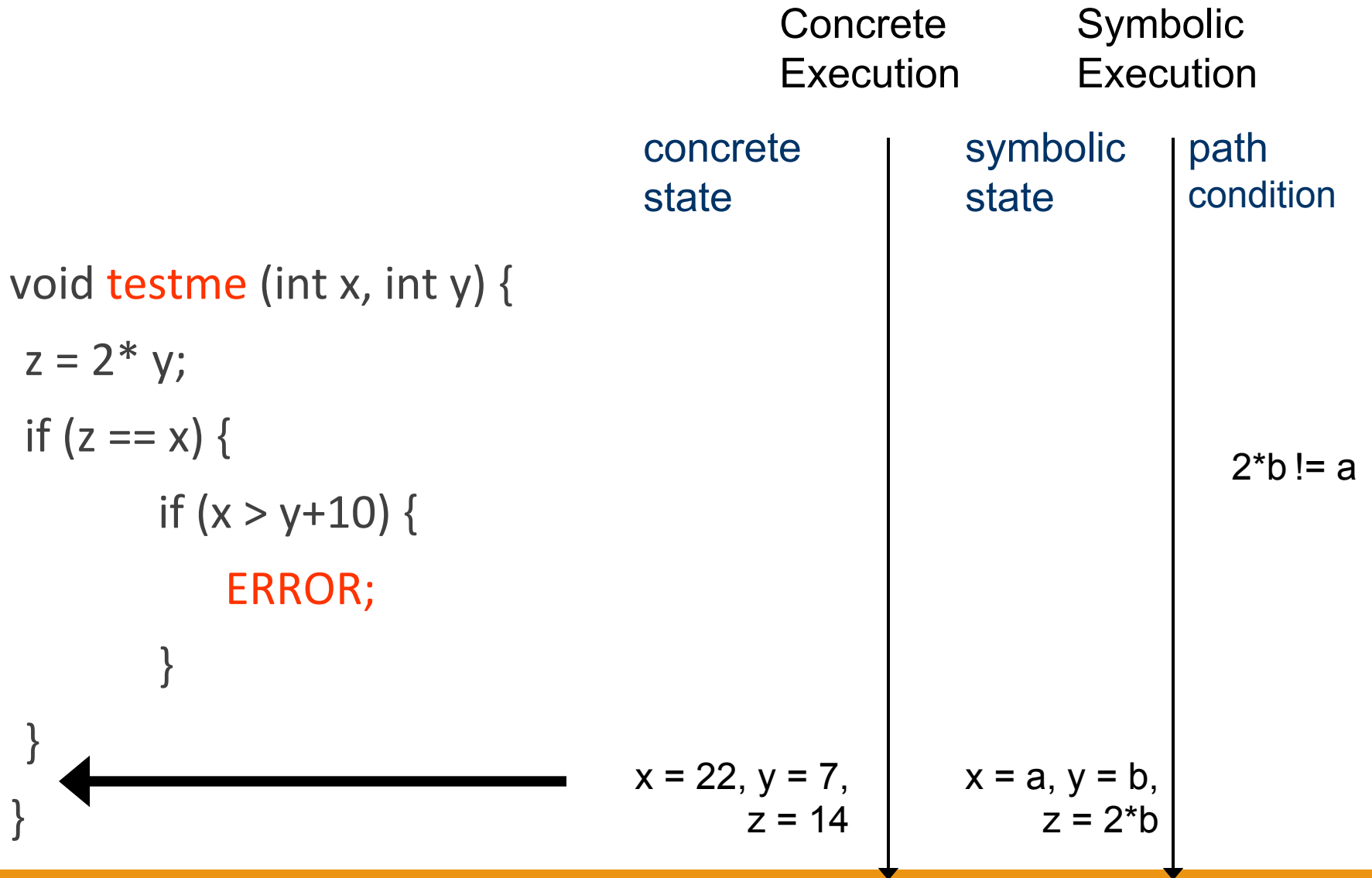
# Concolic execution example



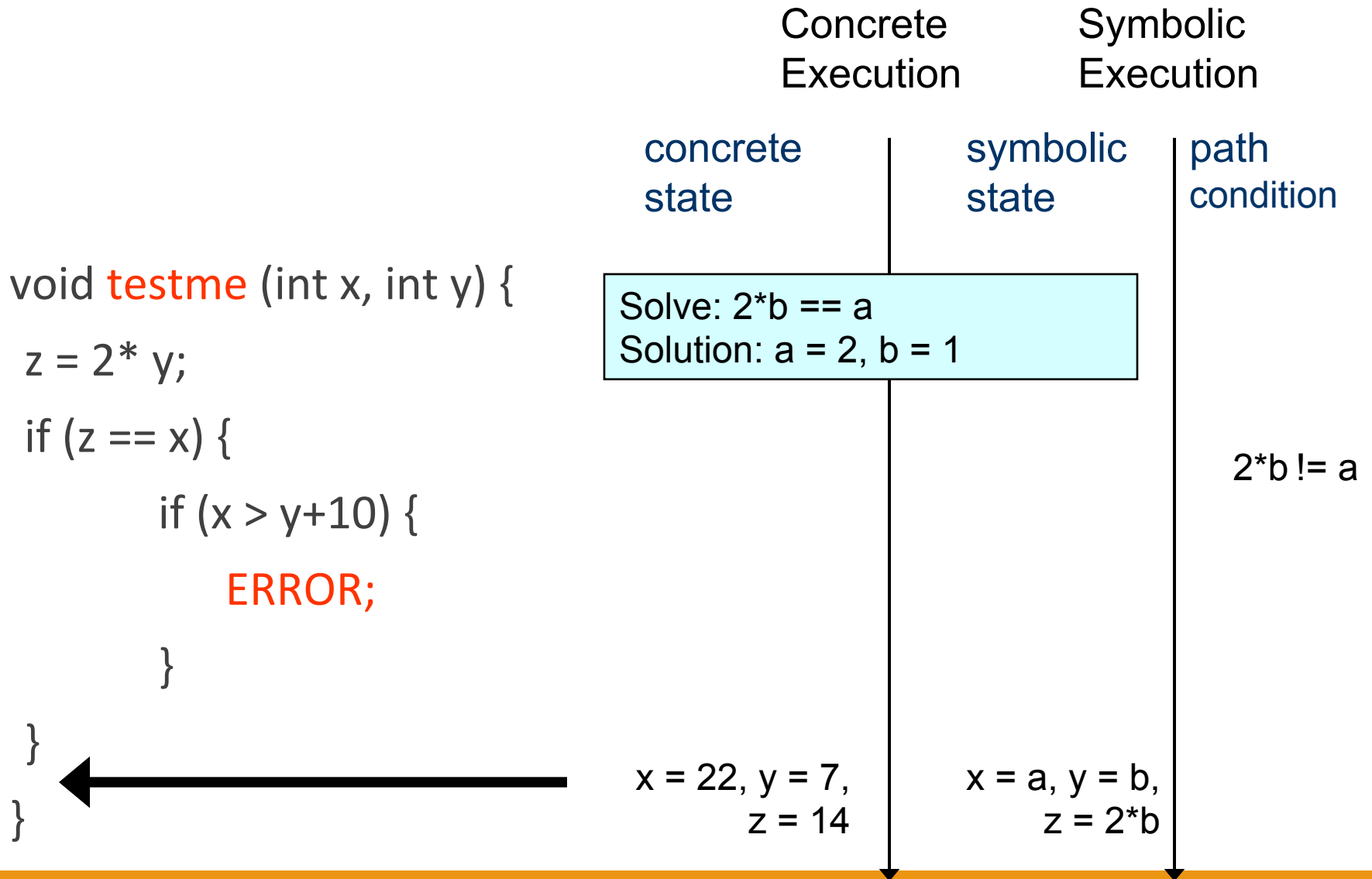
# Concolic execution example



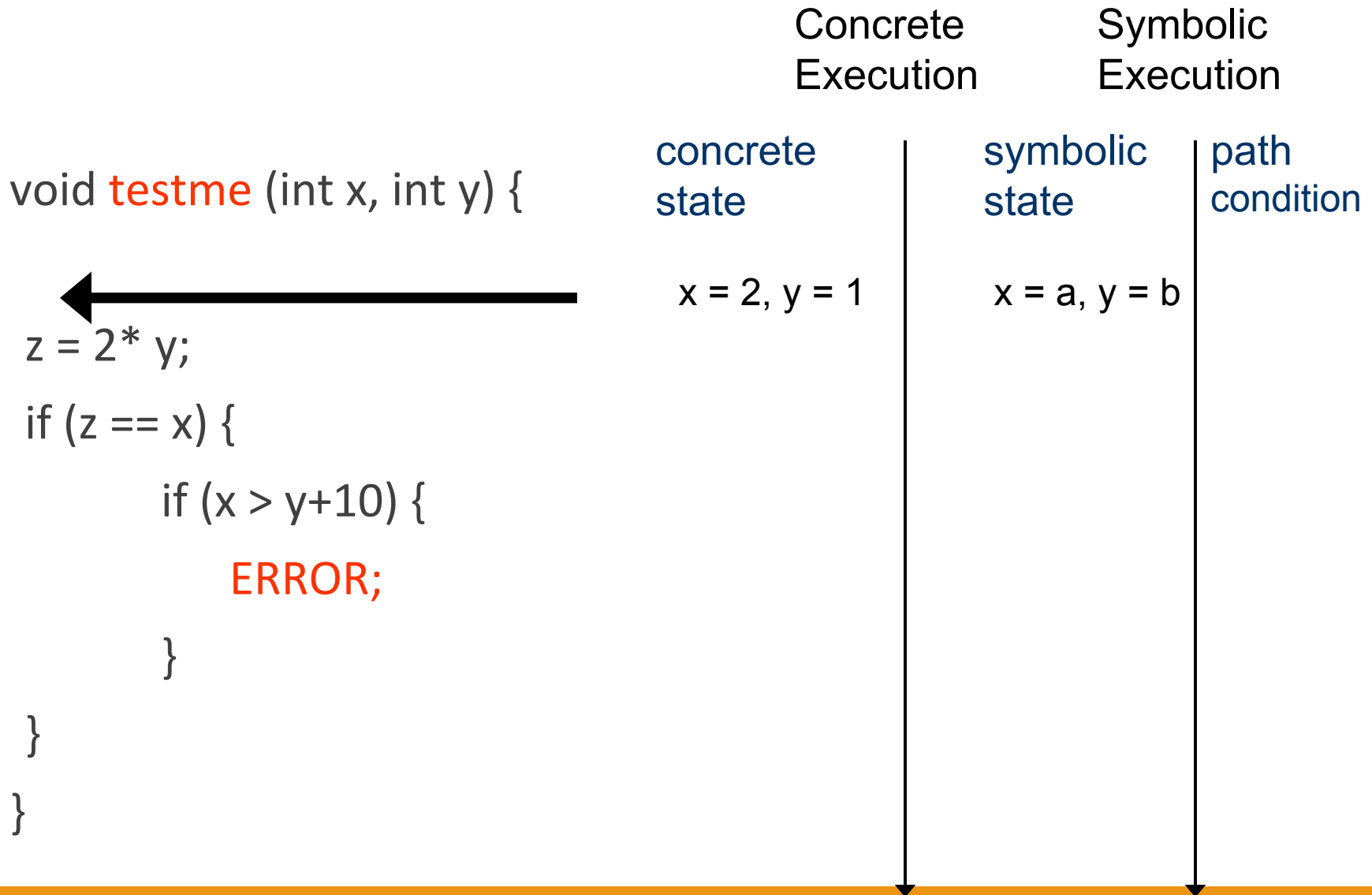
# Concolic execution example



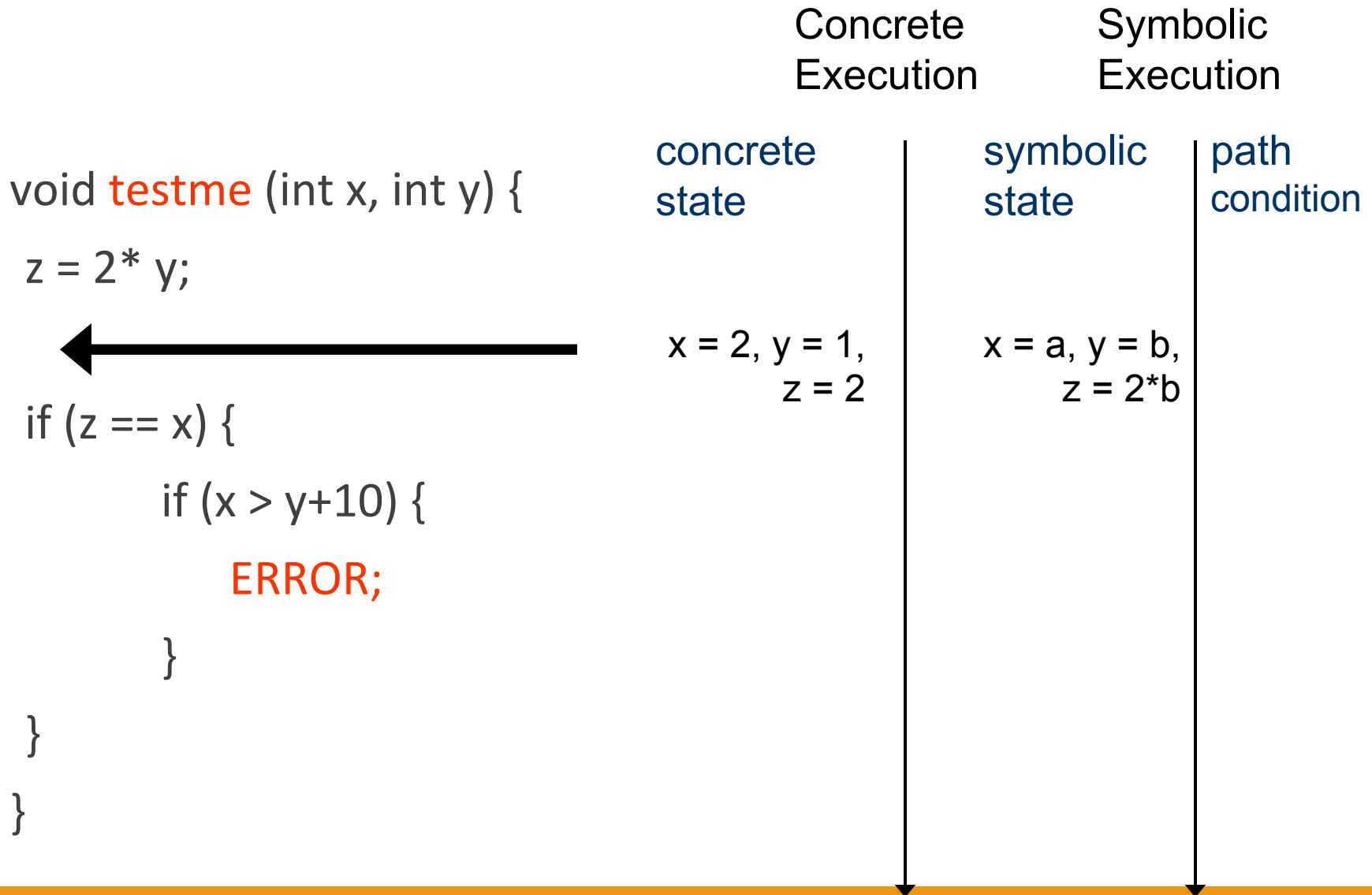
# Concolic execution example



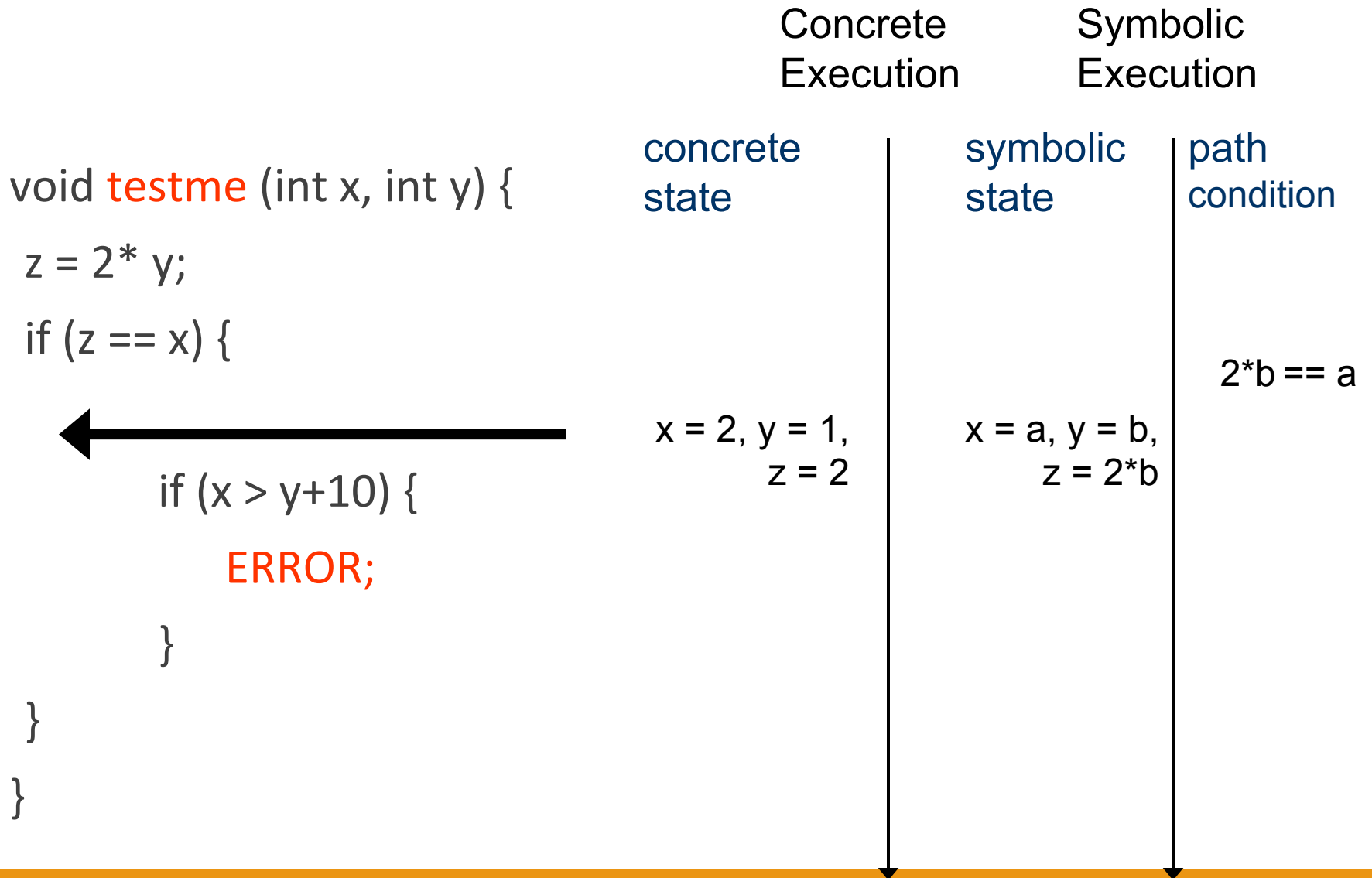
# Concolic execution example



# Concolic execution example

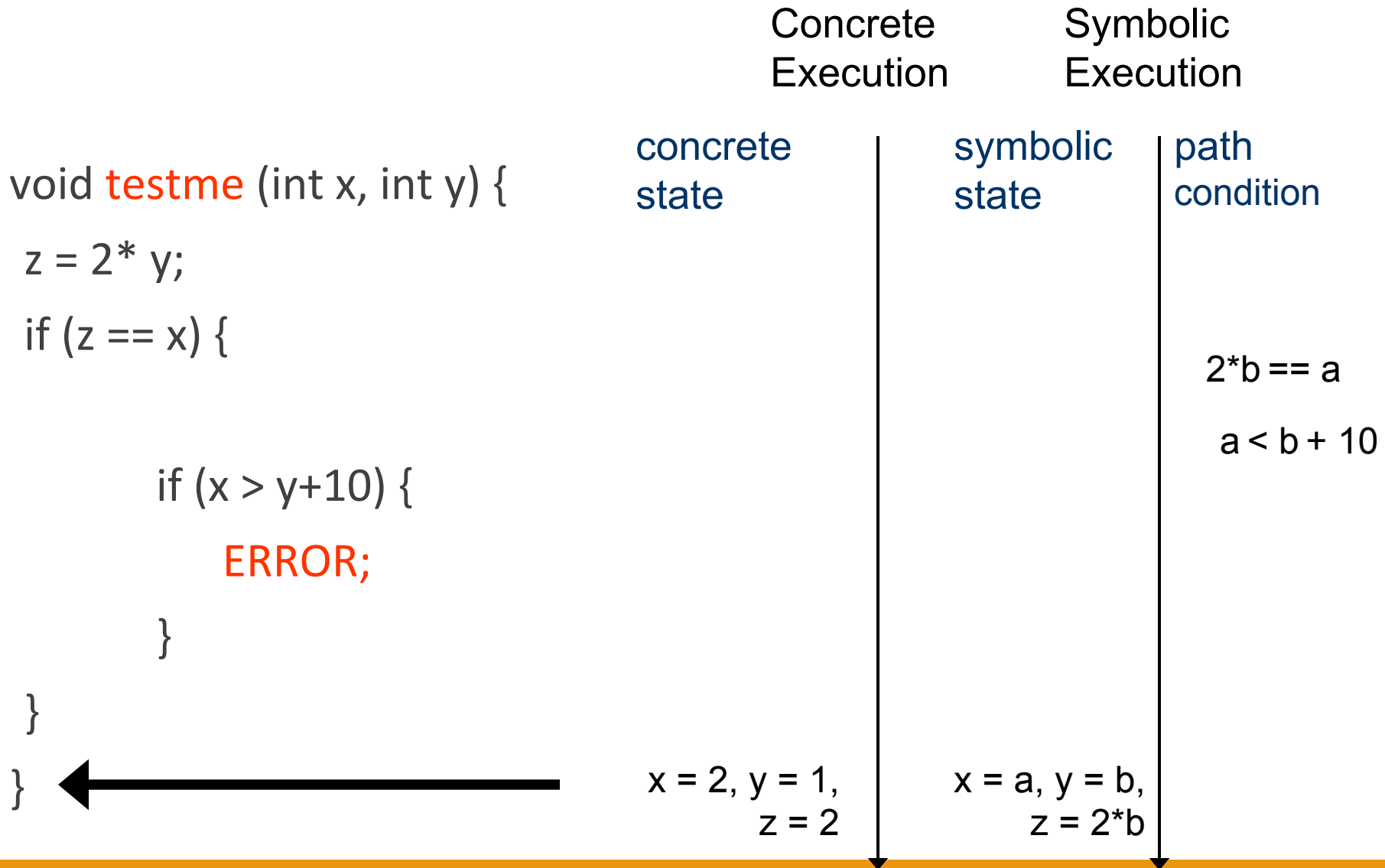


# Concolic execution example

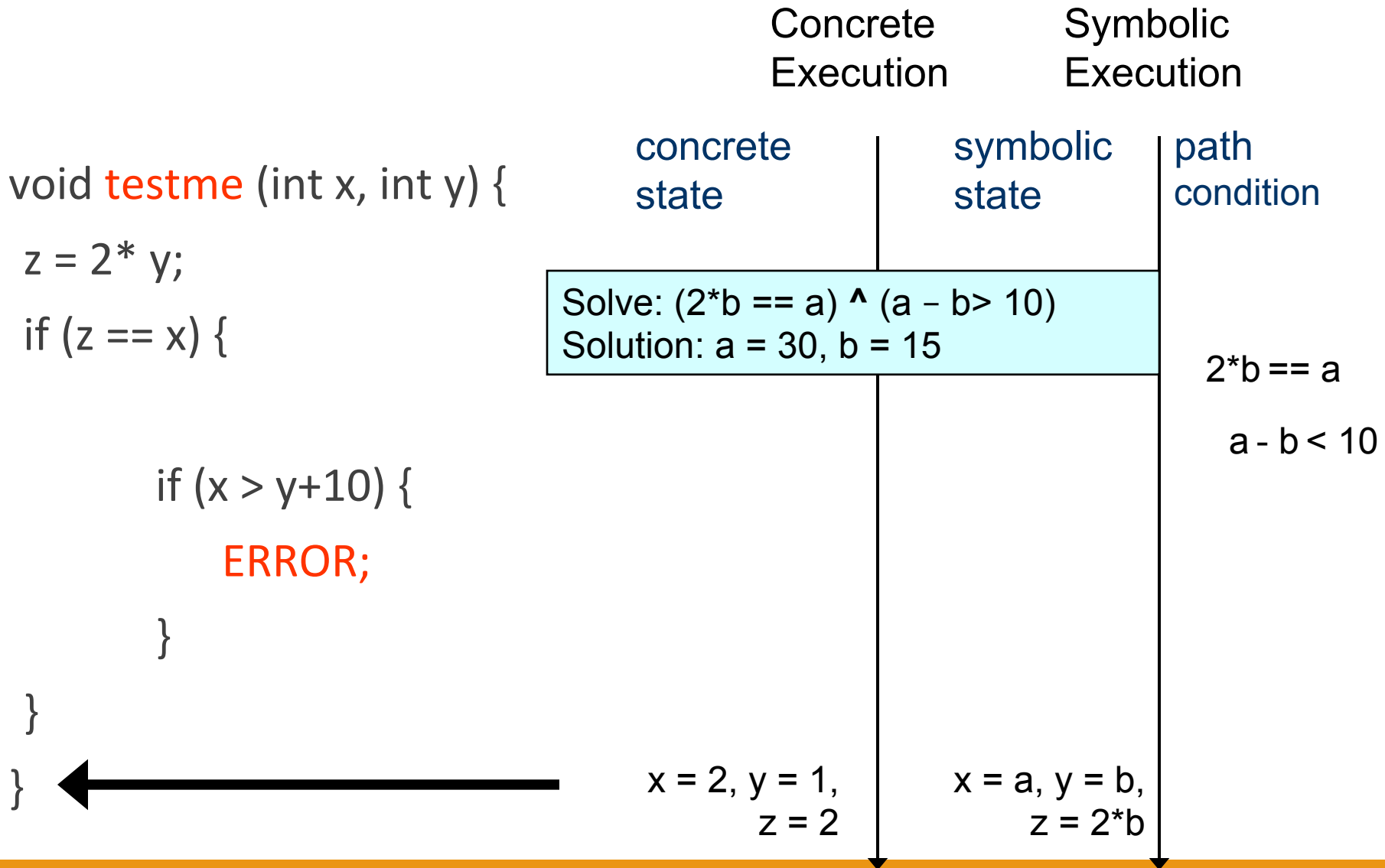




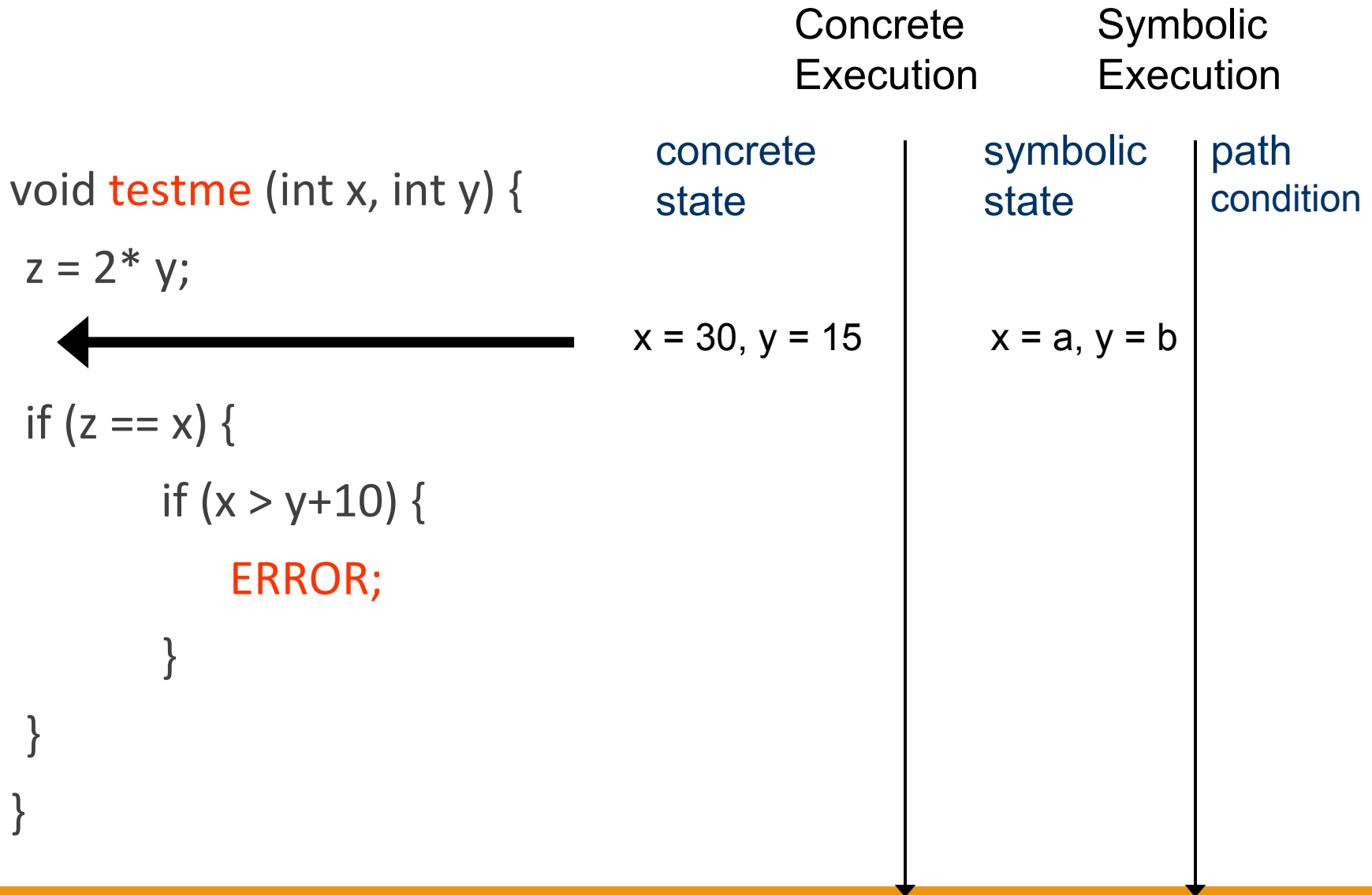
# Concolic execution example



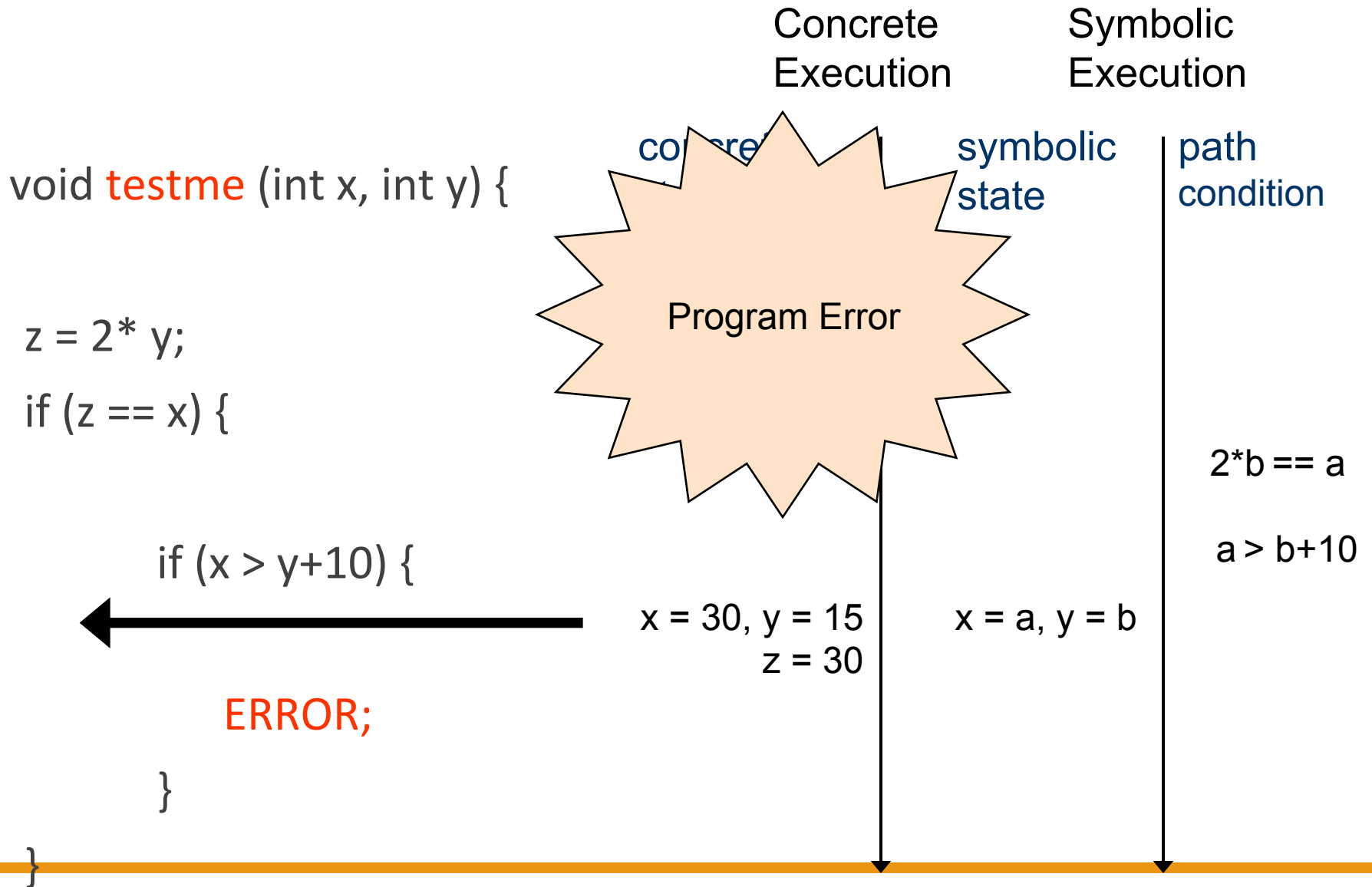
# Concolic execution example



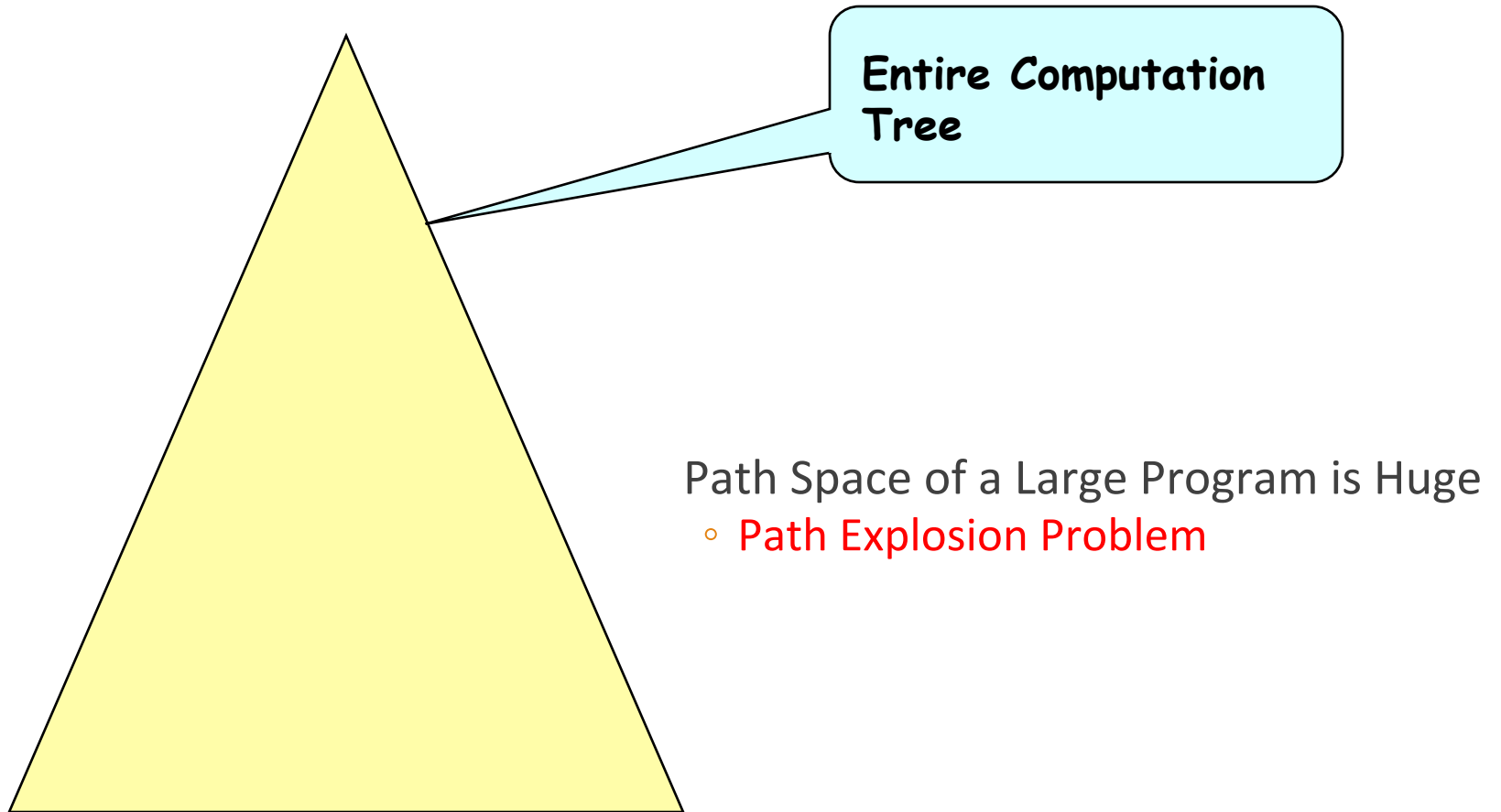
# Concolic execution example



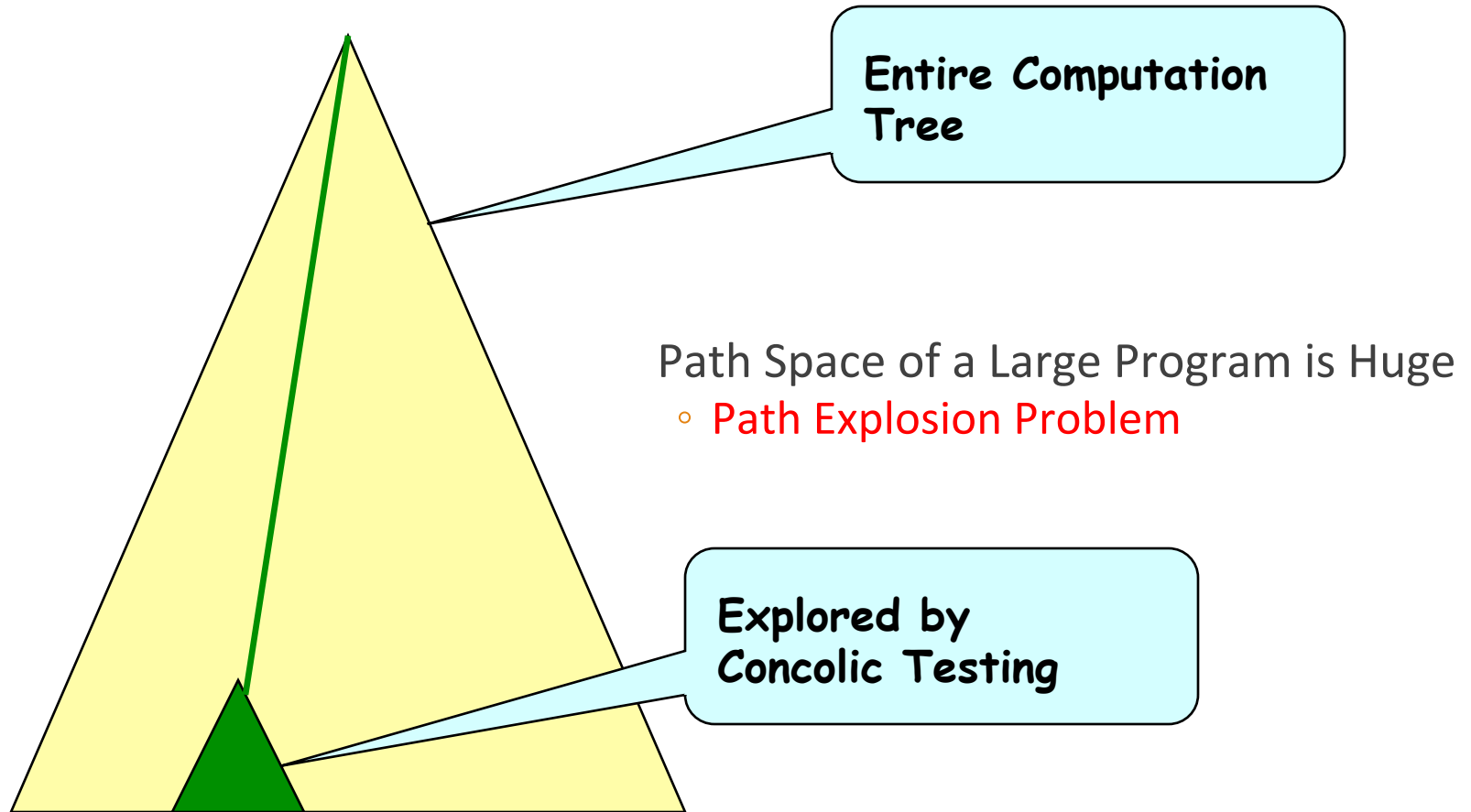
# Concolic execution example



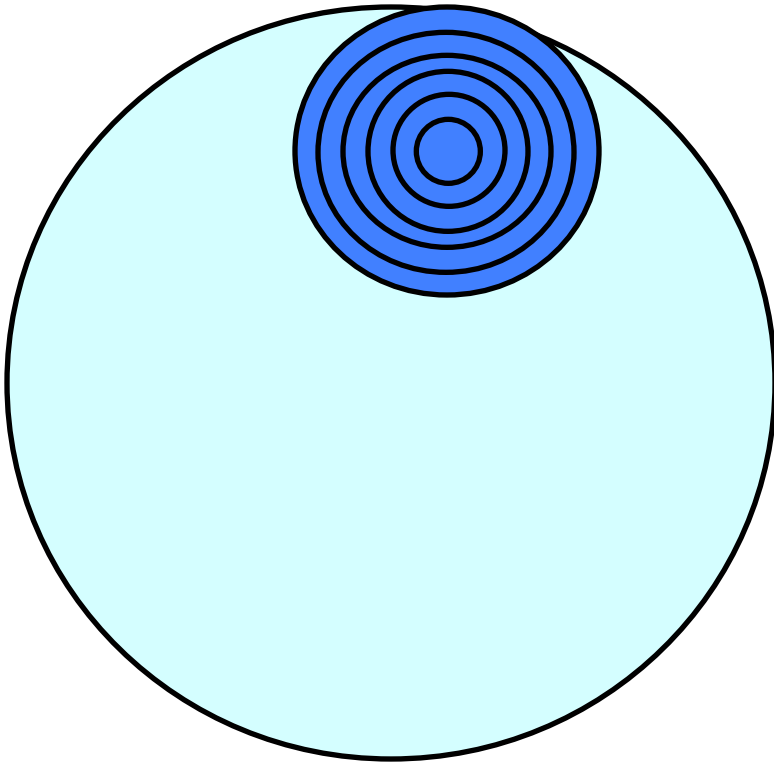
# Limitations



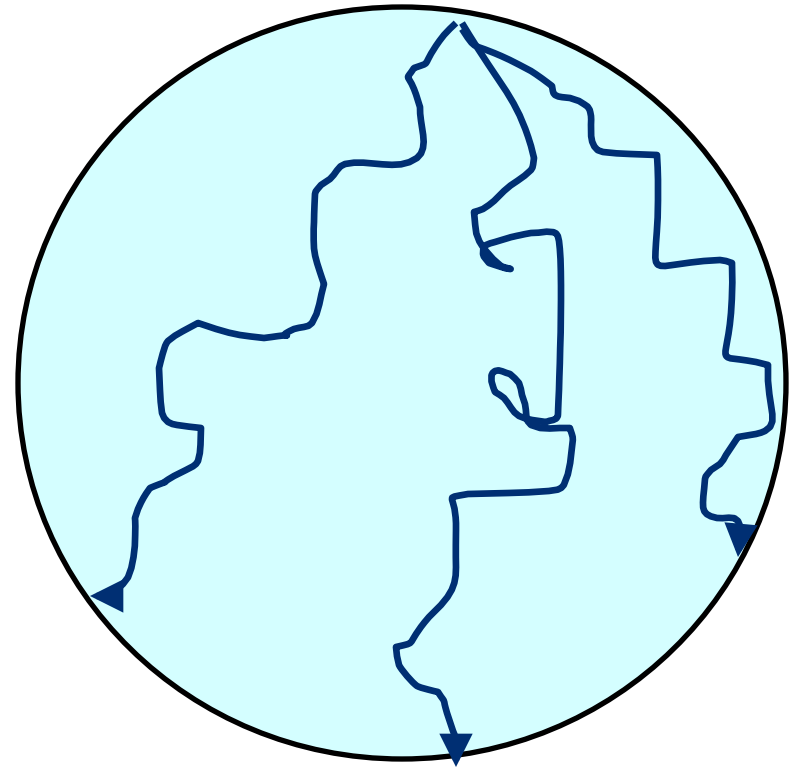
# Limitations



# Limitations: a comparative view



**Concolic: Broad, shallow**



**Random: Narrow, deep**

# Limitations: Example

```
Example ( ) {  
1: state = 0;  
2: while(1) {  
3:  s = input();  
4:  c = input();  
5:  if(c=='.' && state==0)  
        state=1;  
6:  else if(c=='\n' && state==1)  
        state=2;  
7:  else if (s[0]=='I' &&  
        s[1]=='C' &&  
        s[2]=='S' &&  
        s[3]=='E' &&  
        state==2) {  
                                COVER_ME;  
        }  
    }  
}
```

- Want to hit **COVER\_ME**
- **input()** denotes external input
- Can be hit on an input sequence  
    s = "ICSE"  
    c : '.' '\n'

Similar code in

- Text editors (vi)
- Parsers (lexer)
- Event-driven programs (GUI)



# Limitations: Example

```
Example ( ) {  
1: state = 0;  
2: while(1) {  
3:  s = input();  
4:  c = input();  
5:  if(c=='.' && state==0)  
        state=1;  
6:  else if(c=='\n' && state==1)  
        state=2;  
7:  else if (s[0]=='I' &&  
        s[1]=='C' &&  
        s[2]=='S' &&  
        s[3]=='E' &&  
        state==2) {  
                                COVER_ME;  
        }  
    }  
}
```

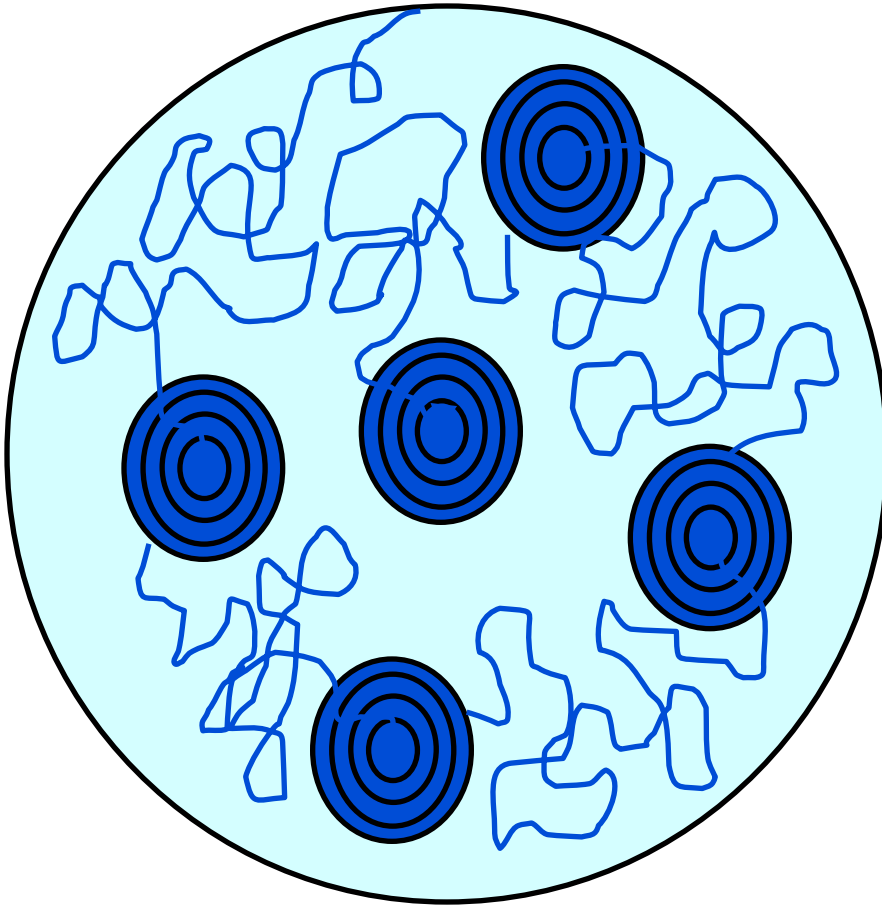
- **Pure random testing** can get to state = 2

But difficult to get 'ICSE' as a Sequence

Probability  $1/(2^8)^6 \gg 3 \times 10^{-15}$

- Conversely, **concolic testing** can generate 'ICSE' but explores many paths to get to state = 2

# Hybrid concolic testing



```
while (not required coverage) {
```

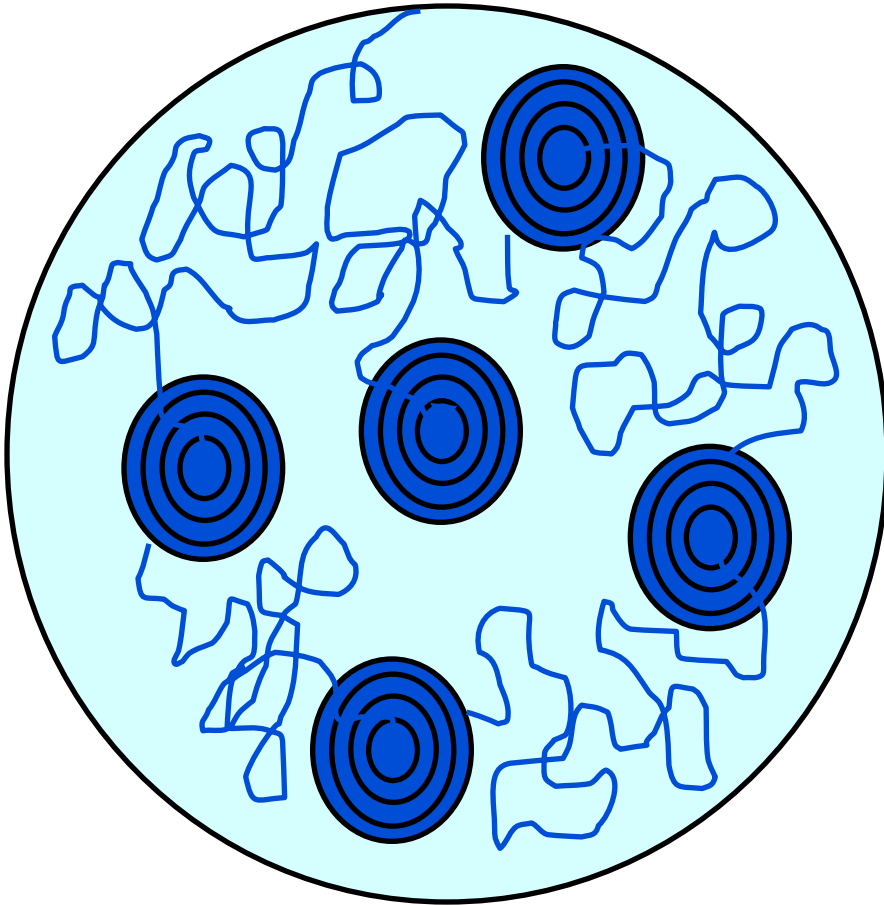
```
    while (not saturation)  
        perform random testing;
```

```
    Checkpoint;  
    while (not increase in coverage)  
        perform concolic testing;  
    Restore;
```

```
}
```

**Interleave Random Testing and  
Concolic Testing to increase coverage**

# Hybrid Concolic Testing



```
while (not required coverage) {
```

```
    while (not saturation)  
        perform random testing;
```

```
    Checkpoint;  
    while (not increase in coverage)  
        perform concolic testing;  
    Restore;
```

```
} Interleave Random Testing and  
   Concolic Testing to increase coverage
```

Deep, broad search  
**Hybrid** Search

# Hybrid Concolic Testing

```
Example ( ) {  
  1: state = 0;  
  2: while(1) {  
  3:   s = input();  
  4:   c = input();  
  5:   if(c=='.' && state==0)  
        state=1;  
  6:   else if(c=='\n' && state==1)  
        state=2;  
  7:   else if (s[0]=='I' &&  
        s[1]=='C' &&  
        s[2]=='S' &&  
        s[3]=='E' &&  
        state==2) {  
                                COVER_ME::  
        }  
  }  
}
```

## Random Phase

- '\$', '&', '-', '6', ':', '%', '^', '\n', 'x', '~' ...
- Saturates after many (~10000) iterations
- In less than 1 second
- **COVER\_ME** is not reached

# Hybrid Concolic Testing

```
Example ( ) {  
  1: state = 0;  
  2: while(1) {  
  3:  s = input();  
  4:  c = input();  
  5:  if(c=='.' && state==0)  
      state=1;  
  6:  else if(c=='\n' && state==1)  
      state=2;  
  7:  else if (s[0]=='I' &&  
      s[1]=='C' &&  
      s[2]=='S' &&  
      s[3]=='E' &&  
      state==2) {  
      COVER_ME;;  
  }  
  }  
}
```

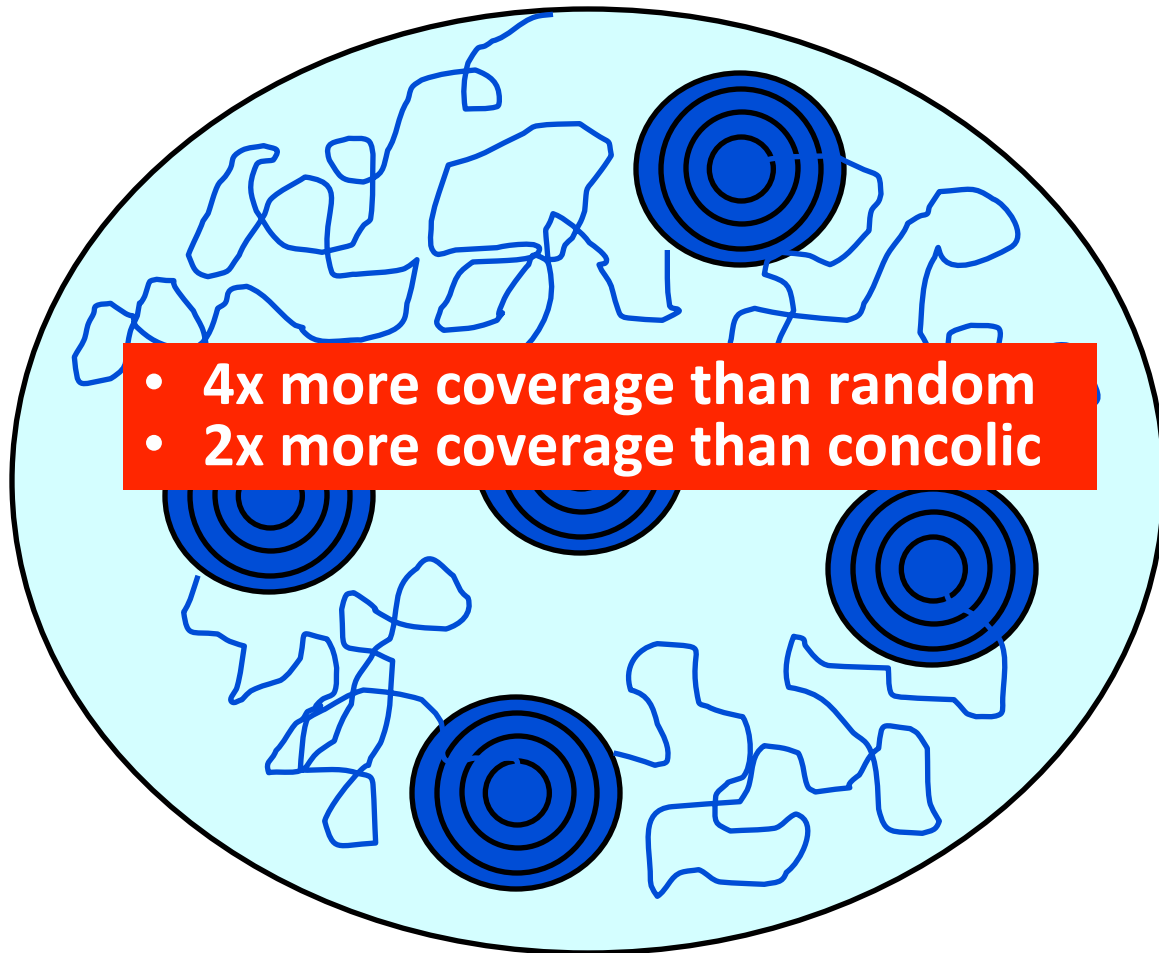
## Random Phase

- '\$', '&', '-', '6', ':', '%', '^', '\n', 'x', '~' ...
- Saturates after many (~10000) iterations
- In less than 1 second
- **COVER\_ME** is not reached

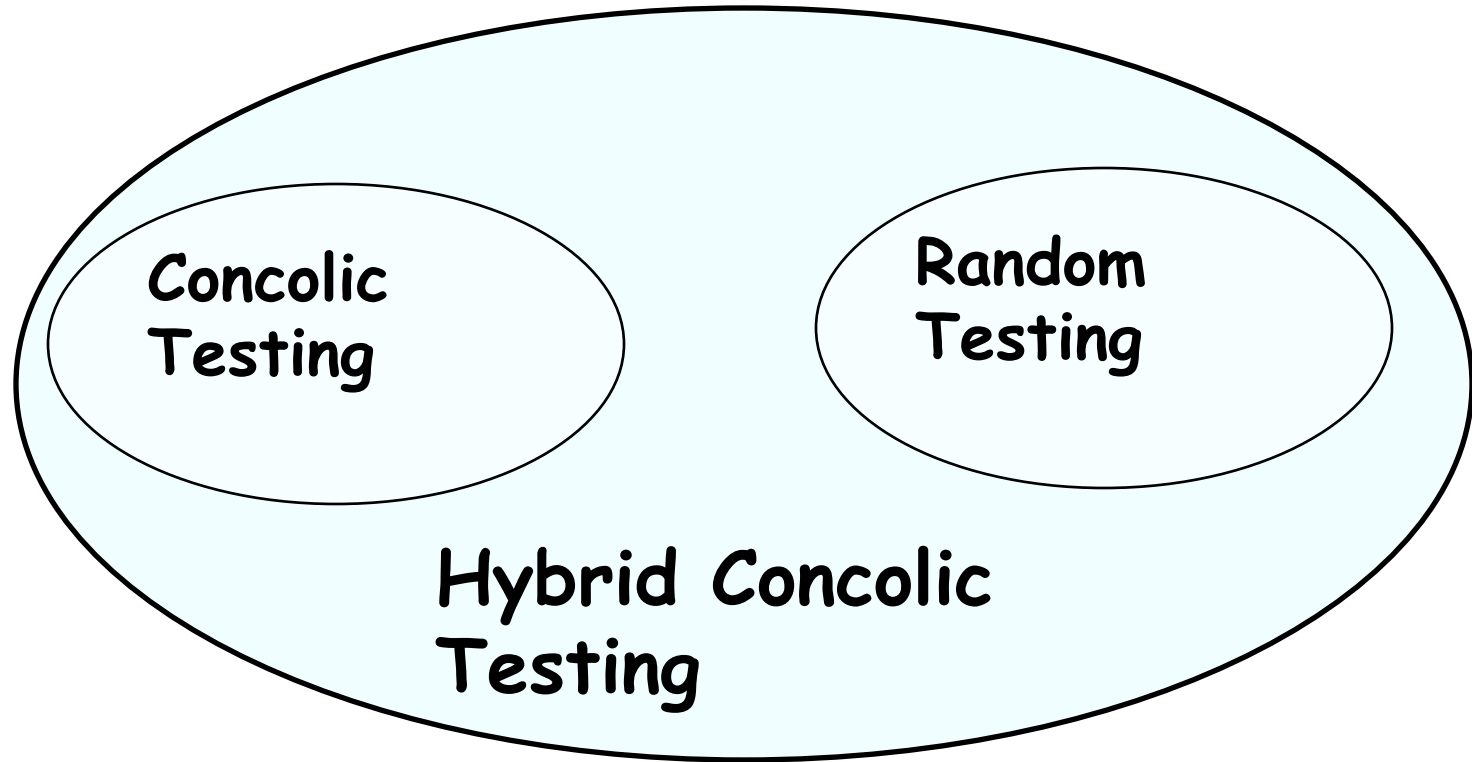
## Concolic Phase

- s[0]='I', s[1]='C', s[2]='S', s[3]='E'
- Reaches **COVER\_ME**

# Hybrid Concolic Testing



# Summary



# Further reading

[Symbolic execution and program testing](#) - James King

[KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs](#) - Cadar et. al.

[Symbolic Execution for Software Testing: Three Decades Later](#) - Cadar and Sen

[DART: Directed Automated Random Testing](#) - Godefroid et. al.

[CUTE: A Concolic Unit Testing Engine for C](#) - Sen et. al.

---