



Learning Abstractions for Program Synthesis

Xinyu Wang^{1(✉)}, Greg Anderson^{1(✉)}, Isil Dillig^{1(✉)}, and K. L. McMillan^{2(✉)}

¹ University of Texas, Austin, USA
`{xwang,ganderso,isil}@cs.utexas.edu`

² Microsoft Research, Redmond, USA
`kenmcmil@microsoft.com`

Abstract. Many example-guided program synthesis techniques use *abstractions* to prune the search space. While abstraction-based synthesis has proven to be very powerful, a domain expert needs to provide a suitable abstract domain, together with the abstract transformers of each DSL construct. However, coming up with useful abstractions can be non-trivial, as it requires both domain expertise and knowledge about the synthesizer. In this paper, we propose a new technique for learning abstractions that are useful for instantiating a general synthesis framework in a new domain. Given a DSL and a small set of training problems, our method uses *tree interpolation* to infer reusable predicate templates that speed up synthesis in a given domain. Our method also learns suitable abstract transformers by solving a certain kind of second-order constraint solving problem in a data-driven way. We have implemented the proposed method in a tool called ATLAS and evaluate it in the context of the BLAZE meta-synthesizer. Our evaluation shows that (a) ATLAS can learn useful abstract domains and transformers from few training problems, and (b) the abstractions learned by ATLAS allow BLAZE to achieve significantly better results compared to manually-crafted abstractions.

1 Introduction

Program synthesis is a powerful technique for automatically generating programs from high-level specifications, such as input-output examples. Due to its myriad use cases across a wide range of application domains (e.g., spreadsheet automation [1–3], data science [4–6], cryptography [7,8], improving programming productivity [9–11]), program synthesis has received widespread attention from the research community in recent years.

Because program synthesis is, in essence, a very difficult search problem, many recent solutions prune the search space by utilizing *program abstractions* [4, 12–16]. For example, state-of-the-art synthesis tools, such as BLAZE [14], MORPHEUS [4] and Scythe [16], symbolically execute (partial) programs over some abstract domain and reject those programs whose abstract behavior is inconsistent with the given specification. Because many programs share the same behavior in terms of their abstract semantics, the use of abstractions allows these synthesis tools to significantly reduce the search space.

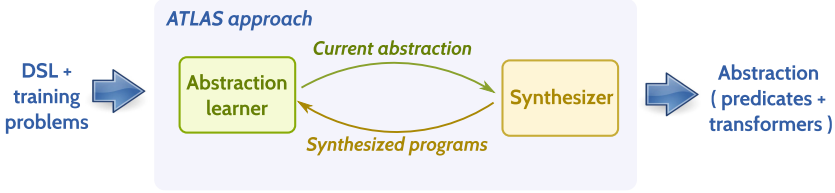


Fig. 1. Schematic overview of our approach.

While the abstraction-guided synthesis paradigm has proven to be quite powerful, a down-side of such techniques is that they require a domain expert to manually come up with a suitable abstract domain and write abstract transformers for each DSL construct. For instance, the BLAZE synthesis framework [14] expects a domain expert to manually specify a universe of predicate templates, together with sound abstract transformers for every DSL construct. Unfortunately, this process is not only time-consuming but also requires significant insight about the application domain as well as the internal workings of the synthesizer.

In this paper, we propose a novel technique for automatically learning domain-specific abstractions that are useful for instantiating an example-guided synthesis framework in a new domain. Given a DSL and a training set of synthesis problems (i.e., input-output examples), our method learns a useful abstract domain in the form of predicate templates and infers sound abstract transformers for each DSL construct. In addition to eliminating the significant manual effort required from a domain expert, the abstractions learned by our method often outperform manually-crafted ones in terms of their benefit to synthesizer performance.

The workflow of our approach, henceforth called ATLAS¹, is shown schematically in Fig. 1. Since ATLAS is meant to be used as an *off-line* training step for a general-purpose programming-by-example (PBE) system, it takes as input a DSL as well as a set of synthesis problems \mathcal{E} that can be used for training purposes. Given these inputs, our method enters a refinement loop where an *Abstraction Learner* component discovers a sequence of increasingly precise abstract domains $\mathcal{A}_1, \dots, \mathcal{A}_n$, and their corresponding abstract transformers $\mathcal{T}_1, \dots, \mathcal{T}_n$, in order to help the *Abstraction-Guided Synthesizer* (AGS) solve all training problems. While the AGS can reject many incorrect solutions using an abstract domain \mathcal{A}_i , it might still return some incorrect solutions due to the insufficiency of \mathcal{A}_i . Thus, whenever the AGS returns an incorrect solution to any training problem, the Abstraction Learner discovers a more precise abstract domain and automatically synthesizes the corresponding abstract transformers. Upon termination of the algorithm, the final abstract domain \mathcal{A}_n and transformers \mathcal{T}_n are sufficient for the AGS to correctly solve *all* training problems. Furthermore, because our method learns *general* abstractions in the form of

¹ ATLAS stands for AuTomated Learning of AbStractions.

predicate templates, the learned abstractions are expected to be useful for solving many *other* synthesis problems beyond those in the training set.

From a technical perspective, the Abstraction Learner uses two key ideas, namely *tree interpolation* and *data-driven constraint solving*, for learning useful abstract domains and transformers respectively. Specifically, given an incorrect program \mathcal{P} that cannot be refuted by the AGS using the current abstract domain \mathcal{A}_i , the Abstraction Learner generates a tree interpolant \mathcal{I}_i that serves as a proof of \mathcal{P} 's incorrectness and constructs a new abstract domain \mathcal{A}_{i+1} by extracting templates from the predicates used in \mathcal{I}_i . The Abstraction Learner also synthesizes the corresponding abstract transformers for \mathcal{A}_{i+1} by setting up a *second-order constraint solving* problem where the goal is to find the unknown relationship between symbolic constants used in the predicate templates. Our method solves this problem in a data-driven way by sampling input-output examples for DSL operators and ultimately reduces the transformer learning problem to solving a system of linear equations.

We have implemented these ideas in a tool called ATLAS and evaluate it in the context of the BLAZE program synthesis framework [14]. Our evaluation shows that the proposed technique eliminates the manual effort involved in designing useful abstractions. More surprisingly, our evaluation also shows that the abstractions generated by ATLAS outperform manually-crafted ones in terms of the performance of the BLAZE synthesizer in two different application domains.

To summarize, this paper makes the following key contributions:

- We describe a method for learning abstractions (domains/transformers) that are useful for instantiating program synthesis frameworks in new domains.
- We show how tree interpolation can be used for learning abstract domains (i.e., predicate templates) from a few training problems.
- We describe a method for automatically synthesizing transformers for a given abstract domain under certain assumptions. Our method is guaranteed to find the unique best transformer if one exists.
- We implement our method in a tool called ATLAS and experimentally evaluate it in the context of the BLAZE synthesis framework. Our results demonstrate that the abstractions discovered by ATLAS outperform manually-written ones used for evaluating BLAZE in two application domains.

2 Illustrative Example

Suppose that we wish to use the BLAZE meta-synthesizer to automate the class of string transformations considered by FlashFill [1] and BlinkFill [17]. In the original version of the BLAZE framework [14], a domain expert needs to come up with a universe of suitable predicate templates as well as abstract transformers for each DSL construct. We will now illustrate how ATLAS automates this process, given a suitable DSL and its semantics (e.g., the one used in [17]).

In order to use ATLAS, one needs to provide a set of synthesis problems \mathcal{E} (i.e., input-output examples) that will be used in the training process. Specifically, let us consider the three synthesis problems given below:

$$\mathcal{E} = \left\{ \begin{array}{l} \mathcal{E}_1 : \{ \text{"CAV"} \mapsto \text{"CAV2018"}, \text{"SAS"} \mapsto \text{"SAS2018"}, \text{"FSE"} \mapsto \text{"FSE2018"} \}, \\ \mathcal{E}_2 : \{ \text{"510.220.5586"} \mapsto \text{"510-220-5586"} \}, \\ \mathcal{E}_3 : \{ \text{"\\Company\\Code\\index.html"} \mapsto \text{"\\Company\\Code\\"}, \\ \text{"\\Company\\Docs\\Spec\\specs.html"} \mapsto \text{"\\Company\\Docs\\Spec\\"} \} \end{array} \right\}.$$

In order to construct the abstract domain \mathcal{A} and transformers \mathcal{T} , ATLAS starts with the trivial abstract domain $\mathcal{A}_0 = \{\top\}$ and transformers \mathcal{T}_0 , defined as $\llbracket F(\top, \cdot, \top) \rrbracket^\# = \top$ for each DSL construct F . Using this abstraction, ATLAS invokes BLAZE to find a program \mathcal{P}_0 that satisfies specification \mathcal{E}_1 under the current abstraction $(\mathcal{A}_0, \mathcal{T}_0)$. However, since the program \mathcal{P}_0 returned by BLAZE is incorrect with respect to the concrete semantics, ATLAS tries to find a more precise abstraction that allows BLAZE to succeed.

Towards this goal, ATLAS enters a refinement loop that culminates in the discovery of the abstract domain $\mathcal{A}_1 = \{\top, \text{len}(\boxed{\alpha}) = c, \text{len}(\boxed{\alpha}) \neq c\}$, where α denotes a variable and c is an integer constant. In other words, \mathcal{A}_1 tracks equality and inequality constraints on the length of strings. After learning these predicate templates, ATLAS also synthesizes the corresponding abstract transformers \mathcal{T}_1 . In particular, for each DSL construct, ATLAS learns one abstract transformer for each combination of predicate templates used in \mathcal{A}_1 . For instance, for the `Concat` operator which returns the concatenation y of two strings x_1, x_2 , ATLAS synthesizes the following abstract transformers, where \star denotes any predicate:

$$\mathcal{T}_1 = \left\{ \begin{array}{l} \llbracket \text{Concat}(\top, \star) \rrbracket^\# = \top \\ \llbracket \text{Concat}(\star, \top) \rrbracket^\# = \top \\ \llbracket \text{Concat}(\text{len}(x_1) \neq c_1, \text{len}(x_2) \neq c_2) \rrbracket^\# = \top \\ \llbracket \text{Concat}(\text{len}(x_1) = c_1, \text{len}(x_2) = c_2) \rrbracket^\# = (\text{len}(y) = c_1 + c_2) \\ \llbracket \text{Concat}(\text{len}(x_1) = c_1, \text{len}(x_2) \neq c_2) \rrbracket^\# = (\text{len}(y) \neq c_1 + c_2) \\ \llbracket \text{Concat}(\text{len}(x_1) \neq c_1, \text{len}(x_2) = c_2) \rrbracket^\# = (\text{len}(y) \neq c_1 + c_2) \end{array} \right\}.$$

Since the AGS can successfully solve \mathcal{E}_1 using $(\mathcal{A}_1, \mathcal{T}_1)$, ATLAS now moves on to the next training problem.

For synthesis problem \mathcal{E}_2 , the current abstraction $(\mathcal{A}_1, \mathcal{T}_1)$ is *not* sufficient for BLAZE to discover the correct program. After processing \mathcal{E}_2 , ATLAS refines the abstract domain to the following set of predicate templates:

$$\mathcal{A}_2 = \{ \top, \text{len}(\boxed{\alpha}) = c, \text{len}(\boxed{\alpha}) \neq c, \text{charAt}(\boxed{\alpha}, i) = c, \text{charAt}(\boxed{\alpha}, i) \neq c \}.$$

Observe that ATLAS has discovered two additional predicate templates that track positions of characters in the string. ATLAS also learns the corresponding abstract transformers \mathcal{T}_2 for \mathcal{A}_2 .

Moving on to the final training problem \mathcal{E}_3 , BLAZE can already successfully solve it using $(\mathcal{A}_2, \mathcal{T}_2)$; thus, ATLAS terminates with this abstraction.

3 Overall Abstraction Learning Algorithm

Our top-level algorithm for learning abstractions, called LEARNABSTRACTIONS, is shown in Fig. 2. The algorithm takes two inputs, namely a domain-specific

```

1: procedure LEARNABSTRACTIONS( $\mathcal{L}, \mathcal{E}$ )
   input: Domain-specific language  $\mathcal{L}$  and a set of training problems  $\mathcal{E}$ .
   output: Abstract domain  $\mathcal{A}$  and transformers  $\mathcal{T}$ .
2:    $\mathcal{A} \leftarrow \{ \top \};$  ▷ Initialization.
3:    $\mathcal{T} \leftarrow \{ \llbracket F(\top, \cdot, \top) \rrbracket^\sharp = \top \mid F \in \text{Constructs}(\mathcal{L}) \};$ 
4:   for  $i \leftarrow 1, \dots, |\mathcal{E}|$  do
5:     while true do ▷ Refinement loop.
6:        $\mathcal{P} \leftarrow \text{Synthesize}(\mathcal{L}, \mathcal{E}_i, \mathcal{A}, \mathcal{T});$  ▷ Invoke AGS.
7:       if  $\mathcal{P} = \text{null}$  then break;
8:       if  $\text{IsCorrect}(\mathcal{P}, \mathcal{E}_i)$  then break;
9:        $\mathcal{A} \leftarrow \mathcal{A} \cup \text{LEARNABSTRACTDOMAIN}(\mathcal{P}, \mathcal{E}_i);$ 
10:       $\mathcal{T} \leftarrow \text{LEARNTRANSFORMERS}(\mathcal{L}, \mathcal{A});$ 
11:   return  $(\mathcal{A}, \mathcal{T});$ 

```

Fig. 2. Overall learning algorithm. **Constructs** gives the DSL constructs in \mathcal{L} .

language \mathcal{L} (both syntax and semantics) as well as a set of training problems \mathcal{E} , where each problem is specified as a *set* of input-output examples \mathcal{E}_i . The output of our algorithm is a pair $(\mathcal{A}, \mathcal{T})$, where \mathcal{A} is an abstract domain represented by a set of predicate templates and \mathcal{T} is the corresponding abstract transformers.

At a high-level, the LEARNABSTRACTIONS procedure starts with the most imprecise abstraction (just consisting of \top) and incrementally improves the precision of the abstract domain \mathcal{A} whenever the AGS fails to synthesize the correct program using \mathcal{A} . Specifically, the outer loop (lines 4–10) considers each training instance \mathcal{E}_i and performs a fixed-point computation (lines 5–10) that terminates when the current abstract domain \mathcal{A} is good enough to solve problem \mathcal{E}_i . Thus, upon termination, the learned abstract domain \mathcal{A} is sufficiently precise for the AGS to solve all training problems \mathcal{E} .

Specifically, in order to find an abstraction that is sufficient for solving \mathcal{E}_i , our algorithm invokes the AGS with the current abstract domain \mathcal{A} and corresponding transformers \mathcal{T} (line 6). We assume that **Synthesize** returns a program \mathcal{P} that is consistent with \mathcal{E}_i under abstraction $(\mathcal{A}, \mathcal{T})$. That is, symbolically executing \mathcal{P} (according to \mathcal{T}) on inputs $\mathcal{E}_i^{\text{in}}$ yields abstract values φ that are consistent with the outputs $\mathcal{E}_i^{\text{out}}$ (i.e., $\forall j. \mathcal{E}_{ij}^{\text{out}} \in \gamma(\varphi_j)$). However, while \mathcal{P} is guaranteed to be consistent with \mathcal{E}_i under the abstract semantics, it may not satisfy \mathcal{E}_i under the concrete semantics. We refer to such a program \mathcal{P} as *spurious*.

Thus, whenever the call to **IsCorrect** fails at line 8, we invoke the LEARNABSTRACTDOMAIN procedure (line 9) to learn additional predicate templates that are later added to \mathcal{A} . Since the refinement of \mathcal{A} necessitates the synthesis of new transformers, we then call LEARNTRANSFORMERS (line 10) to learn a new \mathcal{T} . The new abstraction is guaranteed to rule out the spurious program \mathcal{P} as long as there is a unique best transformer of each DSL construct for domain \mathcal{A} .

4 Learning Abstract Domain Using Tree Interpolation

In this section, we present the `LEARNABSTRACTDOMAIN` procedure: Given a spurious program \mathcal{P} and a synthesis problem \mathcal{E} that \mathcal{P} does not solve, our goal is to find new predicate templates \mathcal{A}' to add to the abstract domain \mathcal{A} such that the Abstraction-Guided Synthesizer no longer returns \mathcal{P} as a valid solution to the synthesis problem \mathcal{E} . Our key insight is that we can mine for such useful predicate templates by constructing a *tree interpolation* problem. In what follows, we first review tree interpolants (based on [18]) and then explain how we use this concept to find useful predicate templates.

Definition 1 (Tree interpolation problem). A *tree interpolation problem* $T = (V, r, P, L)$ is a directed labeled tree, where V is a finite set of nodes, $r \in V$ is the root, $P : (V \setminus \{r\}) \mapsto V$ is a function that maps children nodes to their parents, and $L : V \mapsto \mathbb{F}$ is a labeling function that maps nodes to formulas from a set \mathbb{F} of first-order formulas such that $\bigwedge_{v \in V} L(v)$ is unsatisfiable.

In other words, a tree interpolation problem is defined by a tree T where each node is labeled with a formula and the conjunction of these formulas is unsatisfiable. In what follows, we write $Desc(v)$ to denote the set of all descendants of node v , including v itself, and we write $NonDesc(v)$ to denote all nodes other than those in $Desc(v)$ (i.e., $V \setminus Desc(v)$). Also, given a set of nodes V' , we write $L(V')$ to denote the set of all formulas labeling nodes in V' .

Given a tree interpolation problem T , a *tree interpolant* \mathcal{I} is an annotation from every node in V to a formula such that the label of the root node is *false* and the label of an internal node v is entailed by the conjunction of annotations of its children nodes. More formally, a tree interpolant is defined as follows:

Definition 2 (Tree interpolant). Given a tree interpolation problem $T = (V, r, P, L)$, a tree interpolant for T is a function $\mathcal{I} : V \mapsto \mathbb{F}$ that satisfies the following conditions:

1. $\mathcal{I}(r) = \text{false}$;
2. For each $v \in V$: $\left(\left(\bigwedge_{P(c_i)=v} \mathcal{I}(c_i) \right) \wedge L(v) \right) \Rightarrow \mathcal{I}(v)$;
3. For each $v \in V$: $Vars(\mathcal{I}(v)) \subseteq Vars(L(Desc(v))) \cap Vars(L(NonDesc(v)))$.

Intuitively, the first condition ensures that \mathcal{I} establishes the unsatisfiability of formulas in T , and the second condition states that \mathcal{I} is a valid annotation. As standard in Craig interpolation [19, 20], the third condition stipulates a “shared vocabulary” condition by ensuring that the annotation at each node v refers to the common variables between the descendants and non-descendants of v .

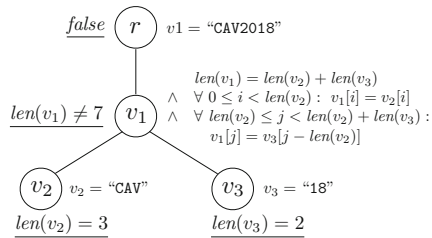


Fig. 3. A tree interpolation problem and a tree interpolant (underlined).

```

1: procedure LEARNABSTRACTDOMAIN( $\mathcal{P}, \mathcal{E}$ )
   input: Program  $\mathcal{P}$  that does not solve problem  $\mathcal{E}$  (set of examples).
   output: Set of predicate templates  $\mathcal{A}'$ .
2:    $\mathcal{A}' \leftarrow \emptyset$ ;
3:   for each  $(e_{in}, e_{out}) \in \mathcal{E}$  do
4:     if  $\llbracket \mathcal{P} \rrbracket_{e_{in}} \neq e_{out}$  then
5:        $T \leftarrow \text{CONSTRUCTTREE}(\mathcal{P}, e_{in}, e_{out})$ ;
6:        $\mathcal{I} \leftarrow \text{FindTreelp}(T)$ ;
7:       for each  $v \in \text{Nodes}(T) \setminus \{r\}$  do
8:          $\mathcal{A}' \leftarrow \mathcal{A}' \cup \{\text{MakeSymbolic}(\mathcal{I}(v))\}$ ;
9:   return  $\mathcal{A}'$ ;

```

Fig. 4. Algorithm for learning abstract domain using tree interpolation.

Example 1. Consider the tree interpolation problem $T = (V, r, P, L)$ in Fig. 3, where $L(v)$ is shown to the right of each node v . A tree interpolant \mathcal{I} for this problem maps each node to the corresponding underlined formula. For instance, we have $\mathcal{I}(v_1) = (\text{len}(v_1) \neq 7)$. It is easy to confirm that \mathcal{I} is a valid interpolant according to Definition 2.

To see how tree interpolation is useful for learning predicates, suppose that the spurious program \mathcal{P} is represented as an abstract syntax tree (AST), where each non-leaf node is labeled with the axiomatic semantics of the corresponding DSL construct. Now, since \mathcal{P} does not satisfy the given input-output example (e_{in}, e_{out}) , we are able to use this information to construct a labeled tree where the conjunction of labels is unsatisfiable. Our key idea is to mine useful predicate templates from the formulas used in the resulting tree interpolant.

With this intuition in mind, let us consider the LEARNABSTRACTDOMAIN procedure shown in Fig. 4: The algorithm uses a procedure called CONSTRUCTTREE to generate a tree interpolation problem T for each input-output example (e_{in}, e_{out}) ² that program \mathcal{P} does not satisfy (line 5). Specifically, letting Π denote the AST representation of \mathcal{P} , we construct $T = (V, r, P, L)$ as follows:

- V consists of all AST nodes in Π as well as a “dummy” node d .
- The root r of T is the dummy node d .
- P is a function that maps children AST nodes to their parents and maps the root AST node to the dummy node d .
- L maps each node $v \in V$ to a formula as follows:

$$L(v) = \begin{cases} v' = e_{out} & v \text{ is the dummy root node with child } v'. \\ v = e_{in} & v \text{ is a leaf representing program input } e_{in}. \\ v = c & v \text{ is a leaf representing constant } c. \\ \phi_F[v'/x, v/y] & v \text{ represents DSL operator } F \text{ with axiomatic semantics } \phi_F(x, y) \text{ and } v' \text{ represents children of } v. \end{cases}$$

² Without loss of generality, we assume that programs take a single input x , as we can always represent multiple inputs as a list.

Essentially, the `CONSTRUCTTREE` procedure labels any leaf node representing the program input with the input example e_{in} and the root node with the output example e_{out} . All other internal nodes are labeled with the axiomatic semantics of the corresponding DSL operator (modulo renaming).³ Observe that the formula $\bigwedge_{v \in V} L(v)$ is guaranteed to be unsatisfiable since \mathcal{P} does not satisfy the I/O example (e_{in}, e_{out}) ; thus, we can obtain a tree interpolant for T .

Example 2. Consider program $\mathcal{P} : \text{Concat}(x, \text{"18"})$ which concatenates constant string “18” to input x . Figure 3 shows the result of invoking `CONSTRUCTTREE` for \mathcal{P} and input-output example (“CAV”, “CAV2018”). As mentioned in Example 1, the tree interpolant \mathcal{I} for this problem is indicated with the underlined formulas.

Since the tree interpolant \mathcal{I} effectively establishes the incorrectness of program \mathcal{P} , the predicates used in \mathcal{I} serve as useful abstract values that the synthesizer (AGS) should consider during the synthesis task. Towards this goal, the `LEARNABSTRACTDOMAIN` algorithm iterates over each predicate used in \mathcal{I} (lines 7–8 in Fig. 4) and converts it to a suitable template by replacing the constants and variables used in $\mathcal{I}(v)$ with symbolic names (or “holes”). Because the original predicates used in \mathcal{I} may be too specific for the current input-output example, extracting templates from the interpolant allows our method to learn reusable abstract domains.

Example 3. Given the tree interpolant \mathcal{I} from Example 1, `LEARNABSTRACTDOMAIN` extracts two predicate templates, namely, $\text{len}(\boxed{\alpha}) = c$ and $\text{len}(\boxed{\alpha}) \neq c$.

5 Synthesis of Abstract Transformers

In this section, we turn our attention to the `LEARNTRANSFORMERS` procedure for synthesizing abstract transformers \mathcal{T} for a given abstract domain \mathcal{A} . Following presentation in prior work [14], we consider abstract transformers that are described using equations of the following form:

$$\llbracket F(\chi_1(x_1, \mathbf{c}_1), \dots, \chi_n(x_n, \mathbf{c}_n)) \rrbracket^\# = \bigwedge_{1 \leq j \leq m} \chi'_j(y, \mathbf{f}_j(\mathbf{c})) \quad (1)$$

Here, F is a DSL construct, χ_i, χ'_j are predicate templates⁴, x_i is the i 'th input of F , y is F 's output, $\mathbf{c}_1, \dots, \mathbf{c}_n$ are vectors of *symbolic* constants, and \mathbf{f}_j denotes a vector of *affine functions* over $\mathbf{c} = \mathbf{c}_1, \dots, \mathbf{c}_n$. Intuitively, given concrete predicates describing the inputs to F , the transformer returns concrete predicates describing the output. Given such a transformer τ , let $\text{Outputs}(\tau)$ be the set of pairs (χ'_j, \mathbf{f}_j) in Eq. 1.

³ Here, we assume access to the DSL's axiomatic semantics. If this is not the case (i.e., we are only given the DSL's operational semantics), we can still annotate each node as $v = c$ where c denotes the output of the partial program rooted at node v when executed on e_{in} . However, this may affect the quality of the resulting interpolant.

⁴ We assume that χ'_1, \dots, χ'_m are distinct.


```

1: procedure LEARNTRANSFORMERS( $\mathcal{L}, \mathcal{A}$ )
   input: DSL  $\mathcal{L}$  and abstract domain  $\mathcal{A}$ .
   output: A set of transformers  $\mathcal{T}$  for constructs in  $\mathcal{L}$  and abstract domain  $\mathcal{A}$ .

2:   for each  $F \in \text{Constructs}(\mathcal{L})$  do
3:     for  $(\chi_1, \dots, \chi_n) \in \mathcal{A}^n$  do
4:        $\varphi \leftarrow \top$ ;  $\triangleright \varphi$  is output of transformer.
5:       for  $\chi'_j \in \mathcal{A}$  do
6:          $E \leftarrow \text{GENERATEEXAMPLES}(\phi_F, \chi'_j, \chi_1, \dots, \chi_n)$ ;
7:          $\mathbf{f}_j \leftarrow \text{Solve}(E)$ ;
8:         if  $\mathbf{f}_j \neq \text{null} \wedge \text{Valid}(\Lambda[\mathbf{f}_j])$  then  $\varphi \leftarrow (\varphi \wedge \chi'_j(y, \mathbf{f}_j(\mathbf{c}_1, \dots, \mathbf{c}_n)))$ 
9:        $\mathcal{T} \leftarrow \mathcal{T} \cup \{ \llbracket F(\chi_1(x_1, \mathbf{c}_1), \dots, \chi_n(x_n, \mathbf{c}_n)) \rrbracket^\# = \varphi \}$ ;
10:  return  $\mathcal{T}$ ;

```

Fig. 5. Algorithm for synthesizing abstract transformers. ϕ_F at line 6 denotes the axiomatic semantics of DSL construct F . Formula Λ at line 8 refers to Eq. 5.

We define the soundness of a transformer τ for DSL operator F with respect to F 's axiomatic semantics ϕ_F . In particular, we say that the abstract transformer from Eq. 1 is *sound* if the following implication is valid:

$$\left(\phi_F(\mathbf{x}, y) \wedge \bigwedge_{1 \leq i \leq n} \chi_i(x_i, \mathbf{c}_i) \right) \Rightarrow \bigwedge_{1 \leq j \leq m} \chi'_j(y, \mathbf{f}_j(\mathbf{c})) \quad (2)$$

That is, the transformer for F is sound if the (symbolic) output predicate is indeed implied by the (symbolic) input predicates according to F 's semantics.

Our key observation is that the problem of learning sound transformers can be reduced to solving the following *second-order constraint solving* problem:

$$\exists \mathbf{f}. \forall \mathbf{V}. \left(\left(\phi_F(\mathbf{x}, y) \wedge \bigwedge_{1 \leq i \leq n} \chi_i(x_i, \mathbf{c}_i) \right) \Rightarrow \bigwedge_{1 \leq j \leq m} \chi'_j(y, \mathbf{f}_j(\mathbf{c})) \right) \quad (3)$$

where $\mathbf{f} = \mathbf{f}_1, \dots, \mathbf{f}_m$ and \mathbf{V} includes all variables and functions from Eq. 2 other than \mathbf{f} . In other words, the goal of this constraint solving problem is to find interpretations of the unknown functions \mathbf{f} that make Eq. 2 valid. Our key insight is to solve this problem in a *data-driven* way by exploiting the fact that each unknown function $\mathbf{f}_{j,k}$ is affine.

Towards this goal, we first express each affine function $\mathbf{f}_{j,k}(\mathbf{c})$ as follows:

$$\mathbf{f}_{j,k}(\mathbf{c}) = p_{j,k,1} \cdot c_1 + \dots + p_{j,k,|\mathbf{c}|} \cdot c_{|\mathbf{c}|} + p_{j,k,|\mathbf{c}|+1}$$

where each $p_{j,k,l}$ corresponds to an unknown integer constant that we would like to learn. Now, arranging the coefficients of functions $\mathbf{f}_{j,1}, \dots, \mathbf{f}_{j,|\mathbf{f}_j|}$ in \mathbf{f}_j into a $|\mathbf{f}_j| \times (|\mathbf{c}| + 1)$ matrix P_j , we can represent $\mathbf{f}_j(\mathbf{c})$ in the following way:

$$\mathbf{f}_j(\mathbf{c})^\top = \underbrace{\begin{bmatrix} f_{j,1}(\mathbf{c}) \\ \vdots \\ f_{j,|\mathbf{f}_j|}(\mathbf{c}) \end{bmatrix}}_{\mathbf{c}'^\top} = \underbrace{\begin{bmatrix} p_{j,1,1} & \cdots & p_{j,1,|\mathbf{c}|+1} \\ \vdots & & \vdots \\ p_{j,|\mathbf{f}_j|,1} & \cdots & p_{j,|\mathbf{f}_j|,|\mathbf{c}|+1} \end{bmatrix}}_{P_j} \underbrace{\begin{bmatrix} c_1 \\ \vdots \\ c_{|\mathbf{c}|} \\ 1 \end{bmatrix}}_{\mathbf{c}^\dagger} \quad (4)$$

where \mathbf{c}^\dagger is \mathbf{c}^\top appended with the constant 1.

Given this representation, it is easy to see that the problem of synthesizing the unknown functions $\mathbf{f}_1, \dots, \mathbf{f}_m$ from Eq. 2 boils down to finding the unknown matrices P_1, \dots, P_m such that each P_j makes the following implication valid:

$$\Delta \equiv \left(\left((\mathbf{c}'^\top = P_j \mathbf{c}^\dagger) \wedge \phi_F(\mathbf{x}, y) \wedge \bigwedge_{1 \leq i \leq n} \chi_i(x_i, \mathbf{c}_i) \right) \Rightarrow \chi'_j(y, \mathbf{c}'_j) \right) \quad (5)$$

Our key idea is to infer these unknown matrices P_1, \dots, P_m in a data-driven way by generating input-output examples of the form $[i_1, \dots, i_{|\mathbf{c}|}] \mapsto [o_1, \dots, o_{|\mathbf{f}_j|}]$ for each \mathbf{f}_j . In other words, \mathbf{i} and \mathbf{o} correspond to instantiations of \mathbf{c} and $\mathbf{f}_j(\mathbf{c})$ respectively. Given sufficiently many such examples for every \mathbf{f}_j , we can then reduce the problem of learning each unknown matrix P_j to the problem of solving a system of linear equations.

Based on this intuition, the LEARNTRANSFORMERS procedure from Fig. 5 describes our algorithm for learning abstract transformers \mathcal{T} for a given abstract domain \mathcal{A} . At a high-level, our algorithm synthesizes one abstract transformer for each DSL construct F and n argument predicate templates χ_1, \dots, χ_n . In particular, given F and χ_1, \dots, χ_n , the algorithm constructs the “return value” of the transformer as:

$$\varphi = \bigwedge_{1 \leq j \leq m} \chi'_j(y, \mathbf{f}_j(\mathbf{c}))$$

where \mathbf{f}_j is the inferred affine function for each predicate template χ'_j .

The key part of our LEARNTRANSFORMERS procedure is the inner loop (lines 5–8) for inferring each of these \mathbf{f}_j 's. Specifically, given an output predicate template χ'_j , our algorithm first generates a set of input-output examples E of the form $[p_1, \dots, p_n] \mapsto p_0$ such that $\llbracket F(p_1, \dots, p_n) \rrbracket^\sharp = p_0$ is a sound (albeit overly specific) transformer. Essentially, each p_i is a concrete instantiation of a predicate template, so the examples E generated at line 6 of the algorithm can be viewed as sound input-output examples for the general symbolic transformer given in Eq. 1. (We will describe the GENERATEEXAMPLES procedure in Sect. 5.1).

Once we generate these examples E , the next step of the algorithm is to learn the unknown coefficients of matrix P_j from Eq. 5 by solving a system of linear equations (line 7). Specifically, observe that we can use each input-output example $[p_1, \dots, p_n] \mapsto p_0$ in E to construct one row of Eq. 4. In particular, we can directly extract $\mathbf{c} = c_1, \dots, c_n$ from p_1, \dots, p_n and the corresponding value of $\mathbf{f}_j(\mathbf{c})$ from p_0 . Since we have one instantiation of Eq. 4 for each of the input-output examples in E , the problem of inferring matrix P_j now reduces to solving a system of linear equations of the form $AP_j^T = B$ where A is a $|E| \times (|\mathbf{c}| + 1)$ (input) matrix and B is a $|E| \times |\mathbf{f}_j|$ (output) matrix. Thus, a solution to the

```

1: procedure GENERATEEXAMPLES( $\phi_F, \chi_0, \dots, \chi_n$ )
   input: Semantics  $\phi_F$  of operator  $F$  and templates  $\chi_0, \dots, \chi_n$  for output and inputs.
   output: A set of valid input-output examples  $E$  for DSL construct  $F$ .
2:    $E \leftarrow \emptyset$ ;
3:   while  $\neg \text{FullRank}(E)$  do
4:     Draw  $(s_1, \dots, s_n)$  randomly from distribution  $R_F$  over  $\text{Domain}(F)$ ;
5:      $s_0 \leftarrow \llbracket F(s_1, \dots, s_n) \rrbracket$ ;
6:      $(A_0, \dots, A_n) \leftarrow \text{Abstract}(s_0, \chi_0, \dots, s_n, \chi_n)$ ;
7:     for each  $(p_0, \dots, p_n) \in A_0 \times \dots \times A_n$  do
8:       if  $\text{Valid}(\bigwedge_{1 \leq i \leq n} p_i \wedge \phi_F \Rightarrow p_0)$  then  $E \leftarrow E \cup \{[p_1, \dots, p_n] \mapsto p_0\}$ ;
9:   return  $E$ ;

```

Fig. 6. Example generation for learning abstract transformers.

equation $AP_j^T = B$ generated from E corresponds to a candidate solution for matrix P_j , which in turn uniquely defines \mathbf{f}_j .

Observe that the call to **Solve** at line 7 may return *null* if no affine function exists. Furthermore, any *non-null* \mathbf{f}_j returned by **Solve** is just a *candidate* solution and may not satisfy Eq. 5. For example, this situation can arise if we do not have sufficiently many examples in E and end up discovering an affine function that is “over-fitted” to the examples. Thus, the validity check at line 8 of the algorithm ensures that the learned transformers are actually sound.

5.1 Example Generation

In our discussion so far, we assumed an oracle that is capable of generating valid input-output examples for a given transformer. We now explain our GENERATEEXAMPLES procedure from Fig. 6 that essentially implements this oracle. In a nutshell, the goal of GENERATEEXAMPLES is to synthesize input-output examples of the form $[p_1, \dots, p_n] \mapsto p_0$ such that $\llbracket F(p_1, \dots, p_n) \rrbracket^\sharp = p_0$ is sound where each p_i is a concrete predicate (rather than symbolic).

Going into more detail, GENERATEEXAMPLES takes as input the semantics ϕ_F of DSL construct F for which we want to learn a transformer for as well as the input predicate templates χ_1, \dots, χ_n and output predicate template χ_0 that are supposed to be used in the transformer. For any example $[p_1, \dots, p_n] \mapsto p_0$ synthesized by GENERATEEXAMPLES, each concrete predicate p_i is an instantiation of the predicate template χ_i where the symbolic constants used in χ_i are substituted with *concrete* values.

Conceptually, the GENERATEEXAMPLES algorithm proceeds as follows: First, it generates *concrete* input-output examples $[s_1, \dots, s_n] \mapsto s_0$ by evaluating F on randomly-generated inputs s_1, \dots, s_n (lines 4–5). Now, for each concrete I/O example $[s_1, \dots, s_n] \mapsto s_0$, we generate a set of *abstract* I/O examples of the form $[p_1, \dots, p_n] \mapsto p_0$ (line 6). Specifically, we assume that the return value (A_0, \dots, A_n) of **Abstract** at line 6 satisfies the following properties for every $p_i \in A_i$:

- p_i is an instantiation of template χ_i .
- p_i is a sound over-approximation of s_i (i.e., $s_i \in \gamma(p_i)$).
- For any other p'_i satisfying the above two conditions, p'_i is not logically stronger than p_i .

In other words, we assume that **Abstract** returns a set of “best” sound abstractions of (s_0, \cdot, s_n) under predicate templates (χ_0, \cdot, χ_n) .

Next, given abstractions (A_0, \cdot, A_n) for (s_0, \cdot, s_n) , we consider each candidate abstract example of the form $[p_1, \cdot, p_n] \mapsto p_0$ where $p_i \in A_i$. Even though each p_i is a sound abstraction of s_i , the example $[p_1, \cdot, p_n] \mapsto p_0$ may not be valid according to the semantics of operator F . Thus, the validity check at line 8 ensures that each example added to E is in fact valid.

Example 4. Given abstract domain $\mathcal{A} = \{\text{len}(\boxed{\alpha}) = c\}$, suppose we want to learn an abstract transformer τ for the **Concat** operator of the following form:

$$\llbracket \text{Concat}(\text{len}(x_1) = c_1, \text{len}(x_2) = c_2) \rrbracket^\sharp = (\text{len}(y) = f([c_1, c_2]))$$

We learn the affine function f used in the transformer by first generating a set E of I/O examples for f (line 6 in **LEARNTRANSFORMERS**). In particular, **GENERATEEXAMPLES** generates concrete input values for **Concat** at random and obtains the corresponding output values by executing **Concat** on the input values. For instance, it may generate $s_1 = \text{“abc”}$ and $s_2 = \text{“de”}$ as inputs, and obtain $s_0 = \text{“abcde”}$ as output. Then, it abstracts these values under the given templates. In this case, we have an abstract example with $p_1 = (\text{len}(x_1) = 3)$, $p_2 = (\text{len}(x_2) = 2)$ and $p_0 = (\text{len}(y) = 5)$. Since $[p_1, p_2] \mapsto p_0$ is a valid example, it is added in E (line 8 in **GENERATEEXAMPLES**). At this point, E is not yet full rank, so the algorithm keeps generating more examples. Suppose it generates two more valid examples $(\text{len}(x_1) = 1, \text{len}(x_2) = 4) \mapsto (\text{len}(y) = 5)$ and $(\text{len}(x_1) = 6, \text{len}(x_2) = 4) \mapsto (\text{len}(y) = 10)$. Now E is full rank, so **LEARNTRANSFORMERS** computes f by solving the following system of linear equations:

$$\begin{bmatrix} 3 & 2 & 1 \\ 1 & 4 & 1 \\ 6 & 4 & 1 \end{bmatrix} P^T = \begin{bmatrix} 5 \\ 5 \\ 10 \end{bmatrix} \xrightarrow{\text{Solve}} P = [1 \ 1 \ 0]$$

Here, P corresponds to the function $f([c_1, c_2]) = c_1 + c_2$, and this function defines the sound transformer: $\llbracket \text{Concat}(\text{len}(x_1) = c_1, \text{len}(x_2) = c_2) \rrbracket^\sharp = (\text{len}(y) = c_1 + c_2)$ which is added to \mathcal{T} at line 9 in **LEARNTRANSFORMERS**.

6 Soundness and Completeness

In this section we present theorems stating some of the soundness, completeness, and termination guarantees of our approach. All proofs can be found in the extended version of this paper [21].

Theorem 1 (Soundness of LEARNTRANSFORMERS). *Let \mathcal{T} be the set of transformers returned by **LEARNTRANSFORMERS**. Then, every $\tau \in \mathcal{T}$ is sound according to Eq. 2.*

The remaining theorems are predicated on the assumptions that for each DSL construct F and input predicate templates χ_1, \dots, χ_n (i) there exists a unique best abstract transformer and (ii) the strongest transformer expressible in Eq. 2 is logically equivalent to the unique best transformer. Thus, before stating these theorems, we first state what we mean by a *unique best abstract transformer*.

Definition 3 (Unique best function). *Consider a family of transformers of the shape $\llbracket F(\chi_1(x_1, \mathbf{c}_1), \dots, \chi_n(x_n, \mathbf{c}_n)) \rrbracket^\# = \chi'(y, \star)$. We say that \mathbf{f} is the unique best function for $(F, \chi_1, \dots, \chi_n, \chi')$ if (a) replacing \star with \mathbf{f} yields a sound transformer, and (b) replacing \star with any other \mathbf{f}' yields a transformer that is either unsound or strictly worse (i.e., $\chi'(y, \mathbf{f}) \Rightarrow \chi'(y, \mathbf{f}')$ and $\chi'(y, \mathbf{f}') \not\Rightarrow \chi'(y, \mathbf{f})$).*

We now define unique best transformer in terms of unique best function:

Definition 4 (Unique best transformer). *Let F be a DSL construct and let $(\chi_1, \dots, \chi_n) \in \mathcal{A}^n$ be the input templates for F . We say that the abstract transformer τ is a unique best transformer for F, χ_1, \dots, χ_n if (a) τ is sound, and (b) for any predicate template $\chi \in \mathcal{A}$, we have $(\chi, \mathbf{f}) \in \text{Outputs}(\tau)$ if and only if \mathbf{f} is a unique best function for $(F, \chi_1, \dots, \chi_n, \chi)$ for some affine \mathbf{f} .*

Definition 5 (Complete sampling oracle). *Let F be a construct, \mathcal{A} an abstract domain, and R_F a probability distribution over $\text{DOMAIN}(F)$ with finite support S . Further, for any input predicate templates χ_1, \dots, χ_n and output predicate template χ_0 in \mathcal{A} admitting a unique best function \mathbf{f} , let $C(\chi_0, \dots, \chi_n)$ be the set of tuples (c_0, \dots, c_n) such that $(\chi_0(y, c_0), \chi_1(x_1, c_1), \dots, \chi_n(x_n, c_n)) \in A_0 \times \dots \times A_n$ and $c_0 = \mathbf{f}(c_1, \dots, c_n)$, where $A_0 \times \dots \times A_n = \text{ABSTRACT}(s_0, \chi_0, \dots, s_n, \chi_n)$ and $(s_1, \dots, s_n) \in S$ and $s_0 = \llbracket F(s_1, \dots, s_n) \rrbracket$. The distribution R_F is a complete sampling oracle if $C(\chi_0, \dots, \chi_n)$ has full rank for all χ_0, \dots, χ_n .*

The following theorem states that `LEARNTRANSFORMERS` is guaranteed to synthesize the best transformer if a unique one exists:

Theorem 2 (Completeness of `LEARNTRANSFORMERS`). *Given an abstract domain \mathcal{A} and a complete sampling oracle R_F for \mathcal{A} , `LEARNTRANSFORMERS` terminates. Further, let \mathcal{T} be the set of transformers returned and let τ be the unique best transformer for DSL construct F and input predicate templates $\chi_1, \dots, \chi_n \in \mathcal{A}^n$. Then we have $\tau \in \mathcal{T}$.*

Using this completeness (modulo unique best transformer) result, we can now state the termination guarantees of our `LEARNABSTRACTIONS` algorithm:

Theorem 3 (Termination of `LEARNABSTRACTIONS`). *Given a complete sampling oracle R_F for every abstract domain and the unique best transformer assumption, if there exists a solution for every problem $\mathcal{E}_i \in \mathcal{E}$, then `LEARNABSTRACTIONS` terminates.*

7 Implementation and Evaluation

We have implemented the proposed method as a new tool called ATLAS, which is written in Java. ATLAS takes as input a set of training problems, an Abstraction-Guided Synthesizer (AGS), and a DSL and returns an abstract domain (in the form of predicate templates) and the corresponding transformers. Internally, ATLAS uses the Z3 theorem prover [22] to compute tree interpolants and the JLinAlg linear algebra library [23] to solve linear equations.

To assess the usefulness of ATLAS, we conduct an experimental evaluation in which our goal is to answer the following two questions:

1. How does ATLAS perform during training? That is, how many training problems does it require and how long does training take?
2. How useful are the abstractions learned by ATLAS in the context of synthesis?

7.1 Abstraction Learning

To answer our first question, we use ATLAS to automatically learn abstractions for two application domains: (i) string manipulations and (ii) matrix transformations. We provide ATLAS with the DSLs used in [14] and employ BLAZE as the underlying Abstraction-Guided Synthesizer. Axiomatic semantics for each DSL construct were given in the theory of equality with uninterpreted functions.

Training Set Information. For the string domain, our training set consists of exactly the four problems used as motivating examples in the BlinkFill paper [17]. Specifically, each training problem consists of 4–6 examples that demonstrate the desired string transformation. For the matrix domain, our training set consists of four (randomly selected) synthesis problems taken from online forums. Since almost all online posts contain a single input-output example, each training problem includes one example illustrating the desired matrix transformation.

Main Results. Our main results are summarized in Fig. 7. The main take-away message is that ATLAS can learn abstractions quite efficiently and does not require a large training set. For example, ATLAS learns 5 predicate templates and 30 abstract transformers for the string domain in a total of 10.2 s. Interestingly, ATLAS does not need all the training problems to infer these four predicates and converges to the final abstraction after just processing the first training instance. Furthermore, for the first training instance, it takes ATLAS 4 iterations in the learning loop (lines 5–10 from Fig. 2) before it converges to the final abstraction. Since this abstraction is sufficient, it takes just one iteration for each following training problem to synthesize a correct program.

Looking at the right side of Fig. 7, we also observe similar results for the matrix domain. In particular, ATLAS learns 10 predicate templates and 59 abstract transformers in a total of 22.5 s. Furthermore, ATLAS converges to the final abstract domain after processing the first three problems⁵ and the number of iterations for each training instance is also quite small (ranging from 1 to 3).

⁵ The learned abstractions can be found in the extended version of this paper [21].

	\mathcal{A} \mathcal{T} Iters.			Running time (sec)			
				T_{AGS}	$T_{\mathcal{A}}$	$T_{\mathcal{T}}$	T_{total}
\mathcal{E}_1	5	30	4	0.6	0.2	0.2	1.0
\mathcal{E}_2	5	30	1	4.9	0	0	4.9
\mathcal{E}_3	5	30	1	0.2	0	0	0.2
\mathcal{E}_4	5	30	1	4.1	0	0	4.1
Total	5	30	7	9.8	0.2	0.2	10.2

String domain

	\mathcal{A} \mathcal{T} Iters.			Running time (sec)			
				T_{AGS}	$T_{\mathcal{A}}$	$T_{\mathcal{T}}$	T_{total}
\mathcal{E}_1	8	45	3	2.9	0.7	0.5	4.1
\mathcal{E}_2	8	45	1	2.8	0	0	2.8
\mathcal{E}_3	10	59	2	0.5	0.3	0.2	1.0
\mathcal{E}_4	10	59	1	14.6	0	0	14.6
Total	10	59	7	20.8	1.0	0.7	22.5

Matrix domain

Fig. 7. Training results. $|\mathcal{A}|$, $|\mathcal{T}|$, Iters denote the number of predicate templates, abstract transformers, and iterations taken per training instance (lines 5–10 from Fig. 2), respectively. T_{AGS} , $T_{\mathcal{A}}$, $T_{\mathcal{T}}$ denote the times for invoking the synthesizer (AGS), learning the abstract domain, and learning the abstract transformers, respectively. T_{total} shows the total training time in seconds.

	Original BLAZE [†] benchmarks				Additional benchmarks				All benchmarks		
	#Solved		Running time improvement		#Solved		Running time improvement		Time (sec)	Running time improvement	
	BLAZE [*]	BLAZE [†]	max.	avg.	BLAZE [*]	BLAZE [†]	max.	avg.	avg.	max.	avg.
String	93	91	15.7×	2.1×	40	40	56×	22.3×	2.8	56×	8.3×
Matrix	39	39	6.1×	3.1×	20	19	83×	21.5×	5.0	83×	9.2×

Fig. 8. Improvement of BLAZE^{*} over BLAZE[†] on string and matrix benchmarks.

7.2 Evaluating the Usefulness of Learned Abstractions

To answer our second question, we integrated the abstractions synthesized by ATLAS into the BLAZE meta-synthesizer. In the remainder of this section, we refer to all instantiations of BLAZE using the ATLAS-generated abstractions as BLAZE^{*}. To assess how useful the automatically generated abstractions are, we compare BLAZE^{*} against BLAZE[†], which refers to the manually-constructed instantiations of BLAZE described in [14].

Benchmark Information. For the string domain, our benchmark suite consists of (1) *all* 108 string transformation benchmarks that were used to evaluate BLAZE[†] and (2) 40 additional challenging problems that are collected from online forums which involve manipulating file paths, URLs, etc. The number of examples for each benchmark ranges from 1 to 400, with a median of 7 examples. For the matrix domain, our benchmark set includes (1) *all* 39 matrix transformation benchmarks in the BLAZE[†] benchmark suite and (2) 20 additional challenging problems collected from online forums. *We emphasize that the set of benchmarks used for evaluating BLAZE^{*} are completely disjoint from the set of synthesis problems used for training ATLAS.*

Experimental Setup. We evaluate BLAZE^{*} and BLAZE[†] using the same DSLs from the BLAZE paper [14]. For each benchmark, we provide the same set of input-output examples to BLAZE^{*} and BLAZE[†], and use a time limit of 20 min per synthesis task.

Main Results. Our main evaluation results are summarized in Fig. 8. The key observation is that BLAZE^{*} consistently improves upon BLAZE[†] for both string and matrix transformations. In particular, BLAZE^{*} not only solves more benchmarks than BLAZE[†] for both domains, but also achieves about an order of magnitude speed-up on average for the common benchmarks that both tools can solve. Specifically, for the string domain, BLAZE^{*} solves 133 (out of 148) benchmarks within an average of 2.8 s and achieves an average 8.3 \times speed-up over BLAZE[†]. For the matrix domain, we also observe a very similar result where BLAZE^{*} leads to an overall speed-up of 9.2 \times on average.

In summary, this experiment confirms that the abstractions discovered by ATLAS are indeed useful and that they outperform manually-crafted abstractions despite eliminating human effort.

8 Related Work

To our knowledge, this paper is the first one to automatically learn abstract domains and transformers that are useful for program synthesis. We also believe it is the first to apply interpolation to program synthesis, although interpolation has been used to synthesize other artifacts such as circuits [24] and strategies for infinite games [25]. In what follows, we briefly survey existing work related to program synthesis, abstraction learning, and abstract transformer computations.

Program Synthesis. Our work is intended to complement example-guided program synthesis techniques that utilize program abstractions to prune the search space [4, 14–16]. For example, SIMPL [15] uses abstract interpretation to speed up search-based synthesis and applies this technique to the generation of imperative programs for introductory programming assignments. Similarly, SCYTHE [16] and MORPHEUS [4] perform enumeration over program sketches and use abstractions to reject sketches that do not have any valid completion. Somewhat different from these techniques, BLAZE constructs a finite tree automaton that accepts all programs whose behavior is consistent with the specification according to the DSL’s abstract semantics. We believe that the method described in this paper can be useful to all such abstraction-guided synthesizers.

Abstraction Refinement. In verification, as opposed to synthesis, there have been many works that use Craig interpolants to refine abstractions [20, 26, 27]. Typically, these techniques generalize the interpolants to abstract domains by extracting a vocabulary of predicates, but they do not generalize by adding parameters to form templates. In our case, this is essential because interpolants

derived from fixed input values are too specific to be directly useful. Moreover, we *reuse* the resulting abstractions for subsequent synthesis problems. In verification, this would be analogous to re-using an abstraction from one property or program to the next. It is conceivable that template-based generalization could be applied in verification to facilitate such reuse.

Abstract Transformers. Many verification techniques use logical abstract domains [28–32]. Some of these, following Yorsh, *et al.* [33] use sampling with a decision procedure to evaluate the abstract transformer [34]. Interpolation has also been used to compile efficient symbolic abstract transformers [35]. However, these techniques are restricted to finite domains or domains of finite height to allow convergence. Here, we use infinite parameterized domains to obtain better generalization; hence, the abstract transformer computation is more challenging. Nonetheless, the approach might also be applicable in verification.

9 Limitations

While this paper takes a first step towards automatically inferring useful abstractions for synthesis, our proposed method has the following limitations:

Shapes of Transformers. Following prior work [14], our algorithm assumes that abstract transformers have the shape given in Eq. 1. We additionally assume that constants \mathbf{c} used in predicate templates are numeric values and that functions in Eq. 1 are affine. This assumption holds in several domains considered in prior work [4, 14] and allows us to develop an efficient learning algorithm that reduces the problem to solving a system of linear equations.

DSL Semantics. Our method requires the DSL designer to provide the DSL’s logical semantics. We believe that giving logical semantics is much easier than coming up with useful abstractions, as it does not require insights about the internal workings of the synthesizer. Furthermore, our technique could, in principle, also work without logical specifications although the learned abstract domain may not be as effective (see Footnote 3 in Sect. 4) and the synthesized transformers would not be provably sound.

UBT Assumption. Our completeness and termination theorems are predicated on the *unique best transformer (UBT)* assumption. While this assumption holds in our evaluation, it may not hold in general. However, as mentioned in Sect. 6, we can always guarantee termination by including the concrete predicates used in the interpolant \mathcal{I} in addition to the symbolic templates extracted from \mathcal{I} .

10 Conclusion

We proposed a new technique for automatically instantiating abstraction-guided synthesis frameworks in new domains. Given a DSL and a few training prob-

lems, our method automatically discovers a useful abstract domain and the corresponding transformers for each DSL construct. From a technical perspective, our method uses tree interpolation to extract reusable templates from failed synthesis attempts and automatically synthesizes unique best transformers if they exist. We have incorporated the proposed approach into the BLAZE meta-synthesizer and show that the abstractions discovered by ATLAS are very useful.

While we have applied the proposed technique to program synthesis, we believe that some of the ideas introduced here are more broadly applicable. For instance, the idea of extracting reusable predicate templates from interpolants and synthesizing transformers in a data-driven way could also be useful in the context of program verification.

References

1. Gulwani, S.: Automating string processing in spreadsheets using input-output examples. In: *Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL*, pp. 317–330. ACM (2011)
2. Singh, R., Gulwani, S.: Transforming spreadsheet data types using examples. In: *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL*, pp. 343–356. ACM (2016)
3. Wang, X., Gulwani, S., Singh, R.: FIDEX: filtering spreadsheet data using examples. In: *OOPSLA*, pp. 195–213. ACM (2016)
4. Feng, Y., Martins, R., Van Geffen, J., Dillig, I., Chaudhuri, S.: Component-based synthesis of table consolidation and transformation tasks from examples. In: *PLDI*, pp. 422–436. ACM (2017)
5. Wang, X., Dillig, I., Singh, R.: Synthesis of data completion scripts using finite tree automata. *Proc. ACM Program. Lang.* 1(OOPSLA), 62:1–62:26 (2017)
6. Yaghmazadeh, N., Wang, X., Dillig, I.: Automated migration of hierarchical data to relational tables using programming-by-example. In: *Proceedings of the VLDB Endowment* (2018)
7. Gascón, A., Tiwari, A., Carmer, B., Mathur, U.: Look for the proof to find the program: decorated-component-based program synthesis. In: Majumdar, R., Kunčák, V. (eds.) *CAV 2017. LNCS*, vol. 10427, pp. 86–103. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-63390-9_5
8. Tiwari, A., Gascón, A., Dutertre, B.: Program synthesis using dual interpretation. In: Felty, A.P., Middeldorp, A. (eds.) *CADE 2015. LNCS (LNAI)*, vol. 9195, pp. 482–497. Springer, Cham (2015). https://doi.org/10.1007/978-3-319-21401-6_33
9. Feng, Y., Martins, R., Wang, Y., Dillig, I., Reps, T.W.: Component-based synthesis for complex APIs. In: *POPL*, vol. 52, pp. 599–612. ACM (2017)
10. Gvero, T., Kuncak, V., Kuraj, I., Piskac, R.: Complete completion using types and weights. In: *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI*, pp. 27–38. ACM (2013)
11. Mandelin, D., Xu, L., Bodík, R., Kimelman, D.: Jungloid mining: helping to navigate the API jungle. In: *Proceedings of the 26th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI*, pp. 48–61. ACM (2005)
12. Feser, J.K., Chaudhuri, S., Dillig, I.: Synthesizing data structure transformations from input-output examples. In: *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI*, pp. 229–239. ACM (2015)

13. Polikarpova, N., Kuraj, I., Solar-Lezama, A.: Program synthesis from polymorphic refinement types. In: Proceedings of the 37th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI, pp. 522–538. ACM (2016)
14. Wang, X., Dillig, I., Singh, R.: Program synthesis using abstraction refinement, vol. 2, pp. 63:1–63:30. ACM (2017)
15. So, S., Oh, H.: Synthesizing imperative programs from examples guided by static analysis. In: Ranzato, F. (ed.) SAS 2017. LNCS, vol. 10422, pp. 364–381. Springer, Cham (2017). https://doi.org/10.1007/978-3-319-66706-5_18
16. Wang, C., Cheung, A., Bodik, R.: Synthesizing highly expressive SQL queries from input-output examples. In: Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI, pp. 452–466. ACM (2017)
17. Singh, R.: BlinkFill: semi-supervised programming by example for syntactic string transformations. *Proc. VLDB Endow.* **9**(10), 816–827 (2016)
18. Blanc, R., Gupta, A., Kovács, L., Kragl, B.: Tree interpolation in vampire. In: McMillan, K., Middeldorp, A., Voronkov, A. (eds.) LPAR 2013. LNCS, vol. 8312, pp. 173–181. Springer, Heidelberg (2013). https://doi.org/10.1007/978-3-642-45221-5_13
19. McMillan, K.L.: Applications of craig interpolants in model checking. In: Halbwachs, N., Zuck, L.D. (eds.) TACAS 2005. LNCS, vol. 3440, pp. 1–12. Springer, Heidelberg (2005). https://doi.org/10.1007/978-3-540-31980-1_1
20. McMillan, K.L.: Interpolation and SAT-based model checking. In: Hunt, W.A., Somenzi, F. (eds.) CAV 2003. LNCS, vol. 2725, pp. 1–13. Springer, Heidelberg (2003). https://doi.org/10.1007/978-3-540-45069-6_1
21. Wang, X., Dillig, I., Singh, R.: Learning Abstractions for Program Synthesis. arXiv preprint [arXiv:1804.04152](https://arxiv.org/abs/1804.04152) (2018)
22. Z3. <https://github.com/Z3Prover/z3>
23. Keilhauer, A., Levy, S., Lochbihler, A., Ökmen, S., Thimm, G., Würzebesser, C.: JLinAlg: a java-library for linear algebra without rounding errors. Technical report (2003–2010). <http://jlinalg.sourceforge.net/>
24. Bloem, R., Egly, U., Klampfl, P., Könighofer, R., Lonsing, F.: Sat-based methods for circuit synthesis. In: Formal Methods in Computer-Aided Design, FMCAD 2014, 21–24 October 2014, Lausanne, Switzerland, pp. 31–34. IEEE (2014)
25. Farzan, A., Kincaid, Z.: Strategy synthesis for linear arithmetic games. *Proc. ACM Program. Lang.* **2**(POPL), 61 (2017)
26. Beyer, D., Henzinger, T.A., Jhala, R., Majumdar, R.: The software model checker BLAST. *Int. J. Softw. Tools Technol. Transf.* **9**(5–6), 505–525 (2007)
27. Albarghouthi, A., Li, Y., Gurfinkel, A., Chechik, M.: UFO: a framework for abstraction- and interpolation-based software verification. In: Madhusudan, P., Seshia, S.A. (eds.) CAV 2012. LNCS, vol. 7358, pp. 672–678. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-31424-7_48
28. Lev-Ami, T., Manevich, R., Sagiv, M.: TVLA: a system for generating abstract interpreters. In: Jacquart, R. (ed.) Building the Information Society. IIFIP, vol. 156, pp. 367–375. Springer, Boston, MA (2004). https://doi.org/10.1007/978-1-4020-8157-6_28
29. Lev-Ami, T., Sagiv, M.: TVLA: a system for implementing static analyses. In: Palsberg, J. (ed.) SAS 2000. LNCS, vol. 1824, pp. 280–301. Springer, Heidelberg (2000). https://doi.org/10.1007/978-3-540-45099-3_15

30. Pnueli, A., Ruah, S., Zuck, L.: Automatic deductive verification with invisible invariants. In: Margaria, T., Yi, W. (eds.) TACAS 2001. LNCS, vol. 2031, pp. 82–97. Springer, Heidelberg (2001). https://doi.org/10.1007/3-540-45319-9_7
31. Lahiri, S.K., Bryant, R.E.: Constructing quantified invariants via predicate abstraction. In: Steffen, B., Levi, G. (eds.) VMCAI 2004. LNCS, vol. 2937, pp. 267–281. Springer, Heidelberg (2004)
32. Reps, T., Thakur, A.: Automating abstract interpretation. In: Jobstmann, B., Leino, K.R.M. (eds.) VMCAI 2016. LNCS, vol. 9583, pp. 3–40. Springer, Heidelberg (2016). https://doi.org/10.1007/978-3-662-49122-5_1
33. Reps, T., Sagiv, M., Yorsh, G.: Symbolic implementation of the best transformer. In: Steffen, B., Levi, G. (eds.) VMCAI 2004. LNCS, vol. 2937, pp. 252–266. Springer, Heidelberg (2004)
34. Thakur, A., Reps, T.: A method for symbolic computation of abstract operations. In: Madhusudan, P., Seshia, S.A. (eds.) CAV 2012. LNCS, vol. 7358, pp. 174–192. Springer, Heidelberg (2012). https://doi.org/10.1007/978-3-642-31424-7_17
35. Jhala, R., McMillan, K.L.: Interpolant-based transition relation approximation. In: Etessami, K., Rajamani, S.K. (eds.) CAV 2005. LNCS, vol. 3576, pp. 39–51. Springer, Heidelberg (2005). https://doi.org/10.1007/11513988_6

Open Access This chapter is licensed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

