

Interactive Program Synthesis

Vu Le

Microsoft Corporation
levu@microsoft.com

Mohammad Raza

Microsoft Corporation
moraza@microsoft.com

Daniel Perelman

Microsoft Corporation
danpere@microsoft.com

Abhishek Udupa

Microsoft Corporation
abudup@microsoft.com

Oleksandr Polozov

University of Washington
polozov@cs.washington.edu

Sumit Gulwani

Microsoft Corporation
sumitg@microsoft.com

Abstract

Program synthesis from incomplete specifications (e.g. input-output examples) has gained popularity and found real-world applications, primarily due to its ease-of-use. Since this technology is often used in an interactive setting, efficiency and correctness are often the key user expectations from a system based on such technologies. Ensuring efficiency is challenging, since the highly combinatorial nature of program synthesis algorithms does not fit in a 1–2 second response expectation of a user-facing system. Meeting correctness expectations is also difficult, given that the specifications provided are incomplete, and that the users of such systems are typically non-programmers.

In this paper, we describe how *interactivity* can be leveraged to develop efficient synthesis algorithms, as well as to decrease the cognitive burden that a user endures trying to ensure that the system produces the desired program. We build a formal model of user interaction along three dimensions: *incremental algorithm*, *step-based problem formulation*, and *feedback-based intent refinement*. We then illustrate the effectiveness of each of these forms of interactivity with respect to synthesis performance and correctness on a set of real-world case studies.

1. Introduction

Program synthesis is the task of generating a program in an underlying *domain-specific language* (DSL) from an intent specification provided by a user [3]. When the user in question is a non-programmer, the specification method must be concise and easy to provide without any programming expertise. The *programming by examples* (PBE) paradigm,¹ where the user intent is specified by the means of input-output examples or constraints, satisfies this requirement ideally. Although examples are succinct and easy for users to provide, they form an *under-specification* on the behavior of the desired program, adding inherent ambiguity into the problem definition.

Thanks to its ease of use, PBE has been effectively applied to many real-world scenarios in mass-market deployments. Two prominent examples are FlashFill [4] and FlashExtract [9]. FlashFill is a technology for automating repetitive string transformations, released as a feature in Microsoft Excel 2013. FlashExtract is a technology for extracting hierarchical data from semi-structured text files, released for log analytics in Microsoft Operations Management Suite and as the ConvertFrom-String cmdlet in Windows PowerShell.

The large-scale and continued adoption of PBE techniques is contingent on ensuring (a) the performance of the synthesizer, and (b) the correctness of the synthesized program. Responsiveness is critical to making PBE technologies usable. Users are willing to interact with the system in many rounds providing constraints iteratively, but any wait time exceeding 1–2 seconds per round leads to a frustrating experience. Deployed systems like FlashFill and FlashExtract ensure performance by restricting the expressiveness

of the underlying DSL or by bounding the execution time of the synthesizer. Such restrictions limit the applicability and growth of these technologies: when the underlying DSL is enriched to meet the users’ demands for capturing a larger class of tasks, the performance of the synthesizer starts degrading.

The correctness of the synthesized program is critical to building trust in PBE systems. In the past, intent ambiguity in PBE has been primarily handled by imposing a sophisticated ranking on the DSL [14]. While ranking goes a long way in avoiding undesirable interpretations of the user’s intent, it is not a complete solution. For example, FlashFill in Excel is designed to cater to users that care not about the program but about its behavior on the small number of input rows in the spreadsheet. Such users can simply eye-ball the outputs of the synthesized program and provide another example if they are incorrect. However, this becomes much more cumbersome (or impossible) with a larger spreadsheet.²

We have observed that inspecting the synthesized program directly also does not establish enough confidence in it even if the user knows programming. Two main reasons for this are (i) program readability,³ and (ii) the users’ uncertainty in the desired intent due to hypothetical unseen corner cases in the data. This feedback was consistent with a recent user study [10], which illustrated that users find it less useful and approachable to inspect the synthesized programs but would prefer more interactive models to converge to the desired program and to be confident of its correctness.

Due to ambiguity of intent in PBE, the standard user interaction model in this setting is for the user to provide constraints *iteratively* until the user is satisfied with the synthesized program or its behavior on the known inputs. However, most work in this area, including FlashFill and FlashExtract, has not been formally modeled as an iterative process. In this paper, we propose an interactive formulation of program synthesis that leverages the inherent iterative nature of synthesis from under-specifications.

Interactivity

We present *interactivity* as the solution to addressing the performance and the correctness challenges associated with PBE. Our inspiration to make program synthesis interactive comes from the standard world of programming. In programming, interactivity manifests in at least three key dimensions:

Incremental: A programmer writes a function by iteratively refining it, for instance, as in test-driven development.

² John Walkenbach, famous for his Excel textbooks, labeled FlashFill as a “controversial” feature. He wrote: “It’s a great concept, but it can also lead to lots of bad data. I think many users will look at a few “flash filled” cells, and just assume that it worked. But my preliminary tests leads me to this conclusion: Be very careful.” [16]

³ Stephen Owen, a certified MVP (“Most Valued Professional”) in Microsoft technologies, said the following of a program synthesized by FlashExtract: “If you can understand this, you’re a better person than I am.” [11]

¹ In this article, whenever applicable, we use PBE to more generally denote programming using under-specifications.

Step-based: A programmer splits the task of writing a program into simpler steps of writing various functions for individual sub-tasks one by one.

Feedback-based: A programmer may use various tools, such as program analysis or test coverage measurement, to obtain actionable feedback on code quality, correctness, and performance.

We propose integrating these dimensions into the interactive process of program synthesis.

Incremental synthesis The standard PBE model requires the user to refine her intent in iterative rounds by providing additional constraints on the current candidate program. The standard approach has been to re-run the synthesizer afresh with the conjunction of the original constraints and the new constraints. In this paper, we describe an alternative technique, which makes the synthesizer *incremental* and provides significant performance benefits.

Most PBE techniques use a data structure called *version space algebra* (VSA) [12] to succinctly represent and compute the set of programs in the underlying DSL that are consistent with the user-provided constraints. Our key idea is to observe that a VSA is simply an *Abstract Syntax Tree (AST) based representation of a language*, and it can be translated into a sub-DSL of the original DSL. The new round of synthesis is then carried out over this new sub-DSL using only the new constraints.

Step-based synthesis In many PBE domains, there exist well-defined sub-computations that can be exposed to the user intuitively as steps. This helps with both synthesizer performance and ensuring correctness of the synthesized program.

For instance, the programs in the FlashFill DSL apply a conditional logic to compute different string transformations based upon the input string. When the user provides a few input-output examples, the FlashFill synthesizer conjectures which of those examples will be addressed by the same branch of the top-level conditional. This decision is based on a few examples and can be incorrect. Furthermore, for a complicated task that requires many examples, the conditional learning algorithm in FlashFill does not scale [8, 17].

The conditional logic is a sub-computation that can be naturally exposed to the user as a clustering of input rows. The user can drive this task by providing examples of rows that should be in the same cluster. Our key idea to enable step-based synthesis is to *allow associating user constraints with named subexpressions* in the DSL.

Feedback-based synthesis In the standard PBE model, the user is responsible for providing additional constraints in each iterative round of synthesis. In each iteration, the user chooses what additional constraints to provide, based purely on the behavior of the synthesized program, but without any directed feedback from the synthesizer. However, the synthesizer has knowledge about the specific ambiguities in the constraints provided by the user w.r.t. the underlying DSL. Our key idea is to translate this knowledge into *natural queries for the user*, whose responses *form the next set of additional constraints*. This provides two key benefits:

- The feedback from the tool can help continue the progress towards the desired program by pointing out remaining ambiguities. Otherwise, the user may stop earlier than intended.
- The new set of constraints constructed from the user’s response can accelerate the progress towards the desired program by resolving the ambiguity faster (relative to the constraints that the user would have provided otherwise). To this end, we introduce a novel automatic component in the conventional PBE model, called *the hypothesizer*. It proactively analyzes the current set of candidate programs and constructs a set of queries that, if answered by the user, would best resolve the ambiguities.

We give illustrations of different kinds of queries for real-world PBE domains that can be easily answered by the user with little cognitive load and that help reduce the ambiguity significantly. To avoid asking too many queries, we associate queries with a *disambiguation score* that represents the benefit of asking that query for ambiguity resolution. We generate feedback only if this score exceeds a certain threshold. Our experimental results illustrate that our strategy for choosing the disambiguation score and the threshold leads to few false positives and almost no false negatives.

Contributions This paper makes the following contributions:

- We formally define the general problem of interactive program synthesis, extending upon the CEGIS [5, 15], SyGuS [1], and FlashMeta [12] formalisms.
- We propose an approach for incremental synthesis that leverages the observation that VSAs can be viewed as DSLs. We present experimental results on its performance benefits.
- We show how to model step-based interaction as part of the general program synthesis paradigm. We present experimental analysis of its effectiveness in improving synthesis convergence.
- We present various feedback strategies for helping users in refining their intent. We present qualitative and quantitative results on the effectiveness of these strategies in ensuring that the synthesized programs are correct.

Structure This paper is structured as follows. §2 provides some background on the problem of PBE, and introduces the FlashFill, FlashExtract, and FlashSplit DSLs, which are used as running examples throughout the paper. §3 gives a high-level overview of our interactive synthesis formulation. The next 3 sections formally describe each problem dimension: incremental (§4), step-based (§5), and feedback-based (§6) synthesis. §7 presents our evaluation for each dimension. Finally, §8 reviews related work, and §9 concludes.

2. Background

In this section, we introduce three DSLs that are used as case studies in the paper, and provide some background on inductive synthesis.

2.1 Domain-Specific Language

We follow the formalism of FlashMeta (a.k.a. the PROSE framework) [12]. A synthesis problem is defined for a given *domain-specific language* (DSL) \mathcal{L} . A DSL is specified as a context-free grammar (CFG), with each nonterminal symbol N defined through a set of *rules*. Each rule is an application of an *operator* to some symbols of \mathcal{L} . All symbols and operators are *typed*. If $N := F(N_1, \dots, N_k)$ is a grammar rule and $N : \tau$, then the output type of F must be τ . A DSL has a designated *output symbol* $\text{output}(\mathcal{L})$, which is a start nonterminal in the CFG of \mathcal{L} .

Every (sub-)program P rooted at a symbol $N : \tau$ in \mathcal{L} maps an *input state*⁴ σ to a value of type τ . A state is a mapping of free variables $\text{FV}(P)$ to their bound values. Variables in a DSL are introduced by *let* definitions and λ -functions. The output symbol has a single free variable—an *input symbol* $\text{input}(\mathcal{L})$ of the DSL. For brevity, we use the notation $N[x := v]$ for “*let* $x = v$ *in* N ”.

Every operator F in a DSL \mathcal{L} has some executable semantics. Many operators are generic, and typically reused across different DSLs (e.g. Filter and Map list combinators). Others are domain-specific, and defined only for a given DSL. Operators are assumed to be deterministic and pure, modulo unobservable side effects.

FlashFill DSL Figure 1 shows an excerpt from the definition of FlashFill DSL, which transforms a list of strings (i.e., a spreadsheet

⁴ DSLs in PBE typically do not involve mutation, so an input σ is technically an *environment*, not a *state*. We keep the term “state” for historical reasons.

```

Language FlashFill;

@output string start := e | std.ITE(cond, e, start);
string e := f | Concat(f, e);
string f := ConstStr(w)
    | let string x = std.Kth(vs, k) in sub;

string sub           := SubStr(x, pp);
tuple<int, int> pp    := std.Pair(pos, pos);
int pos             := AbsPos(x, k) | RegPos(x, rr, k);
tuple<Regex, Regex> rr := std.Pair(r, r);

bool cond := let string s = std.Kth(vs, k) in b;
@extern[std.text.match] bool b; // FV(b) = {s: string}
@input string[] vs; string w; int k; Regex r;

```

Figure 1: FlashFill DSL \mathcal{L}_{FF} for string transformations in spreadsheets [4]. Each program rooted at *start* takes an input a spreadsheet row *vs* and performs a chain of if-elseif matches on some cells of *vs*. The expression in the chosen ITE branch returns a concatenation of constants and input substrings.

```

Language FlashExtract.Sequence;

@output StringRegion[] seq :=
    std.Map(λx ⇒ std.Pair(pos, pos), lines) // LinesMap
    | std.Map(λt ⇒ let string x = GetSuffix(d, t) in
        std.Pair(t, pos), posSeq) // StartSeqMap
    | std.Map(λt ⇒ let string x = GetPrefix(d, t) in
        std.Pair(pos, t), posSeq); // EndSeqMap

int[] posSeq := std.FilterInt(i0, k, rrSeq);
int[] rrSeq := RegexMatches(d, rr);
StringRegion[] lines := std.FilterInt(i0, k, fllLines);
StringRegion[] fllLines := std.Filter(λs ⇒ b, allLines);
StringRegion[] allLines := SplitLines(d);

@extern[std.text.match] bool b; // FV(b) = {s: string}
@extern[FlashFill] int pos; // FV(pos) = {x: string}
@input StringRegion d; int i0; int k;

```

Figure 2: FlashExtract DSL \mathcal{L}_{ES} for selecting a sequence of spans in a textual document *d* [9]. Each program rooted at *seq* is either a LinesMap program (split *d* into lines and select a span in each line), a StartSeqMap program (select a sequence of starting positions of the spans and map each to its corresponding ending position), or a EndSeqMap program (select a sequence of ending positions of the spans and map each to its corresponding starting position). This DSL references position extraction logic *pos* from \mathcal{L}_{FF} (Figure 1) and string predicates *b* from the standard library.

```

Language FlashExtract;

@output TreeNode E := struct | arr;
ObjectNode struct := Struct(E1, ..., En) | Prop(id, Er);
ArrayNode arr := Seq(id, Es) | std.Map(λd ⇒ E, Es);

StringRegion[] Es := @extern[FlashExtract.Sequence] seq;
StringRegion Er := @extern[FlashFill] sub[x := d];
@input StringRegion d; string id;

```

Figure 3: FlashExtract meta-DSL \mathcal{L}_{FE} for extraction of a dataset from a textual document *d* [9]. Each program rooted at *E* builds an object tree of sequences and structs, extracted from *d*. They are extracted using sequence selection logic from \mathcal{L}_{ES} (Figure 2) and substring selection logic from \mathcal{L}_{FF} (Figure 1), respectively. The leaves of the tree are field programs: region or sequence extractions. Each field program is marked with a unique ID.

```

Language FlashSplit;

@output StringRegion[] fields := SplitByDelimiters(v, d);
StringRegion[] d := LookAround(v, c, rr) | Union(d, d);
StringRegion[] c := ExactMatches(v, s);
    | IncludeWhitespace(v, s);
tuple<Regex, Regex> rr := std.Pair(r, r);

@input StringRegion v; string s; Regex r;

```

Figure 4: FlashSplit DSL \mathcal{L}_{FS} for splitting an input record *v* into a sequence of fields separated by delimiters *d*. Each delimiter is determined by a LookAround operator, which represents a constant string match *c* in *v* that also matches regular expressions *r*₁ and *r*₂ on the left and right side respectively (similar to position extraction logic from \mathcal{L}_{FF} in Figure 1). Constant string matches are either exact matches of a string, or including any surrounding whitespace.

```

191.128.19.55 - [09/Jun/2016:18:05:33 -0800] "GET /checks.txt HTTP/1.1" 200 617 "www.facebook.com/"
174.13.04.3 - [09/Jun/2016:19:43:23 -0800] "GET /images/pic.png HTTP/1.1" 403 18 "-"
192.16.201.109 - [10/Jun/2016:06:10:03 -0800] "GET /pdf/document.pdf HTTP/1.1" 204 23 "https://www.microsoft.com/en-us/windows"
11.0.4.50 - [10/Jun/2016:16:10:02 -0800] "GET /index.html HTTP/1.1" 234 74 "https://www.yahoo.com/"
191.169.12.13 - [13/Jun/2016:11:30:02 -0800] "GET /index.html HTTP/1.1" 505 75 "https://www.yahoo.com/"
172.18.0.102 - [13/Jun/2016:16:12:34 -0800] "GET /logs/access-log HTTP/1.1" 234 32 "https://google.com/"
192.19.2.100 - [11/Jun/2016:17:32:36 -0800] "GET /index.html HTTP/1.1" 500 15003 "-"
50.129.2.78 - [11/Jun/2016:17:45:38 -0800] "GET /data/2/4 HTTP/1.1" 130 933 "-"
171.19.3.12 - [11/Jun/2016:22:12:01 -0800] "GET /data/ HTTP/1.1" 254 5265 "https://www.aol.com"
191.168.125.112 - [11/Jun/2016:23:12:52 -0800] "GET /pictures/pic2.png HTTP/1.1" 306 502340 "https://www.google.com/"
174.26.0.223 - [12/Jun/2016:16:13:04 -0800] "GET /images/pic4.gif HTTP/1.1" 800 6876 "https://www.google.com/"
175.16.0.24 - [12/Jun/2016:17:29:33 -0800] "GET /index.html HTTP/1.1" 206 56254 "-"
196.168.29.105 - [12/Jun/2016:18:33:11 -0800] "GET /styles.css HTTP/1.1" 354 1346 "https://www.bing.com/"
101.22.54.38 - [13/Jun/2016:20:32:43 -0800] "GET /js/scripts.js HTTP/1.1" 304 3472 "https://www.yahoo.com/"

```

```

191.128.19.55 -- 09/Jun/2016:18:05:33 0800 -- GET /checks.txt HTTP/1.1 200 617 www.facebook.com/
174.13.04.3 -- 09/Jun/2016:19:43:23 0800 -- GET /images/pic.png HTTP/1.1 403 18 -
192.16.201.109 -- 10/Jun/2016:06:10:03 0800 -- GET /pdf/document.pdf HTTP/1.1 204 23 https://www.microsoft.com/en-us/windows
11.0.4.50 -- 10/Jun/2016:16:10:02 0800 -- GET /index.html HTTP/1.1 234 74 https://www.yahoo.com/
191.169.12.13 -- 11/Jun/2016:11:30:02 0800 -- GET /index.html HTTP/1.1 505 75 https://www.yahoo.com/
172.18.0.102 -- 13/Jun/2016:16:12:34 0800 -- GET /logs/access-log HTTP/1.1 234 32 https://google.com/
192.19.2.100 -- 11/Jun/2016:17:32:36 0800 -- GET /index.html HTTP/1.1 500 15003 -
50.129.2.78 -- 11/Jun/2016:17:45:38 0800 -- GET /data/2/4 HTTP/1.1 130 933 -
171.19.3.12 -- 11/Jun/2016:22:12:01 0800 -- GET /data/ HTTP/1.1 254 5265 https://www.aol.com
191.168.125.112 -- 11/Jun/2016:23:12:52 0800 -- GET /pictures/pic2.png HTTP/1.1 306 502340 https://www.google.com/
174.26.0.223 -- 12/Jun/2016:16:13:04 0800 -- GET /images/pic4.gif HTTP/1.1 800 6876 https://www.google.com/
175.16.0.24 -- 12/Jun/2016:17:29:33 0800 -- GET /index.html HTTP/1.1 206 56254 -
196.168.29.105 -- 12/Jun/2016:18:33:11 0800 -- GET /styles.css HTTP/1.1 354 1346 https://www.bing.com/
101.22.54.38 -- 13/Jun/2016:20:32:43 0800 -- GET /js/scripts.js HTTP/1.1 304 3472 https://www.yahoo.com/

```

Figure 5: A sample splitting of a log file from a web server into fields. Fields are delimited by various delimiters, such as “- - [” and “] ”. Note how “/” is used as a delimiter between some fields, but also occurs as a non-delimiter in other fields such as the URLs.

row) into an output string. In this DSL, non-prefixed operators such as AbsPos and Concat are user-defined, while namespace-prefixed operators such as std.Kth are defined in the standard library of PROSE. Symbols marked as @extern reference symbols in other DSLs or in the standard library.

FlashExtract DSL Figure 3 shows the definition of FlashExtract DSL. A FlashExtract program extracts from a text document *d* a hierarchical tree that consists of nested structs and sequences. A leaf of this tree is either a sequence *seq* of regions (i.e., StringRegions) or a region *sub* within a parent node. We refer to both *seq* and *sub* as a *field* in the program. All fields in \mathcal{L}_{FE} are associated with some IDs. Each program corresponds to a *schema* that defines the structure of the output tree. The logic for extracting a sequence of spans within a given region is separated into a sub-DSL \mathcal{L}_{ES} (Figure 2). The logic for extracting a subregion within a region is imported from \mathcal{L}_{FF} .

FlashSplit DSL Figure 4 shows the definition of the FlashSplit DSL. FlashSplit is a PBE system for field splitting in text-based data sources (e.g. log files). Such files often contains data rows with a large number of fields, separated by arbitrary delimiters. Figure 5 shows a sample splitting of a web server log.

A FlashSplit program splits an input data record into a sequence of field values separated by delimiters. In practice, a string that is used as a delimiter between some fields also occurs inside other field values. To address this, FlashSplit supports *contextual delimiters*, which match constant strings that occur between certain regular expression patterns on the left and right.

2.2 Inductive Synthesis Problem

An *inductive synthesis* problem refers to synthesis of a *program set* $\tilde{N} \subset \mathcal{L}$ that is consistent with a given *inductive specification* φ . An inductive specification (or simply “spec”) is a collection of input-output constraints $\{\sigma_i \rightsquigarrow \psi_i\}_{i=1}^n$. Each *constraint* ψ_i is a unary Boolean predicate over the *output* of the desired program on the corresponding input state σ_i . The simplest form of ψ is a concrete output value; in this case, φ is a collection of *input-output examples*. In this work, we use PBE and inductive synthesis interchangeably.

A program P *satisfies* a spec φ (written $P \models \varphi$) iff it satisfies all constraints $\sigma_i \rightsquigarrow \psi_i$ in φ . A program P satisfies an input-output constraint $\sigma \rightsquigarrow \psi$ iff its output $\llbracket P \rrbracket \sigma$ on the given input state σ satisfies the corresponding constraint predicate ψ , i.e. if $\psi(\llbracket P \rrbracket \sigma)$ is true. A program set \tilde{N} is said to be *valid* w.r.t. a spec φ (written $\tilde{N} \models \varphi$) iff all programs $P \in \tilde{N}$ satisfy φ .

A *synthesis algorithm*, given a symbol $N \in \mathcal{L}$ and a spec φ , learns a *valid* program set \tilde{N} of programs rooted at N . We denote the corresponding problem definition as $\text{Learn}(N, \varphi)$. By definition, the algorithm must be *sound*—every program in \tilde{N} must satisfy φ . A synthesis algorithm is said to be *complete* if it learns *all possible* programs in \mathcal{L} that satisfy φ . Since \mathcal{L} is usually infinite because of unrestricted constants, completeness is usually defined w.r.t. some finitization of \mathcal{L} (possibly dependent on a given synthesis problem).

Example 1. A typical spec for FlashFill synthesis is shown below:

$$\varphi = \left\{ \begin{array}{l} \{vs \mapsto [“323-708-7700”]\} \rightsquigarrow “(323) 708-7700” \\ \{vs \mapsto [“555.988.0139”]\} \rightsquigarrow “(555) 988.0139” \end{array} \right.$$

It consists of 2 constraints ψ_1 and ψ_2 . Each constraint ψ_i is an *example constraint*: it maps a single input state σ_i to the desired FlashFill output string o_i . In each input state σ_i , the input variable vs of \mathcal{L}_{FF} is bound to a single input spreadsheet cell.

Example 2. A typical spec for FlashExtract sequence synthesis is shown below:

$$\varphi = \{d \mapsto [“Brazil 23 21\n Bulgaria 35 32\n”]\} \rightsquigarrow [“Brazil”, \dots]$$

It contains a single *prefix constraint* ψ with an input document d in the state σ and a *prefix* of the desired sequence of selections.

2.3 Version Space Algebra

A *version space algebra* (VSA) is a data structure for efficient storage of candidate programs in deductive synthesis. Since deductive synthesis typically works with large program sets (up to 10^{50} programs), it requires a special data structure to represent them in polynomial space and perform polynomial-time set operations. We refer the reader to [12, §4] for a detailed overview of VSAs; this section provides only a brief background.

Definition 1 (Version space algebra). *Let N be a symbol in a DSL \mathcal{L} . A *version space algebra* is a representation for a set \tilde{N} of programs rooted at N . The grammar of VSAs is:*

$$\tilde{N} := \{P_1, \dots, P_k\} \mid \mathbf{U}(\tilde{N}_1, \dots, \tilde{N}_k) \mid F_{\bowtie}(\tilde{N}_1, \dots, \tilde{N}_k)$$

where F is any k -ary operator in \mathcal{L} , and P_j are some programs in \mathcal{L} . The semantics of VSA as a set of programs is given as follows:

$$\begin{array}{ll} P \in \{P_1, \dots, P_k\} & \text{if } \exists j: P = P_j \\ P \in \mathbf{U}(\tilde{N}_1, \dots, \tilde{N}_k) & \text{if } \exists j: P \in \tilde{N}_j \\ P \in F_{\bowtie}(\tilde{N}_1, \dots, \tilde{N}_k) & \text{if } P = F(P_1, \dots, P_k) \wedge \forall j: P_j \in \tilde{N}_j \end{array}$$

Intuitively, a VSA is a DAG where each node represents a set of programs. *Leaf nodes* contain explicit enumerations of programs; they are composed into larger sets by two possible VSA constructor nodes. *Union nodes* represent a set union of their constituent VSAs.

Join nodes represent a cross-product of their constituent VSAs, with an associated operator F applied to all combinations of parameter programs from the cross-product.

VSAs support multiple efficient set operations. In PBE we typically make use of: intersection $\tilde{N}_1 \cap \tilde{N}_2$, clustering based on program outputs on a given input \tilde{N}/σ , ranking w.r.t. a scoring function $\text{Top}_h(\tilde{N}, k)$, and projection (filtering) onto a subset of programs satisfying a given spec $\tilde{N} \upharpoonright \varphi$.

2.4 Backpropagation

In the FlashMeta formalism, the main synthesis algorithm typically employed for PBE is *backpropagation*, or *deductive synthesis*. It follows the grammar of \mathcal{L} *top-down*, applying the principle of divide-and-conquer. At each step, it reduces the synthesis problem $\text{Learn}(N, \varphi)$ to simpler subproblems: either on parameters of the symbol N , or on subexpressions of the spec φ .

Deductive synthesis makes use of small domain-specific procedures called *witness functions*. They *backpropagate* constraints on a program $F(N_1, \dots, N_k)$ to deduced constraints on its subexpressions N_1, \dots, N_k . Two kinds of witness functions exist: conditional and non-conditional. Non-conditional witness functions take as input a spec φ on F and transform it into a spec on their respective parameter N_i . Conditional witness functions take as input a spec φ on F with a bound value v_j of some other parameter N_j , and transform φ into a spec on parameter N_i under the assumption that $\llbracket N_j \rrbracket \sigma = v_j$. Intuitively, conditional witness functions introduce branching in the top-down search process of deductive synthesis. They split the search space into disjoint partitions based on possible outputs of a target subprogram rooted at N_j , and then continue with synthesis in each partition independently.

At a high level, deductive synthesis solves a problem $\text{Learn}(N, \varphi)$ via a combination of 3 problem reduction techniques (see [12, §5] for a detailed presentation):

Split on alternatives: If $N := N_1 \mid N_2$, then the algorithm solves the subproblems $\text{Learn}(N_1, \varphi)$ and $\text{Learn}(N_2, \varphi)$, and takes a union of results: $\tilde{N} = \tilde{N}_1 \mathbf{U} \tilde{N}_2$.

Backpropagation of witnesses: If $N := F(N')$, then the algorithm invokes the corresponding *witness function* $\omega_{N'}$ for N' in F . The witness function transforms the spec φ into a necessary (and often sufficient) spec φ' on N' . The algorithm then solves the subproblem $\text{Learn}(N', \varphi')$.

If $\omega_{N'}$ is precise (i.e., φ' is sufficient), then any valid program $P' \in \tilde{N}'$, when used as an argument F , produces a valid program $F(P') \models \varphi$. Thus, the program set $F_{\bowtie}(\tilde{N}')$ is a valid solution for φ . If $\omega_{N'}$ is imprecise (i.e., φ' is only necessary), then the algorithm returns a projection $F_{\bowtie}(\tilde{N}') \upharpoonright \varphi$.

Split on conditional execution: It is often impossible to construct a precise witness function for a parameter program under all possible spec conditions. However, it is usually possible assuming additional restrictions on *other* parameters of the same operator. If $N := F(N_1, N_2)$, often F permits two backpropagation procedures: a simple witness function $\omega_1(\varphi)$ for the first parameter, and a *conditional* witness function $\omega_2(\varphi \mid \llbracket N_1 \rrbracket \sigma = v)$ for the second parameter. The function ω_2 produces a spec φ_2 for N_2 that is necessary (and often sufficient) to satisfy φ under the *assumption* that the program chosen for N_1 evaluates to v .

In such situation, the algorithm first invokes ω_1 and solves the produced subproblem $\text{Learn}(N_1, \varphi_1)$. It then *clusters* \tilde{N}_1 on the given inputs, and splits the search into independent branches, one per cluster (i.e. one per each possible output of programs from \tilde{N}_1). Within each branch, it operates under the assumption that all programs in a cluster $\tilde{N}_i \subset \tilde{N}_1$ produce the same concrete

value v_i as an output. Given that value, the function ω_2 produces a spec φ_{2i} for N_2 , and the algorithm solves the subproblem $\text{Learn}(N_2, \varphi_{2i})$. The final result is a union of solution sets over all branches: $\tilde{N} = \bigcup_i F_{\boxtimes}(\tilde{N}_{1i}, \tilde{N}_{2i})$, valid by construction.

3. Overview

A typical workflow of a PBE session with an end user follows a flowchart shown in Figure 6. The synthesis system is parameterized with (i) a DSL \mathcal{L} , which defines a search space for target programs, and (ii) a ranking function h , which resolves ambiguity between multiple program candidates. The user communicates her intent to the system in the form of input-output examples (or, more generally, constraints) φ for the desired program. The system performs a search in \mathcal{L} for the subset of programs that are consistent with φ , ranks them w.r.t. h , and returns top-ranked candidate program(s) to the user. The user inspects the program(s), and, if it does not match her desired behavior, refines the spec φ by introducing additional examples (or constraints), without any help from the synthesizer. The system now searches for a subset of programs in \mathcal{L} that are consistent with the new spec φ' . This cycle continues until either (a) the user is satisfied with the current program, or (b) the system discovers that the current spec is unsatisfiable in \mathcal{L} .

Figure 6 and similar flowcharts in synthesis literature implement different variants of *counterexample-guided inductive synthesis* (CEGIS) [15], a common inductive synthesis technique. While effective in many applications, we found this workflow lacking in several aspects when applied on a *mass-market industrial scale* with an end user playing the role of an *oracle*. We outline our observations and solutions to associated problems below.

Incrementality In conventional CEGIS, the learner uses the refined spec φ' at each iteration to synthesize a new valid program set $\tilde{N} \models \varphi'$. All the information accumulated in the synthesis session is contained in φ' as a conjunction of provided examples, counterexamples, and more general constraints. Typically in PBE, the learner takes it all into account by solving a fresh synthesis problem, searching in the DSL \mathcal{L} for programs that are consistent with all constraints in φ' . As the size of φ' grows with each iteration, this synthesis problem becomes more complex, thereby slowing down the search process [6].

Our key observation here is that *the program set \tilde{N}_i learned at iteration i can be transformed into a new DSL \mathcal{L}' , which will become the search space for synthesis in iteration $i + 1$ instead of \mathcal{L}* . Notice that every refined spec φ' imposes an *additional* restriction on the desired program. Thus, an i^{th} program set \tilde{N}_i must be a subset of the program set \tilde{N}_{i-1} , learned at the previous iteration.

We developed an efficient procedure for transforming a VSA \tilde{N}_i into a DSL definition \mathcal{L}_{i+1} , which replaces \mathcal{L} in the next iteration of synthesis. This replacement achieves two significant speedups. First, the size of \mathcal{L}_i is monotonically decreasing, and quickly becomes many orders of magnitude smaller than \mathcal{L} . Second, at each iteration i we are only searching for programs consistent with the latest introduced constraint ψ_i , since the DSL \mathcal{L}_i by construction only contains programs that satisfy previous constraints $\psi_1, \dots, \psi_{i-1}$.

Step-based formulation In conventional CEGIS, the user is limited to providing constraints on the overall behavior of the desired program. Apart from the current candidate program P , the user does not have any insight into the learner and its configuration (i.e., \mathcal{L} and h). Thus, the user’s guidance is limited to *counterexamples* to the candidate program, which the learner includes in the refined spec φ' . As the number of refining iterations grows, so does the user’s frustration, since she cannot influence the debugging experience without any understanding of the learner.

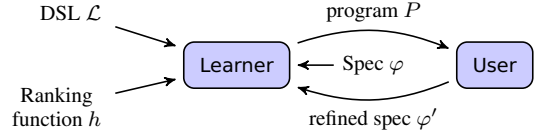


Figure 6: Learner-user communication in conventional PBE.

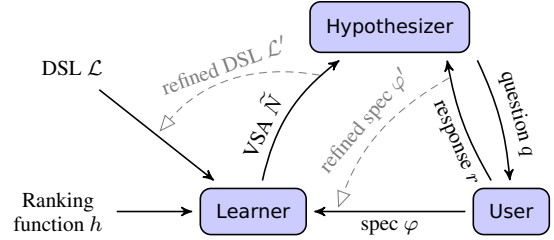


Figure 7: Learner-user communication in interactive PBE.

Many synthesis tasks have a notion of *sub-tasks*, on which the user can easily and naturally specify individual constraints. One way to model this is to support synthesis of various sub-tasks as independent synthesis tasks in respective sub-DSLs. In that case management of the composition of those sub-tasks lies with the application layer on top of the synthesis sub-systems. This is non-trivial, sophisticates the application logic, and is often implemented in an application-specific manner disallowing code re-use.

In our formalism, we model such interaction by allowing the user to provide constraints on *named subexpressions* in a *compound DSL* instead of just top-level constraints. These named subexpressions correspond to aforementioned sub-tasks (or “steps”) in the interaction process. Analogously to bottom-up programming, we allow the user to define building blocks of the target program *step-by-step* first, before assembling them into a larger expression.

In addition to simplifying synthesis application development, step-based formalism also improves the debugging experience during learner-user interaction. In conventional CEGIS, when the current candidate program is incorrect, the user has to analyze the behavior of the whole program, come up with a counterexample, and communicate it to the learner, starting a new iteration. In contrast, with the step-based formulation she can focus on learning individual named subexpressions in the program first. In addition to reducing the scope of required reasoning, it also allows the user to “lock” individual subexpressions as correct. The synthesizer can then leverage knowledge about their behavior when learning other subexpressions of the desired program, thereby completing the entire synthesis task in fewer iterations. In our evaluation (§7.2) we have verified that the step-based formulation significantly reduces the number of examples for many real-world synthesis tasks.

Feedback-based interaction Discovering counterexamples is relatively easy when the CEGIS oracle is a modern SMT solver (for domains that can be modeled in an SMT theory) and when the full spec is known. However, an end user serving as the oracle often does not have a clear understanding of the full spec, and may suffer significant cognitive load with this task. At every iteration, the current candidate program set \tilde{N}_i contains thousands of ambiguous programs, and humans struggle with reasoning about possible ambiguities in intent specification. In contrast, the synthesis system can analyze ambiguities in \tilde{N}_i and derive the most efficient way to resolve them by proactively soliciting concrete knowledge from the user. This observation introduces a third important actor in the CEGIS flowchart, which we call *the hypothesizer*.

Formally, any *constraint type* used in PBE (e.g. example, prefix, output type) states a *property* on a subset of the DSL. Given a program set \tilde{N}_i , the hypothesizer deduces possible properties that best disambiguate among programs in \tilde{N}_i . Any such property is convertible to a Boolean or multiple-choice question q , which the hypothesizer asks the user. Any response r for q is convertible to a concrete constraint ψ , which begins a new iteration of synthesis.

Such *feedback-based* interaction has several major benefits. First, it reduces the cognitive load on the user: instead of analyzing the program’s behavior, she only answers concrete questions. Second, it significantly reduces the number of synthesis iterations thanks to the hypothesizer’s insight into the program set \tilde{N} and its choice of disambiguating questions. Finally, feedback and proactiveness increases the user’s confidence and trust in the system.

Interactive Synthesis

Figure 7 shows a typical workflow in *interactive program synthesis*. As before, the learner is parameterized with a DSL $\mathcal{L} = \mathcal{L}_0$ and a ranking function h . Suppose the initial user-provided spec $\varphi = \varphi_1$ consists of a single constraint ψ_1 .

The learner synthesizes a valid program set $\tilde{N}_1 \models \psi_1$. This set becomes the search space \mathcal{L}_1 for the next iteration of synthesis. To refine the spec for the next iteration, the learner either waits for new constraints from the user, or proactively invokes the hypothesizer. The hypothesizer analyzes the set \tilde{N}_1 and generates the best disambiguating question q_1 for the user. After the user answers it with a response r_1 , the hypothesizer translates it to a constraint ψ_2 , appends it to the spec, and invokes the next synthesis iteration with the refined spec φ_2 . The learner synthesizes a valid program set $\tilde{N}_2 \models \varphi_2$ from the DSL \mathcal{L}_1 , and the cycle continues until convergence or unsatisfiability.

At each iteration, the user can change the context of the synthesis and specify constraints either on the overall program P , or on some named subexpression P' in the program. The learner then continues synthesis for this subexpression. From this moment, any iteration that changes a candidate program for P' also triggers an incremental relearning of any subexpressions in P that are dependant on P' .

The challenges of conventional CEGIS described above motivate the need for modeling the interactive workflow in a first-class manner in the program synthesis formalism. Building on the commonly used formulations of CEGIS [5, 15] and SyGuS [1], we extend their problem definition to incorporate learner-user interaction.

Problem 1 (Interactive Program Synthesis). Let \mathcal{L} be a DSL, and N be a symbol in \mathcal{L} . Let \mathcal{A} be an *inductive synthesis algorithm* for \mathcal{L} , which solves problems of type $\text{Learn}(N, \varphi)$ where φ is an *inductive spec* on a program rooted at N . The specs φ are chosen from a fixed class of *supported spec types* Φ . The result of $\text{Learn}(N, \varphi)$ is some set \tilde{N} of programs rooted at N that are consistent with φ .

Let φ^* be a spec on the output symbol of \mathcal{L} , called a *task spec*. A φ^* -driven interactive program synthesis process is a finite series of 4-tuples $\langle N_0, \varphi_0, \tilde{N}_0, \Sigma_0 \rangle, \dots, \langle N_m, \varphi_m, \tilde{N}_m, \Sigma_m \rangle$, where

- Each N_i is a nonterminal in \mathcal{L} ,
- Each φ_i is a spec on N_i ,
- Each \tilde{N}_i is some set of programs rooted at N_i s.t. $\tilde{N}_i \models \varphi_i$,
- Each Σ_i is an **interaction state**, explained below,

which satisfies the following axioms for any program $P \in \mathcal{L}$:

- A.** $(P \models \varphi^*) \Rightarrow (P \models \varphi_i)$ for any $0 \leq i \leq m$;
- B.** $(P \models \varphi_j) \Rightarrow (P \models \varphi_i)$ for any $0 \leq i < j \leq m$ s.t. $N_i = N_j$.

We say that the process is **converging** iff the top-ranked program of the last program set in the process satisfies the task spec:

$$P^* = \text{Top}_h(\tilde{N}_m, 1) \models \varphi^*$$

and the process is **failing** iff the last program set is empty: $\tilde{N}_m = \emptyset$.

An **interactive synthesis algorithm** $\hat{\mathcal{A}}$ is a procedure (parameterized by \mathcal{L} , \mathcal{A} , and h) that solves the following problem:

$$\text{Learnter: } \begin{cases} \langle N_0, \varphi_0, \perp \rangle \mapsto \langle \tilde{N}_0, \Sigma_0 \rangle \\ \langle N_i, \varphi_i, \Sigma_{i-1} \rangle \mapsto \langle \tilde{N}_i, \Sigma_i \rangle, \quad i > 0 \end{cases}$$

In other words, at each iteration i the algorithm receives the i^{th} learning task $\langle N_i, \varphi_i \rangle$ and its own interaction state Σ_{i-1} from the previous iteration. The type and content of Σ_i is unspecified and can be implemented by $\hat{\mathcal{A}}$ arbitrarily.

Definition 2. We say that an *interactive synthesis algorithm* $\hat{\mathcal{A}}$ is **complete** iff for any task spec φ^* :

- If $\exists P \in \mathcal{L}$ s.t. $P \models \varphi^*$ then $\hat{\mathcal{A}}$ eventually converges for any φ^* -driven interactive synthesis process.
- Otherwise, $\hat{\mathcal{A}}$ eventually fails for any φ^* -driven interactive synthesis process.

The notion of an interactive synthesis process formally models a typical learner-user interaction where φ^* describes the desired program. The general nature of definitions in Problem 1 allows many different implementations for $\hat{\mathcal{A}}$. In addition to completeness, different implementations (and choices for the state Σ) strive to satisfy different *performance objectives*, such as:

- Number of interaction rounds (e.g. examples) m ,
- The total amount of information communicated by the user,
- Cumulative execution time of all $m + 1$ learning calls.

In the rest of this paper, we present several specific instantiations of interactive synthesis algorithms that optimize these objectives.

4. Incremental Synthesis

Our incremental synthesis algorithm is based on two key ideas: **(a)** translation of VSAs as DSLs, and **(b)** local resolution of different constraint types by memoization of intermediate subproblems.

VSA as a DSL Recall that a VSA \tilde{N} is a DAG-like program set representation with two kinds of constructor nodes (unions and joins) and one kind of leaf nodes (explicit sets). Our key observation here is that this representation is in fact simply an *AST-based representation* of a sub-DSL $\mathcal{L}' \subset \mathcal{L}$. More specifically, the DAG of \tilde{N} is isomorphic to a context-free grammar of a subset of \mathcal{L} .

Figure 8 shows an algorithm for translating \tilde{N} into a grammar of \mathcal{L}' . It performs an isomorphic graph translation, converting VSA unions (**U**) into CFG alternatives ($N := N_1 \mid N_2$), VSA joins (F_{\bowtie}) into CFG operator productions ($N := F(\dots)$), and explicit program sets into CFG terminals annotated with their possible values.

Note that as a subset of \mathcal{L} , the new DSL \mathcal{L}' does not introduce any new operators. Thus, all witness functions for \mathcal{L} are still applicable for synthesis in \mathcal{L}' . Moreover, \mathcal{L}' is finite and its terminals are annotated with explicit sets of permitted values, which allows fast learning of constants for any spec type simply via set filtering.

Constraint resolution Deductive synthesis relies on existence of witness functions, which backpropagate constraints top-down through the grammar. Every witness function is defined for a particular *constraint type* ψ that it is able to decompose. While some generic operators allow efficient backpropagation procedures for common constraint types (see, e.g. [2, 9]), most witness functions are domain-specific.

We have identified a useful set of constraint types that occur in various PBE domains and often permit efficient witness functions. We broadly classify these constraints in 3 categories depending on their *descriptive power*, with different incremental synthesis techniques required for each category.

```

function VSAToDSL(VSA  $\tilde{N}$ )
1: Let  $V$  be a set of fresh nonterminals, one per each non-leaf node in  $\tilde{N}$ 
2: Let  $\Sigma$  be a set of fresh terminals, one per each leaf node in  $\tilde{N}$ 
3: // We write  $\text{sym}(\tilde{N}') \in V \cup \Sigma$  to denote the corresponding fresh
   // symbol for a node  $\tilde{N}'$  from  $\tilde{N}$ 
4: Productions  $R \leftarrow \emptyset$ 
5: // Create "symbol := symbol" productions for all union nodes
6: for all union nodes  $\tilde{N}' = \mathbf{U}(\tilde{N}_1, \dots, \tilde{N}_k)$  in  $\tilde{N}$  do
7:    $R \leftarrow R \cup \{\text{sym}(\tilde{N}') := \text{sym}(\tilde{N}_i) \mid i = 1 \dots k\}$ 
8: // Create operator productions for all join nodes
9: for all join nodes  $\tilde{N}' = F_{\bowtie}(\tilde{N}_1, \dots, \tilde{N}_k)$  in  $\tilde{N}$  do
10:   $R \leftarrow R \cup \{\text{sym}(\tilde{N}') := F(\text{sym}(\tilde{N}_1), \dots, \text{sym}(\tilde{N}_k))\}$ 
11: // Annotate terminal symbols with values extracted from leaf nodes
12: for all leaf nodes  $\tilde{N}' = \{P_1, \dots, P_k\}$  in  $\tilde{N}$  do
13:  Annotate in  $\Sigma$  that  $\text{sym}(\tilde{N}') \in \{P_1, \dots, P_k\}$ 
14: return the context-free grammar  $G = \langle V, \Sigma, R, \text{sym}(\tilde{N}) \rangle$ 

```

Figure 8: An algorithm for translating a VSA \tilde{N} of programs in a DSL \mathcal{L} into an isomorphic grammar for a sub-DSL $\mathcal{L}' \subset \mathcal{L}$.

Definitive constraints: constructively *define* a subset of the DSL by the means of backpropagation through witness functions. In other words, properties of the program’s output that they describe are narrow enough to enable deductive reasoning. For instance:

- *Example constraint:* “output = v ”,
- *Membership constraint:* “output $\in \{v_1, v_2, v_3\}$ ”,
- *Prefix constraint:* “output = $[v_1, v_2, \dots]$ ”,
- *Subset/subsequence constraint:* “output $\sqsupseteq [v_1, v_2, v_3]$ ”.

Locally refining constraints: do not define a DSL subset on their own, but can be used to *refine* an existing program set in a witness function for some DSL operator(s). For instance:

- *Datatype constraint:* “output: τ ”. Eliminates all top-level programs rooted at any type-incompatible DSL operators.
- *Provenance constraint:* describes the desired construction method for some parts of the output. For example, in Flash-Fill it may take form “substring $[i : j]$ of the output example o_k is extracted from location ℓ of the corresponding input σ_k ”, or “substring $[i : j]$ of the output example o_k is a date value, formatted as “YYYY-MM-DD””. Allows simple domain-specific elimination of invalid subprograms.
- *Relevance constraint:* marks inputs or parts of the input as *required* or *irrelevant*. Eliminates all programs that do not use any required parts or reference any irrelevant parts.

Globally refining constraints: do not define a DSL subset on their own and do not permit any efficient local refining logic in witness functions. They can only be satisfied by filtering an existing program set on the topmost level of the DSL (i.e., by projecting the set on the constraint). For instance:

- *Negative example constraint:* “output $\neq v$ ”,
- *Negative membership constraint:* “output $\not\in v$ ”.

Given a program set \tilde{N} and a new constraint ψ , we filter \tilde{N} w.r.t. ψ differently depending on the category of ψ .

- If ψ is definitive, it seeds a new full round of deductive synthesis, which may narrow down the set unpredictably. We convert the set \tilde{N} into an isomorphic DSL $\text{VSAToDSL}(\tilde{N})$, and use it as a search space for a new synthesis round with a spec consisting of a single constraint ψ .

- If ψ is locally refining, it is only relevant to select witness functions. Suppose these witness functions backpropagate specs for the operator $F(N_1, \dots, N_k)$.

We first identify occurrences of F in \tilde{N} . Each node of kind $F_{\bowtie}(\tilde{N}_1, \dots, \tilde{N}_k)$ in \tilde{N} has been constructed during the top-down grammar traversal in a previous iteration of deductive synthesis as a solution to some intermediate synthesis subproblem $\text{Learn}(F(N_1, \dots, N_k), \varphi_F)$. To enable incrementality, we keep references to all intermediate specs φ_F that were produced at this level in the previous synthesis iteration.

We now repeat the learning for F on all retained subproblems, but with their previous specs φ_F conjoined with the new constraint ψ . The witness functions for F take ψ into account and produce potentially more refined specs for N_1, \dots, N_k . These new specs are definitive (since they were produced by witness functions), and thus initiate new rounds of incremental synthesis on $\text{VSAToDSL}(\tilde{N}_1), \dots, \text{VSAToDSL}(\tilde{N}_k)$. The results replace $\tilde{N}_1, \dots, \tilde{N}_k$ in \tilde{N} without affecting the rest of the set.

- If ψ is globally refining, it cannot be efficiently resolved by inspecting \tilde{N} or invoking witness functions. Thus, we have to compute the projection $\tilde{N}|_{\psi}$ at the top level. An efficient implementation of the projection operation computes the clustering $\tilde{N}/_{\sigma}$ on an input σ from the constraint ψ . All programs in the same cluster \tilde{N}_k produce the same output v_k on σ . If the constraint ψ only references the output of the desired program (as all globally refining constraints do), then all programs in \tilde{N}_k either satisfy ψ or not. Thus, we simply take a union of clusters where the corresponding outputs v_k satisfy ψ .

We note that for many globally refining constraints this operation is trivial once the clustering is computed. For example, a negative example constraint “output $\neq v$ ” eliminates at most 1 cluster—the one where $v_k = v$, if one exists.

This incremental learning algorithm can be expressed as an instance of interactive synthesis formalism from [Problem 1](#). For that, set Σ_i to be a tuple of \tilde{N}_i (required to become the search space at the next iteration) and all intermediate specs produced by witness functions (required to resolve locally refining constraints).

Theorem 1. *If the underlying one-shot learning algorithm \mathcal{A} for \mathcal{L} is complete, then (a) the incremental learning $\tilde{\mathcal{A}}$ is complete, and (b) at each iteration i the result \tilde{N}_i of incremental learning is equal to the result of cumulative one-shot learning $\text{Learn}_{\mathcal{A}}(N_i, \varphi_0 \wedge \dots \wedge \varphi_i)$.*

Proof. Omitted for lack of space. Intuitively, the theorem follows from the fact that \tilde{N}_i are monotonically non-increasing and from the axioms of interactive synthesis in [Problem 1](#). \square

5. Step-based Synthesis

In this section we introduce our formalism for the *step-based program synthesis*. We first start with a motivational case study of FlashExtract, then define the step-based interaction process formally, and discuss its applicability to our DSL examples in this paper.

5.1 Motivation

Consider the FlashExtract DSL \mathcal{L}_{FE} , presented in [Figure 3](#). A program in \mathcal{L}_{FE} extracts a structured dataset from a given text file. The schema for this constructed dataset may arbitrarily combine sequences, structs, and primitive field extractions. However, the users usually are more inclined to provide examples for a single field and then move on to another field (as opposed to providing examples of tuples). They are more comfortable providing partial

constraints on the individual fields of the desired output, as opposed to burdensome complete examples of extracted objects.

This situation motivated the developers of FlashExtract-based applications to expose constraints on individual fields as the primary UI for end users. They then provided these constraints as individual subproblems for sequence (\mathcal{L}_{ES}) or region (\mathcal{L}_{FF}) extraction. This required the developers to implement their own wrapper code to support the entirety of the learner-user interaction—that is, **(a)** the schema classes for combining sequence and field programs in a compound extraction program, and **(b)** a learning logic to guess the desired schema for a given task. This situation had 3 drawbacks:

1. Implementing (a) and (b) is cumbersome. Moreover, the application-specific implementations usually cannot be shared across different applications.
2. The schema learning logic (b) is non-trivial. In fact, the schema format can be naturally captured using a recursive DSL in our formalism (\mathcal{L}_{FE}), and hence can/should be learned using the synthesis algorithm from the PROSE SDK.
3. When the synthesizer treats multiple fields in the same task as independent problems, it cannot leverage insights learned from one field to improve its learning logic for another field. As a result, the entire task requires more examples to converge.

To address these drawbacks, we define step-based synthesis over compound DSLs as a novel first-class formalism over FlashMeta.

5.2 Problem Definition

Definition 3 (Compound DSL). *A DSL \mathcal{L} is called a **compound DSL** if it includes any **extern nonterminals** N_1, \dots, N_m , which resolve to output symbols of some **sub-DSLs** $\mathcal{L}_1, \dots, \mathcal{L}_m$ respectively.*

Definition 4 (Constraints on named subexpressions). *Given a compound DSL \mathcal{L} , a **named constraint** $\psi^{e: N}$ is specified for an extern nonterminal N in \mathcal{L} . Its meaning is as follows:*

“There exists a subexpression rooted at N in the desired compound program. We mark it with ID e . This subexpression must satisfy the constraint ψ .”

When used in a context of iterative learner-user interaction on a larger task, all named constraints with the same ID e apply to the same subexpression in the desired compound program in \mathcal{L} .

***Named specs** $\varphi^{e: N}$ are defined similarly as conjunctions of named constraints on e .*

At any point during the learner-user interaction, the step-based synthesis algorithm has accumulated the following information in its interaction state Σ_i :

- (i) a spec φ on the compound program,
- (ii) a list of named subexpressions e_1, \dots, e_n , which appear in the compound program, and
- (iii) specs $\varphi^{e_1}, \dots, \varphi^{e_n}$ on these subexpressions.

When the user provides a new named constraint ψ^{e_i} on a subexpression e_i , deductive synthesis must learn a new VSA of compound programs, incorporating the new constraint in the process. It starts the top-down synthesis from φ , as before. When the search process reaches the nonterminal N_i with some deduced spec φ' , the learner compares it with the new constraint ψ^{e_i} (since φ' is known to already be compatible with the previous named spec φ^{e_i}). The witness functions for N_i now must consider the new constraint ψ^{e_i} .

Two options exist at this point: **(a)** we are learning the named subexpression e_i , or **(b)** we are learning an unrelated subexpression, which happens to also start from N_i . Option (b) has already been considered in the previous synthesis iterations, so the learner can just reuse the corresponding VSA. Thus, it must only detect whether option (a) is applicable. To do that, the witness functions for

N_i check if the deduced spec φ' is compatible with the new constraint ψ^{e_i} . If it is, they generate a refined spec φ'' for subsequent learning of e_i in its sub-DSL \mathcal{L}_i , and the VSA for e_i is replaced with the new one.

Re-learning of e_i may impact other subexpressions of the program, which depend on e_i . In the FlashMeta formalism these dependencies are updated automatically, by the means of conditional witness functions. If a VSA learned for a prerequisite subexpression changes, then its clustering will also change, and the new branches will produce new specs for the dependent subexpression.

5.3 Case Studies

Example 3 (FlashExtract). Consider the task of extracting a sequence of customer records, each of which contains a customer name and phone number from the following text file:

```
Carrie Dodson
202-555-0153
Leonard Robledo
945-051-0159
...
```

There are three named subexpressions in the task: the customer record, the customer name, and the phone number. Step-based synthesis allows DSL designers to build different interaction models for FlashExtract with minimal efforts. For instance, one can build a non-step-based model where users provide the compound spec φ (i.e., some nested records with names and phone numbers). In this model, users have to specify the relationship of the subexpressions in the spec φ . The PowerShell cmdlet `ConvertFrom-String` adopts this model because of its command-line user interface. In another model, users only need to provide the named specs iteratively and in a step-based manner for each of the three fields. In particular, users learn each of these in order by increasingly giving some named field instances until the field is identified correctly. Interactive systems such as the one by Mayer et al. [10] adopts this model.

The learner of the first model is simpler than that of the second one because the relationships among the subexpressions (and therefore the structure of the output/program) are given. The learner expands the grammar along the path specified by the structure in φ and learns all subexpressions in φ . In contrast, at any point of time, the learner of the second model has to take into account both the current (partial) program and the new named spec to create a new (partial) program. For instance, while learning for customer name, the learner promotes the sequence *field* of customer records (which is learned in previous step) to a sequence of *structs* which contains a field customer name. Subsequently, the learning of phone number puts a new field phone number into the already constructed struct customer record. Although the step-based, interactive model involves more work (for the DSL designer), it helps to reduce the total number of examples across all fields that the user needs to provide in an extraction task (see §7.2).

Example 4 (FlashFill). Consider the task of normalizing phone numbers into the format “(XXX) XXX-XXXX” as follows:

Input	Output
485-7829	(133) 485-7829
555-0175	(033) 555-0175
555 0122	(033) 555-0122
033 555 6694	(033) 555-6694
...	...

In the non-step-based model, because users provide the whole compound program spec φ , FlashFill has to keep tracks of all possible partitioning of input rows, and for each partitioning, all of its transformation programs. Since the problem is intractable, in practice people usually limit the number of the partitions, which


```

procedure DISAMBIGUATE(Candidate programs  $\tilde{N}$ , current spec  $\varphi$ )
1: Analyze the ambiguities in  $\tilde{N}$  w.r.t.  $\varphi$ .
   Let  $Q$  be a set of questions that may resolve ambiguity in  $\tilde{N}$ 
2:  $q^* \leftarrow \operatorname{argmax}_{q \in Q} \operatorname{ds}(q, \tilde{N}, \varphi)$  // Compare the disambiguation scores
   of all questions
3: if  $\operatorname{ds}(q^*, \tilde{N}, \varphi) < \text{threshold } T$  then
4:   break
5: else
6:   Present the question  $q^*$  to the user
7:   Let  $r$  be the user’s response to  $q$ 
8:   Let  $\psi$  be the response  $r$  converted into a constraint
9:   return  $\psi$  to the learner and invoke a new round of synthesis

```

Figure 9: The hypothesizer’s proactive disambiguation algorithm.

affects the expressiveness of the DSL. The ranking system in this model also requires more efforts because it deals with the space of all satisfying compound programs. For instance, from the examples it is unclear if the last row belongs to the partition that contains “555” (and therefore “033” is a constant string), or it should have its own partition (which takes “033” from the input).

In contrast, step-based model separates the two subproblems as two named subexpressions, and enables users to provide spec for the subproblems. Because this model eliminates the implicit dependency between the two subproblems, it makes the problem more tractable and simplifies the ranking system.

6. Feedback-based Synthesis

In this section, we formally define our proposed learner-user interaction model that leverages proactive feedback in the form of queries to the user. We also present its applications to specific example DSLs of this paper (FlashFill and FlashSplit), and discuss practical issues involved in picking a question, evaluating its effectiveness, and defining a stopping criterion.

6.1 Problem Definition

Let \mathcal{L} be a DSL. Let Ψ be a set of top-level *constraint types* supported by the synthesizer and witness functions for \mathcal{L} . For each constraint type $\psi \in \Psi$ we associate a *descriptive question* q such that a response r for this question directly constitutes an instance of ψ . We denote such a constraint as $\Psi(r)$. Questions q can be Boolean (usually in ternary logic) or multiple-choice. We denote the set of possible responses for q as $R(q)$.

Example 5. An *example constraint* “output = v ” corresponds to a multiple-choice question “Is the desired output on an input σ equal to v_1, v_2, \dots , or v_k ?” A response to this question constitutes an example constraint “output = v_i ” for the chosen i .

A *datatype constraint* “output: τ ” corresponds to a Boolean question “Is the desired program a computation of type τ ? Yes (always), no (never), or unknown (maybe).”

Questions, like constraints, can be domain-specific. In FlashFill, a *relevance constraint* states “an input v_i must/must not appear in the program.” It corresponds to a Boolean question “Should the input v_i be used? Yes (always), no (never), or unknown (maybe).”

As described in §3, we introduce a novel component called the *hypothesizer* in the learner-user interaction model. Figure 9 shows its disambiguation algorithm. Given a VSA \tilde{N} of current candidate programs and the current iteration’s spec φ , the job of the hypothesizer is to analyze \tilde{N} and pick the best question to resolve ambiguities in \tilde{N} . If \tilde{N} has no ambiguities, or if the hypothesizer is not confident in the effectiveness of potential questions, it considers the current candidate program $P^* = \operatorname{Top}_h(\tilde{N}, 1)$ correct and does not ask any questions at this iteration.

Disambiguation score To evaluate a question’s effectiveness, the hypothesizer is parameterized with a *disambiguation score* function $\operatorname{ds}(q, \tilde{N}, \varphi)$. Higher disambiguation scores correspond to more effective questions q —that is, constraints generated by answering q eliminate more incorrect programs from \tilde{N} . Since the hypothesizer cannot predict the user’s response, $\operatorname{ds}(q, \tilde{N}, \varphi)$ must represent *potential effectiveness* of q for any possible outcome.

Disambiguation score functions may be domain-specific or general-purpose. In our evaluation, we found different functions to perform well for different DSLs. In this section, we present one efficient *general-purpose* disambiguation score function, which is independent of the current iteration’s spec φ but takes into account the ranking scores of alternative candidate programs in \tilde{N} .

The ranking-based disambiguation score function prefers a question that promotes higher-ranked programs:

$$\operatorname{ds}_R(q, \tilde{N}, \varphi) \stackrel{\text{def}}{=} \min_{r \in R(q)} \max_{P \in \tilde{N}_r} h(P)$$

where \tilde{N}_r is a set projection $\tilde{N} \upharpoonright_{\Psi(r)}$, and h is a ranking function provided with the DSL. In other words, $\operatorname{ds}_R(q, \tilde{N}, \varphi)$ is higher if every response for the question q leads to a higher-ranked alternative program among the candidates that are consistent with this response.

This disambiguation score can be efficiently evaluated for many constraint types. For instance, calculating $\operatorname{ds}_R(q, \tilde{N}, \varphi)$ for *example constraints* amounts to clustering \tilde{N} and comparing the top-ranked programs across all clusters. Alternatively, we can quickly compute a good approximation to $\operatorname{ds}_R(q, \tilde{N}, \varphi)$ by randomly sampling k programs from the VSA and considering only their outputs.

6.2 Case Studies

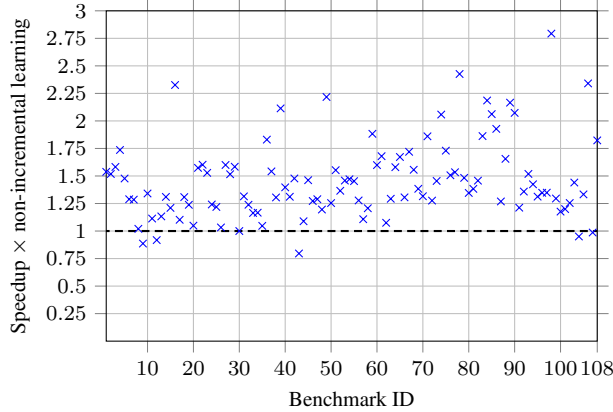
FlashFill Our feedback-driven synthesis for the FlashFill language (Figure 1) uses example constraint questions. In order to efficiently generate these questions, we sample 2000 programs from the VSA and cluster on those, assigning the disambiguation score ds_R as described above. We chose 2000 because empirically it was a good balance between performance and having a high probability of including at least one program from every large cluster.

FlashSplit Our instantiation of the feedback-driven paradigm for the FlashSplit language (Figure 4) intends to resolve ambiguities in learning field-splitting programs with arbitrary delimiters. The FlashSplit’s ranking function favors combinations of delimiters that occur regularly across all input rows and produce a uniform splitting. The synthesis ambiguity lies in choosing a particular combination of consistently aligned delimiters constitutes the desired program. For example, there exists a huge number of natural ways to split the server log data in Figure 5 into fields (e.g. separate the “Date” field into “Day”, “Month”, and “Year”). Examples help to resolve this ambiguity, but each data row may have up to 50 fields; thus, providing even a single complete example of split positions in a row may be too burdensome. We alleviate this user effort in two ways.

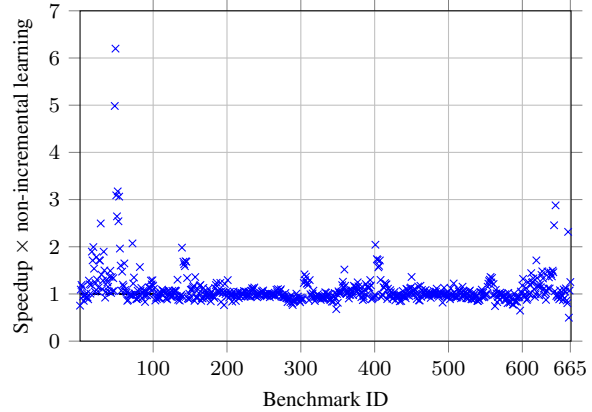
Using constraints other than examples FlashSplit supports a *subset constraint* (as defined in §4), where the user specifies a partial example with only some of the split positions on a given input.

Providing feedback in the form of questions The hypothesizer analyzes the program set \tilde{N} of all candidate delimiter expressions, and guides the user by asking questions about ambiguous split positions. We investigated two kinds of questions to elicit feedback:

Binary position questions. A binary position question $q \in Q_b$ presents a single position in the input row and asks if it is a desired splitting point. The answer “Yes” corresponds to a *positive subset constraint* over desired split positions, and the



(a) Speedup on the FlashFill DSL.



(b) Speedup on the FlashExtract DSL.

Figure 10: Speedups obtained by the incremental synthesis algorithm vs. the non-incremental algorithm. Values higher above the $y = 1$ line (where the runtimes are equal) are better.

answer “No” corresponds to a *negative membership constraint* over desired split positions. The hypothesizer generates such questions when there exist positions that are unique to certain candidate delimiter expressions.

Confirmation questions. A confirmation question $q \in Q_c$ presents a set of positions to the user and asks whether all of these positions are valid splitting points. The answer “Yes” corresponds to a *positive subset constraint* with all presented positions. The answer “No” corresponds to a FALSE constraint, meaning that no program in \mathcal{L}_{FS} can satisfy the user’s constraints. The hypothesizer generates such questions when the user provides an example that can only be satisfied by one delimiter expression in \tilde{N} . In this case, all of the positions determined by this delimiter expression must be part of the output, or else the system can declare failure early on. This saves the user the effort of providing all of these redundant examples individually.

In Section §7.3 we evaluate effectiveness of these question types individually and in combination. For individual question types, since our ranking function h does not distinguish among the top ranked well-aligned delimiters, we use a uniform disambiguation score over all candidate programs for a particular question type.⁵ For a combined system, we found that binary questions perform better than confirmation questions as the number of desired splits grows. This inspires a split-based domain-specific disambiguation score:

$$ds_{FS}(q, \tilde{N}, \varphi) \stackrel{\text{def}}{=} \begin{cases} \text{MinSplits}(\tilde{N}, \varphi) - t & \text{if } q \in Q_b \\ t - \text{MinSplits}(\tilde{N}, \varphi) & \text{if } q \in Q_c \end{cases}$$

where $\text{MinSplits}(\tilde{N}, \varphi)$ is the minimum number of splits produced by any program in \tilde{N} on any input in φ , and t is a predefined threshold. We found $t = 30$ to work well in our experiments.

7. Evaluation

7.1 Incremental Synthesis

We implemented the incremental synthesis algorithm, described in §4, in the PROSE framework. We evaluate the incremental synthesis algorithm in the context of the FlashFill DSL. For this case study, we picked all the benchmarks which required the user to provide two or more examples to learn a correct program, from among the

⁵ h can be improved by statistical analysis of delimiters that more commonly occur in practice. As we found FlashSplit sufficiently efficient for our scenarios in evaluation, we left such analysis for future work.

FlashFill benchmarks. All the data reported in this subsection were obtained by repeating each experiment ten times and averaging the results after discarding outliers.

Figure 10 summarizes the results of our evaluation. Figure 10(a) plots the *speedup* obtained by incremental algorithm over the non-incremental algorithm for each benchmark for the FlashFill DSL. These were computed by dividing the execution time of the non-incremental algorithm by the execution time of the incremental algorithm for each benchmark. We observe that almost all the speedup values are greater than one, with the exceptions being extremely short-running benchmarks as mentioned earlier. Further, the incremental algorithm achieves a geometric mean speedup of 1.42 over the non-incremental algorithm, across the 108 benchmarks considered for the FlashFill DSL.

Figure 10(b) describes the result of a similar experiment with the FlashExtract DSL, where a total of 665 benchmarks were used for evaluation. The performance gains in the case of the FlashExtract benchmarks are more modest on average than in the case of the FlashFill benchmarks. Our investigations revealed that although the pruned VSAs at each iteration consist of fewer programs, the *structure* of these VSAs can be more complex. For example, we observed VSAs where a union operation was performed on several hundred join nodes. Pruning along each of these paths during incremental learning sometimes results in large execution time overheads. We note however, that (a) incremental learning provides speedups in more than half the FlashExtract benchmarks, sometimes over 6X, (b) in most of the cases where incremental learning does not provide a speedup, the execution times are within 10% of the non-incremental algorithm. In fact, slowdowns greater than 10% were observed in just 12.7% of the benchmarks.

We conclude the evaluation of incremental learning by mentioning that over 80% of the FlashFill benchmarks required only two learning iterations, and over 90% of the FlashExtract benchmarks required three or fewer learning iterations. Despite the relatively small number of learning iterations, our evaluation demonstrates that incrementality yields significant performance improvements.

7.2 Step-based Synthesis

We use FlashExtract to evaluate the effectiveness of step-based synthesis. Our benchmarks consist of 100 files collected from help forums, product teams (that have exposed FlashExtract capability as a feature in their products), and their end users. Each file corresponds to an extraction task that extracts several fields into a hierarchical output tree. The number of extracted fields in a task ranges from 1 to 36 fields (mean 5.7, median 4). A field contains from 1 entry (such

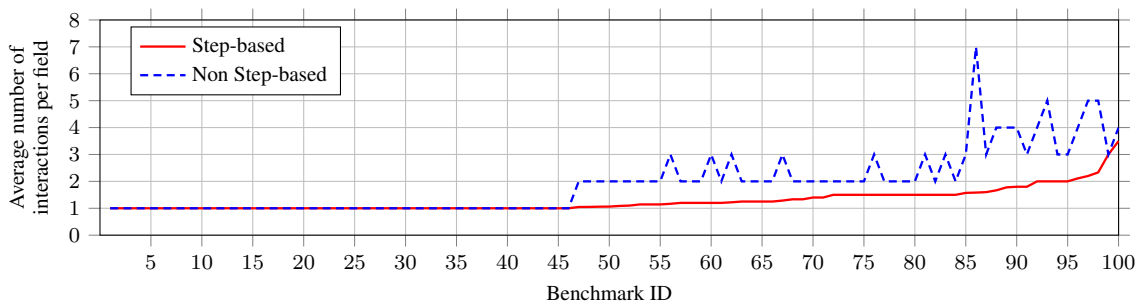


Figure 11: The average number of interactions per field across all benchmarks. Lower is better.

as a customer name in an email) to a few thousand entries (such as the event timestamps in a log file). We simulate user actions in two models: non-step-based and step-based.

Non-step-based FlashExtract The baseline of this evaluation is a non-interactive FlashExtract where the user has to provide examples for all the fields at once. If extraction fails (i.e., the output is not intended), the user needs to provide new examples for each of the *failing* fields. A field fails if its output is not identical to the expected output. The new examples of a field include all output that FlashExtract has correctly identified in the previous interaction and the first discrepancy between the output and the expected output of that field. The extraction succeeds if the executing output tree is identical to the expected output tree.

Step-based FlashExtract The user of this system extracts fields in topological order (i.e., from top-level fields to leaf fields), which is usually also the document order. If the current field fails, the user gives new examples until FlashExtract produces the expected output. The selection of new examples resembles that in the non-interactive setting, which selects the correct prefix of the output and the first discrepancy between the output and the expected output. Once the field extraction succeeds, the user moves to the next field. The whole extraction succeeds if all fields are identified correctly.

Results We refer the step of locating the first output discrepancy and providing new examples for a field as an *interaction*. One of the most important goals of PBE is to reduce the number of interactions to provide better user experience. In fact, some product teams demand that PBE systems should work with only one interaction most of the time for industrial adoption.

Figure 11 shows the average number of interactions per field across all benchmarks, ordered by the number of interactions in step-based FlashExtract. The evaluation shows that step-based FlashExtract requires fewer interactions than non-step-based FlashExtract in more than half of the benchmarks that require more than one example. For benchmarks that require only one interaction, step-based FlashExtract performs similarly to non-step-based FlashExtract because the process is entirely non-interactive. By dividing the extraction task into several steps, step-based FlashExtract can “lock” a field if its extraction has been successful and focus on the remaining fields. The learning of subsequent fields therefore does not have any effects on the previously learned fields. In contrast, non-step-based FlashExtract has to maintain all fields as once. When a field fails, users may have to provide examples for other fields in addition to those for the failing field.

7.3 Feedback-based Synthesis

FlashFill We evaluate the feedback-driven synthesis for FlashFill on a set of 457 text transformation tasks. In the *baseline setting*, the user provides the earliest incorrect row as the next example at each iteration. In the *feedback-driven setting*, the system instead proactively asks the user disambiguating questions on selected input rows until the disambiguation score falls below the threshold T . We set $T = 0.47$ as the mean of the score distribution over our tasks.

We evaluate FlashFill’s feedback on two dimensions: *cognitive burden* and *correctness*. Cognitive burden is defined as the number of rows the user has to *read and verify* in the process. In the baseline setting, it is the number of examples + the number of correct rows before the first discrepancy that are verified at each iteration. In the feedback-driven setting, it is the number of questions answered.

Correctness is a combination of *false positives* and *false negatives*. False positive questions occur when FlashFill keeps asking questions after the program is already correct. False negative questions occur when FlashFill stops before it finds a correct program.

In correctness evaluation, only 2/457 tasks completed incorrectly (i.e., with false negatives). The majority of tasks (342/457) finish with the same number of examples, and the number of false positives in the rest never exceeds 4 (specifically, 90 tasks with 1 false positive, 22 with 2, and 1 task with 4 false positives).

Table 1 compares the cognitive burden of both settings. It shows the histogram distribution of our tasks for each pair of verified row counts in baseline and feedback-driven settings. The baseline setting often requires the user to inspect more rows (that is, the numbers *below* the diagonal in Table 1 are larger than the numbers above it).

FlashSplit We evaluate the feedback-driven synthesis for FlashSplit on a set of 77 splitting tasks on different log files. Figure 12 shows the number of inputs required to complete the task against the number of split fields required by the task, for the following four example-provision strategies:

Baseline Split position examples are provided randomly until the splitting is correct. For each task, we average the number of examples required over 50 different random example orderings. This models the baseline where the system does not ask any questions.

BinaryQ One random example is provided, after which the system keeps asking binary position questions until the correct program is achieved. This strategy is purely system-driven because the user does not provide any examples after the first. She only answers the questions posed by the system.

ConfirmationQ One random example is provided by the user, after which the system poses a confirmation question if one exists. The user then provides another example, and we continue alternating between a user example and a system question until the correct splitting is achieved. This strategy is more evenly balanced between user inputs and system feedback.

CombinedQ The system uses a combination of binary and confirmation questions, using the disambiguation score ds_{F5} from §6.2.

Strategy	Avg. number of inputs
Baseline	8.81
BinaryQ	8.54
ConfirmationQ	6.98
CombinedQ	6.90

In general we see significant improvement with feedback-driven strategies over the baseline, which becomes more drastic with more

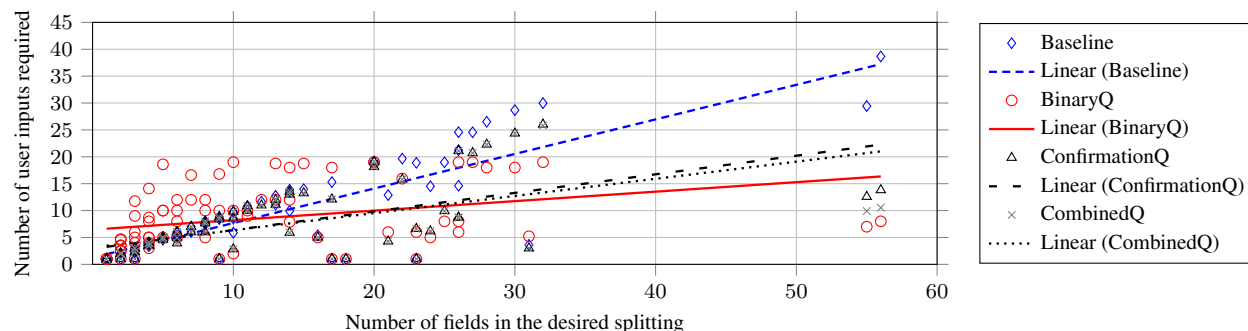


Figure 12: Comparison of effectiveness of different example provision strategies for field splitting tasks from log files. Lower is better.

Baseline	Feedback						
	1	2	3	4	5	6	7
1	241	50	16	0	1	0	0
2	75	30	4	0	0	0	0
3	20	7	2	0	0	0	0
4	0	5	3	0	0	0	0
5	0	1	0	0	0	0	0
6	1	1	0	0	0	0	0

Table 1: Number of rows inspected in the baseline and the feedback-driven settings for FlashFill evaluation.

fields. Although BinaryQ performs much better than the baseline on larger splittings, it performs worse on smaller ones. This is because for smaller splittings, the number of ambiguous programs is much larger than the number of examples requires to perform the task. ConfirmationQ performs better on both small and large splittings. This is because it balances user effort and system feedback: the user’s examples helps reduce the search space while the system’s questions eliminate redundant examples. Finally, with CombinedQ we see further (albeit moderate) improvement on average over the other strategies, illustrating the benefits of a system that uses different question types invoked under different circumstances.

8. Related Work

Conversational Clarification In 2015, Mayer et al. studied different user interaction models that can be applied in PBE to increase the user’s confidence in the learned program and reducing the number of iterations until convergence to the correct program [10]. They compared three interaction models: providing additional examples (positive or negative), presenting a set of candidate programs using an English paraphrasing, and *conversational clarification*, a model that prompts the user with disambiguating questions on a discrepancy between two top-ranked candidate programs. Among them, conversational clarification vastly outperformed the other interaction models in both convergence speed and the users’ confidence.

Mayer et al. established that interaction (in their case, disambiguating questions) in the key to building a user-friendly mass-market PBE system. These results inspired us to give interactive learning first-class treatment in the PBE formalism. In this work, we explore various important dimensions of interactive program synthesis, such as performance of synthesis iterations, impact of different kinds of clarifying questions, and a comprehensive formalism for a step-based synthesis problem in a compound DSL.

Oracle-Guided Inductive Synthesis Jha and Seshia recently developed a novel formalism for inductive synthesis called *oracle-guided inductive synthesis* (OGIS) [5]. It unifies several commonly used approaches such as counterexample-guided inductive synthesis (CEGIS) [15] and distinguishing inputs [7]. In OGIS, an induc-

tive learning engine is parameterized with a concept class (the set of possible programs), and it learns a concept from the class that satisfies a given partial specification by issuing queries to a given oracle. The oracle has access to the complete specification, and is parameterized with the types of queries it is able to answer. Queries and responses range over a finite set of types, including membership queries, witnesses, counterexamples, verification queries, and distinguishing inputs. They also present a theoretical analysis for the CEGIS learner variants, establishing relations between concept classes recognizable by learners under various constraints.

The problem of interactive program synthesis, presented in this work, can be mapped to the OGIS formalism (with the end user playing the role of an oracle). Hence, any theoretical results established by Jha and Seshia for CEGIS automatically hold for the settings of interactive program synthesis where we only issue counterexample queries.

In addition, inspired by our study of mass-market deployment of PBE-based systems, we present further formalism for the user’s interaction with the synthesis system. While the “learner” component in the OGIS formalism is limited to a pre-defined class of queries, our formalism adds a separate modal “hypothesizer” component. Its job is to analyze the current set of candidate programs and to ask the questions that best resolve the ambiguity between programs in the set. The hypothesizer is domain-specific, not learner-specific, and therefore can be refactored out of the learner and reused with different synthesis strategies.

Active Learning In machine learning, *active learning* is a sub-field of semi-supervised learning where the learning algorithm is permitted to issue queries to an oracle (usually a human) [13]. As applied to, e.g., classification problems, a typical query asks for a classification label on a new data point. The goal is to achieve maximum classification accuracy with a minimum number of queries. Research in active learning has focused on two important problems: (a) when to issue a query, and (b) how to pick a data point for it.

This work borrows the ideas of active learning, extends them, and applies them in the domain of program synthesis. In our setting, the issued queries do not necessarily ask for the exact output of the desired program on a given input (an equivalent of “label” in ML), but may also ask for weaker output properties (e.g. verify a candidate element of the output sequence). In all cases, though, our queries are *actionable*: they are convertible to constraints, which automatically trigger a new iteration of synthesis. We also develop a novel approach for picking an input for the query based on its *ambiguity measure* w.r.t. the current set of candidate programs.

9. Conclusion

The standard user interaction model in PBE is for the user to provide constraints in an iterative manner until the user is satisfied with the synthesized program or its behavior on the various inputs. However, this process is far from being interactive:

- The constraints are over the behavior of the entire program. In this paper, we motivated and formalized the notion of associating constraints with sub-expressions of the program.
- The synthesizer is re-run from scratch with the new set of constraints. In this paper, we discussed how to make the synthesizer incremental, leading to a snappier UI experience for the user.
- The refinement of the constraints is a manual process that is guided by the user without any feedback from the synthesizer. In this paper, we discuss various useful feedback mechanisms.

These foundational extensions help address two key challenges of performance and correctness associated with PBE.

References

- [1] R. Alur, R. Bodik, G. Juniwal, M. M. Martin, M. Raghothaman, S. A. Seshia, R. Singh, A. Solar-Lezama, E. Torlak, and A. Udupa. Syntax-guided synthesis. In *FMCAD*, pages 1–17. IEEE, 2013.
- [2] J. K. Feser, S. Chaudhuri, and I. Dillig. Synthesizing data structure transformations from input-output examples. In *ACM SIGPLAN Notices*, volume 50, pages 229–239. ACM, 2015.
- [3] S. Gulwani. Dimensions in program synthesis. In *PPDP*, pages 13–24. ACM, 2010.
- [4] S. Gulwani. Automating string processing in spreadsheets using input-output examples. In *POPL*, volume 46, pages 317–330, 2011.
- [5] S. Jha and S. A. Seshia. A Theory of Formal Synthesis via Inductive Learning. *ArXiv e-prints*, May 2015.
- [6] S. Jha, S. Gulwani, S. A. Seshia, and A. Tiwari. Oracle-guided component-based program synthesis. In *ICSE*, volume 1, pages 215–224. IEEE, 2010.
- [7] S. Jha, S. Gulwani, S. A. Seshia, and A. Tiwari. Oracle-guided component-based program synthesis. In *ICSE*, volume 1, pages 215–224. IEEE, 2010.
- [8] D. Kini and S. Gulwani. FlashNormalize: Programming by examples for text normalization. In *IJCAI*, pages 776–783, 2015.
- [9] V. Le and S. Gulwani. FlashExtract: A framework for data extraction by examples. In *PLDI*, page 55. ACM, 2014.
- [10] M. Mayer, G. Soares, M. Grechkin, V. Le, M. Marron, O. Polozov, R. Singh, B. Zorn, and S. Gulwani. User interaction models for disambiguation in programming by example. In *UIST*, 2015.
- [11] S. Owen. Walkthrough Part Two: Advanced Parsing with ConvertFrom-String. <https://foxdeploy.com/2015/01/13/walkthrough-part-two-advanced-parsing-with-convertfrom-string/>, 2015. [Online; accessed July 4, 2016].
- [12] O. Polozov and S. Gulwani. FlashMeta: A framework for inductive program synthesis. In *OOPSLA*, pages 107–126, 2015.
- [13] B. Settles. Active learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 6(1):1–114, 2012.
- [14] R. Singh and S. Gulwani. Predicting a correct program in programming by example. *CAV*, 2015.
- [15] A. Solar-Lezama. *Program synthesis by sketching*. ProQuest, 2008.
- [16] J. Walkenbach. The research behind Flash Fill. http://spreadsheetpage.com/index.php/blog/the_research_behind_flash_fill/, 2012. [Online; accessed July 4, 2016].
- [17] B. Wu and C. A. Knoblock. An iterative approach to synthesize data transformation programs. In *IJCAI*, 2015.