

Bottom-Up Synthesis of Memory Mutations with Separation Logic

Kasra Ferdowsi ✉ 

University of California San Diego, CA, USA

Hila Peleg ✉ 

Technion, Haifa, Israel

Abstract

Programming-by-Example (PBE) is the paradigm of program synthesis specified via input-output pairs. It is commonly used because examples are easy to provide and collect from the environment. A popular optimization for enumerative synthesis with examples is Observational Equivalence (OE), which groups programs into equivalence classes according to their evaluation on example inputs. Current formulations of OE, however, are severely limited by the assumption that the synthesizer's target language contains only *pure* components with no side-effects, either enforcing this in their target language, or ignoring it, leading to an incorrect enumeration. This limits their ability to use realistic component sets.

We address this limitation by borrowing from Separation Logic, which can compositionally reason about heap mutations. We reformulate PBE using a restricted Separation Logic: Concrete Heap Separation Logic (CHSL), transforming the search for programs into a proof search in CHSL. This lets us perform bottom-up enumerative synthesis without the need for expert-provided annotations or domain-specific inferences, but with three key advantages: we (i) preserve *correctness* in the presence of memory-mutating operations, (ii) *compact* the search space by representing many concrete programs as one under CHSL, and (iii) perform a provably correct OE-reduction.

We present SOBEQ (*Side-effects in OBservational EQuivalence*), a bottom-up enumerative algorithm that, given a PBE task, searches for its CHSL derivation. The SOBEQ algorithm is proved correct with no purity assumptions: we show it is guaranteed to lose no solutions. We also evaluate our implementation of SOBEQ on benchmarks from the literature and online sources, and show that it produces high-quality results quickly.

2012 ACM Subject Classification Software and its engineering → Automatic programming

Keywords and phrases Program synthesis, observational equivalence

Digital Object Identifier 10.4230/LIPIcs.ECOOP.2025.10

Related Version *Full Version*: <https://zenodo.org/records/15271843>

Supplementary Material *Software (Source Code)*: <https://github.com/peleghila/sobeq>

Software (ECOOP 2025 Artifact Evaluation approved artifact):

<https://doi.org/10.4230/DARTS.11.2.2>

Funding This paper was published with the support of the Israel Science Foundation (ISF) #651/23.

Acknowledgements The authors thank Nadia Polikarpova for her help with formalisms, and Ilya Sergey and Yotam M. Y. Feldman for reading many drafts of this work and asking very good questions.

1 Introduction

Program Synthesis is the task of automatically generating a program that satisfies a given specification. A popular specification modality for synthesizing general-purpose programs is input-output examples, where a solution is defined as a program that, evaluated on each provided input, produces the provided output. This is commonly referred to as



© Kasra Ferdowsi and Hila Peleg;

licensed under Creative Commons License CC-BY 4.0

39th European Conference on Object-Oriented Programming (ECOOP 2025).

Editors: Jonathan Aldrich and Alexandra Silva; Article No. 10; pp. 10:1–10:32

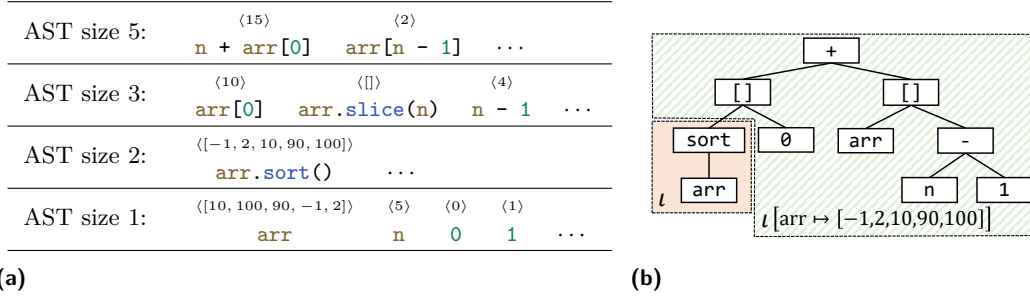
Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany



10:2 Bottom-Up Synthesis of Memory Mutations with Separation Logic



■ **Figure 1** (a) Classical OE bank of programs with their observed evaluation results used to enumerate pure programs. Size 2 contains `arr.sort()` that sorts in-place, making compositions with it incorrect. (b) AST of the target program, divided into the variable states in which each subexpression is evaluated.

Programming-by-Example (PBE) [23, 16, 7, 21, 26]. PBE is notable because, unlike logic-based specifications [49, 56, 29], it does not require specialized knowledge to use, and values for it can be collected from a programming environment. It has been used to develop synthesizers for both end-users [23, 60] and programmers [17, 16, 12].

Consider a synthesis specification for a JavaScript program where, given an array of integers `arr` of length `n`, we want to return the sum of the smallest and largest elements of the array. We might specify this problem with the example:

$$\{\text{arr} \mapsto [10, 100, 90, -1, 2], n \mapsto 5\} \rightarrow 99$$

If we assume the value in `arr` is disposable, e.g., the program may sort it in place, the following is a solution that we would like to synthesize:

```
arr.sort()[0] + arr[n-1]
```

To solve this synthesis task, we can select a set of *components* (operations, library functions, etc.) and search the space of all programs constructed from those components. There are many ways to search, including constructing a representation of the space that can be traversed [15, 39, 23] and encoding the problem for a solver [56, 30]. Such approaches place heavy constraints on which components can be included. But the simplest form of search is to enumerate the space of programs [13, 19, 25, 52] by applying the components, then test resulting programs. Because the space of programs is astronomical, this approach is prohibitively slow when implemented naively. Moreover, most mechanisms to prune the space are domain-specific and restricted by the manual effort for handling each component individually.

Bottom-up Enumeration. One technique that does not require any domain-specific effort is *bottom-up enumeration* with an *Observational Equivalence* (OE) reduction [59, 2]. OE uses components to combine smaller programs into increasingly larger ones. To prune its search space, it maintains a bank of programs that are representatives of the *equivalence classes* of observationally equivalent programs, i.e., programs that, on the given inputs, evaluate to the same outputs.

New programs are obtained by applying components from the component set to smaller sub-expressions already in the bank. For example, in Fig. 1(a) applying `+` to the representative of equivalence classes $\langle 5 \rangle$ and $\langle 10 \rangle$ will yield the program `n + arr[0]`, which will be banked as the representative of equivalence class $\langle 15 \rangle$. The bank lets the enumeration memoize the

evaluation results of each program, and use them when evaluating larger programs, rather than re-evaluating each full program from scratch. The space is pruned by keeping only one program per equivalence class. Despite discarding many programs, this ensures no solution that existed in the un-reduced space is lost. Previous works use variations of this technique in synthesizers for C [59], Python [17, 16], JavaScript [45], Java [22], Selenium [36], and OCaml [2, 40].

Enumeration with heap mutations. The standard assumption of OE is that all components are *pure* and cause no side-effects [40, 17, 60, 59, 2, 36], and any mutations of the state in the target program are performed outside the synthesis task. Other synthesizers make assumptions on the intermediate states [16, 47], reducing the problem to synthesis of pure programs.

This is because OE has a fundamental limitation when handling mutating programs. Notice that JavaScript’s `sort()` works in-place. In the example above, the call to `arr[n-1]` evaluates to a different value if evaluated before `arr.sort()` versus after it. A classical OE enumerator that performs its OE-reduction based on variable values in the inputs would consider `arr[n-1]` to be in the equivalence class $\langle 2 \rangle$. However, unlike its equivalent in languages like Python, `arr.sort()` also returns a self-reference, allowing the result to be composed into larger expressions, as in our target program. If the enumeration itself does not account for this, it will use the wrong memoized value computed using the unsorted value of `arr` in its candidates, yielding a result that does not actually satisfy the specifications. Moreover, the enumeration may see another program whose value is *really* $\langle 2 \rangle$, discard `arr[n-1]`, and never find this solution.

In other words, the classical bottom-up enumeration with OE described above is incorrect with memory mutations. This was also a problem in previous work: some synthesizers ignore the issue, leading to an incorrect enumeration [22, 45]. Peleg et al. [45], for example, allow JavaScript’s in-place `sort()`, breaking their own purity assumption.

The key question addressed by this work, then, is how to correctly enumerate mutating programs without losing solutions while maintaining enough pruning to make the process tractable.

The Problem with Mutations. Classical OE considers only one memory state, the state both before and after evaluation of any expression. A composition with memory mutations, however, must consider a potentially different *after-state*. Fig. 1(b) shows this: sub-programs after a mutation are evaluated in a modified *before-state*, the after-state of the mutation. Since composition by the synthesizer simulates evaluation, composition with a mutating expression requires that the synthesizer to be aware of modified before-states. The synthesizer will need to generate complex expressions, e.g., `arr[n-1]` in Fig. 1(b), in a modified state.

Naively, we could store alongside each program not only its evaluation results, as in classical OE, but also its before- and after-state. While this would correct the equivalence relation and ensure programs are not erroneously discarded, the enumeration would still need to enumerate programs over and over: the program `n-1` that yields the result 4 would be enumerated both when `arr` is in its initial state and its modified state, and, as we do not know what the sequence of mutations in the target program will be (if any), in many other states as well. This is not only costly, as it creates many copies of the program in different equivalence classes, it is unnecessarily wasteful, since `n-1` is independent of the value of `arr`. We need an alternative equivalence reduction that both composes *valid* evaluation sequences, and constructs a *compact* program space.

Our approach. We introduce SOBEQ (*Side-effects in OBservational EQuivalence*), to our knowledge the first bottom-up proof-directed synthesis technique. SOBEQ is a bottom-up enumerative algorithm for expressions proven correct in the presence of mutating components.

A key insight of SOBEQ is that mutating operations only mutate a small portion of the state. If the variable state, usually considered as a variable valuation or a stack, is considered as a heap, then *Separation Logic*'s [51, 44] concept of *local reasoning* can be used to reason about programs in the synthesizer's space. The *Frame Rule* gives us the tools for a *compositional* application of operations, in much the same way that bottom-up synthesis already does for pure programs. Thanks to the Frame Rule, a compact representation of programs suffices, rather than considering the full variable valuation every time.

We define Concrete Heap Separation Logic (CHSL), a representation of the operational semantics of the language in the style of Separation Logic. For example, in CHSL we describe the mutating program `n++` evaluated on the example's input as the triple $\{n \mapsto 5\}n++\{n \mapsto 6; 5\}$ where the precondition $n \mapsto 5$ describes `n`'s initial value, and the postcondition $n \mapsto 6; 5$ describes its value after mutation and its result. Instead of searching for a program, SOBEQ searches for a CHSL derivation from inputs to outputs. Separation Logic has been used in program synthesis before [49, 29], including where information from the environment can be leveraged as a SL specification [20]. But previous work performed the proof search *deductively*, applying manually-crafted inference rules. The concrete heaps of CHSL allow SOBEQ to instead search the proof space bottom-up, composing larger programs from smaller ones, and relying only on program evaluation, like a classical OE search does. CHSL gives SOBEQ two key advantages over classical bottom-up synthesis.

Correct composition via logic. First, SOBEQ's bank stores CHSL triples rather than programs. To build larger programs, SOBEQ also applies each component to smaller triples in the bank, but unlike simple programs, triples can only be combined according to the rules of the logic. This means programs always compose in a way that forms a valid evaluation, and there is no risk of evaluating a program in the wrong state. Moreover, OE can be determined by simply ignoring the program and comparing the specifications surrounding each triple: two programs with the same pre- and postcondition are guaranteed to behave the same, and will also compose the same. This maintains SOBEQ's *correctness* in the presence of mutations.

Local reasoning. Second, thanks to the Frame Rule, CHSL can reason about the local effects of programs: triples do not need to store parts of the state that the program does not touch. For instance, $\{n \mapsto 5\}n++\{n \mapsto 6; 5\}$ does not specify anything about `arr`, because `n++` will behave the same way in all concrete states where the value of `n` is 5. This representation is *compact*, with each triple describing the behavior in many, potentially infinitely many, concrete states. This makes it much easier for SOBEQ to find programs that use variables after a mutation occurs, like our example, compared to the naive solution of enumerating using the full before- and after-states. This local reasoning also allows us to prioritize more general programs: ones that provide the same result and the same mutations for more concrete states. Discarding the less general programs lets us compact the space even further.

Our implementation. We implemented SOBEQ as a synthesizer for JavaScript programs that can create non-pure expressions, including sequences of expressions. We evaluated our synthesizer on 63 benchmarks curated from the literature, from STACKOVERFLOW and from competitive programming website LEETCODE. Compared to state of the art synthesis

$\text{arr} \mapsto [-1, 2, 10, 90, 100] * n \mapsto 5:$	$\{\text{arr} \mapsto [-1, 2, 10, 90, 100] * n \mapsto 5; 100\}$			
	$\text{arr}[n-1]$	\dots		
$\text{arr} \mapsto [10, 100, 90, -1, 2] * n \mapsto 5:$	$\{\text{arr} \mapsto [10, 100, 90, -1, 2] * n \mapsto 5; 2\}$	$\{\text{arr} \mapsto [10, 100, 90, -1, 2] * n \mapsto 5; -1\}$		
	$\text{arr}[n-1]$	$\text{arr}[n-(1+1)]$	\dots	
$\text{arr} \mapsto [-1, 2, 10, 90, 100]:$	$\{\text{arr} \mapsto [-1, 2, 10, 90, 100]; [-1, 2, 10, 90, 100]\}$	$\{\text{arr} \mapsto [-1, 2, 10, 90, 100]; -1\}$		
	arr	$\text{arr}[0]$	\dots	
$\text{arr} \mapsto [10, 100, 90, -1, 2]:$	$\{\text{arr} \mapsto [10, 100, 90, -1, 2]; [10, 100, 90, -1, 2]\}$	$\{\text{arr} \mapsto [10, 100, 90, -1, 2]; 10\}$		
	arr	$\text{arr}[0]$		
	$\{\text{arr} \mapsto [-1, 2, 10, 90, 100]; [-1, 2, 10, 90, 100]\}$	$\{\text{arr} \mapsto [-1, 2, 10, 90, 100]; -1\}$		
	$\text{arr.sort}()$	$\text{arr.sort}()[0]$	\dots	
$n \mapsto 5:$	$\{n \mapsto 5; 5\}$	$\{n \mapsto 5; 4\}$	$\{n \mapsto 6; 5\}$	
	n	$n - 1$	$n++$	\dots
$\text{emp}:$	$\{\text{emp}; 0\}$	$\{\text{emp}; 1\}$	$\{\text{emp}; -1\}$	
	0	1	-1	\dots

■ **Figure 2** The precondition bank used to enumerate the program in Fig. 1(b). Each program is additionally labeled with its postcondition comprising an assertion and its result.

of general mutations, SOBEQ returns more concise and higher quality solutions, with less overfitting and without performing spurious operations, on both tasks that require and do not require mutations. Importantly, SOBEQ is both deterministic, meaning it will more reliably integrate with a larger tool, and proved correct, meaning that it will never discard all solutions to a task while exploring the space.

Contributions. The main contributions of this paper are as follows:

- A formulation of the PBE problem with memory mutations as the search for a derivation in a restricted Separation Logic, CHSL.
- The SOBEQ algorithm, a bottom-up enumerative algorithm that searches for a derivation with a proven-correct OE-reduction.
- An implementation of SOBEQ in a JavaScript synthesizer, and an evaluation on 63 benchmarks curated from the literature, STACKOVERFLOW and LEETCODE, which shows that SOBEQ outperforms state of the art on a variety of synthesis tasks with and without mutations, and highlights the class of tasks for which a naive approach would fail.

2 Overview

In this section, we present SOBEQ through the example introduced in Sec. 1. First, we present SOBEQ’s notation for mutating programs, which is compact compared to the real space of programs. Next, we show how this helps us enumerate programs correctly. Finally, we show how ordering programs by generality helps us prune the space even further.

2.1 Mutating programs

We begin with a reminder of programs and their construction in Observational Equivalence [2, 59], then show how we modify them to correctly reason about side-effects.

Classical OE-reduction and the program bank. Bottom-up enumeration with an OE-reduction uses a component set \mathcal{C} comprising functions, operators, variables, and literals, to search for a program specified by a vector of input-output examples, $\vec{\mathcal{E}}$, where each example is of the form $\iota \rightarrow \omega$. For each program we compute its equivalence class’s label. In classical

OE, we use its evaluation result, $\llbracket p \rrbracket(\iota)$, for each ι in $\vec{\mathcal{E}}$. In our example, the expression `arr` has the *observed behavior* $[10, 100, 90, -1, 2]$ on the provided input, so its equivalence class is $\langle [10, 100, 90, -1, 2] \rangle$.

To enumerate the space of programs, we use a *bank* to store representatives of equivalence classes, as seen in Fig. 1(a). For each function or operator $f \in \mathcal{C}$ with arity k , we collect all k -tuples of programs already in the bank and use them as arguments to apply f . This is usually done with some notion of iterations, e.g., as in Fig. 1(a), by increasing number of AST nodes. When enumerating programs of AST size 5, then, the enumerator considers the component array dereference (`[]`), which has an arity of 2. Among pairs of programs collected from the bank with a total size of 4 will be the program `arr`, representative of the equivalence class $\langle [10, 100, 90, -1, 2] \rangle$, and the program `n - 1`, representative of $\langle 4 \rangle$. Applying array dereference yields `arr[n - 1]` with the observed behavior $\langle 2 \rangle$. Crucially, we do not have to evaluate the full program `arr[n - 1]` to get this value: we can use observed values $[10, 100, 90, -1, 2]$ and 4 and only compute the final step.

When the enumeration finds a program from an equivalence class with no representative in the bank – i.e., no program in the bank has the same label – it is added to the bank to be used when constructing larger programs. This is the case with `arr[n - 1]`, which is added to the bank at size 5. Alternatively, if a representative of the equivalence class already exists in the bank, we discard the new program. If a program’s observed values are equal to the provided outputs, the program is a solution. Since we only use discovered equivalence classes (i.e., programs) to construct larger programs, the effect of OE’s pruning compounds.

The problems with side-effects. We quickly run into issues when \mathcal{C} contains functions or operators that can cause side-effects. The next step in solving our example is to enumerate `arr.sort()[0] + arr[n - 1]`. But, as we are about to see, using the classical approach this program will not be evaluated correctly, and it may even be pruned away.

We already saw how the subprogram $p_1 = \text{arr}[n - 1]$ was enumerated and labeled $\langle 2 \rangle$, and $p_2 = \text{arr.sort}()[0]$ will be enumerated similarly: at size 4, it will be composed from `arr.sort()` (size 2) and `0` (size 1), and added to the bank as a representative of $\langle -1 \rangle$.

The first source of trouble is that OE is a dynamic programming algorithm: when constructing the larger $p_2 + p_1$, the enumerator uses the memoized observed values to compute: $\langle -1 + 2 \rangle = \langle 1 \rangle$. However, since p_2 modifies `arr`, $p_2 + p_1$ *actually* evaluates to 99 on the input ι . Worse, this erroneous value means the synthesizer will consider $p_2 + p_1$, our solution, equivalent to $p_1 + p_2$ so it may be discarded. While a seemingly simple solution is to just evaluate the program in full each time, ensuring labels are correct, this is not only expensive, depending on the operation, but more importantly, insufficient. If the synthesizer sees the constant value 2, e.g., by enumerating `1 + 1`, this would cause the larger – and seemingly equivalent – `arr[n - 1]` to be discarded. OE hinges on the enumeration only discarding a program when it already has a suitable replacement, but `arr[n - 1]` evaluated *after* `arr.sort()` would still have no equivalent program in the bank.

Rethinking the program representation. Our problem was not distinguishing between the *same* program when evaluated using *different* values of `arr`. We also saw in Sec. 1 the need for distinction between programs that return the same value, but their *effect* on variables is different, e.g., `n` and `n++` when evaluated using the same value of `n`.

We replace programs with triples in Concrete Heap Separation Logic (CHSL): the precondition records the variables before its evaluation, and the postcondition comprises an assertion recording the variables after evaluation and the result of the evaluation, like so:

$$\{\text{arr} \mapsto [10, 100, 90, -1, 2]\} \text{arr.sort}() \{\text{arr} \mapsto [-1, 2, 10, 90, 100]; [-1, 2, 10, 90, 100]\}$$

From here on, we name the original array value arr_{orig} and its sorted value arr_{sort} , and use them in our examples and figures for brevity. In this representation, we can tell apart these two triples:

$$\{\mathbf{arr} \mapsto arr_{orig}\} \mathbf{arr} \{ \mathbf{arr} \mapsto arr_{orig}; arr_{orig} \} \quad \{ \mathbf{arr} \mapsto arr_{sort}\} \mathbf{arr} \{ \mathbf{arr} \mapsto arr_{sort}; arr_{sort} \}$$

where the first is \mathbf{arr} evaluated on the original array, returning the original array value, and the second is evaluated on a state where the array has been sorted, returning the sorted value. We use the second triple to construct larger programs evaluated after $\mathbf{sort}()$.

Notice that since all values are concrete, CHSL does not need domain- or library-specific inference rules to connect the pre- and postcondition of the triple, but can rely entirely on evaluation.

Local reasoning. These pre- and postconditions do not encode fully-concrete states, they leave \mathbf{n} unconstrained. This makes our representation of the space more compact: the triple $\{\mathbf{arr} \mapsto arr_{orig}\} \mathbf{arr} . \mathbf{sort}() \{ \mathbf{arr} \mapsto arr_{sort}; arr_{sort} \}$ takes advantage of the fact that only constraining \mathbf{arr} is sufficient to reason about the behavior of this program. We only need to constrain the *footprint* of the program.

CHSL includes the notion of the separating conjunction $*$ to indicate separate parts of the heap (or, in our case, variable store) that can be reasoned about locally. Thanks to aliasing restrictions of the synthesizer's target language detailed in Sec. 3.1, any two variables can be separated with $*$. E.g., the assertion $\mathbf{arr} \mapsto arr_{orig} * \mathbf{n} \mapsto 5$ describes the concrete initial state. Once we replace the programs in our bank with CHSL triples, we can now further borrow from Separation Logic when considering evaluation sequences.

2.2 Enumeration and heaps

A bottom-up enumeration composes larger programs from smaller ones by selecting a sequence of arguments to apply functions or operators f to. In a classical enumeration, our bank contains programs, so if we wish to apply $+$ to two programs in it, this is only a matter of creating the new program $p_1 + p_2$. When executing $p_1 + p_2$ on some variable valuation, the runtime first evaluates p_1 , then p_2 , then their addition. The classical OE enumeration duplicates this by summing the evaluation result banked for p_1 to that banked for p_2 . Without mutations, this is correct for all orderings of arguments to $+$ as all the programs evaluate from the initial state and to the initial state, so their results can always be memoized in the bank.

Now that our bank comprises CHSL triples, however, the enumeration needs to apply f to triples, not programs. Here, we have to be more careful. Let us consider two programs in our space, $\mathbf{arr}[0]$ and $\mathbf{arr} . \mathbf{sort}()[0]$ enumerated with classical OE in Fig. 1(a). In our new bank, they are now the triples:

$$\begin{aligned} tr_1 &= \{\mathbf{arr} \mapsto arr_{orig}\} \mathbf{arr}[0] \{ \mathbf{arr} \mapsto arr_{orig}; 10 \} \\ tr_2 &= \{\mathbf{arr} \mapsto arr_{orig}\} \mathbf{arr} . \mathbf{sort}()[0] \{ \mathbf{arr} \mapsto arr_{sort}; -1 \} \end{aligned}$$

First, we consider what applying $+$ to a pair of triples means: applying it to (tr_1, tr_2) as arguments means trying to create $\mathbf{arr}[0] + \mathbf{arr} . \mathbf{sort}()[0]$, or applying $+$ to the program components of the triples. However, unlike the mutation-free case, not every pair of triples with those two programs can create this application. We expect a sequence of arguments to a component to describe a real evaluation of those arguments.

This is akin to creating a valid composition or sequence in Separation Logic, the postcondition of each argument must be equal to the precondition of the next one in the sequence; this means all mutations are accounted for. For (tr_1, tr_2) , this is not a problem: $\mathbf{arr} \mapsto arr_{orig}$ is

the after-state of tr_1 and the before-state of tr_2 , which means that their sequence mutates the state from the precondition of tr_1 to the postcondition of tr_2 , and they can be used in the application of any component, e.g., $+$, with the values memoized in the triple, so we add the two precomputed results $10 + -1$, and get:

$$\{\text{arr} \mapsto \text{arr}_{orig}\} \text{arr}[0] + \text{arr.sort}()[0] \{\text{arr} \mapsto \text{arr}_{sort}; 9\}$$

However, we cannot compose (tr_2, tr_1) in the same way: $\text{arr.sort}()$ mutated the state, changing the value of arr to arr_{sort} , so tr_1 , whose memoized result was computed where arr is arr_{orig} , cannot follow it. Requiring arguments to a component to form a valid sequence ensures that only valid executions are enumerated by the synthesizer. The enumerator will not try to compose (tr_2, tr_1) under any binary operator, but can – and will – select (tr_1, tr_2) as arguments to $+$. Once a new program is composed, the enumerator invokes the language interpreter to evaluate the composition from the values of its arguments, yielding a triple containing the result and compound effect of the arguments and operation. We denote this invocation the EVAL inference rule.

If we partition our bank of programs by the program’s precondition, as in Fig. 2, then the enumerator can always create correct sequences of arguments by construction, using a triple’s postcondition to fetch all suitable next arguments according to their precondition.

Discovering more preconditions. Let us consider the last subprogram of our target program: $\text{arr}[\text{n} - 1]$. It is going to be added with $\text{arr.sort}()[0]$, which means that for that composition to be allowed, it must match the state after arr is sorted. Specifically, we consider the first argument of the array dereference, which will be needed to enumerate it:

$$tr_3 = \{\text{arr} \mapsto \text{arr}_{sort}\} \text{arr} \{\text{arr} \mapsto \text{arr}_{sort}; \text{arr}_{sort}\}$$

Importantly, this triple will never be enumerated if the enumerator only considers triples starting at a precondition describing the initial state, but without it we will never find the target program.

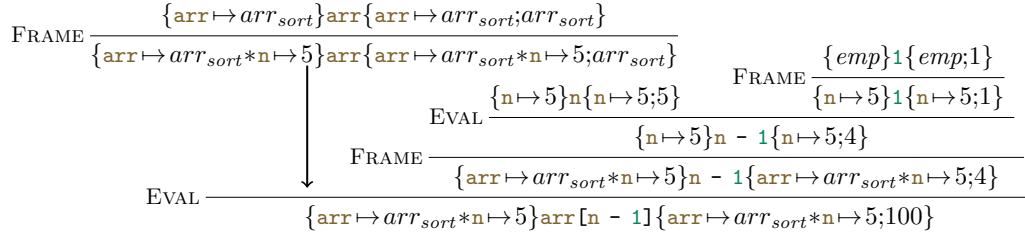
We have to ensure, therefore, that we also enumerate programs that have other preconditions. tr_3 is a necessary building block for such programs.

To this end, whenever a new assertion is discovered by the enumerator, e.g., when a mutating component is evaluated, creating a never-before-seen assertion in the postcondition, as enumerating $\{\text{arr} \mapsto \text{arr}_{orig}\} \text{arr.sort}() \{\text{arr} \mapsto \text{arr}_{sort}; \text{arr}_{sort}\}$ does, we get ready to enumerate programs that *begin* in this new assertion: we create tr_3 , which any program starting at $\text{arr} \mapsto \text{arr}_{sort}$ will need to incorporate, and add it to the bank, which also adds $\text{arr} \mapsto \text{arr}_{sort}$ to the bank as a precondition. The next time the enumerator selects arguments to an operator, tr_3 will be available. This lets the enumerator create programs that begin at $\text{arr} \mapsto \text{arr}_{sort}$.

Programs with different footprints. Now that we have tr_3 , the first argument to the array dereference, let us consider the second:

$$\{\text{n} \mapsto 5\} \text{n} - 1 \{\text{n} \mapsto 5; 4\}$$

We enumerate this program at the initial value of n , i.e., without the need for its precondition to be discovered. However, when the enumerator tries to compose it with tr_3 into a sequence of arguments to $[]$, it runs into a problem: the two triples do not make a valid sequence. The assertion in the postcondition of the first argument, $\text{arr} \mapsto \text{arr}_{sort}$, is different than the precondition of the following argument, $\text{n} \mapsto 5$.



■ **Figure 3** Using FRAME and EVAL to enumerate $\text{arr}[\mathbf{n} - 1]$. Recall arr_{sort} denotes the constant value $[-1, 2, 10, 90, 100]$ for brevity.

However, we clearly see that these two assertions deal with separate parts of the heap, and do not interfere with each other. This is where we draw on Separation Logic’s powerful tool: the FRAME rule. Since these are separate parts of the heap, we can use the separating conjunction to create a unified heap in the pre- and postcondition of both triples. The more constrained triples are now composable, so the enumerator can apply array dereference and EVAL, yielding:

$$\{\text{arr} \mapsto \text{arr}_{\text{sort}} * \mathbf{n} \mapsto 5\} \text{arr} [\mathbf{n} - 1] \{\text{arr} \mapsto \text{arr}_{\text{sort}} * \mathbf{n} \mapsto 5; 100\}$$

The full derivation is shown in Fig. 3. Then similarly, the same extension – now adding $\mathbf{n} \mapsto 5$ to tr_2 and nothing to the second argument – is applied to allow the enumerator to compose the target program. This is an important strength of our technique: the same triple for $\mathbf{n} - 1$ can compose with any value for arr – we will only enumerate it once!

We notice that in this case, this will discover a new assertion not as a postcondition as we did when `sort()` mutated `arr` but as a precondition. Once the program is added to the bank, the enumerator will continue to enumerate more programs with this new assertion as their precondition.

2.3 Equivalence and mutation

So far, we showed how the enumerator composes programs, but we did not show how the OE-reduction works. Now that we have a new representation of programs, we can re-define how each triple is observed and labeled, so that the label vector suits our needs.

In classical OE, the observed value for each example is the program’s evaluation result, but this is now insufficient: while in the precondition $\text{arr} \mapsto [10, 100, 90, -1, 2]$ the program $\text{arr}[1]$ evaluates to 100, if we discard all the following programs that evaluate to 100 we will also discard the $\text{arr}[\mathbf{n} - 1]$ that is evaluated on a sorted array. Likewise, \mathbf{n} and $\mathbf{n}++$ with the same precondition are distinct in that they return the same value given the same precondition, but have a different effect on the state.

Generally, we can see that we can label a triple with its pre- and postcondition, including the result. Once those are identical, any programs that form a valid triple with them would be *interchangeable*, i.e., can be swapped for each other with no change to a larger evaluation. In Fig. 2, the precondition component is expressed by the precondition the triple is stored under in the bank, while the postcondition appears above each program.

Generality of equivalence classes. The equivalence class labels above define an *equivalence relation* for OE. However, in this setting, we can do even better: there are cases where we want to discard programs even though one is not equivalent to the other. Let us consider two triples:

$$tr_4 = \{\mathbf{n} \mapsto 5\} \mathbf{n} - \mathbf{n} \{ \mathbf{n} \mapsto 5; 0 \} \quad tr_5 = \{emp\} 0 \{emp; 0\}$$

where *emp* is the unconstrained heap.

Both triples have the same evaluation result (0), and the same effect (no effect), but different assertions at the pre- and postconditions. But they are not completely different: if we select a *frame axiom* $R = \mathbf{n} \mapsto 5$, we can bridge their difference: $\{emp * R\} 0 \{emp * R; 0\}$ is equivalent to tr_4 . Under this condition, we say that tr_5 is more general than tr_4 : it describes a program that behaves the same over more concrete states. If we let the more general equivalence class *subsume* less general ones, we further compact our representation of the space while remaining correct. While we cannot tractably apply this additional reduction to the entire space, we can at least discard less general programs once more general ones exist in the bank.

Temporary values and variables. This tactic cannot stand on its own: naively employing it means $tr_6 = \{emp\} 5 \{emp; 5\}$ subsumes $tr_7 = \{\mathbf{n} \mapsto 5\} \mathbf{n} \{ \mathbf{n} \mapsto 5; 5 \}$ in the same way. However, in Fig. 2 similar subsumptions do not happen: $\{emp\} -1 \{emp; -1\}$ does not subsume $\{\mathbf{arr} \mapsto arr_sort\} \mathbf{arr}[0] \{ \mathbf{arr} \mapsto arr_sort; -1 \}$ and $\{\mathbf{arr} \mapsto arr_orig * \mathbf{n} \mapsto 5\} \mathbf{arr}[\mathbf{n} - (1+1)] \{ \mathbf{arr} \mapsto arr_orig * \mathbf{n} \mapsto 5; -1 \}$, which are still in the bank. Why?

We notice that if these subsumptions are allowed, the programs will no longer be interchangeable, and some will break entirely. Consider the application of `++`: while $\{\mathbf{n} \mapsto 5\} \mathbf{n}++ \{ \mathbf{n} \mapsto 6; 5 \}$ uses tr_7 as its argument, `5++` using tr_6 does not compile. Likewise, replacing *arr* within `arr.sort()` with the array literal with *arr*'s value would produce a different mutation of the heap, indicating they should be separate programs. The difference, which must be encoded into our representation, is what is *reachable* from state variables, and if so, from *where*.

Our evaluation result therefore comprises two components: the *value* returned by evaluation, and the *location* of that value. tr_3 , then, returns $arr_sort @ arr$, and tr_7 returns $5 @ n$. This ensures the interchangeability within `++` or `sort`: the new program will not only compile, it will affect the same variable.

What, then, should be the location of a temporary variable, i.e., a variable that is not reachable from the stack? We must consider the location's purpose: to separate programs enough to preserve interchangeability when the heap is mutated. Two temporary values, on the other hand, are (inherently) interchangeable, so we do not want to over-separate them which would inflate the number of equivalence classes in the bank. We therefore give all temporary values the *same* location: \perp . This way, tr_4 and tr_5 both return $0 @ \perp$, preserving the subsumption, and tr_6 returns $5 @ \perp$ which separates it from tr_7 .

Fig. 2 omits this for simplicity, but does still separate $\{emp\} -1 \{emp; -1\}$ with result $-1 @ \perp$ from $\{\mathbf{arr} \mapsto arr_sort\} \mathbf{arr}[0] \{ \mathbf{arr} \mapsto arr_sort; -1 \}$ with result $-1 @ arr[0]$ and from $\{\mathbf{arr} \mapsto arr_orig * \mathbf{n} \mapsto 5\} \mathbf{arr}[\mathbf{n} - (1+1)] \{ \mathbf{arr} \mapsto arr_orig * \mathbf{n} \mapsto 5; -1 \}$ with result $-1 @ arr[3]$.

Defining the Solution. We are almost finished: like classical OE, we continue iteratively enumerating programs into the bank until a solution is found. Our final step is to reconsider the definition of a solution to a synthesis specification.

Input variables	Identifiers $x \in \mathcal{V}$	
Literal values	$v \in \mathcal{K}$	
Assertion	$a \in \mathcal{A} ::= emp \mid x \mapsto v * \mathcal{A}$	$\frac{\{\vec{P}_1\}_{p_1}\{\vec{P}_2;\vec{r}_1\} \quad \{\vec{P}_2\}_{p_2}\{\vec{P}_3;\vec{r}_2\} \quad \dots \quad \{\vec{P}_k\}_{p_k}\{\vec{P}_{k+1};\vec{r}_k\}}{c \in \mathcal{C}, \text{arity}(c) = k \quad \forall i. (c(r_1^i, \dots, r_k^i), P_{k+1}^i) \rightarrow (r^i, Q^i)} \text{EVAL}$
Assertions	$\vec{a}, \vec{P}, \vec{Q} \in \mathcal{A} \times \dots \times \mathcal{A}$	
Location	$l \in \mathcal{M} ::= \perp \mid x \mid x[n],$ n is an integer literal	
Result	$r \in \mathcal{R} ::= v @ l$	$\frac{\{\vec{P}\}_p\{\vec{Q};\vec{r}\}}{\{\vec{P} * \vec{R}\}_p\{\vec{Q} * \vec{R};\vec{r}\}} \text{FRAME}$
Results	$\vec{r} \in \mathcal{R} \times \dots \times \mathcal{R}$	
Postcondition	$\{\vec{Q};\vec{r}\}$	

■ **Figure 4** CHSL syntax and inference rules.

When enumerating a program in a non-mutating space with an OE-reduction using the examples' inputs, a solution is a program that is labeled by (i.e., that evaluates to) the provided outputs. In SOBEQ this is insufficient: the enumerator can find a triple returning the correct value in *some* precondition. To be correct when evaluated on the concrete example inputs, we require that a solution's precondition describe the concrete initial state ι . In our example, $\text{arr} \mapsto \text{arr}_{orig} * \mathbf{n} \mapsto 5$ trivially describes ι – it is equal to it – but so does $\mathbf{n} \mapsto 5$.

Thus, to be a solution, a triple must: 1) return the values of the example outputs, and 2) have a precondition that describes ι .

SOBEQ is, to our knowledge, the first proof-directed synthesis algorithm that works bottom-up. It is also the first bottom-up enumerative synthesizer for general mutating components with a correctness guarantee: if a solution exists in the space spanned by \mathcal{C} , SOBEQ will find an equivalent solution.

3 Enumerative Synthesis and Heap Mutations

In this section we define the basic elements with which we construct SOBEQ: programs and their evaluation, and the synthesis task.

3.1 Concrete Heap Separation Logic

In this section, we introduce Concrete Heap Separation Logic (CHSL), a notation for operational semantics in the style of Separation Logic. Its syntax and two inference rules are shown in Fig. 4.

The heap. In CHSL, we use a heap notation to reason about the concrete values assigned to variables from the input. We can reason locally about each individual variable because SOBEQ's target language does not allow assignments, and so as long as nothing is aliased at the input – we can enforce this when generating the synthesis task – each variable is separate from all others and can be separated from them with a separating conjunction, $*$.

CHSL assertions. CHSL reasons about concrete values assigned to variables. As such, a CHSL assertion comprises heaplets of the form $x \mapsto v$ where v is a concrete value.

This also means we can trivially transform a state (partial function) σ into an assertion describing an identical heap s.t. $x_1 \mapsto \sigma(x_1) * \dots * x_n \mapsto \sigma(x_n)$ for all $x_i \in \text{dom}(\sigma)$. We use the function interchangeably with its assertion equivalent, i.e., σ as a shorthand for the assertion above.

10:12 Bottom-Up Synthesis of Memory Mutations with Separation Logic

For a vector of states $\vec{\sigma}$, corresponding to the vector of examples $\vec{\mathcal{E}}$, CHSL likewise uses a vector of assertions. We therefore also define a vector separating conjunction s.t. $\vec{P} * \vec{R} = \langle P_i * R_i \rangle_i$.

Values with location. In the presence of mutations, we must also consider that some values are reachable from our heap. As shown in Sec. 2, this is important for components that modify their arguments (e.g. `arr.sort()`), especially components that return a self-reference and are therefore composable, as in `arr.sort().reverse()`, where both `sort` and `reverse` mutate `arr`.

To support this, a result $r \in \mathcal{R}$ comprises two components: a *value* and that value's *location*, $l \in \mathcal{M}$ that can be a variable, a variable at a *concrete* index, or \perp . We use \perp to indicate a value not reachable from any variable, i.e., a temporary value. We decompose the result as $r = v @ l$. As with variable values, the values in results are concrete.

Triples. A triple in CHSL comprises a precondition $\vec{P} \in \mathcal{A}^k$ (where $|\vec{\mathcal{E}}| = k$), the program p , and a postcondition. Our postconditions have two components: i) an assertion $\vec{Q} \in \mathcal{A}^k$, and ii) the evaluation results of p on \vec{P} . Vectors of assertions and results form a conjunction: $\{\vec{P}\}p\{\vec{Q}; \vec{r}\}$ means all triples $\{P^1\}p\{Q^1; r^1\}, \{P^2\}p\{Q^2; r^2\}, \dots, \{P^k\}p\{Q^k; r^k\}$ hold.

We call the subset of each assertion Q whose values are different from P the *effect* of the program on P . For example, in the triple $\{x \mapsto 1 * y \mapsto 2\}x++ + y\{x \mapsto 2 * y \mapsto 2; 3\}$ the effect is $x \mapsto 2$.

We consider $\{\vec{P}\}p\{\vec{Q}; \vec{r}\}$ to be a *valid* triple if it accurately tracks the evaluation of p on \vec{P} . In other words, assuming a (deterministic) small-step operational semantics for programs the behavior of which is defined by the interpreter of the language: $\rightarrow: (\mathbb{P}, \mathcal{A}) \rightarrow (\mathcal{R}, \mathcal{A})$, we say $\{\vec{P}\}p\{\vec{Q}; \vec{r}\}$ is valid if $\forall 1 \leq i \leq k. (p, P^i) \rightarrow (r^i, Q^i)$.

Because SOBEQ will be *enumerating* triples, it is crucial that it can construct valid triples.

The Eval rule. When constructing a new triple, the EVAL rule denotes the use of the interpreter to ensure the resulting triple is valid. Because an evaluation of an AST node will first evaluate its children in sequence, then use the results to evaluate the node's operation, EVAL requires a sequence of triples. As in the classical Hoare Logic SEQUENCE rule, EVAL requires matching midconditions. However, since we often want to form a sequence from triples where $Q_i \neq P_{i+1}$, we use the FRAME rule to find a larger composed heap where the midconditions do match.

The Frame rule. Generally, all triples in the space contain only their footprint in their pre- and postcondition. As the example in Fig. 3 shows, in this compact representation midconditions may not match. We recall that the frame rule allows us to reason about a larger heap by equally extending the pre- and postcondition with an additional assertion called the *frame axiom*.

Fortunately, in CHSL, frame axioms are very easy to find. To prepare a pair of arguments for EVAL, then, one of two things is true: either for every i , $Q_1^i \cup P_2^i$ is still a partial function (i.e., if the same variable appears in both, it maps to the same value) and R_1^i, R_2^i are found by their difference, or the two triples cannot form a valid sequence.

Combining Eval and Frame. When the SOBEQ enumeration will apply a function or operator f to k child triples, this will combine the two rules. Since not all valid sequences of possible arguments to f have matching midconditions, we apply the FRAME rule to the entire sequence:

given a set of $k = \text{arity}(f)$ triples $tr_i = \{\vec{P}_i\}p_i\{\vec{Q}_i; \vec{r}_i\}$, EVAL is applicable if there exist k frame axioms $\vec{R}_1, \dots, \vec{R}_k$ where $\vec{Q}_1 * \vec{R}_1 = \vec{P}_2 * \vec{R}_2, \vec{Q}_2 * \vec{R}_2 = \vec{P}_3 * \vec{R}_3, \dots, \vec{Q}_{k-1} * \vec{R}_{k-1} = \vec{P}_k * \vec{R}_k$. Some of these \vec{R}_i may be *emp*. Now that the midconditions of the triples match, EVAL can be applied, yielding a triple that evaluates from the new shared precondition $\vec{P}_1 * \vec{R}_1$, and computing the result and precondition for $f(p_1, \dots, p_k)$.

3.2 Problem definition

Since SOBEQ uses CHSL triples to solve a PBE problem, we must re-state the PBE problem for CHSL. In this section we define what it means for a CHSL triple to represent our target program, define a translation from a PBE task, and explain how additional specifications, constraining the program's effects on the state, can also be supported.

The Programming by Example (PBE) task is usually defined as follows: given a vector of input-output examples $\vec{\mathcal{E}} = \langle \iota_i \rightarrow \omega_i \rangle_i$, find a program p over a component set \mathcal{C} where for every input ι_i , evaluating p on ι_i , denoted $\llbracket p \rrbracket(\iota_i)$, is equal to the provided output ω_i . We first want to translate this into a goal of the shape $\{P\}_{-}\{Q; r\}$.

Examples in CHSL. We denote the vector of input states $\vec{I} = \langle \iota_i \mid \iota_i \rightarrow \omega_i \rangle_i$. The direct translation of $\vec{\mathcal{E}}$ to CHSL, then, is $\{\vec{P}\}_{-}\{_;\langle \omega_i @ _ \rangle_i\}$, where $\exists \vec{R}. \vec{P} * \vec{R} = \vec{I}$. In other words, the precondition needs to *describe* the initial states (it is a footprint, so it may not describe all input variables), any postcondition will satisfy our specification, and the correct values can be at any location.

Constraining effect. Users may want to constrain specific mutations as well as the value. We let the user pair with each example any constraints on the effect, e.g., the value of `n` must be 9, or `arr` must end in its initial value $[10, 100, 90, -1, 2]$. To do this, we let the user provide an *effect constraint* $q : \mathcal{V} \rightarrow \mathcal{K}$, constraining the target value for any variable, where if $x \notin \text{dom}(q)$, x is unconstrained.

The SOBEQ task. A PBE user does not need to know about CHSL. Much like the original $\vec{\mathcal{E}}$, they provide a vector $\vec{\Phi}$ where each $\varphi \in \vec{\Phi}$ is a pair of example and effect constraint: $\varphi = (\iota \rightarrow \omega, q)$. If they are unconcerned with effects, this is identical to providing $\vec{\mathcal{E}}$.

We then turn $\vec{\Phi}$ into a synthesis goal in CHSL:

► **Definition 1** (SOBEQ synthesis goal). *Given a specification $\vec{\Phi}$, our goal $\mathcal{G}_{\vec{\Phi}}$ is:*

$$\mathcal{G}_{\vec{\Phi}} = \{\vec{P}\}_{-}\{\langle q_i * _ \rangle_i; \langle \omega_i @ _ \rangle_i\}, \text{ where } \exists \vec{R}. \vec{P} * \vec{R} = \vec{I}$$

The SOBEQ task is to find a triple $\{\vec{\sigma}\}p\{\vec{\sigma}'; \vec{r}\}$ that matches $\mathcal{G}_{\vec{\Phi}}$.

3.3 The program space

Now that we have SOBEQ's specifications, we define the space in which it performs the search. We consider the component set \mathcal{C} as a union of three sets: $\mathcal{V} \cup \text{lits} \cup \mathcal{F}$, where $x \in \mathcal{V}$ are the variables in \vec{I} , $l \in \text{lits}$ are literals in the component set, and $f \in \mathcal{F}$ are the functions and operators that can be applied to subprograms ($\text{arity}(f) > 0$). The full program space of a synthesizer with a component set \mathcal{C} is then every program using the components in \mathcal{C} . We denote this program set $\langle \mathcal{C} \rangle$. This set is infinite and highly redundant, e.g., contains both `x + 1` and `1 + x + 0`, and in our mutating domain, both `arr.sort()` and `arr.sort().sort()`.

Classical Observational Equivalence. In a problem domain with no side-effects, observational equivalence unifies programs with the same evaluation results on the example inputs. This is expressed in the equivalence relation $\equiv_{\vec{\mathcal{E}}}$ that deems two programs p_1, p_2 (observationally) equivalent if $\forall \iota \rightarrow \omega \in \vec{\mathcal{E}}. \llbracket p_1 \rrbracket(\iota) = \llbracket p_2 \rrbracket(\iota)$. The OE-reduced space $\langle \mathcal{C} \rangle^{OE}$ is then defined using $\equiv_{\vec{\mathcal{E}}}$. This idea was first suggested by Udupa et al. [59] and Albarghouthi et al. [2] and is in wide use today.

A classical bottom-up enumeration (CBE), keeps a bank of programs, and the vector of observed results is used to label the program's equivalence class. Only one program with each vector is added to the bank, meaning any observationally equivalent programs are discarded.

A program space with mutations. In the presence of mutations, the space of possible programs no longer comprises only programs evaluated in the initial states \vec{I} . As we saw in Sec. 2, composition in this space requires subprograms that are evaluated in a modified state.

The full program space when \mathcal{F} includes mutations is constructed inductively, building larger valid triples from smaller ones via applications of EVAL, and the FRAME rule when necessary. We denote the possible values for each $x \in \mathcal{V}$ as \mathcal{K}_x . We define the full space of programs $\langle \mathcal{C} \rangle$:

► **Definition 2** (Full program space). *We define the full program space $\langle \mathcal{C} \rangle$ as $\bigcup_{i \geq 0} \langle \mathcal{C} \rangle_{(i)}$ where $\langle \mathcal{C} \rangle_{(i)}$ is defined inductively as follows:*

$$\begin{aligned} \langle \mathcal{C} \rangle_{(0)} &= \{ \{ \text{emp} \} l \{ \text{emp}; \overline{v @ \perp} \} \mid l \in \text{lits}, v \in \mathcal{K} \text{ is the value of } l \} \cup \\ &\quad \{ \{ \vec{\sigma} \} x \{ \vec{\sigma}; \langle v_i @ x \rangle_i \} \mid x \in \mathcal{V}, \vec{\sigma} = \langle x \mapsto v_i \rangle_i, \vec{v} = \mathcal{K}_x \times \dots \times \mathcal{K}_x \} \\ \langle \mathcal{C} \rangle_{(n)} &= \langle \mathcal{C} \rangle_{(n-1)} \cup \{ \{ \vec{P} \} f(p_1, \dots, p_k) \{ \vec{Q}; \vec{r} \} \mid f \in \mathcal{F}, \text{arity}(f) = k, \{ \vec{\sigma}_i \} p_i \{ \vec{\sigma}'_i; r_i \} \in \langle \mathcal{C} \rangle_{(n-1)}, \\ &\quad \exists \vec{R}_1, \dots, \vec{R}_k. \text{EVAL of } f \text{ on } \{ \vec{\sigma}_i * \vec{R}_i \} p_i \{ \vec{\sigma}'_i * \vec{R}_i; \vec{r}_i \} \text{ is } \{ \vec{P} \} f(p_1, \dots, p_k) \{ \vec{Q}; \vec{r} \} \} \end{aligned}$$

Finally, given $\langle \mathcal{C} \rangle$, we can define the OE-reduced space that SOBEQ will enumerate.

3.4 SOBEq: observational equivalence with side effects

As we showed in Sec. 2, classical OE's $\equiv_{\vec{\mathcal{E}}}$ is no longer useful for the program space in Def. 2. We need our equivalence relation to encompass every aspect of the goal φ (Def. 1).

The generalized formalization of OE [45] defined two properties for equivalence, *interchangeability*: two equivalent programs can be used interchangeably within a larger program, and *consistency*: two equivalent programs will either both satisfy or both not satisfy the specification.

To satisfy consistency we must include everything that participates in Def. 1. The *value* returned by the expression, and the postcondition, which includes any constrained effect. We also need the state on which the program will be evaluated to check whether it describes \vec{I} .

As we saw in Sec. 2, to satisfy interchangeability we need the *location* of the value as well: while the literal $[2, 3, 1]$ and a variable x whose valuation is $[2, 3, 1]$ return the same value, they are not interchangeable under a larger program like $?.\text{sort}()$, where they would lead to different postconditions. Moreover, to be able to swap a program in, its postcondition must be the same.

Combined, then, our equivalence relation simply compares the triple excluding the program:

► **Definition 3** (SOBEQ equivalence relation). *The equivalence relation for $\vec{\Phi} = \langle (\iota_i \rightarrow \omega_i, q_i) \rangle_i$ is*

$$\{ \vec{\sigma}_1 \} p_1 \{ \vec{\sigma}'_1; \vec{r}_1 \} \equiv_{\vec{\Phi}} \{ \vec{\sigma}_2 \} p_2 \{ \vec{\sigma}'_2; \vec{r}_2 \} \iff \vec{\sigma}_1 = \vec{\sigma}_2 \wedge \vec{\sigma}'_1 = \vec{\sigma}'_2 \wedge \vec{r}_1 = \vec{r}_2$$

Or in other words, equivalence does not look at the program, only at its specifications.

Subsumption of triples. Finally, we consider the generality ordering of triples and (implicitly) of their equivalence classes, which will let us find more general programs.

In Sec. 2, we considered the two triples $tr_4 = \{\mathbf{n} \mapsto 5\}_{\mathbf{n}} - \mathbf{n}\{\mathbf{n} \mapsto 5; 0@ \perp\}$ and $tr_5 = \{emp\}_{0@ \perp} \{emp; 0@ \perp\}$. While tr_4 and tr_5 are not observationally equivalent under Def. 3, as $emp \neq \mathbf{n} \mapsto 5$, tr_5 is *more general* than tr_4 : it produces the same result ($0@ \perp$) and the same effect (\emptyset) from more concrete states. If tr_5 can absorb tr_4 , we can arrive at a more compact representation of the space. To this end, we formally define this generality: the partial order of programs that differ by only a more general pre- and postcondition.

► **Definition 4** (Generality ordering of triples). *Given a specification $\vec{\Phi}$, we define*

$$\{\vec{P}_1\}_{p_1}\{\vec{Q}_1; \vec{r}_1\} \sqsubseteq_{\vec{\Phi}} \{\vec{P}_2\}_{p_2}\{\vec{Q}_2; \vec{r}_2\} \iff \exists \vec{R}. \{\vec{P}_1\}_{p_1}\{\vec{Q}_1; \vec{r}_1\} \equiv_{\vec{\Phi}} \{\vec{P}_2 * \vec{R}\}_{p_2}\{\vec{Q}_2 * \vec{R}; \vec{r}_2\}$$

I.e., the more general way to achieve the same result and effect is considered larger.

While $\sqsubseteq_{\vec{\Phi}}$ is clearly not an equivalence relation (the existence of a frame axiom is not symmetrical) so it cannot be used to partition the space into classes, we notice that replacing a triple with one that is larger according to $\sqsubseteq_{\vec{\Phi}}$ does preserve the two properties of an OE-relation, *interchangeability* and *consistency*, in one direction. Originally [45], *interchangeability* states that swapping two equivalent programs within a larger program will yield an equivalent larger program, and *consistency* states that two equivalent programs either both satisfy a given specification φ or neither does. Under \sqsubseteq , however, a slightly weaker – but still useful – version holds.

Interchangeability under $\sqsubseteq_{\vec{\Phi}}$ means that given $f \in \mathcal{F}$, $\text{arity}(f) = k$, k triples $tr_j = \{\vec{P}_j\}_{p_j}\{\vec{Q}_j; \vec{r}_j\}$, $1 \leq j \leq k$, and k assertions \vec{R}_j s.t. EVAL can be applied on f and $\{\vec{P}_j * \vec{R}_j\}_{p_j}\{\vec{Q}_j * \vec{R}_j; \vec{r}_j\}$ to infer $\{\vec{P}\}f(p_1, \dots, p_k)\{\vec{Q}; \vec{r}\}$, and a triple $tr'_i = \{\vec{P}'_i\}_{p'_i}\{\vec{Q}'_i; \vec{r}'_i\}$ s.t. $tr_i \sqsubseteq_{\vec{\Phi}} tr'_i$ for some $1 \leq i \leq k$, then:

1. There exists \vec{R}'_i s.t. EVAL can still be applied on f when replacing $\{\vec{P}_i * \vec{R}_i\}_{p_i}\{\vec{Q}_i * \vec{R}_i; \vec{r}_i\}$ with $\{\vec{P}'_i * \vec{R}'_i\}_{p'_i}\{\vec{Q}'_i * \vec{R}'_i; \vec{r}'_i\}$ to infer $\{\vec{P}'\}f(p_1, \dots, p'_i, \dots, p_k)\{\vec{Q}'; \vec{r}'\}$, and
2. $\{\vec{P}\}f(p_1, \dots, p_k)\{\vec{Q}; \vec{r}\} \sqsubseteq_{\vec{\Phi}} \{\vec{P}'\}f(p_1, \dots, p'_i, \dots, p_k)\{\vec{Q}'; \vec{r}'\}$

Consistency under $\sqsubseteq_{\vec{\Phi}}$ means that for goal \mathcal{G} , if $\{\vec{P}_1\}_{p_1}\{\vec{Q}_1; \vec{r}_1\} \sqsubseteq_{\vec{\Phi}} \{\vec{P}_2\}_{p_2}\{\vec{Q}_2; \vec{r}_2\}$, then if there exists \vec{R}_1 s.t. $\{\vec{P}_1 * \vec{R}_1\}_{p_1}\{\vec{Q}_1 * \vec{R}_1; \vec{r}_1\}$ matches \mathcal{G} , then there exists \vec{R}_2 s.t. $\{\vec{P}_2 * \vec{R}_2\}_{p_2}\{\vec{Q}_2 * \vec{R}_2; \vec{r}_2\}$ matches \mathcal{G} . It's important to notice that the \Leftrightarrow of the equivalence version of *consistency* is now \Rightarrow : if $\{\vec{P}_1\}_{p_1}\{\vec{Q}_1; \vec{r}_1\}$ is a solution, so is $\{\vec{P}_2\}_{p_2}\{\vec{Q}_2; \vec{r}_2\}$, but since $\{\vec{P}_2\}_{p_2}\{\vec{Q}_2; \vec{r}_2\}$ is more general \vec{P}_2 might describe \vec{I} even if \vec{P}_1 did not.

► **Lemma 5.** *Interchangeability and consistency are preserved under \sqsubseteq .*

The proof of Lemma 5 is found in the extended version of the paper.

With the program space $\langle \mathcal{C} \rangle$ and the SOBEQ equivalence relation and generality ordering, we are ready to search for a program, i.e., to construct the OE-reduced space via enumeration.

4 The SOBEQ Enumeration

In this section, we describe the construction of the OE-reduced space. We first describe the structure of the space, then the enumeration, and finally how the enumeration can find equivalence classes that, by virtue of their generality, subsume other equivalence classes.

Classical Bottom-Up Enumeration. A bottom-up enumeration consists of applying each $f \in \mathcal{F}$ to k -tuples of child programs from the bank of previously seen programs, where $\text{arity}(f) = k$. Each newly enumerated $f(p_1, \dots, p_k)$ is evaluated, and if an equivalence class labeled by its observed values is not yet represented in the bank, it is added to the bank. Enumeration is done according to some ordering, e.g., height, number of AST nodes, or probability score [7]. This is usually facilitated by the structure of the bank to allow for easier retrieval of relevant child-programs.

The precondition bank. In our setting, we want to bank discovered triples according to their precondition: a bank section \vec{a} contains all discovered programs whose precondition is \vec{a} .

There are two reasons for this: first, when composing a k -tuple of arguments to f , finding a sequence EVAL is enabled for can be done constructively, without considering all possible k -tuples. For a sequence of length $k - 1$ with postcondition assertion \vec{Q}_{k-1} , we can preemptively rule out assertion vectors \vec{a} with different values assigned to a variable appearing in both \vec{Q}_{k-1} and \vec{a} , and will never form a valid sequence. In Algorithm 1 we call this construction COLLECTCHILDREN.

Second, and even more importantly, this lets us track what assertions have been discovered by the enumeration. If the enumeration produces a triple $\{\vec{a}\}p\{\vec{a}'; \vec{r}\}$ where \vec{a} or \vec{a}' have not yet been seen by the enumeration, then a new bank will be initialized for \vec{a} (or \vec{a}' , resp.).

Initializing a new bank. A new assertion \vec{a} can be discovered by the enumeration in one of two ways, as a pre- or a postcondition. For example, the derivation in Fig. 3 combines $\{\mathbf{n} \mapsto 5\}\mathbf{n} - 1\{\mathbf{n} \mapsto 5; 4\}$ with $\{\mathbf{arr} \mapsto \mathbf{arr}_{\text{sort}}\}\mathbf{arr}\{\mathbf{arr} \mapsto \mathbf{arr}_{\text{sort}}; \mathbf{arr}_{\text{sort}}\}$ into $\{\mathbf{arr} \mapsto \mathbf{arr}_{\text{sort}} * \mathbf{n} \mapsto 5\}\mathbf{arr}[\mathbf{n}-1]\{\mathbf{arr} \mapsto \mathbf{arr}_{\text{sort}} * \mathbf{n} \mapsto 5; 100\}$. Both the compound pre- and postcondition might be new for the enumeration, in which case they do not have a section in the bank.

If \vec{a} is discovered as a precondition, we naturally add a section of the bank to bank the newly enumerated program. If \vec{a} was discovered as a postcondition, then it may be a result of an effect, which means some of its atomic heaplets may not be in the bank yet. In this case, we create a section for \vec{a} in the bank, but also add its atomic heaplets to the bank if they are new: we construct $\{a_x\}x\{a_x; r_x\}$ where $a_x = \langle x \mapsto a_i(x) \rangle_i$ and $r_x = \langle a_i(x) \rangle_i$ for each x in \vec{a} , and add it to the bank if a_x is not yet in the bank. In Algorithm 1, this is collectively denoted as INITBANK(\vec{a}).

At the beginning of the enumeration, the bank is trivially initialized with the empty assertion emp for literals and with $\langle x \mapsto \iota_i(x) \rangle_i$ for each variable $x \in \mathcal{V}$.

Enumerating terms. The main enumeration is not very different from the CBE enumeration: the enumerator selects $f \in \mathcal{F}$ with $\text{arity}(f) = k$, and composes as arguments for it all k -tuples of triples $\{\vec{P}_i\}p_i\{\vec{Q}_i; \vec{r}_i\}$ that EVAL can be applied to from triples already in the bank.

The application of EVAL then constructs a new triple $\{\vec{P}\}f(p_1, \dots, p_k)\{\vec{Q}; \vec{r}\}$ that can then be *observed*: we use its pre- and postcondition assertions and its results as the label for its equivalence class. This label is searched for in the bank for precondition \vec{P} . If the bank for \vec{P} has an equivalent program, $\{\vec{P}\}f(p_1, \dots, p_k)\{\vec{Q}; \vec{r}\}$ is discarded, and if no equivalent program was previously seen, $\{\vec{P}\}f(p_1, \dots, p_k)\{\vec{Q}; \vec{r}\}$ is added to the bank for \vec{P} .

We also have at our disposal $\sqsubseteq_{\vec{P}}$ (Def. 4), which means we reduce the space even further.

■ **Algorithm 1** The SOBEQ enumeration.

Input : Set of specifications $\vec{\Phi} = \langle (\iota_i \rightarrow \omega_i, q_i) \rangle_i$, variables \mathcal{V} specified in example inputs
Output : Program p composed of components \mathcal{C} s.t. p that satisfies $\vec{\Phi}$
Globals: Precondition bank $bank$, mapping from assertion to an OE bank

```

1 initialAsserts := {emp}  $\cup \{ \langle x \mapsto \iota_i(x) \rangle_i \mid x \in \mathcal{V} \}$ 
2 foreach  $\vec{a} \in initialAsserts$  do INITBANK( $\vec{a}$ );
3 repeat // until timeout
4   foreach  $f \in \mathcal{F}$ ,  $arity(f) = k$  do
5     foreach  $\{\vec{\sigma}_1 * \vec{R}_1\} p_1 \{\vec{\sigma}'_1 * \vec{R}_1; \vec{r}_1\} \cdots \{\vec{\sigma}_k * \vec{R}_k\} p_k \{\vec{\sigma}'_k * \vec{R}_k; \vec{r}_k\} \in$ 
      COLLECTCHILDREN( $bank, k$ ) do
6        $\{\vec{\sigma}\} p \{\vec{Q}; \vec{r}\} := \text{EVAL}(f, \{\vec{\sigma}_1 * \vec{R}_1\} p_1 \{\vec{\sigma}'_1 * \vec{R}_1; \vec{r}_1\} \cdots \{\vec{\sigma}_k * \vec{R}_k\} p_k \{\vec{\sigma}'_k * \vec{R}_k; \vec{r}_k\})$ 
7       if  $\exists \vec{R}. \{\vec{\sigma} * \vec{R}\} p \{\vec{Q} * \vec{R}; \vec{r}\}$  matches  $\mathcal{G}_{\vec{\Phi}}$  then return  $p$ ;
8       if  $\forall \vec{a} \in keys(bank). \exists \vec{R}. (\vec{a} * \vec{R} = \vec{\sigma} \wedge \vec{R} \subseteq \vec{Q}) \Rightarrow (\vec{a}, \vec{Q} \setminus \vec{R}, \vec{r}) \notin bank(\vec{a})$  then
9         if  $\vec{\sigma} \notin keys(bank)$  then INITBANK( $\vec{\sigma}$ );
10         $bank(\vec{\sigma}) += \{\vec{\sigma}\} p \{\vec{Q}; \vec{r}\}$ 
11      end
12    if  $\vec{Q} \notin keys(bank)$  then INITBANK( $\vec{Q}$ );
13  end
14 end
15 until timeout;
```

Using more general triples. We want to use the $\sqsubseteq_{\vec{\Phi}}$ relation to let more general triples subsume less general ones in the course of enumeration, even when they are not strictly equivalent. However, this can happen in one of two cases: (i) the less general triple being enumerated after the more general triple is in the bank, and (ii) a more general triple being enumerated later.

Handling (i) is simple enough: when enumerating $\{\vec{\sigma}\} p \{\vec{\sigma}'; \vec{r}\}$, even if there is no representative for the equivalence class of $\{\vec{\sigma}\} p \{\vec{\sigma}'; \vec{r}\}$, the enumerator checks whether some $\{\vec{\sigma}_2\} p_2 \{\vec{\sigma}'_2; \vec{r}_2\}$ s.t. $\{\vec{\sigma}\} p \{\vec{\sigma}'; \vec{r}\} \sqsubseteq_{\vec{\Phi}} \{\vec{\sigma}_2\} p_2 \{\vec{\sigma}'_2; \vec{r}_2\}$ already exists in the bank, and if it does, discards $\{\vec{\sigma}\} p \{\vec{\sigma}'; \vec{r}\}$. This is efficiently done by changing the assertions to more general ones and searching the bank.

Handling (ii), however, is a larger feat. *In theory*, once we enumerate the more generic triple $\{\vec{\sigma}\} p \{\vec{\sigma}'; \vec{r}\}$, we can replace all less-general triples $\{\vec{\sigma}_2\} p_2 \{\vec{\sigma}'_2; \vec{r}_2\}$ with $\{\vec{\sigma}\} p \{\vec{\sigma}'; \vec{r}\}$, propagating this change to any compositions that use p_2 . This would be prohibitively costly on its own before the implications of Lemma 5: this replacement can weaken the pre- and postconditions of compositions, changing their place in the bank. This cascade might even weaken the precondition a triple that did not satisfy $\vec{\Phi}$ enough to describe \vec{I} and become a solution. As such, we elect to only incorporate (i) into SOBEQ.

4.1 The enumeration algorithm

We can now compose the complete SOBEQ enumeration, seen in Algorithm 1.

Lines 1–2 of the algorithm initialize the enumeration by calling INITBANK, described above, to create the basic programs for *emp* and for each of the variables' initial values.

Line 4 selects the next component for application; components with an arity of 0 are handled exclusively by INITBANK. Next, on line 5 COLLECTCHILDREN constructs all k -tuples of children that EVAL is enabled for. Since it is a part of the construction process, COLLECTCHILDREN returns each k triples already updated with the frame axioms $\vec{R}_1, \dots, \vec{R}_k$ to a unified composed heap.

Line 6 then applies EVAL to construct the new triple. If the new triple is a solution to the task, it will be returned by line 7. If not, line 8 checks whether an equivalent or subsuming program already exists in the bank via its equivalence class label: if some existing key in the bank can be made equivalent to the current program via a frame axiom. For example, the program $\{\mathbf{n} \mapsto 5\}n - n\{\mathbf{n} \mapsto 5; 0@ \perp\}$ will be enumerated, and will then be tested twice: for $R = \text{emp}$, $(\mathbf{n} \mapsto 5, \mathbf{n} \mapsto 5, 0@ \perp)$ will be searched for in the bank, and also for $R = n \mapsto 5$, $(\text{emp}, \text{emp}, 0@ \perp)$ will be searched and found (for the program 0), so $\{\mathbf{n} \mapsto 5\}n - n\{\mathbf{n} \mapsto 5; 0@ \perp\}$ is discarded. A precondition \vec{a} in the bank is deemed relevant if there exists a frame axiom \vec{R} s.t. \vec{R} both matches \vec{a} to $\vec{\sigma}$ and all values in \vec{R} are unchanged in \vec{Q} (i.e., \vec{R} truly is a frame). Every \vec{a} in the bank lookup uses the relevant equivalence class label, with \vec{a} as its precondition and subtracting the frame axiom from its postcondition assertion. If no triple is found, the new triple is added to the bank. If either the pre- or postcondition assertion do not exist in the bank, INITBANK is called for them to generate triples for their variables.

SOBEQ as a base enumerator. Many synthesizers [16, 45, 36, 5] use a base OE enumeration of expressions to synthesize additional language constructs, e.g., conditionals, higher-order functions, and assignments. These techniques are parameterizable in their base enumerator, and so could leverage SOBEQ to expand their target language.

Correctness of SOBEQ. We show that Algorithm 1 is *correct*: that if a solution p exists in the unreduced space, SOBEQ’s enumeration will find a solution observationally equivalent to p or more general than p .

► **Theorem 6** (Correctness of SOBEQ). *If a valid solution $\{\vec{\sigma}\}p\{\vec{\sigma}'; \vec{r}\}$ to goal $\mathcal{G}_{\vec{\Phi}}$ exists in $\langle\mathcal{C}\rangle$, then there exists $\{\vec{\sigma}''\}p'\{\vec{\sigma}'''; \vec{r}'\}$ that matches $\mathcal{G}_{\vec{\Phi}}$ in the SOBEQ-reduced space, and it is reachable by SOBEQ’s enumeration.*

In the extended version of the paper we prove Theorem 6 in two steps: first we define the (unreduced) space reachable via enumeration and show that any solution will still be in that space, then we apply the reduction of $\equiv_{\vec{\Phi}}$ to the space, defining the *SOBEQ-reduced program space* $\langle\mathcal{C}\rangle^S$, and prove that any solution will have an observationally equivalent or more general under $\sqsubseteq_{\vec{\Phi}}$ solution in $\langle\mathcal{C}\rangle^S$.

5 Implementation

We implemented a SOBEQ synthesizer for JavaScript programs in 3,700 lines of Scala code. The set of functions and operators \mathcal{F} contains 54 components for integers, strings, arrays, and sets, including nine mutating operations: postfix increment for integers, the array functions `sort`, `reverse`, `splice`, `push`, `pop`, `shift`, and the set functions `add` and `delete`. JavaScript strings are immutable, so \mathcal{F} includes no mutating components for strings. \mathcal{F} also includes a sequence operator for any two expressions. We allow a task to provide additional string constants for *lits*, common in benchmarks from the literature [4, 7] or interaction models [17, 16].

Preserving immutability. Effect constraints in $\vec{\Phi}$ can specify that a variable has not changed in the postcondition of the solution. But effect constraints cannot distinguish between, e.g., `x` and `--(++x)` because we only reason about the pre- and postcondition, not the path between them. This is generally desirable as it marks `--(++x)` as redundant, but this will only happen when applying `--`. The user may want `x` to *always* be immutable.

This can be enforced by the enumeration. We add to our task specifications the ability to mark a variable as `immut`. Once the user specifies `x` as immutable, there is no reason to construct `++x` at all: a triple $\{\vec{\sigma}\}p\{\vec{\sigma}'; \vec{r}\}$ where $\langle\sigma_i(x)\rangle_i \neq \langle\sigma'_i(x)\rangle_i$ is discarded immediately, and $\vec{\sigma}'$ is not initialized in the bank, since any programs with the precondition $\vec{\sigma}'$ assume a mutation of `x`.

Goals without an output. Under CHSL it is possible to express a goal that has no specified output, only a specified effect. This is desirable when, e.g., specifying that a value is inserted into a list, but what the insertion returns is irrelevant. If a benchmark has no specified outputs, i.e., $\vec{\Phi} = \langle(\iota_i, q_i)\rangle_i$, we create the synthesis goal $\mathcal{G}_{\vec{\Phi}} = \{\vec{P}\}_-\{\langle q_i * _ \rangle_i; _ \}$. The lack of an output does not change anything in how programs are enumerated and the goal is matched in the same way.

6 Experimental Evaluation

We evaluate SOBEQ on our implementation. All our experiments used a server with Intel Xeon Gold 6338 2 GHz¹ and 128GB of RAM, with the JVM maximum heap size set to 110 GB.

Research Questions. We aim to answer the following research questions:

- RQ1** How does SOBEQ compare to the state-of-the-art in syntax-guided synthesis of heap mutations?
- RQ2** What is the overhead of SOBEQ compared to a classical OE enumeration?
- RQ3** How effective is SOBEQ at pruning the space of programs?
- RQ4** How does SOBEQ compare to enumerating with concrete states?

6.1 Benchmarks and baselines

We evaluated SOBEQ on a set of 63 benchmarks, curated from previous synthesizers [52, 7, 16, 22], programming exercises on LEETCODE.COM, and questions on STACKOVERFLOW.COM. We discarded benchmarks that use types other than strings, arrays and sets, i.e., that would be inherently unsupported by our implementation. As SOBEQ supports a straight-line sequence of expressions, we also discarded benchmarks from Shi et al. [52] and Ferdowsifard et al. [16] if their solution used assignments to intermediate variables, loops, or conditionals.

We split the benchmarks into two sets: \triangleright **May Mutate**: 45 benchmarks that have a non-mutating solution in \mathcal{C} . \triangleright **Must Mutate**: 18 benchmarks that require mutation: benchmarks with a non-empty effect constraint or cannot be solved without mutation in \mathcal{C} . Of these, five are benchmarks from *May* with a constrained effect, ruling out non-mutating solutions. For RQ2, we also create \triangleright **Pure**: a version of *May* where all variables are marked as `immut`.

On average, *May* benchmarks have 2.8 examples (min 1, max 4) and 1.4 variables (min 1, max 3). *Must* benchmarks have 2.7 examples (min 1, max 5), 2.3 variables (min 1, max 4), and 1.2 variables specified in the effect (min 0, max 2). Five *Must* benchmarks have no specified outputs. More information about our benchmarks is provided in the extended version of the paper.

¹ This processor has 32 cores/64 threads, but our implementation is single-threaded.

Selecting a baseline. We considered several synthesizers as baselines for SOBEQ, ruling out synthesizers built for a single example [22], synthesizers that mutate by assigning synthesized pure expressions [16, 17, 59], and tools hinging on Java’s rich type system [15, 61]. FRANGEL [52] is the current state of the art in synthesis of general programs with heap mutations: it is specified with input-output examples, it has the same ability as SOBEQ to declare a variable as immutable, and its stochastic search carries no additional requirements for specification other than the ability to evaluate – just as SOBEQ does. This makes it the most suitable baseline.

FrAngel as a baseline. FRANGEL targets Java collection implementations, which have many more methods than their JavaScript counterparts. This means FRANGEL and SOBEQ can yield solutions that are entirely incomparable. To ensure an apples-to-apples comparison, we created new Java classes that mimic JavaScript types in Java, recreating the same \mathcal{C} used by SOBEQ’s implementation. We then defined our benchmark set over these types.

In addition, in order to force FRANGEL to create a straight-line sequence of expression statements as SOBEQ does, we also removed from its enumeration loops, conditionals, and assignments. Notice that removing choices from its random sampling only improves FRANGEL’s odds of finding a solution; before this additional change, FRANGEL’s run times were considerably worse.

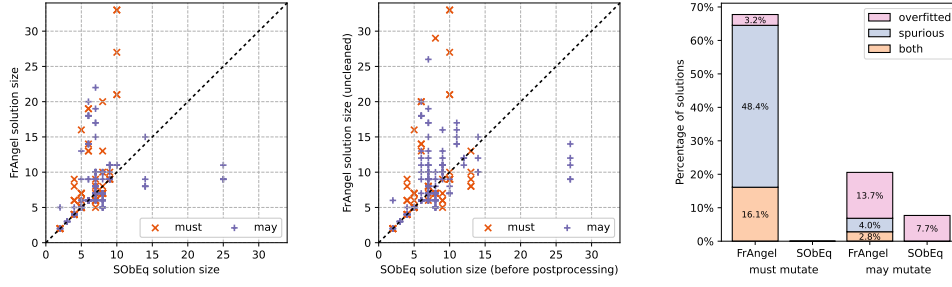
FRANGEL also has the option to set inputs as immutable, which let us encode `immut` variables, and an option to set a target value for an input, which we used for encoding effect constraints in the benchmark. This means FRANGEL can run on every benchmark in SOBEQ’s suite.

Selecting timeouts. We consider two timeouts for our runs. First, we set our timeout at 60 minutes, extending the original 30 minute timeout used to evaluate FRANGEL. However, since we aim SOBEQ to be used in interactive systems, and by other OE-based systems as the base enumerator, we also consider results at an *interactive timeout* or a timeout needed for the synthesizer to run in a user-facing tool and not create a *disruptive interruption* [43]. Building on recent interactive synthesizers, we choose 7 seconds as our interactive timeout.

6.2 RQ1: Comparison to FrAngel

We compare SOBEQ to state of the art FRANGEL, using an equivalent component set as detailed in Sec. 6.1. FRANGEL’s algorithm is stochastic, so we ran FRANGEL five times on each of the 63 benchmarks. We consider a benchmark *solved* if the synthesizer returns a solution that is correct on all examples. While some benchmarks from the literature contained a gold standard solution, we only used those as part of the analysis of result quality. We compare the quality of SOBEQ’s solutions to FRANGEL on three counts: the size of the solution, is it *overfitted*, and does it contain *spurious operations*. These are shown in Fig. 5. In addition, we compare the run times of the two tools at both an interactive 7s and a one-hour timeout.

Size of the solution. While the size of solutions is not, in itself, a good proxy for quality, it is telling in combination with overfitting and spurious operations. SOBEQ postprocessed solutions had an average of 6.2 AST nodes at an interactive timeout and 6.9 at a 1-hour timeout (max 14, 25, resp.) where *Must* benchmarks were generally smaller – 5.7 nodes at the interactive timeout and 6.1 in an hour – compared to *May* benchmarks – 6.4 and 7.1 nodes.



■ **Figure 5** RQ1: Comparing SOBEQ's solutions to FRANGEL's on size (left, center) and quality (right). Solution sizes are per-benchmark: each SOBEQ solution has up to five corresponding FRANGEL values. Above the line means SOBEQ is smaller. We consider both the final results (left) and solutions before postprocessing/cleaning steps (center). Solution quality (right) indicates the percentage of solutions with overfitted or spurious operations, determined by manual inspection (lower is better).

The left and center of Fig. 5 show the AST size for each SOBEQ solution compared to each of FRANGEL's solutions for the same benchmark. Overall, FRANGEL's solutions are much larger than SOBEQ's. There are two benchmarks where SOBEQ's solution is considerably larger than FRANGEL's: *Is All Positive?* (originally from FRANGEL), both SOBEQ's and FRANGEL's solutions were overfitted. In *Get first name from name with comma* (originally from PROBE [7]), FRANGEL returned two overfitted solutions similar to SOBEQ's, both appearing closer to the diagonal line, and three times found a solution closer to the gold standard solution, which is shorter, though never the gold standard solution itself. (The two identical solutions and three identical solutions overlap in Fig. 5.) An interesting trend in SOBEQ's solutions was that the enumeration order favored extracting postfix increment to a sequence. E.g., instead of `buffer[idx++]++`, `idx` SOBEQ returns the larger but functionally identical `buffer[idx]++`, `idx++`, `idx`.

Overfitted solutions. We define overfitted solutions as solutions that do not generalize beyond the provided examples. This was determined by manual inspection of both SOBEQ and FRANGEL solutions, and with reference to task descriptions and gold standard solutions, if they exist.

Two SOBEQ solutions were overfitted, and one additional benchmark was close to a gold standard solution for the task, but since all its inputs were of length 4, it replaced `str.length` in the gold standard solution with 4. Altogether, 3 of SOBEQ's 54 solutions were overfitted (6%). Of 279 FRANGEL solutions (up to five per benchmark, which sometimes differ) 47 were overfitted (17%). As the right of Fig. 5 shows, both tools were more prone to overfitting in *May* benchmarks.

We bring several examples of FRANGEL's overfitted solutions here.

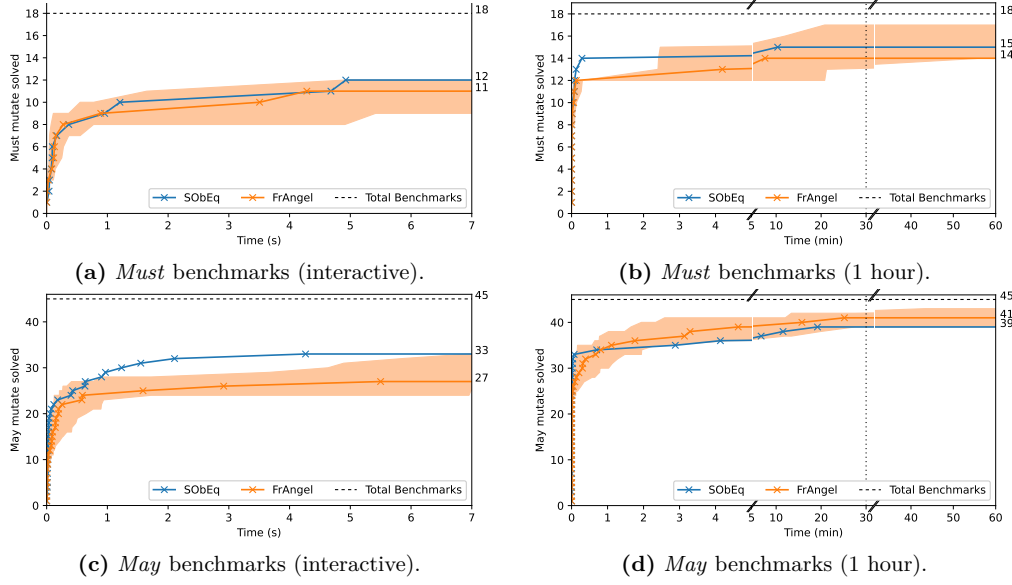
For example in the *Max and Min* benchmark, a simplified version of the example in Sec. 1, all four solutions (one run timed out) are overfitted and *different*. Two of them are:

```
return arr.join(arr.join(Str("")).replaceAll(Str(""), Str("_"))
.concat(arr.join(Str("")).indexOf(arr.slice(Int(0), Int(1)).join(Str(""))));
```

where the result number (max plus min) is generated by `indexOf`, and

```
arr.reverse();
return arr.join(arr.sort().join(Str("")).replaceAll(Str(""), Str("_")).replace(
arr.join(arr.join(Str("")).replaceAll(Str(""), Str("_")), Str("_"))
).length());
```

where it is generated by the length of a string.



■ **Figure 6** RQ1: The number of benchmarks SOBEQ and FRANGEL solve over time. FRANGEL is stochastic: median run for each benchmark shown by line with best and worst runs as range.

We also consider FRANGEL’s results on the three benchmarks that SOBEQ overfitted on: For the *Is All Positive?* benchmark mentioned above for its large solutions, both SOBEQ and FRANGEL overfit to the examples in two different ways. FRANGEL overfit all give solutions, returning `true` only for arrays where the smallest element is 1, and SOBEQ looked for a ‘-’ in an overfitted location in the printed array. For *get from name with comma*, four of five FRANGEL solutions were *not* overfitted. The fifth is:

```
return _arg_0.split(Str(",")).sort().slice(Int(1).subtract(_arg_1),_arg_1).join(Str(""));
```

where the input string is split into words but then sorted before extracting the 1-`_arg_1`th word. This happens to be true because all example inputs have a decreasing lexicographic ordering of the words. In *Negative Index Of*, where SOBEQ string lengths to 4, three of FRANGEL’s solutions are not overfitted, but two rely on string < comparisons to between variables and constants.

Spurious operations. We define spurious operations as terms in the AST that do not contribute to reaching the target values: e.g., sorting an already-sorted list, mutating a temporary value that is then discarded, etc. Of the 279 FRANGEL solutions, 37 had spurious operations (13%). We bring a few examples of these here.

One of the solutions FRANGEL found for the *buffer-add* benchmark was:

```
arr.splice(beg++, 1, toadd).reverse();
return end + 1;
```

Here, `splice` removes a sub-array from `arr`, replaces it with `toadd`, and returns it. The solution then calls `reverse()` on this temporary value, which is immediately discarded.

Likewise, in three of the five solutions to *rotate queue*:

```
mod.equals(inp.push(inp.shift()));
```

where input variable `mod` is compared to the result of the target program (shifting the first element out of the array and pushing it back in at the end).

For *set add-delete*, which asks for an in-place mutation of a the set `seen`, one solution is:

```
seen.delete(s2).compareTo(seen.add(s1).has(null));
```

where both `has` and `compareTo` are useless, computing a result that is not returned.

And, of course, this problem is not disparate from overfitting; the only successful run of FRANGEL for *fill fibonacci window* found this solution:

```
out.reverse().spliceReplace(
  Int(1),
  out.concat(out).length().subtract(
    Str(" ").concat(out.join(Str(" "))).lastIndexOf((new IntArray(Int(1))).join(Str(" ")))
  ).shift();
```

where `reverse().spliceReplace(1, _)` is a consistent, if odd, way to remove the first element and add something after the second element by swapping the elements and then replacing the (now) second one. The second argument to `splice`, which is expected to compute the next fibonacci element, overfits to the examples, and `shifting` on the discarded return value of `splice` is spurious.

SOBEQ, in contrast, had no spurious operations. Some are avoided thanks to OE: FRANGEL can draw programs that perform round trip operations like reversing multiple times (3 of 5 solutions to the example from Sec. 1 do this) or sort a sorted list. Enumerating with OE discards these, as they return the heap to a state it was already at and return the same value.

Other spurious operations are avoided because of the implicit bias of a bottom-up search toward small programs: the program `inp.push(inp.shift())` already satisfies the input, so a bottom-up search would not try to compose it into a larger program adding on spurious operations. Likewise, the enumeration would try to compose `arr.splice(beg++, 1, toadd)` to `end + 1` before trying it with the larger `arr.splice(beg++, 1, toadd).reverse()`.

Runtimes. We compared SOBEQ’s runtimes to FRANGEL’s best, worst and median performance on each benchmark in Fig. 6. In an interactive timeout, SOBEQ outperforms FRANGEL for the majority of the 7 seconds. At the 7s limit, SOBEQ solves 45 benchmarks, 7 benchmarks more than FRANGEL’s median run, whereas FRANGEL solves between 34 and 45 benchmarks. We also notice the main disadvantage of FRANGEL, which is the variance introduced by the stochastic algorithm: a difference of 30% between its best and worst cases. At the 1-hour timeout, SOBEQ solves 54 benchmarks, compared to FRANGEL’s 53–60 (median 55).

Fig. 6 separates benchmarks into sets to ensure performance is not at the expense of one set of benchmarks over the other, important for a general-purpose synthesizer. At an interactive timeout, SOBEQ performs on *Must* benchmarks as well as FRANGEL’s best run, both solving 12 of 18 benchmarks, and solves one more than FRANGEL’s median run of 11. For *May* benchmarks, SOBEQ solves 33 of 45, while FRANGEL solves between 25 and 33 (median 27). After an hour, SOBEQ outperforms FRANGEL’s median run for *Must* benchmarks, solving 15 benchmarks to FRANGEL’s 14 and is within FRANGEL’s range for both benchmark sets.

To answer RQ1, **SOBEQ returns higher quality solutions than FrAngel at comparable runtimes, without the nondeterministic variance in both time and expected solution.**

6.3 RQ2: Overhead of SOBEQ

A more general algorithm inevitably carries an overhead on a (now) special case compared to specialized algorithms for that case. This is often because of additional work by an algorithm making fewer assumptions, but in OE this can also take the form of a more refined equivalence relation, which unifies fewer programs creating a bigger space to explore. We wish to measure this overhead in order to quantify the cost of the changes made to support mutations in SOBEQ.

If \mathcal{C} has no mutating components, SOBEQ would not discover assertions outside the initial preconditions and their compositions. Such a search space is directly comparable to a classical bottom-up enumeration (CBE). This lets us measure the cost of the more

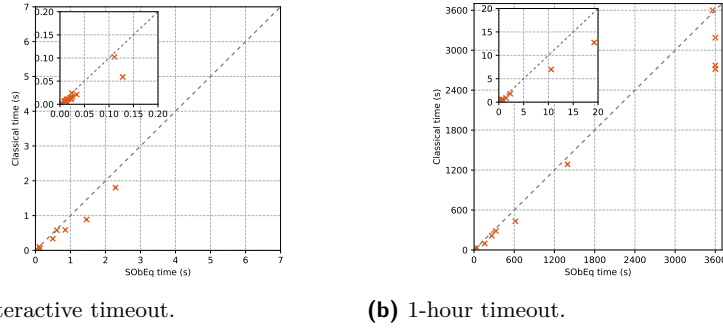


Figure 7 RQ2: Comparison with CBE: time to termination for each *Pure* benchmark with SOBEq compared to a classical bottom-up enumeration. Points at (timeout, timeout) not shown.

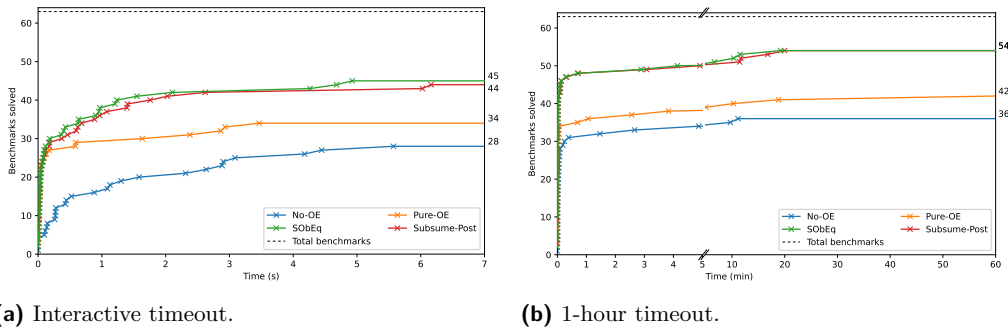


Figure 8 RQ3: Run times of disabling different pruning mechanisms within SOBEq: any OE, OE outside the initial state, and discarding programs according to \sqsubseteq .

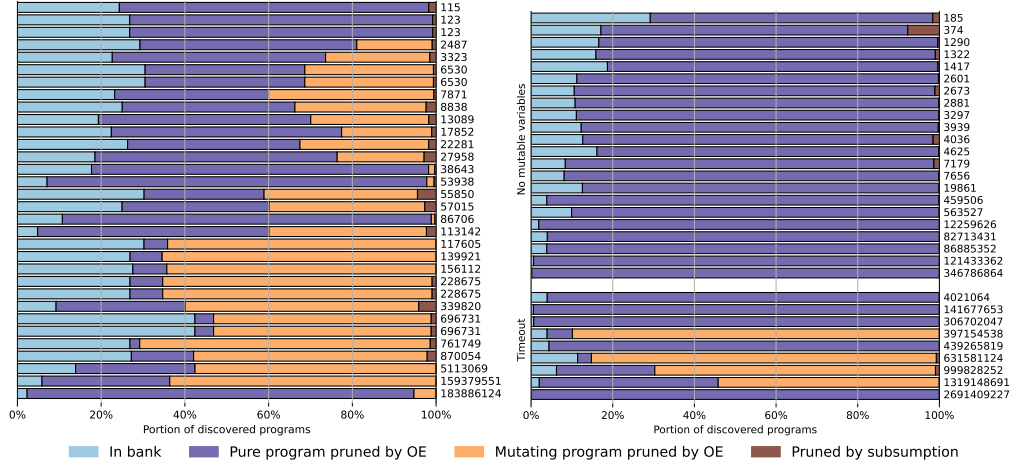
fine-grained equivalence classes caused by using results, not just values, in equivalence class labels, and the additional overhead of creating a compound heap for sub-programs using different variables.

To measure the overhead, we restricted our \mathcal{C} to only non-mutating components, and compared its performance on *Pure* benchmarks to an enumerator that uses only the initial *concrete* state and uses the classical value-based equivalence relation for an OE-reduction. This is equivalent classical OE with a height-based enumeration. By removing the mutating components for this evaluation, we avoid the correctness issues in classical OE, while controlling for performance differences caused by having different grammars.

In this restricted context CBE outperforms SOBEq for most benchmarks (time comparison shown in Fig. 7). In an interactive timeout, CBE solved the benchmarks an average of 1.52 times faster (median 1.5). CBE was also able to solve one additional benchmark, *Max and Min*, within the 7s timeout whereas SOBEq took 10.5s, and it is therefore not shown in Fig. 7. The biggest difference, where CBE was 2x faster or more (up to a maximum of 2.16 times faster), were all in benchmarks with runtimes of under 130ms for both enumerators.

In the full hour, CBE solved the benchmarks an average of 1.48 times faster (median 1.5), and finished with a two-benchmark lead: there are three benchmarks CBE solved that SOBEq did not, and one benchmark that SOBEq solved and CBE did not.

For RQ2, we find that the overhead of SOBEq’s more fine-grained equivalence classes and creating compound heaps is not negligible compared to CBE.



■ **Figure 9** Handling of programs seen in the space of each benchmark (percentages plotted, no. programs in text). Benchmarks where mutation cannot occur and timeouts plotted separately.

6.4 RQ3: SObEq’s effectiveness in pruning

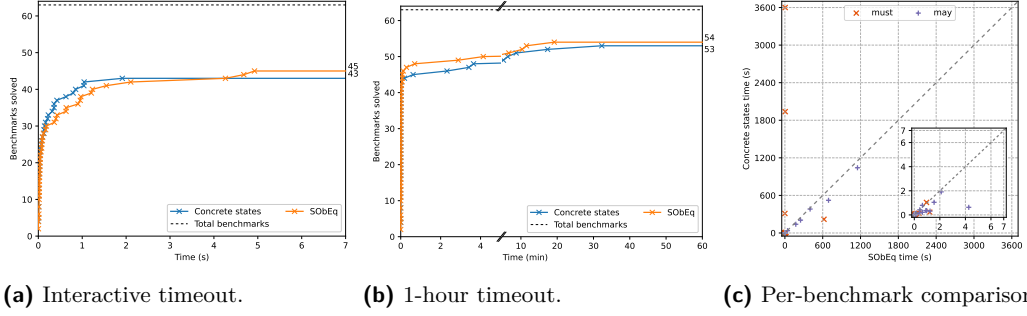
To test SObEq’s effectiveness at pruning its space we compare SObEq to two ablation enumerations: 1. *No-OE*: an enumeration that discards no programs, i.e., enumerates $\langle C \rangle^{reach}$. 2. *Pure-OE*: an enumeration that only performs OE on pure programs, i.e., fully enumerates the space outside preconditions that cover the initial state. Additionally, to measure the effect of subsumption, we also add 3. *Subsume-Post*: an enumeration with the SObEq OE-reduction, but no subsumption during enumeration, only as postprocessing of results. The results are shown in Fig. 8.

It is, of course, unsurprising that *No-OE* performs poorly: with no reduction of the space, few programs are feasible, and this version solves only 28 benchmarks within 7s, diverging completely after 10 minutes. *Pure-OE* is an interesting case: testing whether it should prune a program required computing its effect, an expensive operation that leads *Pure-OE* to be considerably slower. At the 1-hour limit, *Subsume-Post* does as well as SObEq. However, it lags a little behind SObEq for most of its run, having solved one fewer benchmark at 7s.

In summary, **both the application of OE to the entire space and the application of subsumption help SObEq’s performance, particularly in an interactive timeout.**

We also looked at what happens to the discovered programs across runs of SObEq, shown in Fig. 9, specifically to gauge how many programs are discarded and how many are subsumed by a more general program that already exists in the bank. We separate out runs that timeout and benchmarks where all input variables are strings, i.e., no mutation is *possible*. In the latter, the number of mutating programs pruned is trivially 0, but subsumption still takes place.

On (weighted) average, only 3% of programs seen by SObEq are added to the bank (med 14%, min 0.1%, max 42% of the programs). 66% on average are pure programs pruned by equivalence (med 72%, min 2%, max 99.98%), and another 30% are mutating programs pruned by equivalence (med 1.5%, min 0%, max 90%). Finally, an average 0.2% of the programs are pruned by subsumption (med 0.4%, min 0.00002%, max 8%). **While the ratios vary, each of the pruning methods in SObEq is significant for some of the benchmarks.**



■ **Figure 10** RQ4: Comparing an enumeration with preconditions to an enumeration with concrete states. Points at (timeout,timeout) not shown in (c).

6.5 RQ4: Enumerating with preconditions vs. concrete states

Finally, we examine the effect of the *local reasoning* in SOBEQ’s enumeration. Composing partial heaps takes effort, but how costly is it, and does it allow more benchmarks to be solved? To test this, we created a version of SOBEQ where instead of footprints, preconditions are always *full* concrete states. We compare this to SOBEQ’s enumeration in Fig. 10.

For most benchmarks, both enumerations perform similarly with concrete states slightly faster. We expected to see this as an effect of the overhead of composition, but it is also caused by pure programs being banked in one main state, which changes the enumeration order favorably for *May* benchmarks. If we consider only *Must* benchmarks, SOBEQ overtakes the concrete state enumeration after 1.5s and stays ahead for the entire 1-hour run.

Three benchmarks are noteworthy in this experiment. The concrete state enumeration times out on our overview example (benchmark *Max and Min with n*) while SOBEQ solves it in under 5s. Two more, *partition k largest* and *get range loopy while*, SOBEQ solves in 90ms and 4.6s, resp., and enumerating with concrete states takes 5min and 32min, resp. All three share the need to use multiple variables *after* one is mutated, which means subprograms already enumerated with the unmutated variable do not make a valid sequence and need to be re-enumerated.

To conclude, **SOBEQ’s local reasoning is faster for benchmarks that need to mutate, and is not significantly slower for other benchmarks.**

7 Related work

Syntax- and Component-Guided Synthesis. Syntax-Guided Synthesis (SyGuS) [3] is a form of the synthesis problem providing a specification along with a grammar defining its program space. It can also describe a form of synthesis where programs are constructed by syntax rules. The latter is essentially identical to Component-Based Synthesis, where the synthesizer is provided with a list of functions and operations spanning the program space. Both of these formulations are tackled by many synthesis strategies: enumerative [35, 52, 7, 14, 25, 22, 28, 54, 9, 62, 34], SMT [30, 6, 50, 1] and SAT [56, 58] based, by traversing constructed representations of the space such as VSAs [31, 33, 63], version spaces [53] and Petri-nets [15, 24].

Proof-directed synthesizers. A specific form of enumerative synthesis is one that attempts to find a proof for a proof goal, where the proof dictates a program. This idea is not new [57], and some deductive type-driven synthesizers describe their search as a proof

search [24, 42]. Other works' goals are more expressive than types: refinement types [48], relational specifications [32], select/update operators [41], Separation Logic, and even natural language [10]. Importantly, proof searches are overwhelmingly *deductive*, searching the proof space top-down. To our knowledge, SOBEQ is the first synthesis proof search that enumerates the proof space bottom-up.

Synthesis with Separation Logic. Separation Logic proofs are a popular vehicle for proof-directed synthesis. They have been successfully used for parallelizing non-parallel code [8], but are most often used to synthesize heap-manipulating programs. SuSLik [49] performs a deductive search with inferences about the heap to generate programs performing pointer operations. CYPRESS [29] extends this with the ability to synthesize recursive auxiliaries, ROBoSuSLik [11] extends it with read-only borrows, and RusSOL [20] applies it to Rust programs, leveraging Rust's type system to simplify the required specification from formal assertions to functional annotations of Rust functions. All of these accept assertions as their specification, hinge on manually-crafted inference rules, and employ solvers. SOBEQ, in contrast, works on input-output examples and relies only the generic EVAL rule in conjunction with an interpreter.

Observational equivalence. Observational Equivalence [59, 2] is a popular way to reduce the size of the space when enumerating bottom-up. Originally OE works on examples, but RESL [45] generalized it to other specifications on programs as well; this generalization also introduced the notions of *interchangeability* and *consistency* that SOBEQ's correctness relies on. SYMETRIC [18] relaxes OE to observational similarity by choosing representatives via clustering instead of equivalence, then repairing the resulting program. While this compacts the space like SOBEQ's subsumption, correctness of SYMETRIC's solutions hinges on the success of its repair step.

Many of the synthesizers that use OE only synthesize pure programs [46, 7, 45, 2, 60, 40] (we include [45] in this list, despite its misuse of JavaScript's `Array.sort`), others employ synthesis with OE to produce pure expressions that are then used in assignments or other larger mutating expressions [37, 38, 59, 17, 16], creating a top-level program with memory mutations.

ARBORIST [36] and LENS [47] relax the purity assumption in interesting ways. ARBORIST tackles the specific problem of dependent loops (e.g., *fold*) where the composition of smaller programs discovers new inputs and changes the equivalence classes. ARBORIST uses finite tree automata to unify programs without discarding them, which allows it to iteratively break up equivalence classes. ARBORIST finds a solution if its iterative refinement terminates. LENS take advantage of its target domain (straight line assembly code) to treat operations as having no arguments, making sequences the only composition. In this representation, OE values are full register valuations, and synthesis of state-mutating programs reduces to the OE enumeration used for pure programs. Both hinge on the ability to represent equivalence classes *without* discarding any programs using, implementing OE without a bank, which is only tractable under very strict assumptions on the set of components and on values.

Synthesis with heap mutations. The problem of synthesizing mutating snippets has been worked around in many ways. SL proof-directed approaches assume a limited target language of mutations available to the proof. SKETCH [56], TRANSIT [59], SNIPPY [17], SIMPL [55] and COZY [37, 38] assume expressions are pure, then use them in specific mutating contexts validated by other means (e.g., SMT). SyPET [15] and EdSYNTH [61] rely on types for getting programs, and only then evaluate them. Likewise, FRANGEL [52] draws random programs, then evaluates them.

Loopyr’s [16] problem statement bounds the intermediate mutations, and is to enumerate independent paths between these states. Similarly LENS [47] reduces all registers to 3-bit and constructs the full space of intermediate states. COBALT [41] also separates pure expressions from mutations using provided API annotations. TOSHOKAN [27] encodes the specifications for mutating methods as two versions of the function, one for changes to the state and the other the for result, akin to the two parts of SOBEQ’s postcondition.

CODEHINT [22] is, in some ways, the closest to our approach: it defines expressions as equivalent if they have the same value and side effects in the *current state*, or, in our terminology, an equivalence relation on both before- and after-states. How current states are used in enumeration is not defined in the paper, but their implementation uses programs enumerated on the initial state to populate their bank for new after-states, replicating the problem discussed in Sec. 2.1. Moreover, CODEHINT’s focus is its interactivity, and its interaction is designed around a single example, with any additional information only filtering an initial list of programs, not informing the OE enumeration. This is why we ruled it out as a baseline in Sec. 6.1.

SOBEQ sheds assumptions that expressions are pure and that intermediate states are known in advance. Moreover, SOBEQ has the formal guarantee in Theorem 6, unlike the two most general existing solutions: FRANGEL, which randomly draws programs, and CODEHINT, which performs OE only on its first example and mutated states in an unprincipled manner.

Other side effects like file or database writes have been a challenge in program synthesis. RBSYN [25] addresses database writes via a type-guided top-down enumeration from a syntax including *effect holes* that can be filled with annotated mutating methods; expressions not within effect holes have a purity assumption, limiting how effects are composed, unlike SOBEQ that can use mutating operations anywhere. CODEHINT allows users to enable effects on the file system in their enumeration, even though they cannot be undone like memory mutations are.

References

- 1 Alessandro Abate, Cristina David, Pascal Kesseli, Daniel Kroening, and Elizabeth Polgreen. Counterexample guided inductive synthesis modulo theories. In Hana Chockler and Georg Weissenbacher, editors, *Computer Aided Verification*, pages 270–288, Cham, 2018. Springer International Publishing. doi:10.1007/978-3-319-96145-3_15.
- 2 Aws Albarghouthi, Sumit Gulwani, and Zachary Kincaid. Recursive program synthesis. In *International conference on computer aided verification*, pages 934–950. Springer, 2013. doi:10.1007/978-3-642-39799-8_67.
- 3 Rajeev Alur, Rastislav Bodik, Garvit Juniwal, Milo MK Martin, Mukund Raghothaman, Sanjit A Seshia, Rishabh Singh, Armando Solar-Lezama, Emina Torlak, and Abhishek Udupa. Syntax-guided synthesis. *Dependable Software Systems Engineering*, 40:1–25, 2015. doi:10.3233/978-1-61499-495-4-1.
- 4 Rajeev Alur, Dana Fisman, Rishabh Singh, and Armando Solar-Lezama. Sygus-comp 2017: Results and analysis. *arXiv preprint*, 2017. arXiv:1711.11438.
- 5 Rajeev Alur, Arjun Radhakrishna, and Abhishek Udupa. Scaling enumerative program synthesis via divide and conquer. In *Tools and Algorithms for the Construction and Analysis of Systems: 23rd International Conference, TACAS 2017, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2017, Uppsala, Sweden, April 22-29, 2017, Proceedings, Part I 23*, pages 319–336. Springer, 2017. doi:10.1007/978-3-662-54577-5_18.

- 6 Haniel Barbosa, Andrew Reynolds, Daniel Larraz, and Cesare Tinelli. Extending enumerative function synthesis via SMT-driven classification. In *2019 Formal Methods in Computer Aided Design (FMCAD)*, pages 212–220. IEEE, 2019. doi:10.23919/FMCAD.2019.8894267.
- 7 Shraddha Barke, Hila Peleg, and Nadia Polikarpova. Just-in-time learning for bottom-up enumerative synthesis. *Proceedings of the ACM on Programming Languages*, 4(OOPSLA):1–29, 2020. doi:10.1145/3428295.
- 8 Matko Botinčan, Mike Dodds, and Suresh Jagannathan. Proof-directed parallelization synthesis by separation logic. *ACM Trans. Program. Lang. Syst.*, 35(2), July 2013. doi:10.1145/2491522.2491525.
- 9 José Cambronero, Sumit Gulwani, Vu Le, Daniel Perelman, Arjun Radhakrishna, Clint Simon, and Ashish Tiwari. FlashFill++: Scaling programming by example by cutting to the chase. *Proc. ACM Program. Lang.*, 7(POPL), January 2023. doi:10.1145/3571226.
- 10 João Costa Seco, Jonathan Aldrich, Luís Carvalho, Bernardo Toninho, and Carla Ferreira. Derivations with holes for concept-based program synthesis. In *Proceedings of the 2022 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software*, pages 63–79, 2022. doi:10.1145/3563835.3567658.
- 11 Andreea Costea, Amy Zhu, Nadia Polikarpova, and Ilya Sergey. Concise read-only specifications for better synthesis of programs with pointers. In Peter Müller, editor, *Programming Languages and Systems*, pages 141–168, Cham, 2020. Springer International Publishing. doi:10.1007/978-3-030-44914-8_6.
- 12 Ian Drosos, Titus Barik, Philip J Guo, Robert DeLine, and Sumit Gulwani. Wrex: A unified programming-by-example interaction for synthesizing readable code for data scientists. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pages 1–12, 2020. doi:10.1145/3313831.3376442.
- 13 Yu Feng, Ruben Martins, Osbert Bastani, and Isil Dillig. Program synthesis using conflict-driven learning. *ACM SIGPLAN Notices*, 53(4):420–435, 2018. doi:10.1145/3192366.3192382.
- 14 Yu Feng, Ruben Martins, Jacob Van Geffen, Isil Dillig, and Swarat Chaudhuri. Component-based synthesis of table consolidation and transformation tasks from examples. *ACM SIGPLAN Notices*, 52(6):422–436, 2017.
- 15 Yu Feng, Ruben Martins, Yuepeng Wang, Isil Dillig, and Thomas Reps. Component-based synthesis for complex APIs. In *Proceedings of the 44th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2017*, 2017.
- 16 Kasra Ferdowsifard, Shraddha Barke, Hila Peleg, Sorin Lerner, and Nadia Polikarpova. LooPy: Interactive program synthesis with control structures. *Proceedings of the ACM on Programming Languages*, 5(OOPSLA), October 2021. doi:10.1145/3485530.
- 17 Kasra Ferdowsifard, Allen Ordookhanians, Hila Peleg, Sorin Lerner, and Nadia Polikarpova. Small-step live programming by example. In *Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology*, pages 614–626. Association for Computing Machinery, 2020. doi:10.1145/3379337.3415869.
- 18 Jack Feser, Isil Dillig, and Armando Solar-Lezama. Inductive program synthesis guided by observational program similarity. *Proceedings of the ACM on Programming Languages*, 7(OOPSLA2):912–940, 2023. doi:10.1145/3622830.
- 19 John K Feser, Swarat Chaudhuri, and Isil Dillig. Synthesizing data structure transformations from input-output examples. *ACM SIGPLAN Notices*, 50(6):229–239, 2015. doi:10.1145/2737924.2737977.
- 20 Jonáš Fiala, Shachar Itzhaky, Peter Müller, Nadia Polikarpova, and Ilya Sergey. Leveraging Rust types for program synthesis. *Proc. ACM Program. Lang.*, 7(PLDI), June 2023. doi:10.1145/3591278.
- 21 Jonathan Frankle, Peter-Michael Osera, David Walker, and Steve Zdancewic. Example-directed synthesis: a type-theoretic interpretation. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*. ACM, 2016.

- 22 Joel Galenson, Philip Reames, Rastislav Bodik, Björn Hartmann, and Koushik Sen. CodeHint: Dynamic and interactive synthesis of code snippets. In *Proceedings of the 36th International Conference on Software Engineering*, pages 653–663, 2014.
- 23 Sumit Gulwani. Automating string processing in spreadsheets using input-output examples. In *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, POPL ’11, pages 317–330, New York, NY, USA, 2011. Association for Computing Machinery. doi:10.1145/1926385.1926423.
- 24 Zheng Guo, Michael James, David Justo, Jiaxiao Zhou, Ziteng Wang, Ranjit Jhala, and Nadia Polikarpova. Program synthesis by type-guided abstraction refinement. *Proceedings of the ACM on Programming Languages*, 4(POPL):1–28, 2019. doi:10.1145/3371080.
- 25 Sankha Narayan Guria, Jeffrey S Foster, and David Van Horn. RbSyn: type-and effect-guided program synthesis. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, pages 344–358, 2021. doi:10.1145/3453483.3454048.
- 26 Sankha Narayan Guria, Jeffrey S Foster, and David Van Horn. Absynthe: Abstract interpretation-guided synthesis. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1584–1607, 2023. doi:10.1145/3591285.
- 27 Kangjing Huang and Xiaokang Qiu. Bootstrapping library-based synthesis. In *International Static Analysis Symposium*, pages 272–298. Springer, 2022. doi:10.1007/978-3-031-22308-2_13.
- 28 Kangjing Huang, Xiaokang Qiu, Peiyuan Shen, and Yanjun Wang. Reconciling enumerative and deductive program synthesis. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI 2020, pages 1159–1174, New York, NY, USA, 2020. Association for Computing Machinery. doi:10.1145/3385412.3386027.
- 29 Shachar Itzhaky, Hila Peleg, Nadia Polikarpova, Reuben N. S. Rowe, and Ilya Sergey. Cyclic program synthesis. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2021, pages 944–959, New York, NY, USA, 2021. Association for Computing Machinery. doi:10.1145/3453483.3454087.
- 30 Susmit Jha, Sumit Gulwani, Sanjit A Seshia, and Ashish Tiwari. Oracle-guided component-based program synthesis. In *2010 ACM/IEEE 32nd International Conference on Software Engineering*, volume 1, pages 215–224. IEEE, 2010. doi:10.1145/1806799.1806833.
- 31 Ruyi Ji, Jingjing Liang, Yingfei Xiong, Lu Zhang, and Zhenjiang Hu. Question selection for interactive program synthesis. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 1143–1158, 2020. doi:10.1145/3385412.3386025.
- 32 Etienne Kneuss, Ivan Kuraj, Viktor Kuncak, and Philippe Suter. Synthesis modulo recursive functions. In *Proceedings of the 2013 ACM SIGPLAN international conference on Object oriented programming systems languages & applications*, pages 407–426, 2013. doi:10.1145/2509136.2509555.
- 33 Tessa A Lau, Pedro M Domingos, and Daniel S Weld. Version space algebra and its application to programming by demonstration. In *ICML*, pages 527–534, 2000.
- 34 Woosuk Lee. Combining the top-down propagation and bottom-up enumeration for inductive program synthesis. *Proceedings of the ACM on Programming Languages*, 5(POPL):1–28, 2021. doi:10.1145/3434335.
- 35 Woosuk Lee, Kihong Heo, Rajeev Alur, and Mayur Naik. Accelerating search-based program synthesis using learned probabilistic models. *ACM SIGPLAN Notices*, 53(4):436–449, 2018. doi:10.1145/3192366.3192410.
- 36 Xiang Li, Xiangyu Zhou, Rui Dong, Yihong Zhang, and Xinyu Wang. Efficient bottom-up synthesis for programs with local variables. *Proceedings of the ACM on Programming Languages*, 8(POPL), 2024. doi:10.1145/3632894.

- 37 Calvin Loncaric, Michael D Ernst, and Emina Torlak. Generalized data structure synthesis. In *Proceedings of the 40th International Conference on Software Engineering*, pages 958–968, 2018. doi:10.1145/3180155.3180211.
- 38 Calvin Loncaric, Emina Torlak, and Michael D. Ernst. Fast synthesis of fast collections. In *PLDI 2016: Proceedings of the ACM SIGPLAN 2016 Conference on Programming Language Design and Implementation*, pages 355–368, Santa Barbara, CA, USA, June 2016. doi:10.1145/2908080.2908122.
- 39 David Mandelin, Lin Xu, Rastislav Bodík, and Doug Kimelman. Jungloid mining: Helping to navigate the API jungle. In *Proceedings of the 2005 ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '05, pages 48–61, New York, NY, USA, 2005. Association for Computing Machinery. doi:10.1145/1065010.1065018.
- 40 Anders Miltner, Adrian Trejo Nuñez, Ana Brendel, Swarat Chaudhuri, and Isil Dillig. Bottom-up synthesis of recursive functional programs using angelic execution. *Proceedings of the ACM on Programming Languages*, 6(POPL):1–29, 2022. doi:10.1145/3498682.
- 41 Ashish Mishra and Suresh Jagannathan. Specification-guided component-based synthesis from effectful libraries. *Proceedings of the ACM on Programming Languages*, 6(OOPSLA2):616–645, 2022. doi:10.1145/3563310.
- 42 Peter-Michael Osera and Steve Zdancewic. Type-and-example-directed program synthesis. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 619–630, 2015. doi:10.1145/2737924.2738007.
- 43 Antti Oulasvirta and Pertti Saariluoma. Surviving task interruptions: Investigating the implications of long-term working memory theory. *International Journal of Human-Computer Studies*, 64(10):941–961, 2006. doi:10.1016/J.IJHCS.2006.04.006.
- 44 Peter O'Hearn, John Reynolds, and Hongseok Yang. Local reasoning about programs that alter data structures. In *Computer Science Logic: 15th International Workshop, CSL 2001 10th Annual Conference of the EACSL Paris, France, September 10–13, 2001, Proceedings 15*, pages 1–19. Springer, 2001. doi:10.1007/3-540-44802-0_1.
- 45 Hila Peleg, Roi Gabai, Shachar Itzhaky, and Eran Yahav. Programming with a read-eval-synth loop. *Proceedings of the ACM on Programming Languages*, 4(OOPSLA), November 2020. doi:10.1145/3428227.
- 46 Hila Peleg and Nadia Polikarpova. Perfect is the enemy of good: Best-effort program synthesis. In *34th European Conference on Object-Oriented Programming (ECOOP 2020)*. Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2020.
- 47 Phitchaya Mangpo Phothilimthana, Aditya Thakur, Rastislav Bodik, and Dinakar Dhurjati. Scaling up superoptimization. In *Proceedings of the Twenty-First International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 297–310, 2016.
- 48 Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. Program synthesis from polymorphic refinement types. *ACM SIGPLAN Notices*, 51(6):522–538, 2016. doi:10.1145/2908080.2908093.
- 49 Nadia Polikarpova and Ilya Sergey. Structuring the synthesis of heap-manipulating programs. *Proc. ACM Program. Lang.*, 3(POPL), January 2019. doi:10.1145/3290385.
- 50 Andrew Reynolds, Haniel Barbosa, Andres Nötzli, Clark Barrett, and Cesare Tinelli. cvc4sy: smart and fast term enumeration for syntax-guided synthesis. In *International Conference on Computer Aided Verification*, pages 74–83. Springer, 2019. doi:10.1007/978-3-030-25543-5_5.
- 51 John C Reynolds. Separation logic: A logic for shared mutable data structures. In *Proceedings 17th Annual IEEE Symposium on Logic in Computer Science*, pages 55–74. IEEE, 2002. doi:10.1109/LICS.2002.1029817.
- 52 Kensen Shi, Jacob Steinhardt, and Percy Liang. FrAngel: component-based synthesis with control structures. *Proceedings of the ACM on Programming Languages*, 3(POPL):1–29, 2019. doi:10.1145/3290386.

- 53 Xujie Si, Woosuk Lee, Richard Zhang, Aws Albarghouthi, Paraschos Koutris, and Mayur Naik. Syntax-guided synthesis of datalog programs. In *Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pages 515–527, 2018. doi:10.1145/3236024.3236034.
- 54 Calvin Smith and Aws Albarghouthi. MapReduce program synthesis. *Acm Sigplan Notices*, 51(6):326–340, 2016. doi:10.1145/2908080.2908102.
- 55 Sunbeom So and Hakjoo Oh. Synthesizing imperative programs from examples guided by static analysis. In *International Static Analysis Symposium*, pages 364–381. Springer, 2017. doi:10.1007/978-3-319-66706-5_18.
- 56 Armando Solar-Lezama, Liviu Tancau, Rastislav Bodik, Sanjit Seshia, and Vijay Saraswat. Combinatorial sketching for finite programs. In *Proceedings of the 12th international conference on Architectural support for programming languages and operating systems*, pages 404–415, 2006.
- 57 Jamie Stark and Andrew Ireland. Towards automatic imperative program synthesis through proof planning. In *14th IEEE International Conference on Automated Software Engineering*, pages 44–51. IEEE, 1999. doi:10.1109/ASE.1999.802091.
- 58 Yoshiaki Takashima, Ruben Martins, Limin Jia, and Corina S Păsăreanu. SyRust: automatic testing of Rust libraries with semantic-aware program synthesis. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, pages 899–913, 2021.
- 59 Abhishek Udupa, Arun Raghavan, Jyotirmoy V. Deshmukh, Sela Mador-Haim, Milo M.K. Martin, and Rajeev Alur. TRANSIT: Specifying protocols with concolic snippets. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '13*, pages 287–296, New York, NY, USA, 2013. Association for Computing Machinery. doi:10.1145/2491956.2462174.
- 60 Chenglong Wang, Alvin Cheung, and Rastislav Bodik. Synthesizing highly expressive SQL queries from input-output examples. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 452–466. ACM, 2017.
- 61 Zijiang Yang, Jinru Hua, Kaiyuan Wang, and Sarfraz Khurshid. EdSynth: Synthesizing api sequences with conditionals and loops. In *2018 IEEE 11th International Conference on Software Testing, Verification and Validation (ICST)*, pages 161–171. IEEE, 2018. doi:10.1109/ICST.2018.00025.
- 62 Yongho Yoon, Woosuk Lee, and Kwangkeun Yi. Inductive program synthesis via iterative forward-backward abstract interpretation. *Proceedings of the ACM on Programming Languages*, 7(PLDI):1657–1681, 2023. doi:10.1145/3591288.
- 63 Yongwei Yuan, Arjun Radhakrishna, and Roopsha Samanta. Trace-guided inductive synthesis of recursive functional programs. *Proceedings of the ACM on Programming Languages*, 7(PLDI):860–883, 2023. doi:10.1145/3591255.