

Inductive synthesis of structurally recursive functional programs from non-recursive expressions

HANGYEOL CHO 

*Hanyang University,
Department of Computer Science & Engineering,
South Korea
(e-mail: pigon8@hanyang.ac.kr)*

WOOSUK LEE 

*Hanyang University,
Department of Computer Science & Engineering,
South Korea
(e-mail: woosuk@hanyang.ac.kr)*

Abstract

We present a novel approach to synthesizing recursive functional programs from input–output examples. Synthesizing a recursive function is challenging because recursive subexpressions should be constructed while the target function has not been fully defined yet. We address this challenge by using a new technique we call block-based pruning. A block refers to a recursion- and conditional-free expression (i.e., straight-line code) that yields an output from a particular input. We first synthesize as many blocks as possible for each input–output example, and then we explore the space of recursive programs, pruning candidates that are inconsistent with the blocks. Our method is based on an efficient version space learning, thereby effectively dealing with a possibly enormous number of blocks. In addition, we present a method that uses sampled input–output behaviors of library functions to enable a goal-directed search for a recursive program using the library. We have implemented our approach in a system called TRIO and evaluated it on synthesis tasks from prior work and on new tasks. Our experiments show that TRIO significantly outperforms prior work.

1 Introduction

Recent years have witnessed a surge of interest in recursive functional program synthesis (Albarghouthi et al., 2013; Kneuss et al., 2013; Feser et al., 2015; Osera and Zdancewic, 2015; Polikarpova et al., 2016; Lubin et al., 2020; Farzan and Nicolet, 2021; Miltner et al., 2022). In particular, because input–output examples are readily available, inductive synthesis of recursive functional programs has gained a lot of attention, witnessing significant strides. Inductive synthesis problems are typically expressed as a combination of algebraic data types, a library of *external* operators over the data types, and input–output examples that should be satisfied by the target function to be synthesized.

Despite recent advances, synthesizing recursive functional programs from input–output examples is still challenging, mainly due to the following two factors.

- *Recursive calls*: recursive data types often necessitate recursive calls, which are nontrivial to synthesize. That is because we should be able to reason about the target function yet to be defined during the search. As a workaround, previous approaches (Albarghouthi et al., 2013; Osera and Zdancewic, 2015) require the user to provide a *trace-complete specification* where the behaviors of recursive call expressions are part of the specification.¹ However, writing a trace-complete specification is quite unintuitive and difficult even for experts who are familiar with the synthesizers. To overcome this limitation, there have been previous methods, including specification strengthening (Miltner et al., 2022), partial evaluation, and constraint solving (Lubin et al., 2020). However, these approaches have their weaknesses, occasionally suffering from scalability issues even for small programs (see Section 6).
- *External operators*: to synthesize programs that utilize various operators over algebraic data types, synthesizers often require the user to provide a library of external operators. However, there is no general method for accelerating synthesis by exploiting the semantics of such external operators. Previous methods (e.g., Feser et al. (2015)) rely on predefined deductive rules only applicable to a fixed set of combinators (e.g., `map`, `fold`) or resort to naive enumeration. Therefore, the scalability issues worsen in the presence of a targeted library of external operators.

In this paper, we propose a novel approach to the inductive synthesis of recursive functional programs that addresses these challenges.

Our method for handling recursion, which we call *block-based pruning*, is to carry out synthesis in two phases: (1) synthesis of *blocks* satisfying the given examples followed by (2) synthesis of a recursive program. We define a *block* as a recursion- and conditional-free expression (i.e., straight-line program) that yields an output for a particular input, which is called *trace* in the prior work (Summers, 1986; Kitzelmann and Schmid, 2006). For each input–output example, we first synthesize as many blocks satisfying that example as possible. Based on an efficient version space learning, we effectively deal with possibly an enormous number of blocks.² And then, we explore the space of recursive programs top-down, generating incomplete candidate programs with holes (which we call *hypotheses*). For each hypothesis containing recursive calls, we transform it into blocks possibly with holes (which we call *open blocks*) by symbolic evaluation interleaved with concrete evaluation. If the open blocks cannot be the blocks synthesized in the earlier phase by filling the holes, the hypothesis is determined to be *inconsistent* with the blocks and is discarded.

Our method for handling external operators, which we call *library sampling*, is to sample input–output behaviors of library functions and use them for synthesis. This method enables a divide-and-conquer strategy called *top-down propagation* (or top-down

¹ For example, suppose the user tries to provide input–output examples $[] \mapsto 0$ and $[1, 2, 3] \mapsto 3$ to synthesize a function that returns the length of an integer list. To make the specification trace complete, the user should also provide two additional input–output examples: $[2, 3] \mapsto 2$ and $[3] \mapsto 1$.

² For example, a graph of 2470 nodes is used to represent over 7 million blocks (for the `list_rev_append` benchmark in Section 6).

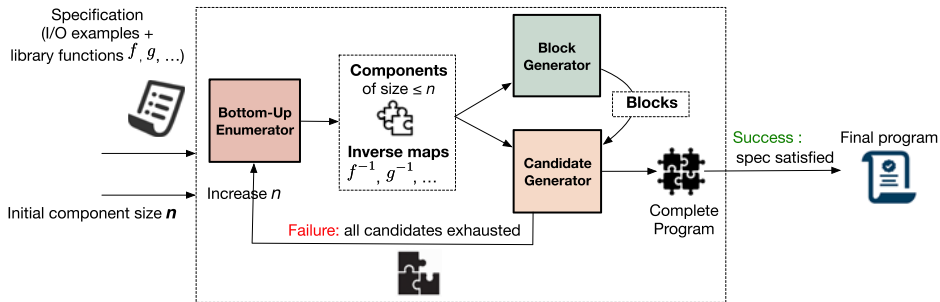


Fig. 1. High-level architecture of our synthesis algorithm.

deductive search) for synthesizing expressions that call *arbitrary* external operators. Top-down propagation hypothesizes an overall structure of the desired program satisfying given input–output examples and then performs deductive reasoning to recursively deduce new examples that should be satisfied by missing subexpressions. For example, suppose we want to synthesize a list manipulating program satisfying an input–output example $[1, 2] \mapsto [3, 4]$. After hypothesizing that the desired program is of form $\text{map}(f, x)$ where x denotes input and f is an unknown function subexpression to be synthesized, we can generate two new input–output examples for f as a synthesis subproblem: $1 \mapsto 3$ and $2 \mapsto 4$. This process is recursively repeated until all subproblems are solved. When hypothesizing the desired program is a function call expression involving external operators, we can use the input–output samples of the library functions to deduce new examples for missing subexpressions. This method is applicable even for black-box libraries.

Figure 1 presents the overall architecture of our synthesis algorithm, inspired by a recently proposed synthesis strategy (Lee, 2021). Our synthesis algorithm consists of three key modules, namely *Bottom-up enumerator*, *Block generator*, and *Candidate generator*:

- **Bottom-up enumerator:** Given synthesis specification comprising input–output examples and usable external operators, and a number n , the *Bottom-up enumerator* module generates two ingredients for the other modules: *components* and *inverse maps*. The components are expressions (of size $\leq n$) that can be used to construct blocks and recursive programs. The inverse maps are finite maps from outputs to inputs of the external operators and derived from input–output samples collected from concrete evaluation of the external operators.
- **Block generator:** Given the two ingredients from *Bottom-up enumerator*, the *Block generator* module generates blocks. For each input–output example, the module generates as many satisfying blocks as possible using the components. The block generation phase can be quickly done by top-down propagation. To control blocks in an enormous amount, we make use of version space representations to efficiently enumerate and store them.
- **Candidate generator:** Given the blocks generated by *Block generator*, the *Candidate generator* module searches for a solution also by top-down propagation. Starting with an empty program, it generates a sequence of hypotheses (i.e., partial programs with holes). During the search, any hypothesis inconsistent with the blocks is discarded early. *Candidate generator* keeps generating hypotheses until it

finds a solution, or all candidates have been explored in the search space. If a solution cannot be found using the current components, the whole process is repeated with the component size n increased by 1, thereby exploring the larger search space in the next iteration.

Our algorithm eventually finds a solution if it exists because *Bottom-up enumerator* will eventually enumerate a solution of finite size. Also, our method does not require the user to provide unintuitive trace-complete specifications.

We implemented our approach in a tool called TRIO. We evaluate TRIO on 65 benchmarks: 45 out of 65 are from prior work (Osera and Zdancewic, 2015), and the others are newly added. We use two types of specifications: (1) input–output examples and (2) reference implementations. Our evaluation results suggest that TRIO is more scalable than prior work on all types of specifications. In particular, our tool can synthesize 100% (65) of the functions from input/output examples and 91% (59) of the functions from reference implementations. We also compare TRIO against simpler variants that do not perform either block-based pruning or library sampling, and we empirically prove the efficacy of the two techniques.

Our contributions are as follows:

- A novel general method for synthesizing recursive programs from input–output examples: We propose a general algorithm for effectively synthesizing recursion and calls to external operators. We believe our method is potentially applicable to other synthesis contexts.
- Confirming the method’s effectiveness in an extensive experimental evaluation: We have conducted an extensive experimental evaluation on synthesis benchmarks from prior work and new benchmarks. Furthermore, we publicly release the implementation of our approach as a tool called TRIO (available at <https://github.com/pslhy/trio>).

Comparison with the previous version. This article is an extension of our previous work (Lee and Cho, 2023). Compared to the previous version, the current article presents a new method for ensuring termination of synthesized programs (Section 4.6). This method enables us to synthesize tail-recursive programs, which are not supported in the previous version. In addition, various optimizations for improving the scalability of the synthesis algorithm (Sections 4.7 and 5.1), which were not discussed in the previous version due to space constraints and a qualitative comparison to a recent work (Yuan et al., 2023) (Section 7) are added. Lastly, the evaluation section (Section 6) is extended with additional benchmarks.³

2 Overview

In this section, we give an overview of our method using the problem of synthesizing a recursive function `mul` for multiplying two natural numbers. The specification for the

³ We have added 5 new tail-recursive benchmarks to the 60 benchmarks used in the previous version to test if our termination check method works well. Also, the tool has been updated to fix some performance bugs and improve the performance.

problem comprises an inductive data type for natural numbers, an external operator `add` for adding two natural numbers that can be used to synthesize `mul`, and input–output examples embedded in a *hypothesis*. A hypothesis is a program that may have placeholders for missing expressions. We will call such a placeholder a *hole* (denoted \square), which is associated with input–output examples that should be satisfied by a subexpression in that position. In the following specification (in OCaml-like syntax),

```
type nat = Z | S of nat
rec add (x : nat, y : nat) : nat = match x with Z -> y | S _ -> S (add
    (S-1(x), y))
rec mul (x : nat, y : nat) : nat =  $\square_{in}$ 
```

\square_{in} is a hole associated with the set of input–output examples $\{(0, 1) \mapsto 0, (1, 2) \mapsto 2, (2, 1) \mapsto 2\}$. S^{-1} denotes a *destructor* which extracts the subcomponent of a constructor application of `S`. Such destructors obviate the need for introducing new variables bound by patterns in match expressions. The following program is a solution.

```
Psol = rec mul (x : nat, y : nat) : nat =
    match x with Z -> Z | S _ -> add (mul (S-1(x), y), y)
```

We will describe how the three modules of our system interact with each other to synthesize the desired program. For brevity, we will often use literals $0, 1, \dots$ as syntactic sugar for the corresponding naturals `Z, S(Z), \dots`.

Component generation and library sampling. *Bottom-up enumerator* first generates the following component pool `C` of expressions whose size is not greater than some user-provided upper bound.

$$C = \{x, y, S^{-1}(x), S^{-1}(S^{-1}(x)), S^{-1}(y), S^{-1}(S^{-1}(y)), Z, \text{add}, (S^{-1}(x), y), \dots\}$$

During bottom-up enumeration, it adopts the existing pruning technique based on *observational equivalence* to avoid maintaining multiple components of the same behaviors with respect to the input examples.⁴ This pruning technique drastically reduces the number of components by removing redundant expressions, which leads to overall performance gains. These components will be used to construct blocks and recursive programs in the following phases.

Next, for each function that a component in `C` may evaluate to, it constructs an inverse map of the function through a method we call library sampling. An inverse map of a function is a finite map from output values to input values of the function. Because `add` is the only function component, we construct the inverse map of `add`, which can be derived from input–output samples of `add`. Such samples can be obtained by evaluating `add` with input values that are not greater than the values in the examples. The reason behind this choice is that we aim to synthesize *structurally decreasing* recursive programs like previous approaches (Frankle et al., 2016; Osera and Zdancewic, 2015; Lubin et al., 2020; Miltner et al., 2022) where arguments of recursive calls are strictly decreasing, and we

⁴ Whenever a new program is enumerated, it is checked if it is “observationally equivalent” to any of the programs already constructed; i.e., those which produce the same outputs on inputs that were given as a specification. If so, the new program is discarded (e.g., $x + x$ is discarded if $2 \times x$ is already enumerated). This is done to avoid enumerating redundant programs.

observe inputs to the target function often flow to the external operators as arguments. Using numbers not greater than the greatest number (2) in the input–output examples (i.e., $(0, 0), (0, 1), \dots, (2, 2)$) as inputs, we evaluate the add function and obtain the inverse map $\text{add}^{-1} = \{0 \mapsto \{(0, 0)\}, 1 \mapsto \{(0, 1), (1, 0)\}, \dots, 4 \mapsto \{(2, 2)\}\}$.

Block generation. Next, for each input–output example, *Block generator* generates a set of blocks satisfying that example. Given an input–output example, it adds all the components in \mathbf{C} satisfying the example to the block set. Irrespective of whether or not any component is added, to find as various blocks as possible, it continues by hypothesizing about the structure of the other possible blocks and deduces new input–output examples that should be satisfied by missing holes. For each hole, it recursively searches for all the blocks that satisfy the hole. For example, consider the second input–output example $(1, 2) \mapsto 2$. *Block generator* first finds all components in \mathbf{C} that satisfy the example. Because the desired output 2 is the value of the second parameter y , y is added to the set of blocks. The search continues by hypothesizing about all possible structures of the other blocks. Suppose it attempts to find blocks of the form $S(\dots)$, generating a hypothesis $S(\square_1)$. The hole \square_1 is associated with $(1, 2) \mapsto 1$ where the output example is obtained by removing the constructor head S from the output example 2. Because the desired output 1 is the value of the first parameter x , $S(x)$ is added to the block set. For other possible blocks in place of \square_1 , it attempts to find blocks of the form $\text{add}(\dots)$. To generate hypotheses involving the external operator, it uses the inverse map of add . Since $\text{add}^{-1}(1) = \{(0, 1), (1, 0)\}$, it generates two hypotheses $S(\text{add}(\square_2, \square_3))$ and $S(\text{add}(\square_3, \square_2))$ where \square_2 and \square_3 are associated with $(1, 2) \mapsto 0$ and $(1, 2) \mapsto 1$, respectively. By finding components satisfying the holes, $S(\text{add}(Z, x))$, $S(\text{add}(x, Z))$, $S(\text{add}(S^{-1}(x), x))$, \dots are added to the block set. It further refines the holes \square_2 and \square_3 by recursively generating other hypotheses in a similar manner to find more blocks.

During the search, hypotheses containing recursive calls and match expressions are not taken into account because resulting blocks should be recursion- and conditional-free. Let us denote \mathbf{B}_i as the set of blocks for the i -th input–output example. We obtain the following blocks.

$$\begin{aligned} (0, 1) \mapsto 0: \quad \mathbf{B}_0 &= \{0, x, \text{add}(0, 0), \text{add}(0, x), \dots\} \\ (1, 2) \mapsto 2: \quad \mathbf{B}_1 &= \{2, y, S(\text{add}(0, x)), \text{add}(S^{-1}(x), y), \dots\} \\ (2, 1) \mapsto 2: \quad \mathbf{B}_2 &= \{x, S(y), \text{add}(x, 0), \text{add}(\text{add}(S^{-2}(x), y), y), \dots\} \end{aligned}$$

Because there are often infinitely many blocks satisfying each example, we limit the maximum number of steps of top-down propagation to ensure the termination of the block generation phase. For example, if we set the maximum number to be 1, in the above example, we would not recursively generate other hypotheses for the holes \square_2 and \square_3 as we already went through one step of top-down propagation. Even though we finitize the search space, there are often still many blocks. To efficiently enumerate and store them, we use a version space representation, which is a data structure that compactly represents a large set of programs (see Section 4.4).

Candidate generation. Equipped with the blocks generated by *Block generator*, *Candidate generator* searches for the desired recursive program by performing top-down

propagation, similar to what *Block generator* does but with a few differences: recursive calls and match expressions are generated, and all the input–output examples are considered at once, in contrast to *Block generator* that only considers one input–output example at a time. Suppose *Candidate generator* hypothesizes that the solution is a match expression with a guessed scrutinee x . Then, it generates the following hypothesis, distributing the input–output examples of \square_{in} into the two different branches.

$$P_0 = \text{rec mul } (x : \text{nat}, y : \text{nat}) : \text{nat} = \text{match } x \text{ with } Z \rightarrow \square_1 \mid S _ \rightarrow \square_2$$

where $\square_1 = \{(0, 1) \mapsto 0\}$ and $\square_2 = \{(1, 2) \mapsto 2, (2, 1) \mapsto 2\}$. Suppose *Candidate generator* fills the hole \square_1 with component x that satisfies the example and moves on to the hole \square_2 , trying to generate a hypothesis of the form $\text{add } (\cdot \cdot \cdot)$ in that position. Similarly to what *Block generator* did, *Candidate generator* uses the inverse map of add to generate new hypotheses. Because two output examples in \square_2 are 2 and there are three inputs of add that lead to the desired output 2 ($\text{add}^{-1}(2) = \{(0, 2), (1, 1), (2, 0)\}$), it deduces $9 (= 3^2)$ new hypotheses. Among them, let us consider the following hypothesis

$$P_1 = \text{rec mul } (x : \text{nat}, y : \text{nat}) : \text{nat} = \text{match } x \text{ with } Z \rightarrow x \mid S _ \rightarrow \text{add } \square_3$$

where $\square_3 = \{(1, 2) \mapsto (0, 2), (2, 1) \mapsto (1, 1)\}$. Observing the desired outputs are tuples of length 2, *Candidate generator* distributes the input–output examples into two new holes, generating the following hypothesis.

$$P_2 = \text{rec mul } (x : \text{nat}, y : \text{nat}) : \text{nat} = \text{match } x \text{ with } Z \rightarrow x \mid S _ \rightarrow \text{add } (\square_4, \square_5)$$

where $\square_4 = \{(1, 2) \mapsto 0, (2, 1) \mapsto 1\}$ and $\square_5 = \{(1, 2) \mapsto 2, (2, 1) \mapsto 1\}$. Suppose now it refines the hole \square_4 by generating a hypothesis of the form $\text{mul } (\cdot \cdot \cdot)$.

$$P_3 = \text{rec mul } (x : \text{nat}, y : \text{nat}) : \text{nat} = \text{match } x \text{ with } Z \rightarrow x \mid S _ \rightarrow \text{add } ((\text{mul } \square_6), \square_5)$$

Because mul is the target function yet to be defined, we cannot deduce examples for \square_6 . In such a case, we try enumerating all the components in \mathbf{C} that can be used as arguments. Recall that we only consider structurally decreasing arguments for recursive calls. For example, $\text{mul } (S^{-1}(x), y)$ is a valid recursive call as the first parameter decreases. By plugging it into the hole, we obtain

$$P_4 = \text{rec mul } (x : \text{nat}, y : \text{nat}) : \text{nat} = \text{match } x \text{ with } Z \rightarrow x \mid S _ \rightarrow \text{add } (\text{mul } (S^{-1}(x), y), \square_5)$$

Whenever a hypothesis containing recursive calls is generated, *Candidate generator* checks the feasibility of the hypothesis. It first performs symbolic evaluation interleaved with concrete evaluation with each input example on the hypothesis to obtain blocks. Our symbolic evaluation obeys the following rules.

- The body of the hypothesis is substituted into every position of a recursive call, and actual parameters are substituted for formal parameters.
- Every scrutinee in a match expression is concretely evaluated with a given input to take a branch.
- Calls to external operators and holes are left unchanged.

We call this process *unfolding*. As a result of unfolding, we obtain an open block, i.e., a block possibly with holes. Let us denote \rightarrow^* as one or more steps of the symbolic evaluation. The followings show how to derive a block B_j from the hypothesis for each j -th input–output example associated with \square_{in} .

```

(0, 1)  $\mapsto$  0:      mul (x, y)
 $\rightarrow$  match x with Z  $\rightarrow$  x | S  $\rightarrow$  add (mul ( $S^{-1}(x)$ , y),  $\square_5$ )
 $\rightarrow$  match 0 with Z  $\rightarrow$  x | S  $\rightarrow$  add (mul ( $S^{-1}(x)$ , y),  $\square_5$ )
 $\rightarrow$  x (=  $B_0$ )
(1, 2)  $\mapsto$  2:      mul (x, y)
 $\rightarrow$  match x with Z  $\rightarrow$  x | S  $\rightarrow$  add (mul ( $S^{-1}(x)$ , y),  $\square_5$ )
 $\rightarrow^*$  add (mul ( $S^{-1}(x)$ , y),  $\square_5$ )
 $\rightarrow$  add (match  $S^{-1}(x)$  with Z  $\rightarrow$   $S^{-1}(x)$  | S  $\rightarrow$  add (mul ( $S^{-2}(x)$ , y),  $\square_5$ ),  $\square_5$ )
 $\rightarrow$  add (match 0 with Z  $\rightarrow$   $S^{-1}(x)$  | S  $\rightarrow$  add (mul ( $S^{-2}(x)$ , y),  $\square_5$ ),  $\square_5$ )
 $\rightarrow$  add ( $S^{-1}(x)$ ,  $\square_5$ ) (=  $B_1$ )
(2, 1)  $\mapsto$  2:      mul (x, y)
 $\rightarrow$  match x with Z  $\rightarrow$  x | S  $\rightarrow$  add (mul ( $S^{-1}(x)$ , y),  $\square_5$ )
 $\rightarrow^*$  add (mul ( $S^{-1}(x)$ , y),  $\square_5$ )
 $\rightarrow$  add (match  $S^{-1}(x)$  with Z  $\rightarrow$   $S^{-1}(x)$  | S  $\rightarrow$  add (mul ( $S^{-2}(x)$ , y),  $\square_5$ ),  $\square_5$ )
 $\rightarrow$  add (match 1 with Z  $\rightarrow$   $S^{-1}(x)$  | S  $\rightarrow$  add (mul ( $S^{-2}(x)$ , y),  $\square_5$ ),  $\square_5$ )
 $\rightarrow^*$  add (add (mul ( $S^{-2}(x)$ , y),  $\square_5$ ),  $\square_5$ )
 $\rightarrow$  add (add (match  $S^{-2}(x)$  with Z  $\rightarrow$   $S^{-2}(x)$  | S  $\rightarrow$  add (mul ( $S^{-3}(x)$ , y),  $\square_5$ ),
     $\square_5$ ),  $\square_5$ )
 $\rightarrow^*$  add (add ( $S^{-2}(x)$ ,  $\square_5$ ),  $\square_5$ ) (=  $B_2$ )

```

where $S^{-2}(x)$ is a shorthand for $S^{-1}(S^{-1}(x))$. Then, for all j , it checks if each block B_j can be identical to another block in \mathbf{B}_j by properly substituting each hole. B_0 , which is x , is identical to x in \mathbf{B}_0 . B_1 , which is $\text{add}(S^{-1}(x), \square_5)$, can be identical to $\text{add}(S^{-1}(x), y)$ in \mathbf{B}_1 . Lastly, B_2 , which is $\text{add}(\text{add}(S^{-2}(x), \square_5), \square_5)$, can be identical to $\text{add}(\text{add}(S^{-2}(x), y), y)$ in \mathbf{B}_2 . This matching process can be efficiently done by traversing the version spaces of the blocks. The fact that the hypothesis can be *unfolded* into blocks satisfying the examples suggests that we may find a solution if we further refine the hypothesis. Thus, P_4 is determined to be feasible. Next, the hole \square_5 can be filled with y , which is a component satisfying the example over the hole, and we find the solution.

Feedback loop for guaranteeing search completeness. Although the block-based pruning presented may be *unsound* in some cases, the overall algorithm eventually finds a solution if it exists. A feasible hypothesis may be mistakenly rejected if *Block generator* misses some satisfying blocks because of its limited search in a finitized space. TRIO uses a feedback loop to avoid such unsound pruning. If a solution cannot be found using a current set of components, TRIO will add larger components into the component pool and repeat the entire process, so that *Block generator* can generate more blocks and hopefully avoid mistakenly rejecting correct hypotheses.

Also, when constructing inverse maps, despite restricting the domain of external functions to be the set of values each of which is not greater than the greatest value in the examples, we do not miss a solution involving external functions. This is because the bottom-up enumerator will eventually enumerate necessary function call expressions of finite size.

$e \in \text{Exp}$	(Expressions)	$\kappa \in \text{Constructors}$	(Constructors)
$v \in \text{Val}$	(Values)	$\square_u \in \text{Val} \xrightarrow{\text{fin}} \text{Val}$	(Input-Output Examples)
$ \begin{aligned} P &::= \text{rec } f(\mathbf{x}) = e \\ e &::= x \mid e_1 e_2 \mid \kappa(e_1, \dots, e_{a(\kappa)}) \mid \kappa^{-1}(e) \mid (e_1, \dots, e_n) \mid e.n \mid \text{rec } f(x) = e \\ &\quad \mid \text{match } e \text{ with } \overline{\kappa_j _ \rightarrow e_j^k} \mid \square_u \\ v &::= \kappa(v_1, \dots, v_k) \mid (v_1, \dots, v_m) \mid \text{rec } f(x) = e \\ \sigma &::= \{ \} \mid \{x \mapsto v\} \cup \sigma \end{aligned} $			

Fig. 2. Our ML-like language.

3 Problem definition

In this section, we define our problem of inductive synthesis of recursive functional programs. We first define an ML-like functional language in which we synthesize programs.

3.1 Language

We consider an idealized functional language similar to the core of ML. Our target language features algebraic data types and recursive functions with the syntax definition depicted in Figure 2. Programs P are recursive functions whose bodies are expressions e . Application is written $e_1 e_2$, κ ranges over data type constructors, $a(\kappa)$ denotes the arity of κ , κ^{-1} denotes a destructor which extracts all the subcomponents of a constructor application of κ as a tuple. An expression $e.n$ projects the n -th component of a tuple. We use ML-style pattern match expressions. We use $\overline{\kappa_j _ \rightarrow e_j^k}$ to denote $\kappa_1 _ \rightarrow e_1 \mid \dots \mid \kappa_k _ \rightarrow e_k$. A hole is written \square_u , where u is the hole name, which we tacitly assume is unique. Each hole is associated with input–output examples, a finite function from input values to output values. Values v are made up of constructor values for data types, tuples, and recursive functions. Recursive functions can be used as input examples when synthesizing higher-order functions but cannot be used as output examples. Environments σ map variables to values. For conciseness, we assume that all functions take a single argument, which does not harm the expressivity of the language since we can represent multiple inputs as a single tuple.

Example 1. The solution program P_{sol} for the overview example in Section 2 is represented as follows in our language.

```

 $P_{\text{sol}} = \text{rec mul } (x : \text{nat} * \text{nat}) : \text{nat} =$ 
  match  $x.1$  with  $Z \rightarrow Z \mid S \_ \rightarrow \text{add } (\text{mul } (S^{-1}(x.1), x.2), x.2)$ 

```

Though we use two parameters x and y in the overview example for better readability, we use a single tuple parameter x in the actual program since our language assumes a single argument for functions.

3.2 Notations

We will use some notations throughout the remaining sections. An *open hypothesis* (resp. open expression) is a program (resp. expression) that contains one or more holes. A *closed*

hypothesis is a program that does not contain any holes. We will use the fixed variables f and x to denote the target function and its formal parameter, respectively. We use $\sigma \vdash P \Rightarrow v$ to denote the standard multistep call-by-value operational semantics of a program P without holes under environment σ . Lastly, we denote the set of all subexpressions of an expression e as $\text{SubExprs}(e)$.

3.3 Problem definition

Given an environment σ that provides definitions of *external* functions and an initial open hypothesis $P_{in} = \text{rec } f(x) = \square_{in}$ where $\square_{in} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}$ represents input-output examples that should be satisfied by the function body, our goal is to find a closed hypothesis of form $P = \text{rec } f(x) = e$ that satisfies the input-output examples of \square_{in} . Formally, $\forall 1 \leq j \leq n. \sigma[f \mapsto \text{rec } f(x) = e, x \mapsto i_j] \vdash P x \Rightarrow o_j$ (denoted $P \models_{\sigma} \square_{in}$). We just use the user-provided external functions without inventing new ones.

4 Algorithm

This section formally describes our algorithm, inspired by the previous methods by Lee (2021) and Feser et al. (2015).

4.1 Overall algorithm

Figure 1 shows the high-level structure of our algorithm. The algorithm takes as input an environment σ that provides definitions of external functions, which we tacitly assume to be globally accessible throughout the algorithm, an initial open hypothesis with input-output examples, and initial component size n . Finally, it returns a program P that satisfies the input-output examples. With initially empty component set C , the main loop of our synthesis procedure (lines 2–28) is repeated until a solution is found. The loop starts by invoking the `COMPONENTGENERATION` procedure (line 3) which takes a current component pool C , the input-output examples \square_{in} , and the target component size n . The procedure generates new components by composing existing components in C . It applies the standard pruning technique based on observational equivalence with respect to the input examples. Expressions with recursive calls to the target function being synthesized can be included in the resulting component pool. Because we cannot evaluate such recursive components as the function is unknown yet, we cannot apply the observational equivalence reduction based on their outputs. Instead, we exploit *functional congruence*, i.e., the same input to the function always results in the same output. For example, we do not maintain both of $f\ 2$ and $f\ (1 + 1)$ in the component pool. Next, the `LibrarySampling` procedure (Section 4.2) is invoked to derive an inverse map for each function expression in the component pool (line 4). Next, the `BlockGen` procedure is invoked to obtain satisfying blocks for each input-output example (line 6). Each block must be recursion- and conditional-free because the target function is unknown yet and conditionals are not necessary when it comes to a single input-output example. Therefore, any components containing recursive calls and conditionals must not be used in blocks. We exclude such components from C (line 5).

Algorithm 1 The TRIO Algorithm

Require: (Global variable) environment σ that provides definitions of external functions
Require: Initial component size n
Require: A hypothesis $P_{in} = \text{rec } f(x) = \square_{in}$ where $\square_{in} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}$
Ensure: A solution program P

```

1:  $C := \emptyset$ 
2: repeat
3:    $C := \text{COMPONENTGENERATION}(C, \square_{in}, n)$ 
4:    $\mathcal{I} := \text{LibrarySampling}(C, \square_{in})$ 
5:    $C_{\text{simple}} := C - \text{RecursiveOrMatchExprs}(C)$ 
6:    $B := \text{BlockGen}(C_{\text{simple}}, \mathcal{I}, \square_{in})$ 
7:    $Q := \{P_{in}\}$ 
8:   while  $Q \neq \emptyset$  do
9:     Remove  $P$  from  $Q$  s.t.  $P$  has minimal cost
10:    if  $\neg \text{Terminate}(P)$  then
11:      continue
12:    end if
13:    if  $\text{Holes}(P) = \emptyset$  then
14:      if  $P \models_{\sigma} \square_{in}$  then return  $P$ 
15:      else continue
16:      end if
17:    end if
18:    Pick a hole  $\square_u$  in  $P$ 
19:    for  $e \in \text{Deduce}(C, \mathcal{I}, \square_u)$  do
20:       $P' := P[e/\square_u]$ 
21:      if  $\text{Holes}(P') = \emptyset$  then insert  $P'$  into  $Q$ 
22:      else if  $\text{BlockConsistent}(P', B, \square_{in})$  then
23:        insert  $P'$  into  $Q$ 
24:      end if
25:    end for
26:  end while
27:   $n := n + 1$ 
28: until false

```

Annotations for Algorithm 1:

- Lines 3–4: *Bottom-Up Enumerator* $\triangleright C : \mathcal{P}(\text{Exp})$ $\triangleright \mathcal{I} : \mathcal{P}(\text{Val} \times \text{Val} \times \text{Val})$
- Lines 5–6: *Block Generator* $\triangleright B : \text{Val} \times \text{Val} \rightarrow \mathcal{P}(\text{Exp})$
- Line 9: \triangleright Using a cost model in Section 5.3
- Lines 13–26: *Candidate Generator*

and provide the reduced component set to the BlockGen procedure (Section 4.4). With inverse maps \mathcal{I} and blocks B , the inner loop (lines 8–26) iteratively processes elements in the priority queue Q . The priority queue Q contains hypotheses (initially only P_{in}) and is sorted according to the cost (Section 5.3) of each hypothesis. In each iteration, we pick a minimum-cost hypothesis P from the queue (line 9). We first check if P is structurally recursive and guaranteed to terminate (line 10) (Section 4.6). If not, we continue to the next hypothesis. If P is closed (line 13) and correct with respect to the top-level input–output examples, P is returned as a solution (line 14). Otherwise, we continue investigating other hypotheses in the queue. If a chosen hypothesis P is open, we pick a hole \square_u in P (line 18). Then, the Deduce procedure (Section 4.3) returns possible replacements for the hole \square_u (line 19). A replacement e for the hole may be a closed expression satisfying the example of \square_u , or an open expression with new holes. For each replacement e , we obtain a new hypothesis P' by replacing the hole with e (line 20). If P' is closed, there are no unknowns left to be synthesized (line 21). Hence, we add P' into the queue, so that its correctness can be checked in the next iterations. If P' is open, we check its consistency with the

blocks **B** before adding it to the queue (line 22) by invoking the BlockConsistent procedure (Section 4.5). If the queue becomes empty before a solution is found, we increase the component size n by 1 (line 27) and restart the main loop.

Our algorithm is sound and complete in that it finds a program correct with respect to the given input–output examples if it exists in the search space.

Theorem 2. *Algorithm 1 finds a solution to a given synthesis problem if it exists.*

Proof Available in the appendix. ■

4.2 Getting inverse maps of external functions by library sampling

This section describes the LibrarySampling procedure that derives a set \mathcal{I} whose each element is a triple (g, v_o, v_i) where g is a function value and v_i and v_o are non-function values, meaning that $\sigma \vdash g \ v_i \Rightarrow v_o$. We will write (g, v_o, v_i) as $g^{-1}(v_o) = v_i$.

Given the component pool **C** and the top-level input–output examples $\square_{in} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}$ as input, we first compute a finite domain $D = \{v \in Val \mid \exists 1 \leq j \leq n. v \sqsubseteq i_j\}$ where \sqsubseteq denotes a well-founded ordering on values. In our implementation, we represent values as abstract syntax trees and use the subtree relation. Using the values in D as inputs, we compute a set of inverse maps as follows:

$$\mathcal{I} = \{g^{-1}(v_o) = v_i \mid g, v_o \in Val, v_i \in D, \exists e \in \mathbf{C}, 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e \Rightarrow g, \sigma \vdash g \ v_i \Rightarrow v_o\}$$

The use of the component pool **C** is for computing inverse maps of functions provided as input examples, which is useful for synthesizing higher-order functions.

Example 3. *Consider the following hypothesis.*

```
type nat = Z | S of nat
rec f (x : (nat -> nat) * nat) : nat =  $\square_{in}$ 
```

where $\square_{in} = \{(rec \ one \ (n) = S(Z), 1) \mapsto 1, (rec \ inc \ (n) = S(n), 0) \mapsto 1\}$. The solution is

```
rec f (x : (nat -> nat) * nat) : nat = x.1 x.2
```

Suppose the component pool $\mathbf{C} = \{x.1, x.2\}$. The domain D is $\{0, 1\}$. We derive $(rec \ one(n) = S(Z))^{-1}(1) = 0 \in \mathcal{I}$ because $\sigma[x \mapsto (rec \ one(n) = S(Z), 1)] \vdash x.1 \Rightarrow rec \ one(n) = S(Z)$ and $\sigma \vdash (rec \ one(n) = S(Z)) \ 0 \Rightarrow 1$. In a similar manner, we conclude $(rec \ inc(n) = S(n))^{-1}(1) = 0$, $(rec \ inc(n) = S(n))^{-1}(2) = 1 \in \mathcal{I}$. In conclusion,

$$\begin{aligned} \mathcal{I} = \{ & (rec \ one(n) = S(Z))^{-1}(1) = 0, (rec \ one(n) = S(Z))^{-1}(1) = 1, \\ & (rec \ inc(n) = S(n))^{-1}(1) = 0, (rec \ inc(n) = S(n))^{-1}(2) = 1\}. \end{aligned}$$

We consider D to be the domain for sampling for the following reason. We permit recursive calls on values that are strictly smaller than the input to ensure that our synthesized programs terminate, and inputs to a function often flow to other functions called inside of it. Therefore, it is likely that \mathcal{I} captures input–output behaviors of library functions that can be observed during the evaluation of the desired program with the user-provided input examples.

$$\begin{array}{c}
\frac{\{e \in \mathbf{C} \mid e \models_{\sigma} \square_u\} \subseteq \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)}{\text{D_COMPONENT}} \\
\\
\frac{\{e \in \mathbf{C} \mid e \text{ contains recursive calls to } f\} \subseteq \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)}{\text{D_REC}} \\
\\
\frac{\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto \kappa(v_{1j}, \dots, v_{kj})\} \quad \forall 1 \leq m \leq k, \square_{u_m} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto v_{mj}\}}{\kappa(\square_{u_1}, \dots, \square_{u_k}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)} \text{D_CTOR} \\
\\
\frac{\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\} \quad \kappa \in \text{Constructors} \quad \square_{u'} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto \kappa(o_j)\}}{\kappa^{-1}(\square_{u'}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)} \text{D_DTOR} \\
\\
\frac{\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\} \quad e \in \mathbf{C} \quad m \in \mathbb{N} \quad \forall 1 \leq j \leq n, \sigma[x \mapsto i_j] \vdash e.m \Rightarrow o_j}{e.m \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)} \text{D_PROJ} \\
\\
\frac{\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto (v_{1j}, \dots, v_{kj})\} \quad \forall 1 \leq m \leq k, \square_{u_m} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto v_{mj}\}}{(\square_{u_1}, \dots, \square_{u_k}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)} \text{D_TUPLE} \\
\\
\frac{\begin{array}{l} \square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\} \quad e \in \mathbf{C} \quad e \text{ does not contain recursive calls to } f \\ \forall 1 \leq m \leq k, \square_{u_m} = \bigcup_{j \in I_m} \{i_j \mapsto o_j\} \text{ where } I_m = \{j \mid 1 \leq j \leq n, \sigma[x \mapsto i_j] \vdash e \Rightarrow \kappa_m(_)\} \end{array}}{\text{match } e \text{ with } \overline{\kappa_i} _ \rightarrow \square_{u_i}^k \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)} \text{D_MATCH} \\
\\
\frac{\begin{array}{l} \square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\} \quad \square_{u_1} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto g_j\} \\ \square_{u_2} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto v_j\} \quad \forall 1 \leq j \leq n, g_j^{-1}(o_j) = v_j \in \mathcal{I} \end{array}}{\square_{u_1} \square_{u_2} \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)} \text{D_EXTCALL}
\end{array}$$

Fig. 3. Inference rules for Deduce.

4.3 The Deduce procedure

We now describe the Deduce procedure that returns a set of expressions that are either closed or open as possible replacements for a given hole \square_u .

$\text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)$ is the smallest set of expressions satisfying the constraints depicted in Figure 3. The D_COMPONENT rule says if we have components in \mathbf{C} that immediately satisfy \square_u , every such component can fill the hole ($e \models_{\sigma} \square_u$ denotes $\forall i \mapsto o \in \square_u. \sigma[x \mapsto i] \vdash e \Rightarrow o$). D_REC indicates that all recursive components in \mathbf{C} are considered potential replacements for the hole. This rule is under an optimistic assumption that any recursive expressions whose semantics is unknown yet may satisfy the hole, although in reality not all recursive expressions can do. Later, any hypothesis containing a recursive call that is determined to be infeasible will be discarded (will be detailed in Section 4.5). D_CTOR generates new examples for arguments of a constructor application. If the example values in the hole \square_u consist of constructor values with a shared constructor κ of arity k , then it creates k new examples constraints over the k arguments of the constructor value. D_DTOR creates a new example for the argument of a destructor value. The new example consists of constructor values with a shared constructor κ where κ can be any constructor. D_PROJ creates closed expressions as replacements for the hole. D_TUPLE is for creating new examples corresponding to arguments that must be synthesized for a tuple expression. The deductive reasoning process is similar to that of D_CTOR. D_MATCH first identifies components that can be used as scrutinees. Then, for each match expression whose scrutinee is such a component, it distributes the given examples in the hole to each branch. Lastly,

$D_EXTCALL$ uses the inverse maps of external functions. For example, for an input-output example $i \mapsto o$, if we can find a triple $g^{-1}(o) = v$ in \mathcal{I} , then we can deduce an example $i \mapsto g$ for the function part and another example $i \mapsto v$ for the argument part.

Comparison with prior work. Compared to the $IREFINE$ rules in MYTH (Osera and Zdancewic, 2015) and the deductive reasoning rules in λ^2 (Feser et al., 2015) that also propagate examples to holes in a top-down manner, the novelty of Deduce lies in the $D_EXTCALL$ rule. In MYTH, new function applications are generated by enumerating all possible combinations of functions and arguments. In contrast, by using the inverse maps in the $D_EXTCALL$ rule, we can expedite the search in a goal-directed manner. The deductive reasoning of λ^2 is only applicable to a fixed set of predefined functions such as filter and map. In contrast, Deduce is applicable to any external function.

Example 4. Consider the overview example in Section 2. Let us denote the target function as f . Suppose we have a hole $\square_u = \{(1, 2) \mapsto 2\}$ and a component pool $\mathbf{C} = \{x, f(S^{-1}(x.1), x.2), 2, add\}$. We can deduce the following constraints by applying the rules in Figure 3.

$$\begin{aligned}
& \text{By } D_COMPONENT, 2 \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\cdot : 2 \in \mathbf{C}, 2 \models_{\sigma} \square_u) \\
& \text{By } D_REC, f(S^{-1}(x.1), x.2) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) \\
& \text{By } D_CTOR, S(\square_{u_1}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\square_{u_1} = \{(1, 2) \mapsto 1\}) \\
& \text{By } D_DTOR, S^{-1}(\square_{u_2}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\square_{u_2} = \{(1, 2) \mapsto 3\}) \\
& \text{By } D_PROJ, x.2 \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\cdot : x \in \mathbf{C}, x.2 \models_{\sigma} \square_u) \\
& \text{By } D_MATCH, \text{match } 2 \text{ with } z \rightarrow \square_{u_3} \mid S_ \rightarrow \square_u \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\square_{u_3} = \emptyset) \\
& (\cdot : 2 \in \mathbf{C}, \sigma[x \mapsto (1, 2)] \vdash 2 \Rightarrow S_)
\end{aligned}$$

By $D_EXTCALL$,

$$\begin{aligned}
& \square_{u_4} \square_{u_5} \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\square_{u_4} = \{(1, 2) \mapsto \text{rec } add \dots\}, \square_{u_5} = \{(1, 2) \mapsto (2, 0)\}) \\
& \square_{u_4} \square_{u_6} \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\square_{u_6} = \{(1, 2) \mapsto (1, 1)\}) \\
& \square_{u_4} \square_{u_7} \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) & (\square_{u_7} = \{(1, 2) \mapsto (0, 2)\})
\end{aligned}$$

$$(\cdot : \{(\text{rec } add \dots)^{-1}(2) = (2, 0), (\text{rec } add \dots)^{-1}(2) = (1, 1), (\text{rec } add \dots)^{-1}(2) = (0, 2)\} \subseteq \mathcal{I})$$

Note that we cannot apply the D_TUPLE rule because \square_u does not contain any tuple. Also, when applying the D_MATCH rule, we cannot use the components x and $f(S^{-1}(x.1), x.2)$ as a scrutinee because neither of them evaluates to a constructor application (in particular, $f(S^{-1}(x.1), x.2)$ cannot evaluate to a concrete value as f is not defined yet). The following is the smallest solution satisfying the constraints over $\text{Deduce}(\mathbf{C}, \square_u)$.

$$\begin{aligned}
& \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u) = \{2, f(S^{-1}(x.1), x.2), S(\square_{u_1}), S^{-1}(\square_{u_2}), x.2, \text{match } 2 \text{ with } z \rightarrow \square_{u_3} \mid S_ \rightarrow \square_u, \\
& \quad \square_{u_4} \square_{u_5}, \square_{u_4} \square_{u_6}, \square_{u_4} \square_{u_7}\}
\end{aligned}$$

The Deduce procedure is *sound* in the following sense.

Definition 5 (Soundness of Deduction). Let \square_u be a set of input-output examples, and let \mathbf{C} and \mathcal{I} be a set of components and a set of inverse maps, respectively. If there exists an expression satisfying \square_u , for every open expression $e \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)$, for every hole $\square_{u'}$ in e , there exists an expression satisfying the hole $\square_{u'}$.

Intuitively, the deduction procedure is sound if there exists a solution to a synthesis task, then there also exists a solution to every synthesis subtask derived from the original synthesis task.

Theorem 6. *Without using the $D_EXTCALL$ rule, the Deduce procedure is sound.*

Proof Available in the appendix. ■

The following example shows that the Deduce procedure can be unsound when using the $D_EXTCALL$ rule.

Example 7. Recall Example 3. Let us denote the first and second input examples in \square_{in} as i_1 and i_2 respectively (i.e., $i_1 = (\text{rec one } (n) = S(Z), 1)$, $i_2 = (\text{rec inc } (n) = S(n), 0)$). We can deduce the following fact by applying the $D_EXTCALL$ rule because $(\text{rec one } \dots)^{-1}(1) = 0 \in \mathcal{I}$.

$$\square_{u_1} \square_{u_2} \in \text{Deduce}(\mathbf{C}, \square_{in}) \quad (\square_{u_1} = \{i_1 \mapsto \text{rec one } \dots, i_2 \mapsto \text{rec one } \dots\}, \square_{u_2} = \{i_1 \mapsto 0, i_2 \mapsto 0\})$$

We cannot synthesize an expression satisfying the hole \square_{u_1} . That is because we cannot synthesize an expression that evaluates to the *one* function under the environment where x is bound to i_2 (the only available function is *inc*). Recall that we do not synthesize any new auxiliary functions.

4.4 Constructing blocks from each input–output example

This section describes the BlockGen procedure for computing satisfying blocks for each input–output example in \square_{in} . The set of blocks is stored in a *version space* (Gulwani, 2011) which is a compact representation of expressions.

We begin with the definition of version spaces.

Definition 8 (Version Space). *A version space is either*

- *A union:* $\bigcup V$ where V is a set of version spaces
- *An expression*
- *An application:* written $(\tilde{e}_1 \tilde{e}_2)$ where \tilde{e}_i are version spaces
- *A tuple:* written $(\tilde{e}_1, \dots, \tilde{e}_k)$ where \tilde{e}_i are version spaces
- *A constructor:* written $\kappa(\tilde{e}_1, \dots, \tilde{e}_k)$ where \tilde{e}_i are version spaces and κ is a constructor
- *A destructor:* written $\kappa^{-1}(\tilde{e})$ where \tilde{e} is a version space and κ is a constructor
- *The empty set, \emptyset*

A version space can be understood as an E-graph where each node represents a set of expressions. Each leaf node represents a single expression, and they are composed into larger sets. The union operator \bigcup symbolizes a nondeterministic choice between multiple expressions, allowing version spaces to compactly represent huge sets of expressions.

$$\begin{array}{c}
\frac{\Box_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}}{\text{BlockGen}(\mathbf{C}, \mathcal{I}, \Box_u) = \{(i_j, o_j) \mapsto \text{Blocks}(\mathbf{C}, \mathcal{I}, i_j \mapsto o_j) \mid 1 \leq j \leq n\}} \text{B_GEN} \\
\\
\frac{\text{SimpleBlocks} = \{e \in \mathbf{C} \mid \sigma[x \mapsto i] \vdash e \Rightarrow o\}}{\text{Blocks}(\mathbf{C}, \mathcal{I}, i \mapsto o) = \bigcup \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, i \mapsto o) \cup \text{SimpleBlocks}} \text{B_GEN_PER_EX} \\
\\
\frac{\kappa(\Box_{u_1}, \dots, \Box_{u_k}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \{i \mapsto o\}) \quad \forall 1 \leq m \leq k. \widetilde{e}_m = \text{Blocks}(\mathbf{C}, \mathcal{I}, \Box_{u_m})}{\kappa(\widetilde{e}_1, \dots, \widetilde{e}_k) \subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, i \mapsto o)} \text{B_CTOR} \\
\\
\frac{\kappa^{-1}(\Box_u) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \{i \mapsto o\}) \quad \widetilde{e} = \text{Blocks}(\mathbf{C}, \mathcal{I}, \Box_u)}{\kappa^{-1}(\widetilde{e}) \subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, i \mapsto o)} \text{B_DTOR} \\
\\
\frac{e.n \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \{i \mapsto o\})}{e.n \in \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, i \mapsto o)} \text{B_PROJ} \\
\\
\frac{(\Box_{u_1}, \dots, \Box_{u_k}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \{i \mapsto o\}) \quad \forall 1 \leq m \leq k. \widetilde{e}_m = \text{Blocks}(\mathbf{C}, \mathcal{I}, \Box_{u_m})}{(\widetilde{e}_1, \dots, \widetilde{e}_k) \subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, i \mapsto o)} \text{B_TUPLE} \\
\\
\frac{\Box_{u_1} \Box_{u_2} \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \{i \mapsto o\}) \quad \forall 1 \leq m \leq 2. \widetilde{e}_m = \text{Blocks}(\mathbf{C}, \mathcal{I}, \Box_{u_m})}{(\widetilde{e}_1 \ \widetilde{e}_2) \subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, i \mapsto o)} \text{B_APP}
\end{array}$$

Fig. 4. Inference rules for BlockGen.

Example 9. The version space ($\text{add}(\bigcup\{x, Z\}, \bigcup\{x, Z\})$) encodes four different expressions: $\text{add}(x, x)$, $\text{add}(x, Z)$, $\text{add}(Z, x)$, and $\text{add}(Z, Z)$.

The set of expressions encoded by a version space is defined as follows:

Definition 10. The set represented by a version space is written $\llbracket \widetilde{e} \rrbracket$ and is defined recursively as

$$\begin{aligned}
\llbracket e \rrbracket &= e \quad (e \text{ is an expression}) & \llbracket \emptyset \rrbracket &= \emptyset & \llbracket \bigcup \mathcal{V} \rrbracket &= \{e \in \llbracket \widetilde{e} \rrbracket \mid \widetilde{e} \in \mathcal{V}\} \\
\llbracket (\widetilde{e}_1, \dots, \widetilde{e}_k) \rrbracket &= \{(e_1, \dots, e_k) \mid \forall 1 \leq j \leq k. e_j \in \llbracket \widetilde{e}_j \rrbracket\} \\
\llbracket \kappa(\widetilde{e}_1, \dots, \widetilde{e}_k) \rrbracket &= \{\kappa(e_1, \dots, e_k) \mid \forall 1 \leq j \leq k. e_j \in \llbracket \widetilde{e}_j \rrbracket\} \\
\llbracket \kappa^{-1}(\widetilde{e}) \rrbracket &= \{\kappa^{-1}(e) \mid e \in \llbracket \widetilde{e} \rrbracket\} & \llbracket (\widetilde{e}_1 \ \widetilde{e}_2) \rrbracket &= \{(e_1 \ e_2) \mid \forall 1 \leq j \leq 2. e_j \in \llbracket \widetilde{e}_j \rrbracket\}
\end{aligned}$$

With this in mind, we are ready to describe how to obtain a version space of blocks. Given a set \mathbf{C} of components and the top-level input–output examples \Box_{in} , the BlockGen procedure computes the smallest version spaces satisfying the constraints in Figure 4. The result maps each input–output example to a version space of satisfying blocks. In the B_GEN_PER_EX rule, $\text{Blocks}(\mathbf{C}, i \mapsto o)$ denotes the version space of blocks satisfying a single input–output example $i \mapsto o$. The other rules depict how to compute a version space for a single example. SimpleBlocks denotes a set of component expressions satisfying that example. CompoundBlocks denotes a set of version spaces each of which is not a single expression. We reuse the Deduce procedure to obtain CompoundBlocks, which can be derived by the other remaining rules B_CTOR, \dots , B_APP. Note that the given set of components \mathbf{C} only includes recursion- and conditional-free expressions (by line 5 in Algorithm 1), and there are no rules for deriving version spaces containing match

expressions and recursive calls, thereby ensuring that the resulting version space represents a set of blocks.

Example 11. Recall $\text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)$ in Example 4. We describe how to obtain a version space of blocks for the input–output example $\square_u = \{(1, 2) \mapsto 2\}$ using the rules in Figure 4. $\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_u)$ is computed as follows: first, by the $B_GEN_PER_EX$ rule,

$$\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_u) = \bigcup \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) \cup \text{SimpleBlocks}$$

where $\text{SimpleBlocks} = \{2\}$ because 2 is the only component satisfying the example. We can deduce the following constraints over $\text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u)$ as follows:

$$\begin{aligned} \text{By } B_CTOR, \mathcal{S}(\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_1})) &\subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) & (\because \mathcal{S}(\square_{u_1}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)) \\ \text{By } B_DTOR, \mathcal{S}^{-1}(\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_2})) &\subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) & (\because \mathcal{S}^{-1}(\square_{u_2}) \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)) \\ \text{By } B_PROJ, x.2 \in \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) & & (\because x.2 \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)) \\ \text{By } B_APP \\ (\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_4}) \text{ Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_5})) &\subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) \\ (\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_4}) \text{ Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_6})) &\subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) \\ (\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_4}) \text{ Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_7})) &\subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u) \end{aligned}$$

where $\square_{u_1}, \dots, \square_{u_7}$ are the ones defined in Example 4. We keep applying the rules to generate constraints over $\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_1}), \dots, \text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_7})$. For example, by the $B_GEN_PER_EX$ rule, $\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_1}) = \bigcup \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_{u_1}) \cup \text{SimpleBlocks}$ where $\text{SimpleBlocks} = \emptyset$ because no component in \mathbf{C} satisfies the example $\square_{u_1} = \{(1, 2) \mapsto 1\}$. Constraints over $\text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_{u_1})$ are generated by applying the rules in a similar manner. For instance, by the B_PROJ rule, a constraint $x.1 \in \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_{u_1})$ will be generated since $x.1 \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \square_{u_1})$. $\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_4})$ will include a version space of a single expression `add` since `add` is a component satisfying the example. $\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_6})$ (where $\square_{u_6} = \{(1, 2) \mapsto (1, 1)\}$) will include a version space of a tuple $(\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_1}), \text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_1}))$ by the B_TUPLE rule.

When generating constraints, a cycle that leads to blocks of infinite length may occur. For example, we may generate the following two constraints: $\mathcal{S}(\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_{u_1})) \subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_u)$ and $\mathcal{S}^{-1}(\text{Blocks}(\mathbf{C}, \mathcal{I}, \square_u)) \subseteq \text{CompoundBlocks}(\mathbf{C}, \mathcal{I}, \square_{u_1})$ that can be used to generate blocks of form $\mathcal{S}(\mathcal{S}^{-1}(\mathcal{S}^{-1}(\dots)))$. In our implementation, we bound the maximum height of version spaces to avoid generating blocks of infinite length.

We can derive the final version space of blocks for \square_u by finding the following smallest version space satisfying the above constraints.

$$\begin{aligned} \text{Blocks}(\mathbf{C}, \mathcal{I}, \square_u) = \\ \bigcup \{2, x.2, \mathcal{S}(x.1), (\text{add}(\bigcup \{x.1, \mathcal{S}^{-1}(x.2), \dots\}, \bigcup \{x.1, \mathcal{S}^{-1}(x.2), \dots\})), \dots\}. \end{aligned}$$

Comparison with prior work. The previous methods for version space construction for synthesis (Gulwani, 2011; Lee, 2021; Polozov and Gulwani, 2015) construct a version space of possible solutions directly. On the other hand, our version space construction is different in that it is used for pruning the search space. In addition, the previous methods rely on inverse semantics (also called witness functions) specialized for operators available in the target language for synthesis. Developers need to manually craft inverse semantics

for each operator. In contrast, our method does not require inverse semantics for arbitrary operators provided as library functions thanks to the use of inverse maps.

4.5 Pruning infeasible hypotheses using blocks

Finally, in this section, we describe the BlockConsistent procedure that prunes infeasible hypotheses using the blocks generated by the BlockGen procedure.

Deriving blocks from a hypothesis by unfolding. We first describe how to obtain blocks from an open hypothesis which we want to determine the feasibility of by a technique we call unfolding. Suppose a currently considered hypothesis is $P = \text{rec } f(x) = e_{\text{body}}$. For each input example i in the top-level input-output example \square_{in} , we perform symbolic evaluation (interleaved with concrete evaluation) over e_{body} and obtain a block, which possibly contains holes. We call such a block possibly with holes (resp. without holes) an *open block* (resp. *closed block*). We formalize our symbolic evaluation via the transition relation $e \rightarrow_{P,i} e'$ induced by the target hypothesis P and the input i . The relation says that the expression e takes a single step to the expression e' . The transition relation is formally defined by the rules in Figure 5. The rules U_CTOR, U_DTOR, U_TUPLE, and U_PROJ perform symbolic evaluation on the arguments of constructor, destructor, tuple, and projection expressions, respectively. U_APP_L and U_APP_R perform symbolic evaluation on the left and right hand sides of applications. The most notable part is the remaining two rules. U_MATCH for match expressions concretely evaluates the scrutinee e of a given match expression with input i . Then, a branch is chosen by the concrete value of the scrutinee. To obtain concrete values of scrutinees, we require scrutinees not to contain recursive calls to the target function, which is unknown yet. Therefore, any hypothesis containing a match expression that pattern matches on a recursive call to the target function (called *inside-out recursion* (Osera, 2015)) will get stuck and thus will be determined to be infeasible. This means *Candidate generator* will never generate programs with inside-out recursion. However, such programs can still be synthesized by our algorithm as *Bottom-up enumerator* will eventually enumerate all programs. U_REC is a special rule for recursive calls. Any recursive call to the target function f is replaced by the body of the function where every occurrence of the parameter x is replaced by the argument expression. Note that there are no transition rules for variables and holes. That is, every variable and hole in the hypothesis remains unchanged.

Given the top-level input-output examples $\square_{in} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}$, the set of open blocks derivable from hypothesis $P = \text{rec } f(x) = e_{\text{body}}$ (denoted B_P) is defined as follows:

$$B_P = \{(i_j, o_j) \mapsto e \mid e_{\text{body}} \rightarrow_{P,i_j}^* e, 1 \leq j \leq n\}.$$

where \rightarrow_{P,i_j}^* is the transitive closure of \rightarrow_{P,i_j} . That is, with each input example, we apply the transition rules till the end to obtain an open block.

Example 12. Recall the hypothesis P_4 in the overview example in Section 2.

$$P_4 = \text{rec } f(x) = \text{match } x.1 \text{ with } Z \rightarrow x.1 \mid S_ \rightarrow \text{add}(f(S^{-1}(x.1), x.2), \square_5)$$

$$\begin{array}{c}
\frac{e_i \rightarrow_{P,i} e'_i \quad 1 \leq i \leq k}{\kappa(e_1, \dots, e_i, \dots, e_k) \rightarrow_{P,i} \kappa(e_1, \dots, e'_i, \dots, e_k)} \text{U_CTOR} \quad \frac{e \rightarrow_{P,i} e'}{\kappa^{-1}(e) \rightarrow_{P,i} \kappa^{-1}(e')} \text{U_DTOR} \\
\\
\frac{e_i \rightarrow_{P,i} e'_i \quad 1 \leq i \leq k}{(e_1, \dots, e_i, \dots, e_k) \rightarrow_{P,i} (e_1, \dots, e'_i, \dots, e_k)} \text{U_TUPLE} \quad \frac{e \rightarrow_{P,i} e'}{e.n \rightarrow_{P,i} e'.n} \text{U_PROJ} \\
\\
\frac{e_1 \rightarrow_{P,i} e'_1}{e_1 e_2 \rightarrow_{P,i} e'_1 e_2} \text{U_APP_L} \quad \frac{e_2 \rightarrow_{P,i} e'_2}{e_1 e_2 \rightarrow_{P,i} e_1 e'_2} \text{U_APP_R} \\
\\
\frac{\sigma[x \mapsto i] \vdash e \Rightarrow \kappa_m _ \quad 1 \leq m \leq k}{\text{match } e \text{ with } \overline{\kappa_j} _ \rightarrow e_j^k \rightarrow_{P,i} e_m} \text{U_MATCH} \quad \frac{}{f e \rightarrow_{P,i} e_{\text{body}}[e/x]} \text{U_REC}
\end{array}$$

Fig. 5. Rules for unfolding (symbolic evaluation interleaved with concrete evaluation) for deriving open blocks from $P = \text{rec } f(x) = e_{\text{body}}$ with input i .

Using the rules in Figure 5, we can derive an open block from P_4 with input example $i = (1, 2)$ as follows:

```

match x.1 with Z -> x.1 | S _ -> add (f (S-1(x.1), x.2), □5)
→P4,i add (f (S-1(x.1), x.2), □5) (By U_MATCH)
→P4,i add (match S-1(x.1) with Z -> S-1(x.1)
| S _ -> add (f (S-2(x.1), x.2), □5), □5) (By U_REC, U_TUPLE, and U_APP_R)
→P4,i add (S-1(x.1), □5) (By U_MATCH, U_TUPLE, and U_APP_R)

```

Checking feasibility of a hypothesis. We check the feasibility of a hypothesis P by checking if it is *block consistent* with respect to the set \mathbf{B} of closed blocks from the BlockGen procedure, which is formally defined as follows:

Definition 13. Given the top-level input-output examples \square_{in} , a hypothesis $P = \text{rec } f(x) = e$ is block consistent with respect to a set \mathbf{B} of blocks if and only if

$$\forall i_j \mapsto o_j \in \square_{in}. B_P(i_j, o_j) \sim \mathbf{B}(i_j, o_j)$$

where \sim is a binary relation over Exp and version spaces, which is defined in Figure 6.

The relation \sim relates an open block to a set of closed blocks. Specifically, for an open block e and a version space of closed blocks \tilde{e} , $e \sim \tilde{e}$ holds if we can obtain an expression in \tilde{e} by properly filling each occurrence of the holes in e . Checking if $e \sim \tilde{e}$ resembles conventional syntactic matching between different expressions but with the following differences. Syntactic matching has as a goal to determine whether two expressions can be made equal by searching for a proper substitution from variables into expressions. For example, $\text{add}(x, Z)$ can be matched with $\text{add}(Z, Z)$ since we can substitute x with Z . On the other hand, in our method, not variables but only holes are targets for substitution. In addition, in contrast to syntactic matching that traverses two expressions, our method simultaneously traverses one expression and a version space to figure out if an open block can be matched with a closed block in the version space.

$$\begin{array}{c}
\frac{\forall 1 \leq m \leq k. e \sim \tilde{e}_m}{e \sim \bigcup \{\tilde{e}_1, \dots, \tilde{e}_k\}} \quad \frac{\forall 1 \leq m \leq k. e_m \sim \tilde{e}_m}{(e_1, \dots, e_k) \sim (\tilde{e}_1, \dots, \tilde{e}_k)} \quad \frac{\forall 1 \leq m \leq k. e_m \sim \tilde{e}_m}{\kappa(e_1, \dots, e_k) \sim \kappa(\tilde{e}_1, \dots, \tilde{e}_k)} \\
\frac{}{\square_u \sim \tilde{e} \neq \emptyset} \quad \frac{e = \kappa^{-1}(e') \quad e' \sim \tilde{e}}{e \sim \kappa^{-1}(\tilde{e})} \quad \frac{e_1 \sim \tilde{e}_1 \quad e_2 \sim \tilde{e}_2}{e_1 e_2 \sim (\tilde{e}_1 \tilde{e}_2)} \quad \frac{e \sim \tilde{e}}{e.n \sim \tilde{e}.n} \quad \frac{}{e \sim \tilde{e} \quad \tilde{e} = e}
\end{array}$$

Fig. 6. Matching rules for checking block consistency.

The first rule in Figure 6 says that any hole can be matched with any expression in \tilde{e} as long as \tilde{e} is not empty. The other rules recursively traverse the version space \tilde{e} of blocks and check if e can be matched with any expression in \tilde{e} .

Finally, the BlockConsistent procedure is defined as follows:

$$\text{BlockConsistent}(P, \mathbf{B}, \square_{in}) = \begin{cases} \text{true} & (\text{if } \forall i_j \mapsto o_j \in \square_{in}. B_P(i_j, o_j) \sim \mathbf{B}(i_j, o_j)) \\ \text{false} & (\text{otherwise}) \end{cases}$$

Example 14. The following derivation tree shows how the open block in Example 12 can be matched with the version space $(add(\bigcup \{Z, S^{-1}(x.1)\}, \bigcup \{x.2, S(x.1)\}))$ using the rules in Figure 6.

$$\begin{array}{c}
\frac{}{S^{-1}(x.1) \sim S^{-1}(x.1)} \\
\frac{}{S^{-1}(x.1) \sim \bigcup \{Z, S^{-1}(x.1)\}} \quad \frac{}{\square_5 \sim \bigcup \{x.2, S(x.1)\}} \\
\frac{}{add \sim add} \quad \frac{(S^{-1}(x.1), \square_5) \sim (\bigcup \{Z, S^{-1}(x.1)\}, \bigcup \{x.2, S(x.1)\})}{add(S^{-1}(x.1), \square_5) \sim (add(\bigcup \{Z, S^{-1}(x.1)\}, \bigcup \{x.2, S(x.1)\}))}
\end{array}$$

As already mentioned in Section 2, the block-based pruning presented may be *unsound*; a valid open hypothesis that can be a solution in the future may be pruned by the block-based pruning. Such a situation may occur if *Block generator* is not able to generate closed blocks for the valid hypothesis due to a lack of components. However, as the component pool grows, such unsoundness may be mitigated.

Please recall that even though the block-based pruning is unsound, our algorithm finds a solution if exists by resorting to *Bottom-up enumerator* that will eventually generate a solution of finite size.

Comparison with prior work. Our rules for unfolding are similar to the evaluation past holes in the Hazel system (Omar et al., 2019), which supports evaluation of incomplete programs for interactive editing of programs. However, we use the rules to prune the search space of recursive programs.

The novelty of our block-based pruning is discussed in Section 7.

4.6 Ensuring termination of synthesized programs

In this section, we describe how to ensure termination of synthesized programs. Through the Terminate procedure on line 10 of Algorithm 1, we check if a chosen candidate program $P = \text{rec } f(x) = e_{\text{body}}$ is guaranteed to terminate. If P is an open hypothesis containing

```

(a) 1: Function Terminate(rec f(x) = ebody)
    2: if RecursiveCalls(ebody) = ∅ then
    3:   return true
    4: else
    5:   return true if ∀ e ∈ RecursiveCalls(ebody). Struct(e, KeyArgs(ebody)) else false
    6: end if
    7: Function Struct(e, K)
    8: if e is of form (e1, ..., em) then
    9:   if K = ∅ then return false
   10:   e' := (ek)k ∈ K
   11:   x' := (x.i)i ∈ K
   12:   return true if e' ⊑ x' else false
   13: else
   14:   return true if e ⊑ x else false
   15: end if
   16: Function KeyArgs(e)
   17: if e is of form match e0 with κj → ejk then
   18:   K := {i ∈ ℕ | x.i ∈ SUBEXPRS(e0)}
   19:   return ⋃j=1k KeyArgs(ej) ∪ K
   20: else
   21:   return ⋃e' ∈ SUBEXPRS(e) KeyArgs(e')
   22: end if

```

The Terminate procedure and its helper functions

$$\begin{array}{c}
 \text{(b) } \boxed{e_1 \sqsubseteq e_2} \\
 \frac{\kappa \in \text{Constructors} \quad \text{ORD_DTOR} \quad \frac{e_1 \sqsubseteq e_2 \vee e_1 = e_2}{e_1.n \sqsubseteq e_2} \quad \text{ORD_PROJ}}{\kappa^{-1}(e) \sqsubseteq e} \\
 \frac{\exists 1 \leq i \leq m. e_i \neq e'_i \quad \forall 1 \leq i \leq m. e_i = e'_i \vee e_i \sqsubseteq e'_i}{(e_1, \dots, e_m) \sqsubseteq (e'_1, \dots, e'_m)} \quad \text{ORD_TUPLE}
 \end{array}$$

Inference rules for the partial order relation \sqsubseteq

Fig. 7. Termination checking procedure.

holes, the Terminate judgment answers if P can be completed to a terminating program. If P is a closed program, the Terminate judgment checks if P is terminating.

The pseudo-code and the inference rules in Figure 7 define our termination checking procedure. Figure 7(a) shows the Terminate procedure and its helper functions. The Terminate procedure takes a target function of the form $\text{rec } f(x) = e_{\text{body}}$ and returns **true** if the program is guaranteed to terminate. If there is no recursive call, the program is guaranteed to terminate (line 3). Otherwise, the program is guaranteed to terminate if all recursive calls are structurally decreasing. This is checked by the Struct function (line 5). For every recursive call $f \ e$ in the body of the target function (denoted $\text{RecursiveCalls}(e_{\text{body}})$), we check if the recursive call is valid. The judgment $\text{Struct}(f \ e, K)$ states that the argument expression e of the recursive call is deemed structurally decreasing where K is a set of indices of the arguments that may have to be structurally decreasing. Let us call such arguments *key arguments*. To see the role of K in ensuring termination of recursive calls, consider the following example.

Example 15. *The following function is a solution to the problem of synthesizing a function that reverses a given list in a tail-recursive manner.*

```

type nat = Z | S of nat
type list = Nil | Cons of nat * list

rec f (x : list * list) : list =
  match x.1 with
  | Nil -> x.2
  | Cons _ -> f (Cons-1(x.1).2, Cons(Cons-1(x.1).1, x.2))

```

The first component of the input tuple is a list to be reversed, and the second component is an accumulator. The function pattern matches on the first component of the input tuple. Therefore, the first component should be structurally decreasing in each recursive call to ensure termination. On the other hand, the second component does not affect the termination of the function. To see if the function is terminating, we keep track of the indices of the key arguments and check if the arguments are structurally decreasing. In this case, the key argument is the first component of the input tuple and the set K is $\{1\}$.

The function `KeyArgs` takes an expression e as input and returns a set of indices of the key arguments. If e is a match expression (line 17), the key arguments are the indices of the arguments that appear in the scrutinee of the match expression (line 18). The reason is as follows: recursive calls are typically made in the branches of match expressions (otherwise, the program would never terminate because of unconditional recursion), and the arguments that appear in the scrutinee of the match expression determine which branch to take, deciding whether recursive calls are made further or not. Therefore, it is likely that the arguments that appear in the scrutinee of the match expression are key arguments. Because match may be nested, key arguments in the branches are collected and merged (line 19). If e is not a match expression, the `KeyArgs` function recursively calls itself on the sub-expressions of e and collects the key arguments by unioning the results (line 21).

Given a set of key arguments K and an expression e that may contain a tuple of arguments or a single argument, the function `Struct(e, K)` checks if the arguments are structurally decreasing. If e is a tuple expression (line 8), it first checks if K is empty (line 9). If K is empty, the function returns `false` because there is no key argument to check (line 9). Otherwise, the function extracts the components of e at the indices in K and the corresponding components of x (line 10 and 11). Here, $(e_k)_{k \in K}$ denotes a tuple of $(e_{k_1}, e_{k_2}, \dots)$ where k_1, k_2, \dots are the indices in K and $(x.i)_{i \in K}$ denotes a tuple of $(x.k_1, x.k_2, \dots)$ where k_1, k_2, \dots are the indices in K . The function then checks if the extracted components are structurally decreasing (line 12) using the partial order relation \sqsubseteq defined in Figure 7(b). The partial order relation \sqsubseteq is defined by the rules `ORD_DTOR`, `ORD_PROJ`, and `ORD_TUPLE`. The rule `ORD_DTOR` states that a destructor expression is structurally smaller than the expression it destructs. The rule `ORD_PROJ` states that a projection expression is structurally smaller than the expression it projects if the expression itself is structurally smaller than the projected expression or the two expressions are equal. The rule `ORD_TUPLE` states that a tuple expression $e = (e_1, \dots, e_m)$ is structurally smaller than another tuple expression $e' = (e'_1, \dots, e'_m)$ if any components of e are structurally smaller than the corresponding components of e' and the rest of the components are equal.

Lastly, if e is not a tuple expression (line 13), the function checks if e is structurally smaller than x using the partial order relation \sqsubset (line 14).

Example 16. Consider the solution in Example 15 where the function body is denoted as e_{body} . Because the scrutinee of the match expression is $x.1$, $\text{KeyArgs}(e_{body}) = \{1\}$. The recursive call in the function body is the one in the second branch of the match expression. By the line 5 of the Terminate procedure and the ORD_TUPLE rule, we check if $\text{Cons}^{-1}(x.1).2 \sqsubset x.1$ because $\text{Cons}^{-1}(x.1).2$ is the first component of the argument tuple in the recursive call and the key argument is the first component of the input tuple. By the ORD_PROJ rule, we should check if $\text{Cons}^{-1}(x.1) \sqsubset x.1$, which is true by the ORD_DTOR rule. Therefore, the function is terminating.

Theorem 17 guarantees that if the Terminate procedure accepts a closed hypothesis, then it is guaranteed to terminate on any input.

Theorem 17. If Terminate accepts P , then P is guaranteed to terminate on any input.

Proof Available in the appendix. ■

Comparison with prior work. The prior work on recursion synthesis (Osera and Zdancewic, 2015; Lubin et al., 2020; Miltner et al., 2022) also ensures termination of synthesized programs. But the termination checking in prior work is simpler than ours, limiting the scope of programs that can be synthesized. For example, BURST cannot synthesize tail-recursive programs and MYTH and SMYTH cannot synthesize tail-recursive programs if the first parameter of the target function is a tail-recursive argument, which is non-decreasing in each recursive call. Our termination checking procedure is more general and can handle such cases. More details on the comparison with prior work are discussed in Section 6.2.

4.7 Optimizations

We describe several optimizations that we use to improve the efficiency of our algorithm.

Normalization and type-based pruning. We also utilize a few standard optimizations in prior work. For ease of presentation, we have presented as if we do not type-check any of the expressions during the search. However, in our implementation, we perform type-based pruning to generate only well-typed expressions, similarly to prior work (Feser et al., 2015; Osera and Zdancewic, 2015). We also generate expressions in β -normal η -long form as done in prior work (Osera and Zdancewic, 2015; Frankle et al., 2016; Lubin et al., 2020). Also, we apply constructor/destructor simplification (Lubin et al., 2020) to avoid generating unnecessarily long programs containing sub-expressions of forms $\kappa(\kappa^{-1}(\dots))$ or $\kappa^{-1}(\kappa(\dots))$ for any constructor κ . Finally, in the COMPONENTGENERATION procedure, we avoid generating unnecessary components. When it comes to generating projection expressions, we do not have to generate components of the form $(e_1, \dots, e_k).n$ because they can be replaced by e_n for $1 \leq n \leq k$. Therefore, we only generate projection

components of the form $e.n$ where e is not a tuple. In addition, we do not generate all possible tuples of a product type. For example, if we have m algebraic data types available in the specification and want to generate tuples of length n , the number of all product types of tuples is m^n . We observe that tuples likely to be used in the target program can be used for constructor applications and function calls. Therefore, only product types that appear in the data type definitions and *input types of external (library) functions and the target function* are considered types of tuples to be generated. Lastly, we do not consider nested recursive calls to the target function of the form $f(\dots f(\dots) \dots)$ because they are unlikely to be useful in practice.

Using another version of the D_EXTCALL rule. As described in Example 7 in Section 4.3, the D_EXTCALL rule in Figure 3 may be unsound (i.e., some of generated holes cannot be filled with any expression). This deduction unsoundness may lead to a scalability issue by generating too many unsatisfiable holes if the number of examples is greater than a certain threshold. In such a case, we use the following two rules instead of the D_EXTCALL rule.

$$\begin{array}{c}
 \frac{\begin{array}{l} \Box_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\} \quad e_1, e_2 \in \mathbf{C} \quad \forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_1 \Rightarrow g_j \\ \forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_2 \Rightarrow v_j \quad \forall 1 \leq j \leq n. g_j^{-1}(o_j) = v_j \in \mathcal{I} \end{array}}{e_1 \ e_2 \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \Box_u)} \quad \text{D_EXTCALL1} \\
 \\
 \frac{\begin{array}{l} \Box_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\} \quad e_1, e_2 \in \mathbf{C} \quad \forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_1 \Rightarrow g_j \\ e_2 \text{ contains recursive calls.} \quad \forall 1 \leq j \leq n. \exists v. g_j^{-1}(o_j) = v \in \mathcal{I} \end{array}}{e_1 \ e_2 \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \Box_u)} \quad \text{D_EXTCALL2}
 \end{array}$$

Both rules use components to generate closed expressions without any holes. The difference between the two rules is in whether a component for the argument part contains recursive calls. The D_EXTCALL2 rule considers every component containing recursive calls to be the potential argument of a library function call. This is based on an optimistic assumption that any recursive expression may be a proper argument (like the D_REC rule in Figure 3). These rules do not generate any unsatisfiable holes, but the search space explored in a single iteration of the main loop (lines 2–28 in Algorithm 1) is smaller than the case where the D_EXTCALL rule is used.

5 Implementation

In this section, we describe various implementation details of our synthesis algorithm.

5.1 Preventing unsafe destructor applications

We do not permit potentially unsafe destructors to be used in any candidate program. For example, suppose the following candidate is explored during the search.

```

rec f (x : nat) : nat =
  match x with Z -> x | S _ -> f (S-1(S-1(x)))
  
```

This program is considered invalid and not generated during the search. In the branch of S , the expression $S^{-1}(S^{-1}(x))$ is not allowed because the pattern match does not guarantee that the input x is of the form $S(S(\dots))$. The only safe destructor usable in the branch of S is $S^{-1}(x)$. To prevent such unsafe destructor applications, $\text{Deduce}(\mathbf{C}, \mathcal{I}, \square_u)$ on line 19 in Algorithm 1 is changed to $\text{Deduce}(\text{RemoveUnsafeComp}(P, \square_u, \mathbf{C}), \mathcal{I}, \square_u)$ where $\text{RemoveUnsafeComp}(P, \square_u, \mathbf{C})$ returns a subset of \mathbf{C} whose components do not contain any unsafe destructor applications at the position of \square_u in P . This filtering prevents components of unsafe destructor applications from being used to fill the hole \square_u in the candidate program P .

5.2 Ensuring termination of block and candidate generation

In our implementation, to guarantee the termination of the BlockGen procedure and the Deduce procedure, we limit the maximum number of subsequent steps of deduction to a certain number.⁵

In other words, to avoid generating infinitely many open hypotheses from a given initial hypothesis, we permit the Deduce procedure to be terminated after a certain number of steps of applications of rules in Figure 3. Because the BlockGen procedure relies on the Deduce procedure, this also guarantees termination of the BlockGen procedure. Note that despite this finitization, the search space is still infinite because there is no limit on the maximum component size (i.e., the component pool will keep growing until a solution is found).

5.3 Program selection

In order to synthesize likely programs, we utilize a cost function: the cost of each candidate program $P = \text{rec } f(x) = e$ is determined by the cost of its body e (denoted $\mathcal{C}(e)$), which is a nonnegative number. Costs of expressions satisfy the following constraints (some cases are omitted):

- $\mathcal{C}(e_1 \ e_2) > \mathcal{C}(e_1) + \mathcal{C}(e_2)$
- $\mathcal{C}(\kappa(e_1, \dots, e_k)) > \sum_{1 \leq j \leq k} \mathcal{C}(e_j)$
- $\mathcal{C}(x) = 0$
- $\mathcal{C}(e.n) = \mathcal{C}(e)$
- $\mathcal{C}(\kappa^{-1}(e)) = \mathcal{C}(e)$

Intuitively, we penalize the use of constructors (thereby constants) and prioritize the use of variables, destructors, and projections. The reason for this is that they are used to extract subcomponents of constructors, which are essentially the same as variables bound by patterns in match expressions. For example, consider the solution P_{sol} of the overview example problem in Section 2.

```
rec mul (x : nat, y : nat) : nat = match x with Z -> Z | S _ -> add (mul
    (S-1(x), y), y)
```

This program can be written as the following program by introducing a new variable x' bound by the pattern of S .

⁵ The Deduce procedure may derive additional holes for a given hole and a sketch with the added holes is added into the queue in line 23 in Algorithm 1. This may cause the algorithm to not terminate if the Deduce procedure keeps generating holes because the queue will never be empty.

```

rec mul (x : nat, y : nat) : nat = match x with Z -> Z | S x' -> add (mul
  (x', y), y)

```

Note that $S^{-1}(x)$ and x' play the same role. Therefore, destructors and projections can be understood as variables in many cases, and prioritizing variables has been proved to be a good heuristic to avoid overfitting (Feser et al., 2015; Gulwani, 2011). Lastly, in case of a tie, we pick a smaller program in terms of AST size. This heuristic has also been popularly used in the majority of previous approaches (Albarghouthi et al., 2013; Feser et al., 2015; Wang et al., 2017; Miltner et al., 2022).

5.4 Checking block consistency

We describe implementation details for improving the pruning power of the BlockConsistent procedure. In Algorithm 1, when choosing a hole in a hypothesis of a match expression, we prefer holes for base cases to ones for inductive cases. By filling holes for base cases first, we can effectively prune infeasible recursive hypotheses. For example, suppose we encounter the following hypothesis while synthesizing mul in Section 2.

```

P1 = rec mul (x : nat, y : nat) : nat = match x with Z -> □1 | S _ -> mul
  (S-1(x), □2)

```

Suppose we first fill the hole □₂ with y, obtaining the following hypothesis.

```

P2 = rec mul (x : nat, y : nat) : nat = match x with Z -> □1 | S _ -> mul
  (S-1(x), y)

```

Note that this hypothesis cannot become a solution no matter what expression we put in the remaining hole. To check block consistency, for every input, we will obtain the open block □₁ as a result of our symbolic evaluation. Although the hypothesis P₂ is infeasible, because a hole can be matched with any expression according to the rules in Figure 6, the hypothesis will be determined to be block consistent with respect to any set of blocks, and will not be pruned.

Now, suppose we first fill the hole □₁ in P₁ with x, obtaining the following hypothesis.

```

P3 = rec mul (x : nat, y : nat) : nat = match x with Z -> x | S _ -> mul
  (S-1(x), □2)

```

Note that this hypothesis also cannot become a solution no matter what expression we put in the remaining hole. For every input, we will obtain a closed block x as a result of our symbolic evaluation. Because x is not included in the blocks for the example $(1, 2) \mapsto 2$, the hypothesis P₃ is determined to be block inconsistent and will be pruned.

5.5 Finding a solution vs. all solutions

We allow the user to choose between finding all possible solutions synthesizable using some set of components and picking the best one, and stopping the search as soon as a solution is found. This allows the user to find a good balance between speed of synthesis and accuracy (i.e., the possibility of generating intended programs). Finding all solutions and picking the best one only requires a slight modification in Algorithm 1 as follows:

even if a solution is found on line 14, the main loop continues to explore the search space until the queue becomes empty. Then, we pick the best one among the multiple solutions found so far. In addition, we can further expedite the process of finding a single solution with a slight modification in Figure 3. Instead of including all components consistent with a given set of input–output examples in the `D_COMPONENT` rule, we just include a single component whose score is the best among the satisfying components.

6 Evaluation

We have implemented our approach in a tool `TRIO`. `TRIO` consists of about 4K lines of OCaml code. We evaluate `TRIO` on synthesis tasks used in prior work and on new tasks collected from an online tutorial. We aim to answer the following research questions:

- **RQ1:** How does `TRIO` perform on various synthesis tasks?
- **RQ2:** How does `TRIO` compare with existing techniques for recursive program synthesis?
- **RQ3:** How effective are block-based pruning and library sampling for accelerating synthesis?

All of our experiments were conducted on a 2.0 GHz Intel Core i5 processor with 16GB of memory running macOS Big Sur. We set the timeout limit to 120 seconds for each synthesis task.

6.1 Experimental setup

Benchmarks. We use 65 recursive functional programs. 45 out of 65 have been used to evaluate prior work (Osera and Zdancewic, 2015; Lubin et al., 2020; Miltner et al., 2022). The remaining 20 programs are from the exercises in the official OCaml online tutorial and their slight variants. The details can be found in Table 1.

For these benchmarks, we consider the following two classes of specifications to evaluate `TRIO` over different specifications.

- **IO:** We use input–output examples written by developers of `SMYTH` (Lubin et al., 2020) and `BURST` (Miltner et al., 2022).
- **Ref:** For 45 benchmarks, we use reference implementations from prior work written by developers of `BURST` (Miltner et al., 2022). For the other 20 benchmarks, we use reference implementations from the official OCaml online and the ones written by us.

Baselines. We compare `TRIO` to state-of-the-art tools for synthesizing recursive functional programs. `SMYTH` (Lubin et al., 2020) performs top-down synthesis from input–output examples. It performs partial evaluation to propagate constraints from partial programs to the remaining holes. `BURST` (Miltner et al., 2022) performs bottom-up

Table 1. List of new 20 benchmarks collected from the exercises in the official OCaml online tutorial (<https://ocaml.org/exercises>) and their variants

	Name	External Operators	Description
Arithmetic	nat_mul	add	Multiplication of two natural numbers.
	nat_sub		Subtraction of two natural numbers.
	nat_fac_tailcall	add, mul	Factorial of a natural number using tail recursion.
	nat_fib_tailcall	add	Fibonacci number using tail recursion.
	expr_boolean	not, and, or	Logical expression evaluator.
	expr	add, mul	Calculator using addition and multiplication operators.
	expr_sub	add, mul, sub	Calculator using addition, subtraction, and multiplication operators.
	expr_div	add, mul, sub, div	Calculator using addition, subtraction, multiplication, and division operator.
Lists	list_dropeven	is_even	Drop the even number(s) in a list.
	list_last2		Return the last two elements in a list.
	list_make		Return the 0 padded list of length n
	list_range	compare	Return the sequence of the given numbers in descending order.
	list_append_tailcall		Append two lists using tail recursion.
	list_length_tailcall		Return the length of a list using tail recursion.
	list_sum_tailcall	add	Return the sum of numbers in a list using tail recursion.
Trees	tree_height	compare, max	Return the height of a tree.
	tree_balanced	compare, max, tree_height, and	Check whether a tree is height-balanced.
	tree_lastleft		Return the node at the left end of the tree.
	tree_notexist	compare, and	Check if a number is in a tree
	tree_sum	add	Return the sum of numbers in a tree.

synthesis from input-output examples or logical specifications. Neither of them requires trace-complete specifications. We aim to confirm the benefits of our bidirectional search strategy by comparing TRIO against the top-down synthesizer SMYTH and the bottom-up synthesizer BURST. The result of the comparison to SMYTH and BURST is presented in all the following sections except for Section 6.6. In addition, we compare TRIO to SYRUP (Yuan et al., 2023) but in a different aspect compared to the above two synthesizers. SYRUP uses version space algebra to avoid overfitting when synthesizing recursive programs. It focuses on minimizing the number of examples required for synthesizing the desired programs rather than the synthesis time. Also, SYRUP is known to be less sensitive to the quality of the input-output examples. Therefore, we compare against SYRUP focusing on the quantity and quality of examples required for synthesizing the desired programs. The result of the comparison to SYRUP is presented in Section 6.6.

6.2 Effectiveness of TRIO for input-output examples

We evaluate TRIO on synthesis problems with **IO** specifications. The initial component size n for TRIO is set to be 6. For each instance, we measure the running time of TRIO and the size of the synthesized program.

The results are summarized in Table. 2. The column “Correct” shows if the synthesized program is the one intended by the user. We manually checked if the synthesized program is semantically equivalent to the known solution for each problem.

TRIO outperforms the other baselines in terms of both the number of solved problems and synthesis time. TRIO can synthesize 65 out of 65 problems, with an average time of 1.5 seconds. On the other hand, BURST and SMYTH can synthesize 54 and 53 problems, with average times of 2.6 and 2.7 seconds, respectively. In addition, TRIO is the fastest tool in 42 problems, whereas BURST and SMYTH are the fastest tools in 14 and 27 problems, respectively.⁶

For every problem, the time taken for synthesizing blocks and constructing inverse maps is negligible (usually less than a second).

We observe BURST and SMYTH occasionally take a large amount of memory, whereas TRIO only requires a small amount of memory. While solving all of the tasks, the peak memory usage of BURST is 5.3GB and that of SMYTH is 1.2GB. On the other hand, TRIO only requires 88 MB even in the worst case. On average, BURST and SMYTH use 647 and 40 MB of memory respectively, whereas TRIO uses 24 MB. Thus, we can conclude TRIO is more memory efficient than the other baselines.

Thanks to the performance gain of TRIO, we can synthesize programs that are hard for the other baselines. The problem `expr_div` is hard in that it requires complex pattern matching involving many external operators. It concerns synthesizing a simple calculator with addition, subtraction, multiplication, and division. The specification is given as follows:

```
type nat = Z | S of nat
type expr = NAT of nat | ADD of expr and expr | SUB of expr and expr
           | MUL of expr and expr | DIV of expr and expr

rec add (x:nat, y:nat) : nat = ... rec sub (x:nat, y:nat) : nat = ...
rec mul (x:nat, y:nat) : nat = ... rec div (x:nat, y:nat) : nat = ...

rec eval (x:nat, y:nat) : nat =  $\square_{in}$ 
```

where \square_{in} is the set of input–output examples $\{\text{NAT } 1 \mapsto 1, \text{ADD}(\text{NAT } 1, \text{NAT } 2) \mapsto 3, \dots\}$, and `add`, `sub`, `mul`, and `div` are the external operators. Finding the following solution is non-trivial because there are extremely many possible combinations of recursive calls, external operators, and case matching.

```
rec eval (x : expr) : nat =
  match x with NAT n -> n
  | ADD (x1, x2) -> add (eval x1, eval x2)
  | SUB (x1, x2) -> sub (eval x1, eval x2)
  | MUL (x1, x2) -> mul (eval x1, eval x2)
  | DIV (x1, x2) -> div (eval x1, eval x2)
```

⁶ If there are ties in a synthesis problem, all tools with the same synthesis time are considered to be the fastest.

Table 2. Results for the **IO** benchmark suite (with 15 easy problems omitted), where “Time” gives synthesis time in seconds, and “Size” shows the size of the synthesized program (measured by number of AST nodes). Synthesis time of the fastest tool for each problem is highlighted in bold.

Benchmark	TRIO			BURST			SMYTH		
	Time(s)	Size	Correct	Time(s)	Size	Correct	Time(s)	Size	Correct
bool_band	0.01	14	✓	0.02	14	✓	0.02	23	✓
list_append	0.05	44	✗	0.15	31	✓	0.02	33	✓
list_compress	0.14	56	✓	0.74	56	✓	Timeout	N/A	N/A
list_drop	0.03	31	✓	2.78	24	✓	0.05	32	✓
list_even_parity	0.02	25	✓	0.05	23	✓	0.03	22	✓
list_filter	0.03	47	✗	0.31	46	✗	0.08	44	✗
list_fold	0.08	38	✗	0.13	43	✗	1.42	49	✗
list_map	0.03	37	✓	0.03	37	✓	0.39	51	✗
list_pairwise_swap	0.03	38	✓	0.32	38	✓	0.02	27	✓
list_rev_append	0.08	35	✗	1.32	22	✓	0.06	20	✓
list_rev_fold	0.01	10	✓	0.03	10	✓	0.04	15	✓
list_rev_snoc	0.05	18	✓	1.43	18	✓	0.02	16	✓
list_rev_tailcall	0.02	37	✗	0.17	43	✗	0.02	33	✓
list_snoc	0.03	36	✓	1.9	36	✓	0.03	37	✓
list_sort_sorted_insert	0.04	18	✓	0.08	18	✓	0.02	16	✓
list_sorted_insert	0.46	60	✓	Timeout	N/A	N/A	1.53	55	✓
list_stutter	0.02	23	✓	4.3	23	✓	0.03	21	✓
list_sum	0.01	10	✓	0.08	10	✓	0.04	10	✓
list_take	0.04	38	✓	Timeout	N/A	N/A	0.04	38	✓
nat_iseven	0.02	16	✓	0.06	16	✓	0.02	13	✓
nat_max	0.19	23	✓	0.27	23	✓	0.07	34	✓
tree_bininsert	1.59	87	✓	3.53	87	✓	Timeout	N/A	N/A
tree_collect_leaves	0.06	27	✓	0.49	27	✓	0.04	24	✓
tree_count_leaves	0.06	22	✓	0.24	22	✓	0.48	25	✓
tree_count_nodes	0.14	22	✓	0.11	22	✓	0.18	20	✓
tree_inorder	0.12	27	✓	8.21	27	✓	0.07	24	✓
tree_map	0.07	47	✗	0.41	49	✓	0.94	61	✓
tree_nodes_at_level	0.83	47	✓	35.52	47	✓	Timeout	N/A	N/A
tree_postorder	0.47	32	✓	3.21	32	✓	Timeout	N/A	N/A
tree_preorder	0.11	27	✓	0.1	27	✓	0.09	24	✓
expr_boolean	6.78	59	✓	0.58	52	✗	Timeout	N/A	N/A
expr	0.67	36	✓	Timeout	N/A	N/A	Timeout	N/A	N/A
expr_sub	3.48	51	✓	Timeout	N/A	N/A	11.46	58	✗
expr_div	26.99	66	✓	Timeout	N/A	N/A	Timeout	N/A	N/A
list_dropeven	0.03	28	✓	Timeout	N/A	N/A	0.05	25	✓
list_last2	0.04	40	✓	0.11	39	✓	0.09	29	✓
list_make	0.02	14	✓	0.78	14	✓	0.01	13	✓
list_range	0.24	60	✓	Timeout	N/A	N/A	Timeout	N/A	Timeout
nat_mul	1.7	27	✓	Timeout	N/A	N/A	109.46	31	✓
nat_sub	0.13	29	✓	8.34	29	✓	0.04	30	✓
tree_balanced	9.48	46	✗	21.12	37	✗	Timeout	N/A	N/A
tree_height	0.51	23	✓	0.22	22	✓	16.51	20	✓
tree_lastleft	0.19	21	✓	0.03	21	✓	0.03	18	✓
tree_notexist	24.2	79	✓	0.79	79	✓	Timeout	N/A	N/A
tree_sum	0.36	28	✓	35.56	28	✓	0.47	25	✓

Table 2. Continued.

Benchmark	TRIO			BURST			SMYTH		
	Time(s)	Size	Correct	Time(s)	Size	Correct	Time(s)	Size	Correct
list_append_tailcall	0.11	51	✗	Timeout	N/A	N/A	0.02	33	✓
list_length_tailcall	0.02	24	✓	0.06	24	✓	0.02	30	✓
list_sum_tailcall	0.04	32	✓	3.33	32	✓	0.04	31	✗
nat_fac_tailcall	2.04	27	✓	Timeout	N/A	N/A	Timeout	N/A	N/A
nat_fib_tailcall	14.01	56	✓	Timeout	N/A	N/A	Timeout	N/A	N/A
# Solved (# correct)	65 (57)			54 (49)			53 (48)		
# Timeout	0			11			12		

However, TRIO can find the solution in 27 seconds.⁷ On the other hand, the other baselines fail to solve all the problems that concern synthesizing calculators (i.e., `expr`, `expr_sub`, `expr_div`).

Analysis of overfitting. We manually inspect the programs synthesized by the three tools to investigate how they are prone to overfitting. 57 out of 65 programs (88%) synthesized by TRIO are the intended ones. 49 out of 54 programs (91%) and 48 out of 53 programs (91%) synthesized by BURST and SMYTH are the intended ones, respectively. Therefore, all the tools are roughly equal in terms of solution quality.

We can mitigate overfitting by making TRIO find all solutions that can be found with a current set of components and choose the best one according to the cost described in Section 5.3. For the 8 problems for which TRIO synthesizes unintended programs, if TRIO is configured to find all solutions and pick the best one, it could find the desired programs for 7 problems except for `list_rev_append` at the cost of overhead ranging from a second to a few minutes. In the case of `list_rev_append`, the specification is not constraining enough for finding the solution. In the experiment with reference implementations based on CEGIS where additional input–output examples can be provided whenever the synthesizer fails, we confirm that TRIO successfully finds the desired solution.

Tail-recursive functions. TRIO can synthesize all of the 6 tail-recursive benchmarks (the benchmarks with the suffix `_tailcall`) thanks to the termination checking mechanism that permits tail-recursive calls. On the other hand, neither of the other baselines can synthesize all of the tail-recursive benchmarks.

BURST cannot synthesize tail-recursive calls because its termination checker is based on the default value order which does not permit tail-recursive calls. For example, `list_rev_tailcall` requires a recursive call on `([2],[1])` for input `([1;2],[1])`, but its value ordering does not consider `([2],[1])` to be strictly smaller than `([1;2],[1])`. However, it produces the correct solution for some of the tail-recursive benchmarks (`list_sum_tailcall` and `list_length_tailcall`) by finding a non-tail-recursive

⁷ In our previous work (Lee and Cho, 2023), we could not synthesize the solution for this problem. The performance improvement is due to the optimization of the implementation of TRIO.

solution that is semantically equivalent to the tail-recursive one. For example, BURST synthesizes the following solution for `list_sum_tailcall`:

```

rec f (x : list * nat) : nat =
  match x.1 with
  | Nil _ -> x.2
  | Cons _ -> add Cons-1(x.1).1 (f (Cons-1(x.1).2, x.2))

```

which is not tail-recursive but semantically equivalent to the tail-recursive solution.

SMYTH employs a check during synthesis to ensure that the argument to a recursive function recursive call is a strict subterm of the parameter to the recursive call. However, all functions in SMYTH are single-parameter functions, and multi-parameter functions are curried. As a result, only recursive calls that are structurally decreasing on the first parameter of multi-parameter (curried) functions are allowed (Lubin, 2020). This restriction limits the scope of programs that can be synthesized by SMYTH. As an evidence, we have tried to put the tail-recursive argument (i.e., the accumulator) as the first parameter for `list_sum_tailcall`. As expected, Smyth fails to synthesize the solution because the first parameter is not strictly decreasing within the timeout limit. Burst also times out for the same reason. However, TRIO can synthesizing the solution. This observation suggests that the termination checking mechanism in TRIO is more flexible than the one in SMYTH.

The overhead of the termination checking mechanism in TRIO is negligible. On average, the termination checking mechanism takes 0.05 seconds. With the exceptions of `tree_notexist` that require 1.5 seconds respectively because of the large number of candidates explored during the search, the termination checking mechanism takes less than 0.1 seconds for all the other benchmarks.

Summary of results. When synthesizing recursive programs from input–output examples, TRIO outperforms state-of-the-art baseline tools in terms of both synthesis time and memory usage. Also, TRIO solves harder synthesis problems beyond the reach of the baselines.

6.3 Effectiveness of TRIO for reference implementations

In this section, we evaluate TRIO on synthesis problems with Ref specifications. We follow the same evaluation procedure as the evaluation of BURST (Miltner et al., 2022) for **Ref** specifications. The authors of BURST integrated BURST and SMYTH into a CEGIS loop and, for each candidate program proposed by each tool, they use the verifier to determine whether the candidate is semantically equivalent to the reference implementation. If not, a new input–output example comprising a counterexample input generated by the verifier and its corresponding output is added.⁸ This process is repeated until the desired program is found.

The goal of this experiment is to confirm how the tools deal with the random examples generated by the verifier, rather than hand-crafted examples.

The results are summarized in Table 3. The column “# Iters” shows the number of CEGIS iterations required until a solution is found. TRIO also outperforms the other

⁸ The authors of BURST use bounded testing instead of verification and manually checked the semantic equivalence between the generated programs and the reference implementation. We use the same method by reusing the artifacts of BURST.

Table 3. Results for the **Ref** benchmark suite where “# Iter” shows the number of CEGIS iterations.

Benchmark	TRIO			BURST			SMYTH		
	Time(s)	Size	# Iter	Time(s)	Size	# Iter	Time(s)	Size	# Iter
bool_band	0.02	14	3	0.02	14	3	0.02	23	3
list_append	0.48	31	6	0.4	31	6	0.66	33	8
list_compress	1.07	56	10	1.53	56	9	Timeout	N/A	N/A
list_drop	0.36	31	5	1.46	31	6	Timeout	N/A	N/A
list_even_parity	0.13	25	5	0.15	23	6	0.16	22	6
list_filter	0.91	56	8	1.13	56	8	Timeout	N/A	N/A
list_fold	0.86	42	6	12.25	42	6	Timeout	N/A	N/A
list_map	0.67	37	6	0.7	37	4	Timeout	N/A	N/A
list_pairwise_swap	0.36	38	7	0.47	38	6	Timeout	N/A	N/A
list_rev_append	0.96	22	5	11.65	22	5	Timeout	N/A	N/A
list_rev_fold	0.76	10	3	0.79	10	2	Timeout	N/A	N/A
list_rev_snoc	1.02	18	5	6.34	18	4	Timeout	N/A	N/A
list_rev_tailcall	Timeout	N/A	N/A	Timeout	N/A	N/A	0.52	33	9
list_snoc	0.35	36	3	0.43	36	3	0.48	36	7
list_sort_sorted_insert	1.33	18	6	1.65	18	5	1.32	16	6
list_sorted_insert	0.9	60	6	1.03	60	6	13.87	69	12
list_stutter	0.36	23	3	0.99	23	3	0.49	21	4
list_sum	0.82	10	2	0.71	10	2	0.8	10	2
list_take	0.38	38	9	15.74	38	9	0.42	38	7
nat_iseven	0.02	16	4	0.02	16	4	0.02	13	4
nat_max	0.29	23	5	4.88	23	5	0.15	34	7
tree_bininsert	7.65	87	7	18.88	87	8	Timeout	N/A	N/A
tree_collect_leaves	5.7	27	4	6.54	27	5	7.8	24	5
tree_count_leaves	5.42	22	4	7.01	22	4	8.27	25	4
tree_count_nodes	5.66	22	4	5.83	22	4	6.83	20	4
tree_inorder	8.09	27	4	12.23	27	5	8.98	24	6
tree_map	6.07	49	7	7.1	49	5	Timeout	N/A	N/A
tree_nodes_at_level	5.88	47	7	Timeout	N/A	N/A	26.21	43	5
tree_postorder	9.75	32	6	10.54	32	7	Timeout	N/A	N/A
tree_preorder	10.18	27	6	39.27	27	5	9.38	24	5
expr_boolean	33.22	59	14	7.7	59	21	Timeout	N/A	N/A
expr	4.36	36	9	Timeout	N/A	N/A	Timeout	N/A	N/A
expr_sub	28	51	12	Timeout	N/A	N/A	Timeout	N/A	N/A
expr_div	Timeout	N/A	N/A	Timeout	N/A	N/A	Timeout	N/A	N/A
list_dropeven	0.43	28	6	0.37	28	6	Timeout	N/A	N/A
list_last2	0.72	40	7	1.9	39	5	1.02	29	6
list_make	0.04	14	3	0.03	14	3	0.03	13	3
list_range	1.04	49	7	118.64	49	6	Timeout	N/A	N/A
nat_mul	0.96	27	6	77.55	27	8	9.01	31	7
nat_sub	0.29	29	6	1.25	29	7	Timeout	N/A	N/A
tree_balanced	81.09	49	9	Timeout	N/A	N/A	Timeout	N/A	N/A
tree_height	8.3	23	5	7.71	22	5	7.79	20	4
tree_lastleft	9.4	21	6	7.92	21	5	8.65	18	5
tree_notexist	11.76	79	8	12.07	79	9	Timeout	N/A	N/A
tree_sum	9.62	28	5	11.62	28	5	9.03	25	5
list_append_tailcall	Timeout	N/A	N/A	Timeout	N/A	N/A	0.55	33	9
list_length_tailcall	0.33	24	4	0.31	24	4	0.45	29	5

Table 3. Continued.

Benchmark	TRIO			BURST			SMYTH		
	Time(s)	Size	# Iter	Time(s)	Size	# Iter	Time(s)	Size	# Iter
list_sum_tailcall	0.67	32	5	2.53	32	6	1.07	34	6
nat_fac_tailcall	Timeout	N/A	N/A	Timeout	N/A	N/A	Timeout	N/A	N/A
nat_fib_tailcall	Timeout	N/A	N/A	Timeout	N/A	N/A	Timeout	N/A	N/A
# Solved	60			56			42		
# Timeout	5			9			23		

baselines in terms of both the number of solved instances and synthesis time. TRIO can synthesize 60 instances, with an average time of 4.5 seconds and an average number of CEGIS iterations of 5.2. BURST can synthesize 56 instances, with an average time of 5.0 seconds and an average number of CEGIS iterations of 5. SMYTH can synthesize 42 instances, with an average time of 3.0 seconds and an average number of CEGIS iterations of 4.7.

Overall these results suggest that TRIO can deal better with random examples generated by the verifier compared to the other baselines.

Failure analysis. The timeout on 5 problems is due to many CEGIS iterations. Because TRIO often synthesizes an overfit solution for these problems, the verifier generates many counterexamples. As the number of input–output examples increases, the time required for each CEGIS iteration increases. This result suggests that TRIO can be improved by adopting a better strategy for avoid overfitting.

Summary of results. Also when synthesizing recursive programs from reference implementations, TRIO outperforms the other baselines in terms of both the number of solved instances and synthesis time. The results suggest the TRIO’s robustness to randomly given examples.

6.4 Ablation study for block-based pruning and library sampling

We now evaluate the effectiveness of the block-based pruning and library sampling techniques used by TRIO. For this purpose, we compare the performance of four variants of TRIO, each using a different combination: TRIO with block-based pruning and library sampling, TRIO^B only with block-based pruning, TRIO^L only with library sampling, and TRIO[–] with both techniques disabled.

Table 4 summarizes the results of this ablation study (more detailed results can be found in cactus plots in Figure 8). For each variant of TRIO, we report the number of solved benchmarks with the **IO** and **Ref** specifications, respectively. In this experiment, we only consider the 20 newly added benchmarks because we realize the other 45 benchmarks from prior work are easy, so that they can be quickly solved by all the variants of TRIO. We conjecture that the reason why even TRIO[–] can solve all of the 45 benchmarks is that it enjoys the benefit of the synergistic combination of top-down and bottom-up search

Table 4. Number of instances that can be solved by four variants of TRIO among **20** newly added benchmarks

	TRIO	TRIO ^B	TRIO ^L	TRIO ⁻⁻
# Solved (IO spec.)	20 (100%)	18 (90%)	15 (75%)	14 (70%)
# Solved (Ref spec.)	16 (80%)	14 (70%)	12 (60%)	11 (55%)

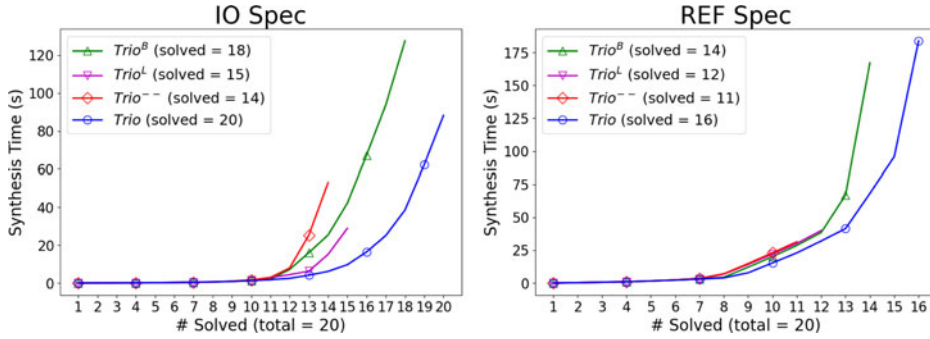


Fig. 8. Comparison of different variants of TRIO.

strategies. As can be seen in the table, TRIO with the two techniques solves more benchmarks than the other three variants. TRIO solved 100% of the new benchmarks with IO specifications, whereas TRIO⁻⁻ could solve 70% of the benchmarks. Such a trend can also be observed in the reference implementation experiment. We notice the efficacy of block-based pruning is higher than that of library sampling because the difference between TRIO⁻⁻ and TRIO^B is more significant than the difference between TRIO⁻⁻ and TRIO^L.

6.5 Benefits of our method compared to prior work

In this section, we analyze why our method outperforms the previous methods. As a representative example, we investigate how the tools work for a simpler version of the expr benchmark where TRIO can quickly find the solution in contrast to the other baselines.

```

type nat = Z | S of nat
type expr = NAT of nat | ADD of expr and expr
rec add (x : nat, y : nat) : nat = ...      rec mul (x : nat, y : nat) :
  nat = ...
rec eval (x : nat, y : nat) : nat =  $\square_{in}$ 

```

where $\square_{in} = \{i_1 \mapsto 1, i_2 \mapsto 4, i_3 \mapsto 7\}$, $i_1 = \text{NAT } 1$, $i_2 = \text{ADD}(\text{NAT } 3, \text{NAT } 1)$, and $i_3 = \text{ADD}(\text{NAT } 4, \text{NAT } 3)$. The solution in our language (depicted in Figure 2) is as follows:

```

rec eval (x : expr) : nat =
  match x with NAT _ -> NAT-1(x) | ADD _ -> add (eval ADD-1(x).1, eval
    ADD-1(x).2)

```

Comparison to SMYTH. Similarly to our method, SMYTH explores the search space by performing top-down propagation. There are two major differences between SMYTH

and our method. First, whenever a hole is filled with some expression, SMYTH “updates” the other remaining holes according to the hole filling, so that the holes are more likely to be filled with the correct expressions. Second, SMYTH solely relies on a top-down search strategy without any bottom-up search. We observe these two differences are the main reasons why SMYTH fails to solve the benchmark. Consider the following hypothesis generated by SMYTH during the search.

$$P_1 = \text{rec eval } (x : \text{expr}) : \text{nat} = \text{match } x \text{ with NAT } _ \rightarrow \square_1 \mid \text{ADD } _ \rightarrow \square_2$$

where $\square_1 = \{i_1 \mapsto 1\}$, $\square_2 = \{i_2 \mapsto 4, i_3 \mapsto 7\}$. SMYTH further refines the hypothesis P_1 by filling the hole \square_2 with a recursive call to `eval`. Like TRIO, SMYTH enumerates structurally-decreasing recursive calls to guarantee the termination of synthesized programs. Suppose the following hypothesis is generated.

$$P_2 = \text{rec eval } (x : \text{expr}) : \text{nat} = \text{match } x \text{ with NAT } _ \rightarrow \square_1 \mid \text{ADD } _ \rightarrow \text{eval ADD}^{-1}(x).1$$

Obviously, the hypothesis P_2 is not desired because it cannot be the solution. However, SMYTH cannot detect this problem for the following reason. As hole \square_2 is filled, SMYTH updates the other remaining hole. SMYTH generates the following hypothesis.

$$P_3 = \text{rec eval } (x : \text{expr}) : \text{nat} = \text{match } x \text{ with NAT } _ \rightarrow \square_3 \mid \text{ADD } _ \rightarrow \text{eval ADD}^{-1}(x).1$$

where $\square_3 = \square_1 \cup \{\text{NAT } 3 \mapsto 4, \text{NAT } 4 \mapsto 7\}$. The additional examples $\{\text{NAT } 3 \mapsto 4, \text{NAT } 4 \mapsto 7\}$ are originated from the original examples $\{i_2 \mapsto 4, i_3 \mapsto 7\}$ and partial evaluation of P_2 with the input examples. SMYTH refines the hole \square_3 by generating the following hypothesis.

$$P_4 = \text{rec eval } (x : \text{expr}) : \text{nat} = \text{match } x \text{ with NAT } _ \rightarrow \text{match S}^{-1}(\text{NAT}^{-1}(x)) \text{ with Z } \rightarrow \square_4 \mid \text{S } _ \rightarrow \square_5 \mid \text{ADD } _ \rightarrow \text{eval ADD}^{-1}(x).1$$

where $\square_4 = \{i_1 \mapsto 1\}$ and $\square_5 = \{\text{NAT } 3 \mapsto 4, \text{NAT } 4 \mapsto 7\}$. SMYTH keeps refining this hypothesis, which is fruitless. In summary, SMYTH’s updating holes by partial evaluation sometimes makes the search more difficult.

On the other hand, TRIO can quickly identify the infeasibility of P_2 as follows: TRIO does not update the other remaining hole \square_1 after filling \square_2 . Then, the hole \square_1 can be easily filled with $\text{NAT}^{-1}(x)$, generating the following program, which can be easily proved to be infeasible by concrete evaluation.

$$P'_3 = \text{rec eval } (x : \text{expr}) : \text{nat} = \text{match } x \text{ with NAT } _ \rightarrow \text{NAT}^{-1}(x) \mid \text{ADD } _ \rightarrow \text{eval ADD}^{-1}(x).1$$

In addition, we note that another source of inefficiency of SMYTH is that it redundantly generates many semantically equivalent hypotheses. For instance, the followings are some of hypotheses generated by SMYTH by filling \square_1 in P_1 with different expressions. The more library functions are usable, the more redundant hypotheses are generated.

$$\begin{array}{ll} \text{match } x \text{ with NAT } _ \rightarrow \text{NAT}^{-1}(x) & \mid \text{ADD } _ \rightarrow \square_2 \\ \text{match } x \text{ with NAT } _ \rightarrow \text{add } (\text{NAT}^{-1}(x), 0) & \mid \text{ADD } _ \rightarrow \square_2 \\ \text{match } x \text{ with NAT } _ \rightarrow 1 & \mid \text{ADD } _ \rightarrow \square_2 \\ \text{match } x \text{ with NAT } _ \rightarrow \text{add } (1, (\text{mul } (0, 0))) & \mid \text{ADD } _ \rightarrow \square_2 \end{array}$$

Note that the first two hypotheses and the last two hypotheses are semantically equivalent respectively. However, TRIO avoids generating such redundant hypotheses because *Bottom-up enumerator* exploits observational equivalence to avoid generating multiple components of the same behaviors.

Comparison to BURST. BURST performs bottom-up synthesis with *angelic execution*.⁹ It first synthesizes a program assuming any recursive calls to the function being synthesized *angelically* behave to make the program correct. Then, it checks if the assumptions made in the previous step are correct. If they are, the solution is found. Otherwise, those assumptions are refuted and never made again by being added to the list of *anti-specifications*. This process is repeated, progressively strengthening the specification of the target function and eventually leading to a solution.

BURST fails to find the solution because it runs into extensive backtracking (i.e., too many steps of specification strengthening). Specifically, given the specification of the target function $\text{eval } i_1 = 1 \wedge \text{eval } i_2 = 4 \wedge \text{eval } i_3 = 7$ (which is from the three input-output examples), BURST first enumerates the following candidate program.

```
rec eval (x : expr) : nat = match x with NAT _ -> NAT-1(x) | ADD _ ->
  eval ADD-1(x).1
```

Obviously, the above program does not satisfy the second and third input-output examples. However, at this stage, BURST generates a program assuming any recursive call to `eval` can return anything to satisfy the constraints. The above program is generated by assuming $\text{eval } (\text{NAT } 3) = 4 \wedge \text{eval } (\text{NAT } 4) = 7$. BURST then checks if this assumption is correct. It clearly does not because $\text{eval } (\text{NAT } 3) = 3$ and $\text{eval } (\text{NAT } 4) = 4$. Then, it re-attempts synthesis with the following strengthened specification.

$$\text{eval } i_1 = 1 \wedge \text{eval } i_2 = 4 \wedge \text{eval } i_3 = 7 \wedge \text{eval } (\text{NAT } 3) = 4 \wedge \text{eval } (\text{NAT } 4) = 7.$$

After the search within a bounded space, BURST fails to find a program that satisfies the strengthened specification. It concludes that the assumption made in the previous step is incorrect. Then, it adds the negation of the assumption (i.e., $\neg(\text{eval } (\text{NAT } 3) = 4 \wedge \text{eval } (\text{NAT } 4) = 7)$) into the list of anti-specifications, searches for a program that does not violate the anti-specifications and generates the following program.

```
rec eval (x : expr) : nat =
  match x with NAT _ -> NAT-1(x) | ADD _ -> S (S (eval ADD-1(x).1))
```

Obviously, this program also does not satisfy the original specification. However, it does not violate the anti-specification because the above program is correct assuming $\text{eval } (\text{NAT } 3) = 2 \wedge \text{eval } (\text{NAT } 4) = 5$. Then, it re-attempts synthesis with the following strengthened specification.

$$\text{eval } i_1 = 1 \wedge \text{eval } i_2 = 4 \wedge \text{eval } i_3 = 7 \wedge \neg(\text{eval } (\text{NAT } 3) = 4 \wedge \text{eval } (\text{NAT } 4) = 7) \\ \wedge \text{eval } (\text{NAT } 3) = 2 \wedge \text{eval } (\text{NAT } 4) = 5$$

⁹ As a side note, this angelic execution-based method is agnostic to whether or not the underlying synthesis is top-down or bottom-up (see Section 5 of Miltner et al. (2022)). We explain BURST as a bottom-up synthesis tool just because it is how the BURST tool is currently implemented.

Table 5. Comparison of TRIO and SYRUP on the 43 benchmarks with random input-output examples. Each row represents the results for a different number of examples. “Succ. Rate” gives the success rate of each tool. “Avg Time” shows the average synthesis time for successful trials. “#T/O” denotes the number of time-outs

#Examples	TRIO			SYRUP		
	Succ. Rate(%)	Avg Time(s)	#T/O	Succ. Rate(%)	Avg Time(s)	#T/O
1	19.53	0.07	0	22.56	0.20	11
2	24.42	0.08	0	27.21	0.56	19
3	43.02	0.40	8	45.58	4.11	40
4	56.51	1.29	15	56.98	2.44	72
5	63.49	1.72	14	56.28	6.16	107
6	67.21	2.36	13	57.21	5.80	123
7	70.70	3.42	19	54.65	12.03	158
8	75.35	2.87	18	52.09	13.16	180

Again, after a bounded search, BURST fails to find a program that satisfies the strengthened specification. It concludes that the assumption made in the previous step is incorrect. Then, BURST refutes the assumption, increasing the list of anti-specifications. In this manner, BURST initially overapproximates the specification of the target function and then refines it by repeatedly adding anti-specifications. However, because the space of possible anti-specifications is too large, this method is not effective.

6.6 Sensitivity to the quantity and quality of examples

In this section, we compare TRIO with SYRUP (Yuan et al., 2023) in terms of the number of input-output examples (chosen randomly) required for synthesizing the desired programs. The goal of this experiment is to confirm how the quantity and quality of examples affects the performance of TRIO by comparing it with SYRUP which aims to synthesize generalizable programs from a few examples by leveraging a version space algebra.

Setup. We use the same benchmark set as the one used for evaluating SYRUP, which consists of 43 benchmarks. This set is a subset of the benchmarks used in our previous experiments with the 20 newly added benchmarks and two trivial benchmarks (bool_always_true and bool_always_false) excluded.¹⁰ For each benchmark, we generate 10 sets of random input-output examples with sizes ranging from 1 to 8. Every set includes the *base case* for the recursive programming task (i.e., the example(s) with the smallest input) because SYRUP is known to perform better with the base case according to the paper of SYRUP.

Table 5 summarizes the results. For each size of example set, we report the success rate, average synthesis time, and the number of timeouts for each tool. The success rate is defined to be the ratio of the number of successful synthesis trials (i.e., the number of

¹⁰ Extending the benchmark set to include the 20 newly added benchmarks is not easy because SYRUP requires all library functions to be first-order and monomorphic SMT functions that can be interpreted by Z3.

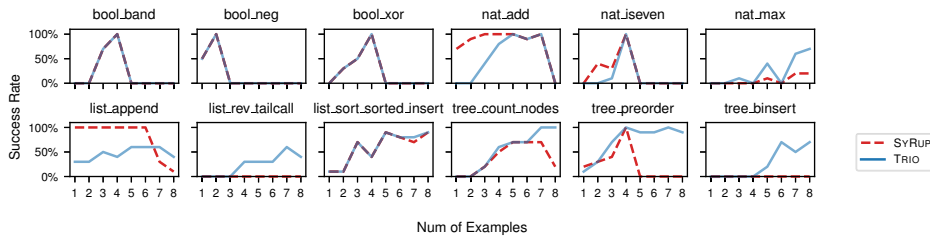


Fig. 9. Success rates of TRIO and SYRUP for 12 chosen benchmarks for different numbers of examples (1–8). The x-axis label indicates the number of examples, and the y-axis label indicates the success rate. The plots for the other 31 benchmarks are available in the appendix.

desired programs found) to the number of trials (i.e., the number of example sets multiplied by the number of benchmarks, which is 430 in this case). SYRUP performs better than TRIO in terms of success rate when the number of examples is small (1–4). In terms of efficiency, TRIO is faster than SYRUP as shown in the average synthesis time and the number of timeouts. However, when the number of examples increases (5–8), TRIO consistently outperforms SYRUP in both success rate and efficiency. We observe that SYRUP suffers from the scalability issue as the number of examples grows. SYRUP’s lower success rate is due to the fact that it often times out before finding the desired program when the number of examples is large. This is because SYRUP performs computationally expensive version space intersections, which become more expensive as the number of examples increases. This performance degradation can be observed in terms of memory usage as well. The average memory usage of SYRUP is 29MB when handling a single example and scales up to 131.6MB when handling 8 examples. In contrast, TRIO’s memory usage is much lower, starting at 16.4MB for 1 example and growing gradually to 22.1MB with 8 examples.

Figure 9 shows more detailed results for 12 chosen benchmarks. We present the success rate of each tool for each benchmark. The x-axis label in each plot indicates the number of examples, and the y-axis label indicates the success rate. There are three cases: (1) the two tools show similar success rates (*bool_band*, *bool_neg*, *bool_xor*, and *list_sort_sorted_insert*), (2) TRIO consistently outperforms SYRUP irrespective of the number of examples (*nat_max*, *list_rev_tailcall*, *tree_bininsert*, *tree_count_nodes*, and *tree_preorder*), and (3) SYRUP outperforms TRIO when the number of examples is small but SYRUP is comparable to or worse than TRIO when the number of examples is large (*nat_add*, *nat_iseven*, and *list_append*). The lower success rates of SYRUP is due to the version space intersection, which becomes more expensive as the number of examples increases. In contrast, TRIO’s performance remains stable or even improves as the number of examples grows because more examples can resolve the ambiguity in the search space. In conclusion, SYRUP cannot enjoy the benefits of more examples that can resolve the ambiguity in the search space because of the computational cost of version space intersection.

Summary of results. When given a small number of examples less than 5, SYRUP finds the desired programs more frequently than TRIO thanks to its version space algebra-based approach. However, SYRUP’s performance degrades as the number of

examples increases, leading to more timeouts and lower success rates. In contrast, TRIO's performance remains stable or even improves as the number of examples grows.

7 Related work

We divide the prior work related to our paper into three categories: (1) synthesis of functional recursive programs, (2) version-space-based synthesis, and (3) bidirectional search-based synthesis. We elaborate on these categories of work. For a broader survey of program synthesis, we refer the reader to Gulwani et al. (2017).

Synthesis of recursive programs. There is a large body of work on the synthesis of functional recursive programs. Various approaches have been proposed to synthesize functional recursive programs from input–output examples (Osera and Zdancewic, 2015; Feser et al., 2015; Lubin et al., 2020), refinement types (Polikarpova et al., 2016), logical specifications (Kneuss et al., 2013; Itzhaky et al., 2021), and a reference implementation with desired type invariants (Farzan and Nicolet, 2021). In the following, we will mainly focus on the prior work of inductive synthesis of functional recursive programs.

THESYS (Summers, 1986) and its reincarnation IGOR2 (Kitzelmann and Schmid, 2006) are similar to ours in the sense that they stage synthesis into (1) non-recursive program synthesis and (2) recursive program synthesis. They first synthesize non-recursive programs for the given example by a top-down search. Then, by identifying syntactic patterns, these systems “fold” the synthesized non-recursive programs into a recursive one. Similarly, CYPRESS (Itzhaky et al., 2021), which is for synthesizing recursive programs from separation logic specifications, also generates a satisfying straight-line program, then folds it into a generalized recursive one. Contrary to these systems, instead of exploring the space of possible foldings, which is prohibitively large in our case, we prune the search space of recursive programs by “unfolding” each candidate into a non-recursive program; we check if it can be one of the non-recursive programs synthesized earlier.

The *recursion-free approximation* in SYNDUCE (Farzan and Nicolet, 2021) is related to our block-based pruning. SYNDUCE is a system for synthesizing a recursive program from a reference implementation and type invariants. It also synthesizes recursive programs from non-recursive programs. It eliminates recursion in a given specification by replacing each recursive call with a variable, synthesizes a satisfying non-recursive program, and then changes the variables back to their corresponding recursive calls. This method differs from ours in that we do not directly construct a recursive solution from a non-recursive one. Instead, we prune the search space of recursive programs using non-recursive ones.

MYTH (Osera and Zdancewic, 2015) and λ^2 (Feser et al., 2015) pioneered the idea of top-down deductive search for functional recursive programs, which hypothesizes the overall structure of a program and then tries to synthesize the subcomponents. The major shortcoming of MYTH is the requirement for trace-complete specifications that our system does not need. The major shortcoming of λ^2 is that it only applies deductive reasoning to a fixed set of primitive list and tree combinators such as `filter` and `map`. Our deductive reasoning is not limited to a certain set of operators but can be applied to any usable external operators thanks to the use of inverse maps.

SMYTH (Lubin et al., 2020) and BURST (Miltner et al., 2022) are recently proposed systems for recursive program synthesis that do not require trace-complete specifications. SMYTH explores the search space top-down and generates partial programs with holes. To alleviate the trace-completeness requirement, for each partial program, SMYTH performs partial evaluation to propagate example constraints over the entire program into holes in it. However, as already shown in Section 6.5, SMYTH occasionally runs into the problem of continuously refining infeasible candidates by updating constraints over holes. On the other hand, our tool can quickly identify infeasible candidates, significantly outperforming SMYTH as already shown in Section 6.5. BURST performs bottom-up synthesis with *angelic execution* as already explained in Section 6.5. BURST inherits scalability issues of the prior bottom-up strategies where a goal-directed search in top-down strategies is missing. In contrast, our method combines top-down and bottom-up synthesis to overcome the limitation of bottom-up synthesis. In addition, BURST may run into the problem of *extensive backtracking* as explained in Section 6.5, whereas we do not have such an issue since we explore the full search space of recursive programs without any refinement process.

CONTATA (Miltner et al., 2024) is a recently proposed extension of BURST. CONTATA aims to synthesize recursive functional programs from *relational specifications*. As relational specifications do not constrain input–output behavior of individual functions but rather the relationship between multiple functions represented by logical formulas, CONTATA tackles a more challenging problem than ours. The major difference between CONTATA and BURST is that CONTATA is free from the problem of extensive backtracking because it does not overapproximate the specification of the target function. CONTATA discards infeasible candidates by checking their consistency with respect to the relational specifications.

Eguchi et al. (2018) have proposed a technique for synthesizing both a functional program and recursive helper functions from refinement types. Their method infers specifications of recursive helper functions by trying with a number of predefined templates. Our work focuses on synthesizing the target function when library functions are given. We expect our work can be combined with their work to synthesize the target function with a mixture of known and unknown library functions.

PARA (Hong and Aiken, 2024) is a recently proposed system for synthesizing recursive functional programs from input–output examples. Instead of general recursive programs, PARA targets paramorphisms. A paramorphism is a generalization of catamorphism (like fold) that only provides the recursive result. It retains both of the recursive result and the original input at each recursive step. Based on the observation that a broad range of recursive functions in practice can be expressed as paramorphisms, PARA constrains the search space to paramorphisms. By leveraging the structure of paramorphisms and a stochastic search strategy, PARA has been shown to outperform prior work on recursive program synthesis. However, not all recursive functions can be expressed as paramorphisms. For example, the McCarthy 91 function is not definable by a single paramorphism. Our techniques can be used to synthesize general recursive functions in principle (if the termination checker can handle them).

Version space-based synthesis. To efficiently represent the set of all programs correct with respect to a given specification, the prior version space approaches to synthesis

use a space-efficient data structure. FLASHFILL (Gulwani, 2011) first used e-graphs like version space representations to efficiently represent the set of all correct programs and choose the best one among them. This method is generalized in the FLASHMETA (Polozov and Gulwani, 2015) framework and its instantiations (Le and Gulwani, 2014; Kini and Gulwani, 2015; Rolim et al., 2017) have shown successful applications of the version space approach to synthesis in various domains. These methods construct version spaces top-down as we do in our system. There have been also previous methods that construct version spaces in a bottom-up fashion. Finite tree automata (FTAs) have been used to represent version spaces of functional programs (Wang et al., 2017; Miltner et al., 2022). In particular, BURST (Miltner et al., 2022) uses FTAs to represent the version space of *recursive* functional programs.

The major difference between our method and these previous methods is that we do not use the version space representation directly for finding a solution. Instead, we construct the version space of non-recursive programs to prune the search space of recursive programs.

DREAMCODER (Ellis et al., 2021) also indirectly uses version space representations for synthesis. DREAMCODER stores a large number of possible refactorings to each training program into version space representations. Those refactorings expose common sub-expressions that correspond to library functions, which DREAMCODER can use for other synthesis tasks. In contrast to DREAMCODER, we construct and use version spaces within a single synthesis task rather than across different tasks.

SYRUP (Yuan et al., 2023) uses version space algebra to avoid overfitting when synthesizing recursive functional programs from input–output examples. It uses pairs of recursive programs and execution traces that capture chains of recursive calls in the program in order to prioritize generalizable programs. SYRUP has a different goal from our work: SYRUP aims to minimize the number of examples required to synthesize a desired recursive program while we aim to improve the efficiency of synthesis.

Combining top-down and bottom-up search for recursive program synthesis. The idea of combining top-down and bottom-up synthesis often appears in prior work on recursive program synthesis. λ^2 (Feser et al., 2015) enumerates open hypotheses (i.e., partial programs with holes) by a top-down deductive search and closed hypotheses by a bottom-up search. Such closed hypotheses are used to fill holes in open hypotheses. MYTH (Osera and Zdancewic, 2015) also enumerates expressions bottom-up up to a certain size, and uses them during a top-down deductive search.

Our work is different from these methods in that (1) we use bottom-up enumeration for collecting not only sub-expressions but also inverse maps that enable top-down propagation for arbitrary external operators and (2) we use a combination of top-down and bottom-up synthesis not only for finding a recursive solution but also for finding all non-recursive blocks.

8 Conclusion

We have presented a new technique for synthesizing recursive functional programs from input–output examples. Our approach differs from prior work in that we first synthesize satisfying blocks (straight-line programs) for each input–output example, and then we

prune the space of recursive programs by removing candidates that are inconsistent with the blocks. Additionally, we propose a technique we call library sampling, which accelerates deductive reasoning over a library by using sampled input–output behaviors of library functions. We have implemented our algorithm in a tool called TRIO. Our comparison against the state-of-the-art synthesizers shows that TRIO advances the state of the art of inductive synthesis of recursive functional programs.

Acknowledgments

We thank the reviewers for their insightful comments that helped us improve the paper. This work was supported by the National Research Foundation of Korea (NRF) grant funded by the Korea government (MSIT) (No. 2021R1A5A1021944) and Institute for Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korea government (MSIT) (Nos. 2022-0-00995, RS-2024-00341722).

Conflicts of interest

The authors report no conflict of interest.

Data availability statement

The artifact is available at Zenodo: <https://doi.org/10.5281/zenodo.15690878>.

References

- Albarghouthi, A., Gulwani, S. & Kincaid, Z. (2013) Recursive program synthesis. In Proceedings of the 25th International Conference on Computer Aided Verification, CAV 2013, Saint Petersburg, Russia, July 13–19, 2013, Proceedings. Springer-Verlag, Berlin, Heidelberg, pp. 934–950.
- Eguchi, S., Kobayashi, N. & Tsukada, T. (2018) Automated synthesis of functional programs with auxiliary functions. In Programming Languages and Systems – 16th Asian Symposium, APLAS 2018, Wellington, New Zealand, December 2–6, 2018, Proceedings. Cham: Springer International Publishing, Cham, pp. 223–241.
- Ellis, K., Wong, C., Nye, M., Sablé-Meyer, M., Morales, L., Hewitt, L., Cary, L., Solar-Lezama, A. & Tenenbaum, J. B. (2021) Dreamcoder: Bootstrapping inductive program synthesis with wake-sleep library learning. In Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation. New York, NY, USA. Association for Computing Machinery, pp. 835–850.
- Farzan, A. & Nicolet, V. (2021) Counterexample-guided partial bounding for recursive function synthesis. In Computer Aided Verification – 33rd International Conference, CAV 2021, Virtual Event, July 20–23, 2021, Proceedings, Part I. Springer, pp. 832–855.
- Feser, J. K., Chaudhuri, S. & Dillig, I. (2015) Synthesizing data structure transformations from input-output examples. In Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation. New York, NY, USA. Association for Computing Machinery, pp. 229–239.
- Frankle, J., Osera, P.-M., Walker, D. & Zdancewic, S. (2016) Example-directed synthesis: A type-theoretic interpretation. In Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. New York, NY, USA. Association for Computing Machinery, pp. 802–815.
- Gulwani, S. (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. New York, NY, USA. Association for Computing Machinery, pp. 317–330.

- Gulwani, S., Polozov, A. & Singh, R. (2017) *Program Synthesis*. Foundations and Trends® in Programming Languages 4, 1–2, 1–119. Now Publishers Inc., Boston.
- Hong, Q. & Aiken, A. (2024) Recursive program synthesis using paramorphisms. *Proc. ACM Program. Lang.* **8**(PLDI), 102–125.
- Itzhaky, S., Peleg, H., Polikarpova, N., Rowe, R. N. S. & Sergey, I. (2021) Cyclic program synthesis. In Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation. New York, NY, USA. Association for Computing Machinery, pp. 944–959.
- Kini, D. & Gulwani, S. (2015) Flashnormalize: Programming by examples for text normalization. In Proceedings of the 24th International Conference on Artificial Intelligence (IJCAI'15), Buenos Aires, Argentina. AAAI Press, Palo Alto, CA, USA, pp. 776–783.
- Kitzelmann, E. & Schmid, U. (2006) Inductive synthesis of functional programs: An explanation based generalization approach. *J. Mach. Learn. Res.* **7**(15), 429–454.
- Kneuss, E., Kuraj, I., Kuncak, V. & Suter, P. (2013) Synthesis modulo recursive functions. In Proceedings of the 2013 ACM SIGPLAN International Conference on Object Oriented Programming Systems Languages & Applications. New York, NY, USA. Association for Computing Machinery, pp. 407–426.
- Le, V. & Gulwani, S. (2014) Flashextract: A framework for data extraction by examples. In Proceedings of the 35th ACM SIGPLAN Conference on Programming Language Design and Implementation. New York, NY, USA. Association for Computing Machinery, pp. 542–553.
- Lee, W. (2021) Combining the top-down propagation and bottom-up enumeration for inductive program synthesis. *Proc. ACM Program. Lang.* **5**(POPL), 1–28.
- Lee, W. & Cho, H. (2023) Inductive synthesis of structurally recursive functional programs from non-recursive expressions. *Proc. ACM Program. Lang.* **7**(POPL), 2048–2078.
- Lubin, J. (2020) *Forging Smyth: The Implementation of Program Sketching with Live Bidirectional Evaluation*. Bachelor's thesis, Department of Computer Science, University of Chicago, Chicago, Illinois.
- Lubin, J., Collins, N., Omar, C. & Chugh, R. (2020) Program sketching with live bidirectional evaluation. *Proc. ACM Program. Lang.* **4**(ICFP), 1–29.
- Miltner, A., Nuñez, A. T., Brendel, A., Chaudhuri, S. & Dillig, I. (2022) Bottom-up synthesis of recursive functional programs using angelic execution. *Proc. ACM Program. Lang.* **6**(POPL), 1–29.
- Miltner, A., Wang, Z., Chaudhuri, S. & Dillig, I. (2024) Relational synthesis of recursive programs via constraint annotated tree automata. In Computer Aided Verification. Cham. Springer Nature Switzerland, pp. 41–63.
- Omar, C., Voysey, I., Chugh, R. & Hammer, M. A. (2019) Live functional programming with typed holes. *Proc. ACM Program. Lang.* **3**(POPL), 1–32.
- Osera, P.-M. (2015) *Program Synthesis with Types*. PhD Dissertation, Computer and Information Science, University of Pennsylvania, Philadelphia, Pennsylvania.
- Osera, P.-M. & Zdancewic, S. (2015) Type-and-example-directed program synthesis. *ACM SIGPLAN Not.* **50**(6), 619–630.
- Polikarpova, N., Kuraj, I. & Solar-Lezama, A. (2016) Program synthesis from polymorphic refinement types. In Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation. New York, NY, USA. Association for Computing Machinery, pp. 522–538.
- Polozov, O. & Gulwani, S. (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. New York, NY, USA. Association for Computing Machinery, pp. 107–126.
- Rolim, R., Soares, G., D'Antoni, L., Polozov, O., Gulwani, S., Gheyi, R., Suzuki, R. & Hartmann, B. (2017) Learning syntactic program transformations from examples. In Proceedings of the 39th International Conference on Software Engineering. IEEE Press, Piscataway, NJ, USA, pp. 404–415.

- Summers, P. D. (1986) A methodology for lisp program construction from examples. In *Readings in Artificial Intelligence and Software Engineering*, Rich, C. & Waters, R. C. (eds). Morgan Kaufmann, San Francisco, CA, USA, pp. 309–316. Available at: <https://www.sciencedirect.com/science/article/pii/B9780934613125500288>.
- Wang, X., Dillig, I. & Singh, R. (2017) Program synthesis using abstraction refinement. *Proc. ACM Program. Lang.* **2**(POPL), 1–30.
- Yuan, Y., Radhakrishna, A. & Samanta, R. (2023) Trace-guided inductive synthesis of recursive functional programs. *Proc. ACM Program. Lang.* **7**(PLDI), 860–883.

1 Proofs

We first prove the soundness of the overall synthesis algorithm in Section 4.1 and the Deduce procedure in Section 4.3. Then, we prove the soundness of the termination check in Section 4.6.

Theorem 2. *Algorithm 1 finds a solution to a given synthesis problem if it exists.*

Proof Suppose that a program P_{sol} of size k is a solution to the synthesis problem. The target component size n keeps increasing until P_{sol} is found because of line 27 in Algorithm 1. If n becomes k , P_{sol} , or a program of size k that is observationally equivalent to P_{sol} , will be included in **C** by the COMPONENTGENERATION procedure (line 3). By the D_COMPONENT rule in Figure 3, the solution will be included in the queue \mathcal{Q} (lines 19–21) and will be returned as a solution (line 14). ■

Theorem 6. *Without using the D_EXTCALL rule, the Deduce procedure is sound.*

Proof Suppose **C**, \mathcal{I} , and \Box_u are provided to the Deduce procedure as input, and there exists an expression e_u satisfying \Box_u .

We prove the theorem by contradiction. We will show that a contradiction occurs if there exists an open expression $e \in \text{Deduce}(\mathbf{C}, \mathcal{I}, \Box_u)$ such that there is a hole in e that cannot be satisfied by any expression. The Deduce procedure can generate an open expression only via the rules D_CTOR, D_DTOR, D_TUPLE, and D_MATCH (because we assume the D_EXTCALL rule is not used).

- Case 1: $e = \kappa(\Box_{u_1}, \dots, \Box_{u_k})$. We will show that there exists a destructor application satisfying a hole. Because e must have been generated by the D_CTOR rule, $\Box_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto \kappa(v_{1j}, \dots, v_{kj})\}$ and for $1 \leq m \leq k$, $\Box_{u_m} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto v_{mj}\}$. From the assumption there exists an expression e_u satisfying \Box_u ,

$$\forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_u \Rightarrow \kappa(v_{1j}, \dots, v_{kj}).$$

By the standard semantics of the language,

$$\forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash \kappa^{-1}(e_u) \Rightarrow (v_{1j}, \dots, v_{kj})$$

and

$$\forall 1 \leq j \leq n, 1 \leq m \leq k. \sigma[x \mapsto i_j] \vdash \kappa^{-1}(e_u).m \Rightarrow v_{mj}$$

which can be rewritten as

$$\forall 1 \leq m \leq k. \kappa^{-1}(e_u).m \models_{\sigma} \square_{u_m}.$$

This contradicts the assumption that there is a hole in e that cannot be satisfied by any expression.

- Case 2: $e = \kappa^{-1}(\square_{u'})$. We will show that there exists a constructor application satisfying a hole. Because e must have been generated by the D_DTOR rule, $\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}$ and for $\square_{u'} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto \kappa(o_j)\}$. From the assumption there exists an expression e_u satisfying \square_u ,

$$\forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_u \Rightarrow o_j.$$

By the standard semantics of the language,

$$\forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash \kappa(e_u) \Rightarrow \kappa(o_j)$$

which can be rewritten as

$$\kappa(e_u) \models_{\sigma} \square_{u'}.$$

This contradicts the assumption that there is a hole in e that cannot be satisfied by any expression.

- Case 3: $e = (\square_{u_1}, \dots, \square_{u_k})$. We will show that there exists a projection satisfying a hole. Because e must have been generated by the D_TUPLE rule, $\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto (v_{1j}, \dots, v_{kj})\}$ and for $1 \leq m \leq k$, $\square_{u_m} = \bigcup_{1 \leq j \leq n} \{i_j \mapsto v_{mj}\}$. From the assumption there exists an expression e_u satisfying \square_u ,

$$\forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_u \Rightarrow (v_{1j}, \dots, v_{kj}).$$

By the standard semantics of the language,

$$\forall 1 \leq j \leq n, 1 \leq m \leq k. \sigma[x \mapsto i_j] \vdash (e_u).m \Rightarrow v_{mj}$$

which can be rewritten as

$$\forall 1 \leq m \leq k. (e_u).m \models_{\sigma} \square_{u_m}.$$

This contradicts the assumption that there is a hole in e that cannot be satisfied by any expression.

- Case 4: $e = \text{match } e' \text{ with } \overline{\kappa_i _ \rightarrow \square_{u_i}}^k$. We will show that there exists an arbitrary expression satisfying a hole. Because e must have been generated by the D_MATCH rule, $\square_u = \bigcup_{1 \leq j \leq n} \{i_j \mapsto o_j\}$, $e' \in C$, $\forall 1 \leq m \leq k$. $\square_{u_m} = \bigcup_{j \in I_m} \{i_j \mapsto o_j\}$ where $I_m = \{j \mid 1 \leq j \leq n, \sigma[x \mapsto i_j] \vdash e \Rightarrow \kappa_m(_)\}$. From the assumption there exists an expression e_u satisfying \square_u ,

$$\forall 1 \leq j \leq n. \sigma[x \mapsto i_j] \vdash e_u \Rightarrow o_j.$$

Because $\forall 1 \leq m \leq k. I_m \subseteq \{j \mid 1 \leq j \leq n\}$,

$$\forall 1 \leq m \leq k, j \in I_m. \sigma[x \mapsto i_j] \vdash e_u \Rightarrow o_j$$

which can be rewritten as

$$\forall 1 \leq m \leq k. e_u \models_{\sigma} \square_{u_m}.$$

This contradicts the assumption that there is a hole in e that cannot be satisfied by any expression. ■

Now we prove the soundness of the termination check (the Terminate procedure in Algorithm 7a).

Let $<$ be the standard subterm relation which is well-founded. The following lemma relates the subterm relation $<$ to the partial order relation \sqsubset (defined in Figure 7b) used in the termination check.

Lemma 18. *If $e_1 \sqsubset e_2$, then $\sigma \vdash e_1 \Rightarrow v_1$ and $\sigma \vdash e_2 \Rightarrow v_2$ implies $v_1 < v_2$.*

Proof By the definition of \sqsubset , there are three cases:

1. $e_1 = \kappa^{-1}(e_2)$ for some constructor κ ,
2. $e_1 = e_2.n$ where $e_1 \sqsubset e_2$ or $e_1 = e_2$ for some $n \in \mathbb{N}$,
3. $e_1 = (e_{1,1}, \dots, e_{1,m})$ and $e_2 = (e_{2,1}, \dots, e_{2,m})$ for some $m \in \mathbb{N}$.

The first case is for the base case of the induction, and the other two cases are for the inductive step. We will show the lemma holds for each case.

(1) By the standard semantics of the language, if $\sigma \vdash e_2 \Rightarrow \kappa(v'_1, \dots, v'_k)$ for some values v'_1, \dots, v'_k , then $\sigma \vdash e_1 \Rightarrow (v'_1, \dots, v'_k)$. Therefore, $v_1 = (v'_1, \dots, v'_k)$ and $v_2 = \kappa(v'_1, \dots, v'_k)$, and $v_1 < v_2$ holds by the definition of the subterm relation.

(2) Suppose $\sigma \vdash e_1 \Rightarrow (v_{1,1}, \dots, v_{1,k})$ for some values $v_{1,1}, \dots, v_{1,k}$, and $\sigma \vdash e_2 \Rightarrow (v_{2,1}, \dots, v_{2,k})$ for some values $v_{2,1}, \dots, v_{2,k}$ (i.e., $v_1 = (v_{1,1}, \dots, v_{1,k})$ and $v_2 = (v_{2,1}, \dots, v_{2,k})$). By the standard semantics of the language, $\sigma \vdash e_1.n \Rightarrow v_{1,n}$ where $1 \leq n \leq k$. If $e_1 = e_2$, then $v_{1,n} = v_{2,n}$. By the definition of the subterm relation, $v_1 < v_2$ holds since v_1 is a component of v_2 . If $e_1 \sqsubset e_2$, then by the induction hypothesis, $(v_{1,1}, \dots, v_{1,k}) < (v_{2,1}, \dots, v_{2,k})$ holds. Because $v_{1,n} < (v_{1,1}, \dots, v_{1,k})$ by the definition of the subterm relation, $v_{1,n} < (v_{2,1}, \dots, v_{2,k})$ holds by the transitivity of the subterm relation.

(3) Because $e_1 \sqsubset e_2$, there exists $1 \leq i \leq m$ such that $e_{1,i} \sqsubset e_{2,i}$ and all other $e_{1,j}$ are equal to $e_{2,j}$ or $e_{1,j} \sqsubset e_{2,j}$ for $j \neq i$.

Let $v_{1,k}$ be the evaluation result of $e_{1,k}$ and $v_{2,k}$ be the evaluation result of $e_{2,k}$ for $1 \leq k \leq m$. In other words, $v_1 = (v_{1,1}, \dots, v_{1,m})$ and $v_2 = (v_{2,1}, \dots, v_{2,m})$.

By the induction hypothesis, $\sigma \vdash e_{1,i} \Rightarrow v_{1,i}$ and $\sigma \vdash e_{2,i} \Rightarrow v_{2,i}$ implies $v_{1,i} < v_{2,i}$. Since $v_{1,j} = v_{2,j}$ for $j \neq i$, $v_1 < v_2$ holds by the definition of the subterm relation. ■

The following two lemmas are used to prove the soundness of the termination check.

Lemma 19. *For a given program $P = \text{rec } f(x) = e_{\text{body}}$, if $e \sqsubset x$ is true for every recursive call $f e$ in P , then P is guaranteed to terminate on any input.*

Proof Proof by contradiction. Suppose P does not terminate on some input i . Then, an infinite sequence of recursive calls will be generated:

$$f e_1, f e_2, f e_3, \dots$$

where e_i are the argument expressions of the recursive calls. When the first recursive call $f\ e_1$ is called, $e_1 \sqsubset x$ by assumption. Let $\sigma[x \mapsto i] \vdash e_1 \Rightarrow v_1$ for some value v_1 . By Lemma 18, $v_1 < i$. When $f\ e_2$ is called, $e_2 \sqsubset x$ by assumption. Let $\sigma[x \mapsto v_1] \vdash e_2 \Rightarrow v_2$ for some value v_2 . By Lemma 18, $v_2 < v_1$. In this way, the infinite sequence of recursive calls will generate an infinite sequence of values

$$i > v_1 > v_2 > v_3 > \dots$$

However, since $<$ is a well-founded relation, no such infinite sequence of values can exist. Therefore, the assumption that P does not terminate on some input i is false. Therefore, P is guaranteed to terminate on any input. ■

Lemma 20. *For a given program $P = \text{rec } f(x) = e_{\text{body}}$ where the input of P is a tuple of length m , if $(e_k)_{k \in K} \sqsubset (x.i)_{i \in K}$ where K is a non-empty set of indices is true for every recursive call $f(e_1, \dots, e_m)$ in P , then P is guaranteed to terminate on any input.*

Proof Proof by contradiction. Suppose P does not terminate on some input (i_1, \dots, i_m) . Then, an infinite sequence of recursive calls will be generated:

$$f(e_{1,1}, \dots, e_{1,m}), f(e_{2,1}, \dots, e_{2,m}), f(e_{3,1}, \dots, e_{3,m}), \dots$$

where $e_{i,j}$ comprise the argument expressions of the recursive calls. When the first recursive call $f(e_{1,1}, \dots, e_{1,m})$ is called, $(e_{1,k})_{k \in K} \sqsubset (x.i)_{i \in K}$ by assumption. Let $\sigma[x \mapsto (i_1, \dots, i_m)] \vdash e_{1,j} \Rightarrow v_{1,j}$ for some values $v_{1,j}$ where $1 \leq j \leq m$. By Lemma 18, $v_{1,k} < i_k$ for $k \in K$.

When $f(e_{2,1}, \dots, e_{2,m})$ is called, $(e_{2,k})_{k \in K} \sqsubset (x.i)_{i \in K}$ by assumption. Let $\sigma[x \mapsto (v_{1,1}, \dots, v_{1,m})] \vdash e_{2,j} \Rightarrow v_{2,j}$ for some value $v_{2,j}$ where $1 \leq j \leq m$. By Lemma 18, $v_{2,k} < v_{1,k}$ for $k \in K$. In this way, the infinite sequence of recursive calls will generate infinite sequences of values

$$i_k > v_{1,k} > v_{2,k} > \dots$$

for every $k \in K$. However, since $<$ is a well-founded relation, no such infinite sequences of values can exist. Therefore, the assumption that P does not terminate on some input (i_1, \dots, i_m) is false. Therefore, P is guaranteed to terminate on any input. ■

The following theorem shows that the termination check is sound.

Theorem 17. *If Terminate accepts P , then P is guaranteed to terminate on any input.*

Proof If P does not contain any recursive calls, then $\text{Terminate}(P)$ is true by line 3 of Algorithm 7a.

Otherwise, if P contains recursive calls, $\text{Terminate}(P)$ is true if and only if for every recursive call $f\ e$ in P , $\text{Struct}(f\ e, \text{KeyArgs}(e_{\text{body}}))$ is true.

There are two cases for $\text{Struct}(f\ e, \text{KeyArgs}(e_{\text{body}}))$ to be true.

First, if e is not a tuple, $e \sqsubset x$ by line 14. By Lemma 19, P is guaranteed to terminate on any input.

Second, if e is a tuple, $e' \sqsubset x'$ is true where e' and x' are defined in lines 10 and 11. No matter what K is, K is not empty by line 9. By Lemma 20, P is guaranteed to terminate on any input.

Therefore, if Terminate accepts P , then P is guaranteed to terminate on any input. ■

2 Evaluation

In this section, we add the evaluation results omitted in the main paper due to space constraints. Tables 6 and 7 show the results for the 15 easy problems in the **IO** benchmark suite and the **Ref** benchmark suite, respectively. Figure 10 shows the full results of Figure 9.

Table 6. Results for the 15 easy problems in the **IO** benchmark suite, where “Time” gives synthesis time in seconds, and “Size” shows the size of the synthesized program (measured by number of AST nodes). Synthesis time of the fastest tool for each problem is highlighted in bold.

Benchmark	TRIO			BURST			SMYTH		
	Time(s)	Size	Correct	Time(s)	Size	Correct	Time(s)	Size	Correct
bool_always_false	0.09	4	✓	0.02	4	✓	0.02	4	✓
bool_always_true	0.01	4	✓	0.02	4	✓	0.02	4	✓
bool_bor	0.01	14	✓	0.02	14	✓	0.02	23	✓
bool_impl	0.01	13	✓	0.02	13	✓	0.02	23	✓
bool_neg	0.01	8	✓	0.02	8	✓	0.02	8	✓
bool_xor	0.02	20	✓	0.02	20	✓	0.02	27	✓
list_concat	0.04	19	✓	0.03	19	✓	0.03	17	✓
list_hd	0.02	10	✓	0.02	10	✓	0.02	9	✓
list_inc	0.03	18	✓	0.02	18	✓	0.03	12	✓
list_last	0.03	22	✓	0.02	22	✓	0.02	19	✓
list_length	0.02	13	✓	0.02	13	✓	0.02	12	✓
list_nth	0.04	33	✓	1.99	33	✓	0.07	34	✓
list_tl	0.06	9	✓	0.02	9	✓	0.02	9	✓
nat_add	0.11	22	✓	0.04	22	✓	0.02	27	✓
nat_pred	0.01	7	✓	0.02	7	✓	0.02	7	✓
# Solved (# correct)	15 (15)			15 (15)			15 (15)		
# Timeout	0			0			0		

Table 7. Results for the 15 easy problems in the **Ref** benchmark suite where “# Iter” shows the number of CEGIS iterations

Benchmark	TRIO			BURST			SMYTH		
	Time(s)	Size	# Iter	Time(s)	Size	# Iter	Time(s)	Size	# Iter
bool_always_false	0.02	4	1	0.01	4	0	0.02	4	0
bool_always_true	0.01	4	0	0.01	4	1	0.02	4	1
bool_bor	0.02	14	3	0.02	14	4	0.03	23	4
bool_impl	0.02	13	2	0.02	13	3	0.02	23	3
bool_neg	0.02	8	2	0.01	8	2	0.02	8	2
bool_xor	0.02	20	3	0.02	20	3	0.04	27	4
list_concat	1.06	19	5	1.19	19	4	1.14	17	4
list_hd	0.19	10	2	0.21	10	2	0.22	9	2
list_inc	0.42	18	4	0.49	18	3	0.63	12	2
list_last	0.38	22	4	0.47	22	3	0.43	19	5
list_length	0.43	13	4	0.35	13	4	0.36	12	3
list_nth	0.38	33	6	0.45	33	7	0.45	34	5
list_tl	0.18	9	2	0.21	9	2	0.23	9	2
nat_add	0.15	22	5	0.08	22	6	0.07	27	6
nat_pred	0.02	7	2	0.01	7	2	0.02	7	2
# Solved	15			15			15		
# Timeout	0			0			0		

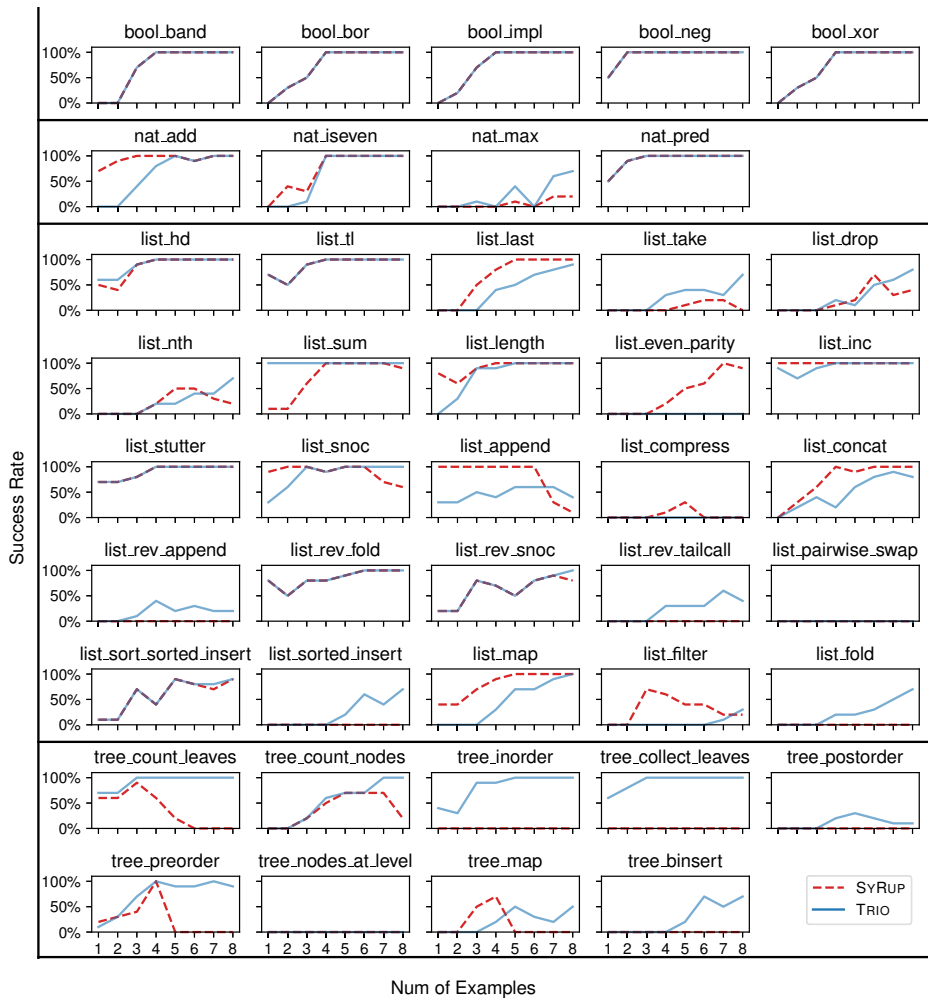


Fig. 10. Full results of Figure 9. The x-axis represents the number of examples, and the y-axis represents the success rate. The empty plot indicates that both tools failed to synthesize a program within the time limit.