# Oracle-Guided Component-Based Program Synthesis

Reporter: Song Ao-qi

# Authors

**Susmit Jha**

Technical Director, NuSCI Research Group,

Computer Science Laboratory, SRI International

**Sumit Gulwani**

Research Manager/Principal Researcher @ Microsoft, Redmond

Affiliate Faculty Member @ Univ. of Washington

# Authors

**Sanjit A. Seshia**

Cadence Founders Chair Professor

Group in Logic and the Methodology of Science  University of California, Berkeley

**Ashish Tiwari**

Principal Researcher @ Microsoft, Redmond

# CONTENT

1. From program synthesis to search problem

   - Problem Definition

   - Encoding Programs

   - Constrains

   - Algorithm

2. Solve the search problem

3. Illustration on Running Example
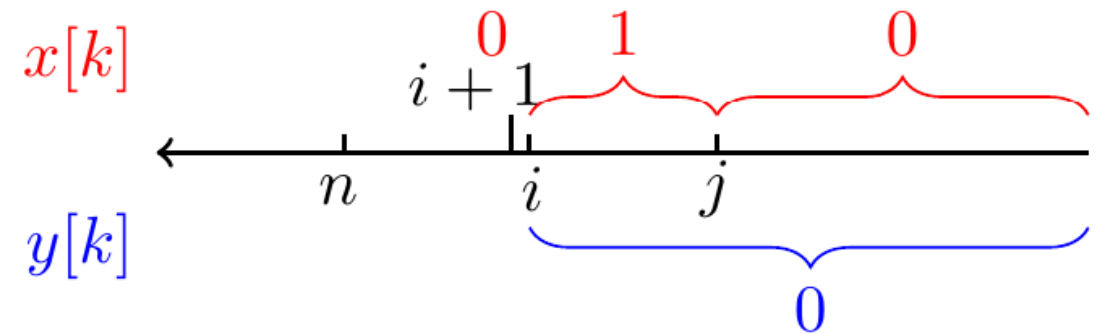
# Two examples

## ☐ Bit Manipulation

- **Given** a bit-vector x

- **Construct** a new bit-vector y that corresponds to x with the right most string of contiguous 1s turned off

$$\exists i, j. \{ 0 \leq i, j < n \wedge (\forall k. j \leq k \leq i \implies x[k] = 1)$$
$$\wedge (\forall k. 0 \leq k < j \implies x[k] = 0)$$
$$\wedge (x[i + 1] = 0 \vee i = n - 1)$$
$$\wedge (\forall k. i < k < n \implies x[k] = y[k])$$
$$\wedge (\forall k. 0 \leq k \leq i \implies y[k] = 0) \}$$

**specification**
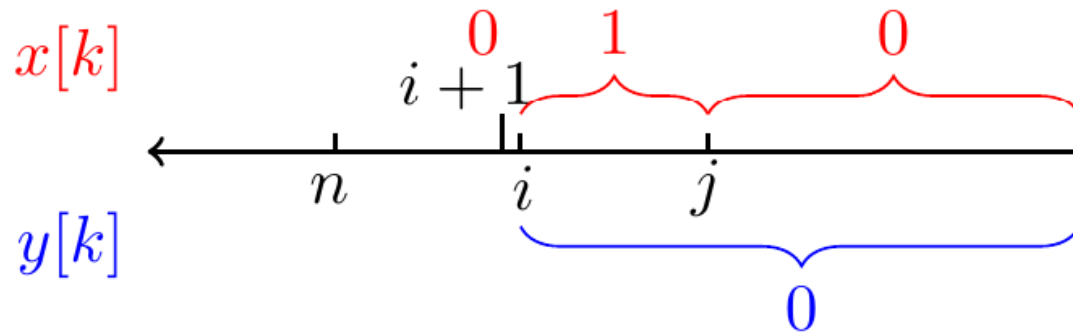


**sketch**

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Two examples

☐ **Bit Manipulation**



```
1 turnOffRightMostOneBitString (x)
2 { t₁ = x-1; t₂ = (x | t₁); t₃ = t₂+1;
3    t₄ = (t₃ & x); return t₄; }
```

- Oracle-guide: input-output example

- Component-based：a given set of base components

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Two examples

☐ **Deobfuscation**

```
 1 genStringObs(int input)
 2 {
 3    a1=1, a2=0; b1=1, b2=0; c1=0; c2=0;
 4    if (input == 0)
 5       { a1 = 0; a2 = 0; b1=0; b2 = 0; }
 6    else if (input == 1)
 7       { c1=0; c2 = 1; }
 8    else if (input == 2)
 9       { a1 =1; a2 = 0; c1=1; c2=1; }
10    else if (input == 3)
11       { b1 = 0; b2 = 0; c1=1; c2=1; }
12    else return NULL;
13    c = 2*c1 + c2;
14    if(c == 1) { return rot13("EPCG GB, 7"); }
15    else
16    if (c == 2) {
17       if (input * (input-1) % 2 == 0)
18          return rot13("EPCG GB", 7);
19       else
20          return rot13("RUYB", 4);
21    }
22    else {
23       if (b1 ⊕ b2) return rot13("ZNVY SEBZ",9)
24       else if ((a1 ⊕ a2) = (b1 ⊕ b2))
25             return rot13("RUYB", 4);
26       else return rot13("QNGN", 4);
27    }
28 }
```

```
29 rot13(char *buf, int sz)
30 {
31    char *buf1 = malloc((sz+1) * sizeof(char));
32    char a;
33    while (a =~ *buf)
34    {
35       *buf1 = (~a-1/(~((a | 32))/13*2-11)*13);
36       buf++; buf1++;
37    }
38    return buf1;
39 }
```

```
 1 genString(int input)
 2 { if(input == 0) return "EHLO";
 3    else if (input == 1) return "RCPT TO";
 4    else if (input == 2) return "MAIL FROM";
 5    else if (input == 3) return "DATA";
 6    else return NULL;
 7 }
```

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Problem Definition

**Base components (subprograms) → candidate program → desired program**

**Straight-line program (loop-free)**

$$P(\vec{I}):$$
$$O_{\pi_1} := f_{\pi_1}(\vec{V}_{\pi_1}); \quad \ldots \; ; \; O_{\pi_N} := f_{\pi_N}(\vec{V}_{\pi_N});$$
$$\textbf{return } O_{\pi_N};$$

**Assumptions**

- program P is using **all components** from the library

- program P is using each base component **only once**.

**Given**

**Oracle**:

- Validation oracle: validates the correctness of a candidate program

- I/O oracle        : generates input/output examples

**Specification** for base components (library)

$$\{\langle \vec{I}_i, O_i, \phi_i(\vec{I}_i, O_i)\rangle \mid i = 1, \ldots, N\}$$

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Encoding Programs

☐ **Location coding**

line number

| | |
|---|---|
| 0 | input1 |
| 1 | input2 |
| ... | |
| $\lvert \vec{I} \rvert - 1$ | input$\lvert \vec{I} \rvert - 1$ |
| $\lvert \vec{I} \rvert$ | the first line of program $P$ |
| ... | |
| $\lvert \vec{I} \rvert + N - 1$ | the last line of program $P$ |

**Coding of desired program/candidate program**

**Assumptions**

- All components have exactly one output program

- All inputs and outputs have the same type

**Location variables L**

$$L := \{ l_x \mid x \in \mathbf{P} \cup \mathbf{R} \}$$

$$\mathbf{P} := \bigcup_{i=1}^{N} \vec{I_i} \qquad \mathbf{R} := \bigcup_{i=1}^{N} \{O_i\} = \{O_1, \ldots, O_N\}$$

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Well-formedness constraints

□ **Syntactic constraints**

$$\psi_{\mathtt{wfp}}(L) \overset{\mathrm{def}}{=} \bigwedge_{x \in \mathbf{P}} (0 \leq l_x < M) \;\wedge\; \bigwedge_{x \in \mathbf{R}} (|\vec{I}| \leq l_x < M)$$

$$\wedge\; \psi_{\mathtt{cons}}(L) \;\wedge\; \psi_{\mathtt{acyc}}(L)$$

**Consistency constraints:**
$$\psi_{\mathtt{cons}} \overset{\mathrm{def}}{=} \bigwedge_{x,y \in \mathbf{R}, x \not\equiv y} (l_x \neq l_y)$$

**Acyclicity constraints:**
$$\psi_{\mathtt{acyc}} \overset{\mathrm{def}}{=} \bigwedge_{i=1}^{N} \bigwedge_{x \in \vec{I}_i, y \equiv O_i} l_x < l_y$$

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Well-formedness constraints

□ **semantic constraints**

$$\phi_{\texttt{func}}(L, \vec{I}, O) \stackrel{\text{def}}{=} \exists \mathbf{P}, \mathbf{R} \; \psi_{\texttt{wfp}}(L) \; \wedge \; \phi_{\texttt{lib}}(\mathbf{P}, \mathbf{R})$$
$$\wedge \; \psi_{\texttt{conn}}(L, \vec{I}, O, \mathbf{P}, \mathbf{R})$$

**Semantics of the base components:**
$$\phi_{\texttt{lib}}(\mathbf{P}, \mathbf{R}) \stackrel{\text{def}}{=} \left( \bigwedge_{i=1}^{N} \phi_i(\vec{I}_i, O_i) \right)$$

**Dataflow semantics:**
$$\psi_{\texttt{conn}}(L, \vec{I}, O, \mathbf{P}, \mathbf{R}) \stackrel{\text{def}}{=} \bigwedge_{x,y \in \mathbf{P} \cup \mathbf{R} \cup \vec{I} \cup \{O\}} (l_x = l_y \; \Rightarrow \; x = y)$$

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Well-formedness constraints

□ **Phase summary**

**Syntactic constrains:**

- Ordering of location variables [acyclicity]

- Uniqueness of location variables [consistency]

**Semantic constrains:**

- Relation between inputs and outputs of each component (implies connections between different components)

- "boundary conditions"

**Objective function:**

$$\phi_{\mathtt{func}}(L, \vec{I}, O) \overset{\mathrm{def}}{=} \exists \mathbf{P}, \mathbf{R} \; \psi_{\mathtt{wfp}}(L) \; \wedge \; \phi_{\mathtt{lib}}(\mathbf{P}, \mathbf{R})$$
$$\wedge \; \psi_{\mathtt{conn}}(L, \vec{I}, O, \mathbf{P}, \mathbf{R})$$

$$L := \{ l_x \; | \; x \in \mathbf{P} \cup \mathbf{R} \}$$

$$\phi_{func}(L, (\vec{I}, O)) == loss(\theta, (x, y))$$

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# I/O-behavioral Constraint &  Distinguishing Constraint

## ☐ I/O-behavioral Constraint

$$\text{Behave}_E(L) \overset{\text{def}}{=} \bigwedge_{(\alpha_j, \beta_j) \in E} \phi_{\text{func}}(L, \alpha_j, \beta_j)$$

where E is a set of input-output examples.

## ☐ Distinguishing Constraint

**target:** find a unique solution

**way:** data augmentation, but generative ones

$$\text{Distinct}_{E,L}(\vec{I}) \overset{\text{def}}{=} \exists L', O, O' \ \text{Behave}_E(L') \wedge \phi_{\text{func}}(L, \vec{I}, O)$$
$$\wedge \ \phi_{\text{func}}(L', \vec{I}, O') \ \wedge \ O \neq O'$$

~~How do we know which of O and O' is the true output? Or both are false.~~

~~Can not halt~~

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Algorithm

```
IterativeSynthesis():
```

1  `// Input:  Set of base components used in`
2  `// construction of Behave_E and Distinct_{E,L}`
3  `// Output:  Candidate Program`
4  $E := \{(\alpha_0, \mathcal{I}(\alpha_0))\}$ `//` $\alpha_0$ `is an arbitrary value for` $\vec{I}$
5  `while (1) {`
6  $\quad L := $ `T-SAT(Behave`$_E(L)$`);`
7  $\quad$ `if (`$L == \perp$`) return "Components insufficient";`
8  $\quad \alpha := $ `T-SAT(Distinct`$_{E,L}(\vec{I})$`);`
9  $\quad$ `if (`$\alpha == \perp$`) {`
10  $\quad\quad P := $ `Lval2Prog(`$L$`);`
11  $\quad\quad$ `if (`$\mathcal{V}(P)$`) return` $P$`;`
12  $\quad\quad$ `else return "Components insufficient"; }`
13  $E := E \cup \{\alpha, \mathcal{I}(\alpha)\};$ `}`

- T-SAT
- Lval2Prog()
- I/O oracle
- Validation oracle

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Oracle & components

## Bit manipulation

- Base component: standard library; extended library (user)

- I/O oracle: user

- Validation oracle: semantically unique candidate program to be correct program in practice; user

## Deobfuscation

- Base component: standard library; extended library (user)

- I/O oracle: obfuscated program

- Validation oracle: a program equivalence checking tool; user

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Optimization

**SMT solvers:** z3(python)

**Sampling:**

- Sampling Uniformly at Random

- Sampling With Bias (a priori)

```
ConstrainedRandomInput:
1  // cnt is a global variable initialized to 0
2  // K is a parameter (number of rightmost bits to se
3  if (cnt < 2^K) {
4     α := sample(Inputs);
5     α := Set rightmost K bits of α to cnt;
6     cnt := cnt + 1; }
7  else α := T-SAT(Distinct_{E,L}(I⃗));
```

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Illustration on Running Example

**A given set of components:** bit-wise logical operations; basic arithmetic operations

**I/O oracle:** the user

Iteration 1 Given the input/output pair (01011, 01000)

    [Behave_E] Candidate program 1: (x + 1) & (x − 1)

    [Distinct_E,L] Candidate program 2: (x + 1) & x

          with distinguishing input 00000

    [I/O oracle] (00000, 00000)

    E = {(01011, 01000), (00000, 00000)}

Iteration 2

    [Behave_E] Candidate program 1: −(¬x) & x

    ……

(((x − 1)|x) + 1)&x

```
IterativeSynthesis():
1    // Input:  Set of base components used in
2    // construction of Behave_E and Distinct_E,L
3    // Output:  Candidate Program
4    E := {(α_0, I(α_0))} // α_0 is an arbitrary value for Ī
5    while (1) {
6        L := T-SAT(Behave_E(L));
7        if (L == ⊥) return "Components insufficient";
8        α := T-SAT(Distinct_E,L(Ī));
9        if (α == ⊥) {
10           P := Lval2Prog(L);
11           if (V(P)) return P;
12           else return "Components insufficient"; }
13       E := E ∪ {α, I(α)}; }
```

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224

# Conclusion & Discussion

**Conclusion**

1.  Encoding program: Modeling Program Synthesis as a Search Problem

2.  Program synthesis based on oracle-guided learning from examples and SMT solvers.

**Discussion**

1.  Learning & reasoning

2.  Modeling:

    *   location coding, functional coding

    *   Discover more constrains

S. Jha, S. Gulwani, S. A. Seshia and A. Tiwari, "Oracle-guided component-based program synthesis," 2010 ACM/IEEE 32nd International Conference on Software Engineering, Cape Town, South Africa, 2010, pp. 215-224