



# Trace-Guided Inductive Synthesis of Recursive Functional Programs

YONGWEI YUAN, Purdue University, USA

ARJUN RADHAKRISHNA, Microsoft, USA

ROOPSHA SAMANTA, Purdue University, USA

We propose a novel trace-guided approach to tackle the challenges of ambiguity and generalization in synthesis of recursive functional programs from input-output examples. Our approach augments the search space of programs with recursion traces consisting of recursive subcalls of the programs. Our method is based on a new version space algebra (VSA) for succinct representation and efficient manipulation of pairs of recursion traces and programs that are consistent with each other. We have implemented this approach in a tool called SyRUP and evaluated it on benchmarks from prior work. Our evaluation demonstrates that SyRUP not only requires fewer examples to achieve a certain success rate than existing synthesizers, but is also less sensitive to the quality of the examples.

CCS Concepts: • **Software and its engineering** → **Automatic programming**; **Functional languages**; **Recursion**; **Programming by example**.

Additional Key Words and Phrases: Program Synthesis, Recursive Functional Programs

## ACM Reference Format:

Yongwei Yuan, Arjun Radhakrishna, and Roopsha Samanta. 2023. Trace-Guided Inductive Synthesis of Recursive Functional Programs. *Proc. ACM Program. Lang.* 7, PLDI, Article 141 (June 2023), 24 pages. <https://doi.org/10.1145/3591255>

## 1 INTRODUCTION

The last decade has witnessed significant advances in synthesis of recursive functional programs [Albarghouthi et al. 2013; Farzan and Nicolet 2021; Feser et al. 2015; Itzhaky et al. 2021; Kneuss et al. 2013; Lee and Cho 2023; Lubin et al. 2020; Miltner et al. 2022; Osera and Zdancewic 2015; Polikarpova et al. 2016], and, in particular, in inductive synthesis of recursive programs from input-output (I/O) examples [Albarghouthi et al. 2013; Feser et al. 2015; Lubin et al. 2020; Miltner et al. 2022; Osera and Zdancewic 2015]. The latter problem, however, is far from solved and continues to pose several challenges that are best illustrated through an example. Let us consider a synthesis task where the goal is a recursive program that adds two natural numbers and is consistent with an I/O example  $(2, 1) \rightsquigarrow 3$ . In Fig. 1, we list several candidate programs in the search space of a typical synthesizer for functional programs, along with the corresponding *recursion traces*; informally, recursion traces capture sequences of recursive subcalls of a program induced by an I/O example. While all programs are consistent with the example  $(2, 1) \rightsquigarrow 3$ , notice that only the programs annotated with a  $\checkmark$  are correct for all possible program inputs. So, how would the typical synthesizer resolve the ambiguity between these programs and pick one that generalizes to other inputs?

Authors' addresses: Yongwei Yuan, yuan311@purdue.edu, Purdue University, USA; Arjun Radhakrishna, Microsoft, USA, arradha@microsoft.com; Roopsha Samanta, Purdue University, USA, roopsha@purdue.edu.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2023 Copyright held by the owner/author(s).

2475-1421/2023/6-ART141

<https://doi.org/10.1145/3591255>

Recursion Trace	Program	
$(2,1) \rightsquigarrow 3$ $\downarrow$ $(1,1) \rightsquigarrow 2$ $\downarrow$ $(0,1) \rightsquigarrow 1$	<pre>let rec nat_add m n =   match m with     Zero -&gt; n     Succ m' -&gt;     Succ (nat_add m' n)</pre>	✓
$(2,1) \rightsquigarrow 3$ $\downarrow$ $(1,1) \rightsquigarrow 3$ $\downarrow$ $(0,1) \rightsquigarrow 3$	<pre>let rec nat_add1 m n =   match m with     Zero -&gt; Succ (Succ n)     Succ m' -&gt;     nat_add1 m' n</pre>	✗
$(2,1) \rightsquigarrow 3$ $\downarrow$ $(1,0) \rightsquigarrow 1$	<pre>let rec nat_add2 m n =   match m, n with     Zero, Zero -&gt; Zero     Succ m', Zero -&gt; m     Zero, Succ n' -&gt; n     Succ m', Succ n' -&gt;     Succ (Succ (nat_add2 m' n'))</pre>	✓
$(2,1) \rightsquigarrow 3$	<pre>let rec nat_add3 m n =   Succ (Succ n)</pre>	✗

Fig. 1. Consistent trace-program pairs for an I/O example  $(2, 1) \rightsquigarrow 3$ .

- The synthesizer may use everyone’s favorite inductive bias—*Occam’s razor*—to pick the “smallest” program among the candidates [Barke et al. 2020; Gulwani 2011; Wang et al. 2017]. In our example, however, program `nat_add3` is the smallest and is obviously incorrect; further, programs `nat_add` and `nat_add1` are about the same size and only one of them is correct.
- Or, the synthesizer may require the user to provide a *complete* recursion trace, say  $((2, 1) \rightsquigarrow 3 \rightarrow (1, 1) \rightsquigarrow 2 \rightarrow (0, 1) \rightsquigarrow 1)$ , and use the recursion trace to find a program, such as `nat_add`, consistent with the trace [Albarghouthi et al. 2013; Feser et al. 2015; Osera and Zdancewic 2015]. This is an instance of programming by demonstration [Lau et al. 2003] and trace-based synthesis [Chugh et al. 2016], and helps the synthesizer by expecting more information about the behavior of the target program. Unfortunately, requiring users to provide a complete recursion trace, or, essentially, a *trace-complete* set of examples can be prohibitive.

In this paper, we propose a new *trace-guided* synthesis approach to tackle the challenges of trace-completeness, generalization, and ambiguity-resolution in inductive synthesis of recursive programs. Our approach augments the search space of programs with recursion traces, thereby using recursion traces as an inductive bias without burdening users. Given a set of I/O examples, our method *jointly* explores the space of recursion traces and programs that are consistent with each other, keeping track of consistent trace-program pairs in a version space [Lau et al. 2003; Mitchell 1982], and using a ranking function over such trace-program pairs to pick the most likely program.

The key technical contribution of this work is a new version space algebra (VSA) for succinct representation and efficient manipulation of consistent trace-program pairs. Our VSA has several important properties. First, it allows us to explore the space of recursion traces and programs *incrementally*. In each step, we leverage *angelic semantics* [Bodik et al. 2010; Floyd 1967] to track possible valuations of recursive subcalls that a consistent program might make in the *next* recursion step. Second, it allows us to construct a version space of consistent trace-program pairs *compositionally*. The VSA operators enable us to compose the version space for the current recursion step with version spaces for subsequent recursion steps through a clever manipulation that involves *consing* recursion traces and *unifying* sets of programs. Last, but not the least, all VSA operators are designed to be *consistency-preserving*, thereby eliminating the need for backtracking.

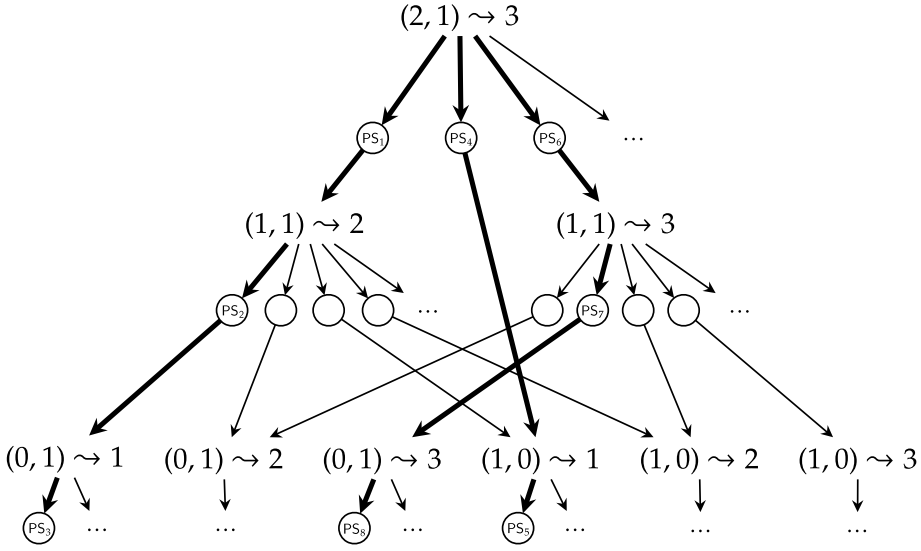


Fig. 2. An AND/OR graph representing a trace-indexed VS for I/O example  $(2, 1) \rightsquigarrow 3$ . Each  $\bigcirc$  is an AND node representing a set of programs and each  $i \rightsquigarrow o$  is an OR node representing a synthesis problem.

We have implemented our approach in a tool called SyRUP<sup>1</sup> and evaluated it on 43 programming tasks collected from prior work [Osera and Zdanczewicz 2015]. Our experiments extensively evaluate SyRUP’s ability to generalize across different classes of I/O example sets. SyRUP not only requires *fewer* examples to achieve a certain *success rate* than the two state-of-the-art inductive synthesizers SMYTH [Lubin et al. 2020] and BURST [Miltner et al. 2022], but is also *less sensitive* to the *quality* of the examples.

To summarize, our paper makes the following contributions.

- (1) We present a new trace-guided procedure for inductive synthesis of recursive programs to tackle the challenges of trace-completeness, generalization, and ambiguity-resolution. Our procedure is sound and complete modulo a notion of observational equivalence (Sec. 5).
- (2) We define a new VSA for succinct representation and efficient manipulation of the space of recursion traces and programs that are consistent with each other (Sec. 4).
- (3) We develop a tool SyRUP and extensively evaluate it on synthesis tasks. We demonstrate that SyRUP outperforms prior work in multiple aspects (Sec. 6, Sec. 7).

## 2 TRACE-GUIDED SYNTHESIS: ILLUSTRATIVE OVERVIEW

In this section, we illustrate key aspects of our trace-guided synthesis algorithm using our running example of synthesis from one I/O example  $(2, 1) \rightsquigarrow 3$ .

**Recursion Traces.** Given an I/O example and a recursive program, a *recursion trace* (or, in short, *trace*) is a tree that witnesses the execution of the program on the given example. In particular, the nodes  $i \rightsquigarrow o$  of a trace correspond to recursive subcalls with input  $i$  and output  $o$ , where the root node represents the given I/O example and the descendants of a node represent subsequent recursive subcalls made in the program execution. A trace is *linear* iff each node has at most one child. Given an I/O example  $i \rightsquigarrow o$ , we say that a *trace-program pair*  $(Tr, P)$  is *consistent* iff the nodes

<sup>1</sup>SyRUP stands for “Synthesizer for Recursive Programs”.

of the trace  $Tr$  witness the recursive calls necessary to complete the execution of  $P$  on input  $i$  as described above.

For instance, for the example  $(2, 1) \rightsquigarrow 3$ , each row in Fig. 1 presents a consistent trace-program pair, where each trace is linear.

**Trace-Indexed Version Space.** Our approach leverages (recursion) traces for inductive synthesis of recursive programs using a *trace-indexed version space* (VS). Given an example  $i \rightsquigarrow o$ , a trace-indexed VS consists of pairs  $(Tr, PS)$  of traces and sets of programs such that for each program  $P$  in  $PS$ , the trace-program pair  $(Tr, P)$  is consistent.

To help illustrate how our method constructs a trace-indexed VS, we use a rooted AND/OR graph as shown in Fig. 2. Each  $(Tr, PS)$  pair in the trace-indexed VS for example  $(2, 1) \rightsquigarrow 3$  corresponds to a rooted tree from the root OR node to leaf AND nodes in the AND/OR graph. For instance, `nat_add` and its trace in the first row of Fig. 1 corresponds to the leftmost path in the AND/OR graph (indicated using thick edges), and so on. Note that because each trace in this example is linear, the rooted tree corresponding to each  $(Tr, PS)$  pair in the trace-indexed VS is simply a path.

We now illustrate how our method constructs the trace-indexed VS for the leftmost path in Fig. 2. Each OR node  $i \rightsquigarrow o$  represents a new synthesis problem, defined by an example  $i \rightsquigarrow o$ , with the root OR node representing the original synthesis problem defined by  $(2, 1) \rightsquigarrow 3$ . The children of an OR node  $i \rightsquigarrow o$  are AND nodes (denoted by circles) and the children of an AND node are OR nodes. Each AND node represents a program set containing programs that are consistent with the parent OR-node, *assuming* the children OR-nodes. For instance, each program  $P$  in the program set  $PS_1$  is consistent with the original example  $(2, 1) \rightsquigarrow 3$  given the assumption (that  $P$  is also consistent with)  $(1, 1) \rightsquigarrow 2$ . An example of a program in the program set  $PS_1$  is as shown below:

```
let rec nat_add4 m n =
  match m with
  | Zero -> ??
  | Succ m' -> Succ (nat_add4 m' n)
```

where `??` can be substituted with any well-typed expression. This program belongs to program set  $PS_1$  because given the assumption `nat_add 1 1 = 2`, we have `nat_add 2 1` equals `Succ (nat_add 1 1)` equals 3, as expected.

Similarly, the program set  $PS_2$  contains programs that are consistent with  $(1, 1) \rightsquigarrow 2$ , given the assumption  $(0, 1) \rightsquigarrow 1$ . In particular, the program set  $PS_2$  also contains the partial program shown above; given the assumption `nat_add 0 1 = 1`, we have `nat_add 1 1` equals `Succ (nat_add 0 1)` equals 2, as expected.

Finally, the program set  $PS_3$  contains programs that are consistent with  $(0, 1) \rightsquigarrow 1$ , and includes the following program that does not involve any recursive subcalls:

```
let rec nat_add5 m n =
  match m with
  | Zero -> n
  | Succ m' -> ??
```

where `??` can be substituted by any well-typed expression as before.

The  $(Tr, PS)$  pair for the leftmost path can now be constructed by “combining” the assumptions into the trace  $(2, 1) \rightsquigarrow 3 \rightarrow (1, 1) \rightsquigarrow 2 \rightarrow (0, 1) \rightsquigarrow 1$ , and by “intersecting” the program sets  $PS_1$ ,  $PS_2$ , and  $PS_3$ . Notice that this intersection contains the correct program `nat_add` in Fig. 1.

By exhaustively exploring other rooted trees in the AND/OR graph, our method constructs a trace-indexed VS that consists of  $(Tr, PS)$  pairs corresponding to all rooted trees in the AND/OR graph by *splicing together assumptions and program sets*. In particular, our method uses operators defined in our version space algebra to do these manipulations efficiently.

*Inferring Assumptions.* The exploration of an AND/OR graph, and thus the construction of a trace-indexed VS, is driven by a procedure for inferring assumptions about the target program. This involves conjecturing new I/O examples that the target program is expected to be consistent with, based on the synthesis problem defined by the current OR node. Let's revisit the assumption  $\text{nat\_add4 } 1 \ 1 = 2$  from earlier. To tackle the synthesis problem defined by  $(2, 1) \rightsquigarrow 3$ , we not only search for candidate programs like  $\text{nat\_add4}$ , but also infer the output of the recursive subcall  $\text{nat\_add4 } 1 \ 1$  such that  $\text{Succ}(\text{nat\_add4 } 1 \ 1)$  equals to 3. In this case, the inference is trivially done by “deconstructing”  $3 \equiv \text{Succ}(\text{Succ}(\text{Succ}(\text{Zero})))$ . In general, the inference may be more involved. Consider an alternative candidate program for the synthesis problem:

```
let rec nat_add5 m n =
  match m with
  | Zero -> ??
  | Succ m' -> plus_one (nat_add5 m' n)
```

Here,  $\text{plus\_one}$  is a background function with the same semantics as the constructor  $\text{Succ}$ . The inverse semantics of background functions, particularly recursive ones, are not always readily available. Hence, we symbolically encode the semantics of  $\text{plus\_one}$  in an SMT solver, and query the solver for the output  $\alpha_{(1,1)}$  of the recursive subcall on  $(1, 1)$  such that  $\text{plus\_one } \alpha_{(1,1)} = 3$ .

**Trace-Based Ranking Function.** Given a trace-indexed VS for an I/O example, how do we pick programs that generalize better? As argued in Sec. 1, simply choosing a program with the “smallest” size may not always suffice. Fortunately, the recursion traces in the trace-indexed VS encode rich *semantic* features of the corresponding programs. For instance, one may easily identify the recursion scheme used in a program from the shape of its recursion trace; if we simply prefer programs with longer traces, then the program  $\text{nat\_add}$  ranks higher than the smaller non-recursive program  $\text{nat\_add3}$  in Fig. 1.

### 3 PROBLEM DEFINITION AND SOLUTION OVERVIEW

program	P	::=	$\text{fix } f(x_{\text{in}}).e$
expression	$e, g$	::=	$x \mid e_1 (e_2) \mid C() \mid C(e)$ $\mid \text{match } x \text{ with } \{C_1(x_1) \rightarrow e_1, \dots, C_n(x_n) \rightarrow e_n\}$ $\mid \text{cond } \{g_1 \rightarrow e_1, \dots, g_n \rightarrow e_n\}$
value	$i, o, v$	::=	$\text{true} \mid \text{false} \mid C() \mid C(v)$

$\frac{e_2 \Downarrow v_2 \quad [x_{\text{in}} := v_2, f := (\text{fix } f(x_{\text{in}}).e_1)] e_1 \Downarrow v}{(\text{fix } f(x_{\text{in}}).e_1) (e_2) \Downarrow v} \text{EREC}$	$\frac{e \Downarrow v}{C(e) \Downarrow C(v)} \text{ECTOR}$
$\frac{g \Downarrow \text{true} \quad e \Downarrow v}{\text{cond } \{\dots g \rightarrow e \dots\} \Downarrow v} \text{ECOND}$	$\frac{e' \Downarrow C(v') \quad [x := v'] e \Downarrow v}{\text{match } e' \text{ with } \{\dots C(x) \rightarrow e \dots\} \Downarrow v} \text{EMATCH}$

Fig. 3. tinyML: Syntax and Semantics.

**Algorithm 1** Top-Level Synthesis Algorithm**Input:** I/O examples  $\mathcal{E}$ **Input:** Layered search space  $(\mathcal{L}_j)_{j \in \mathbb{N}}$ **Output:** Program  $P$ 

```

1: procedure SYNTHESIZE( $\mathcal{E}$ ,  $(\mathcal{L}_j)_{j \in \mathbb{N}}$ )
2:   for  $j \in \mathbb{N}$  do
3:      $\Psi \leftarrow \text{LEARN}(\mathcal{E}, \mathcal{L}_j)$ 
4:     if  $\Psi \neq \emptyset$  then
5:        $P \leftarrow \text{RANK}(\Psi)$ 
6:   return  $P$ 

```

Given a set  $\mathcal{E}$  of I/O examples, our goal is to synthesize a recursive functional program  $P$  of the form `fix  $f(x_{\text{in}}).e$`  in a search space defined by a DSL. In what follows, we present our target DSL, formalize our problem statement, and present our top-level synthesis algorithm.

Our target DSL `tinyML` is a functional ML-like programming language with recursive functions, algebraic data types, and pattern matching (see Fig. 3). The syntax of expressions is mostly standard, with variables  $x$ , function applications  $e_1(e_2)$ , and constructors  $C$ . Both function applications and constructors may take multiple arguments (omitted in Fig. 3). To ensure that programs terminate, recursive subcalls are only permitted on values that are strictly smaller than the input. Expressions include a restricted form of `match`, where scrutinees can only be variables, and patterns are required to be non-overlapping and exhaustive, i.e., there is exactly one branch for each constructor of the data type being matched. Finally, expressions include conditional expressions `cond` consisting of a set of *guarded* expressions, where guards are essentially Boolean expressions that do not permit recursive subcalls to  $f$ . The guards in a `cond` expression are required to be mutually-exclusive and exhaustive; this ensures that `cond` expressions are equivalent in expressiveness to the more commonly used `if` expressions. The program semantics is standard.

**Definition 3.1** (Program Correctness). Given a program  $P$  in `tinyML`,  $P$  is consistent w.r.t. a set  $\mathcal{E}$  of I/O examples, denoted  $P \models \mathcal{E}$ , if and only if  $P(i) \Downarrow o$  for each I/O example  $i \mapsto o \in \mathcal{E}$ .

Let us denote by  $\mathcal{L}$  the set of programs generated by `tinyML`.

**Definition 3.2** (Problem Definition). Given a set  $\mathcal{E}$  of I/O examples, find a program  $P$  such that  $P \in \mathcal{L}$  and  $P \models \mathcal{E}$ .

Our top-level synthesis procedure is presented in Algo. 1. Because  $\mathcal{L}$  is unbounded, we follow common practice in program synthesis literature to both bound the search space and prioritize search. Thus, we assume that  $\mathcal{L}$  is *layered* into a sequence of finite subspaces  $(\mathcal{L}_1, \mathcal{L}_2, \dots)$ , where each subspace  $\mathcal{L}_j$  is obtained from  $\mathcal{L}$  by constraining the value of some suitable program parameter(s), and the ordering of subspaces encodes preference based on parameter values. For instance, a popular layering strategy is to use *program size* to bound the search space and prefer *smaller* programs by ordering subspaces by increasing size.

**Definition 3.3** (Layered Search Space<sup>2</sup>). Search space  $\mathcal{L}$  is defined to be *layered* by a sequence  $(\mathcal{L}_1, \mathcal{L}_2, \dots)$  of finite subsets  $\mathcal{L}_j$  of  $\mathcal{L}$ , denoted  $\mathcal{L} = (\mathcal{L}_j)_{j \in \mathbb{N}}$ , iff  $\mathcal{L} = \bigcup_{j \in \mathbb{N}} \mathcal{L}_j$ .

Given a set  $\mathcal{E}$  of examples and a layered search space  $\mathcal{L} = (\mathcal{L}_j)_{j \in \mathbb{N}}$ , procedure SYNTHESIZE (1) iterates through each finite subspace  $\mathcal{L}_j$  of  $\mathcal{L}$ , (2) learns a (representation of a) set  $\Psi$  of programs

<sup>2</sup>Notice that, unlike conventional notions of layering, say based on program size, we do not require  $\mathcal{L}_{j+1} \supseteq \mathcal{L}_j$ . As we shall see in Sec. 5 and Sec. 7, our relaxed notion of layering offers more flexibility in exploring the search space  $\mathcal{L}$ .

in subspace  $\mathcal{L}_j$  such that  $P \models \mathcal{E}$  for each  $P \in \Psi$ , and (3) returns a program in  $\Psi$  that maximizes a ranking function on  $\mathcal{L}$ . In any iteration, if  $\Psi$  is empty, the procedure moves onto the next subspace.

Most existing solutions [Barke et al. 2020; Lee 2021; Lubin et al. 2020; Miltner et al. 2022] for the program synthesis problem in Def. 3.3 are instances of Algo. 1 and iteratively increase the upper bound on some program parameter (such as AST size or depth) to systematically expand the search space. The key distinguishing elements lie in the design of the layered search space, the LEARN procedure, and the ranking function.

In what follows, we formalize our version space algebra to represent and manipulate sets of programs in Sec. 4. We then present the definition of our layered search space along with the LEARN and RANK procedures for our program representation in Sec. 5 and Sec. 6 respectively.

## 4 TRACE-INDEXED EXTENDED VERSION SPACES

In this section, we present the definitions required to describe our synthesis algorithm.

### 4.1 Trace-Witnessed Semantics

**Frames and Recursion Traces.** We use *frames* to represent the behavior within the current instance of a call to the recursive function  $f$  and *recursion traces* to represent the behavior of the full tree of recursive sub-calls. A frame  $Fr$  is denoted as  $(\Delta \leadsto v) \rightarrow \text{Asm}$  where: (a)  $\Delta$  is a sequence of substitutions  $x_j := v_j$ , called valuation context, with  $[\Delta] e$  being the result of applying substitutions in  $e$ , (b)  $v$  is a target value, and (c)  $\text{Asm} = [i_1 \leadsto o_1 \ \dots \ i_n \leadsto o_n]$  is an *assumption sequence* of I/O examples  $i_j \leadsto o_j$  where we also have that  $i_j = i_k \implies o_j = o_k$ . A recursion trace  $Tr$  is given by a tree  $(i \leadsto o) \rightarrow Trs$  where: (a) the root  $\text{Root}(Tr)$  is an I/O example  $(i \leadsto o)$ , and (b) the (possibly empty) sequence of recursion traces  $Trs = [Tr_1 \ \dots \ Tr_n]$  are the children. We assume that the set of I/O examples in a recursion trace is consistent, i.e., no two nodes are labeled with examples having the same input but differing outputs.

For brevity, frame  $((x_{in} := i) \leadsto o) \rightarrow \text{Asm}$  is written as  $(i \leadsto o) \rightarrow \text{Asm}$  for the rest of the paper.

**Example 4.1.** Consider the tinyML program  $P \equiv (\text{fix } f(x_{in}).e)$  where  $e \equiv \text{match } x_{in} \text{ with } \{\text{nil} \rightarrow 0, \text{cons}(h, t) \rightarrow h + f(t)\}$ . On the input  $i = \text{cons}(5, \text{cons}(10, \text{nil}))$ , the program  $P$  returns 15. The invocation to  $f(i)$  makes one direct recursive sub-call (i.e.,  $f(\text{cons}(10, \text{nil}))$ ), which in turn makes another recursive sub-call  $f(\text{nil})$ . This behavior is represented by the frame  $Fr = (i \leadsto 15) \rightarrow [\text{cons}(10, \text{nil}) \leadsto 10]$  and the recursion trace  $Tr = (i \leadsto 15) \rightarrow [( \text{cons}(10, \text{nil}) \leadsto 10 ) \rightarrow [(\text{nil} \leadsto 0) \rightarrow []]]$ . Here, the frame  $Fr$  captures the local behavior, i.e., that calling  $f$  on  $i$  returns 15 while making one direct recursive sub-call characterized by  $\text{cons}(10, \text{nil}) \leadsto 10$ ; while the recursion trace  $Tr$  characterizes the full nested and thus global recursion behavior.  $\square$

**Example 4.2.** Consider a standard recursive implementation  $e^*$  for computing Fibonacci Numbers:  $e^* \equiv \text{cond } \{x_{in} = 0 \rightarrow 1, x_{in} = 1 \rightarrow 1, x_{in} > 1 \rightarrow f(x_{in} - 2) + f(x_{in} - 1)\}$ . Similarly, the recursive sub-call  $(\text{fix } f(x_{in}).e^*) (3)$  is characterized by the frame  $(3 \leadsto 3) \rightarrow [1 \leadsto 1, 2 \leadsto 2]$  and the trace  $Tr^* = (3 \leadsto 3) \rightarrow [(1 \leadsto 1) \rightarrow [], (2 \leadsto 2) \rightarrow [(0 \leadsto 1) \rightarrow [], (1 \leadsto 1) \rightarrow []]]$ . Note that  $Tr^*$  (unlike  $Tr$  from Example 4.1) is branching, i.e., that a node has more than 1 child corresponding to multiple direct recursive sub-calls.  $\square$

**Angelic Semantics and Trace-Witnessed Semantics.** The angelic semantics and trace-witnessed semantics formalize the idea of frames and traces representing local and global recursion behavior as seen in Example 4.1. The rules for angelic and trace-witnessed semantics are shown in Fig. 4. Intuitively, the angelic semantics evaluate a term left-to-right in a standard manner, except that when a recursive function call  $f(e)$  is encountered, its output value is obtained from the assumption sequence  $\text{Asm}$ . Formally, the judgment  $e \Downarrow_{\text{Asm}}^{\text{Asm}'} v$  asserts that: (a) evaluating the expression  $e$  requires  $n = |\text{Asm}| - |\text{Asm}'|$  direct recursive calls, (b) the arguments to the recursive calls and their



$$\begin{array}{c}
\frac{e \Downarrow v}{e \Downarrow_{\text{Asm}}^{\text{Asm}} v} \text{AEBASE} \qquad \frac{g \Downarrow_{\text{Asm}}^{\text{Asm}}, \text{true} \quad e \Downarrow_{\text{Asm}'}^{\text{Asm}'}, v}{\text{cond } \{\dots g \rightarrow e \dots\} \Downarrow_{\text{Asm}''}^{\text{Asm}'', v} v} \text{AECOND} \\
\\
\frac{e \Downarrow_{\text{Asm}'}^{\text{Asm}'}, v}{C(e) \Downarrow_{\text{Asm}}^{\text{Asm}}, C(v)} \text{AECTOR} \qquad \frac{e \Downarrow_{\text{Asm}'}^{\text{Asm}'}, C(v) \quad [x := v] e \Downarrow_{\text{Asm}''}^{\text{Asm}'', v}}{\text{match } e \text{ with } \{\dots C(x) \rightarrow e \dots\} \Downarrow_{\text{Asm}''}^{\text{Asm}'', v} v} \text{AEMATCH} \\
\\
\frac{e \Downarrow_{\text{Asm}'}^{\text{Asm}'}, i \quad \text{Asm}' = [i \rightsquigarrow o] \sqcup \text{Asm}''}{f(e) \Downarrow_{\text{Asm}''}^{\text{Asm}'', o}} \text{AEREC} \qquad \frac{\text{Asm} = [\dots \text{Root}(\text{Tr}_j) \dots] \quad [x := i] e \Downarrow_{\square}^{\text{Asm}} o \quad \forall j. (\text{fix } f(x).e) \models \text{Tr}_j}{(\text{fix } f(x).e) \models ((i \rightsquigarrow o) \rightarrow [\dots \text{Tr}_j \dots])} \text{TEREC}
\end{array}$$

Fig. 4. Angelic semantics and trace-witnessed semantics of a recursive program  $P = \text{fix } f(x).e$ .

return values form the first  $n$  I/O examples in  $\text{Asm}$ , and (c)  $\text{Asm}'$  is the suffix of  $\text{Asm}$  excluding the first  $n$  I/O examples. And we write  $\text{Asm} \setminus \text{Asm}'$  to denote the assumption sequence obtained from removing the suffix  $\text{Asm}'$  from  $\text{Asm}$ .

The trace-witnessed semantics build on the angelic semantics. Intuitively, Rule **TERec** combines the local behavior (the 2<sup>nd</sup> premise stating that  $f$  angelically evaluates  $i$  to  $o$  making assumptions  $\text{Asm}$ ) and the global behavior (the 3<sup>rd</sup> premise stating that for each assumption  $i_j \rightsquigarrow o_j$  in  $\text{Asm}$ , there is a trace  $\text{Tr}_j$  witnessing it).

**Example 4.3.** Consider the program  $P$ , frame  $\text{Fr} = (i \rightsquigarrow 15) \rightarrow \text{Asm}$ , and  $\text{Tr}$  from Example 4.1. Using the angelic semantics, we can derive  $e \Downarrow_{\square}^{\text{Asm}} 15$ , i.e., assuming  $f(\text{cons}(10, \text{nil}))$  evaluates to 10, we have that  $e$  evaluates to 15. Note that the value of the recursive call is not obtained by executing  $e$  on the arguments, but from the assumption list  $[\text{cons}(10, \text{nil}) \rightsquigarrow 10]$ . In fact, we also have  $e' \Downarrow_{\square}^{[\text{cons}(10, \text{nil}) \rightsquigarrow 10]} 15$  for  $e' \equiv \text{match } x_{\text{in}} \text{ with } \{\text{nil} \rightarrow 100, \text{cons}(h, t) \rightarrow h + f(t)\}$ . But we cannot derive  $(\text{fix } f(x_{\text{in}}).e') \models \text{Tr}'$  for any  $\text{Tr}' = (\text{cons}(5, \text{cons}(10, \text{nil})) \rightsquigarrow 15) \rightarrow \text{Tr}_s$ , while we can derive  $(\text{fix } f(x_{\text{in}}).e) \models \text{Tr}$  for  $\text{Tr}$  defined in Example 4.1.  $\square$

Thm. 4.4 establishes the equivalence of the trace-witnessed semantics and standard semantics of programs in tinyML with the help of Lem. 4.5 and Lem. 4.6.

**Theorem 4.4** (Soundness and Completeness of Trace-Witnessed Semantics). Given an I/O example  $i \rightsquigarrow o$ , we have that  $P(i) \Downarrow o$  if and only if  $P \models \text{Tr}$  for some trace  $\text{Tr} = (i \rightsquigarrow o) \rightarrow \text{Tr}_s$ .

**PROOF.** Soundness follows by rule induction on  $P \models \text{Tr}$ . The only case that applies is **TERec**. We have  $P = (\text{fix } f(x_{\text{in}}).e)$ ,  $\text{Tr} = (i \rightsquigarrow o) \rightarrow [\dots \text{Tr}_j \dots]$ ,  $[x_{\text{in}} := i] e \Downarrow_{\square}^{\text{Asm}} o$ , and  $(\text{fix } f(x_{\text{in}}).e) \models \text{Tr}_j$  for all  $j$ . Let's say  $\text{Root}(\text{Tr}_j) = i_j \rightsquigarrow o_j$  for all  $j$ , and  $\text{Asm} = [\dots i_j \rightsquigarrow o_j \dots]$ . By the inductive hypothesis,  $P(i_j) \Downarrow o_j$ . Then by Lem. 4.5 (see below),  $P(i) \Downarrow o$ . Lastly, completeness follows directly from Lem. 4.6 (see below).  $\square$

**Lemma 4.5.** If  $e \Downarrow_{\text{Asm}'}^{\text{Asm}}, v$  and  $P(i) \Downarrow o$  for all  $i \rightsquigarrow o \in \text{Asm} \setminus \text{Asm}'$  then  $e \Downarrow v$ .

**PROOF.** By rule induction on  $e \Downarrow_{\text{Asm}'}^{\text{Asm}}, v$ . Since the assumptions used in evaluation of sub-expression(s) have to be subsequences of  $\text{Asm} \setminus \text{Asm}'$ , applying the inductive hypothesis gives us the evaluation of sub-expression(s) in standard semantics. Applying rules from standard semantics, we reconstruct the evaluation of  $e$  in standard semantics.  $\square$



**Lemma 4.6.** If  $e \Downarrow v$  then  $\begin{cases} e_1 \Downarrow v_1 \text{ and } P \models (v_1 \rightsquigarrow v) \rightarrow \text{Trs} \text{ for some Trs} & \text{if } e = P(e_1) \\ e \Downarrow_{\text{Asm}_2}^{\text{Asm}_1} v \text{ for some Asm}_2 \text{ for arbitrary Asm}_1 \text{ and} & \\ P \models (i_j \rightsquigarrow o_j) \rightarrow \text{Trs}_j \text{ for some Trs}_j & \text{otherwise} \\ \text{for all } i_j \rightsquigarrow o_j \in \text{Asm}_1 \setminus \text{Asm}_2 & \end{cases}$

PROOF. By rule induction on  $e \Downarrow v$  and case analysis on whether the sub-expression(s) in  $e$  is a recursive subcall or not, we apply the inductive hypothesis accordingly and reconstruct the evaluation of  $e$  in trace-witnessed semantics along with witnessing Asm and/or Tr.  $\square$

## 4.2 Extended Version Spaces

Version space algebras (VSAs) are a way to compactly represent and manipulate solution sets to synthesis problems [Gulwani 2011; Lau et al. 2003; Polozov and Gulwani 2015]. In this section, we extend this representation with a new kind of version space operator called *choice* which allows us to more efficiently represent solution sets that involve conditional operators.

An *extended version space* is given by a term specified by the syntax below.

$$\text{PS} ::= \top \mid \{e_1, \dots, e_n\} \mid \cup (\text{PS}_1, \dots, \text{PS}_n) \mid \bowtie_{\text{op}} (\text{PS}_1, \dots, \text{PS}_n) \mid \Pi (\text{PS}_1, \dots, \text{PS}_n)$$

Here,  $\text{PS}_j$  are extended version spaces,  $e_j$  are tinyML expressions, and  $\text{op}$  is a function, a constructor, or a `match` from tinyML. The first four operators (top, list, union, and join) are standard version space algebra operators and have standard semantics. More formally, given any tinyML expression  $e$ , we have that: (a)  $e \in \top$ , (b)  $e \in \{e_1, \dots, e_n\}$  if  $\exists j. e = e_j$ , (c)  $e \in \cup (\text{PS}_1, \dots, \text{PS}_n)$  if  $\exists j. e \in \text{PS}_j$ , and (d)  $e \in \bowtie_{\text{op}} (\text{PS}_1, \dots, \text{PS}_n)$  if  $e$  is of the form  $\text{op}(e_1, \dots, e_n)$  and  $\forall j. e_j \in \text{PS}_j$ . Note that these semantics do not cover conditional terms, i.e., those involving the operator `cond`.

To understand how conditional terms are handled and how they interact with the choice operator  $\Pi(\cdot)$ , we first define the notion of *supporting valuations*. Formally, each extended version space  $\text{PS}$  is implicitly annotated with a set of valuations denoted by  $\text{Support}(\text{PS})$ . Intuitively, in our synthesis procedures from Sec. 5, if  $\text{PS}$  represents the solution set to a synthesis task, any program  $P$  that behaves equivalent to some program  $P^* \in \text{PS}$  on all valuations in  $\text{Support}(\text{PS})$  is also a solution. Our procedures ensure that a version space rooted at a non-choice operator has the same support as its sub-spaces, while the support of a version space rooted at a choice operator is the *disjoint* union of the support of its sub-spaces.

**Example 4.7.** Consider the version spaces  $\text{PS}_1 \equiv \{x_{\text{in}}, 2 \times x_{\text{in}}\}$  with support  $\text{Support}(\text{PS}_1) = \{x_{\text{in}} := 1, x_{\text{in}} := 2\}$  and  $\text{PS}_2 \equiv \{-x_{\text{in}}\}$  with support  $\text{Support}(\text{PS}_2) = \{x_{\text{in}} := -1, x_{\text{in}} := 0\}$ . These two version spaces are disjoint in the classical sense. However, the program  $P^* \equiv \text{cond } \{x_{\text{in}} > 0 \rightarrow x_{\text{in}}, x_{\text{in}} \leq 0 \rightarrow -x_{\text{in}}\}$  behaves like a program from both  $\text{PS}_1$  and  $\text{PS}_2$  on their respective supports.  $\square$

Thus, the choice operator  $\Pi(\cdot)$  in *extended version spaces* enables the inclusion of programs like  $P^*$  in  $\Pi(\text{PS}_1, \text{PS}_2)$ . Formally, given a tinyML expression  $e \equiv \text{cond } \{\dots, g_k \rightarrow e_k, \dots\}$  and an extended version space  $\text{PS} \equiv \Pi(\dots, \text{PS}_j, \dots)$ , we have  $e \in \text{PS}$  if

$$\forall j. \exists k. (e_k \in \text{PS}_j \wedge \forall \Delta \in \text{Support}(\text{PS}_j). [\Delta] g_k \Downarrow \text{true})$$

Combining the rule with the standard rules from the previous paragraph gives the full semantics of extended version spaces.

**Example 4.8.** Continuing from Example 4.7, the program  $P^*$  is contained in  $\Pi(\text{PS}_1, \text{PS}_2)$ , i.e.,  $P^* \in \Pi(\text{PS}_1, \text{PS}_2)$ . This is because for both  $\text{PS}_1$  and  $\text{PS}_2$ , there exists a unique matching branch in  $P^*$  whose guard is satisfied by any valuations in  $\text{Support}(\text{PS}_1)$  and  $\text{Support}(\text{PS}_2)$  respectively, i.e.,  $(x_{\text{in}} > 0)$  is satisfied by  $\{x_{\text{in}} := 1, x_{\text{in}} := 2\}$ , and  $(x_{\text{in}} \leq 0)$  is satisfied by  $\{x_{\text{in}} := -1, x_{\text{in}} := 0\}$ .  $\square$

We now define the extended version of the version space intersection operation. The intersection of two extended version spaces is defined by the following rules (modulo symmetry):

- ①  $\top \cap \text{PS} = \text{PS}$
- ②  $\{\text{PS}\} \cap \{\text{PS}\} = \{\text{PS}\}$
- ③  $\{\text{PS}_1, \dots, \text{PS}_n\} \cap \text{PS} = \{\text{PS}_1 \cap \text{PS}, \dots, \text{PS}_n \cap \text{PS}\}$
- ④  $\bigcup (\text{PS}_1, \dots, \text{PS}_n) \cap \text{PS} = \bigcup (\text{PS}_1 \cap \text{PS}, \dots, \text{PS}_n \cap \text{PS})$
- ⑤  $\bowtie_{\text{op}} (\text{PS}'_1, \dots, \text{PS}'_k) \cap \bowtie_{\text{op}} (\text{PS}''_1, \dots, \text{PS}''_k) = \bowtie_{\text{op}} (\text{PS}'_1 \cap \text{PS}''_1, \dots, \text{PS}'_k \cap \text{PS}''_k)$
- ⑥  $\text{PS}_1 \cap \text{PS}_2 = \Pi (\text{PS}_1, \text{PS}_2, \bigcup (\text{PS}_1, \text{PS}_2))$  if rules ①-⑤ do not apply

Rules ①-② are the base cases. Rules ③-⑤ push intersections through the list, union, and join operators whenever possible. This is akin to how intersection is defined on standard version spaces. However, when these rules fail, we have a special rule ⑦ that produces a choice node  $\Pi (\text{PS}_1, \text{PS}_2, \bigcup (\text{PS}_1, \text{PS}_2))$ . Here, if  $S_1$  and  $S_2$  are the supports of the inputs  $\text{PS}_1$  and  $\text{PS}_2$ , the 3 operands of the choice operator have disjoint supports  $S_1 \setminus S_2$ ,  $S_2 \setminus S_1$ , and  $S_2 \cap S_1$ , respectively. Thus, the intersection contains conditional programs (e.g.  $P^*$ ) that “behave” like  $\text{PS}_1$  on  $S_1 \setminus S_2$ , like  $\text{PS}_2$  on  $S_2 \setminus S_1$ , and like one of  $\text{PS}_1$  and  $\text{PS}_2$  on  $S_1 \cap S_2$ . We use this notion of *observational equivalence* to formalize the soundness and completeness of extended version space intersection in the following theorem. The notation  $P_1 \equiv_S P_2$  denotes that  $P_1$  always evaluates to the same output as  $P_2$  for the valuations in support  $S$ .

**Theorem 4.9** (Observational Soundness and Completeness of Extended Version Space Intersection). Given extended version spaces  $\text{PS}_1$  and  $\text{PS}_2$  with corresponding supports  $\text{Support}(\text{PS}_1)$  and  $\text{Support}(\text{PS}_2)$ , their intersection  $\text{PS}_1 \cap \text{PS}_2$  (as defined by the above rules) satisfies: (a)  $P \in \text{PS}_1 \cap \text{PS}_2 \implies (P \in \text{PS}_1 \wedge P \in \text{PS}_2) \vee (\forall j \in \{1, 2\}. \forall P_j \in \text{PS}_j. P_j \equiv_{\text{Support}(\text{PS}_j)} P)$  (b)  $P \in \text{PS}_1 \wedge P \in \text{PS}_2 \implies P \in \text{PS}_1 \cap \text{PS}_2$ , and (c)  $\text{Support}(\text{PS}_1 \cap \text{PS}_2) = \text{Support}(\text{PS}_1) \cup \text{Support}(\text{PS}_2)$ .

### 4.3 Indexed Extended Version Spaces

An *indexed extended version space (indexed VS)*  $\Psi = \{\iota_1 \mapsto \text{PS}_1, \dots, \iota_n \mapsto \text{PS}_n\}$  is a mapping from distinct indices  $\iota_j$  to extended version spaces  $\text{PS}_j$ . In general, the indices represent some semantic information related to the expressions  $e \in \text{PS}$ . More specifically, fixing the type of indices to be one of frames, recursion traces, or recursion trace sequences yields frame-indexed, trace-indexed or traces-indexed VSes, respectively. These are usually denoted using the symbols  $\Psi^{\text{fr}}$ ,  $\Psi^{\text{tr}}$ , and  $\Psi^{\text{trs}}$ , respectively.

We define notions of consistency for indexed version spaces that state that the programs in the version space “match” the semantics of the index. Formally,

- (a) a frame-indexed VS  $\Psi^{\text{fr}}$  is consistent if  $\forall ((\Delta \leadsto v) \rightarrow \text{Asm}) \mapsto \text{PS} \in \Psi^{\text{fr}}. \forall e \in \text{PS}. [\Delta]e \Downarrow_{\square}^{\text{Asm}} v$ ;
- (b) a trace-indexed VS  $\Psi^{\text{tr}}$  is consistent if  $\forall \text{Tr} \mapsto \text{PS} \in \Psi^{\text{tr}}. \forall e \in \text{PS}. (\text{fix } f(x).e) \models \text{Tr}$ ; and
- (c) a traces-indexed VS  $\Psi^{\text{trs}}$  is consistent if  $\forall \text{Trs} \mapsto \text{PS} \in \Psi^{\text{trs}}. \forall e \in \text{PS}, \text{Tr} \in \text{Trs}. (\text{fix } f(x).e) \models \text{Tr}$ .

Now, we define operations on trace-indexed VSes, traces-indexed VSes, and frame-indexed VSes in Fig. 5, which preserve the consistency of indexed extended version spaces before and after the operations. Note that one may easily generalize definitions in Fig. 5 to handle more than 2 arguments.

- **CONCAT operation.** The CONCAT operation takes  $n$  consistent trace-indexed VSes  $\Psi_j^{\text{tr}}$  whose indices are all distinct, and returns a trace-indexed VS containing all  $\text{Tr} \mapsto \text{PS}$  mappings from each  $\Psi_j^{\text{tr}}$ . The resulting trace-indexed VS is trivially consistent.
- **PROD operation.** The PROD operation takes a consistent traces-indexed VS  $\Psi^{\text{trs}}$ , and a consistent trace-indexed VS  $\Psi^{\text{tr}}$ . For each pair of their mappings, we compute VS intersection while

$$\begin{aligned}
\text{CONCAT}(\Psi_1^{\text{tr}}, \dots, \Psi_n^{\text{tr}}) &\stackrel{\text{def}}{=} \{ \text{Tr} \mapsto \text{PS} \mid \exists 1 \leq j \leq n, (\text{Tr} \mapsto \text{PS}) \in \Psi_j^{\text{tr}} \} \quad \text{where all Tr are distinct} \\
\text{PROD}(\Psi^{\text{trs}}, \Psi^{\text{tr}}) &\stackrel{\text{def}}{=} \{ \text{Trs} \sqcup [\text{Tr}] \mapsto \text{PS}_1 \cap \text{PS}_2 \mid (\text{Trs} \mapsto \text{PS}_1) \in \Psi^{\text{trs}}, (\text{Tr} \mapsto \text{PS}_2) \in \Psi^{\text{tr}} \} \\
\text{CONS}(\text{Fr} \mapsto \text{PS}, \Psi^{\text{trs}}) &\stackrel{\text{def}}{=} \{ ((i \rightsquigarrow o) \rightarrow \text{Trs}) \mapsto \text{PS} \cap \text{PS}' \mid (\text{Trs} \mapsto \text{PS}') \in \Psi^{\text{trs}} \} \\
&\quad \text{where Fr} = (i \rightsquigarrow o) \rightarrow \text{Asm} \text{ and } \forall (\text{Trs} \mapsto \text{PS}) \in \Psi^{\text{trs}}, \text{Roots}(\text{Trs}) = \text{Asm}
\end{aligned}$$

$$\begin{aligned}
\text{MERGE}(\Psi_1^{\text{fr}}, \Psi_2^{\text{fr}}) &\stackrel{\text{def}}{=} \{ \text{Fr} \mapsto \cup (\text{PS}_1, \text{PS}_2) \mid \text{Fr} \mapsto \text{PS}_1 \in \Psi_1^{\text{fr}}, \text{Fr} \mapsto \text{PS}_2 \in \Psi_2^{\text{fr}} \} \\
&\quad \cup \{ \text{Fr} \mapsto \text{PS}_1 \in \Psi_1^{\text{fr}} \mid \text{Fr} \mapsto \text{PS}_2 \notin \Psi_2^{\text{fr}} \} \\
&\quad \cup \{ \text{Fr} \mapsto \text{PS}_2 \in \Psi_2^{\text{fr}} \mid \text{Fr} \mapsto \text{PS}_1 \notin \Psi_1^{\text{fr}} \} \\
\text{JOIN}(\text{op}, (\Delta \rightsquigarrow v), \Psi_1^{\text{fr}}, \Psi_2^{\text{fr}}) &\stackrel{\text{def}}{=} \text{MERGE}(\{ ((\Delta \rightsquigarrow v) \rightarrow \text{Asm}_1 \sqcup \text{Asm}_2) \mapsto \bowtie_{\text{op}} (\text{PS}_1, \text{PS}_2) \mid \\
&\quad ((\Delta_j \rightsquigarrow v_j) \rightarrow \text{Asm}_j) \mapsto \text{PS}_j \in \Psi_j^{\text{fr}}, j \in \{1, 2\} \} \} \\
&\quad \text{where } \Delta_1 \rightsquigarrow v_1 \text{ and } \Delta_2 \rightsquigarrow v_2 \text{ are sufficient specification of } \Delta \rightsquigarrow v
\end{aligned}$$

Fig. 5. Operations that compose indexed extended version spaces.

concatenating traces. The consistency of the result follows from the observational soundness of VS intersection.

- *CONS operation.* The CONS operation takes a consistent frame-VS mapping  $\text{Fr} \mapsto \text{PS}$ , and a consistent traces-indexed VS  $\Psi^{\text{trs}}$ . Then, it constructs a trace-indexed VS by Consing Fr onto Trs, the indices of  $\Psi^{\text{trs}}$ , and taking the intersection of two VSes consistent with Fr and Trs. The consistency of the resulting  $\Psi^{\text{trs}}$  follows from Rule [TERec](#).
- *MERGE operation.* The MERGE operation takes two consistent frame-indexed VSes  $\Psi_1^{\text{fr}}$  and  $\Psi_2^{\text{fr}}$ , and merges them by taking the union of  $\text{PS}_1$  and  $\text{PS}_2$  consistent with the same Fr. The consistency of the resulting  $\Psi^{\text{fr}}$  follows from the semantics of  $\cup(\cdot)$ .
- *JOIN operation.* Like join nodes in VSes, the JOIN operation is parametrized by an n-arity operator op. It also takes a synthesis problem specified by  $\Delta \rightsquigarrow v$  along with n consistent frame-indexed VSes  $\Psi_j^{\text{fr}}$ , denoting the solutions to synthesis sub-problems  $\Delta_j \rightsquigarrow v_j$  for op's arguments. Then it reconstructs a frame-indexed VS for  $\Delta \rightsquigarrow v$  by concatenating assumptions in  $\Psi_j^{\text{fr}}$  while joining the corresponding VSes with  $\bowtie_{\text{op}}(\cdot)$ , whose consistency follows from Rules [AERec](#), [AECtor](#), and [AEMatch](#).

## 5 SYNTHESIS ALGORITHM

In this section, we present our instantiation of the top-level synthesis procedure from [Algo. 1](#). We first define our layering  $(\mathcal{L}_j)_{j \in \mathbb{N}}$  of the search space  $\mathcal{L}$ . We then present a trace-guided LEARN procedure to learn a traces-indexed VS from the user-provided I/O examples over a finite subspace  $\mathcal{L}_j$  of this layered search space. We defer to [Sec. 6](#) the specifics of our trace-based ranking function over traces-indexed VSes which is used to select the program returned by [Algo. 1](#). The discussion focuses on the body  $e$  of the program `fix  $f(x_{\text{in}}).e$`  to be synthesized.

$$\begin{array}{c}
\frac{\text{isAtomic}(e) \quad \text{ASTSize}(e) \leq \mu(\text{atomicSize})}{e \vdash \mu} \text{SATOMIC} \quad \frac{\mu(\text{ctorDepth}) > 0 \quad \mu' = \mu[\text{ctorDepth}--] \quad e' \vdash \mu'}{C(e') \vdash \mu} \text{SCTOR} \\
\\
\frac{\begin{array}{c} \mu(\text{matchDepth}) > 0 \quad \mu_{\text{base}} = \mu[\text{matchDepth}--] \quad \mu_{\text{rec}} = \mu_{\text{base}}[\text{ctorDepth}++] \\ \forall j. \text{arity}(C_j) = 0 \Rightarrow e_j \vdash \mu_{\text{base}} \quad \forall j. \text{arity}(C_j) > 0 \Rightarrow e_j \vdash \mu_{\text{rec}} \end{array}}{\text{match } x \text{ with } \{\dots C_j(x_j) \rightarrow e_j \dots\} \vdash \mu} \text{SMATCH} \\
\\
\frac{\mu(\text{condWidth}) \geq n \quad \forall j. \text{isAtomic}(g_j) \quad \forall j. \text{ASTSize}(g_j) \leq \mu(\text{atomicSize}) \quad \forall j. e_j \vdash \mu}{\text{cond } \{\dots g_j \rightarrow e_j \dots\} \vdash \mu} \text{SCOND}
\end{array}$$

Fig. 6. Satisfaction of a program  $P$  (its function body) with respect to a  $\mu$ -value.

### 5.1 Layered Search Space

In contrast to commonly used size-based search space layering, we layer our search space  $\mathcal{L}$  in a more *artful* way to assimilate domain-specific *structural attributes* or *shapes* of programs. Specifically, we provide (1) a tuple  $\mu$  of parameters of a program's AST, (2) a characterization of the finite subspace  $\mathcal{L}_\mu$  of  $\mathcal{L}$  defined by each  $\mu$ -value, and (3) an ordering over  $\mu$ -values that yields a sequence  $(\mathcal{L}_{\mu_j})_{j \in \mathbb{N}}$ . We describe the parametrization  $\mu$  in Fig. 6 by defining judgment  $P \vdash \mu$ , denoting that program  $P$  satisfies the tuple  $\mu$  of parameter values. Then the subspace  $\mathcal{L}_\mu$  is characterized by the  $\mu$ -value in the sense that  $P \in \mathcal{L}_\mu$  iff  $P \in \mathcal{L}$  and  $P \vdash \mu$ . The specific ordering strategy used by us is a low-level heuristic and is hence deferred to Sec. 6.

Now we motivate each parameter in  $\mu$ , and describe judgment  $P \vdash \mu$ . (1) *matchDepth* specifies the maximum depth of pattern matching (Rule **SMATCH**), and thus helps bias against consuming too much data from input at a time. (2) *ctorDepth* specifies the maximum depth of constructors (Rule **SCTOR**) without pattern matching, and is set to be 0 by default. As we want to favor the usage of constructor when data is consumed from input via pattern matching,<sup>3</sup> the third premise in Rule **SMATCH** increments *ctorDepth* for *match* branches with patterns of non-zero arity. (3) *atomicSize* specifies the maximum size of atomic expressions (**IsAtomic**), namely function applications, constructors, variables, and their compounds, and bounds the size of enumeration in effect. (4) *condWidth* specifies the maximum width of *cond* expressions (Rule **SCOND**), and thus limits the complexity of branching logic in the programs.

**Example 5.1.** Consider the following two semantically equivalent implementations of the Fibonacci number function: (1)  $e^* \equiv \text{cond } \{x_{\text{in}} = 0 \rightarrow 1, x_{\text{in}} = 1 \rightarrow 1, x_{\text{in}} > 1 \rightarrow f((x_{\text{in}} - 2)) + f((x_{\text{in}} - 1))\}$ ; (2)  $e^\dagger \equiv \text{match } x_{\text{in}} \text{ with } \{\text{zero} \rightarrow 1, \text{succ}(n) \rightarrow \text{match } n \text{ with } \{\text{zero} \rightarrow 1, \text{succ}(m) \rightarrow f(m) + f(n)\}\}$ . Subspace  $\mathcal{L}_\mu$  contains  $e^*$  when  $\mu(\text{atomicSize})$  is at least 9 and  $\mu(\text{condWidth})$  is at least 3, and contains  $e^\dagger$  when  $\mu(\text{atomicSize})$  is at least 5 and  $\mu(\text{matchDepth})$  is at least 2.

### 5.2 Trace-Guided Synthesis

We present a trace-guided synthesis procedure **LEARN** (Algo. 2) which iterates through the set  $\mathcal{E}$  of user-provided I/O examples to generate a traces-indexed VS  $\Psi^{\text{trs}}$  of programs in subspace  $\mathcal{L}_\mu$ , along with their recursion traces, that is consistent with  $\mathcal{E}$ .<sup>4</sup> This traces-indexed VS is constructed by iteratively computing the product (**PROD**, Line 9) over the trace-indexed VS  $\Psi^{\text{tr}}$  of programs consistent with each I/O example  $i \rightsquigarrow o$ .

<sup>3</sup>In contrast, “corecursive calls” may only take place under guarding constructors to ensure productivity [Coquand 1994].

<sup>4</sup>For brevity, we omit the synthesis context, which may include definition of background functions.

**Algorithm 2** Trace-guided synthesis**Input:** I/O examples  $\mathcal{E}$ **Input:** Search space  $\mathcal{L}_\mu$ **Output:** Traces-indexed VS  $\Psi^{\text{trs}}$ 


---

```

1: procedure LEARN( $\mathcal{E}, \mathcal{L}_\mu$ )
2:    $\Psi^{\text{trs}} \leftarrow \{[] \mapsto \top\}$ 
3:   for all  $i \leadsto o \in \mathcal{E}$  do
4:      $\Psi^{\text{tr}} \leftarrow \emptyset$ 
5:      $\Psi^{\text{fr}} \leftarrow \text{LEARNANGELIC}(i \leadsto o, \mathcal{L}_\mu)$ 
6:     for all  $(\text{Fr} \mapsto \text{PS}) \in \Psi^{\text{fr}}$  do ▷  $\text{Fr} = (i \leadsto o) \rightarrow \text{Asm}$  for some Asm
7:        $\Psi_{\text{temp}}^{\text{tr}} \leftarrow \text{CONS}(\text{Fr} \mapsto \text{PS}, \text{LEARN}(\text{Asm}, \mathcal{L}_\mu))$ 
8:        $\Psi^{\text{tr}} \leftarrow \text{CONCAT}(\Psi^{\text{tr}}, \Psi_{\text{temp}}^{\text{tr}})$ 
9:    $\Psi^{\text{trs}} \leftarrow \text{PROD}(\Psi^{\text{trs}}, \Psi^{\text{tr}})$ 
10:  return  $\Psi^{\text{trs}}$ 

```

---

The interesting part of LEARN then lies in Lines 4–8, where LEARN constructs  $\Psi^{\text{tr}}$  for I/O example  $i \leadsto o$ .<sup>5</sup> First, LEARN uses a sub-procedure LEARNANGELIC (Algo. 3) to collect angelic solutions, along with the corresponding execution frames, in a frame-indexed VS  $\Psi^{\text{fr}}$ . Then, for each frame-VS mapping  $\text{Fr} \mapsto \text{PS}$  in  $\Psi^{\text{fr}}$ , LEARN recursively calls itself with the assumptions Asm from Fr as the new synthesis task (Line 7), and uses the output traces-indexed VS of this recursive call and the frame-VS mapping to iteratively build up  $\Psi^{\text{tr}}$  (Line 7 and Line 8). Recall that the CONS operation on Line 7 not only computes an intersection between the expressions in PS and in the output of LEARN, but also splices together Fr and the traces Trs in the output of LEARN into a new trace. The CONCAT operation on Line 8 collects the trace-indexed VSes for each mapping into one trace-indexed VS for the current I/O example.

**Example 5.2.** The user provides an example  $\text{Ex} = \text{cons}(1, \text{cons}(2, \text{cons}(3, \text{nil}))) \leadsto 6$ , and intends to have the synthesizer find a program that produces the sum of elements in the input list. LEARN calls LEARNANGELIC on Ex, and has it return a frame-indexed VS  $\Psi^{\text{fr}} = \{\dots, (\text{Ex} \rightarrow \text{Asm}) \mapsto \text{PS}, \dots\}$ . Say  $\text{Asm} = [\text{cons}(2, \text{cons}(3, \text{nil})) \leadsto 5]$  and PS contains  $e^* \equiv \text{match } x_{\text{in}} \text{ with } \{\text{nil} \rightarrow 0, \text{cons}(\text{hd}, \text{tl}) \rightarrow \text{hd} + f(\text{tl})\}$ . Recursively calling LEARN on Asm may produce a traces-indexed VS  $\Psi^{\text{trs}} = \{\dots, [\text{Tr}] \mapsto \text{PS}', \dots\}$  where Tr is a linear trace  $(\text{cons}(2, \text{cons}(3, \text{nil})) \leadsto 5) \rightarrow [(\text{cons}(3, \text{nil}) \leadsto 3) \rightarrow [(\text{nil} \leadsto 0) \rightarrow []]]$  and  $\text{PS}'$  also contains  $e^*$ .

As a result, the traces-indexed VS produced by calling LEARN on the user-provided Ex is of the form  $\{\dots, [\text{Ex} \rightarrow [\text{Tr}]] \mapsto \text{PS} \cap \text{PS}', \dots\}$ . And the desired implementation  $e^*$ , of course, is contained in  $\text{PS} \cap \text{PS}'$ .  $\square$

### 5.3 Angelic Synthesis

Given a synthesis problem specified by an example  $\Delta \leadsto v$  and a bounded search space  $\mathcal{L}_\mu$ , the synthesis procedure LEARNANGELIC (Algo. 3) constructs a frame-indexed VS  $\Psi^{\text{fr}}$  of expressions in  $\mathcal{L}_\mu$ , along with the assumptions that make them satisfy  $\Delta \leadsto v$ . The procedure is based on top-down enumerative search and lightweight deductive reasoning. In particular, the procedure uses deduction for expressions with top-level operators that enable deductive reasoning (Lines 3–5) and uses enumeration for atomic expressions (Line 6). At each step, the MERGE operation is used to

<sup>5</sup>In our implementation, trace-indexed VSes are memoized.

**Algorithm 3** Angelic synthesis

<b>Input:</b> Valuation context $\Delta$	8: <b>procedure</b> LEARNOP( $\text{op}, \Delta \rightsquigarrow v, \mathcal{L}_\mu$ )
<b>Input:</b> Goal value $v$	9: <b>for</b> $j \leftarrow 1, \dots, \text{ARITY}(\text{op})$ <b>do</b>
<b>Input:</b> Search space $\mathcal{L}_\mu$	10: $(\Delta_j \rightsquigarrow v_j) \leftarrow \text{DEDUCE}_{\text{op},j}(\Delta \rightsquigarrow v)$
<b>Output:</b> Frame-indexed VS $\Psi^{\text{fr}}$	11: $\mathcal{L}_{\mu_j} \leftarrow \text{WALKDOWN}_{\text{op},j}(\mathcal{L}_\mu)$
	12: $\Psi_j^{\text{fr}} \leftarrow \text{LEARNANGELIC}(\Delta_j \rightsquigarrow v_j, \mathcal{L}_{\mu_j})$
	13: <b>return</b> JOIN( $\text{op}, \Delta \rightsquigarrow v, \dots, \Psi_j^{\text{fr}}, \dots$ )
1: <b>procedure</b> LEARNANGELIC( $\Delta \rightsquigarrow v, \mathcal{L}_\mu$ )	14: <b>procedure</b> LEARNENUM( $\Delta \rightsquigarrow v, \mathcal{L}_\mu$ )
2: $\Psi^{\text{fr}} \leftarrow \emptyset$	15: $\Psi^{\text{fr}} \leftarrow \emptyset$
3: <b>for all</b> $\text{op} \in \text{DEDUCIBLEOP}(\mathcal{L}_\mu)$ <b>do</b>	16: <b>for all</b> $e \in \text{ATOMICEXP}(\mathcal{L}_\mu)$ <b>do</b>
4: $\Psi_{\text{temp}}^{\text{fr}} \leftarrow \text{LEARNOP}(\text{op}, \Delta \rightsquigarrow v, \mathcal{L}_\mu)$	17: <b>for all</b> $\text{Asm} \in \text{INFERASM}([\Delta] e, v)$ <b>do</b>
5: $\Psi^{\text{fr}} \leftarrow \text{MERGE}(\Psi^{\text{fr}}, \Psi_{\text{temp}}^{\text{fr}})$	18: $\Psi^{\text{fr}} \leftarrow \text{MERGE}(\Psi^{\text{fr}},$
6: $\Psi_{\text{enum}}^{\text{fr}} \leftarrow \text{LEARNENUM}(\Delta \rightsquigarrow v, \mathcal{L}_\mu)$	19: $\{((\Delta \rightsquigarrow v) \rightarrow \text{Asm}) \mapsto \{e\}\})$
7: <b>return</b> MERGE( $\Psi^{\text{fr}}, \Psi_{\text{enum}}^{\text{fr}}$ )	20: <b>return</b> $\Psi^{\text{fr}}$

combine frame-indexed VSes by grouping the corresponding angelic solutions by their assumptions (Line 5, Line 6).

The procedure LEARNOP implements a general top-down algorithm that propagates the angelic synthesis problem, specified by example  $\Delta \rightsquigarrow v$  and search space  $\mathcal{L}_\mu$ , into the arguments of an operator  $\text{op}$  in tinyML. The JOIN operation then combines the angelic solutions for each argument by joining the respective version spaces and concatenating the respective assumptions.

The only two “deducible” operators in tinyML are  $C(\cdot)$  and **match**  $x$  **with**  $\{\dots C_j(x_j) \rightarrow \dots\}$ . And thus, we present DEDUCE procedures for them as follows. Notice that, in both cases, at most one example per argument is deduced. In particular, only one branch of **match** will get an example<sup>6</sup>.

$$\text{DEDUCE}_{\text{ctor},j}(\Delta \rightsquigarrow v) = \Delta \rightsquigarrow v_j \quad \text{for all } j, \text{ where } v = C(\dots, v_j, \dots)$$

$$\text{DEDUCE}_{\text{match}(x),j}(\Delta \rightsquigarrow v) = \begin{cases} (\Delta, x_j := v') \rightsquigarrow v & \text{if } [\Delta] x = C_j(v') \\ \emptyset & \text{otherwise} \end{cases}$$

The procedure LEARNENUM iterates through all atomic expressions  $e$  such that  $\text{ASTSize}(e) \leq \mu(\text{atomicSize})$ , and collects, in a frame-indexed VS, assumptions that make  $e$  angelically evaluate to the goal value  $v$ .

The procedure INFERASM is responsible for identifying all assumptions  $\text{Asm}$  such that  $e \Downarrow_{\text{[]}}^{\text{Asm}} v$  (Fig. 4). The key idea in the procedure is to identify the desired outputs for recursive subcalls in  $e$  that make  $e$  evaluate to  $v$  using an SMT solver (Z3 in our case). First, we substitute  $\Delta$  into  $e$  so that  $f$  is the only free variable in  $e$  if any; for variables denoting background functions, we substitute them with the Z3 encoding of their functional definitions. Then we substitute all recursive subcalls  $f(i_j)$  in  $e$  with a fresh symbol  $\alpha_{i_j}$  denoting their yet-to-be-determined outputs (making sure that the same symbol is used for recursive subcalls sharing the argument  $i_j$ ).<sup>7</sup> This leaves us a formula  $\epsilon$  such that  $e \equiv [\alpha_{i_j} := f(i_j)] \epsilon$ . Repeatedly, we query Z3 for satisfying assignments  $\alpha_{i_j} \leftrightarrow o_j$  of formula  $\phi \equiv (\epsilon = v)$ , and add conjuncts, that negate the seen assignments, to  $\phi$  until  $\phi$  becomes unsatisfiable. Each set of satisfying assignments  $\alpha_{i_j} \leftrightarrow o_j$  becomes a set of consistent assumptions about  $f$ , i.e.,  $f(i_j) = o_j$ .

<sup>6</sup>The other branches get no example, and we use a placeholder  $\top$  for them to indicate that any expression of type  $\tau$  works.

<sup>7</sup> $i_j$  has to be a value rather than a nested recursive subcalls since we only permit structural recursion.



**Example 5.3.** Recall in Example 5.2, LEARN passes an example  $\text{cons}(1, \text{cons}(2, \text{cons}(3, \text{nil}))) \rightsquigarrow 6$  to LEARNANGELIC. Calling LEARNANGELIC on this Ex may produce a frame-indexed VS  $\Psi^{\text{fr}} = \{(\text{Ex} \rightarrow \text{Asm}_1) \mapsto \text{PS}_1, (\text{Ex} \rightarrow \text{Asm}_2) \mapsto \text{PS}_2, \dots\}$  where:

- $\text{Asm}_1 = [\text{cons}(2, \text{cons}(3, \text{nil})) \rightsquigarrow 5]$  and  $\text{PS}_1$  contains `match  $x_{\text{in}}$  with  $\{\text{nil} \rightarrow 0, \text{cons}(\text{hd}, \text{tl}) \rightarrow \text{hd} + f(\text{tl})\}$` , and
- $\text{Asm}_2 = [\text{cons}(2, \text{cons}(3, \text{nil})) \rightsquigarrow 6]$  and  $\text{PS}_2$  contains `match  $x_{\text{in}}$  with  $\{\text{nil} \rightarrow 1, \text{cons}(\text{hd}, \text{tl}) \rightarrow \text{hd} \times f(\text{tl})\}$` .

In the first case, assuming  $f([2, 3]) = 5$ ,  $\text{head}(x_{\text{in}}) + f(\text{tail}(x_{\text{in}}))$  conforms to Ex and similarly, assuming  $f([2, 3]) = 5$ ,  $\text{head}(x_{\text{in}}) \times f(\text{tail}(x_{\text{in}}))$  conforms to Ex. Note that, though the two programs mentioned here happen to truly satisfy Ex,  $\text{PS}_1$  and  $\text{PS}_2$  may contain programs that only *angelically* satisfy Ex. `match  $x_{\text{in}}$  with  $\{\text{nil} \rightarrow 1, \text{cons}(\text{hd}, \text{tl}) \rightarrow \text{hd} + f(\text{tl})\}$`  is such an angelic candidate program that will be pruned away during VS intersection on Line 7 of Algo. 2.  $\square$

## 5.4 Soundness and Completeness

Since INFERASM returns all assumptions Asm such that  $e \Downarrow_{[]}^{\text{Asm}} v$ , and operations MERGE and JOIN are consistency preserving (as discussed in Sec. 4.3), frame-indexed VSes returned by LEARNANGELIC are consistent by construction. As a result, LEARNANGELIC is sound (Thm. 5.4).

**Theorem 5.4** (Soundness of LEARNANGELIC). For all  $\text{Fr} \mapsto \text{PS} \in \text{LEARNANGELIC}(\text{Ex}, \mathcal{L}_\mu)$  and for all  $e \in \text{PS}$ , we have  $e \in \mathcal{L}_\mu$  and  $[\Delta] e \Downarrow_{[]}^{\text{Asm}} v$  where  $\text{Fr} = (\Delta \rightsquigarrow v) \rightarrow \text{Asm}$ .

Since operations CONCAT, PROD, and CONS are consistency preserving (as discussed in Sec. 4.3), the soundness of LEARN (Thm. 5.5) follows directly from Thm. 5.4.

**Theorem 5.5** (Soundness of LEARN). For all  $\text{Trs} \mapsto \text{PS} \in \text{LEARN}(\mathcal{E}, \mathcal{L}_\mu)$  and for all  $e \in \text{PS}$ , we have that  $e \in \mathcal{L}_\mu$  and  $(\text{fix } f(x).e) \models \text{Trs}$ .

Now, we discuss the completeness of procedures LEARNANGELIC and LEARN. Careful readers may have noticed that neither procedure LEARNDEDUCE nor procedure LEARNENUM considers `cond` expressions. This is because neither deduction nor enumeration scales well for synthesis of `cond` expressions. The return value of LEARNANGELIC, *implicitly* and thus efficiently, carries `cond` expressions by only storing their branches, which may later be unified into choice nodes during VS intersection (Line 7 and Line 9)!

As a result, the completeness of LEARN and LEARNANGELIC holds modulo conditional expressions. That is, when an expression  $e$  uses a conditional expression when it is not necessary, it is not guaranteed that it is present in the return value of LEARNANGELIC or LEARN — only that another expression  $e'$  (possibly conditional free) that is observationally equivalent to  $e$  is present.

**Theorem 5.6** (Observational Completeness of LEARNANGELIC). If  $e \in \mathcal{L}_\mu$  and  $[\Delta] e \Downarrow_{[]}^{\text{Asm}} v$ , then there exists  $((\Delta \rightsquigarrow v) \rightarrow \text{Asm}) \mapsto \text{PS} \in \text{LEARNANGELIC}(\Delta \rightsquigarrow v, \mathcal{L}_\mu)$  and  $e' \in \text{PS}$  such that  $e' \equiv_{\text{Support}(\text{PS})} e$ .

**Theorem 5.7** (Observational Completeness of LEARN). If  $e \in \mathcal{L}_\mu$  and  $(\text{fix } f(x).e) \models \text{Trs}$ , then there exists  $\text{Trs} \mapsto \text{PS} \in \text{LEARN}(\text{Roots}(\text{Trs}), \mathcal{L}_\mu)$  and  $e' \in \text{PS}$  such that  $e' \equiv_{\text{Support}(\text{PS})} e$ .

## 6 IMPLEMENTATION

We implement a type-directed synthesizer, called SYRUP, in OCaml from scratch. Instead of parsing user-provided definitions, SYRUP incorporates inbuilt definitions of algebraic data types and background functions along with their Z3 encodings. In this section, we will expand on important implementation details and heuristics omitted in previous sections.



## 6.1 Parameterized Search Space Layering

In Sec. 5.1, we have seen that the search space parametrization enables us to concentrate on programs with a specific *shape* rather than a certain *size* in each synthesis iteration.

We now describe how SyRUP changes parameters across synthesis iterations, thereby prioritizing the search for programs of certain shapes over others. Parameters `atomicSize`, `matchDepth`, and `condWidth` are successively increased in lexicographic order; first `atomicSize` is incremented up to an upper bound; if a solution isn't found, then `matchDepth` is incremented, while `atomicSize` is reset to its initial value and incremented again up to an upper bound; and so on. Similar methods have been used in prior work [Lubin et al. 2020], and are referred to as “staging”.

Finally, we use initial values of the parameters that are sufficiently large to encompass interesting programs in the search space, but not excessively large to prevent the inclusion of improbable programs in the search space. Specifically, we initially let `atomicSize` = 3 to include function application, `matchDepth` = 1 since pattern matching is ubiquitous in functional programming, and `condWidth` = 1 to avoid overfitting.

## 6.2 Program Ranking Function

Program rank plays an important role in picking programs that generalize to unseen examples as we will see in Sec. 7.4. From the traces-indexed VS  $\Psi^{\text{trs}}$  returned by Algo. 2, SyRUP picks  $\text{Trs} \mapsto \text{PS}$  that maximizes the number of unique assumptions in  $\text{Trs}$ , favoring candidate programs that exploit more input information via recursion. Then, from  $\text{PS}$ , SyRUP picks a program with the smallest AST depth.

In the presence of choice nodes in the VSeS in  $\Psi^{\text{trs}}$ , SyRUP proceeds slightly differently. In particular, SyRUP first attempts to select a program using the above ranking function without any match expressions in its conditional guards. If this fails, SyRUP then attempts to select a program using the above ranking function with at most 1 match expression in its guards.

## 6.3 Tracking of Recursion Traces

Instead of tracking complete recursion traces, we choose to track an abstraction of recursion traces, specifically, the unique recursive sub-calls subsumed. Recall that we track supporting valuation contexts for our extended version spaces; the supports of top-level VS nodes exactly corresponds to the input of recursive sub-calls made by the programs contained. For example, a join node that contains the program in the first row of Fig. 1 has the following supports  $\{[m := 2, n := 1], [m := 1, n := 1], [m := 0, n := 1]\}$ . While this abstraction is sufficient for our benchmarks and evaluation, we defer the exploration of more precise abstractions to future work.

## 7 EVALUATION

In our experimental evaluation of SyRUP's performance, we consider the following research questions.

**RQ1** How does SyRUP perform on various synthesis tasks compared to existing techniques?

**RQ2** How sensitive is SyRUP's performance to the quality of the I/O examples?

**RQ3** How does SyRUP's layered search space and ranking function impact its overall performance?

All experiments are performed in a single thread on a department-wide shared Linux server equipped with two 2.90GHz Intel Xeon E5-2690 2.90GHz 8-core processors and 192GB of RAM, with a time limit of 120 seconds.

## 7.1 Experimental Setup

**Synthesis Tasks.** A synthesis task consists of a programming task  $T$  (e.g., write a program that adds two natural numbers) and a set of I/O examples e.g.  $\{(1, 1) \leadsto 2, (2, 1) \leadsto 3\}$ . We use a suite of 43 functional programming tasks that have been used to evaluate prior work [Lubin et al. 2020; Osera and Zdancewic 2015]. Of these, 32 tasks require recursion. For each programming task  $T$ , we spawn numerous synthesis tasks by generating I/O examples using the following two strategies. These strategies simulate the expertise of different users, allowing us to assess the performance of SyRUP in diverse scenarios.

- (1) **Expert Examples:** Each programming task in our suite was originally developed with a corresponding *trace-complete* set of I/O examples [Osera and Zdancewic 2015]. Similar to evaluations in prior work (SMYTH, BURST), we evaluate our synthesizer’s performance on subsets of these trace-complete example sets. Because such subsets are aware of the (recursive) semantics of the desired program to some degree, we regard them as *expert examples*. However, evaluations in prior work only use a *single, manually crafted subset* of the original trace-complete example set for each programming task. To effectively examine how well a synthesizer generalizes to unseen inputs given expert examples, we argue that it is important to evaluate across different subsets of these trace-complete example sets. In our experiments, we use subsets of different sizes, ranging from 1 to the total number of available expert examples, sampling up to 10 subsets for each size.
- (2) **Random Examples:** Expert examples, while not necessarily trace-complete, serve as a reasonable approximation of a user who is familiar with recursive programs and synthesis tools. To evaluate SyRUP’s performance on examples provided by non-experts, we generate random inputs and execute reference implementations on them, obtaining 10 example sets of sizes ranging from 1 to 6. To ensure that these example sets can be provided by a human user, we bound the number of constructors allowed in the input values, following a similar approach to Lubin et al. [2020]. For inputs with a function type, we randomly select from the functions used in expert examples.

**Success Rate.** To evaluate a synthesizer’s performance on a programming task  $T$  under various scenarios, we spawn synthesis tasks associated with  $T$  using I/O example sets with different attributes (identified by example generation strategies and number of examples). Then for a set  $\mathbb{E}_T$  of example sets  $\mathcal{E}$  with the same attribute, we compute the synthesizer’s *success rate*  $S_{\mathbb{E}_T}$  as follows:

$$S_{\mathbb{E}_T} = \frac{|\{\mathcal{E} \in \mathbb{E}_T \mid \text{given } \mathcal{E}, \text{ the synthesizer generates a correct program for } T\}|}{|\mathbb{E}_T|}$$

**Baselines.** We compare SyRUP against the two state-of-the-art tools, SMYTH [Lubin et al. 2020] and BURST [Miltner et al. 2022], for synthesis of recursive functional programs from examples. Neither SMYTH nor BURST requires trace-complete specifications. While SMYTH is based on top-down propagation and partial evaluation, BURST performs bottom-up search using a tree automata-based program representation.

We do not consider other synthesizers for recursive functional programs such as LEON [Kneuss et al. 2013],  $\lambda^2$  [Feser et al. 2015], and SYNQUID [Polikarpova et al. 2016] in our evaluation. Both LEON and SYNQUID are less effective in handling example-based specifications, and, in particular, ones that are not trace-complete [Lubin et al. 2020]. While  $\lambda^2$  can handle examples that are not trace-complete, the technique heavily relies on the use of language-level recursion schemes in the form of higher-order combinators such as `map` and `fold`.

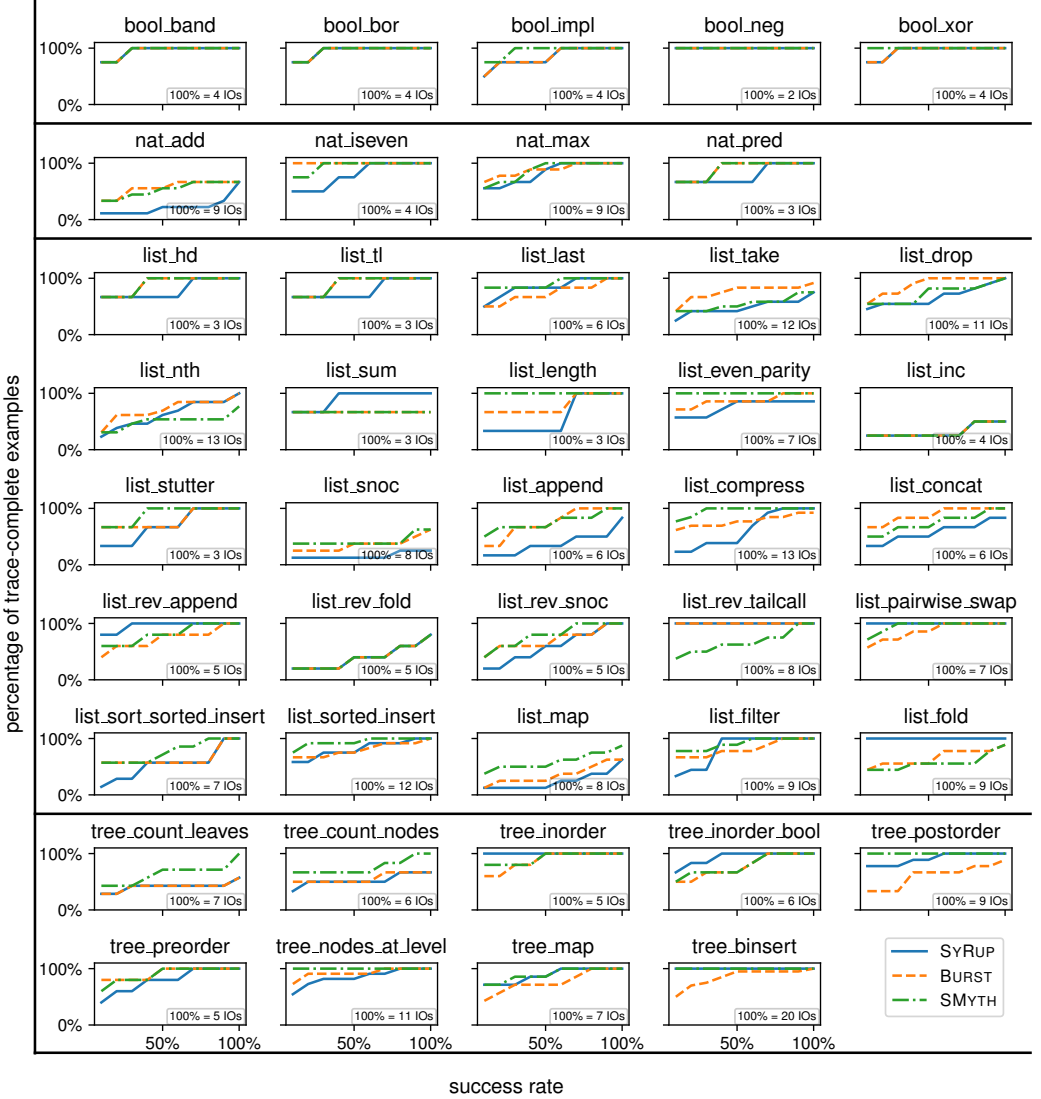


Fig. 7. Percentage of trace-complete examples needed to achieve specific success rates for each programming task. The x-axis label in each sub-figure indicates the size of the trace-complete example set for the corresponding programming task.

## 7.2 SyRUP’s Performance on Expert Examples

To evaluate SyRUP’s performance in the absence of trace-complete examples, we spawn synthesis tasks using sets of expert examples of different sizes for each programming task, and compute the respective success rates as discussed earlier. In Fig. 7, we present the *percentages of trace-complete examples* needed to achieve specific success rates (ranging from 10%, 20%, ..., 100%) for each programming task<sup>8</sup>.

<sup>8</sup>tree\_inorder\_bool corresponds to tree\_collect\_leaves in prior work. We correct the name to reflect its true semantics.

For all 43 programming tasks, SyRUP over-all needs less or similar number of examples to achieve the same success rate. For instance, to achieve 50% success rate, SyRUP needs less examples on 16 tasks than both BURST and SMYTH, and only needs more examples on 5 tasks than both of them (`list_sum`, `list_rev_append`, `list_filter`, `list_fold`, and `tree_inorder_bool`). This is better demonstrated in Fig. 8, which shows the percentage of trace-complete examples needed on average to achieve specific success rates. For instance, to achieve 50% success rate, SyRUP needs 70%, BURST needs 76%, and SMYTH needs 81% of trace-complete examples on average.

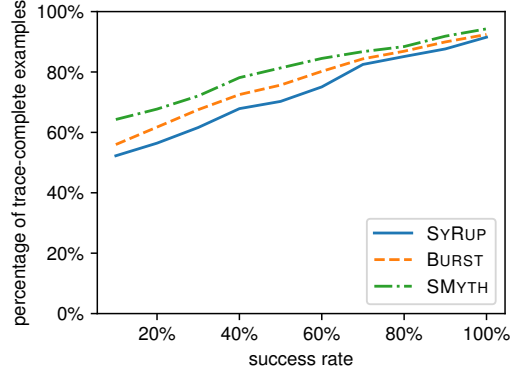


Fig. 8. Average percentage of trace-complete examples needed to achieve specific success rates.

### 7.3 SyRUP's Sensitivity to the Quality of Examples

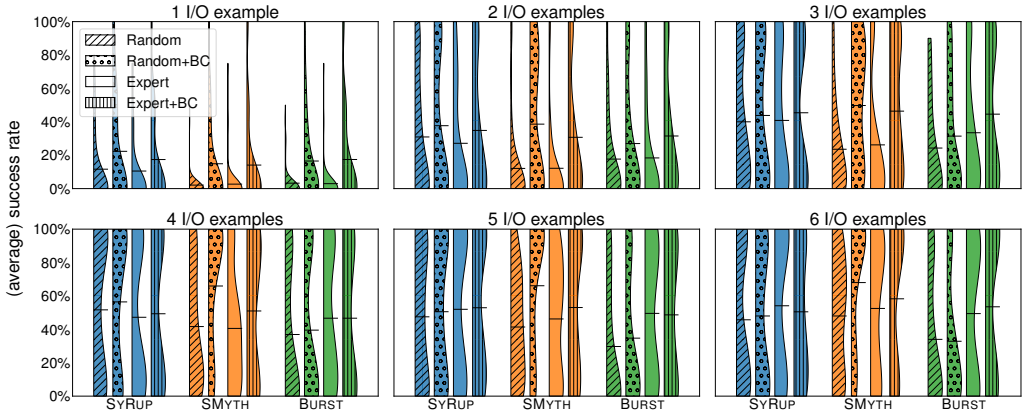


Fig. 9. Average success rate on example sets of size 1 to 6.

To assess the sensitivity of SyRUP to the quality of examples, we spawn synthesis tasks on example sets generated in different ways. Then, we present the average (the horizontal bars) and the variation (the violin plots) of success rate across all 43 programming tasks in Fig. 9. In addition to the two strategies for example generation mentioned already, we differentiate between example sets that *must* or *may* contain the *base case* for the recursive programming tasks (i.e. the example(s) with the smallest input). This yields 4 classes of example sets as shown in the figure (Random, Random + BC, Expert, Expert + BC). We repeat the experiment for example sets of sizes ranging from 1 to 6.

We observe that: (1) the average success rate of SyRUP across all programming tasks exhibits less variation compared to BURST and SMYTH when provided with the same number of I/O examples with different attributes; (2) the distribution of success rates of SyRUP across all programming tasks demonstrates a consistent trend regardless of the attribute of the given I/O examples (the shapes of the violin plots within each subplot are similar for SyRUP); and (3) with one exception

(discussed in Sec. 7.5), SyRUP overall has a higher or comparable success rate. Additionally, the specific attributes of example sets used highlights that: (1) SyRUP gains little improvement in success rate from the inclusion of base case examples, and thus does not fully rely on knowledge about the base case to synthesize the desired program; (2) providing random examples or expert examples has little impact on SyRUP's success rate, and thus SyRUP does not rely on users' mental model of the implementation to synthesize the desired program.

#### 7.4 Ablation Study

Recall Algo. 1. To evaluate the impact of SyRUP's design choice of search space layering and ranking function on its overall performance, we perform an ablation study using commonly-used variants for each design choice.

**Layered Search Space.** First, we implement a version of SyRUP, called SyRUP<sup>+</sup> that layers the search space based on AST height (similar to BURST) and rerun our experiment for expert examples. Unfortunately, SyRUP<sup>+</sup> times out very often for 33 out of 43 programming tasks.

**Trace-based Ranking Function.** We next implement a version of SyRUP, called SyRUP\*, that uses the same layered search space as SyRUP but uses a different ranking function.

In particular, SyRUP\* selects a program with the minimum AST height from the constructed traces-indexed VS. We rerun our experiment for expert examples. Since our ranking function greatly biases towards programs with longer recursion traces, the results shown in Fig. 10 focus on the 32 recursive tasks. As can be seen, both SyRUP and SyRUP\* need less examples to achieve the same success rate on average than BURST and SMYTH, with SyRUP slightly outperforming SyRUP\*. After inspecting the performance difference on individual tasks manually (omitted in the paper to save space), we observe that for 8 out of 10 tasks on which SyRUP outperforms (needs less examples to achieve the same success rate) SyRUP\*, SyRUP also outperforms BURST and SMYTH. Therefore, we conclude that, both trace-guided synthesis and trace-based ranking function contribute to SyRUP's performance on expert examples.

#### 7.5 Discussion

**Evaluation Summary.** Our evaluation demonstrates SyRUP's ability to effectively tackle the challenges of generalization, trace-completeness, and ambiguity-resolution across diverse programming tasks and I/O example sets. In particular, the evaluation shows that SyRUP requires fewer (trace-complete) examples than its counterparts to achieve the same success rates, *and* is significantly less sensitive to example quality. Thus, SyRUP's trace-guided approach utilizing recursion traces as an inductive bias can benefit a variety of synthesis tasks for recursive programs.

**Synthesis Time and Scalability.** The synthesis time of SyRUP is comparable to that of SMYTH and BURST and is typically within 5 seconds or less when given a few examples with exceptions on a few relatively difficult tasks, which is often accompanied by time-outs. For instance, given 10 sets of expert examples of size 6, SyRUP times out at least once on 6 tasks, and BURST times out at least once on 9 tasks, out of 43 programming tasks. On the other hand, SMYTH, while rarely timing out, might fail to return a solution due to the absence of a base case example.

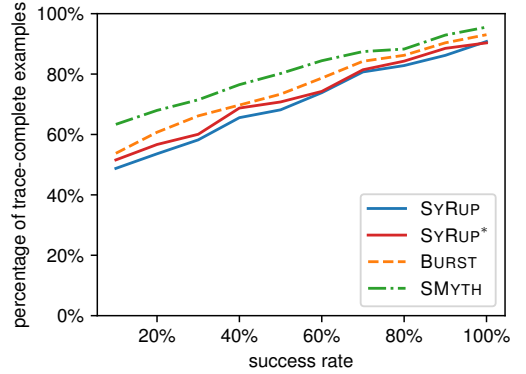


Fig. 10. Average percentage of trace-complete examples needed to achieve certain success rates across 32 recursive programming tasks.

It is interesting to note that both SYRUP and BURST start to time out more often when the number of examples increases or the quality of examples decreases (if we consider expert examples as high quality), while SMYTH's runtime is less affected. However, this is expected as both SYRUP and BURST predominantly rely on computationally expensive version space (tree automata for BURST) intersections. Consequently, we see in Sec. 7.3 that the increase in success rates of SYRUP and BURST slows down as the number of user-provided examples increases, especially when the examples are randomly generated. While we plan to explore optimizations for version space intersection, it is helpful to note that scalability is not an issue for the current benchmarks and the intended use-case of SYRUP only expects users to provide a small number of examples.

**Failure Analysis and Limitations.** SYRUP fails to find a correct implementation for `list_rev_tailcall` because we only permit recursive subcalls on values that are strictly smaller than the input argument—this is also the case in prior work with the same assumption [Lubin et al. 2020; Miltner et al. 2022; Osera and Zdancewic 2015]. Furthermore, due to efficiency considerations, we limit scrutinees in `match` expressions to be variables. As a result, SYRUP fails on `list_pairwise_swap`, whose correct implementation relies on *inside-out recursion* [Osera and Zdancewic 2015]. Allowing scrutinees to make recursive subcalls would require enumerating and checking the consistency of candidate programs as a whole since our algorithm relies on the concrete value of the scrutinee to deduce examples for branches in `match` expressions. Lastly, our implementation of the INFERASM decision procedure assumes that all background functions are first-order and monomorphic, as it is based on the Z3 Solver. While this limitation does not cause any issues with the current set of benchmarks, alternative or augmented implementations of INFERASM that overcome this limitation are worth exploring in future work.

## 8 RELATED WORK

**Synthesis of Recursive Programs.** Synthesis of recursive programs has witnessed great progress since Summers [1977], especially in the past decade. There are techniques that synthesize recursive programs from logical specifications [Itzhaky et al. 2021; Kneuss et al. 2013; Kuncak et al. 2010], refinement types [Polikarpova et al. 2016], or reference implementations [Farzan et al. 2022; Farzan and Nicolet 2021]. Here, we focus on synthesis of recursive programs from I/O examples. To enable efficient synthesis, initial attempts rely on either (1) comprehensive (trace-complete) examples, including IGOR II [Hofmann 2010], ESCHER [Albarghouthi et al. 2013], and MYTH [Osera and Zdancewic 2015], or (2) built-in recursion schemes, such as recursive higher-order functions used by  $\lambda^2$  [Feser et al. 2015] and BIG $\lambda$  [Smith and Albarghouthi 2016].

More recently, SMYTH [Lubin et al. 2020] and BURST [Miltner et al. 2022] have relaxed the requirement. Specifically, SMYTH devises a technique called *un-evaluation* and uses it in combination with best-effort forward evaluation [Omar et al. 2019] for top-down propagation of examples. BURST leverages angelic semantics to identify an over-approximating set of satisfying programs, picks the smallest program along with witnesses to the angelic satisfaction, and refines the program set with the witnesses, until a truly satisfying program is found. However, both SMYTH and BURST can miss programs that generalize better, and, in particular, have more suitable recursion traces. SYRUP extends the idea of angelic semantics to explore all possible traces in a bounded search space, and presents a novel VSA to represent and manipulate consistent trace-program pairs.

**Representation-based Synthesis.** The use of VSAs to learn and store sets of programs dates back to [Lau et al. 2003; Mitchell 1982]. Since FLASHFILL [Gulwani 2011] introduced a compact representation similar to e-graphs [Downey et al. 1980], multiple efforts have used similar representations in different problem settings [Barowy et al. 2015; Le and Gulwani 2014; Lee 2021; Polozov and Gulwani 2015; Rolim et al. 2017; Wang et al. 2016]. Another line of work uses finite tree automata (FTA) as a version space representation [Koppel et al. 2022; Miltner et al. 2022; Wang et al. 2017, 2018a,b]. Our



work presents a richer representation augmented with recursion traces and conditional semantics. Generalization to other semantic features, compact representation of indexed extended version spaces, and adaption to FTAs, are all exciting opportunities for future work.

**Condition Abduction/Unification.** Conditional expressions have been widely used to improve the scalability of synthesis using a divide-and-conquer approach [Albarghouthi et al. 2013; Alur et al. 2017; Gulwani 2011; Ji et al. 2021; Lee 2021; Shi et al. 2019; Wang et al. 2017]. Most of these approaches introduce conditionals in an *eager* way. Leveraging the conditional semantics of extended version space, SyRUP introduces conditionals in an *on-demand* (like ESCHER [Albarghouthi et al. 2013]), and *delayed* (like FRANGEL [Shi et al. 2019]) manner. A potential direction for future work is extension of the conditional semantics of extended version spaces using unification mechanisms that go beyond conditionals [Alur et al. 2015].

**Search Space Layering.** Existing synthesizers typically use a bounded search space and increase the bound (e.g., on program ASTs or intermediate values) when the current search space is exhausted. While one can construct a larger search space from the smaller ones, the search space still grows exponentially with respect to the bound [Barke et al. 2020]. Inspired by MYTH [Osera and Zdancewicz 2015] and SMYTH [Lubin et al. 2020], SyRUP parameterizes functional programs in a domain-specific and structural way. The parameters allow us to layer the search space more finely so that the search space not only grows more slowly, but also focuses on programs of different shapes. This means that our technique is likely to benefit from parallelizing the search with different threads focusing on programs of different shapes like Bornholt et al. [2016]; Padhi et al. [2019].

**Program Ranking.** Existing synthesis techniques, especially those based on VSAs, use a monotonic function defined over program ASTs to pick a “smallest” program. Smaller programs may not always generalize better as we have seen in Fig. 1. To the best of our knowledge, the only work that considers semantic features for ranking programs is Ellis and Gulwani [2017]. However, their approach relies on post-synthesis computation to rank candidate programs. In contrast, we collect and compose semantic features *during* synthesis. Our notion of indexed extended version space also opens up opportunities for exploration of other semantic features with different levels of abstraction.

## 9 DATA-AVAILABILITY STATEMENT

Our implementation along with the code required for reproduction, is available at Yuan et al. [2023].

## ACKNOWLEDGMENTS

We thank our anonymous reviewers and our shepherd, Hila Peleg, for their helpful feedback, and Patrick LaFontaine for his detailed comments on the initial draft of the paper. We thank Anders Miltner and Justin Lubin for answering our many questions about BURST and SMYTH. This research was partially supported by the National Science Foundation under Grant No. 1846327. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the National Science Foundation.

## REFERENCES

- Aws Albarghouthi, Sumit Gulwani, and Zachary Kincaid. 2013. Recursive Program Synthesis. In *Computer Aided Verification (Lecture Notes in Computer Science)*, Natasha Sharygina and Helmut Veith (Eds.). Springer, Berlin, Heidelberg, 934–950. [https://doi.org/10.1007/978-3-642-39799-8\\_67](https://doi.org/10.1007/978-3-642-39799-8_67)
- Rajeev Alur, Pavol Cerný, and Arjun Radhakrishna. 2015. Synthesis Through Unification. In *Computer Aided Verification (Lecture Notes in Computer Science)*, Daniel Kroening and Corina S. Păsăreanu (Eds.). Springer International Publishing, Cham, 163–179. [https://doi.org/10.1007/978-3-319-21668-3\\_10](https://doi.org/10.1007/978-3-319-21668-3_10)
- Rajeev Alur, Arjun Radhakrishna, and Abhishek Udupa. 2017. Scaling Enumerative Program Synthesis via Divide and Conquer. In *Tools and Algorithms for the Construction and Analysis of Systems*, Axel Legay and Tiziana Margaria (Eds.). Vol. 10205. Springer Berlin Heidelberg, Berlin, Heidelberg, 319–336. [https://doi.org/10.1007/978-3-662-54577-5\\_18](https://doi.org/10.1007/978-3-662-54577-5_18)



- Shraddha Barke, Hila Peleg, and Nadia Polikarpova. 2020. Just-in-Time Learning for Bottom-up Enumerative Synthesis. *Proceedings of the ACM on Programming Languages* 4, OOPSLA (Nov. 2020), 1–29. <https://doi.org/10.1145/3428295>
- Daniel W. Barowy, Sumit Gulwani, Ted Hart, and Benjamin Zorn. 2015. FlashRelate: Extracting Relational Data from Semi-Structured Spreadsheets Using Examples. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '15)*. Association for Computing Machinery, New York, NY, USA, 218–228. <https://doi.org/10.1145/2737924.2737952>
- Rastislav Bodik, Satish Chandra, Joel Galenson, Doug Kimelman, Nicholas Tung, Shaon Barman, and Casey Rodarmor. 2010. Programming with Angelic Nondeterminism. In *Proceedings of the 37th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '10)*. Association for Computing Machinery, New York, NY, USA, 339–352. <https://doi.org/10.1145/1706299.1706339>
- James Bornholt, Emina Torlak, Dan Grossman, and Luis Ceze. 2016. Optimizing Synthesis with Metasketches. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '16)*. Association for Computing Machinery, New York, NY, USA, 775–788. <https://doi.org/10.1145/2837614.2837666>
- Ravi Chugh, Brian Hempel, Mitchell Spradlin, and Jacob Albers. 2016. Programmatic and Direct Manipulation, Together at Last. *ACM SIGPLAN Notices* 51, 6 (June 2016), 341–354. <https://doi.org/10.1145/2980983.2908103>
- Thierry Coquand. 1994. Infinite Objects in Type Theory. In *Types for Proofs and Programs (Lecture Notes in Computer Science)*, Henk Barendregt and Tobias Nipkow (Eds.). Springer, Berlin, Heidelberg, 62–78. [https://doi.org/10.1007/3-540-58085-9\\_72](https://doi.org/10.1007/3-540-58085-9_72)
- Peter J. Downey, Ravi Sethi, and Robert Endre Tarjan. 1980. Variations on the Common Subexpression Problem. *J. ACM* 27, 4 (Oct. 1980), 758–771. <https://doi.org/10.1145/322217.322228>
- Kevin Ellis and Sumit Gulwani. 2017. Learning to Learn Programs from Examples: Going Beyond Program Structure. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence Organization, Melbourne, Australia, 1638–1645. <https://doi.org/10.24963/ijcai.2017/227>
- Azadeh Farzan, Danya Lette, and Victor Nicolet. 2022. Recursion Synthesis with Unrealizability Witnesses. In *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI 2022)*. Association for Computing Machinery, New York, NY, USA, 244–259. <https://doi.org/10.1145/3519939.3523726>
- Azadeh Farzan and Victor Nicolet. 2021. Counterexample-Guided Partial Bounding for Recursive Function Synthesis. In *Computer Aided Verification (Lecture Notes in Computer Science)*, Alexandra Silva and K. Rustan M. Leino (Eds.). Springer International Publishing, Cham, 832–855. [https://doi.org/10.1007/978-3-030-81685-8\\_39](https://doi.org/10.1007/978-3-030-81685-8_39)
- John K. Feser, Swarat Chaudhuri, and Isil Dillig. 2015. Synthesizing Data Structure Transformations from Input-Output Examples. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '15)*. Association for Computing Machinery, New York, NY, USA, 229–239. <https://doi.org/10.1145/2737924.2737977>
- Robert W. Floyd. 1967. Nondeterministic Algorithms. *J. ACM* 14, 4 (Oct. 1967), 636–644. <https://doi.org/10.1145/321420.321422>
- Sumit Gulwani. 2011. Automating String Processing in Spreadsheets Using Input-Output Examples. In *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '11)*. Association for Computing Machinery, New York, NY, USA, 317–330. <https://doi.org/10.1145/1926385.1926423>
- Martin Hofmann. 2010. IGOR2 - an Analytical Inductive Functional Programming System: Tool Demo. In *Proceedings of the 2010 ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation (PEPM '10)*. Association for Computing Machinery, New York, NY, USA, 29–32. <https://doi.org/10.1145/1706356.1706364>
- Shachar Itzhaky, Hila Peleg, Nadia Polikarpova, Reuben N. S. Rowe, and Ilya Sergey. 2021. Cyclic Program Synthesis. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation (PLDI 2021)*. Association for Computing Machinery, New York, NY, USA, 944–959. <https://doi.org/10.1145/3453483.3454087>
- Ruyi Ji, Jingtao Xia, Yingfei Xiong, and Zhenjiang Hu. 2021. Generalizable Synthesis through Unification. *Proceedings of the ACM on Programming Languages* 5, OOPSLA (Oct. 2021), 1–28. <https://doi.org/10.1145/3485544>
- Etienne Kneuss, Ivan Kuraj, Viktor Kuncak, and Philippe Suter. 2013. Synthesis modulo Recursive Functions. In *Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications*. ACM, Indianapolis Indiana USA, 407–426. <https://doi.org/10.1145/2509136.2509555>
- James Koppel, Zheng Guo, Edsko de Vries, Armando Solar-Lezama, and Nadia Polikarpova. 2022. Searching Entangled Program Spaces. *Proceedings of the ACM on Programming Languages* 6, ICFP (Aug. 2022), 91:23–91:51. <https://doi.org/10.1145/3547622>
- Viktor Kuncak, Mikael Mayer, Ruzica Piskac, and Philippe Suter. 2010. Complete Functional Synthesis. In *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '10)*. Association for Computing Machinery, New York, NY, USA, 316–329. <https://doi.org/10.1145/1806596.1806632>
- Tessa Lau, Steven A. Wolfman, Pedro Domingos, and Daniel S. Weld. 2003. Programming by Demonstration Using Version Space Algebra. *Machine Learning* 53, 1 (Oct. 2003), 111–156. <https://doi.org/10.1023/A:1025671410623>

- Vu Le and Sumit Gulwani. 2014. FlashExtract: A Framework for Data Extraction by Examples. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '14)*. Association for Computing Machinery, New York, NY, USA, 542–553. <https://doi.org/10.1145/2594291.2594333>
- Woosuk Lee. 2021. Combining the Top-down Propagation and Bottom-up Enumeration for Inductive Program Synthesis. *Proceedings of the ACM on Programming Languages* 5, POPL (Jan. 2021), 1–28. <https://doi.org/10.1145/3434335>
- Woosuk Lee and Hangeol Cho. 2023. Inductive Synthesis of Structurally Recursive Functional Programs from Non-recursive Expressions. *Proceedings of the ACM on Programming Languages* 7, POPL (Jan. 2023), 70:2048–70:2078. <https://doi.org/10.1145/3571263>
- Justin Lubin, Nick Collins, Cyrus Omar, and Ravi Chugh. 2020. Program Sketching with Live Bidirectional Evaluation. *Proceedings of the ACM on Programming Languages* 4, ICFP (Aug. 2020), 109:1–109:29. <https://doi.org/10.1145/3408991>
- Anders Miltner, Adrian Trejo Nuñez, Ana Brendel, Swarat Chaudhuri, and Isil Dillig. 2022. Bottom-up Synthesis of Recursive Functional Programs Using Angelic Execution. *Proceedings of the ACM on Programming Languages* 6, POPL (Jan. 2022), 21:1–21:29. <https://doi.org/10.1145/3498682>
- Tom M. Mitchell. 1982. Generalization as Search. *Artificial Intelligence* 18, 2 (March 1982), 203–226. [https://doi.org/10.1016/0004-3702\(82\)90040-6](https://doi.org/10.1016/0004-3702(82)90040-6)
- Cyrus Omar, Ian Voysey, Ravi Chugh, and Matthew A. Hammer. 2019. Live Functional Programming with Typed Holes. *Proceedings of the ACM on Programming Languages* 3, POPL (Jan. 2019), 1–32. <https://doi.org/10.1145/3290327>
- Peter-Michael Osera and Steve Zdancewicz. 2015. Type-and-Example-Directed Program Synthesis. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '15)*. Association for Computing Machinery, New York, NY, USA, 619–630. <https://doi.org/10.1145/2737924.2738007>
- Saswat Padhi, Todd Millstein, Aditya Nori, and Rahul Sharma. 2019. Overfitting in Synthesis: Theory and Practice. In *Computer Aided Verification (Lecture Notes in Computer Science)*, Isil Dillig and Serdar Tasiran (Eds.). Springer International Publishing, Cham, 315–334. [https://doi.org/10.1007/978-3-030-25540-4\\_17](https://doi.org/10.1007/978-3-030-25540-4_17)
- Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program Synthesis from Polymorphic Refinement Types. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '16)*. Association for Computing Machinery, New York, NY, USA, 522–538. <https://doi.org/10.1145/2908080.2908093>
- Oleksandr Polozov and Sumit Gulwani. 2015. FlashMeta: A Framework for Inductive Program Synthesis. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*. ACM, Pittsburgh PA USA, 107–126. <https://doi.org/10.1145/2814270.2814310>
- Reudismam Rolim, Gustavo Soares, Loris D'Antoni, Oleksandr Polozov, Sumit Gulwani, Rohit Gheyi, Ryo Suzuki, and Björn Hartmann. 2017. Learning Syntactic Program Transformations from Examples. In *2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE)*. 404–415. <https://doi.org/10.1109/ICSE.2017.44>
- Kensen Shi, Jacob Steinhardt, and Percy Liang. 2019. FrAngel: Component-Based Synthesis with Control Structures. *Proceedings of the ACM on Programming Languages* 3, POPL (Jan. 2019), 1–29. <https://doi.org/10.1145/3290386>
- Calvin Smith and Aws Albarghouthi. 2016. MapReduce Program Synthesis. In *Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI '16)*. Association for Computing Machinery, New York, NY, USA, 326–340. <https://doi.org/10.1145/2908080.2908102>
- Phillip D. Summers. 1977. A Methodology for LISP Program Construction from Examples. *J. ACM* 24, 1 (Jan. 1977), 161–175. <https://doi.org/10.1145/321992.322002>
- Xinyu Wang, Isil Dillig, and Rishabh Singh. 2017. Synthesis of Data Completion Scripts Using Finite Tree Automata. *Proceedings of the ACM on Programming Languages* 1, OOPSLA (Oct. 2017), 1–26. <https://doi.org/10.1145/3133886>
- Xinyu Wang, Isil Dillig, and Rishabh Singh. 2018a. Program Synthesis Using Abstraction Refinement. *Proceedings of the ACM on Programming Languages* 2, POPL (Jan. 2018), 1–30. <https://doi.org/10.1145/3158151>
- Xinyu Wang, Sumit Gulwani, and Rishabh Singh. 2016. FIDEX: Filtering Spreadsheet Data Using Examples. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA 2016)*. Association for Computing Machinery, New York, NY, USA, 195–213. <https://doi.org/10.1145/2983990.2984030>
- Yuepeng Wang, Xinyu Wang, and Isil Dillig. 2018b. Relational Program Synthesis. *Proceedings of the ACM on Programming Languages* 2, OOPSLA (Oct. 2018), 1–27. <https://doi.org/10.1145/3276525>
- Yongwei Yuan, Arjun Radhakrishna, and Roopsha Samanta. 2023. Artifact for "Trace-Guided Inductive Synthesis of Recursive Functional Programs". <https://doi.org/10.5281/zenodo.7812616>

Received 2022-11-10; accepted 2023-03-31