

Lecture 19: Program Synthesis

17-355/17-665/17-819: Program Analysis

Rohan Padhye

March 31, 2022

* Course materials developed with Jonathan Aldrich and Claire Le Goues
With slide inspiration gratitude to Emina Torlak and Ras Bodik

Program Synthesis Overview

- A mathematical characterization of program synthesis: prove that $\exists P . \forall x . \varphi(x, P(x))$
- In constructive logic, the witness to the proof of this statement is a program P that satisfies property φ for all input values x

Program Synthesis Overview

- A mathematical characterization of program synthesis: prove that

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

- In constructive logic, the witness to the proof of this statement is a program P that satisfies property φ for all input values x
- What could the inferred program P be?
 - Historically, a protocol, interpreter, classifier, compression algorithm, scheduling policy, cache coherence policy, ...
- How is property φ expressed?
 - Historically, as a formula, a reference implementation, input/output pairs, traces, demonstrations, a sketch, ...

Exercise: specify $P_{max}(list)$

- Specify a program $P_{max}(l)$ that finds the maximum number in a list l . How many different ways can you do it?

Expressing User Intent

- How do we constrain the program to be synthesized?
 - Express what we know about the problem and/or solution
 - Usually incomplete
- Two forms of specification can constrain synthesis
 - Observable behavior: input/output relations, executable specification, safety property
 - Structural properties: constraints on internal computation, such as a sketch, template, assertions about structure (e.g. number of iterations)

The Search Space of Programs

- Constraining the search space can help make synthesis feasible
 - Subset of a real programming language?
 - Grammar for combining fixed set of operators and control structures?
 - DSL?
 - Logic?

Two approaches to searching for programs

- Deductive synthesis
 - Maps a high-level specification to an implementation, using a theorem prover
 - Efficient, provably correct
 - Require complete specifications, sufficient axiomatization of the domain
 - Can be as complicated as writing the program itself!
 - Used for e.g. controllers
 - A lot like compilation!
- Inductive synthesis
 - Takes a partial, perhaps multi-modal specification and constructs a program that satisfies it
 - Flexible in specification requirements, require no axioms
 - May be less efficient, weaker guarantees on correctness/optimalty
 - Search techniques: brute-force, probabilistic, genetic programming, logical reasoning
 - Major current focus of research

Inductive synthesis

Find a program correct on a set of inputs and hope (or verify) that it's correct on other inputs.

A **partial program** syntactically defines the candidate space.

Inductive synthesis search phrased as a **constraint problem**.

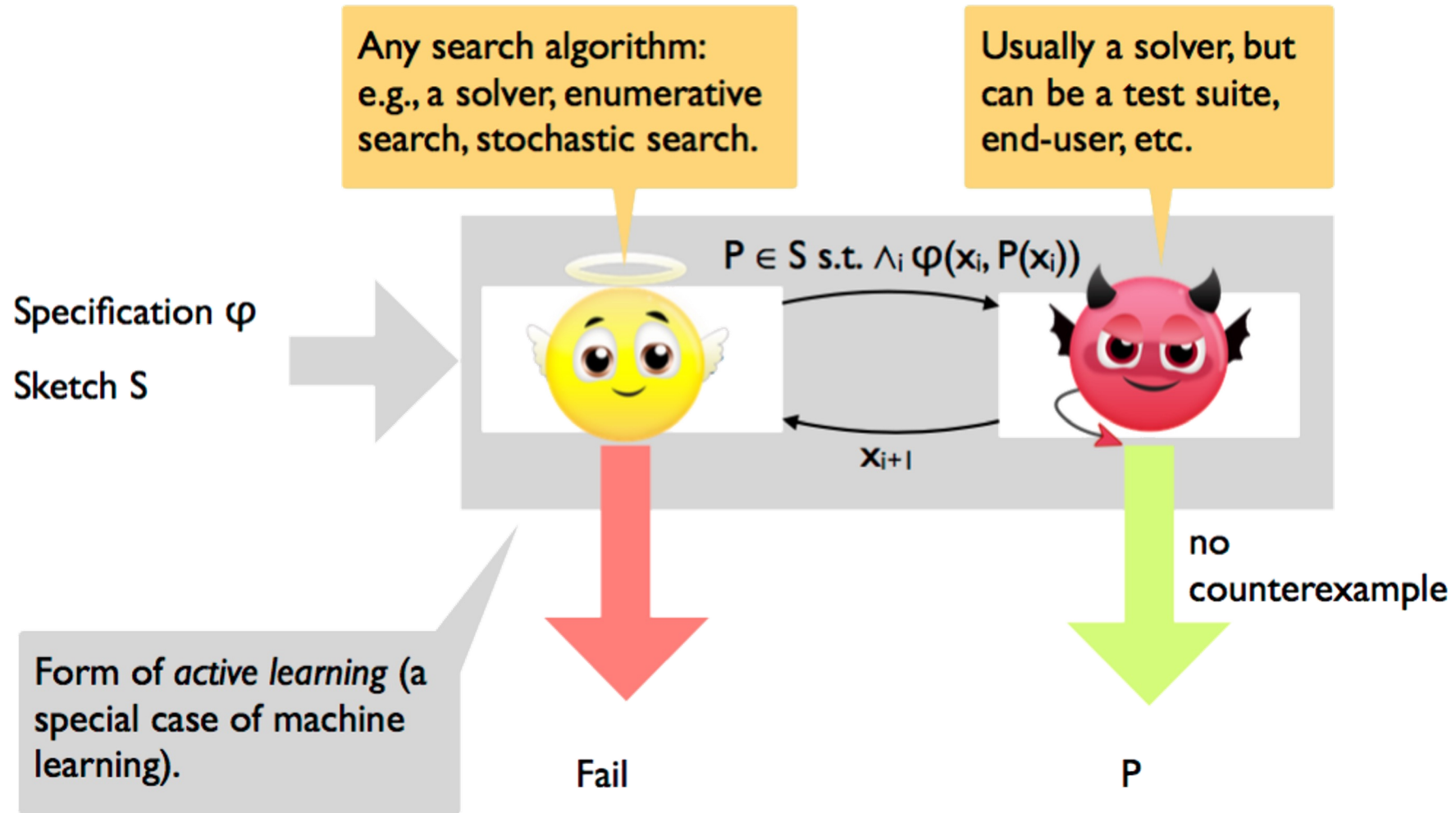
Program found by (symbolic) interpretation of a (space of) candidates, not by deriving the candidate.

So, to find a program, we need only an interpreter, not a sufficient set of derivation axioms.

Exercise: validate $P_{max}(list)$

- Given a candidate program $P_{max}(l)$ that finds the maximum number in a list l , how can you check if it is correct?

Overview of CEGIS



Partial or multi-modal specification of the desired program

`reg6 * 4 + 1`

Solves $\exists P . \varphi(x_1, P(x_1)) \wedge \dots \wedge \varphi(x_n, P(x_n))$ for representative inputs x_1, \dots, x_n

A program P from the given space of candidates that satisfies φ on all (usually bounded) inputs

`s4addl(reg6, 1)`

CEGIS: Counterexample-guided Inductive Synthesis
[Solar-Lezama et al., ASPLOS 06]

`expr :=
const | reg6 |
s4addl(expr, expr) |
...`

A syntactic *sketch* describing the shape of the desired program; defines the space of candidate programs to search. Can be tuned for performance.

Sketching intuition

Extend the language with two constructs

```
spec: int foo (int x) {  
    return x + x;  
}
```

$\phi(x, y): y = \mathbf{foo}(x)$

```
sketch: int bar (int x) implements foo {  
    return x << ??;  
}
```

?? substituted with an
int constant meeting ϕ

```
result: int bar (int x) implements foo {  
    return x << 1;  
}
```

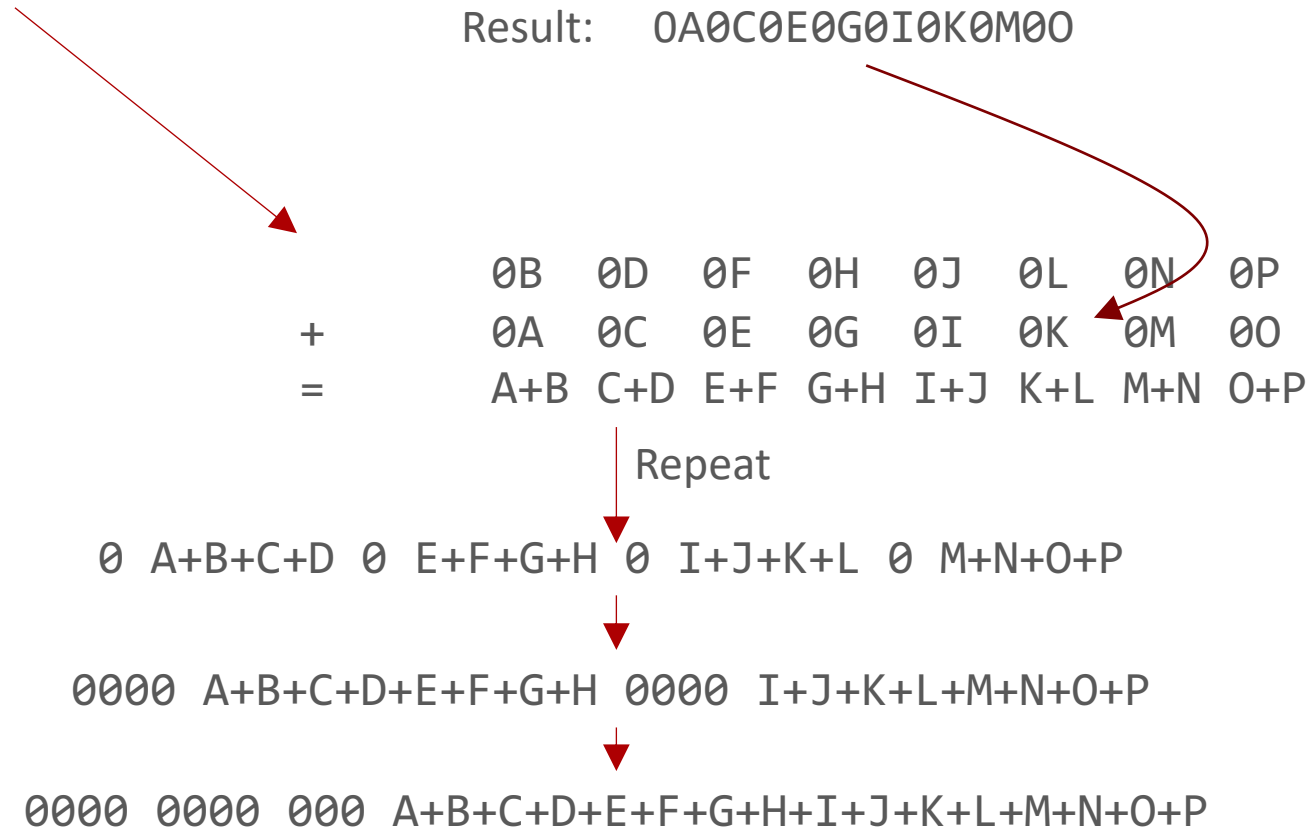
EXAMPLE: BIT COUNTING

```
1. bit[W] countBits(bit[W] x)
2. {
3.     int count = 0;
4.     for (int i = 0; i < W; i++) {
5.         if (x[i]) count++;
6.     }
7.     return count;
8. }
```

Intuition

Bit string: ABCDEFGHIJKLMNOP
Mask: 0101010101010101
Result: 0B0D0F0H0J0L0N0P

Bit string: ABCDEFGHIJKLMNOP
Bits >> 1: 0ABCDEFGHIJKLMNO
Mask: 0101010101010101
Result: 0A0C0E0G0I0K0M00



```
1.  bit[W] countSketched(bit[W] x)
2.      implements countBits {
3.      loop (??) {
4.          x = (x & ??) +
5.              ((x >> ??) & ?? );
6.      }
7.      return x;
8.  }
```



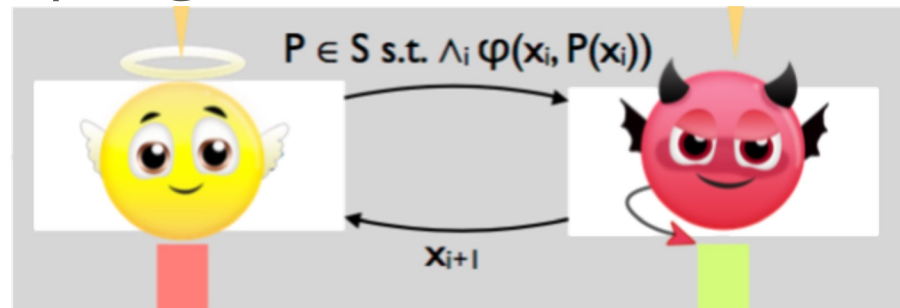
```
1.  bit[W] countSketched(bit[W] x)
2.  {
3.    x = (x & 0x5555) +
4.        ((x >> 1) & 0x5555);
5.    x = (x & 0x3333) +
6.        ((x >> 2) & 0x3333);
7.    x = (x & 0x0077) +
8.        ((x >> 4) & 0x0077);
9.    x = (x & 0x000F) +
10.       ((x >> 4) & 0x000F);
11.  return x;
12. }
```

Oracle-Guided Inductive Synthesis

1. Generalize CEGIS (counterexample-guided inductive synthesis)
 - From sketches to arbitrary programs
2. Synthesize programs from components

CEGIS: A Mathematical View

- Let's formalize Counterexample-Guided Inductive Synthesis (CEGIS)
- Consider a formalization of synthesizing a *max* function for lists
$$\exists P_{max} \forall l, m : P_{max}(l) = m \Rightarrow (m \in l) \wedge (\forall x \in l : m \geq x)$$
- CEGIS iterates between synthesis from examples and counterexample generation



- How do we generate a counterexample?

Counterexample generation, formalized

- Let's say we have a candidate program P_{max} . Does it meet the spec?

- Here's how that can be formalized:

$$\forall l, m : P_{max}(l) = m \Rightarrow (m \in l) \wedge (\forall x \in l : m \geq x)$$

- By De Morgan's Law, this is equivalent to disproving the negation:

$$\exists l, m : (P_{max}(l) = m) \wedge (m \notin l \vee \exists x \in l : m < x)$$

- This finds a list l and a corresponding incorrect output m
- Let's tweak this to generate the correct output, m^* :

$$\exists l, m^* : (P_{max}(l) \neq m^*) \wedge (m^* \in l) \wedge (\forall x \in l : m^* \geq x)$$

- We can use this to help generate the next version of P_{max}

Oracle-Guided Component-Based Program Synthesis (from examples)

- Goal: given a set of N components f_1, \dots, f_N and a set of T input/output pairs $\langle \alpha_0, \beta_0 \rangle \dots \langle \alpha_T, \beta_T \rangle$, synthesize a function f such that: $\forall i \in [0, T]: f(\alpha_i) = \beta_i$.
- We search for programs of a particular form:

Put inputs in variables

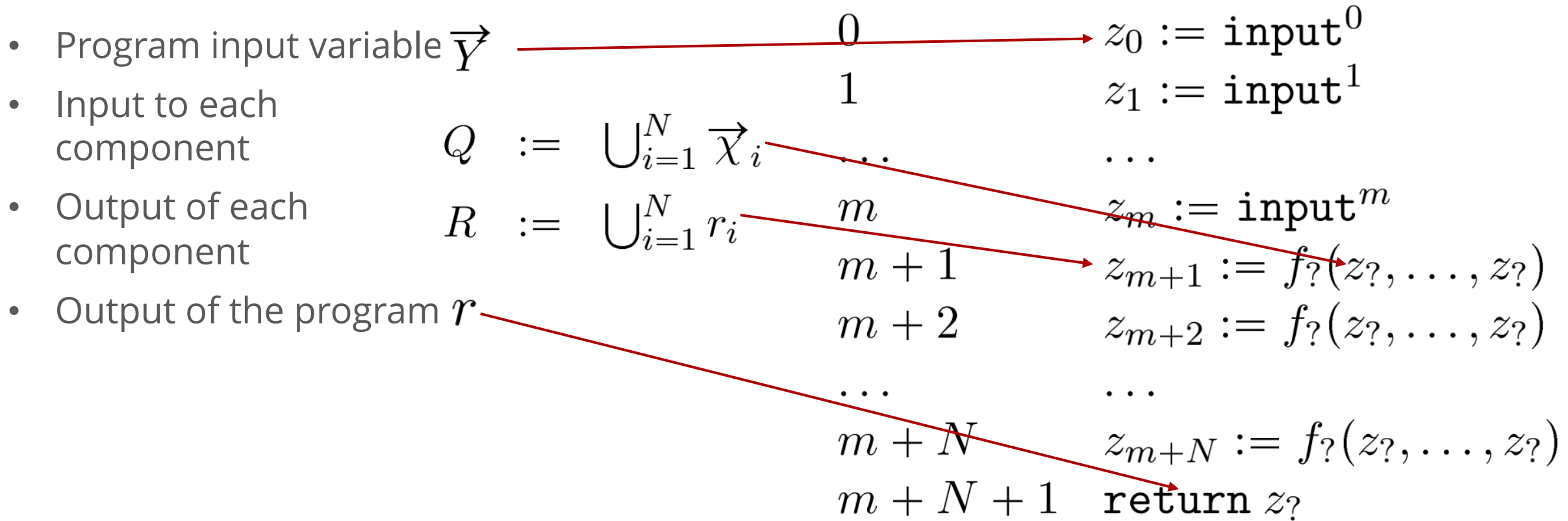
```
0      z0 := input0
1      z1 := input1
...
m      zm := inputm
m + 1  zm+1 := f?(z?, ..., z?)
m + 2  zm+2 := f?(z?, ..., z?)
...
m + N  zm+N := f?(z?, ..., z?)
m + N + 1  return z?
```

Compute N functions,
each of which has
arguments

Choices: fill in the ?s

- What order are the functions in?
- What variables are passed to each function?
- What variable is returned?

The program is defined by a set of variables



Program variables are specified by location variables

- Location variable l_x specifies where x is defined
- L is the set of location variables

$$L := \{l_x \mid x \in Q \cup R \cup \vec{Y} \cup r\}$$

(again: component inputs, component results, program inputs, and program result)

0	$z_0 := \text{input}^0$
1	$z_1 := \text{input}^1$
...	...
m	$z_m := \text{input}^m$
$m + 1$	$z_{m+1} := f^?(z^?, \dots, z^?)$
$m + 2$	$z_{m+2} := f^?(z^?, \dots, z^?)$
...	...
$m + N$	$z_{m+N} := f^?(z^?, \dots, z^?)$
$m + N + 1$	return $z^?$

Example of Location Variables

- Imagine we have one input and one component, +

- Here's a sample program:
0 $z_0 := \text{input}^0$
1 $z_1 := z_0 + z_0$
2 **return** z_1

- This can be specified by the location variables

$$\{l_{r_+} \mapsto 1, l_{\chi_+^1} \mapsto 0, l_{\chi_+^2} \mapsto 0, l_r \mapsto 1, l_Y \mapsto 0\}$$

Practice with Location Variable Encodings

Assume two components, $*$ and \ll , each of which takes two inputs and produces a single output. Provide a map which assigns values to location variables that describe the following straight-line code.

For your reference, the variables are: $\vec{Y} \ r \ \vec{X}_i \ r_i$

$z_0 = \text{input}_0$

$z_1 = \text{input}_1$

$z_2 = z_0 \ll z_1$ *// component \ll*

$z_3 = z_2 * z_2$ *// component $*$*

return z_2

Well-formedness constraints on the generated program

- Component inputs come from locations $0 \dots M$
 - $M = \text{number of inputs } |\vec{Y}| + \text{number of functions } N$

$$\bigwedge_{x \in Q} (0 \leq l_x < M)$$

- Component outputs defined after program inputs

$$\bigwedge_{x \in R} (|\vec{Y}| \leq l_x < M)$$

- One component per line

$$\bigwedge_{x, y \in R, x \neq y} (l_x \neq l_y)$$

- Component inputs are defined before use

$$\bigwedge_{i=1}^N \bigwedge_{x \in \vec{X}_i} l_x < l_{r_i}$$

0

1

...

m

$m + 1$

$m + 2$

...

$m + N$

$m + N + 1$

$z_0 := \text{input}^0$

$z_1 := \text{input}^1$

...

$z_m := \text{input}^m$

$z_{m+1} := f?(z?, \dots, z?)$

$z_{m+2} := f?(z?, \dots, z?)$

...

$z_{m+N} := f?(z?, \dots, z?)$

return $z?$

Functionality constraints

- Variables defined at the same location are the same (have the same value)
 - Basically: define value flow from definition to use $\bigwedge_{x,y \in Q \cup R \cup \vec{Y} \cup \{r\}} (l_x = l_y \Rightarrow x = y)$
- The program inputs and outputs match a test case pair
 - We repeat this for all test cases $(\alpha = \vec{Y}) \wedge (\beta = r)$
- Functional components obey their specification $(\bigwedge_{i=1}^N \phi_i(\vec{x}_i, r_i))$

Component-Based Synthesis, Overall

- We conjoin the well-formedness and functionality constraints into one big formula
- We have an SMT solver solve that formula
- The result is a witness, assigning integer values to each location variable
 - We can then convert the witness into a program
 - Line i of the program:

$$z_i = f_j(z_{\sigma_1}, \dots, z_{\sigma_\eta}) \text{ when } l_{r_j} == i \text{ and } \bigwedge_{k=1}^{\eta} (l_{\chi_j^k} == \sigma_k)$$

- We can then put this into a CEGIS loop:

