
Learning Natural Programs from a Few Examples in Real-Time

Nagarajan Natarajan¹ Danny Simmons² Naren Datha¹ Prateek Jain¹ Sumit Gulwani²

¹Microsoft Research, India ²Microsoft Corporation, Redmond

Abstract

Programming by examples (PBE) is a rapidly growing subfield of AI, that aims to synthesize user-intended programs using input-output examples from the task. As users can provide only a few I/O examples, capturing user-intent accurately and ranking user-intended programs over other programs is challenging even in the simplest of the domains. Commercially deployed PBE systems often require years of engineering effort and domain expertise to devise ranking heuristics for real-time synthesis of accurate programs. But such heuristics may not cater to new domains, or even to a different segment of users from the same domain. In this work, we develop a novel, real-time, ML-based program ranking algorithm that enables synthesis of natural, user-intended, personalized programs. We make two key technical contributions: 1) a new technique to embed programs in a vector space making them amenable to ML-formulations, 2) a novel formulation that interleaves program search with ranking, enabling real-time synthesis of accurate user-intended programs. We implement our solution in the state-of-the-art PROSE framework. The proposed approach learns the intended program with just *one* I/O example in a variety of real-world string/date/number manipulation tasks, and outperforms state-of-the-art neural synthesis methods along multiple metrics.

1 Introduction

Programming by examples (PBE) is an important and emerging subfield of AI (Parisotto et al., 2016; Balog et al., 2017; Devlin et al., 2017; Bunel et al., 2018; Kalyan et al., 2018), where a user-intended program is synthesized automatically with the help of a few input-output examples

Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics (AISTATS) 2019, Naha, Okinawa, Japan. PMLR: Volume 89. Copyright 2019 by the author(s).

Input	Output
Missing page numbers, 1993	1993
64-67, 1995	?
2002 (1-27)	?

Table 1: I/O spec provided to a PBE system. Goal is to find a program that is: a) *consistent* (maps the first input into the corresponding output), b) *generalizable* or accurate (computes desired output on last two inputs). While millions of programs in the DSL in Figure 1 are *consistent*, only a handful of them *generalize* well to *unseen* inputs.

(I/O specification, or spec for short). A large fraction of computer users are not experts in programming, and synthesizing programs automatically enables them to be more productive. Table 1 describes a typical PBE task.

PBE is essentially a needle-in-haystack problem where the goal is to search for a *consistent* program (i.e. one that satisfies given I/O spec) in a certain *Domain Specific Language* (DSL) that might contain infinitely many programs. The problem becomes significantly more difficult due to user-centric focus of the systems — the PBE system has to be real-time and should be able to synthesize non-trivial programs; and often in under-specified situations as one cannot expect a user to provide a large number of I/O examples. Unfortunately, these requirements are somewhat contradictory. That is, if the DSL is rich and can support complicated programs, then a small number of I/O examples might not be able to uniquely identify a program in the DSL. For example, for the specification in Table 1, we can generate several consistent programs in the DSL of Figure 1, such as “extract the first number” or “extract the last token”. However, human programmers are typically able to figure out the correct program using a few I/O examples. So, the key question is: *can we synthesize rich user-intended programs using a small number of I/O examples, in real-time?*

Starting with the FlashFill PBE system (Gulwani, 2011) that was commercially deployed in MS Excel (PCWorld, 2012), there has been tremendous progress in this domain over the past few years. Typical PBE techniques search for a program in a carefully-designed DSL and can be categorized into: a) symbolic deduction based techniques (Polo-

zov and Gulwani, 2015; Gulwani et al., 2017; Alur et al., 2017; Le et al., 2017), b) neural computation based techniques (Parisotto et al., 2016; Balog et al., 2017; Devlin et al., 2017; Bunel et al., 2018; Kalyan et al., 2018).

Most neural synthesis (Parisotto et al., 2016; Balog et al., 2017; Devlin et al., 2017; Bunel et al., 2018) models are trained on synthetic data and hence in general, do not capture user-intended programs with a small number of I/O examples. In contrast, symbolic computation based PBE systems handcode the structure of programs and domain knowledge tightly leading to significantly more accurate programs in certain cases. However, manual engineering of the system makes it challenging to extend the solution for even slightly different scenario or a new domain.

Our work alleviates concerns with both the approaches by carefully combining ML techniques with the symbolic search techniques well-understood by the PL community.

(1) Our first contribution addresses the fundamental question of embedding heterogeneous programs/expressions in a vector-space which can make programs more amenable to standard learning techniques (Section 3). In the context of PBE, a few key learning tasks that are enabled by program embeddings are: clustering programs/expressions (Padhi et al., 2018), predicting correct programs (Singh and Gulwani, 2015; Ellis and Gulwani, 2017), and ranking programs (Polozov and Gulwani, 2015; Balog et al., 2017). Defining program embedding is challenging because programs are inherently *recursive* and can be composed of heterogeneous sub-expressions. Furthermore, semantically different programs can often behave equivalently on a given I/O spec, so the embedding should take I/O spec into account as well.

(2) We show how the proposed embedding can be leveraged for learning to rank programs, a crucial component of PBE systems. However, we cannot apply standard ranking techniques — we cannot even enumerate all the candidate programs to rank as there can be millions of consistent programs. So, we need to interleave synthesis and ranking for real-time synthesis, which in turn requires comparing heterogeneous programs, subprograms, expressions, etc. The problem is further complicated by unavailability of supervision for such intermediate subprograms, and by biased training data that the bootstrapping process induces. We propose three novel and successively refined formulations to address the above mentioned challenges (Section 4).

(3) Finally, we integrate our ranking solution with the state-of-the-art PROgram Synthesis using Examples, PROSE (2015) framework. In particular, we show that on real-world data wrangling tasks, the proposed ranking approach outperforms baselines, as well as state-of-the-art neural-synthesis approaches significantly. Our solution is competitive wrt. the ranker tuned over two expert-years that currently ships in Microsoft products (MS Excel, PowerShell, Azure ML). Using just one I/O example, our method

```
@start program := tr | If(cond) Then(tr)
                    Else(program);
bool cond := Matches(input, r);
string tr := atom | Concat(atom, tr);
string atom := ConstStr(s) | let string x
                    : input in SubStr(x, pp) | input;
Tuple<int, int> pp := Pair(pos, pos) |
                    RegexOccurrence(x, r, k);
int pos := AbsPos(x, k);
@input string input; string s; int k;
Regex r; //Terminals
```

Figure 1: An illustrative subset of the FlashFill DSL (Gulwani, 2011). A program takes a string *input*, and returns a string, a concatenation of *atoms*. The operators are self-explanatory. See Appendix B for the full DSL.

synthesizes a desired program for about 67% of the tasks while baselines are successful only in at most 44%.

2 Background

In this section, we define the PBE problem formally, introduce various aspects of PBE systems and terminology/nomenclature used in the rest of the paper.

The goal of a PBE system is to generate *user-intended* program(s) where the user intention is specified using input-output examples (I/O spec): $\zeta = \{\sigma_i \mapsto \psi_i\}_{i=1}^m \cup \{\sigma_i\}_{i=m+1}^n$. σ_i is the i -th example’s input and ψ_i is the corresponding output (when available). Unlabeled inputs are often available and can be used for doing simple validation checks on synthesized programs (See Remark 1).

Typically, PBE systems restrict the search for a program to a domain specific language \mathcal{L} that is powerful enough for solving critical tasks in a certain domain, but is still restrictive and structured enough for efficient program synthesis. A DSL \mathcal{L} is represented as a context-free grammar (CFG) consisting of *terminal* symbols T , *non-terminal* symbols N , *rules* that govern how non-terminals are expanded, and *operators* $F : (N \cup T)^* \rightarrow N$ that make the production rules. As an example, consider the popular *FlashFill* DSL meant for data wrangling tasks in spreadsheets (Gulwani, 2011; Polozov and Gulwani, 2015). The core DSL is captured in Table 1 (Appendix B has the full DSL).

A program or an expression $\mathcal{L} \ni P : \sigma \rightarrow \psi$ is a structured entity with precise syntax and semantics defined by the DSL.

Remark 1 (Unlabeled inputs). *Using unlabeled inputs (i.e. $\{\sigma_i\}_{i=m+1}^n$) can be often helpful in characterizing program behavior; for example, programs that map many of the unlabeled inputs to nulls or empty strings can be indicative of unintended behavior.*

For a PBE system to be usable in an interactive setting, it

should satisfy three key requirements:

- (R1) be *consistent* (see Definition 1), i.e., return program(s) that satisfy the user-provided I/O spec,
- (R2) be *generalizable*, i.e., the synthesized program(s) should give desired output on *unseen* inputs; for severely underspecified problems (say $m = 1$ I/O example) there can be millions of consistent programs (see Table 1), and
- (R3) be *real-time*, i.e., the synthesized generalizable programs on consumer-class devices.

Definition 1 (Consistent Program). *A program $P \in \mathcal{L}$ is “consistent” on a given input-output specification $\{\sigma_i, \psi_i\}_{i=1}^m$, if $P(\sigma_i) = \psi_i$, for $i = 1, 2, \dots, m$. Otherwise, P is inconsistent.*

While consistency (R1) is essentially a search problem, (R2) is more critical and interesting from a machine learning perspective — often there can be millions of programs that satisfy (R1), but the user would find most of the consistent programs unusable because they do not generalize to new inputs. It is not possible to formally specify “naturalness” of programs with symbolic logic. Typically, (R2) is addressed by means of a ranking function that can help choose the “best” program from possibly many consistent programs. One way to address this is to first synthesize all the consistent programs, and *then* rank them (Ellis and Gulwani, 2017). Unfortunately, the naive approach cannot be done in real-time — it can take hours to even enumerate the consistent programs, thus contradicting (R3). State-of-the-art neural-network based synthesis approaches are trained on synthetic datasets/programs, so they fail to capture the structure in the domain. As a result, neural synthesis approaches suffer in the quality of synthesized programs, especially for underspecified synthesis tasks (See Section 5).

It is therefore crucial to look at the search and the ranking problem as a whole (i.e., (R1)-(R3)). Successful, commercially-deployed PBE systems (Gulwani, 2011; Gulwani et al., 2015; Alur et al., 2013) use symbolic logic and deductive synthesis techniques to efficiently address (R1) and (R3). In particular, the symbolic PBE systems use a top-down deductive synthesis strategy based on the divide-and-conquer paradigm. Here, the search problem for a given I/O spec is reduced into smaller subproblems with suitably modified specs¹. For e.g., the synthesis problem $\zeta = \{\text{“New York”} \mapsto \text{“NY”}\}$ is broken down into finding a set of subprograms \mathcal{P}_1 with spec $\zeta_1 = \{\text{“New York”} \mapsto \text{“N”}\}$ and a set of subprograms \mathcal{P}_2 with spec $\zeta_2 = \{\text{“New York”} \mapsto \text{“Y”}\}$, i.e., programs in \mathcal{P}_1 generating “N” and those in \mathcal{P}_2 generating “Y”. Then the final program set is given by $\mathcal{P} = \{\text{Concat}(P_1, P_2) \text{ s.t. } P_1 \in \mathcal{P}_1, P_2 \in \mathcal{P}_2\}$. Each of the synthesis subproblems is solved recursively using the same strategy.

¹it is beyond the scope of the paper to describe how spec for the subproblems are obtained. See Polozov and Gulwani (2015) for details of search. The key idea is to leverage *inverse semantics* of the involved operators.

However, the aforementioned PBE systems rely on heuristics for (R2), i.e. ranking (Polozov and Gulwani, 2015; Rolim et al., 2017; Wang et al., 2017) (such as choosing smaller programs/expressions over larger ones). Simple heuristics may result in bad failures even in simple cases. For illustration, consider the data formatting task with just one I/O example: $\{\text{“[CCC-0001”} \mapsto \text{“[CCC-0001”}\}$. Adopting naive heuristics such as “prefer programs with fewer constants” or “prefer shorter programs” leads to the incorrect program: `Concat(input, ConstStr("]"))`, which would fail on an already formatted input, say “[CCC-002]”. On the other hand, developing carefully-tuned ranking heuristics often takes one to two expert-years; and requires continual effort to keep up with domain changes, let alone scaling to new domains. Also, it can be challenging to *personalize* the heuristics to user segments with unique biases/preferences.

The primary goal of our work is to develop an ML-based ranking solution for real-time synthesis of natural programs. Programs are difficult objects to analyse/rank, so we need to be able to embed them in a suitable feature space. To this end, we first address the problem of embedding *heterogeneous* programs/expressions in a common vector space. Defining an embedding that handles the heterogeneity is non-trivial, and it turns out that we need to *learn* the embeddings themselves. Existing embedding techniques (Ellis and Gulwani, 2017) do not work because they are defined for homogeneous programs. We address the embedding challenges and our solution in Section 3. Subsequently, we consider the problem of doing program ranking and search jointly. A priori, it is unclear how to set up/formulate the machine learning problem, or what loss function to optimize. Ranking programs/expressions is challenging for multiple reasons: 1) classical ranking techniques (Liu et al., 2009) do not work, as we do not even have a clean supervised dataset to begin with, and 2) search for user-intended programs is a sequential decision making problem, therefore a mistake at any point in the search may be irrevocable; this necessitates a novel ranking formulation that admits *interleaved* search and ranking during synthesis. We address these challenges and propose ranking solutions in Section 4.

3 Program-Spec Embedding

Informally, the problem is to find a representation for programs/expressions $P \in \mathcal{L}$ together with the I/O spec ζ , such that the embedding captures syntactic and semantic structure (defined by DSL), as well as behavioral properties (defined by I/O spec). Defining a feature vector for programs/expressions that captures the complex structure/properties is not obvious. Simple techniques like using the abstract syntax tree (AST) directly do not suffice. Programs with very similar ASTs can differ arbitrarily in their semantics. Consider two programs from the *FlashFill*

DSL for the task in Table 1, $P_1 = \text{let } x : \text{input in SubStr}(x, \text{RegexOccurrence}(x, \text{"Number"}, 1))$ and $P_2 = \text{let } x : \text{input in SubStr}(x, \text{RegexOccurrence}(x, \text{"Number"}, -1))$; P_1 and P_2 have identical ASTs but different semantics (extracting the first number vs the last number in the input). On the other hand, two programs with very different ASTs can produce identical outputs on given inputs.

(1) It is crucial to embed I/O spec along with the program/expression itself. The utility of a program can vary drastically based on the I/O spec. For e.g., the program $P = \text{let } x : \text{input in SubStr}(x, \text{Pair}(1, 3))$ has the outcome of extracting first three digits of SSN in $\zeta_1 = \{\text{"123-45-6789"} \mapsto \text{"123"}, \text{"555-21-9012"} \mapsto \text{"555"}\}$ vs an undesirable outcome of extracting first three letters of name in $\zeta_2 = \{\text{"Joe Smith"} \mapsto \text{"Joe"}\}$. So the embedding must be defined on the tuple (P, ζ) rather than just P .

(2) The embedding should facilitate comparisons between expressions and programs of different sizes, types and complexities. For e.g., we want the expressions $\text{Concat}(\text{Concat}(\text{ConstStr}(\text{"@"}), \text{ConstStr}(\text{"gmail"})), \text{ConstStr}(\text{".com"}))$ and $\text{ConstStr}(\text{"@gmail.com"})$ to yield similar representations. This is highly non-trivial; existing embedding techniques do not impose/satisfy such a requirement.

(3) Programs are compositional, e.g. $\text{Concat}(\text{Concat}(P_1, P_2), P_3)$. We want the embedding to be *recursive*, thereby preserving the compositional structure. The embedding of a program should respect and conform to the embeddings of its subprograms/expressions.

Often domain knowledge can help us define features for individual operators in the DSL. Concretely, let $\Phi_{Op}(P, \zeta) \in \mathbb{R}^{d_{Op}}$ be the set of given d_{Op} features for an operator Op . For e.g., for the `Concat` operator, the length of its prefix string argument is a feature (note that the feature may depend on the spec ζ). See Appendix B.3 for features in FlashFill DSL.

Define the dimensionality d to be $d = \sum_{Op \in CFG(\mathcal{L})} d_{Op}$.

Definition 2 (Program-Spec embedding). *For any given program/expression $P \in \mathcal{L}$, operator features $\Phi_{Op} \in \mathbb{R}^{d_{Op}}$ for all operators $Op \in CFG(\mathcal{L})$, and I/O spec ζ , we want an embedding $\Phi(P; \zeta) \in \mathbb{R}^d$ that satisfies the aforementioned three requirements.*

To handle the recursive nature of programs (in the requirement (3) above), and the grammar itself, we critically exploit the fact that \mathcal{L} is represented as an unambiguous grammar that has a unique parse $\mathcal{T}(P)$ for P . Let $Op(P)$ be the operator at the top of $\mathcal{T}(P)$, and let $\mathcal{C}(P)$ denote the immediate children nodes of P in $\mathcal{T}(P)$. We obtain embedding for P by combining the given features for the top operator in $\mathcal{T}(P)$ with a weighted combination of embeddings of each child node of P in $\mathcal{T}(P)$. We define embedding

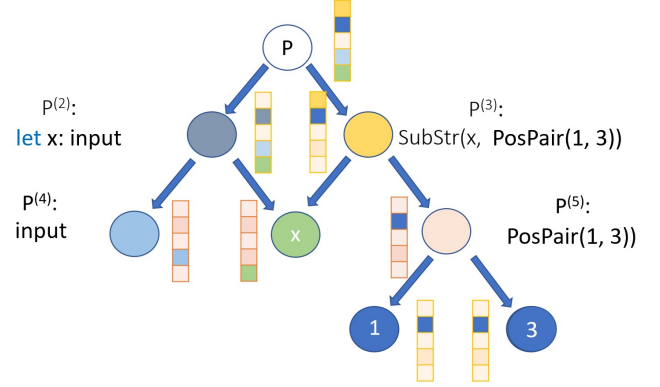


Figure 2: Parse-tree and embedding for the program: $\text{let } x : \text{input in SubStr}(x, \text{PosPair}(1, 3))$. The types of nodes (variables/operators) are color-coded.

$\Phi(P)$ of P recursively as:

$$\Phi(P; \zeta) = \Phi_{Op(P)}(P, \zeta; w) + \sum_{P' \in \mathcal{C}(P)} w(P') \Phi(P', \zeta_{P'}; w), \quad (1)$$

where $\Phi_{Op(P)}(P; \zeta)$ are the given features for the root operator $Op(P)$ of P , $\zeta_{P'}$ is the spec for subprogram P' defined as $\{\sigma_i \mapsto P'(\sigma_i) \mid \sigma_i \in \zeta\}$, and $w(P')$ is the weight assigned to the operator at child P' , i.e. $w(P') := w(Op(P'))$, in the parse tree of P (see Figure 2). Thus, in addition to the given features, the embeddings are characterized by children operator weights $w(P')$ as well, i.e., $w \in \mathbb{R}^{|\mathcal{L}_{Op}|}$ where $|\mathcal{L}_{Op}|$ is the number of operators in \mathcal{L} .

Remark 2. *Note that although the definition of the embedding is recursive, we can compute it once weights w are fixed. Observe that the leaf nodes in $\mathcal{T}(P)$ have only the given features and hence the embeddings are well-defined and immediately obtained; thus, the embedding for the program P can be computed efficiently in a bottom-up fashion.*

Thus we have a homogeneous embedding $\Phi(P, \zeta; w)$ of the program P in the same d -dimensional space as that of its constituent expressions. The weights $w(\cdot)$ can be learned based on the end task that the embedding will be used for, addressed in Section 4.

Remark 3. *Our program embedding technique is also an independent technical contribution, as it enables key learning tasks such as clustering programs/expressions (Padhi et al., 2018), predicting correct programs (Singh and Gulwani, 2015; Ellis and Gulwani, 2017), and ranking programs in code-completion task (Balog et al., 2017).*

4 Program Ranking

The goal of program ranking is to learn a ranking function s that provides the highest score to user-intended programs; and to facilitate synthesis of a user-intended program from a few I/O examples. However, as mentioned in Section 2, a

standard approach (Ellis and Gulwani, 2017) of generating all the *consistent* programs and then ranking them using standard formulations (Liu et al., 2009) is not feasible for real-time systems.

Instead, a key motivating observation for our solution is that the search process of the synthesis algorithm partitions the program generation into multiple *smaller* program synthesis sub-problems. So, the ranking algorithm should be able to generate “correct” subprograms for each of the smaller sub-problems as well; we call a program *correct* if it produces the desired output on unseen inputs as well.

That is, say a program $P = Op(P_1, \dots, P_r), \{P_j \in \mathcal{P}_j\}$ is generated for specification ζ with operator Op in the DSL \mathcal{L} . Each $P_j \in \mathcal{P}_j$, $1 \leq j \leq r$, is in turn generated by solving a smaller PBE problem with “refined” specification ζ_j (discussed briefly in Section 2). Now, we require the ranking function s to be such that it not only scores P higher than other programs $P' \in \mathcal{L}$ for specification ζ but it also scores each P_j above other programs $P'_j \in \mathcal{L}$ for specification ζ_j . That is the ranking function is *monotonic*.

Definition 3 (Program ranking). *Let $\zeta = \{\sigma_i \mapsto \psi_i\}_{i=1}^m \cup \{\sigma_i\}_{i=m+1}^n$ denote I/O spec given to the PBE system. We want to learn a ranking function $s : \mathbb{R}^d \rightarrow \mathbb{R}$ as well as the embedding function $\Phi(\cdot, \cdot)$ such that below hold:*

1. **Correctness:** $s(\Phi(P; \zeta)) > s(\Phi(P'; \zeta))$, for **correct programs** $P \in \mathcal{L}$ and **incorrect programs** $P' \in \mathcal{L}$.
2. **Monotonicity:** Let $\mathcal{P}_1, \dots, \mathcal{P}_r$ denote the top- K programs returned for each subproblem with specification ζ_j , $1 \leq j \leq r$ and let the final set of programs be $\mathcal{P} = \{P, s.t. P = Op(P_1, \dots, P_r), P_j \in \mathcal{P}_j\}$. Then, the following holds:

$$\forall P \in \mathcal{P}, \forall P' \in \mathcal{L} \setminus \mathcal{P}, s(\Phi(P; \zeta)) \geq s(\Phi(P'; \zeta)) \Rightarrow \forall P_j \in \mathcal{P}_j, \forall P'_j \in \mathcal{L} \setminus \mathcal{P}_j, s(\Phi(P_j; \zeta_j)) \geq s(\Phi(P'_j; \zeta_j)).$$

To learn the ranking function, we use a benchmark of real-world programming tasks that should capture the typical user-intent. Each task has a set of input-output examples; while we provide a small number of them for synthesizing the program, the remaining I/O examples are used for testing if a synthesized program succeeds on the task. Designing such a function requires further solving the following two key challenges:

- (1) **Biased training data:** Learning a ranking function requires generating data from the PBE system itself (by applying it to a few tasks in the benchmark). To bootstrap and to generate training data, we supply the PBE system with a baseline ranker s_0 (e.g., a ranker that generates random scores); generated training data is used to learn a new ranker s_1 . When we deploy s_1 , the distribution of the subprograms generated itself changes based on s_1 ’s rankings, hence the accuracy can be arbitrarily poor as s_1 was trained on data generated from s_0 .
- (2) **Distant supervision:** Though the ranking function s

Algorithm 1 Algorithm for training ML-PROSE.

```

function ML-PROSE( $\mathcal{L}, \theta_0, T = \{\zeta_i, i \in [|T|]\}, \Gamma$ )
1:  $w(P')_0 = 1$  for all  $P' \in CFG(\mathcal{L})$ 
2: for all  $0 \leq \tau \leq \Gamma$  do
3:    $\mathcal{P}_j = Synthesis(h_\tau, \zeta_j), 1 \leq j \leq |T|$ , Synthesized programs by applying  $s_\tau$  to spec  $\zeta_j$ 
4:   Assign  $y_P = 1$  for each correct  $P \in \mathcal{P}_j, \forall j$ 
5:   Assign  $y_P = -1$  for each incorrect  $P \in \mathcal{P}_j, \forall j$ 
6:    $\theta = \theta_\tau, w = w_\tau$ 
7:   while not converged do
8:     Compute  $\Phi(P, \zeta_j; w_\tau)$  using (1),  $P \in \mathcal{P}_j, \forall j$ 
9:     Update  $\theta$  by solving (3) with fixed  $w$ 
10:    Update  $w$  by solving (3) with fixed  $\theta$  and  $s(P, \zeta; w)$  computed recursively using (2)
11:     $w_{\tau+1} = w, \theta_{\tau+1} = \theta$ 
12: return  $s_\Gamma = (w_\Gamma, \theta_\Gamma)$ 

```

is applied to rank smaller subprograms as well as the final programs, the feedback (correctness label) is available only for final programs; i.e., we can apply the final set of synthesized programs on unseen inputs to measure their accuracy, but we cannot get similar feedback for their subprograms.

4.1 Learning to Rank Programs/Subprograms

In this section, we describe three methods to generate ranking problems; successive methods capture the problem structure better and try to address the above mentioned challenges more directly. In this work, we focus on linear scoring functions over the embedding Φ (parameterized by w) defined in (1) (see Remark 4 for discussion on non-linear functions), i.e., the score $s(P, \zeta; w)$ for program P with spec ζ is given by: $s(P, \zeta; w) := \theta^T \Phi(P, \zeta; w)$, θ are the weights. If $P = Op(P_1, \dots, P_r)$ then,

$$s(P, \zeta; w) = \theta_{Op}^T \Phi_{Op}(P, \zeta; w) + \sum_j w(P_j) s(P_j, \zeta_j; w), \quad (2)$$

is a recursive scoring (ranking) function as desired, where θ_{Op} is the projection of θ onto given features for operator Op . Note that $w(P_j) \geq 0$ is a necessary condition for satisfying monotonicity (Definition 3). Now, we want to learn weights $w(\cdot) \geq 0$ as well as θ in (2) such that the ranking problem in Definition 3 is feasible and can be solved accurately. For a DSL \mathcal{L} , let $T = \{\zeta^1, \zeta^2, \dots, \zeta^{|T|}\}$ denote a set of tasks, where each task corresponds to an I/O spec $\zeta^\tau = \{\sigma_i^\tau \mapsto \psi_i^\tau\}_{i=1}^{m_\tau} \cup \{\sigma_i^\tau\}_{i=m_\tau+1}^{n_\tau}$. For task ζ^τ , let \mathcal{P}_{ζ^τ} denote the set of programs synthesized. It is *always* possible to generate at least one correct program for *offline* training tasks by providing sufficiently many I/O examples (as search returns only consistent programs). Note that correctness of a program (if it produces the desired output on all unseen inputs as well) can be easily determined

for *training data*. Let $y(P) = 1$ if $P \in P_{\zeta^\tau}$ is correct for task ζ^τ , else $y(P) = -1$.

(I) Basic formulation (ML-PROSE): In the first formulation, we avoid the challenges mentioned in the previous section by starting with a random ranker and by comparing only the final programs. That is, the goal is to learn a scoring function that ranks any correct program above all incorrect programs, i.e. $\theta^T \Phi(P_a, \zeta; w) > \theta^T \Phi(P_b, \zeta; w)$, for programs $P_a, P_b \in \mathcal{P}(\zeta^\tau)$ generated for a task τ , such that $y(P_a) = 1$ and $y(P_b) = -1$. More generally, we want to penalize the difference between their scores using a suitable loss function ℓ . The corresponding optimization problem is written as:

$$\begin{aligned} \min_{\theta, w} \quad & \sum_{\tau=1}^{|T|} \sum_{\substack{P \in \mathcal{P}_{\zeta^\tau}, \\ y(P)=1}} \sum_{\substack{P' \in \mathcal{P}_{\zeta^\tau}, \\ y(P')=-1}} \ell(s(P, \zeta^\tau; w) - s(P', \zeta^\tau; w)) \\ & + C_1 \|\theta\|_2^2 + C_2 \sum_{Op \in CFG(\mathcal{L})} w_{Op}^2, \quad \text{s.t. } w_{Op} \geq 0, \forall Op, \end{aligned} \quad (3)$$

where $s(\cdot)$ is defined in Equation (2) and the loss function $\ell(a)$ penalizes negative a ; we use standard hinge loss for ℓ in our experiments. We solve the above given problem by alternating over θ and w ; note that each of the sub-problems for θ and w is individually convex and easy to optimize.

Remark 4 (Non-linearity). *We can capture non-linearity in the ranking model by generating polynomial features for the local features Φ_{Op} . This enables learning complex scoring functions like the one in Figure 4 (See Appendix).*

(II) Handling distant supervision (ML-PROSE-SubPRG): The above formulation ignores the fact that subprograms are generated by solving smaller synthesis problems. So, even if the scoring function s is accurate for final programs, it can be arbitrarily poor for the subprograms. We alleviate this issue partially by sampling subprograms in the training data to solve Problem (3). We use a baseline ranker to generate both the final programs as well as the subprograms and include a sample from the subprograms in Problem (3). We address the issue of distant supervision by fixing ‘‘correctness’’ of a subprogram P' as follows: $y(P') = 1$ if P' appears as part of at least one correct program for a given task, or else we assign $y(P') = -1$. Table 2 clearly shows that the ranking function can be improved significantly by inclusion of subprograms when solving (3).

(III) MinMax formulation: Problem (3) does not directly address the crucial requirement of deductive program synthesis (even if we include sampled subprograms as in ML-PROSE-SubPRG) — we want *all* the subprograms of a *given* correct program to be ranked correctly during synthesis. Furthermore, it suffices to rank *any one* correct program above *all* incorrect programs for a given task. For the

subprograms of given correct program P for task ζ , i.e., $P_j \in \mathcal{T}(P)$, let $\zeta_j, j = 1, 2, \dots, |\mathcal{T}(P)|$ denote their respective subproblem specification. Let \mathcal{P}_{ζ_j} denote the set of all programs in \mathcal{L} that satisfy ζ_j , for each j . Of course, it is impossible to enumerate the entire set, but we can sample many such subprograms for each subproblem specification. We determine the label for the subprograms in \mathcal{P}_{ζ_j} as before (+1 if the subprogram is part of at least one correct program for the task, or else -1). Define the loss on a correct program P as the max over the losses of all the comparisons during its synthesis:

$$\Delta(P) = \max_{P_j \in \mathcal{T}(P)} \max_{\substack{P'_j \in \mathcal{P}_{\zeta_j} \\ y(P'_j)=-1}} \ell(s(P_j, \zeta_j; w) - s(P'_j, \zeta_j; w)). \quad (4)$$

We solve:

$$\min_{\theta, w \geq 0} \sum_{\tau=1}^{|T|} \min_{\substack{P \in \mathcal{P}_{\zeta^\tau}, \\ y(P)=1}} \Delta(P) + C_1 \|\theta\|_2^2 + C_2 \sum_{Op} w_{Op}^2. \quad (5)$$

The above optimization problem is non-convex even in θ , however, we can still define sub-gradient for the problem. In particular, we implement stochastic sub-gradient descent method for this problem using the widely-used Tensorflow framework (www.tensorflow.org).

4.2 Iterative Training

The above formulations still do not address the biased training data challenge. In fact, even if we have a good ranker s_0 to bootstrap with, the bias of baseline ranker still persists. To alleviate this concern, we use an iterative scheme to ensure that the train-test distribution for our ranking function matches while we improve the ranking function itself. Using a base ranker we synthesize programs for the training tasks, sample programs and solve problem (3) (or (5), for training the MinMax model). We then deploy the learned ranker in the PBE system, synthesize (possibly different set of) programs for the tasks, sample programs afresh again to re-learn the ranking model, and repeat. The iterative procedure is described in Algorithm 1 (and in Figure 3 of Appendix) and is able to handle the biased training data issue effectively. In order to ensure *smooth* refinement of s , we combine data from a few recent iterations $\mathcal{D}_1 \cup \mathcal{D}_2 \dots \cup \mathcal{D}_\Gamma$; here \mathcal{D}_τ is the training dataset generated using s_τ . This also helps us avoid poor local minima and helps the ranker converge to a reasonable stationary point.

5 Experiments

We have implemented our learning approach in the PROSE (2015) framework, which is the state-of-the-art PBE system for data wrangling tasks, and is publicly available for academic use.

Benchmark tasks. We use 740 real-world string/date/time manipulation tasks obtained from Polozov and Gulwani (2015). Each task in the benchmark consists of a list of input strings and their corresponding outputs (See Appendix C). The available number of I/O examples per task varies from two to a few hundreds. We use 100 tasks for training, and the remaining 640 for testing. Permuting the order of I/O examples in each *training* task, and varying the spec size m , we get several variants of a single training task. Results are reported on the 640 test tasks used as-is from the benchmark.

Performance Metrics. We want the PBE system to get an intended program in top- K . We report results for $K = 1$ (ACC@1) as well as $K = 10$ (ACC@10).

Initial Ranking Model. In our experiments, we use the ranking function that prefers shorter programs as the initial ranking function in Algorithm 1. Natural programs tend to be terse and often short, so this is a reasonable starting point. Here, $s(P; \zeta) = \frac{1}{|\mathcal{C}(P)|} \sum_{P' \in \mathcal{C}(P)} s(P'; \zeta) - 1$.

Training Data. At each iteration of Algorithm 1, we take the top-1000 programs for each task generated with the ranking model of the previous iteration. This ensures we have good mix of correct and incorrect programs to sample from. With 100 tasks (and their variants) in the training set, sampling about 40 correct and 40 incorrect programs from each task results in about 1.2M data points in total per iteration for learning the ML-PROSE model (3). To train the MinMax model, we sample 50 correct programs from each training task, to compute the inner min in (5); on average there are about 20 subprograms per program corresponding to the outer max in (4) and about 20 (incorrect) subprograms corresponding to the inner max in (4); this sampling strategy leads to about 2M training data points. ACC@1 for training tasks flattens after about 5 iterations, as shown in Figure 5 (in Appendix A); so we use the ranking function at the end of 6 iterations to report results on test data.

5.1 Results on FlashFill benchmark

Compared methods. Our three proposed ranking algorithms are (i) ML-PROSE where we use only top-level programs for training, (ii) ML-PROSE-SubPRG where we use both programs as well as subprograms for learning the ranking model in the objective (3), and (iii) MinMax model that uses the more directed objective in (5). We compare our methods against four baseline ranking functions: (i) RANDOM ranking function where each weight $\theta_i \sim \text{Uniform}([-1, 1])$; (ii) the initial ranking model outlined earlier, that prefers shorter programs, which we call SHORTEST-PROGRAM (Wang et al., 2017; Osera and Zdanczewicz, 2015); here, we discard trivial ConstStr programs (which is by definition the shortest program, when the I/O spec has only one example), (iii) a ranking score model that prefers fewer and shorter constants, which we

call FEWER-CONSTANTS; good constants like delimiters tend to be short, so this is a reasonable heuristic; (iv) combining the ranking models of SHORTEST-PROGRAM and FEWER-CONSTANTS (i.e. prefer programs that are short *as well as* with fewer, shorter constants).

Accuracy. The results for accuracy at top-1 and at top-10 for the different methods are presented in Table 2 (columns 1-4). The best performing method in terms of ACC@1 is ML-PROSE-SubPRG, which retrieves the intended program at the top in 67% test tasks, *using just one I/O example*. Note that the hand-designed PROSE ranker (that comes with PROSE (2015) SDK, and is shipped as part of Microsoft products including MS Excel, Powershell, and Azure ML), tuned using the entire benchmark, i.e. training *as well as* test tasks, achieves 0.72 top-1 accuracy with $m = 1$. However, its top-10 accuracies are comparable to ML-PROSE-SubPRG. In terms of ACC@10, the MinMax model is the clear winner, in both $m = 1$ and $m = 2$ cases; this suggests that the $\Delta(P)$ loss (4) effectively captures the synthesis-time “competitions” among potential subprograms. Another important takeaway from the results is that the synthesis problem becomes significantly easier with $m = 2$ compared to $m = 1$. This is evident from observing the lift in performance of all the baseline methods, especially the fourth one.

Synthesis time. Our ranking models are competitive compared to the optimized PROSE ranker in terms of synthesis times (i.e. elapsed CPU time to synthesize top-1 program for a given I/O spec). See Figure 5 in Appendix A.

5.2 Comparison to state-of-the-art ML methods

Two important neural synthesis techniques in PBE context are the **RobustFill** framework (Devlin et al., 2017) and the **DeepCoder** framework (Balog et al., 2017). For fair comparison, we conduct experiments on a simpler DSL that Devlin et al. (2017) use. In particular, we use 73 tasks from the FlashFill benchmark, which is an *exact* subset of our 640 test tasks, on which the results are reported in Kalyan et al. (2018). We summarize the results in Table 1 of Kalyan et al. (2018) as well as present comparisons to our method in Table 3 of Appendix A. We find that even the SHORTEST-PROGRAM baseline achieves 32% ACC@1 with $m = 1$, about 7% better than RobustFill with $m = 1$, on the *exact* 73 tasks. The simple baseline performs reasonably well because, in this subset of tasks, 2 or 3 I/O examples are sufficient for the search strategy to find *consistent programs that also generalize very well*; on the other hand, RobustFill cannot even guarantee consistent programs. Our ranker ML-PROSE-SubPRG performs the best on the 73 tasks, achieving 70% ACC@1 with $m = 1$. We exclude comparisons to Menon et al. (2013) as it requires additional information beyond I/O spec for synthesis.

RANKING METHOD	ACC@1		ACC@10	
	$m = 1$	$m = 2$	$m = 1$	$m = 2$
RANDOM	0.22	0.60	0.38	0.67
(A) SHORTEST PROGRAM	0.37	0.69	0.49	0.80
(B) FEWER CONSTANTS	0.38	0.60	0.59	0.80
(A) and (B)	0.44	0.72	0.60	0.87
ML-PROSE	0.63	0.78	0.73	0.87
ML-PROSE-SubPRG	0.67	0.83	0.75	0.89
MinMax	0.65	0.81	0.79	0.92

Table 2: Performance on the FlashFill benchmark. The number of I/O examples given to the PBE system for each of the 640 test tasks is m . The proposed methods, especially ML-PROSE-SubPRG and MinMax, perform significantly better than the baselines. The expert-designed ranker, currently shipped as part of several Microsoft products, tuned using training as well as test tasks, gets 0.72 ACC@1 with $m = 1$, and 0.85 with $m = 2$; though its ACC@10 is worse than MinMax.

5.3 Personalization

A single ranking function may not cater to all types of users, even within the same domain. A significant advantage of our ranking solution is that we can re-train the scoring model in order to capture the unique biases/preferences for different user segments. For e.g, geography often determines date/time formats; we want the ranking function to prefer the default formatting style for the specific user locale, unless additional I/O examples overrule the assumed preferences. One simple and effective way to capture these biases is to repeat the task, on which the ranker deviates from the desired behavior, multiple times (or equivalently, weigh the loss associated with this task higher). Below, we present two scenarios for personalized ranking.

Rounding Numbers. Say we want to induce the following preference for rounding a number: “Nearest” > “TowardsZero” > “Down”. The preference that our method (using MinMax ranking formulation) learns from the training data is “TowardsZero” > “Nearest” > “Down” (See Figure 7, Appendix D). Learning this preferential order from the randomly sampled training data is likely because in many number transformation tasks where “Nearest” rounding operation applies, “TowardsZero” also leads to correct programs (and “Down” is the least representative rounding operation in the entire benchmark). By replicating three training tasks that induce the preferred rounding behavior 10 times and re-training, the (MinMax) ranking model learns “Nearest” as the most-preferred rounding operation (See Figure 8 and Example 1, Appendix D).

Formatting Dates. In many tasks, the intended output format is ambiguous unless one looks at several I/O examples. Say, some users prefer “m/d” to “M/dd” (2/3 vs 02/03 for 3rd Feb) for date, or “h:mm:ss” to “hh:mm:ss” for time. Our (MinMax) ranker learns a preference towards “mm/dd” and “hh:mm:ss” formats which are representative of the training data. By replicating 2 tasks that induce the desired formatting behavior in the training data and re-

training, the ranking model learns the desired formatting preferences (See Example 2, Appendix D).

Remark 5 (Maintenance and Debugging). *The personalization scenarios above also imply another significant advantage of our ML-based ranking solution over neural synthesis approaches — transparency. It is crucial for an ML-based PBE system to be maintainable and debuggable.*

6 Related Work

As mentioned in Section 1, there are two lines of work on program synthesis, symbolic and ML/neural-synthesis based approaches. For symbolic techniques, Gulwani (2010) and Gulwani et al. (2017) provide extensive surveys. State-of-the-art neural program synthesis techniques have already been mentioned/discussed earlier. See Gulwani and Jain (2017) for recent results that are at the intersection of ML and PL. Statistical learning techniques for PBE have also received some attention. Ellis and Gulwani (2017) try to improve the accuracy of existing PROSE implementations, by learning to re-rank the top K consistent programs for given I/O spec, assuming a “good” ranking function is already in place unlike our approach. The statistical learning framework of Menon et al. (2013) employs a log-linear model for inferring likelihood of consistent programs from a probabilistic CFG. In addition to I/O spec, it also needs “clues” to be able to narrow down the rules to consider for enumeration, so that synthesis time is not prohibitive. Singh and Gulwani (2015) learn a ranking function (that prefers generalizable programs) using only top-level programs but apply the learned function recursively to rank subprograms during synthesis; their method has not been implemented in a PBE system to demonstrate real-time synthesis. Raychev et al. (2016) focus on the synthesis setting where one has access to many, and potentially noisy, I/O examples. Christakopoulou and Kalai (2017) specify intent through a “glass-box” scoring program that evaluates candidate programs; they do not use any I/O spec.

References

- Alur, R., Bodík, R., Juniwal, G., Martin, M. M. K., Raghthaman, M., Seshia, S. A., Singh, R., Solar-Lezama, A., Torlak, E., and Udupa, A. (2013). Syntax-guided synthesis. In *Formal Methods in Computer-Aided Design (FMCAD)*, pages 1–8.
- Alur, R., Radhakrishna, A., and Udupa, A. (2017). Scaling enumerative program synthesis via divide and conquer. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 319–336. Springer.
- Balog, M., Gaunt, A. L., Brockschmidt, M., Nowozin, S., and Tarlow, D. (2017). Deepcoder: Learning to write programs. *International Conference on Learning Representations (ICLR)*.
- Bunel, R., Hausknecht, M., Devlin, J., Singh, R., and Kohli, P. (2018). Leveraging grammar and reinforcement learning for neural program synthesis. In *International Conference on Learning Representations*.
- Christakopoulou, K. and Kalai, A. T. (2017). Glass-box program synthesis: A machine learning approach. *arXiv preprint arXiv:1709.08669*.
- Devlin, J., Uesato, J., Bhupatiraju, S., Singh, R., Mohamed, A.-r., and Kohli, P. (2017). Robustfill: Neural program learning under noisy i/o. In *International Conference on Machine Learning*, pages 990–998.
- Ellis, K. and Gulwani, S. (2017). Learning to learn programs from examples: Going beyond program structure. *IJCAI*.
- Gulwani, S. (2010). Dimensions in program synthesis. In *Proceedings of the 12th international ACM SIGPLAN symposium on Principles and practice of declarative programming*, pages 13–24. ACM.
- Gulwani, S. (2011). Automating string processing in spreadsheets using input-output examples. In *ACM SIGPLAN Notices*, volume 46, pages 317–330. ACM.
- Gulwani, S., Hernández-Orallo, J., Kitzelmann, E., Muggleton, S. H., Schmid, U., and Zorn, B. (2015). Inductive programming meets the real world. *Communications of the ACM*, 58(11):90–99.
- Gulwani, S. and Jain, P. (2017). Programming by examples: PL meets ML. In *Asian Symposium on Programming Languages and Systems*, pages 3–20. Springer.
- Gulwani, S., Polozov, O., and Singh, R. (2017). Program synthesis. *Foundations and Trends® in Programming Languages*, 4(1-2):1–119.
- Kalyan, A., Mohta, A., Polozov, O., Batra, D., Jain, P., and Gulwani, S. (2018). Neural-guided deductive search for real-time program synthesis from examples. *International Conference on Learning Representations (ICLR)*.
- Le, X.-B. D., Chu, D.-H., Lo, D., Le Goues, C., and Visser, W. (2017). S3: syntax-and semantic-guided repair synthesis via programming by examples. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*, pages 593–604. ACM.
- Liu, T.-Y. et al. (2009). Learning to rank for information retrieval. *Foundations and Trends® in Information Retrieval*, 3(3):225–331.
- Menon, A., Tamuz, O., Gulwani, S., Lampson, B., and Kalai, A. (2013). A machine learning framework for programming by example. In *International Conference on Machine Learning*, pages 187–195.
- Osera, P.-M. and Zdancewic, S. (2015). Type-and-example-directed program synthesis. In *ACM SIGPLAN Notices*, volume 50, pages 619–630. ACM.
- Padhi, S., Jain, P., Perelman, D., Polozov, O., Gulwani, S., and Millstein, T. D. (2018). Flashprofile: a framework for synthesizing data profiles. *PACMPL*, 2(OOPSLA):150:1–150:28.
- Parisotto, E., Mohamed, A.-r., Singh, R., Li, L., Zhou, D., and Kohli, P. (2016). Neuro-symbolic program synthesis. In *International Conference on Learning Representations (ICLR)*.
- PCWorld (2012). Microsoft Office 2013 Preview: Hands On.
- Polozov, O. and Gulwani, S. (2015). Flashmeta: A framework for inductive program synthesis. *ACM SIGPLAN Notices*, 50(10):107–126.
- PROSE (2015). Microsoft SDK.
- Raychev, V., Bielik, P., Vechev, M., and Krause, A. (2016). Learning programs from noisy data. In *ACM SIGPLAN Notices*, volume 51, pages 761–774. ACM.
- Rolim, R., Soares, G., D’Antoni, L., Polozov, O., Gulwani, S., Gheyi, R., Suzuki, R., and Hartmann, B. (2017). Learning syntactic program transformations from examples. In *Proceedings of the 39th International Conference on Software Engineering*, pages 404–415. IEEE Press.
- Singh, R. and Gulwani, S. (2015). Predicting a correct program in programming by example. In *International Conference on Computer Aided Verification*, pages 398–414. Springer.
- Wang, X., Dillig, I., and Singh, R. (2017). Program synthesis using abstraction refinement. *Proceedings of the ACM on Programming Languages*, 2(POPL):63.