

# Invariant Generation via Genetic Programming

Soonho Kong

`soon@ropas.snu.ac.kr`

12 Mar 2010

ROPAS Show&Tell

# Problem:

## Inductive Invariant Generation

For the annotated loop

$\{\delta\}$  while  $\rho$  do  $S$  end  $\{\epsilon\}$

Find an invariant  $I$  satisfying the following conditions:

- (A)  $\delta \Rightarrow \iota$  (  $\iota$  holds when entering the loop)
- (B)  $\iota \wedge \rho \Rightarrow Pre(\iota, S)$  (  $\iota$  holds at each iteration)
- (C)  $\iota \wedge \neg\rho \Rightarrow \epsilon$  (  $\iota$  gives  $\epsilon$  after leaving the loop)

$$\underbrace{\delta}_{\bar{l} \text{ strongest under-approximation of an invariant}} \Rightarrow I \Rightarrow \underbrace{\epsilon \vee \rho}_{\bar{l} \text{ weakest over-approximation of an invariant}}$$

Precondition

$\{ \text{phase} = \text{F} \wedge \text{success} = \text{F} \wedge \text{give\_up} = \text{F} \wedge \text{cutoff} = 0 \wedge \text{count} = 0 \}$

Loop Body

```
1 while  $\neg(\text{success} \vee \text{give\_up})$  do
2    $\text{entered\_phase} := \text{F};$ 
3   if  $\neg\text{phase}$  then
4     if  $\text{cutoff} = 0$  then  $\text{cutoff} := 1;$ 
5     else if  $\text{cutoff} = 1 \wedge \text{maxcost} > 1$  then  $\text{cutoff} := \text{maxcost};$ 
6       else  $\text{phase} := \text{T}; \text{entered\_phase} := \text{T}; \text{cutoff} := 1000;$ 
7     if  $\text{cutoff} = \text{maxcost} \wedge \neg\text{search}$  then  $\text{give\_up} := \text{T};$ 
8   else
9      $\text{count} := \text{count} + 1;$ 
10    if  $\text{count} > \text{words}$  then  $\text{give\_up} := \text{T};$ 
11    if  $\text{entered\_phase}$  then  $\text{count} := 1;$ 
12     $\text{linkages} := \text{nondet};$ 
13    if  $\text{linkages} > 5000$  then  $\text{linkages} := 5000;$ 
14     $\text{canonical} := 0; \text{valid} := 0;$ 
15    if  $\text{linkages} \neq 0$  then
16       $\text{valid} := \text{nondet};$  assume  $0 \leq \text{valid} \wedge \text{valid} \leq \text{linkages};$ 
17       $\text{canonical} := \text{linkages};$ 
18    if  $\text{valid} > 0$  then  $\text{success} := \text{T};$ 
19 end
```

Postcondition

$\{ (\text{valid} > 0 \vee \text{count} > \text{words} \vee (\text{cutoff} = \text{maxcost} \wedge \neg\text{search})) \wedge$   
 $\text{valid} \leq \text{linkages} \wedge \text{canonical} = \text{linkages} \wedge \text{linkages} \leq 5000 \}$

**Fig. 3.** A Sample Loop in SPEC2000 Benchmark PARSER  
with 20 Atomic Propositions (Building Blocks)

```
{ phase = F ∧ success = F ∧ give_up = F ∧ cutoff = 0 ∧ count = 0 }
```

```
1 while ¬(success ∨ give_up) do
2   entered_phase := F;
3   if ¬phase then
4     if cutoff = 0 then cutoff := 1;
5     else if cutoff = 1 then
```

An Invariant ✓

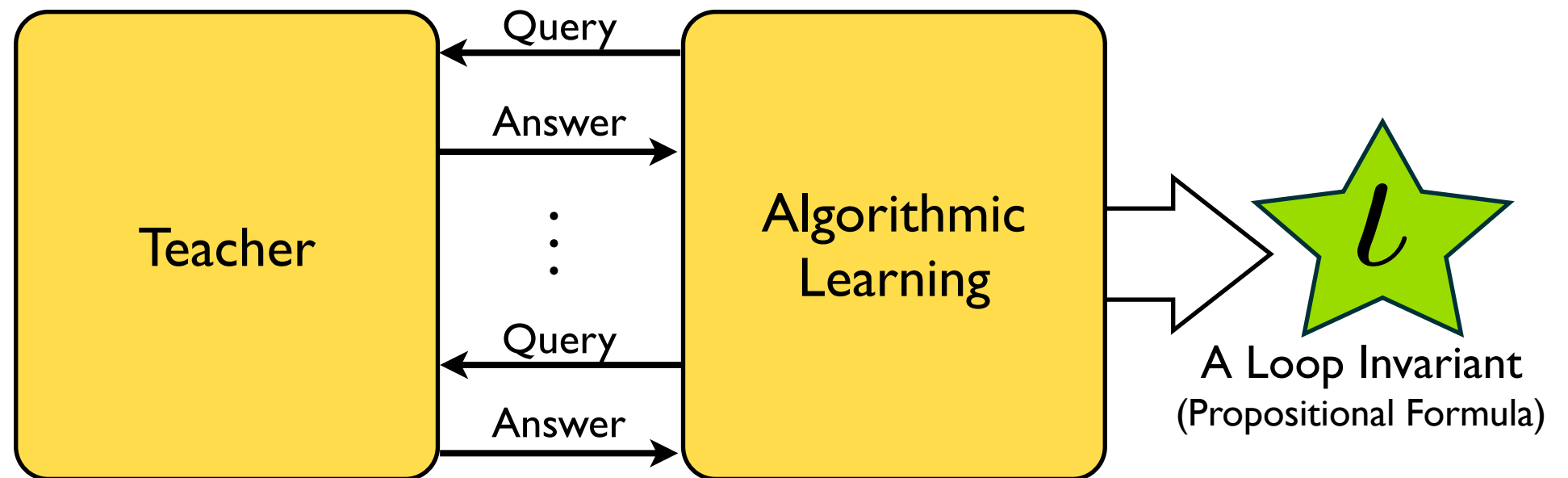
$success \Rightarrow (valid \leq linkages \wedge linkages \leq 5000 \wedge canonical = linkages) \wedge$   
 $success \Rightarrow (\neg search \vee count > words \vee valid \neq 0) \wedge$   
 $success \Rightarrow (count > words \vee cutoff = maxcost \vee (canonical \neq 0 \wedge valid \neq 0 \wedge linkages \neq 0)) \wedge$   
 $give\_up \Rightarrow ((valid = 0 \wedge linkages = 0 \wedge canonical = linkages) \vee$   
 $(canonical \neq 0 \wedge valid \leq linkages \wedge linkages \leq 5000 \wedge canonical = linkages)) \wedge$   
 $give\_up \Rightarrow (cutoff = maxcost \vee count > words \vee$   
 $(canonical \neq 0 \wedge valid \neq 0 \wedge linkages \neq 0)) \wedge$   
 $give\_up \Rightarrow (\neg search \vee count > words \vee valid \neq 0)$

```
16   canonical := nondet; assume 0 ≤ valid ∧ valid ≤ linkages;
17   canonical := linkages;
18   if valid > 0 then success := T;
19 end
```

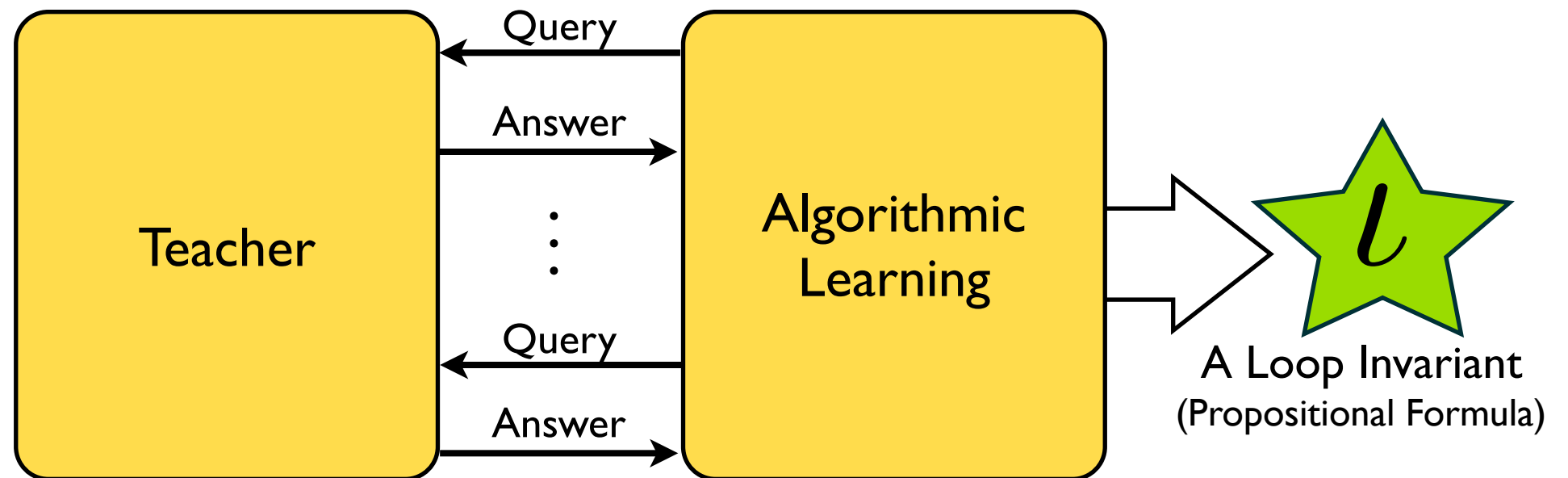
```
{ (valid > 0 ∨ count > words ∨ (cutoff = maxcost ∧ ¬search)) ∧
  valid ≤ linkages ∧ canonical = linkages ∧ linkages ≤ 5000 }
```

**Fig. 3.** A Sample Loop in SPEC2000 Benchmark PARSER  
with 20 Atomic Propositions (Building Blocks)

# Observations on Algorithmic Learning Approach



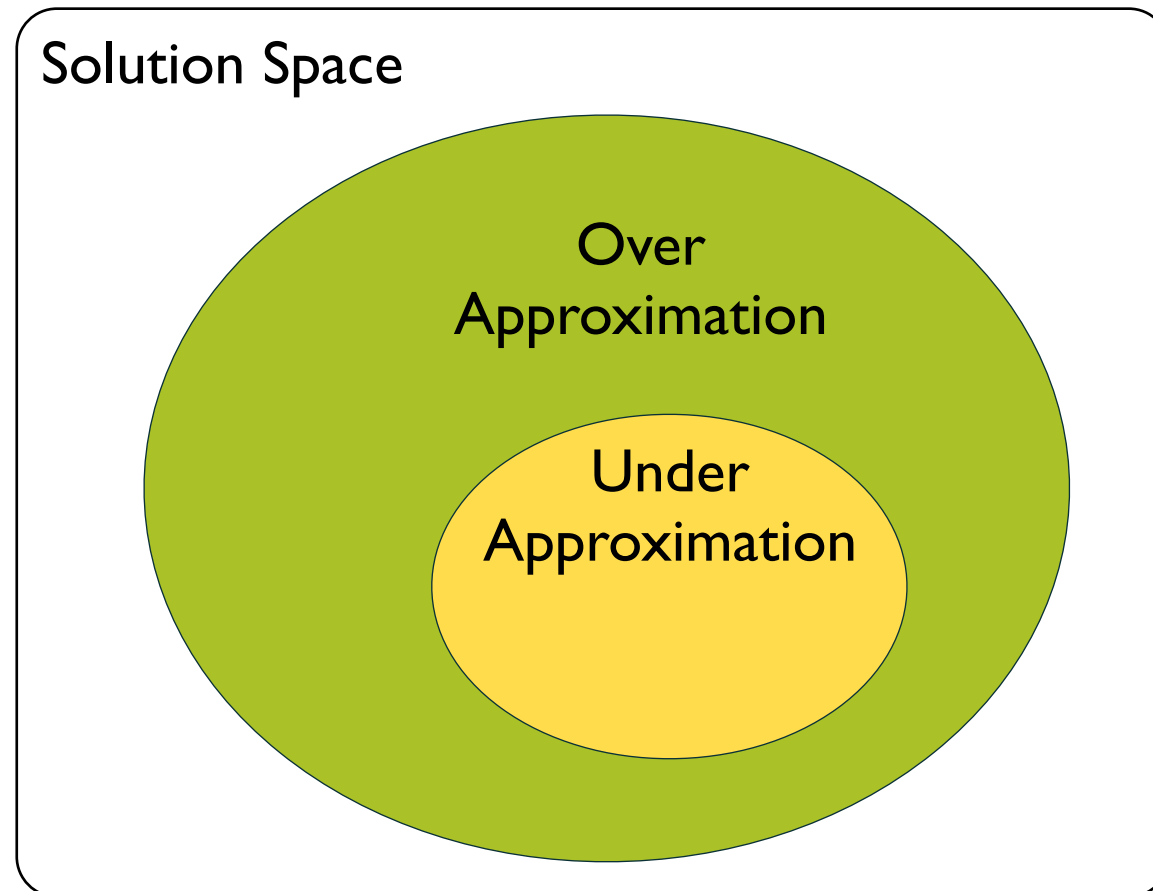
# Observations on Algorithmic Learning Approach



## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.

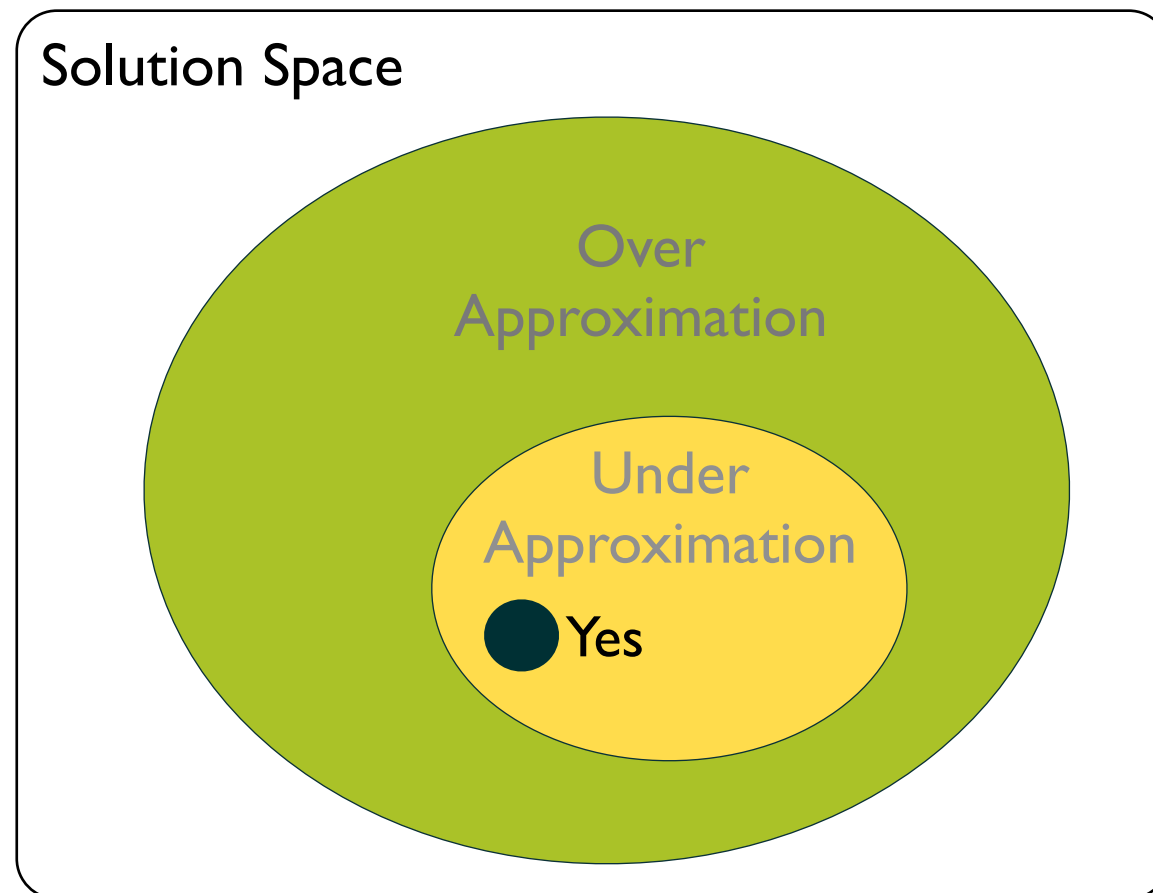
# Observations on Algorithmic Learning Approach



## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.

# Observations on Algorithmic Learning Approach

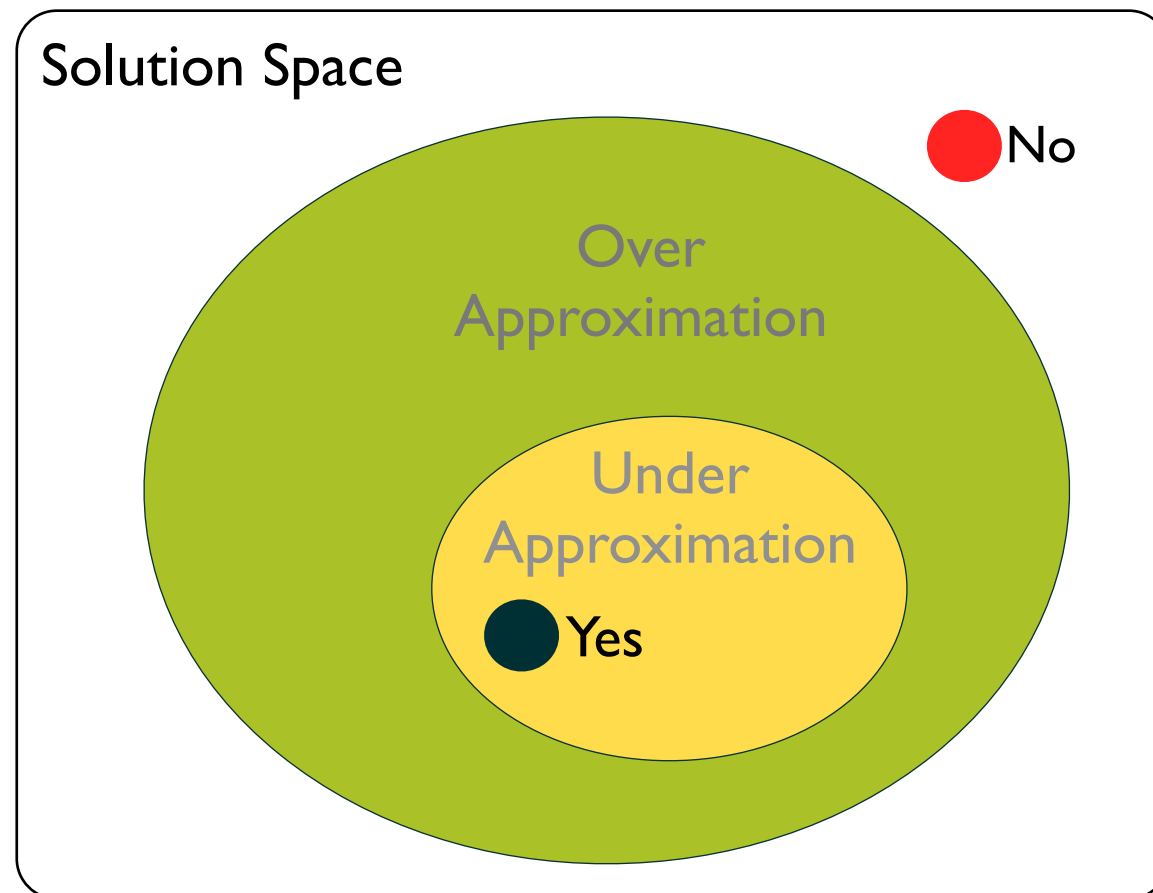


## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.



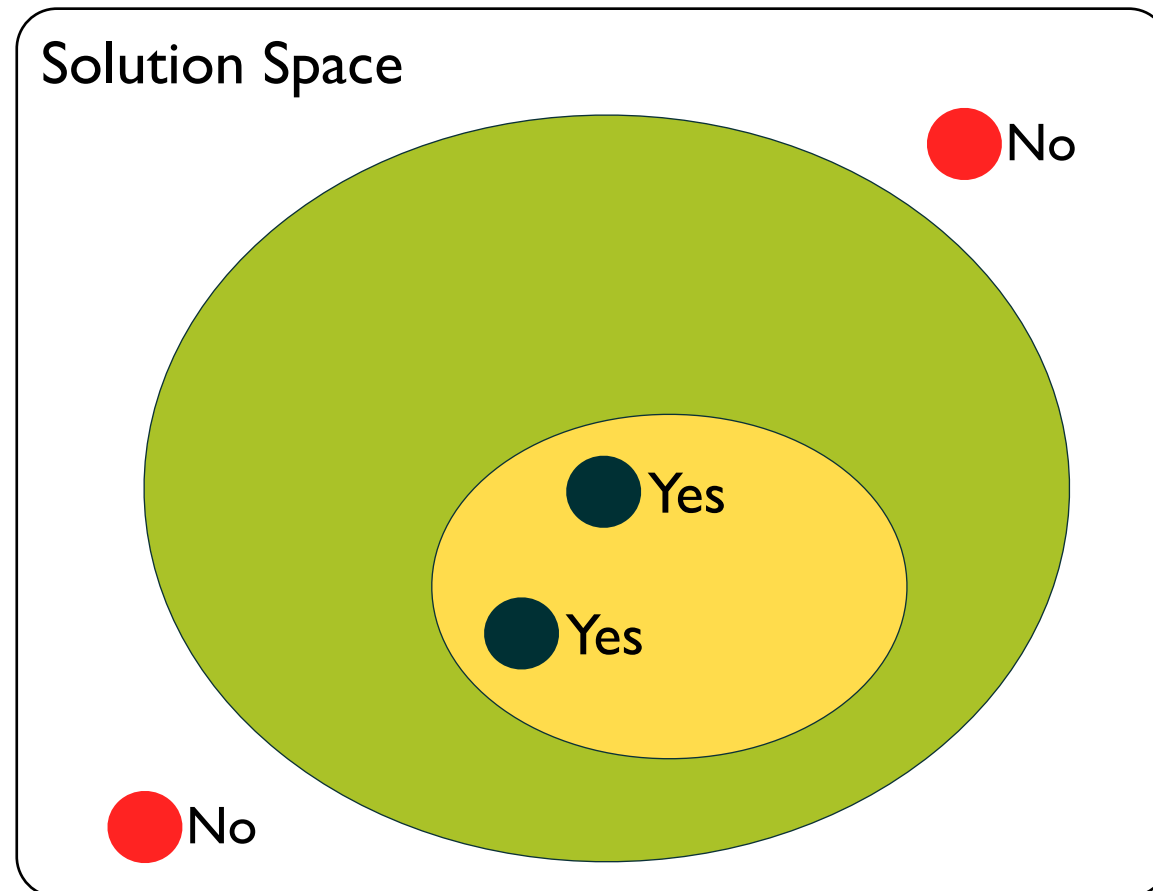
# Observations on Algorithmic Learning Approach



## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.

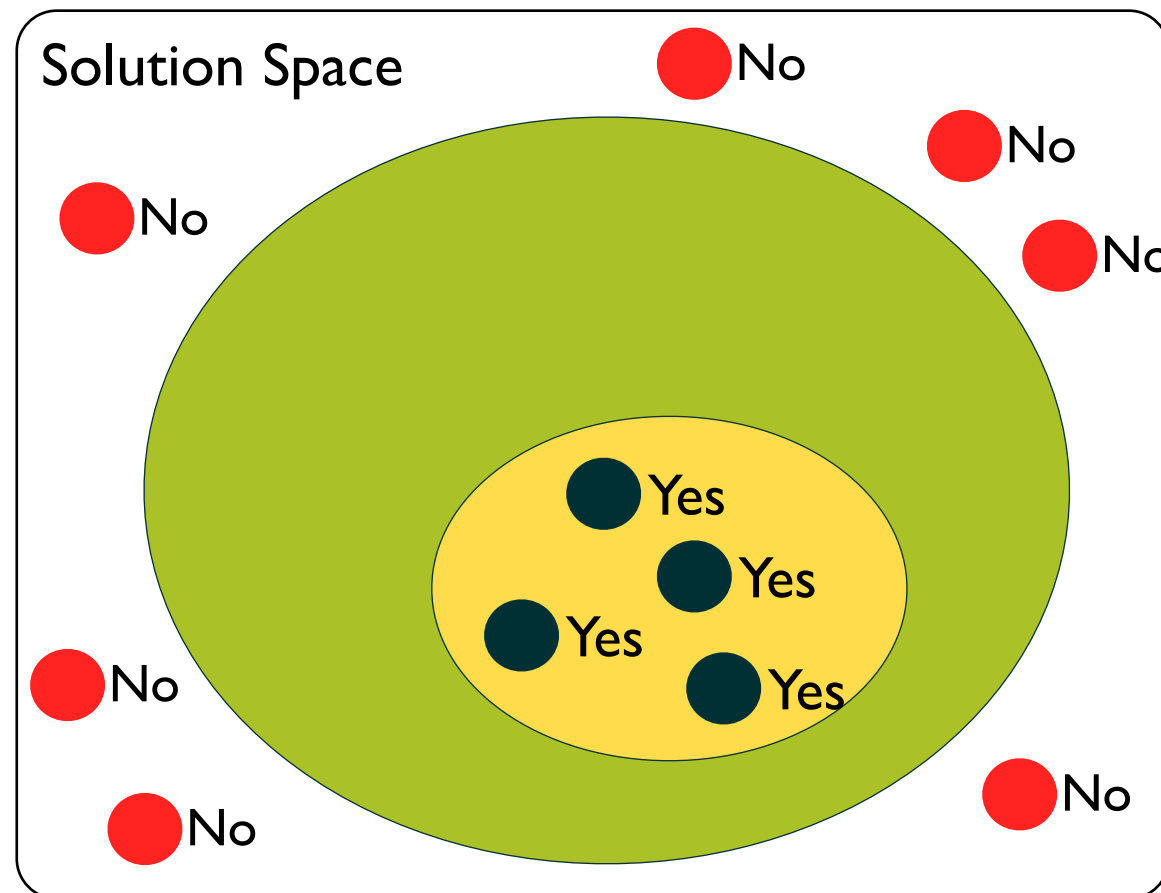
# Observations on Algorithmic Learning Approach



## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.

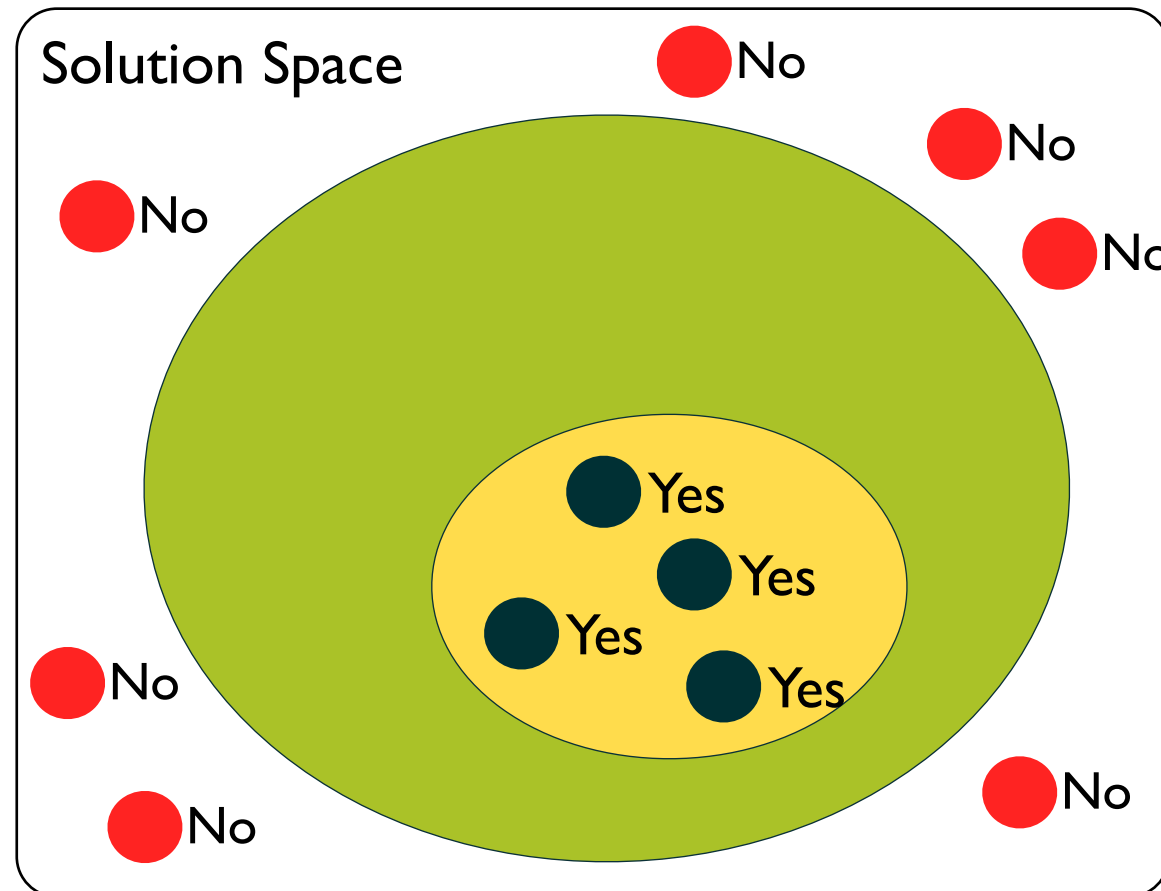
# Observations on Algorithmic Learning Approach



## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.

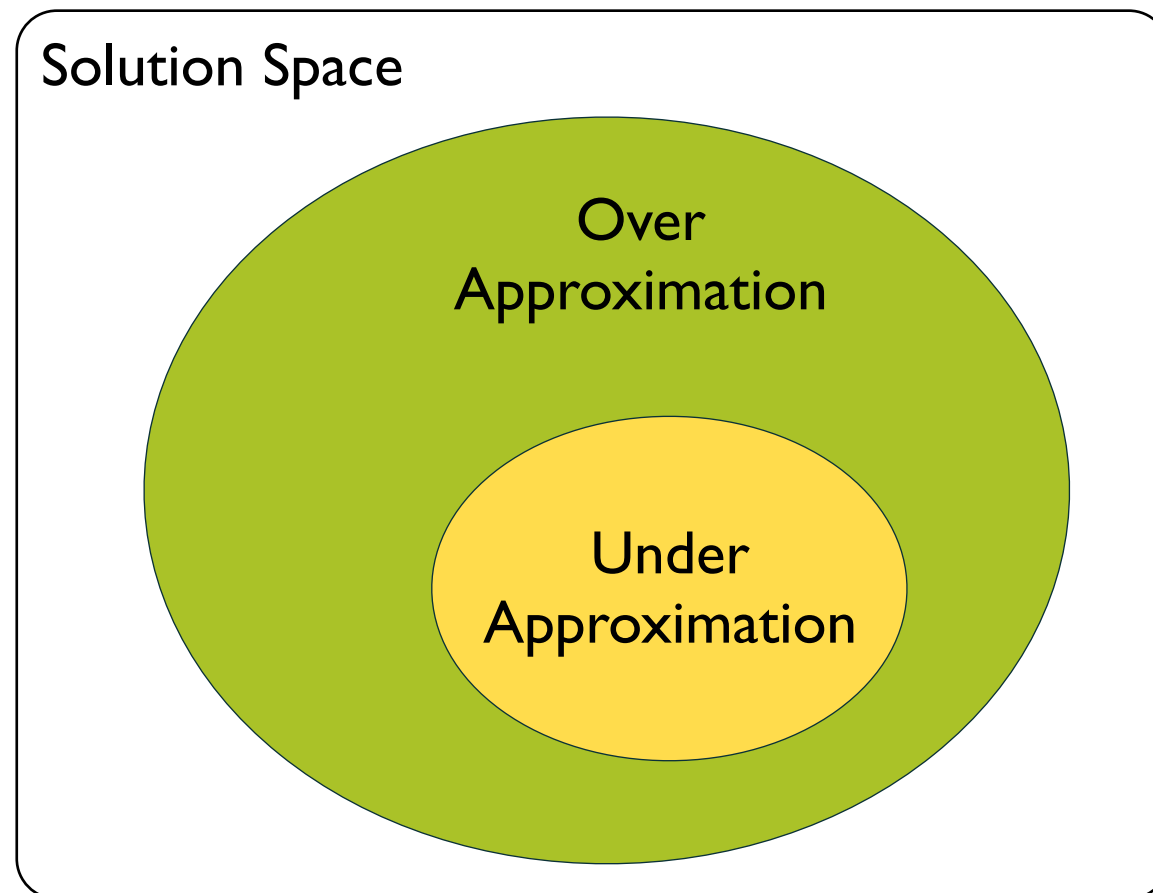
# Observations on Algorithmic Learning Approach



## Problem #1: **Active Learning**

Invariant approximations **are not directly** utilized.

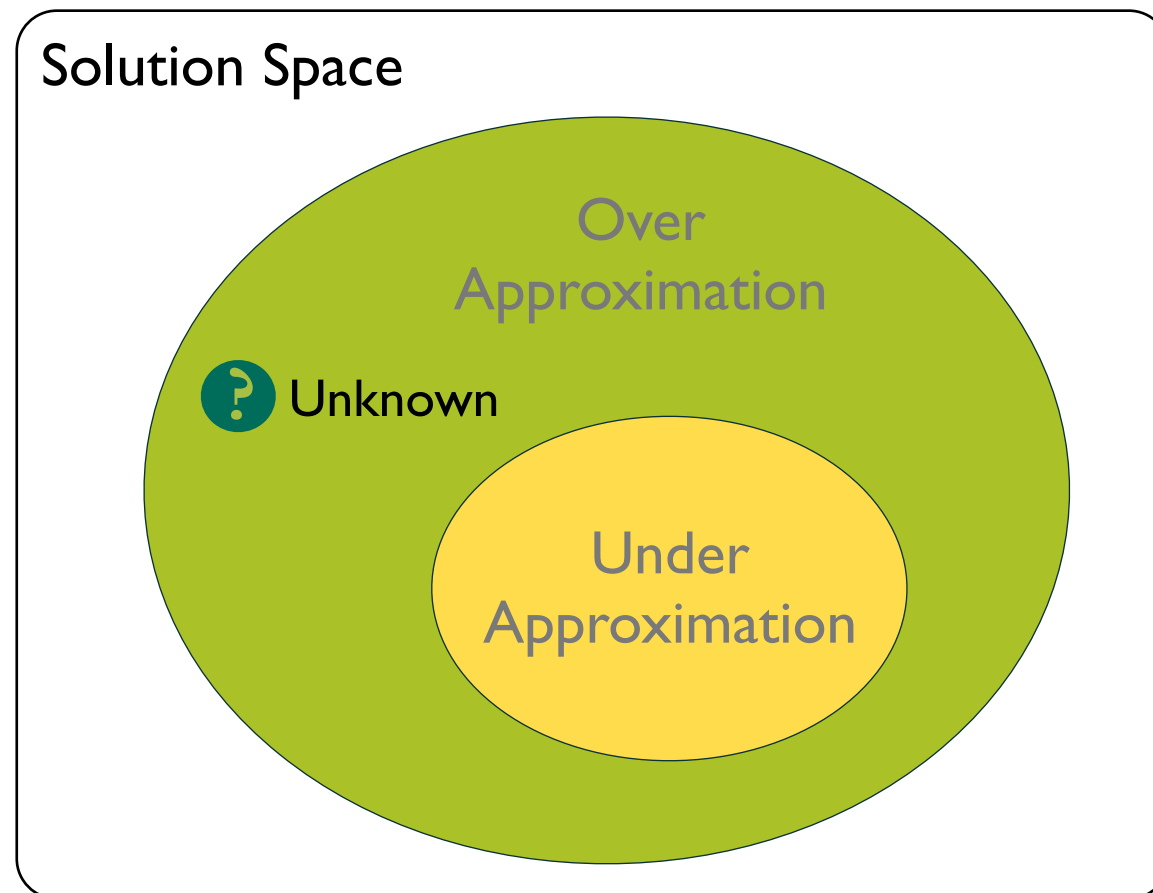
# Observations on Algorithmic Learning Approach



Problem #2: **Naive Randomized Mechanism**

Teacher **solely relies on** naive randomized mechanism when approximations are not helpful.

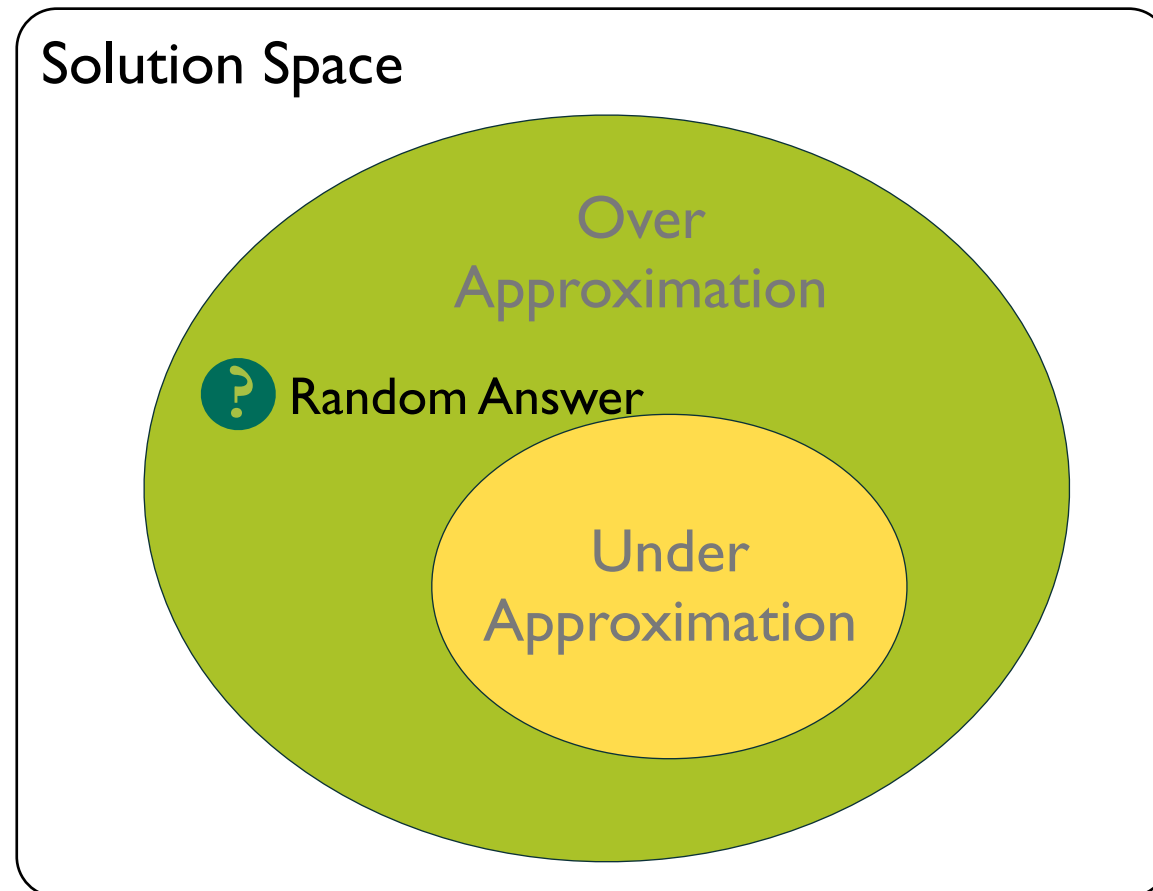
# Observations on Algorithmic Learning Approach



Problem #2: **Naive Randomized Mechanism**

Teacher **solely relies on** naive randomized mechanism when approximations are not helpful.

# Observations on Algorithmic Learning Approach

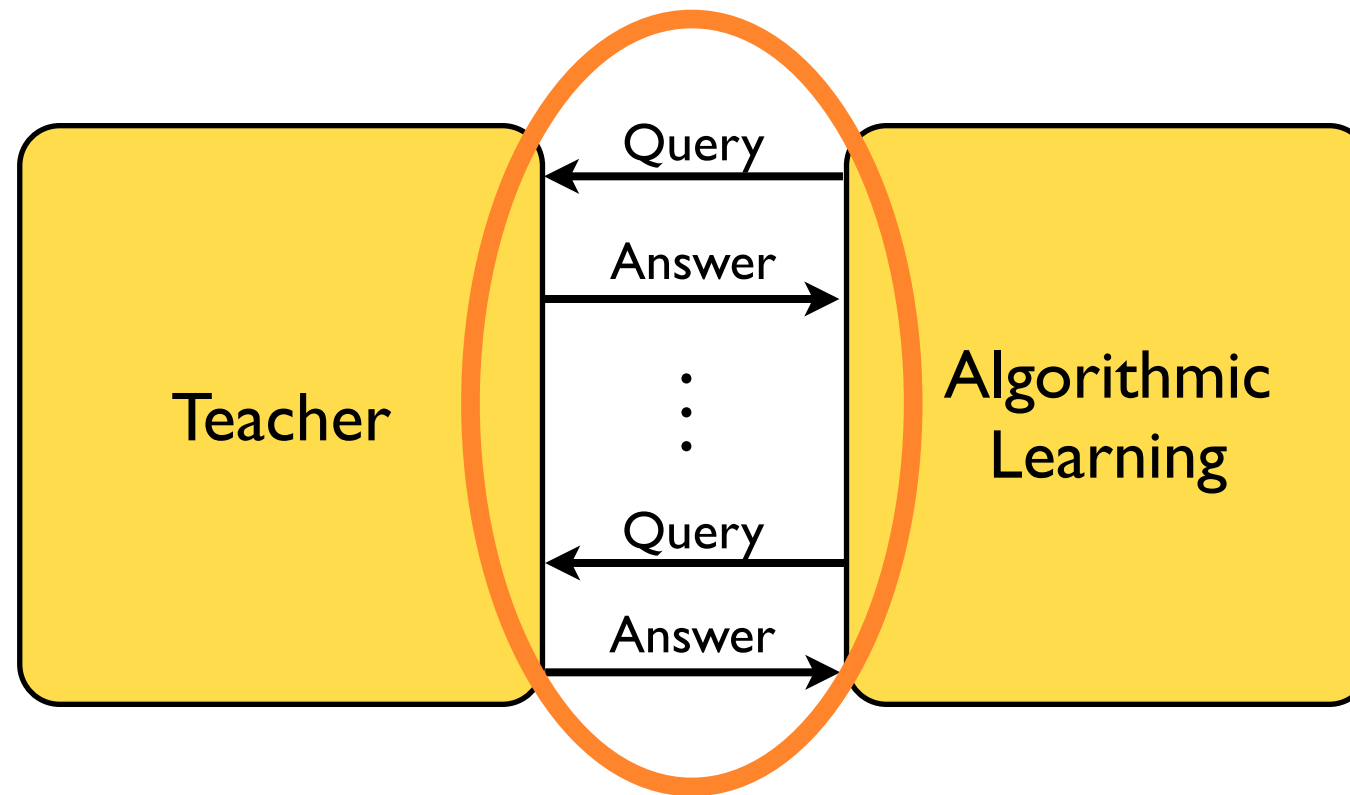


Problem #2: **Naive Randomized Mechanism**

Teacher **solely relies on** naive randomized mechanism when approximations are not helpful.



# Observations on Algorithmic Learning Approach

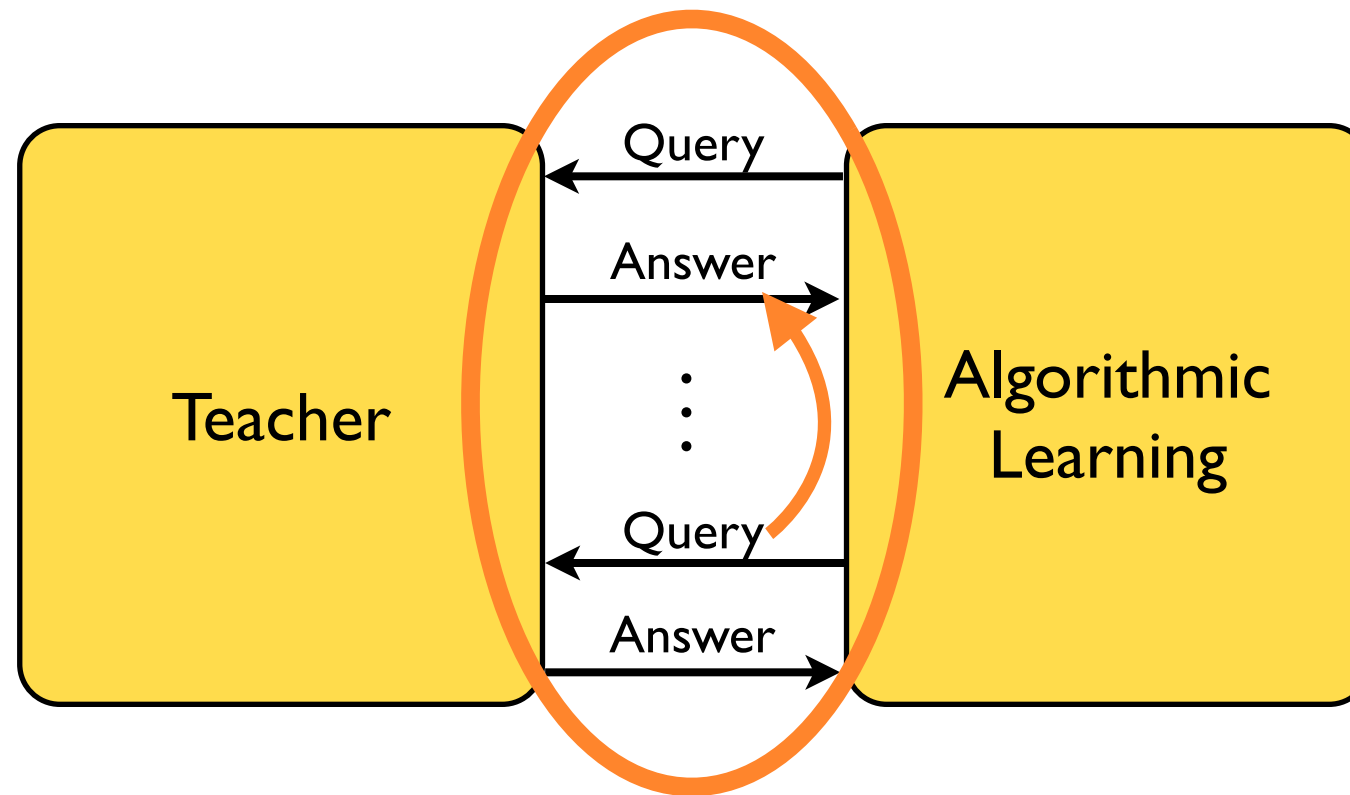


Problem #3: **Serialized Query-Answer Process**

Queries are **dependent** on the previous query/  
answer pairs and make parallelization **difficult**.



# Observations on Algorithmic Learning Approach



Problem #3: **Serialized Query-Answer Process**

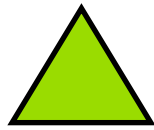
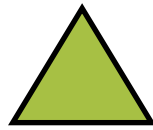
Queries are **dependent** on the previous query/  
answer pairs and make parallelization **difficult**.

# Genetic Programming

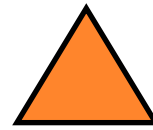
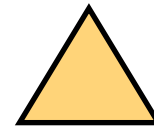


# Genetic Programming

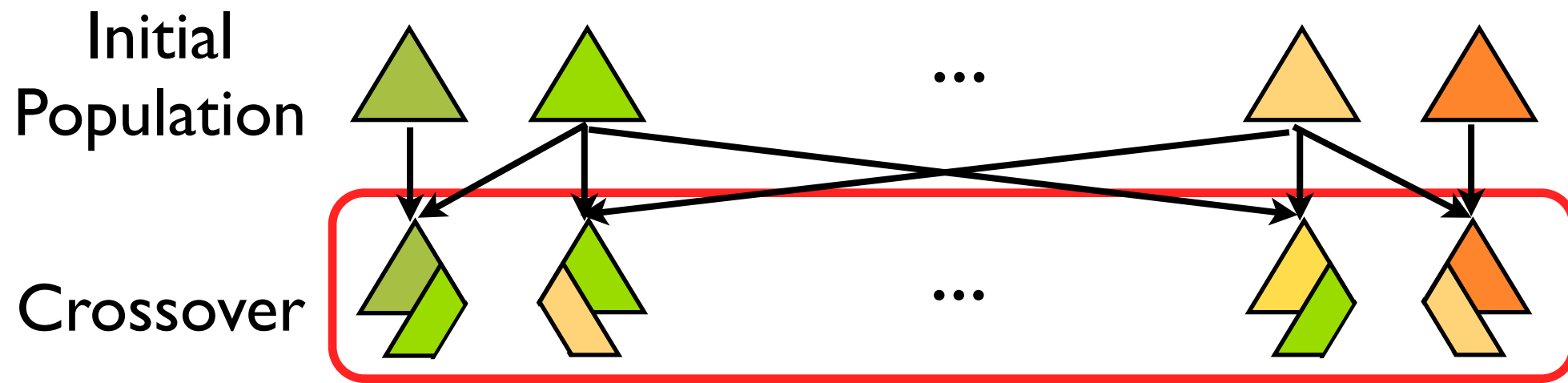
Initial  
Population



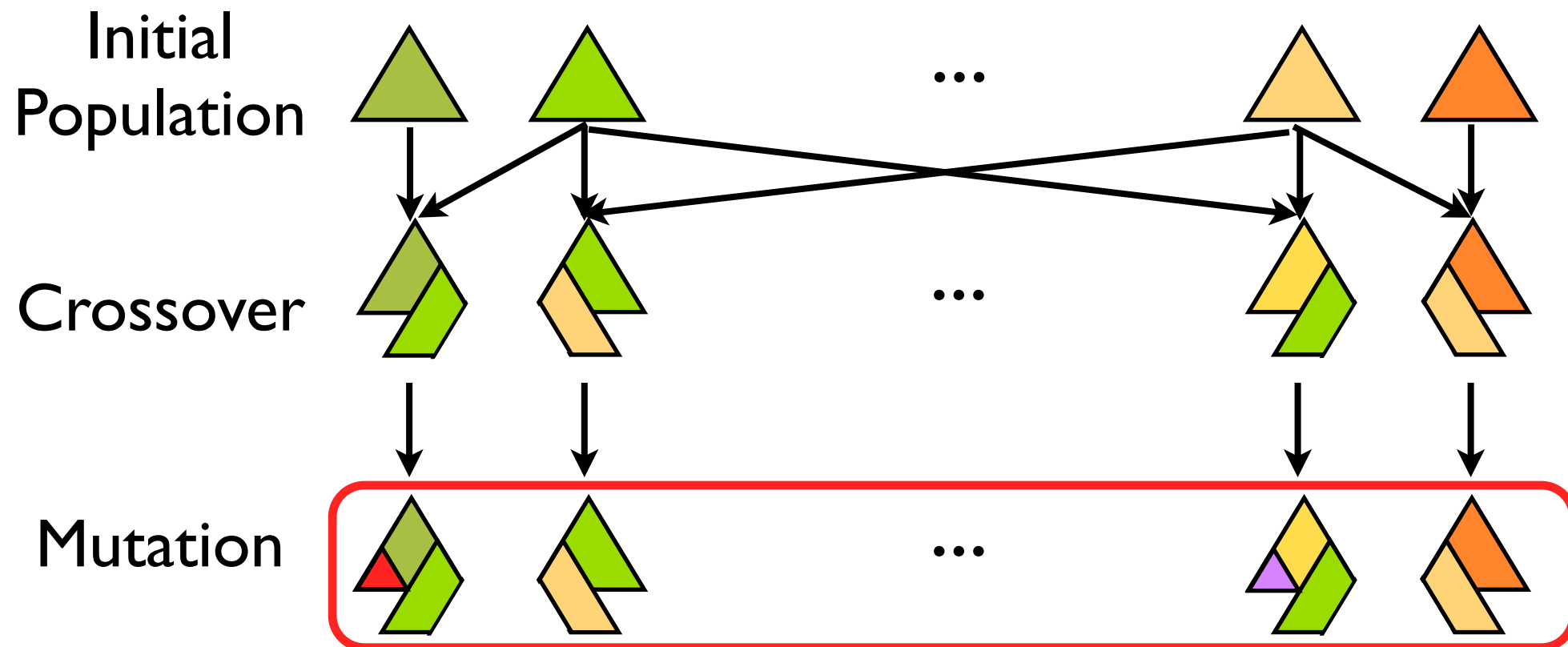
...



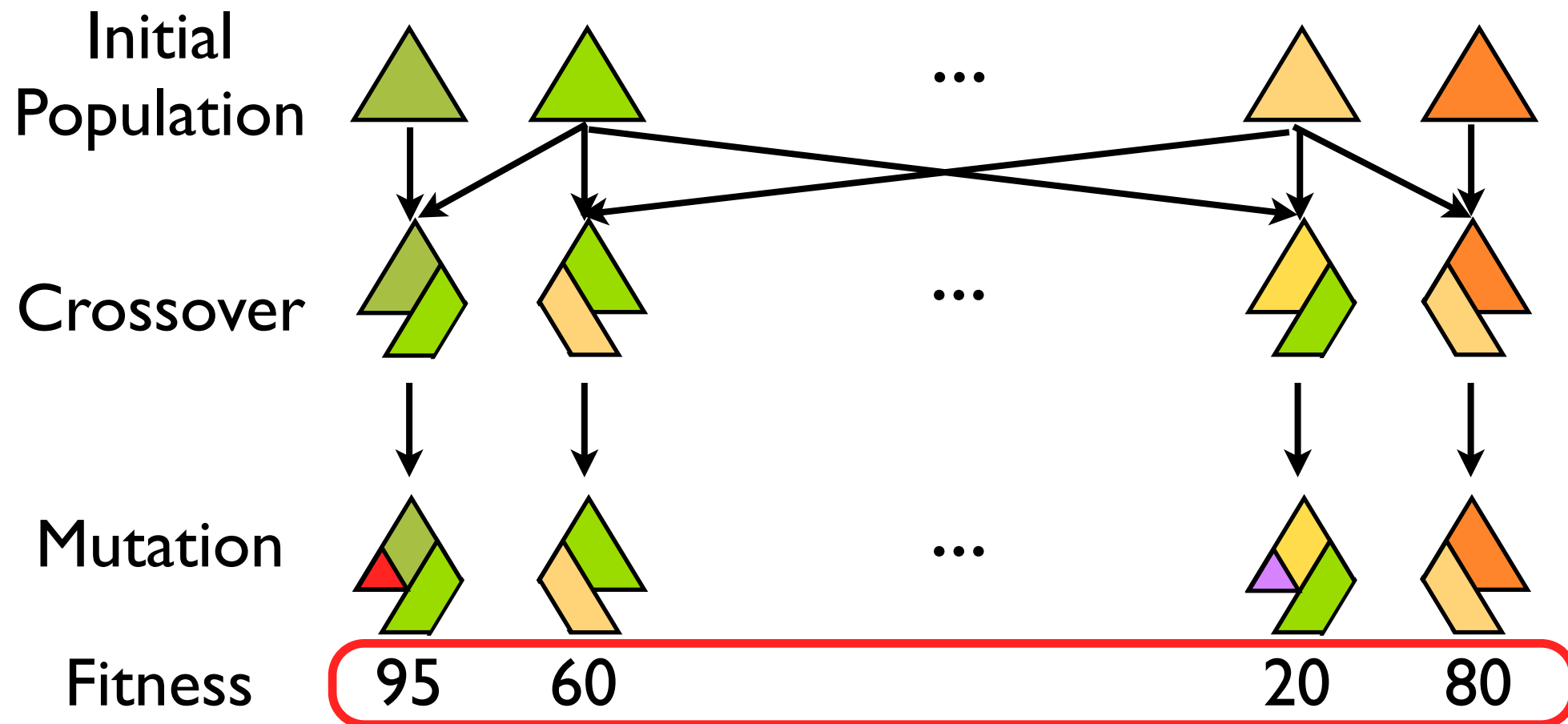
# Genetic Programming



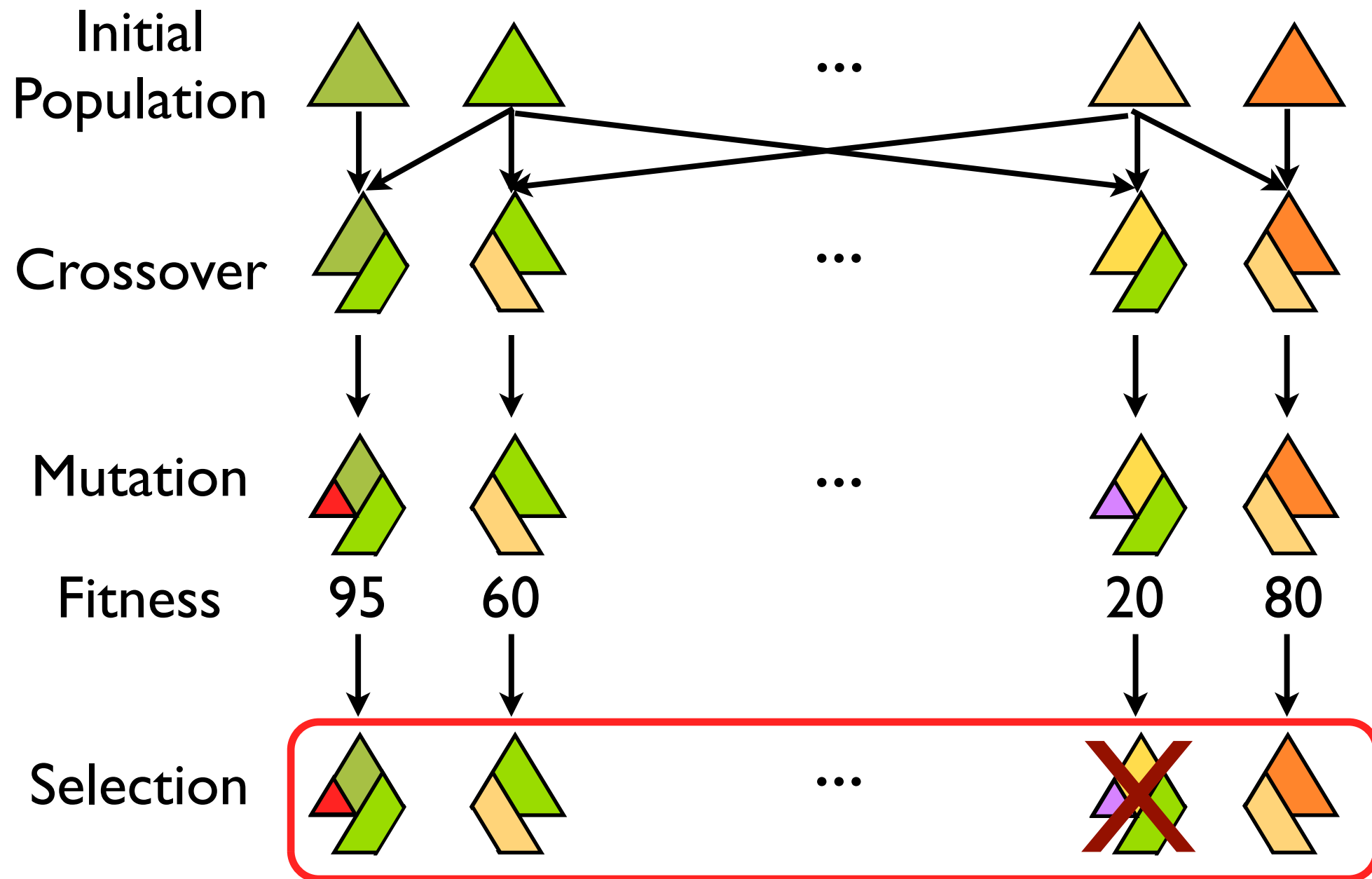
# Genetic Programming



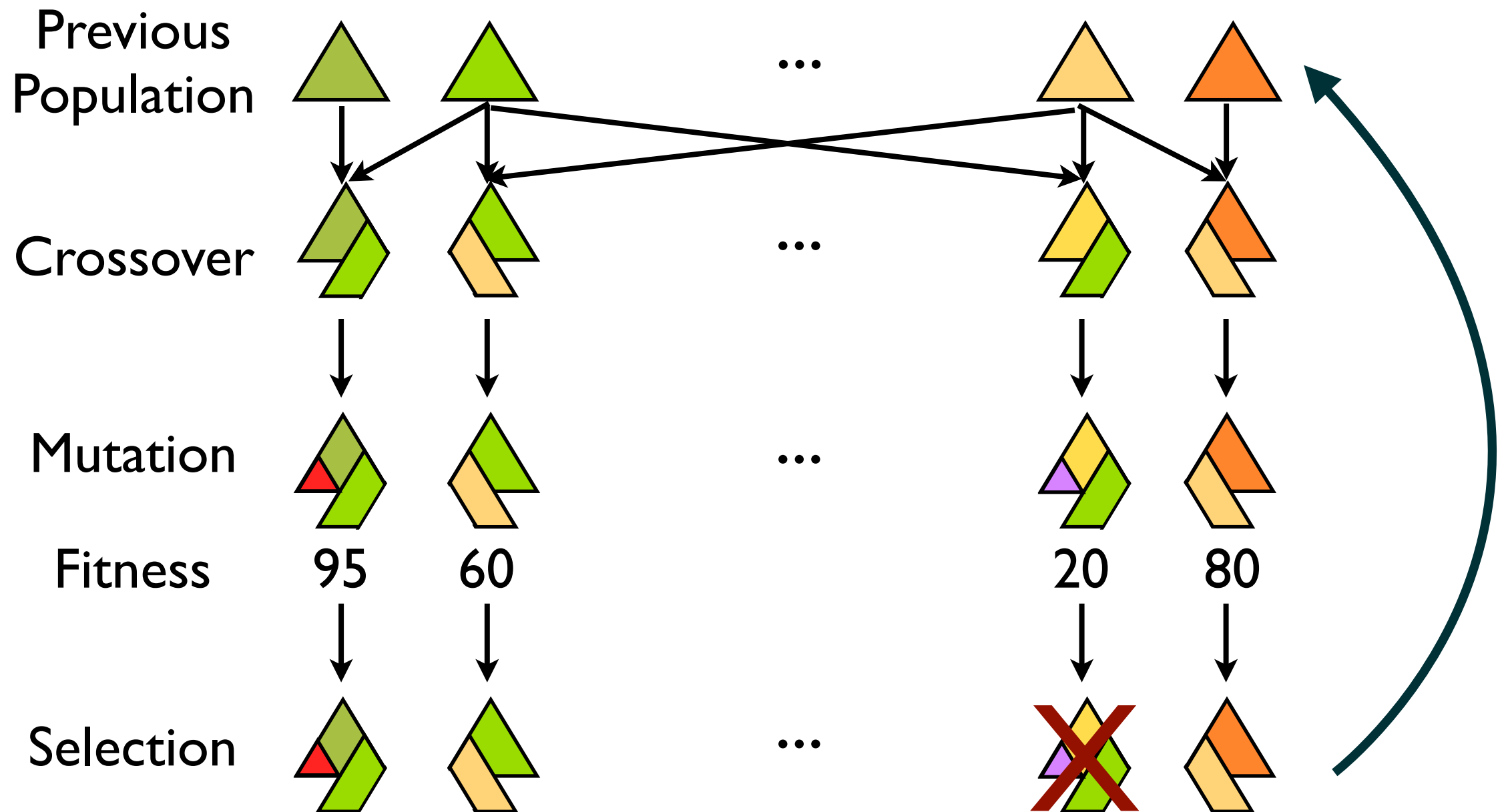
# Genetic Programming



# Genetic Programming



# Genetic Programming

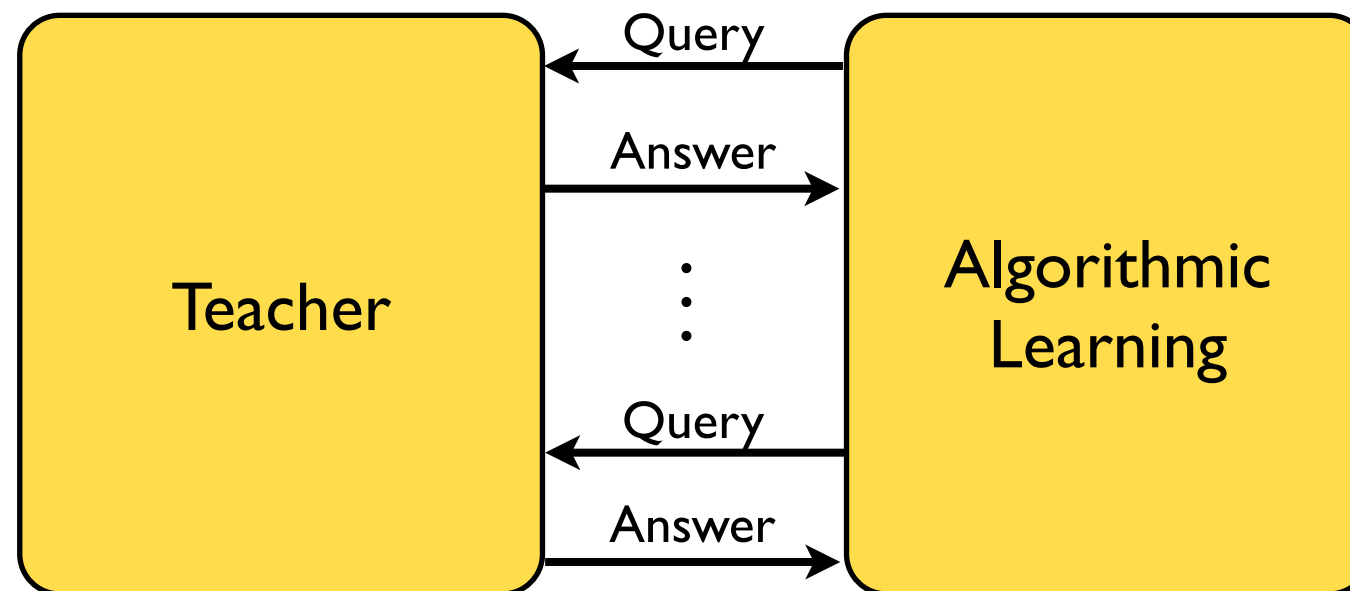




# Genetic Programming as a Solution

## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.



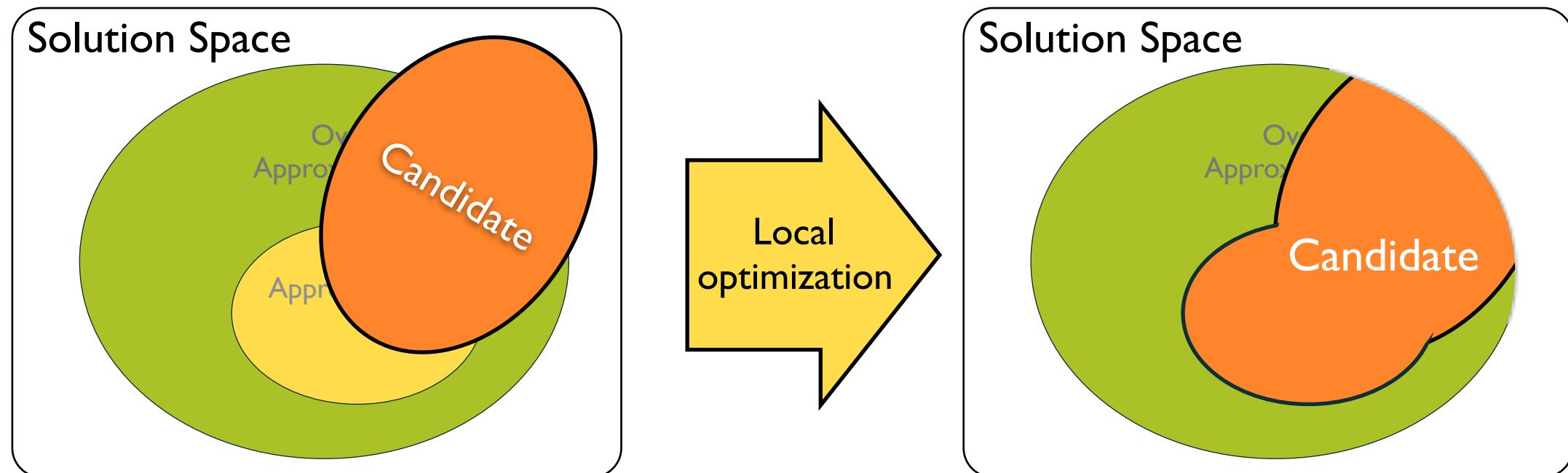
# Genetic Programming as a Solution

## Problem #1: **Active Learning**

Teacher can guide the learner **only through** the form of answers.

## **Solution:** **Active Teaching**

In genetic programming, we can manipulate the formula directly:  $c \Rightarrow (c \vee \underline{l}) \wedge \bar{l}$



# Genetic Programming as a Solution

## Problem #2: **Naive Randomized Mechanism**

Teacher **solely relies on** naive randomized mechanism when approximations are not helpful.

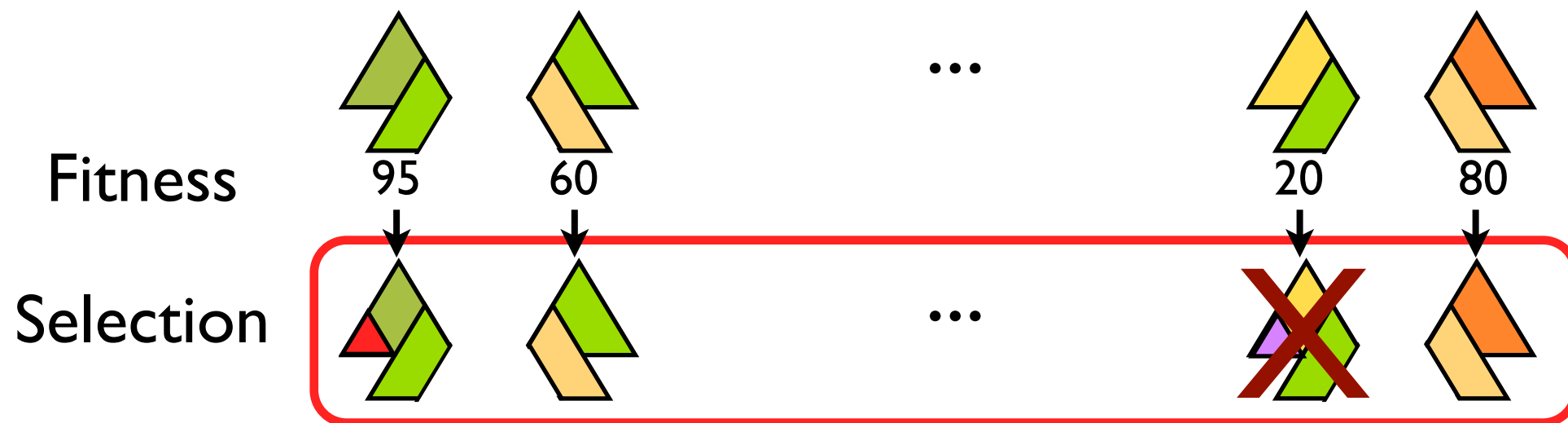
# Genetic Programming as a Solution

## Problem #2: Naive Randomized Mechanism

Teacher **solely relies on** naive randomized mechanism when approximations are not helpful.

## Solution: Natural Selection

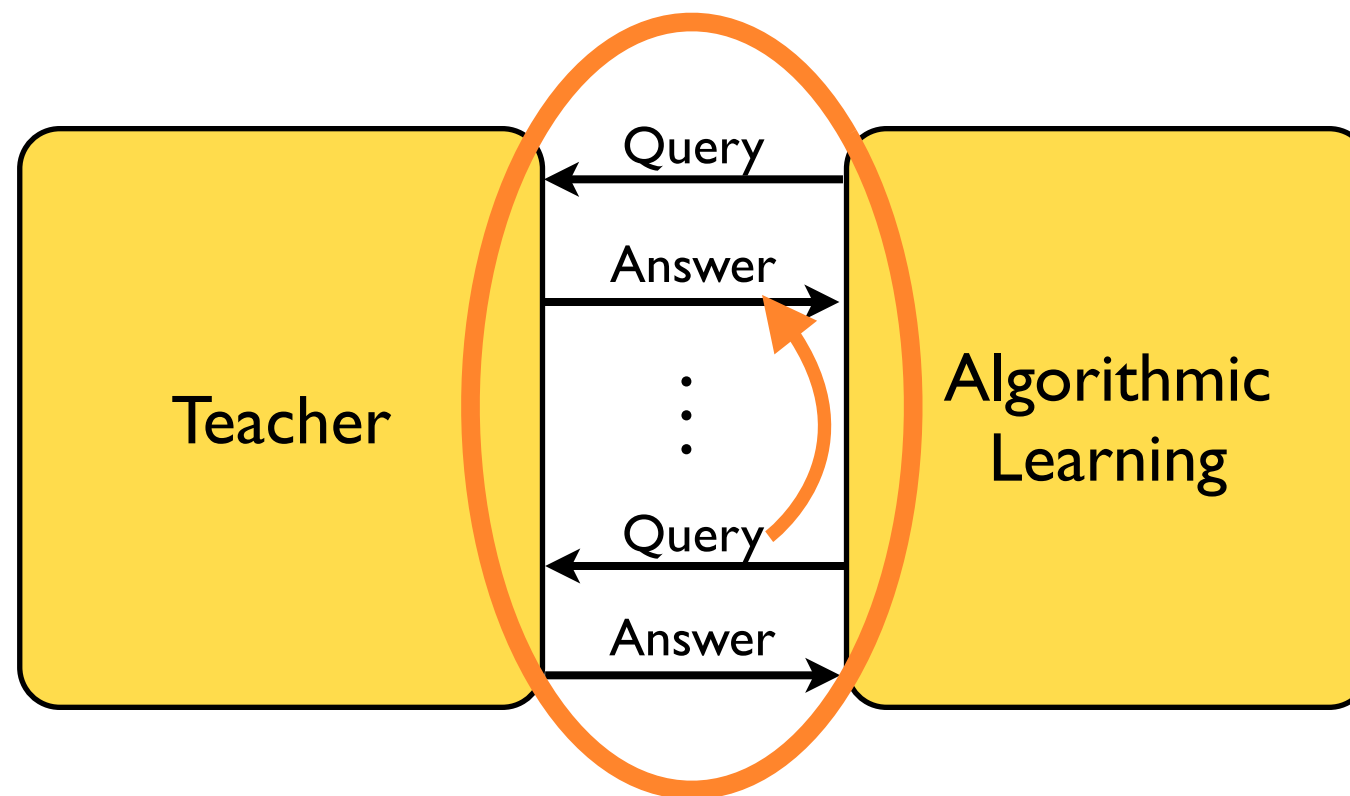
Genetic programming **favors** better solutions.



# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

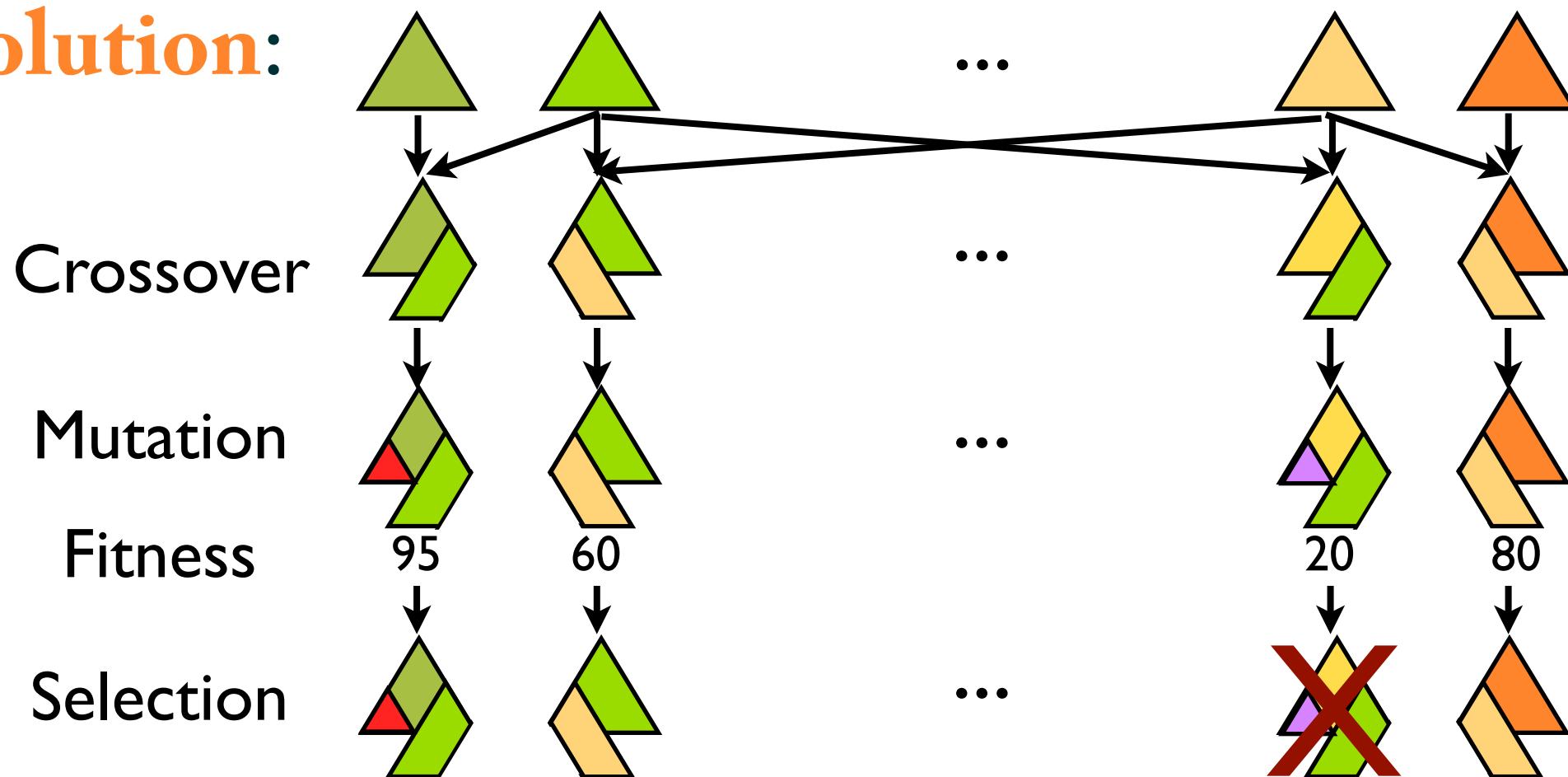


# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

### Solution:

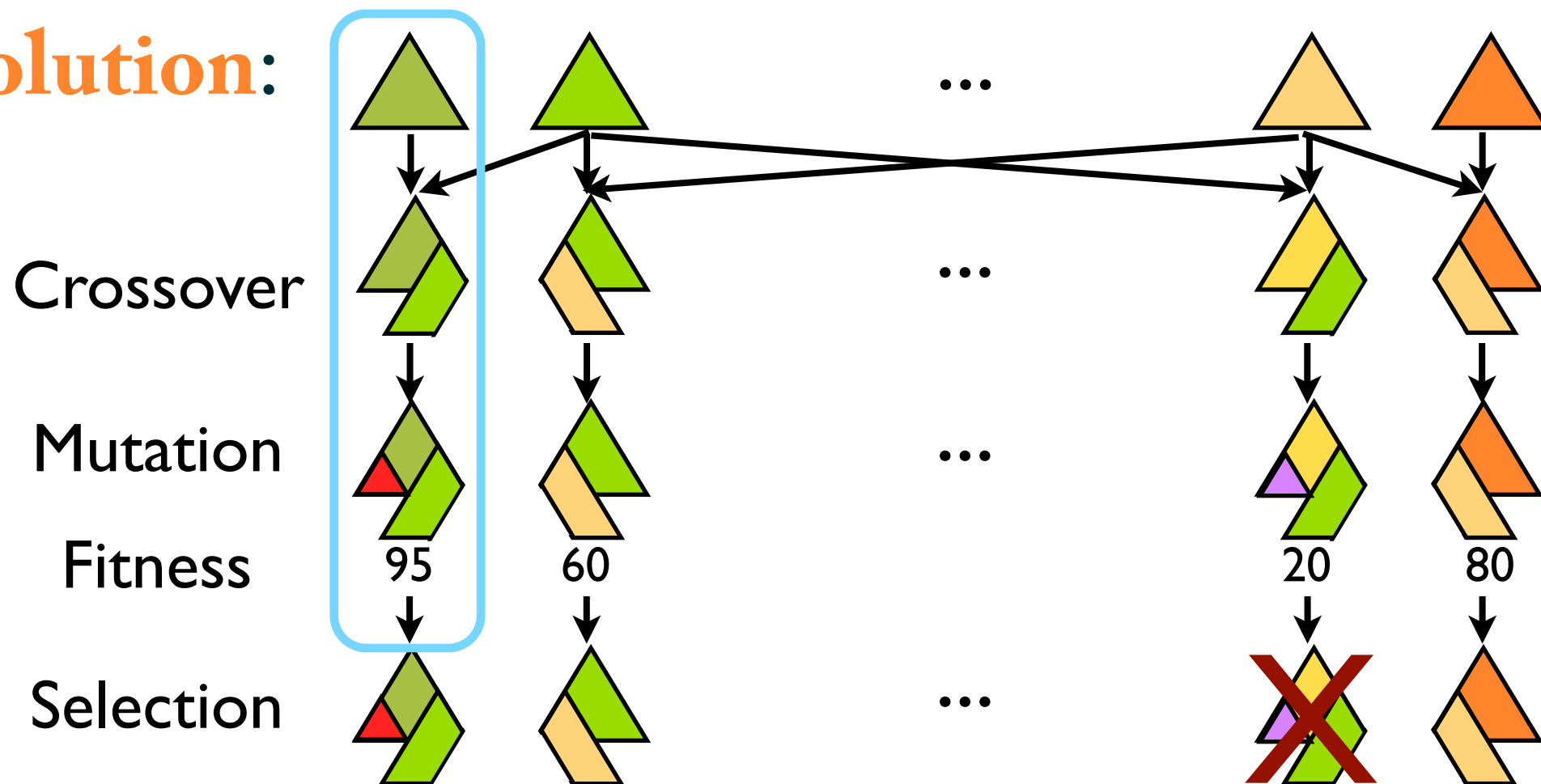


# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

### Solution:

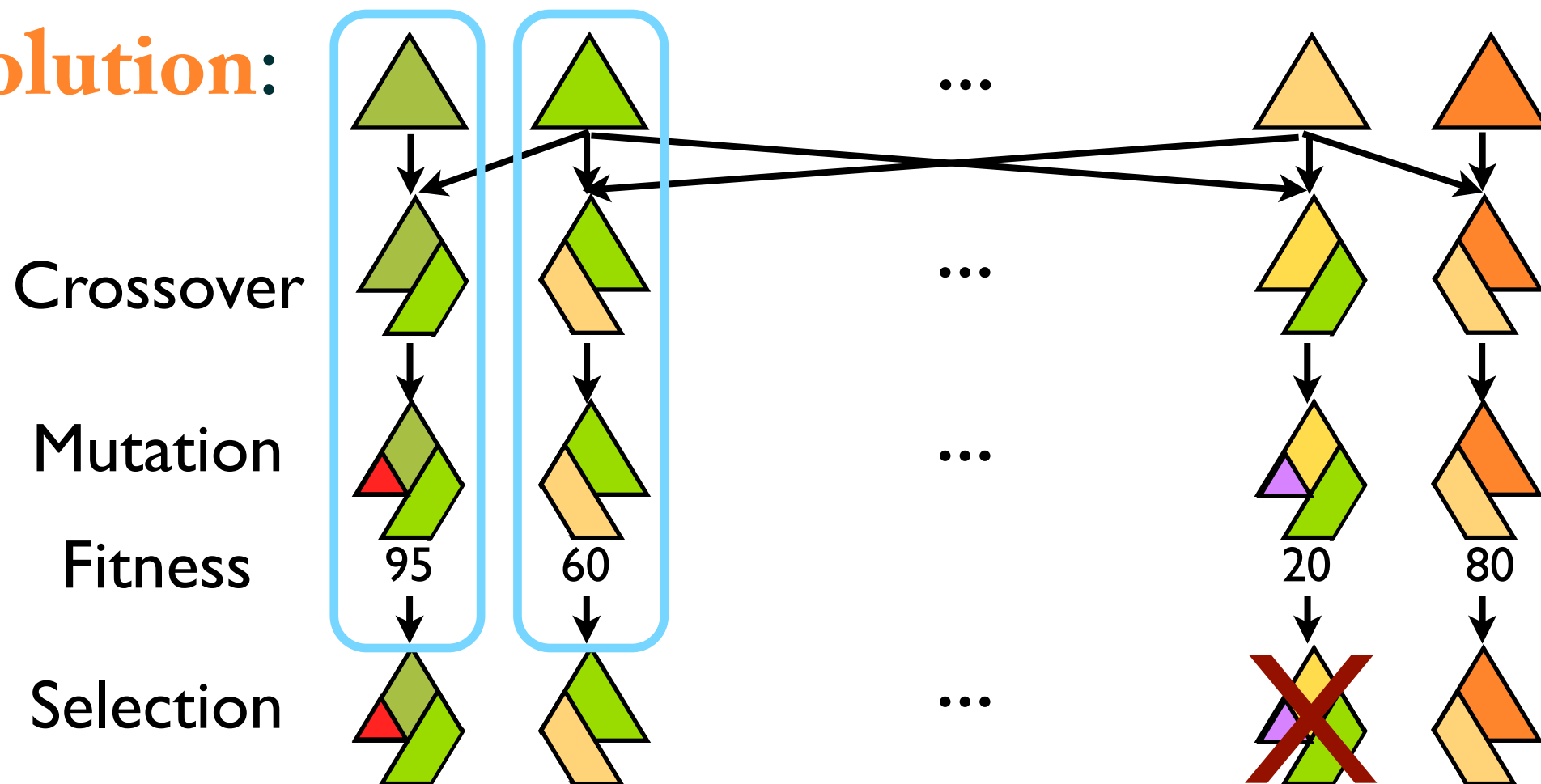


# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

### Solution:



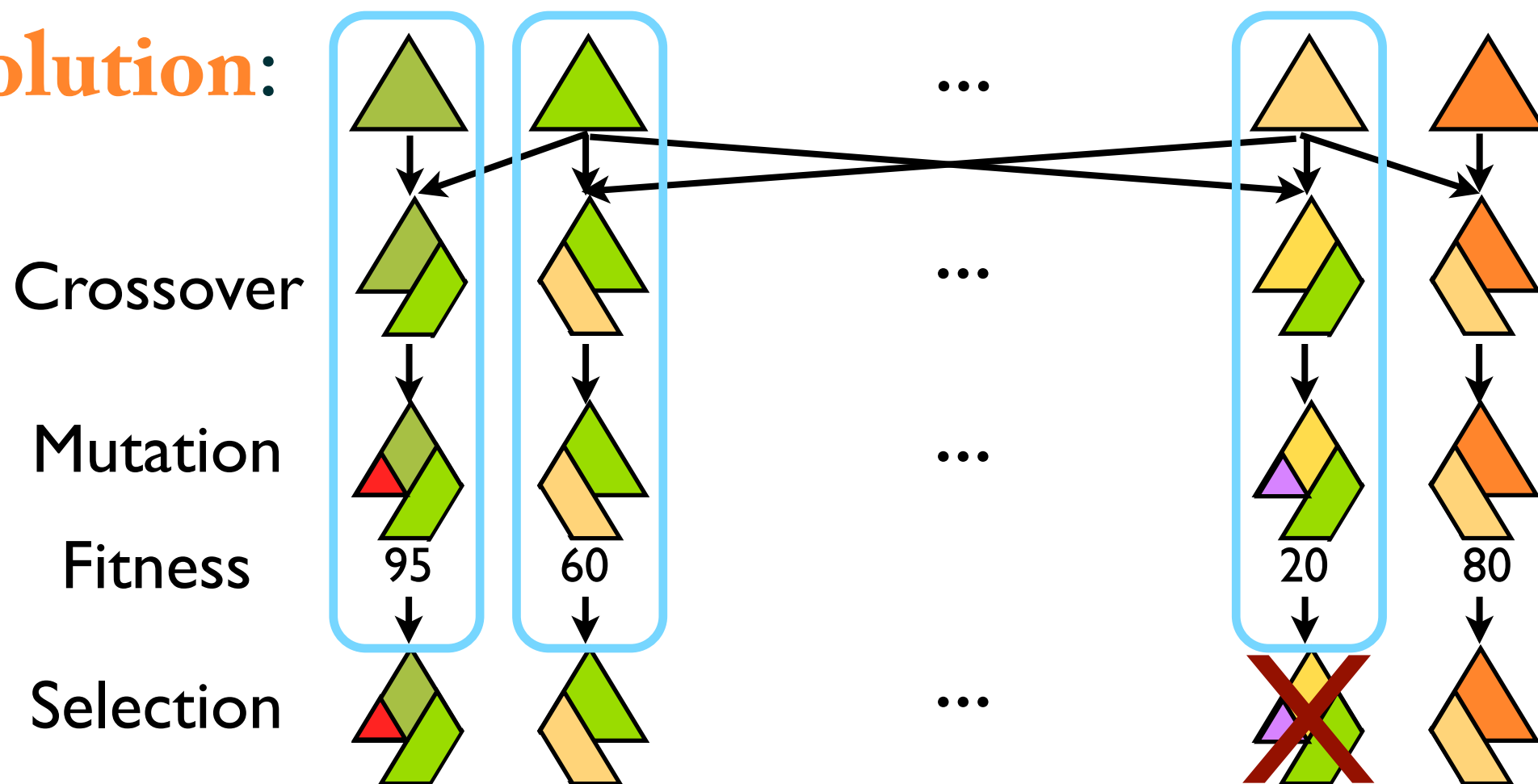


# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

### Solution:

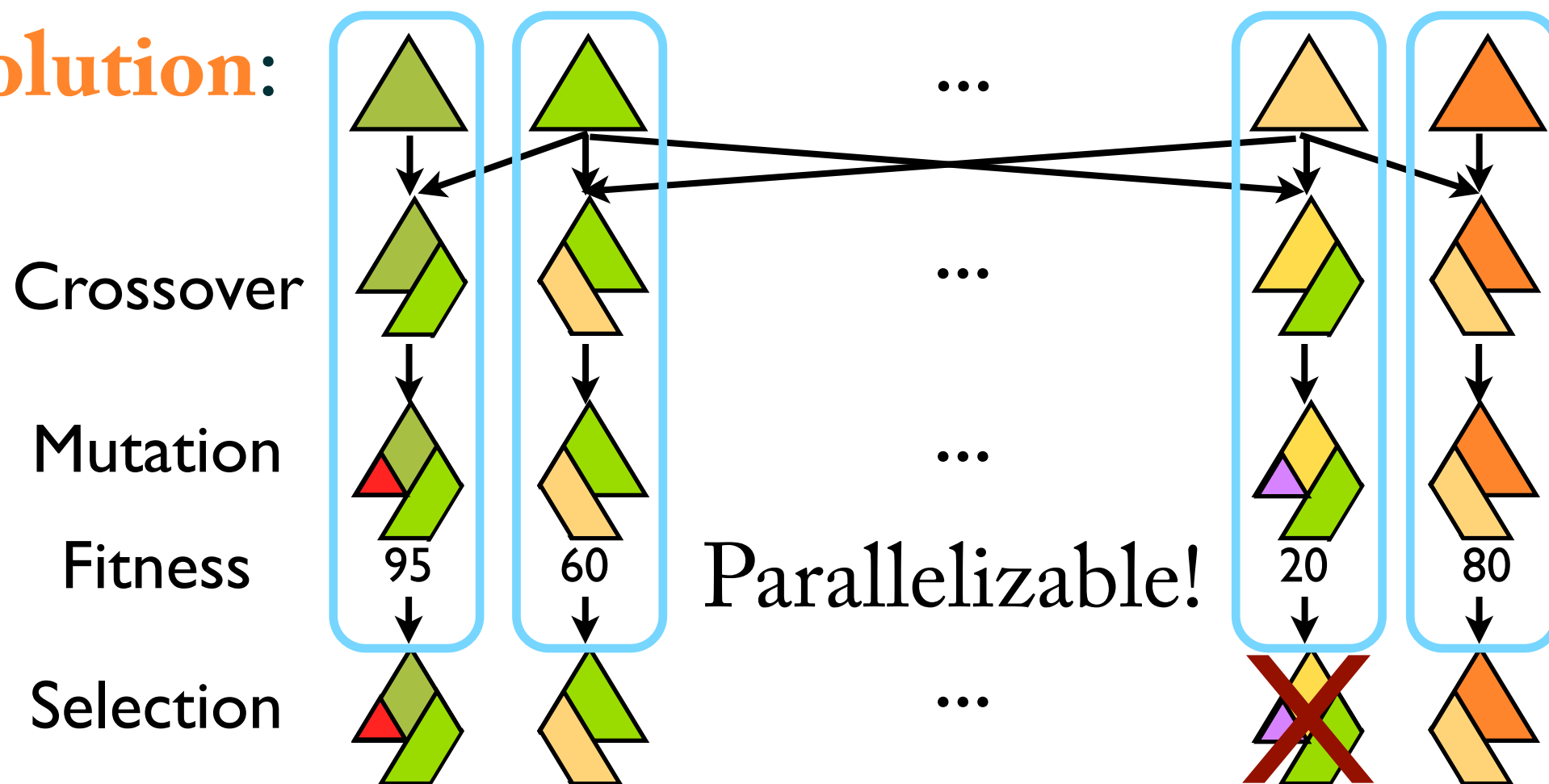


# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

### Solution:

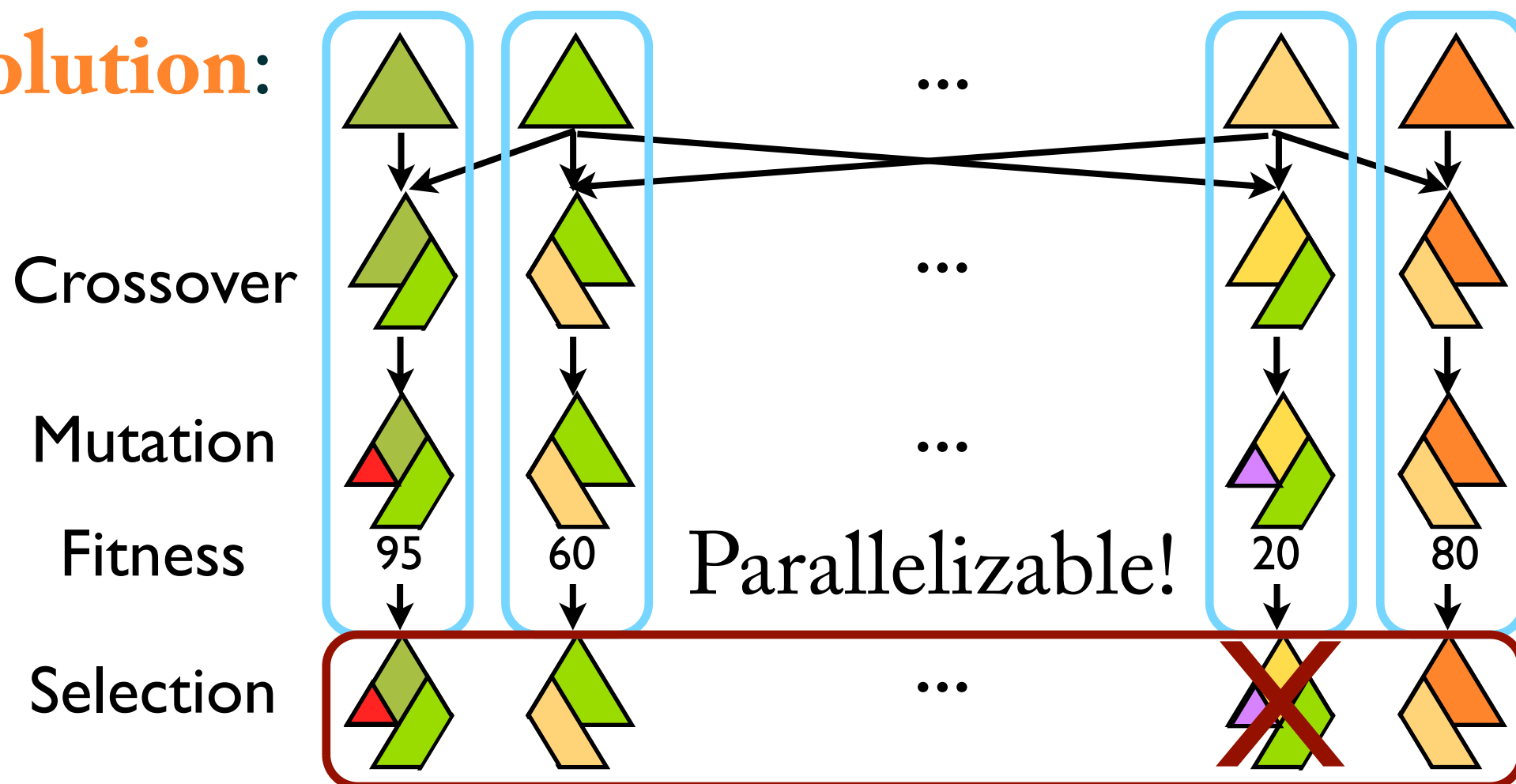


# Genetic Programming as a Solution

## Problem #3: Serialized Query-Answer Process

Queries are **dependent** on the previous query/answer pairs and make parallelization **difficult**.

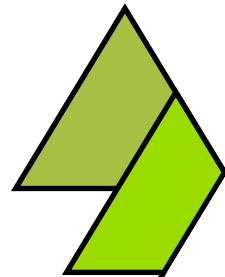
### Solution:



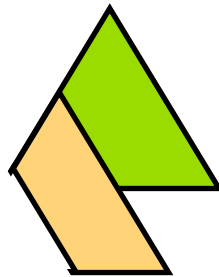
# Key Issue: How to define Fitness Function?



Fitness

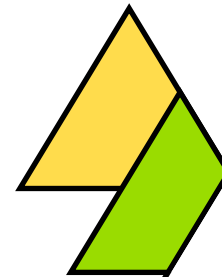


95

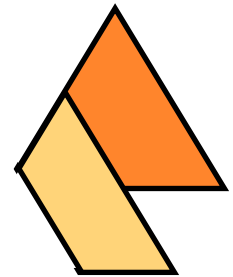


60

...



20



80

# Role of Fitness Function

For the annotated loop

$\{\delta\}$  while  $\rho$  do  $S$  end  $\{\epsilon\}$

Find an invariant  $I$  satisfying the following conditions:

- (A)  $\underline{\iota} \Rightarrow \iota$  (  $\iota$  holds when entering the loop)
- (B)  $\iota \wedge \rho \Rightarrow Pre(\iota, S)$  (  $\iota$  holds at each iteration)
- (C)  $\iota \Rightarrow \bar{\iota}$  (  $\iota$  gives  $\epsilon$  after leaving the loop)

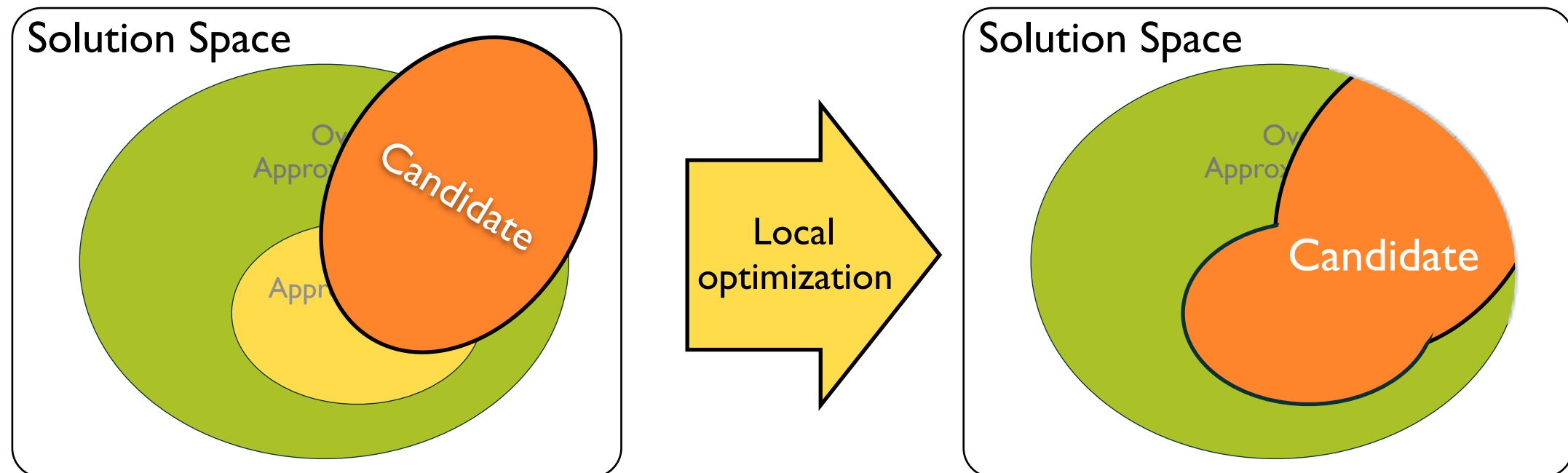
# Role of Fitness Function

For the annotated loop

$$\{\delta\} \text{ while } \rho \text{ do } S \text{ end } \{\epsilon\}$$

Find an invariant  $I$  satisfying the following conditions:

- (A)  $\underline{\iota} \Rightarrow \iota$  (  $\iota$  holds when entering the loop)
- (B)  $\iota \wedge \rho \Rightarrow Pre(\iota, S)$  (  $\iota$  holds at each iteration)
- (C)  $\iota \Rightarrow \bar{\iota}$  (  $\iota$  gives  $\epsilon$  after leaving the loop)



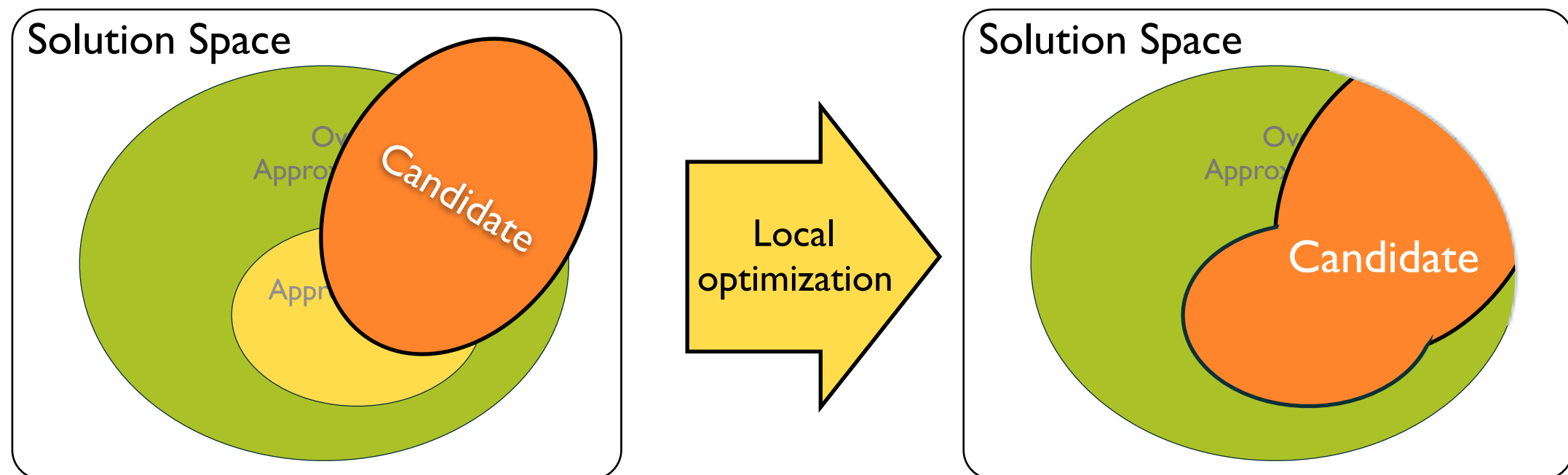
# Role of Fitness Function

For the annotated loop

$$\{\delta\} \text{ while } \rho \text{ do } S \text{ end } \{\epsilon\}$$

Find an invariant  $I$  satisfying the following conditions:

- ✓ (A)  $\underline{\iota} \Rightarrow \iota$  (  $\iota$  holds when entering the loop)
- (B)  $\iota \wedge \rho \Rightarrow Pre(\iota, S)$  (  $\iota$  holds at each iteration)
- ✓ (C)  $\iota \Rightarrow \bar{\iota}$  (  $\iota$  gives  $\epsilon$  after leaving the loop)



# Role of Fitness Function

For the annotated loop

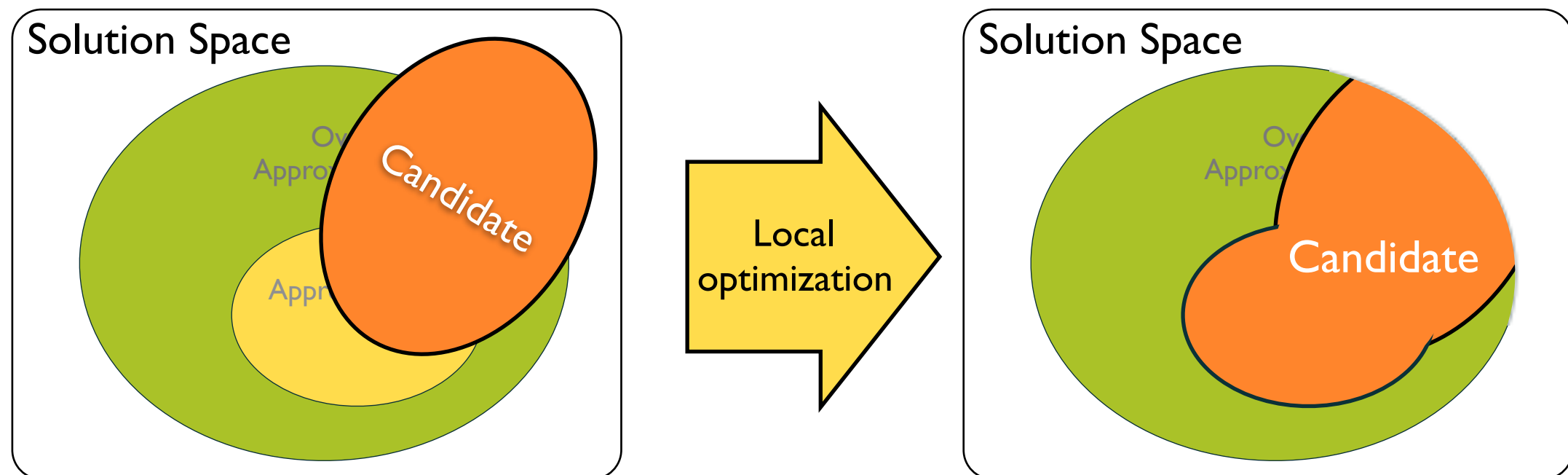
$\{\delta\}$  while  $\rho$  do  $S$  end  $\{\epsilon\}$

Find an invariant  $I$  satisfying the following conditions:

✓ (A)  $\underline{\iota} \Rightarrow \iota$  (  $\iota$  holds when entering the loop)

(B)  $\iota \wedge \rho \Rightarrow Pre(\iota, S)$  (  $\iota$  holds at each iteration)

✓ (C)  $\iota \Rightarrow \bar{\iota}$  (  $\iota$  gives  $\epsilon$  after leaving the loop)

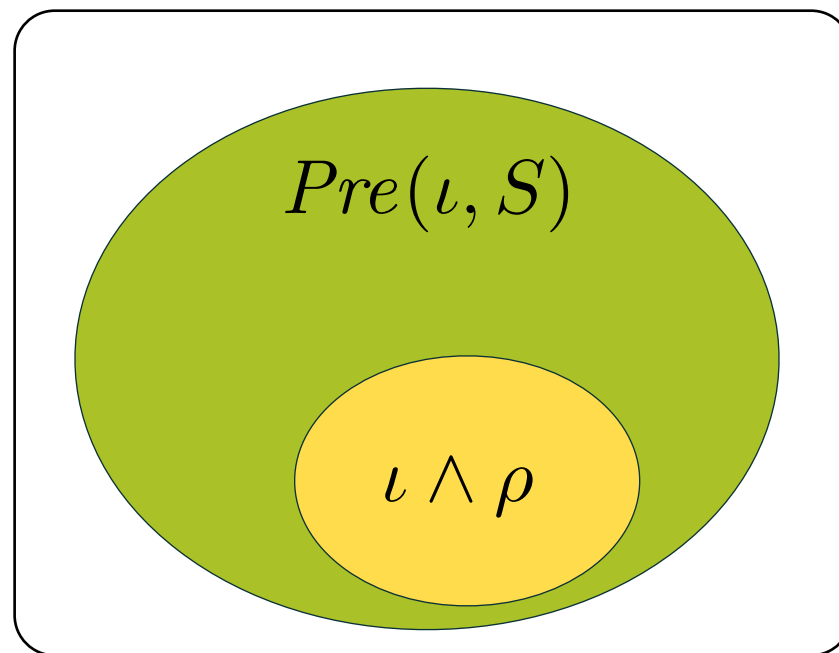




# Role of Fitness Function

(B)  $\iota \wedge \rho \Rightarrow Pre(\iota, S)$  ( $\iota$  holds at each iteration)

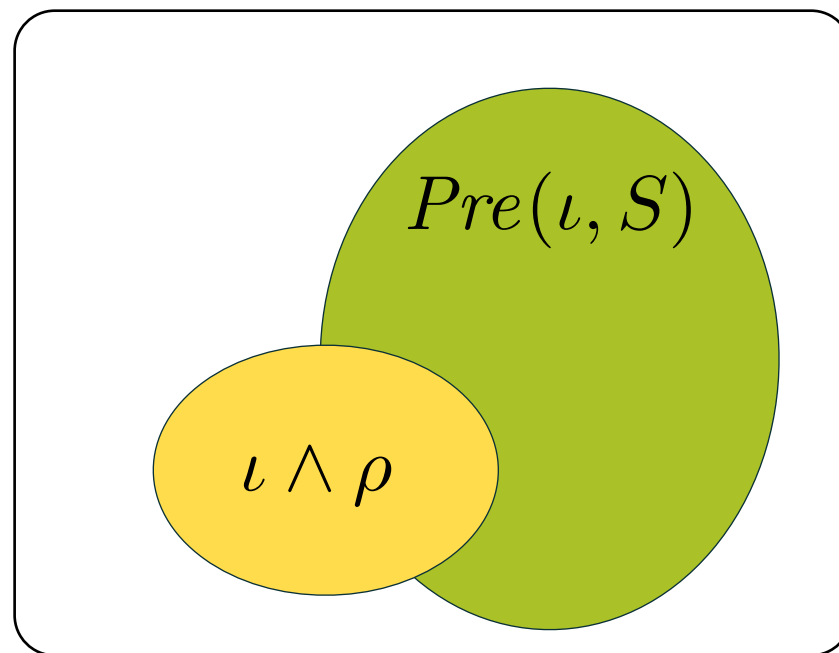
If (B) holds, then we find an invariant.



# Role of Fitness Function

$$(B) \quad \iota \wedge \rho \Rightarrow Pre(\iota, S) \quad (\iota \text{ holds at each iteration})$$

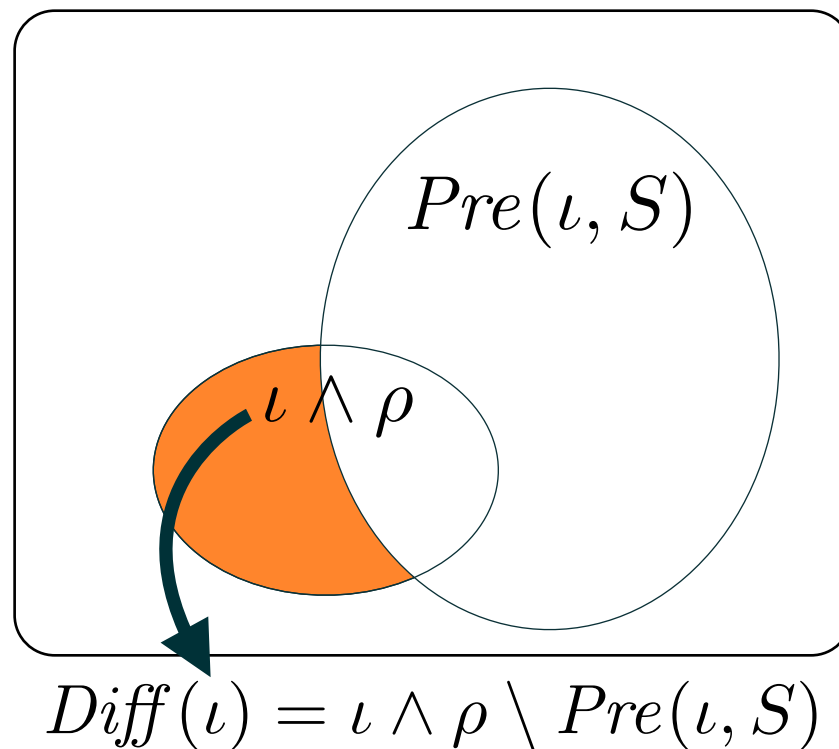
If (B) does not hold,  
we need to quantify how bad the candidate is.



# Role of Fitness Function

$$(B) \quad \iota \wedge \rho \Rightarrow Pre(\iota, S) \quad (\iota \text{ holds at each iteration})$$

If (B) does not hold,  
we need to quantify how bad the candidate is.



Fitness Function  
Property

$$Diff(\iota_1) \subseteq Diff(\iota_2) \Rightarrow Fit(\iota_1) \leq Fit(\iota_2)$$

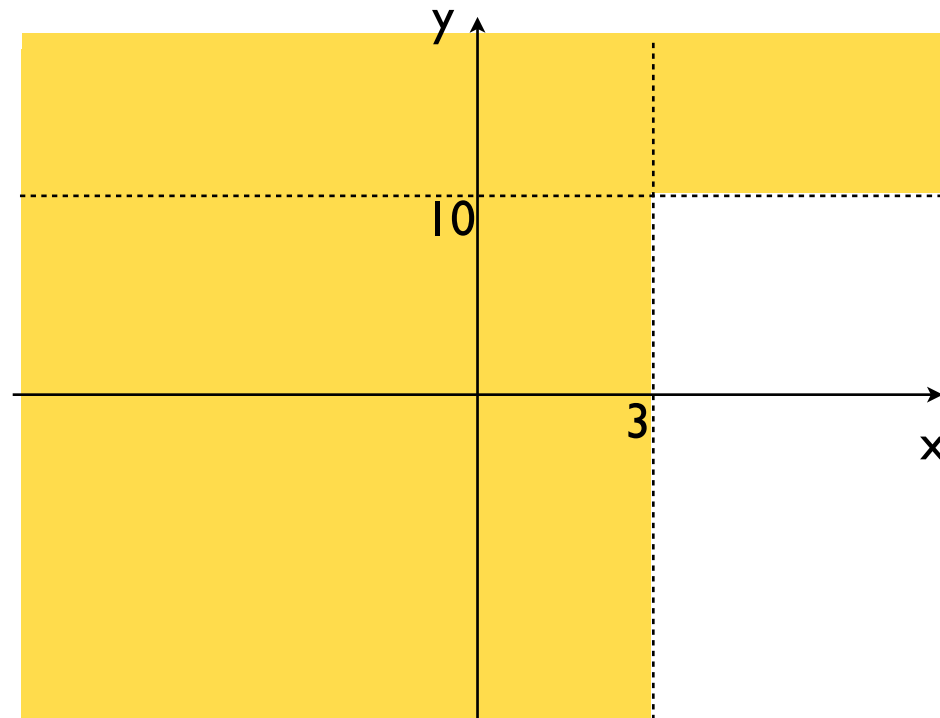
# Infeasible Fitness Function

Count the number of models in **concrete domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$Fit(\iota) = |\{(2, 11), (10, 15), (0, 0), \dots\}|$$



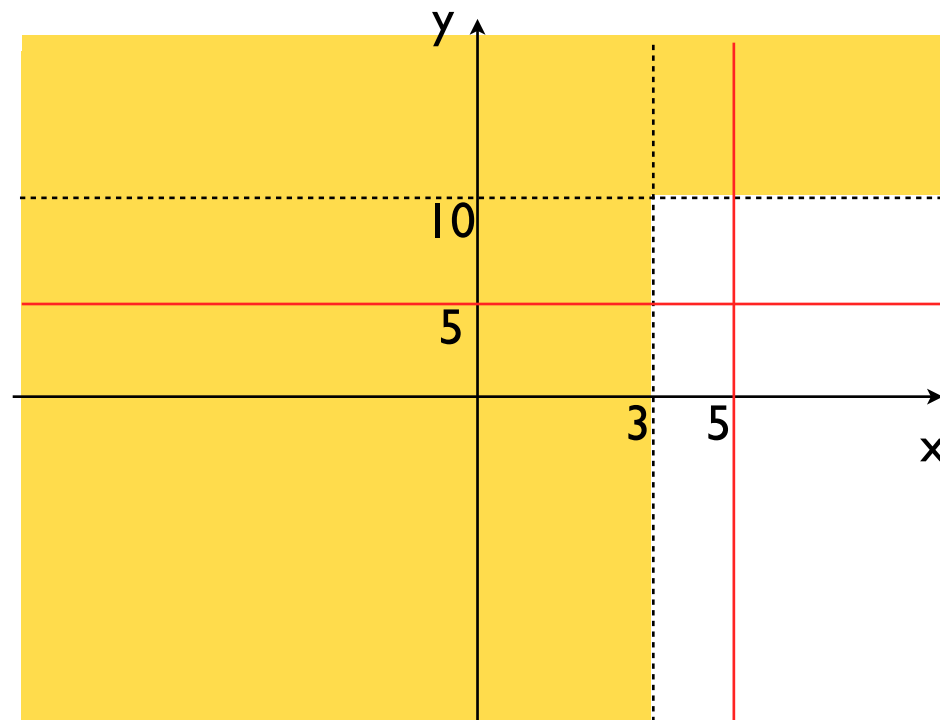
Can **not** count them all in general.

# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$



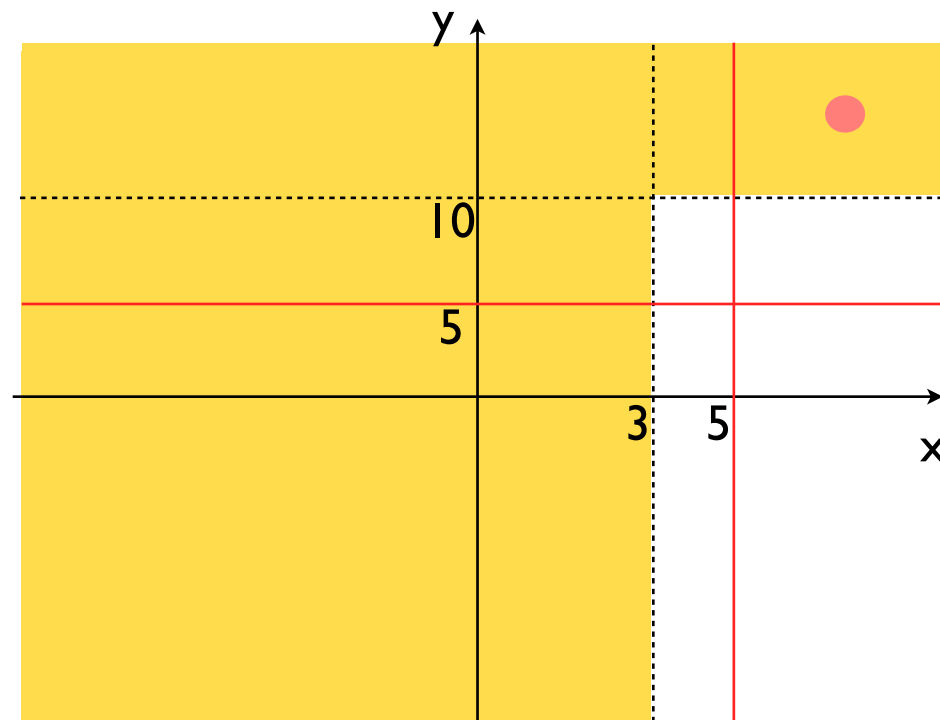
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$\text{SMT: } (10, 15) \models Diff(\iota)$$



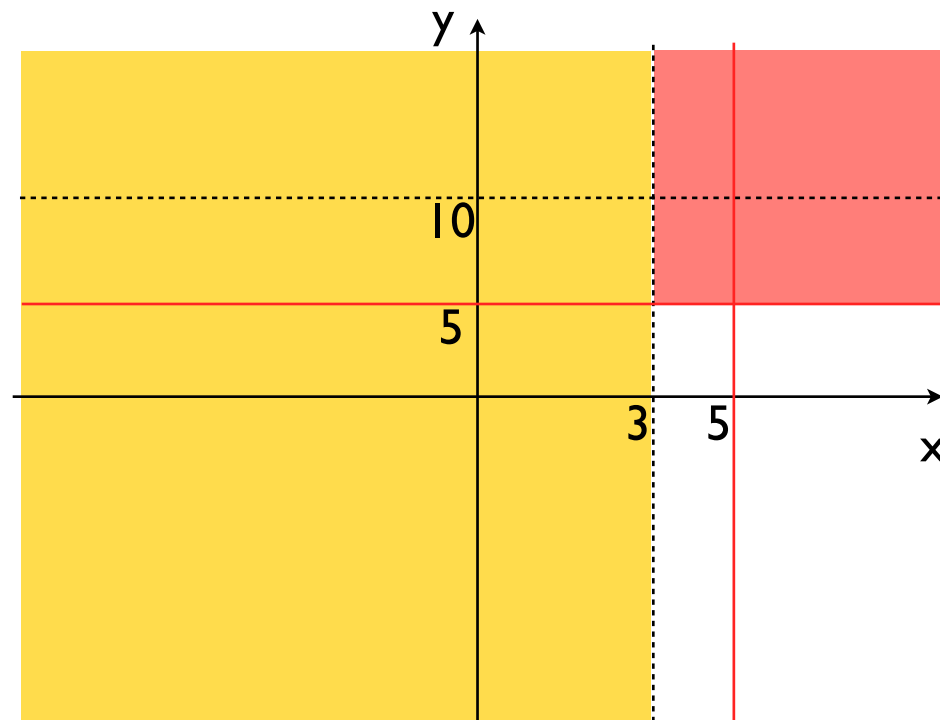
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$\text{SMT: } (10, 15) \models Diff(\iota)$$



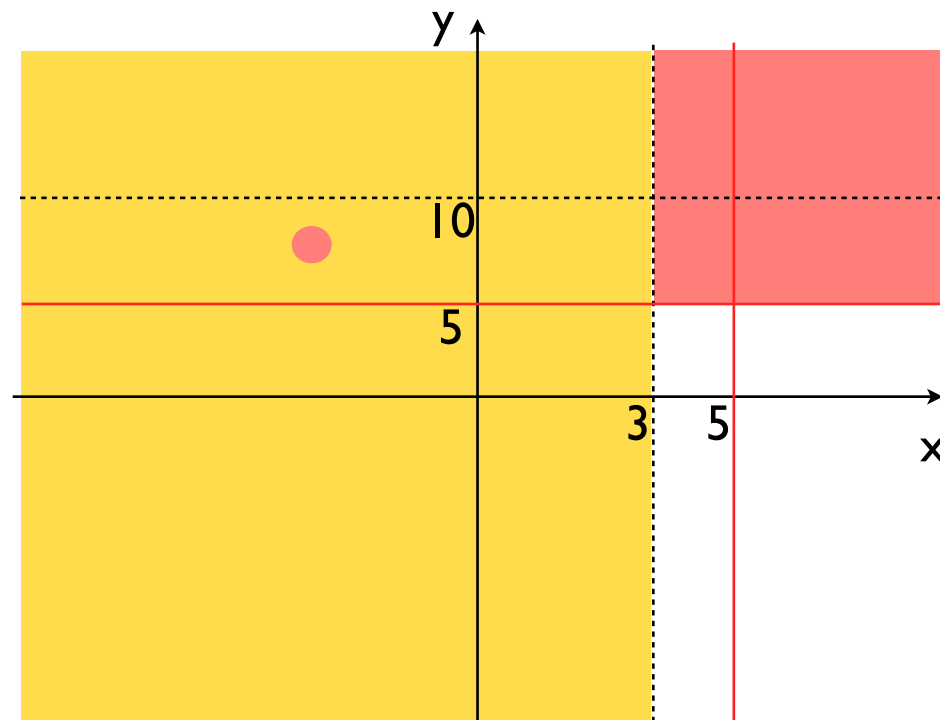
# Feasible Fitness Function

Count the number of models in **abstract domain**.

### Example

$$Diff(\iota) = x < 3 \vee y > 10$$

SMT:  $(-3, 7) \models Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)]$





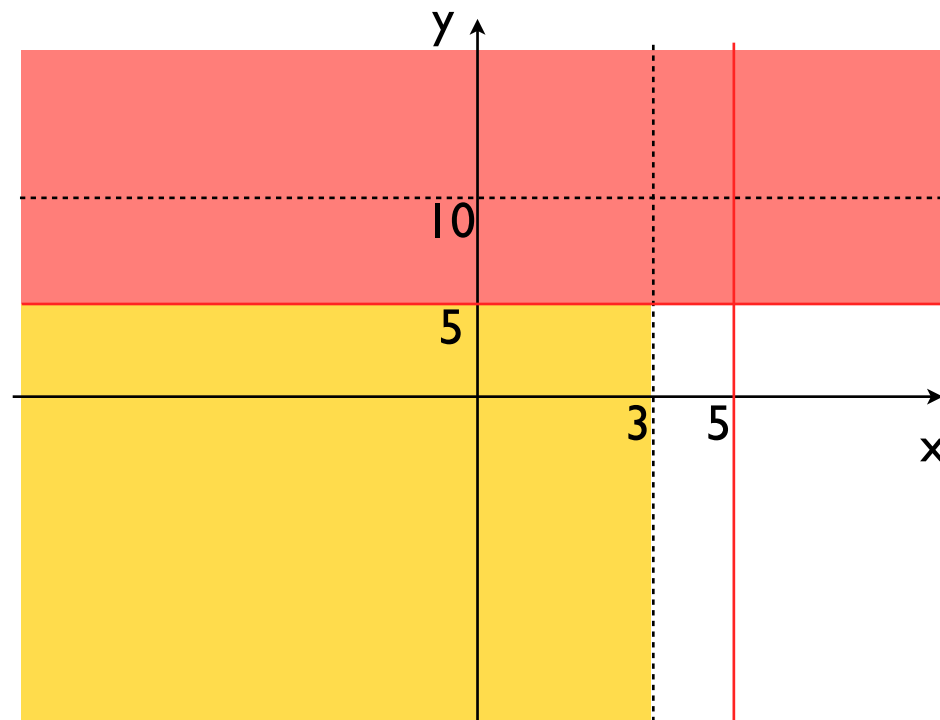
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$\text{SMT: } (-3, 7) \models Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)]$$



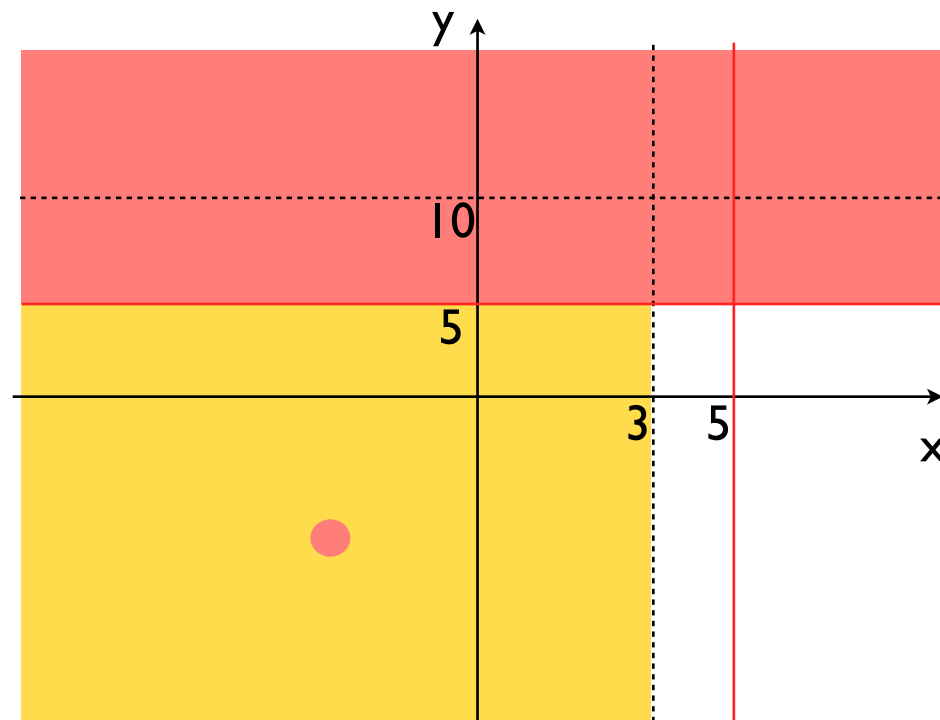
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

SMT:  $(-5, -5) \models Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y > 5)]$



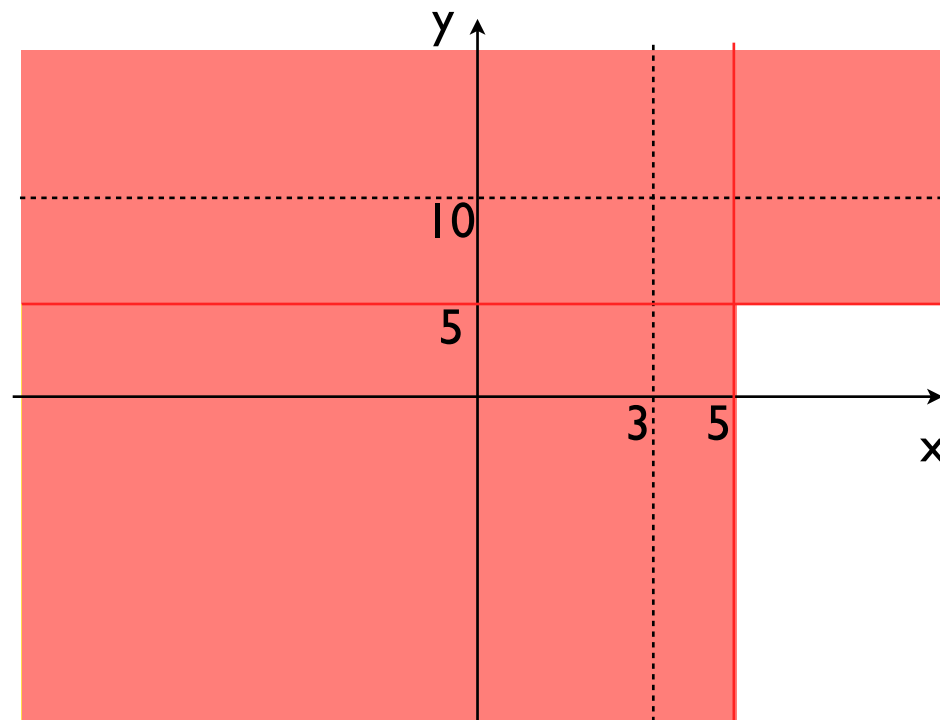
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

SMT:  $(-5, -5) \models Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y > 5)]$



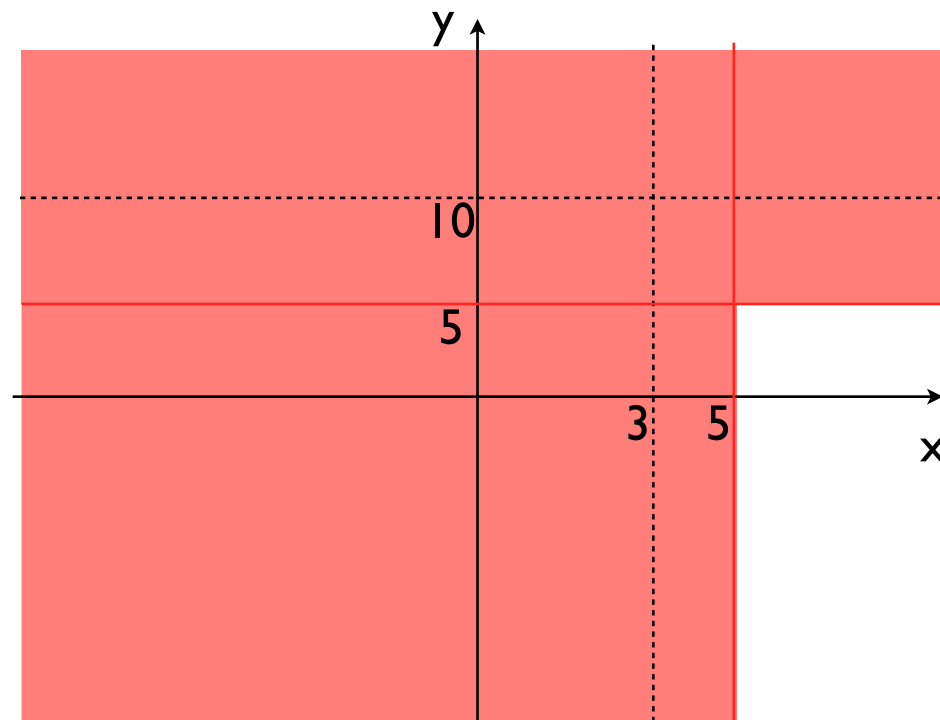
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$SMT(Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y \leq 5)]) \\ = \text{UNSAT}$$



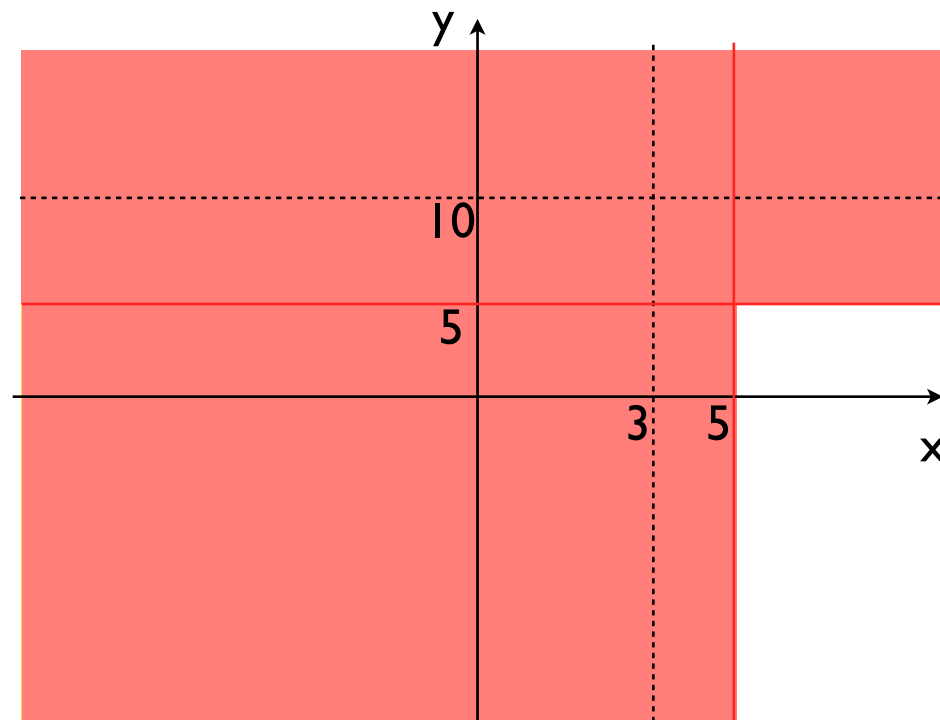
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$SMT(Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y \leq 5)]) \\ = \text{UNSAT}$$



$$\widehat{Fit}(\iota) = 3$$

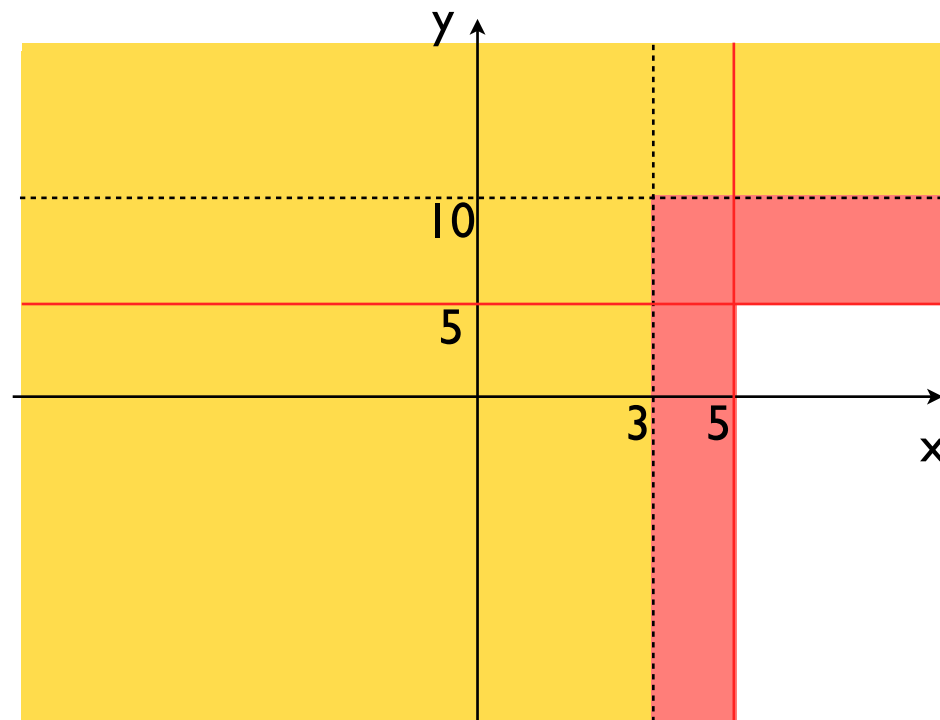
# Feasible Fitness Function

Count the number of models in **abstract domain**.

Example

$$Diff(\iota) = x < 3 \vee y > 10$$

$$SMT(Diff(\iota) \wedge \neg[(x > 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y \leq 5)]) \\ = \text{UNSAT}$$



Fitness Function  
Property ✓

$$Diff(\iota_1) \subseteq Diff(\iota_2) \Rightarrow Fit(\iota_1) \leq Fit(\iota_2)$$

# Complexity of Fitness Function

Number of SMT calls could raise up to  $2^n$   
if we have  $n$  atomic propositions in predicate  
abstraction.

# Complexity of Fitness Function

Number of SMT calls could raise up to  $2^n$  if we have  $n$  atomic propositions in predicate abstraction.

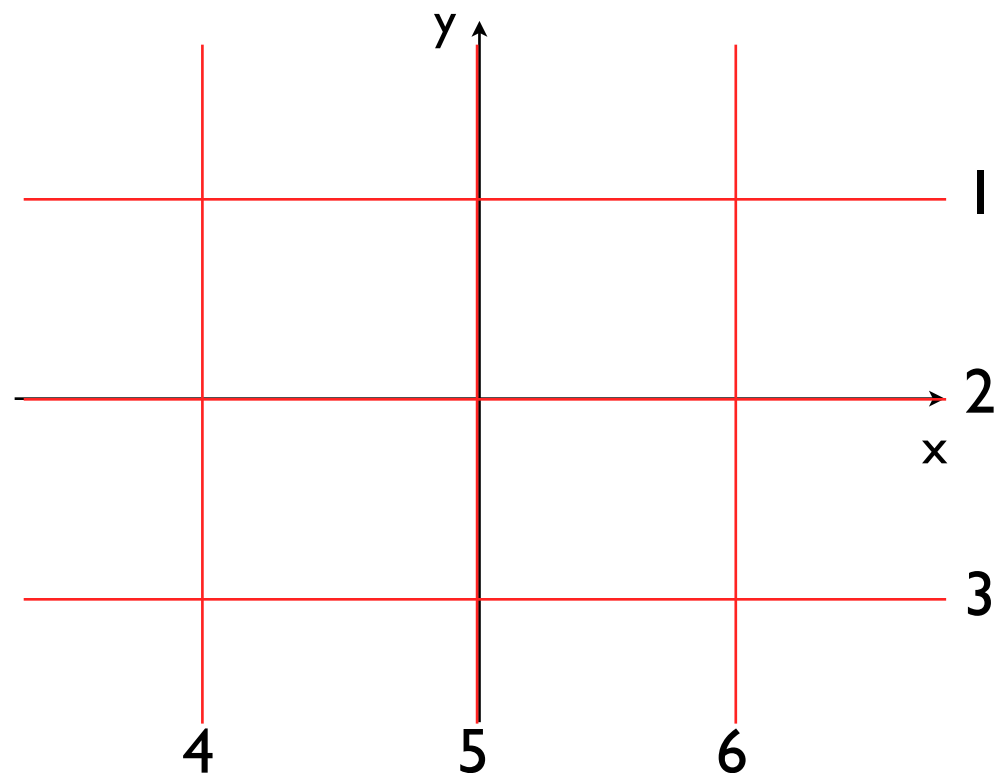
In practice, that is **not** the case because of dependencies among atomic propositions.



# Complexity of Fitness Function

Number of SMT calls could raise up to  $2^n$  if we have  $n$  atomic propositions in predicate abstraction.

In practice, that is **not** the case because of dependencies among atomic propositions.



# Complexity of Fitness Function

Number of SMT calls could raise up to  $2^n$  if we have  $n$  atomic propositions in predicate abstraction.

In practice, that is **not** the case because of dependencies among atomic propositions.

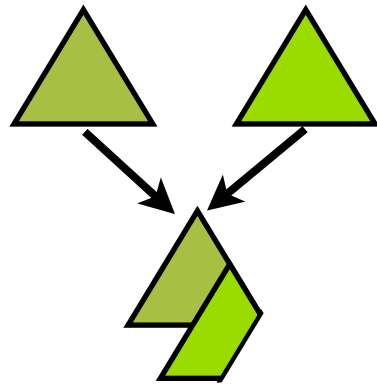
1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

$$16 \leq 2^6 = 64$$

# Issue: Formula Size

Formula size can **grow infinitely**.

## 1. Crossover operation



## 2. Fitness Function

SMT:  $(10, 15) \models \text{Diff}(\iota)$

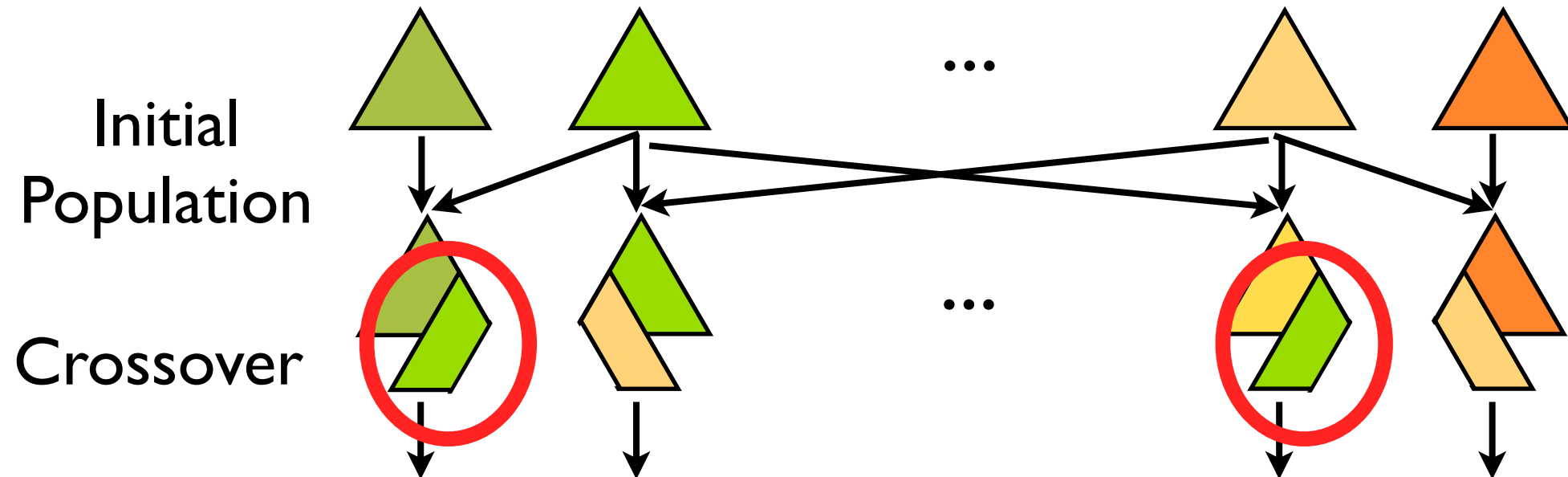
$(-3, 7) \models \text{Diff}(\iota) \wedge \neg[(x > 5) \wedge (y > 5)]$

$(-5, -5) \models \text{Diff}(\iota) \wedge \neg[(x > 5) \wedge (y > 5)] \wedge \neg[(x \leq 5) \wedge (y > 5)]$

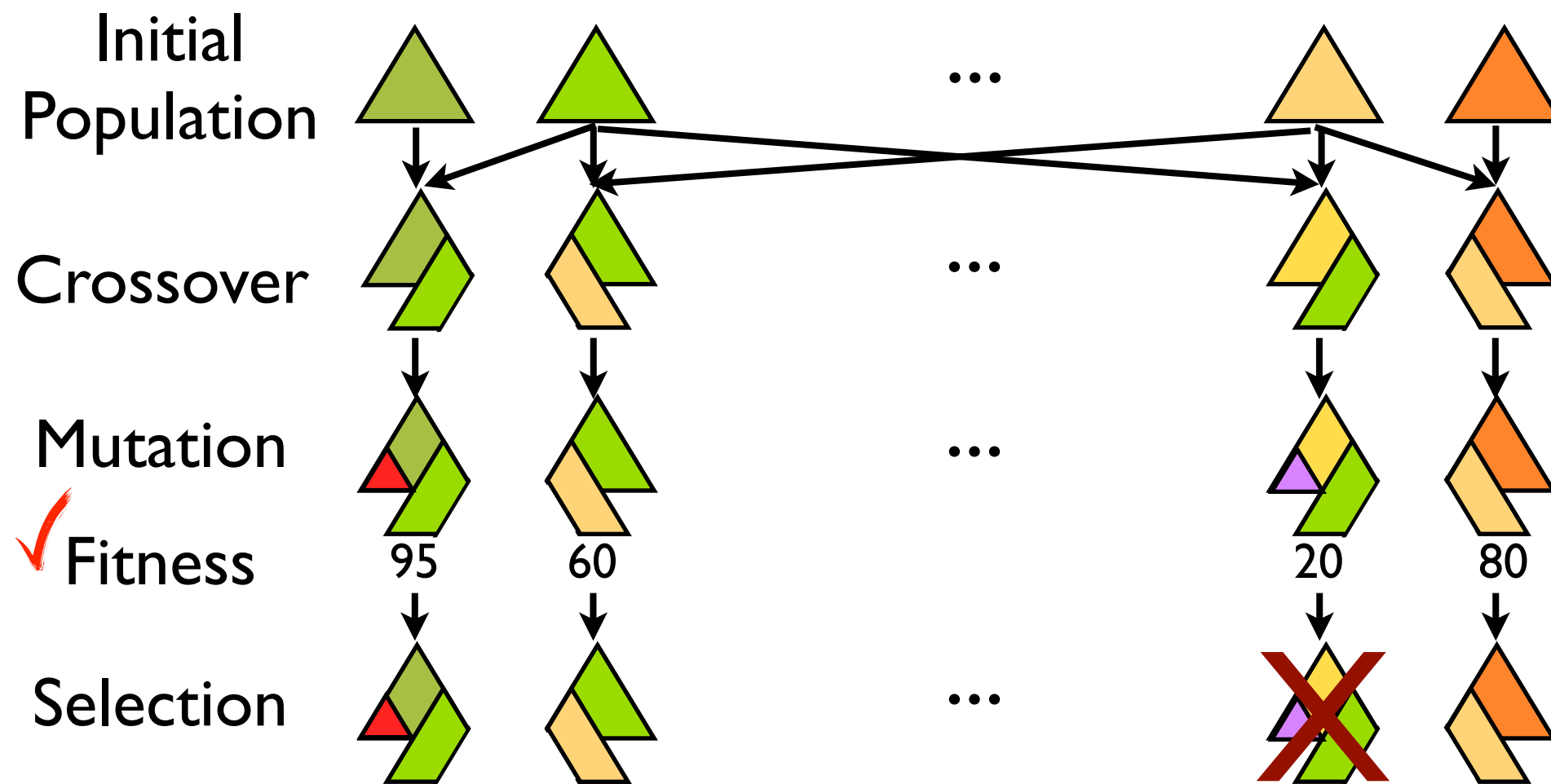
**Solution:** Simplify formulae.

# Progress

1. Design and implement abstract data type for formula to **share subformulae**  
(Heejae Shin & Wonchan Lee)



# Plan



Thank you