

1 Top-Down Synthesis For Library Learning

2 ANONYMOUS AUTHOR(S)

3 This paper introduces *corpus-guided top-down synthesis* as a mechanism for synthesizing library functions that
 4 capture common functionality from a corpus of programs in a domain specific language (DSL). The algorithm
 5 builds abstractions directly from initial DSL primitives, using syntactic pattern matching of intermediate
 6 abstractions to intelligently prune the search space and guide the algorithm towards abstractions that maximally
 7 capture shared structures in the corpus. We present an implementation of the approach in a tool called STITCH
 8 and evaluate it against the state-of-the-art deductive library learning algorithm from DreamCoder [Ellis et al.
 9 2021]. Our evaluation shows that STITCH is 3-4 orders of magnitude faster and uses 2 orders of magnitude
 10 less memory while maintaining comparable or better library quality (as measured by compressivity). We also
 11 demonstrate STITCH’s scalability on corpora containing hundreds of complex programs, which are intractable
 12 with prior deductive approaches, and show empirically that it is robust to terminating the search procedure
 13 early—further allowing it to scale to challenging datasets by means of early stopping.

14
 15 CCS Concepts: • Software and its engineering → Automatic programming.

16 Additional Key Words and Phrases: Program Synthesis, Library Learning, Abstraction Learning

17 **ACM Reference Format:**

18 Anonymous Author(s). 2022. Top-Down Synthesis For Library Learning. *Proc. ACM Program. Lang.* 1, CONF,
 19 Article 1 (June 2022), 29 pages.

20
 21 **1 INTRODUCTION**

22 One way programmers manage complexity is by hiding functionality behind functional abstractions.
 23 For example, consider the graphics programs at the bottom of Fig. 1A. Each uses a generic set of
 24 drawing primitives and renders a technical schematic of a hardware component, shown to the left
 25 of each program. Faced with the task of writing *more* of these rendering programs, an experienced
 26 human programmer likely would not continue using these low-level primitives. Instead, they would
 27 introduce new functional abstractions, like the one at the bottom of Fig. 1B which renders a regular
 28 polygon given a size and number of sides. Useful abstractions like these allow more concise and
 29 legible programs to render the existing schematics. More importantly, well-written abstractions
 30 should generalize, making it easier to write new graphics programs for similar graphics tasks.

31 Recently, the *program synthesis* community has introduced new approaches that can mimic this
 32 process, automatically building a library of functional abstractions in order to tackle more complex
 33 synthesis problems [Dechter et al. 2013; Ellis et al. 2020, 2021; Lázaro-Gredilla et al. 2019; Shin
 34 et al. 2019]. One popular approach to library learning is to search for common tree fragments
 35 across a corpus of programs, which can be introduced as new abstractions [Dechter et al. 2013;
 36 Lázaro-Gredilla et al. 2019; Shin et al. 2019]. Ellis et al. 2021, however, introduces an algorithm that
 37 reasons about variable bindings to abstract out well-formed *functions* instead of just tree fragments.
 38 While it produces impressive results, the system in Ellis et al. 2021 takes a *deductive* approach to
 39 library learning that is difficult to scale to larger datasets of longer and deeper input programs. This
 40 deductive approach attempts to refactor existing programs to expose shared structure, and must
 41 represent and evaluate a large space of proposed refactorings to identify common functionality
 42 across the corpus. Prior, work such as Ellis et al. 2020, 2021, approaches this challenge by combining
 43 a dynamic bottom-up approach to refactoring with version spaces to more efficiently search over
 44 the refactored programs. However, these deductive approaches face daunting memory and search
 45 requirements as the corpus scales in size and complexity.

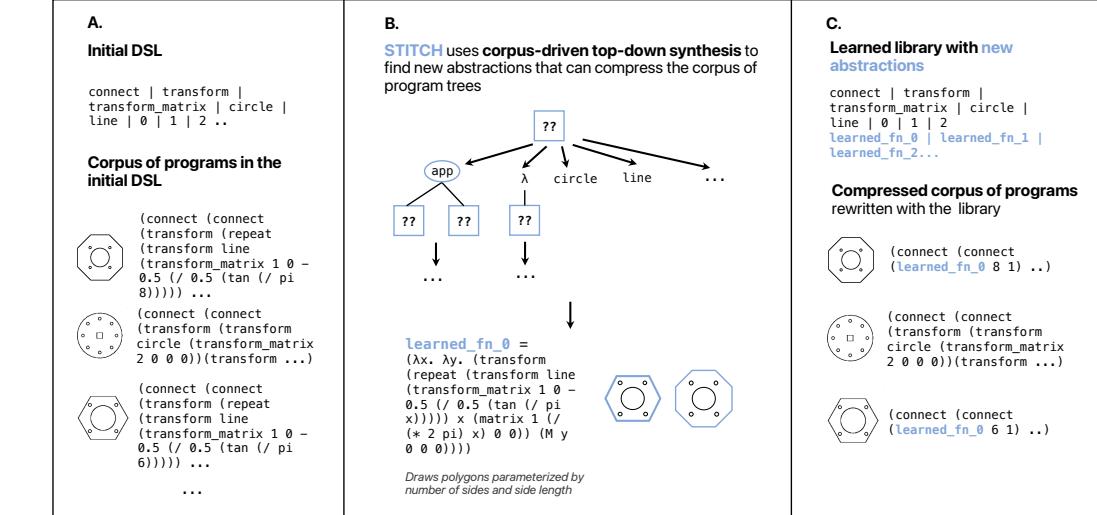


Fig. 1. (A) Given an initial DSL of lower-level primitives and a corpus of programs written in the initial DSL, (B) STITCH uses a fast and memory efficient *corpus-guided top-down search* algorithm to construct function abstractions which maximally capture shared structure across the corpus of programs. (C) STITCH automatically rewrites the initial corpus using the learned library functions.

This paper introduces an alternate approach to library learning: instead of taking a *deductive* approach based on *refactoring the corpus*, we *directly synthesize abstractions*. We call this approach *corpus-guided top-down synthesis*, and it is based on the insight that when applied to the task of synthesizing abstractions, top-down search can be guided precisely towards discovering shared abstractions over a set of existing training programs — at every step of the search, syntactic comparisons between a partially constructed abstraction and the set of training programs can be used to strongly constrain the search space and direct the search towards abstractions that capture the greatest degree of shared syntactic structure.

We implement this approach in STITCH¹, a corpus-guided top-down library learning tool written in Rust (Fig. 1). We evaluate STITCH through a series of experiments (Section 6), and find that: STITCH learns libraries of comparable quality to those found by the algorithm of Ellis et al. 2021 on their iterative bootstrapped library learning task, while being 3-4 orders of magnitude faster and using 2 orders of magnitude less memory (Section 6.1); STITCH learns high quality libraries within seconds to single-digit minutes when run on corpora containing a few hundred programs with mean lengths between 76–189 symbols (sourced from Wong et al. 2022), while even the simplest of these corpora lies beyond the reach of the algorithm of Ellis et al. 2021 (Section 6.2); and that STITCH degrades gracefully when resources are constrained (Section 6.3). We also perform ablation studies to expose the relative impact of different optimizations in STITCH (Section 6.4). Finally, we show that STITCH is complementary to deductive rewrite approaches (Section 6.5).

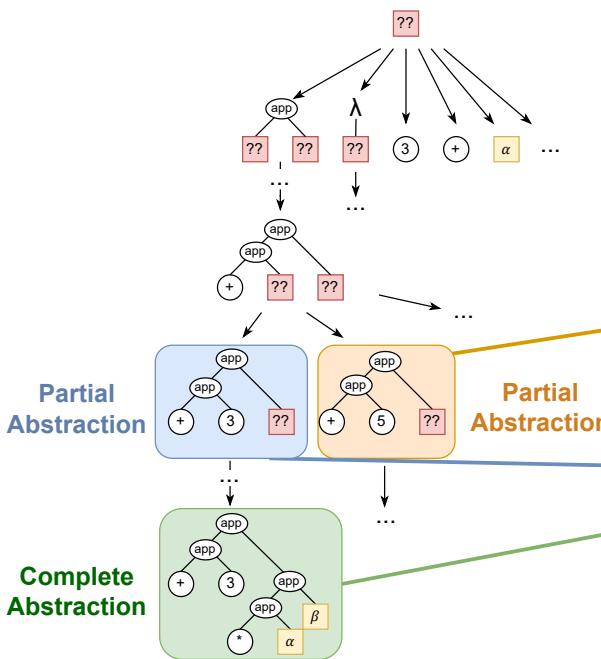
In summary, our paper makes the following contributions:

- (1) **Corpus-guided top-down synthesis (CTS):** A novel, strongly-guided branch-and-bound algorithm for synthesizing function abstractions (Sec. 3).

¹Anonymized version of the STITCH tool is provided in supplement: [supplement/anonymized_stitch.zip](#).

99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121

A.



B.

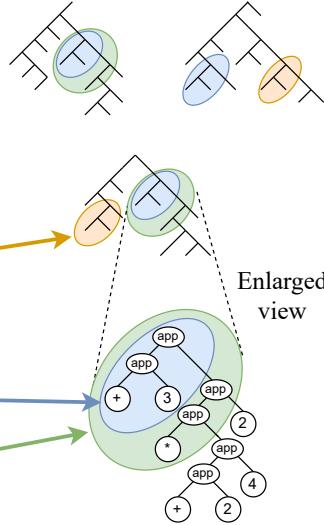


Fig. 2. (A) Schematic of the *corpus-guided top-down search* search. Partial abstractions can contain *holes* indicating unfinished subtrees, denoted ??, while complete abstractions do not contain holes. **(B)** Partial and complete abstractions *match* at locations in the corpus. The blue partial abstraction can be expanded into the green complete abstraction so the green abstraction matches in a *subset* of places that the blue abstraction does. For complete abstractions, *matching* indicates that the subtree can be rewritten to use the abstraction.

(2) **CTS for program compression:** An instantiation of the CTS framework for utility functions favoring abstractions that compress a corpus of programs (Sec. 4).

(3) **STRICH**: A parallel Rust implementation of CTS for compression that achieves 3-4 orders of magnitude of speed and memory improvements over prior work, and the analysis of its performance, scaling, and optimizations through several experiments (Sec. 6).

2 OVERVIEW

In this section, we build intuition for the algorithmic insights that power STITCH. As a running example, we focus on learning a single abstraction from the following corpus of programs:

```

 $\lambda x. + 3 (* (+ 2 4) 2)$ 
 $\lambda xs. \text{map } (\lambda x. + 3 (* 4 (+ 3 x))) \ xs$  (1)
 $\lambda x. * 2 (+ 3 (* x (+ 2 1)))$ 

```

The notion of what a *good* abstraction is depends on the application, so our algorithm is generic over the utility function that we seek to maximize. Following prior work [Dechter et al. 2013; Ellis et al. 2021; Shin et al. 2019] we focus on *compression* as a utility function: a good abstraction is one which minimizes the size of the corpus when rewritten with the abstraction. The utility function used by STITCH is detailed in Section 4 and corresponds exactly to the compression objective, but

at a high level the function seeks to maximize the product of *the size of the abstraction* and the *number of locations* where the abstraction can be used. This product balances two key features of a highly compressive abstraction: the abstraction should be general enough that it applies in many locations, but specific enough that it captures a lot of structure at each location.

The optimal abstraction that maximizes our utility in this example is:

$$\text{fn0} = \lambda\alpha. \lambda\beta. (+ 3 (* \alpha \beta)) \quad (2)$$

And the shared structure that this abstraction captures is highlighted in blue in Eq. (1). When rewritten to use this abstraction, the size of the resulting programs is minimized:

$$\begin{aligned} & \lambda x. \text{fn0} (+ 2 4) 2 \\ & \lambda xs. \text{map} (\lambda x. \text{fn0} 4 (+ 3 x)) xs \\ & \lambda x. * 2 (\text{fn0} x (+ 2 1)) \end{aligned} \quad (3)$$

STITCH synthesizes the optimal abstraction directly through *top-down search*. Top-down search methods construct program syntax trees by iteratively refining partially-completed programs that have unfinished *holes*, until a complete program that meets the specification is produced [Balog et al. 2016; Ellis et al. 2021; Feser et al. 2015; Nye et al. 2021; Polikarpova et al. 2016; Shah et al. 2020]. This kind of search is often made efficient by identifying branches of the search tree that can be pruned away because you can efficiently prove that none of the programs in that branch can be correct. We aim to apply a similar idea to the problem of synthesizing good abstractions; the idea is to explore the space of functions in the same top-down way, but in search of a function that maximizes the utility measure. The key observation is that this new objective affords even more aggressive pruning opportunities than the traditional correctness objective, allowing us to synthesize optimal abstractions very efficiently. We call this approach *corpus-guided top-down search* (CTS).

Corpus-guided top-down search. Like other top-down synthesis approaches, our algorithm explores the space of abstractions by repeatedly refining abstractions with holes, as illustrated in Fig. 2A. We call abstractions with holes *partial abstractions* in contrast to *complete abstractions* which have no holes. Our top-down algorithm searches over abstraction *bodies*, so to synthesize the optimal abstraction fn0 in our running example, we synthesize its body: $(+ 3 (* \alpha \beta))$.

We say that a partial abstraction *matches* at a subtree in the corpus if there's a possible assignment to the holes and arguments that yields the subtree. For example, consider the subtree $(+ 3 (* (+ 2 4) 2))$ from the first program in our running example, shown as a syntax tree at the bottom of figure 2B. The partial abstraction $(+ 3 ??)$ shown in blue matches here with the hole $?? = (* (+ 2 4) 2)$, and the complete abstraction $(+ 3 (* \alpha \beta))$ matches here with $\alpha = (+ 2 4)$ and $\beta = 2$. For complete abstractions, matching corresponds to being able to use the abstraction to rewrite this subtree, resulting in compression. We refer to the set of subtrees at which a complete or partial abstraction matches as the set of *match locations*.

In traditional top-down synthesis, a branch of search can be safely pruned by proving that a program satisfying the specification cannot exist in that branch. In CTS, we can safely prune a branch of search if we can prove that it cannot contain the optimal abstraction. One way to prove this is by computing an upper bound on the utility of abstractions in the branch, and discarding the branch if we have previously found an abstraction with a utility that is higher than this bound. One way to compute this upper bound is to overapproximate the set of match locations, then upper bound the compressive utility gain from each match location.

Our key observation to compute this bound is that during search, expanding a hole in a partial abstraction yields an abstraction that is more precise and thus matches at a *subset* of the locations

197
198

(A) Learning abstractions

(B) Rewriting with abstractions

199
200

Primitives

shapes
c (circle)
r (rectangle)
l (line)algebra
+ (add)
- (subtract)
* (multiply)
/ (divide)trigonometry
cos (cosine)
sin (sin)
pi (pi)utility
C (connect)
T (transform)
M (matrix)
repeatBase
DSL

(C (C (T (repeat (T 1 (N 1 0 -0.5 (/ 0.5 (* pi 8)))) 5 (M 1 (/ (* 2 pi) 8) 0 0)) (W 4 0 0 0)) (T (T c (M 2 0 0 0) 0) (M 1 0 0 0)) (repeat (T (T c (M 0.5 0 0 0)) (M 1 0 0 0) 1.5 (cos (/ pi 4)))) (* 1.5 (sin (/ pi 4)))))) 6 (M 1 (/ (* 2 pi) 6) 0 0)))

201
202

Scaled N-gon

fn_0 2 6

fn_0 4 6

fn_0 2 8

fn_0 4 8

+ fn_0

(C (C (fn_0 4 8) (T (T c (M 2 0 0 0)) (M 1 0 0 0))) (repeat (T (T c (M 0.5 0 0 0)) (M 1 0 0 0) (* 1.5 (cos (/ pi 4)))) (* 1.5 (sin (/ pi 4)))))) 6 (M 1 (/ (* 2 pi) 6) 0 0)))

203
204

Scaled circle

fn_1 2 0 5

fn_2 2 1

fn_2 2 2

fn_2 4

+ fn_2

(C (C (fn_0 4 8) (fn_2 1)) (fn_1 c 1.5 6 (M 1 (/ (* 2 pi) 6) 0 0)))

205
206

Scaled circle nested inside generic shape

fn_7 1 (fn_2 4)

fn_7 2 (fn_2 4)

fn_7 1 (fn_0 4 6)

fn_7 2 (fn_0 4 8)

+ fn_7

(fn_7 1 (fn_0 4 8) (fn_6 1.5 c 6))

207
208

Scaled circle nested inside Length-4 N-gon

fn_13 14 6 1

fn_14 6 2

fn_14 8 1

fn_14 8 2

+ fn_14

(fn_14 8 1 (fn_10 c 6))

209
210211
212213
214

Fig. 3. Visualization of STITCH library learning on the nuts-bolts subdomain from Wong et al. 2022. (A) From the base DSL primitives (top row), STITCH iteratively discovers a series of abstractions that compress programs in the domain. Arrows demonstrate how abstractions from selected iterations build on one another to achieve increasingly higher-level behaviors. (B) Rewriting a single item from the domain with the cumulative benefit of discovered abstractions yields increasingly compact expressions. Colors indicate correspondence between object parts and program fragments: orange = outer octagon, green = ring of six circles, purple = inner circle.

222

that the original matched. This *subset observation* is depicted in Fig. 2 where the refinement of the partial abstraction $(+ 3 ??)$ (blue) to the goal $(+ 3 (* \alpha \beta))$ (green) results in a larger abstraction that matches at a subset of the locations (match locations shown at the top of Fig. 2B). An important consequence of this is that the match locations of a partial abstraction serves as an over-approximation of the match locations of any abstraction in this branch of top-down search.

To upper bound the compressive utility gained by rewriting at a match location, notice that the most compressive abstraction that matches at a subtree is the constant abstraction corresponding to the subtree itself. For example, the largest possible abstraction at $(+ 3 (* (+ 2 4) 2))$ is just $(+ 3 (* (+ 2 4) 2))$. In other words, the size of an expression is a strict upper bound on how much compression can be achieved. Thus we get the upper bound for some partial abstraction $A_{??}$:

$$U_{\text{upperbound}}(A_{??}) = \sum_{e \in \text{matches}(A_{??})} \text{size}(e) \quad (4)$$

233
234

Equipped with this bound, our algorithm performs a branch-and-bound style top-down search, where at each step of search it discards all partial abstractions that have utility upper bounds that are less than the utility of the best complete abstraction found so far.

Our full algorithm presented in Section 3 has some additional complexity. It handles rewriting in the presence of variables soundly, it uses an exact utility function that accounts for additional compression gained by using the same variable in multiple places, and it handles situations where match locations overlap and are mutually exclusive. We also introduce two other important forms of pruning while maintaining the optimality of the algorithm, and we use the upper bound as a heuristic to prioritize more promising branches of search first.

245

Building up abstraction libraries. The top-down search algorithm described above yields a single abstraction. However, we can easily run this algorithm for multiple iterations on a corpus of programs to build up an entire *library of abstractions*. Fig. 3 illustrates the power of this kind of iterative library learning, which interleaves compression and rewriting. At each iteration, STITCH discovers a single abstraction that is used to rewrite the entire corpus of programs. Successive iterations therefore yield abstractions that build hierarchically on one another, achieving increasingly higher-level behaviors. As the library grows to contain richer and more complex abstractions, individual programs shrink into compact expressions.

Structure of the paper. In the subsequent sections, we formalize our corpus-guided top-down search algorithm (Section 3), its application to the problem of compression (Section 4), and how it may be layered on top of data structures such as version spaces (Section 5). We then report experimental results (Section 6) showcasing both diverse library learning settings as well as ablations of STITCH’s search mechanisms. Finally, we conclude by situating STITCH within the landscape of related (Section 7) and future (Section 8) work in the areas of library learning and program synthesis.

3 CORPUS-GUIDED TOP-DOWN SEARCH

We first provide the definitions necessary to understand our problem then introduce our algorithm.

Grammar. Our algorithm operates on lambda-calculus expressions with variables represented through de Bruijn indices [de Bruijn 1972]; expressions come from a context-free grammar of the form $e ::= \lambda. e' | e' e'' | \$i | t$, where $t \in \mathcal{G}_{\text{sym}}$ refers to the set of built-in primitives in the domain-specific language. For example, in an arithmetic domain \mathcal{G}_{sym} would include operators like $+$ and constants like 3. We say that e' is a subexpression of e if $e = C[e']$, where C is a context as defined in the contextual semantics of Felleisen and Hieb [1992]. An expression is *closed* if it has no free variables, in which case the expression is also a *program*. A set of programs \mathcal{P} is a *corpus*.

Abstraction. Given a grammar \mathcal{G} , we define the *abstraction grammar* \mathcal{G}_A as \mathcal{G} extended to include *abstraction variables*, denoted by Greek letters α, β , etc. Formally, $A ::= \lambda. A' | A' A'' | \$i | t | \alpha$ where $t \in \mathcal{G}_{\text{sym}}$. A term A from this grammar represents the *body* of an abstraction; e.g. the abstraction $\lambda\alpha. \lambda\beta. (+ \alpha \beta)$ is simply represented by the term $(+ \alpha \beta)$ from the language of \mathcal{G}_A .

Partial abstraction. A *partial abstraction* $A_{??}$ is an abstraction that can additionally include *holes*. A *hole*, denoted by $??$, represents an unfinished subtree of the abstraction. Thus, $A_{??} ::= A | ?? | \lambda. A'_{??} | A'_{??} A''_{??}$. Any abstraction is thus also a partial abstraction. Given a grammar \mathcal{G} , the grammar of partial abstractions is denoted by $\mathcal{G}_{A_{??}}$. Each hole in a partial abstraction A can be referred to by a *unique* index, which can be explicitly written as $??_i$.

Lambda-aware unification. We introduce *lambda-aware unification* as a modification of traditional unification adapted to our algorithm. LAMBDAUNIFY(A, e) returns a map from abstraction variables and holes to expressions $[\alpha_i \rightarrow e'_i, \dots, ??_j \rightarrow u'_j, \dots]$ such that

$$(\lambda\alpha_i. \dots \lambda??_j. \dots A) e'_i \dots u'_j \dots = e \quad (5)$$

through beta reduction. If no such mapping exists, LAMBDAUNIFY fails. The key challenge in defining LAMBDAUNIFY(A, e) is that the expression e may be deep inside a program written using de Bruijn indices, so the function must perform some index arithmetic in order to generate its output mapping.

The definition of LAMBDAUNIFY is presented in Fig. 4 (left). In U-APP, MERGE(l_1, l_2) merges two $[\alpha_i \rightarrow e_i]$ mappings to create a new mapping that includes all bindings from l_1 and l_2 , and fails

| | |
|--|--|
| 295 U-ABSVAR $\frac{}{\text{LAMBDAUNIFY}(\alpha, e) \rightsquigarrow [\alpha \rightarrow e]}$ | |
| 296 | |
| 297 U-HOLE $\frac{}{\text{LAMBDAUNIFY}(??_i, e) \rightsquigarrow [??_i \rightarrow e]}$ | $\downarrow e = \downarrow_0 e$ |
| 298 | |
| 299 | $\downarrow_d \lambda.b = \lambda. \downarrow_{d+1} b$ |
| 300 | $\downarrow_d (fx) = (\downarrow_d f \downarrow_d x)$ |
| 301 U-APP $\frac{l = \text{MERGE}(l_1, l_2)}{\text{LAMBDAUNIFY}((A_1 A_2), (e_1 e_2)) \rightsquigarrow l}$ | $\downarrow_d \$i = \begin{cases} \$i, & \text{if } i < d \\ \$i-1, & \text{if } i > d \\ \&(i-1), & \text{if } i = d \end{cases}$ |
| 302 | |
| 303 | $\downarrow_d \&i = \&(i-1),$ |
| 304 | $\downarrow_d t = t, \text{ if } t \text{ is a primitive (i.e. } t \in \mathcal{G}_{\text{sym}}\text{)}$ |
| 305 U-LAM $\frac{l = \text{DOWNSHIFTALL}(l')}{\text{LAMBDAUNIFY}((\lambda. A), (\lambda. e)) \rightsquigarrow l}$ | |
| 306 | |
| 307 U-SAME $\frac{}{\text{LAMBDAUNIFY}(e, e) \rightsquigarrow []}$ | |
| 308 | |

Fig. 4. (Left) Inference rules for lambda-aware unification. (Right) Definition of the downshifting operator.

312 if the same abstraction variable α maps to different expressions in l_1 and l_2 . The rule U-SAME
 313 applies only when the abstraction argument to LAMBDAUNIFY is an expression (i.e. it contains no
 314 holes or abstraction variables) and is syntactically identical to the expression it is being unified
 315 with. In U-LAM, DOWNSHIFTALL returns a new mapping with all abstracted expressions e and hole
 316 expressions u downshifted by 1 using the \downarrow operator presented in Figure 4 (right), which has been
 317 modified from the traditional de Bruijn index operator.

318 This operator recursively downshifts all free variables in an expression by 1, extended to switch
 319 to $\&i$ variables when a downshift would otherwise convert a free variable to a (different, incorrect)
 320 bound variable. If an abstracted expression e contains $\&i$ variables at this point then the procedure
 321 fails, as $\&i$ variables represent references to lambdas present within the body of the abstraction.
 322 Notably, the beta reduction in Eq. (5) must account for hole expressions u_j that contain $\&i$ indices.
 323 We thus use a modified form of beta reduction in which $\&i$ indices are *always* upshifted, and at the
 324 very end of beta reduction we replace each $\&i$ with $\$i$. With this modification of beta reduction,
 325 any ampersands with negative indices will ultimately be shifted back to positive indices during
 326 reduction.

327 To understand the definition of the \downarrow operator and the motivation for the $\&i$ indices, consider
 328 what happens when you run.

$$\text{LAMBDAUNIFY}(\lambda. f ??_0, \lambda. e) \quad (6)$$

331 By Eq. (5), the goal is to produce a mapping $[??_0 \rightarrow u'_0]$ such than when replacing $??_0$ with u'_0 it
 332 produces $(\lambda. e)$. In other words, $(\lambda. \lambda. f \$1) u'_0 = \lambda. e$. Now, suppose

$$\text{LAMBDAUNIFY}(f ??_0, e) \rightsquigarrow [??_0 \rightarrow \lambda. \$i] \quad (7)$$

333 There are three possibilities for i that need to be considered, corresponding to the 3 cases in the
 334 definition of $\downarrow_d \$i$. In the first case, when $i = 0$, it means that $\$i$ in Eq. (7) is bound to the lambda
 335 in the return mapping. In this case, DOWNSHIFTALL in U-LAM should not do anything because if
 336 we can replace $??_0$ with $(\lambda. \$0)$ in $(f ??_0)$ to produce e , the same replacement in $(\lambda. f ??_0)$ will
 337 produce $(\lambda. e)$.

338 The second case is when $i \geq 2$; this means that e has some variables that were defined outside
 339 of it and they have been captured by $??_0$. For example, if $i = 2$ it means that $(\lambda. f \$0)(\lambda. \$2) = e$,
 340 which implies that $e = f (\lambda. \$2)$. In this case, the solution to Eq. (6) is $[??_0 \rightarrow \lambda. \$1]$, since

344 $(\lambda. \lambda. f \$1)(\lambda. \$1) = \lambda. e$; in other words, when computing LAMBDAUNIFY($\lambda. f ??_0, \lambda. e$) from
345 Eq. (7), \$i has to be downshifted to account for the fact that it will be substituted inside an additional
346 lambda.

347 The third case, when $i = 1$, is more problematic. In this case, we have that $(\lambda. f \$0)(\lambda. \$1) = e$, which implies that $e = f (\lambda. \$1)$. This creates a problem, since there is no i such that $(\lambda. \lambda. f \$1)(\lambda. \$i) = \lambda. e$; the $\$i$ needs to refer to the lambda surrounding e . The downshift operator addresses this through the special treatment of the & index.

351 **Match locations.** The set of *match locations* of a partial abstraction $A_{??}$ in a corpus \mathcal{P} , denoted
352 MATCHES($\mathcal{P}, A_{??}$), is the maximal set of context-expression pairs $\{(C_1, e_1), (C_2, e_2), \dots, (C_k, e_k)\}$
353 such that $\forall k. p = C_k [e_k]$ and $p \in \mathcal{P}$ and LAMBDAUNIFY($e_k, A_{??}$) succeeds. In Section 3.1 we explain
354 how we maintain this set of matches incrementally.
355

356 **Rewrite Strategies.** A corpus \mathcal{P} can be *rewritten* to use a complete abstraction A as follows. We introduce a new terminal symbol t_A into the symbol grammar \mathcal{G}_{sym} to represent the abstraction, and consider it semantically equivalent to $(\lambda \alpha_0. \dots \lambda \alpha_k. A)$. We then re-express \mathcal{P} in terms of t_A as follows. We can replace the match location $(C, e) \in \text{MATCHES}(\mathcal{P}, A)$ with t_A applied with the argument assignments $[\alpha_i \rightarrow e_i] = \text{LAMBDAUNIFY}(A, e)$ to yield a semantically equivalent expression, $C[(t_A e_0 \dots e_k)]$. Note that since this is a complete abstraction, it contains no holes.

362 One complication is that rewriting at one match location may *preclude* rewriting at another match location if they overlap with each other. A *rewrite strategy* \mathcal{R} is a procedure for selecting a subset of $\text{MATCHES}(\mathcal{P}, A)$ to rewrite at, in which no matches are mutually exclusive. We refer to this set as REWRITELocations $_{\mathcal{R}}(\mathcal{P}, A)$. We refer to rewriting a corpus \mathcal{P} under an abstraction A with rewrite strategy \mathcal{R} by REWRITE $_{\mathcal{R}}(\mathcal{P}, A)$. We refer to rewriting at only a single particular match location m with abstraction A by REWRITEONE(m, A).

368 **Utility.** CTS optimizes a user-defined utility function $U_{\mathcal{P}, \mathcal{R}}(A)$ which scores an abstraction A given a corpus \mathcal{P} and rewrite strategy \mathcal{R} . There are no strict constraints on the form or properties of the utility function. While in Section 4 we will focus on compressive utility functions, our framework does not generally require this.

373 We say that a rewrite strategy is *optimal* with respect to a utility function if the rewrite strategy chooses locations to rewrite at such that the utility is maximized. A naïve optimal rewrite strategy exhaustively checks the utility of rewriting with each of the $2^{|\text{MATCHES}(\mathcal{P}, A)|}$ possible subsets of $\text{MATCHES}(\mathcal{P}, A)$, however depending on the specific utility function computing an optimal or approximately optimal strategy may be significantly more computationally tractable. In Section 4 we detail the rewrite strategy used by STITCH, which is an optimal, linear-time rewrite strategy for compressive utility functions based on dynamic programming.
379

380 3.1 Algorithm

381 Given a corpus \mathcal{P} , rewrite strategy \mathcal{R} , and utility function $U_{\mathcal{P}, \mathcal{R}}(A)$, the objective of CTS is to find the abstraction A that maximizes the utility $U_{\mathcal{P}, \mathcal{R}}(A)$.

384 **Expansion.** We construct the body of a partial abstraction $A_{??}$ through a series of top-down expansions starting from the trivial partial abstraction $??$. Given a partial abstraction $A_{??}$ we can expand a hole in this partial abstraction using any production rule from the partial abstraction grammar $\mathcal{G}_{A_{??}}$, yielding a new abstraction $A'_{??}$. We denote this expansion by $A_{??} \rightarrow A'_{??}$. In Figure 2A our overall top-down search is depicted as a series of expansions. We say that a complete abstraction A can be *derived* from a partial abstraction $A_{??}$, denoted $A_{??} \rightarrow^* A$, if A is in the transitive reflexive closure of the expansion operation, i.e. there exists a sequence of expansions from $A_{??}$ to A . Any abstraction can be derived from $??$.

A naïve approach. The goal of our algorithm is to find the maximum-utility abstraction. One simple, inefficient approach to this is to simply enumerate the *entire* space of abstractions through top down synthesis and return the one with the highest utility. This naïve approach maintains a queue of partial abstractions initialized with just ?? . At each step of the algorithm it pops an abstraction from the queue, chooses a hole in it and expands that hole using each possible production rule. Whenever an expansion produces a new partial abstraction it pushes it to the queue, and when it produces a complete abstraction it calculates its utility and updates the best abstraction found so far. Since an abstraction cannot match a program that is smaller than it, the algorithm can stop expanding when all expansions would lead to abstractions larger than the largest program in the corpus. This algorithm will enumerate all possible abstractions exactly once each.

Introducing subsumption pruning. While the naïve approach will find the optimal abstraction, it is extremely inefficient. We can improve on this by allowing the algorithm to choose to *prune* a partial abstraction instead of expanding it, in which case it simply discards the abstraction and chooses another from the queue. Of course, to maintain optimality we must be certain that we *never* prune the branch of search containing the optimal abstraction. To formalize safe pruning, we introduce the notions of *covering* and *subsumption*.

A complete abstraction A'' *covers* another complete abstraction A' , written $\text{covers}(A'', A')$, if the utility of A'' is greater than that of A' . A partial abstraction $A'_{??}$ is *subsumed* by another partial abstraction $A''_{??}$ if and only if every complete abstraction that is derivable from $A'_{??}$ is covered by a complete abstraction that is derivable from $A''_{??}$.

Formally, $A''_{??}$ subsumes $A'_{??}$, iff $\forall A'. A'_{??} \rightarrow^* A' \implies \exists A''. A''_{??} \rightarrow^* A'' \wedge \text{covers}(A'', A')$.

Equipped with this formalism, we claim that it is safe to prune an abstraction $A'_{??}$ if we know that there exists some $A''_{??}$ which subsumes it. $A''_{??}$ need not have been enumerated yet during search and may even have been pruned; simply proving its existence is enough.

LEMMA 3.1. *Naïve search augmented with subsumption pruning still finds the optimal abstraction.*

PROOF. We proceed with a proof by induction. We seek to prove the equivalent statement that the optimum is never pruned and thus it will be enumerated by the search.

Our inductive hypothesis is that the optimum has not yet been pruned. In the base case this is trivially true, since no pruning has yet taken place. In the inductive step we must prove that a step of pruning maintains the validity of the inductive hypothesis. Recall that a step of pruning will discard some partial abstraction $A'_{??}$ which is subsumed by another abstraction $A''_{??}$. Suppose then for the sake of a contradiction that we did prune the optimum in this step; then the optimum must have been derivable from $A'_{??}$. By subsumption we know that all abstractions derivable from $A'_{??}$ must be covered by some abstraction derivable from $A''_{??}$. Hence, there must exist some abstraction that covers the optimum. However, since the utility of the optimum is greater than or equal to that of all other abstractions, there is by definition *no* such abstraction; we have thus arrived at a contradiction. \square

This lemma ensures the safety of composing together pruning strategies, as long as they all are forms of subsumption pruning, since no instance of subsumption pruning will remove the optimum. Determining which partial abstractions are subsumed by which others is specific to the utility function being used, and in Section 4 we identify two instances of subsumption used when instantiating this framework for a compression-based utility.

Branch-and-bound pruning. We further improve this algorithm to employ branch-and-bound style pruning, in which we bound the maximum utility that can be obtained in a branch of search, and discard the branch if we have previously found a complete abstraction with higher utility

than this bound. Our algorithm is generic over the upper bound function $\bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??})$ which upper bounds the utility of *any complete abstraction* A that can be derived from $A_{??}$.

Formally, the bound must satisfy $\forall A. A_{??} \rightarrow^* A \implies U_{\mathcal{P}, \mathcal{R}}(A) \leq \bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??})$. While this could trivially be satisfied with $\bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??}) = \infty$ for any choice of a utility function, a tighter bound will allow for more pruning. The soundness of branch-and-bound pruning in maintaining the optimality of the solution trivially follows from the fact that the pruned branch only contains complete abstractions that are at most as high-utility as the best abstraction found so far.

While we defer to Section 4 to construct the upper bound used by STITCH, it is worth mentioning a key insight that helps in constructing tight bounds for many utility functions. When bounding the utility of $A_{??}$, it is useful to bound the set of possible locations where rewrites can occur for any abstraction A derived from $A_{??}$. The following lemma provides such a bound:

LEMMA 3.2. *MATCHES($\mathcal{P}, A_{??}$) is an upper bound on the set of locations where rewrites can occur in any A derived from $A_{??}$.*

This follows from the fact that as partial abstractions are expanded they become more precise and thus match at a subset of the locations. Since rewriting only happens at a subset of match locations, MATCHES($\mathcal{P}, A_{??}$) serves as a bound on the set of locations where rewrites can occur.

An important consequence of Lemma 3.2 is that if a partial abstraction matches at zero locations, then all abstractions derived from it will have zero rewriting locations, so no rewriting can occur. Such branches can therefore be safely pruned.

Improving the search order. Finally, without sacrificing the optimality of this algorithm we can heuristically guide the order in which the space of abstractions is explored. The queue of partial abstractions used in top-down search can be replaced with a priority queue ordered by the upper bound. This way, the algorithm will first explore more promising, higher-bound branches of search, narrowing in on the optimal abstraction more quickly. While the algorithm is flexible to which heuristic is used, we find the upper bound heuristic works well in practice.

Efficient incremental matching. When expanding an abstraction $A_{??}$ to a new abstraction $A'_{??}$, it is easy to compute MATCHES($\mathcal{P}, A'_{??}$) since we know it is a subset of MATCHES($\mathcal{P}, A_{??}$). Thus, there is no need to perform matching from scratch against every subtree in the corpus to compute the match locations. Instead, we can simply inspect the relevant subtree at each match location of the original abstraction to see which match locations will be preserved by a given expansion. In fact, except when expanding into an abstraction variable α , the sets of match locations obtained by different expansions of a hole will be disjoint.

When expanding to an abstraction variable α , if α is an existing abstraction variable then this is a situation where the same variable is being used in more than one place, as in the *square* abstraction $(\lambda\alpha. * \alpha \alpha)$. In this case we restrict the match locations to the subset of locations where within the location all instances of α are bound to syntactically identical subtrees.

Additionally, if the user provides a maximum arity limit, then an expansion that causes the abstraction to exceed this limit can be discarded.

Algorithm summary. The CTS algorithm combines all of the aforementioned optimizations in a top-down search algorithm with pruning. In summary, CTS takes as input a corpus \mathcal{P} , a rewrite strategy \mathcal{R} , a utility function $U_{\mathcal{P}, \mathcal{R}}(A)$, and an upper bound function $\bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??})$, and returns the *optimal abstraction* with respect to the utility function. CTS performs a top-down search over partial abstractions and prunes branches of search when partial abstractions match at zero locations or are eliminated by branch-and-bound pruning or subsumption pruning. CTS is made further efficient

491 through the aforementioned *search-order heuristic* and *efficient incremental matching*. None of these
 492 optimizations sacrifice the optimality of the abstraction found.

493 Since the algorithm maintains a best abstraction so-far, it can also be terminated early, making it
 494 an *any-time algorithm*. We also note that this algorithm is amenable to parallelization, as different
 495 workers can process items from the priority queue in parallel, especially since branch-and-bound
 496 pruning remains sound even when best abstraction found so-far is not always up-to-date across
 497 workers. Finally, to learn a library of abstractions, CTS can be run repeatedly, adding one abstraction
 498 at a time to the library and rewriting the corpus with each abstraction as it is learned before running
 499 CTS again. A listing of the full algorithm is provided in Appendix A.²

500 4 APPLYING CORPUS-GUIDED TOP-DOWN SEARCH TO COMPRESSION

501 Having presented the general framework and algorithm of corpus-guided top-down search, we
 502 now instantiate this framework for optimizing a *compression* metric.

503 4.1 Utility

504 In compression we seek to minimize some measure of the size, or more generally the *cost*, of a
 505 corpus of programs after rewriting them with a new abstraction. As in prior work [Ellis et al. 2021]
 506 we penalize large abstractions by including abstraction size in the utility. The compressive utility
 507 function for a corpus \mathcal{P} and abstraction A is given below.

$$511 \quad U_{\mathcal{P}, \mathcal{R}}(A) \triangleq -\text{cost}(A) + \text{cost}(\mathcal{P}) - \text{cost}(\text{REWRITE}_{\mathcal{R}}(\mathcal{P}, A)) \quad (8)$$

512 Here, $\text{cost}(\cdot)$ is a cost function of the following form:

$$\begin{aligned} 513 \quad \text{cost}(\lambda. e') &= \text{cost}_{\lambda} + \text{cost}(e') \\ 514 \quad \text{cost}(e'e'') &= \text{cost}_{\text{app}} + \text{cost}(e') + \text{cost}(e'') \\ 515 \quad \text{cost}(\$i) &= \text{cost}_{\$i} \\ 516 \quad \text{cost}(t) &= \text{cost}_t(t), \text{ for } t \in \mathcal{G}_{\text{sym}} \\ 517 \quad \text{cost}(\alpha) &= \text{cost}_{\alpha} \end{aligned} \quad (9)$$

518 where cost_{λ} , cost_{app} , $\text{cost}_{\$i}$, and cost_{α} are non-negative constants, and $\text{cost}_t(t)$ is a mapping
 519 from grammar primitives to their (non-negative) costs. Finally, we introduce $\text{cost}^*(\cdot)$ as a version
 520 of $\text{cost}(\cdot)$ where $\text{cost}_{\alpha} = 0$.

521 We can equivalently construct the utility in Eq. (8) by summing over the compression gained
 522 from performing each rewrite separately, as given below.

$$523 \quad U_{\mathcal{P}, \mathcal{R}}(A) = -\text{cost}(A) + \sum_{e \in \text{REWRITELOCATIONS}_{\mathcal{R}}(\mathcal{P}, A)} \text{cost}(e) - \text{cost}(\text{REWRITEONE}(e, A)) \quad (10)$$

524 By reasoning about the way that rewriting transforms a program, this utility can be broken
 525 down even further:

$$526 \quad U_{\mathcal{P}, \mathcal{R}}(A) = -\text{cost}(A) + \sum_{e \in \text{REWRITELOCATIONS}_{\mathcal{R}}(\mathcal{P}, A)} U_{\text{local}}(A, e) \quad (11)$$

527
 528 ²[supplement/appendix.pdf](#)

$$\begin{aligned}
 U_{local}(A, e) = & \underbrace{\text{cost}^*(A)}_{\text{abstraction size}} - (\underbrace{\text{cost}_t(t_A) + \text{cost}_{app} \cdot \text{arity}(A)}_{\text{application utility (negative)}}) + \underbrace{\sum_{[\alpha \rightarrow e'] \in \text{args}(A, e)} (\text{usages}(\alpha) - 1) \cdot \text{cost}(e')}_{\text{multiuse utility}}
 \end{aligned} \tag{12}$$

where t_A is the new primitive corresponding to abstraction A . This form of the utility function can be efficiently computed without explicitly performing any rewrites; and it sheds light on the different sources of compression. The three main terms in this expression are (1) the *abstraction size* that comes from the shared structure that is removed, (2) the negative *application utility* that comes from introducing the new primitive and lambda calculus app nonterminals to apply it to each argument, and (3) the *multiuse utility* which comes from only needing to pass in a single copy of an argument that might be used in more than one place in the body. We emphasize that this form of the utility function is equivalent to the original definition based on rewriting given in Eq. (8).

4.2 Upper bounding the utility

We seek an upper bound function $\bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??})$ such that for any A derived from $A_{??}$, $\bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??}) \geq U_{\mathcal{P}, \mathcal{R}}(A)$. We begin from the decomposition of the utility function given in Section 4.1. Since costs are always non-negative, we can upper bound this by dropping the $-\text{cost}_\alpha(A)$ term as well as the negative term within the sum:

$$U_{\mathcal{P}, \mathcal{R}}(A) \leq \sum_{e \in \text{REWRITELocations}_{\mathcal{R}}(\mathcal{P}, A)} \text{cost}(e) \tag{13}$$

Intuitively, dropping the negative term from the sum is equivalent to assuming we compressed the cost of this location all the way down to cost 0. We can also bound $\text{REWRITELocations}_{\mathcal{R}}(\mathcal{P}, A)$ as $\text{MATCHES}(\mathcal{P}, A_{??})$ using Lemma 3.2, yielding our final upper bound in terms of $A_{??}$:

$$\bar{U}_{\mathcal{P}, \mathcal{R}}(A_{??}) \triangleq \sum_{e \in \text{MATCHES}(\mathcal{P}, A_{??})} \text{cost}(e) \tag{14}$$

4.3 Subsumption pruning

We identify two additional forms of subsumption pruning compatible with a compressive utility. The first is *argument capture*: when a partial abstraction takes the same argument for at least one abstraction variable across all match locations, this abstraction can be discarded. This is because every abstraction derivable is covered by another abstraction which is identical except it has the argument inlined into the body. For example, if an abstraction variable α takes the same argument $(+ 3 5)$ across every match location, there is a superior partial abstraction that simply has $(+ 3 5)$ inlined in its body instead of using an abstraction variable.

The second is *redundant argument elimination*: a partial abstraction can be dropped if it has two abstraction variables that always take the same argument as each other across all match locations, for example if it had variables α and β and $\alpha = \beta = (+ 3 5)$ at one location, $\alpha = \beta = 2$ at another location, and so on for all match locations. This abstraction would then be subsumed by the abstraction that doesn't take β as an argument and instead reuses α in place of β , and can hence be eliminated.

589 4.4 Rewrite strategy

590 STITCH employs a linear-time, optimal rewrite strategy for a compression objective. The goal of an
 591 optimal rewrite strategy is to efficiently select the optimal subset of $\text{MATCHES}(\mathcal{P}, A)$ to perform
 592 rewrites at in order to maximize the utility.

593 The main challenge is that when match locations overlap, the strategy must decide which of
 594 the two to accept. For example, consider the program $(\text{foo} (\text{foo} (\text{foo bar}))$ and the abstraction
 595 $\lambda\alpha. (\text{foo} (\text{foo} \alpha))$. This abstraction matches at the root of the program with $\alpha = (\text{foo bar})$,
 596 resulting in the rewritten program $(t_A (\text{foo bar}))$. However, though the abstraction also matched
 597 at the subtree $(\text{foo} (\text{foo bar}))$ in the original program, this match location is no longer present in
 598 the rewritten program, so only one of the two locations can be chosen by the rewrite strategy.

599 Our approach is a bottom-up dynamic programming algorithm, which begins at the leaves of
 600 the program and proceeds upwards. At each node e , we compute the cumulative utility so far if we
 601 reject a rewrite here (util_R), if we accept a rewrite here (util_A), or if we choose the better of the two
 602 options (util^*):

$$\begin{aligned} 604 \text{util}_R[e] &= \sum_{e' \in \text{children}(e)} \text{util}^*[e'] \\ 605 \\ 606 \text{util}_A[e] &= \begin{cases} 0 & \text{if } e \notin \text{MATCHES}(\mathcal{P}, A) \\ 607 U_{\text{local}}(A, e) + \sum_{e' \in \text{args}(A, e)} \text{util}^*[e'] & \text{otherwise} \end{cases} \\ 608 \\ 609 \text{util}^*[e] &= \max(\text{util}_R[e], \text{util}_A[e]) \end{aligned} \quad (15)$$

612 where $U_{\text{local}}(A, e)$ is the utility gained from a single rewrite as defined in Eq. (12). After calculating
 613 these quantities the rewrite strategy can start from the program root and recurse down the tree,
 614 rewriting at each node where rewriting is optimal (i.e. $\text{util}_A > \text{util}_R$) and then recursing into its
 615 arguments after rewriting. Since all arguments were originally subtrees in the program these
 616 quantities have been calculated for all of them, with the caveat that their de Bruijn indices may
 617 have been shifted. Shifting indices of a subtree does not change whether an abstraction can match
 618 there nor does it change the utility gained from using that abstraction, so this simply requires some
 619 extra bookkeeping to track.

620 This algorithm is optimal by a simple inductive argument - at each step of the dynamic program-
 621 ming problem we can assume that we know the cumulative utility of all children and potential
 622 arguments at this node, so we can use Eq. (15) to calculate the cumulative utility of either rejecting
 623 or accepting the rewrite at this node.

625 5 COMBINING CORPUS-GUIDED TOP-DOWN SEARCH WITH DEDUCTIVE 626 APPROACHES

628 In complex domains, assembling good libraries may involve more than just finding matching code-
 629 templates; sometimes, some initial refactoring is necessary to expose the common structure. For
 630 example, consider learning the abstraction $\lambda\alpha. (* 2 \alpha)$ for doubling integers, given the expressions
 631 $(*) 2 8$, $(*) 7 2$, and $(\text{right-shift} 3 1)$. This is only possible if the system can use the
 632 commutativity of multiplication to rewrite $(*) 7 2$ into $(*) 2 7$, and bitvector properties to
 633 rewrite $(\text{right-shift} 3 1)$ into $(*) 2 3$; such rewrites are not natively supported by CTS.

634 In this section, we discuss how CTS can be combined with refactoring systems based on deductive
 635 rewrites (e.g, e-graphs [Willsey et al. 2021]) to increase its expressivity further. The core idea is
 636 simple: run a rewrite system on the corpus to produce a set of refactorings of the corpus in an

$$\begin{array}{c}
\text{638 RU-AbsVar} \frac{}{\text{REWRITEUNIFY}(\alpha, v) \rightsquigarrow [\alpha \rightarrow v]} \quad \text{RU-HOLE} \frac{}{\text{REWRITEUNIFY}(??_i, v) \rightsquigarrow [??_i \rightarrow v]} \\
\text{639} \\
\text{640 RU-APP} \frac{\text{REWRITEUNIFY}(A_1, v_1) \rightsquigarrow l_1 \quad \text{REWRITEUNIFY}(A_2, v_2) \rightsquigarrow l_2 \quad l = \text{VMERGE}(l_1, l_2)}{\text{REWRITEUNIFY}((A_1 A_2), (v_1 v_2)) \rightsquigarrow l} \\
\text{641} \\
\text{642 RU-LAM} \frac{\text{REWRITEUNIFY}(A, e) \rightsquigarrow l' \quad l = \text{DOWNSHIFTALL}(l')}{\text{REWRITEUNIFY}((\lambda. A), (\lambda. e)) \rightsquigarrow l} \quad \text{RU-SAME} \frac{}{\text{REWRITEUNIFY}(e, e) \rightsquigarrow []} \\
\text{643} \\
\text{644 RU-UNION} \frac{\text{REWRITEUNIFY}(e, v_1) \rightsquigarrow l}{\text{REWRITEUNIFY}(e, v_1 \uplus v_2) \rightsquigarrow l} \\
\text{645} \\
\text{646 VM-SAME} \frac{[\alpha \rightarrow v_1] \in l_1 \quad [\alpha \rightarrow v_2] \in l_2}{[\alpha \rightarrow (v_1 \cap v_2)] \in \text{VMERGE}(l_1, l_2)} \quad \text{VM-DIFF} \frac{[\alpha \rightarrow v] \in l_1 \quad \alpha \notin \text{dom}(l_2)}{[\alpha \rightarrow v] \in \text{VMERGE}(l_1, l_2)} \\
\text{647} \\
\text{648} \\
\text{649} \\
\text{650}
\end{array}$$

Fig. 5. Defining REWRITEUNIFY, which takes as input an abstraction and a version space of possible refactorings, and yields multiple substitutions corresponding to all the ways that the abstraction’s holes and variables can match with programs encoded by the version space.

e-graph or version space, then run CTS over the resulting data structure.³ Intuitively, this will lead to improved performance compared with a purely deductive approach, since the cost of rewriting grows exponentially the more rewrites you make. This is a problem for fully deductive approaches like Ellis et al. 2021 because extracting the abstractions often requires a long chain of rewrites, especially for higher-arity abstractions. However, a small number of rewrites is typically sufficient to expose the underlying commonalities, as it was in the example above; performing only a handful of rewrites and then using CTS to actually extract the library thus avoids the exponential blow-up of computing several rewrites in sequence, while still benefiting from the increase in expressivity afforded by the rewrites.

Version Spaces. To illustrate this hybrid approach, we combine CTS with *version spaces* [Lau et al. 2003; Polozov and Gulwani 2015] representing sets of programs semantically equivalent under the β -reduction rewrite. Version spaces are represented as terms from a grammar obtained by extending the grammar of expressions with the union (\uplus) operator that represents a set of equivalent expressions: $v ::= \emptyset \mid v' \uplus v'' \mid \lambda. v' \mid v' v'' \mid \$i \mid t$. We define a denotation operator, $\llbracket v \rrbracket$, mapping a version space v to a set of terms: for the union operator, $\llbracket v \uplus v' \rrbracket = \llbracket v \rrbracket \cup \llbracket v' \rrbracket$; for applications, $\llbracket v v' \rrbracket = \{e e' : \forall e, e' \in \llbracket v \rrbracket \times \llbracket v' \rrbracket\}$; for lambda abstractions, $\llbracket \lambda v \rrbracket = \{\lambda e : \forall e \in \llbracket v \rrbracket\}$; for the empty set, $\llbracket \emptyset \rrbracket = \emptyset$; and for deBruijn indices and terminals, $\llbracket v \rrbracket = v$. Ellis et al. [2021] give a procedure called $I\beta(e)$, which take as input an expression e and outputs a version space v inverting one step of β -reduction: that is, $e' \rightarrow^\beta e$ iff $e' \in \llbracket I\beta(e) \rrbracket$.

To run CTS on top of this rewriting system requires a generalization of LAMBDAUNIFY: instead of unifying an abstraction with a term, yielding a single substitution, we unify against a *set of terms* (a version space), yielding a set of candidate substitutions. The relation REWRITEUNIFY (Fig. 5) accomplishes this, and allows us to straightforwardly redefine the core STITCH quantities in terms of REWRITEUNIFY. The size of term e after expanding it into $I\beta(e)$ and rewriting it with abstraction A is:

$$\text{cost}(\text{REWRITE}(A, e)) = \min_{\substack{l: \\ \text{REWRITEUNIFY}(A, I\beta(e)) \rightsquigarrow l}} \text{cost}_t(t_A) + \text{cost}_{\text{app}} \cdot \text{arity}(A) + \sum_{v \in \text{range}(l)} \min_{e' \in \llbracket v \rrbracket} \text{cost}(v)$$

³Note that while we have previously presented CTS as operating over syntax trees, the core notions of upper bounds and matching that CTS operates on are not restricted to program trees.

687 And the local utility is then given by:

$$U_{\text{local}}(A, e) = \text{cost}(e) - \text{cost}(\text{REWRITE}(A, e)) \quad (16)$$

690 **Example: Learning map from fold.** Consider learning the higher-order function `map` from two
 691 example programs: doubling a list of numbers (via `(fold (λx.λℓ. (cons (+ x x) ℓ)) nil))`) and
 692 decrementing a list of numbers (via `(fold (λx.λℓ. (cons (- x 1) ℓ)) nil))`). Matching these
 693 programs with the `map` function—`(fold (λx.λℓ. (cons (α x) ℓ)) nil))`—requires re-expressing
 694 them as `(fold (λx.λℓ. (cons ((λz. (+ z z)) x) ℓ)) nil))` and `(fold (λx.λℓ. (cons ((λz.
 695 (- z 1)) x) ℓ)) nil))`, with the substitutions $\alpha \mapsto (\lambda z. (+ z z))$ and $\alpha \mapsto (\lambda z. (- z 1))$
 696 respectively. These two re-expressions are licensed by a β -reduction rewriting step.

6 EXPERIMENTS

698 In this section, we evaluate corpus-guided top-down search for library learning. Specifically, our
 699 evaluation focuses on five hypothesis about the performance of STITCH:
 700

- 702 (1) *STITCH learns libraries of comparable quality to those found by existing deductive library
 703 learning algorithms in prior work, while requiring significantly less resources.* In Section 6.1
 704 we run STITCH on the library learning tasks from [Ellis et al. 2021] and directly compare
 705 STITCH to DreamCoder, the deductive algorithm introduced in that work. We find that
 706 STITCH learns libraries which usually match or exceed the baseline in quality (measured via
 707 a compression metric), while improving the resource efficiency in terms of memory usage
 708 and runtime by 2 and 3-4 orders of magnitude compared to the baseline (respectively).
- 709 (2) *STITCH scales to corpora of programs that contain more and longer programs than would be
 710 tractable with prior work.* In Section 6.2, we evaluate STITCH’s ability to learn libraries within
 711 eight graphics domains from Wong et al. 2022, which are considerably larger and more
 712 complex than have been considered in previous work. We find that STITCH on average obtains
 713 a test set compression ratio of 2.55x-11.57x in 0.28s-134.44s, with a peak memory usage of
 714 11.22MB-715.42MB in these domains. The problems are large enough to time out with the
 715 DreamCoder baseline.
- 716 (3) *STITCH degrades gracefully when resource-constrained.* In Section 6.3 we investigate STITCH’s
 717 performance when run as an *any-time* algorithm, i.e. one that can be terminated early for a
 718 best-effort result if a corpus is too large or there are limits on time or memory. We reuse
 719 the eight graphics domains from Wong et al. 2022 and find that with its heuristic guidance
 720 STITCH converges upon a set of high quality abstractions very early in search, doing so
 721 within 1% of the total search time in 3 out of 8 domains and 10% in all except one.
- 722 (4) *All the elements of STITCH matter.* In Section 6.4, we carry out an ablation study on STITCH and
 723 find the *argument capturing* and *upper bound* pruning methods essential to its performance,
 724 while *redundant argument elimination* also proves useful in certain domains. With all
 725 optimizations disabled, we find that STITCH cannot run in ≤ 90 minutes and ≤ 50 GB of
 726 RAM on any of the domains from Wong et al. 2022.
- 727 (5) *STITCH is complementary to deductive rewrite-based approaches to library learning.* These
 728 prior experiments show the superior runtime performance of STITCH relative to deductive
 729 rewrite systems, but deductive systems have an important advantage over STITCH, the
 730 ability to incorporate arbitrary rewrite rules to expose more commonality among different
 731 programs and in that way discover better libraries. Such deductive approaches are especially
 732 more apt at learning higher-order abstractions. In Section 6.5, we show that this expressivity
 733 gap can be reduced by running STITCH on top of such rewrite systems, allowing it to learn
 734 almost all the same functions they do while still using only <1% of the compute time.

We run all experiments on a machine with two AMD EPYC 7302 processors, 64 CPUs, and 256GB of RAM. We parametrize STITCH’s cost(e) function (as defined in Section 4.1) as follows: $\text{cost}_{\text{app}} = \text{cost}_\lambda = 1$, $\text{cost}_{\$i} = \text{cost}_\alpha = \text{cost}_t(t) = 100$.

6.1 Iterative bootstrapped library learning

Experimental setup. Our first experiment is designed to replicate the experiments in DreamCoder, which is the state-of-the-art in deductive library learning. DreamCoder learns libraries iteratively: the system is initialized with a low-level DSL, and then alternates between synthesizing programs (via a neurally-guided enumerative search) that solve a training corpus of inductive tasks and updating the library of abstractions available to the synthesizer. Traces from the experiments carried out by Ellis et al. 2021 are publicly available⁴ and include all of the intermediate programs that were synthesized as well as the libraries learned from those programs.

In this experiment, we take these traces and evaluate STITCH on each instance where library learning was performed, comparing the quality of the resulting library to the original one found by DreamCoder. We also re-run DreamCoder on these same benchmarks in order to evaluate its resource usage, capturing its runtime and memory usage in the same environment as STITCH.

The library learning algorithm in Ellis et al. 2021 implements a stopping criterion to determine how *many* abstractions to retain on any given set of training programs. In our comparative experiments, we run the DreamCoder baseline first, and then match the number of abstractions learned by STITCH at each iteration to those learned by the baseline under its stopping criterion so that timing comparisons are fair. We record the total time spent both performing abstraction learning and rewriting for both STITCH and DreamCoder.

We replicate experiments on five distinct domains from Ellis et al. 2021:

- **Lists:** a functional programming domain consisting of 108 total inductive tasks;
- **Text:** a string editing domain in the style of FlashFill [Gulwani et al. 2015] consisting of 128 total inductive tasks;
- **LOGO:** a graphics domain consisting of 80 total inductive tasks;
- **Towers:** a block-tower construction domain consisting of 56 total inductive tasks;
- **Physics:** a domain for learning equations corresponding to physical laws from observations of simulated data, consisting of 60 total inductive tasks

Assessing library quality with a compression metric. The standard STITCH configuration optimizes a compression metric that minimizes the size of the programs after being rewritten to use the abstraction. This is a standard metric in program synthesis, since shorter programs are frequently easier to synthesize. Optimizing against this metric is equivalent to maximizing the likelihood of the rewritten programs under a uniform PCFG.

The DreamCoder synthesizer is more sophisticated than simple enumeration; it takes as input a learned typed bigram PCFG and leverages it to synthesize programs more efficiently. When performing compression, it optimizes against this given PCFG in order to find abstractions that will be more profitable for its specific synthesizer. For the purpose of this evaluation, however, we restrict ourselves to the uniform PCFG because the one used by DreamCoder requires programs to be in a particular normal form. Another aspect of DreamCoder relevant to its compression metric is that DreamCoder synthesizes multiple programs that solve the same task and then selects the abstraction that works best on *some* program for a given task. This is expressed formally in the equation below.

$$\text{cost}(A) + \sum_{\substack{p \in \text{task} \\ \text{task}}} \min \text{cost}(\text{REWRITE}(p, A)) \quad (17)$$

⁴https://github.com/mlb2251/compression_benchmark

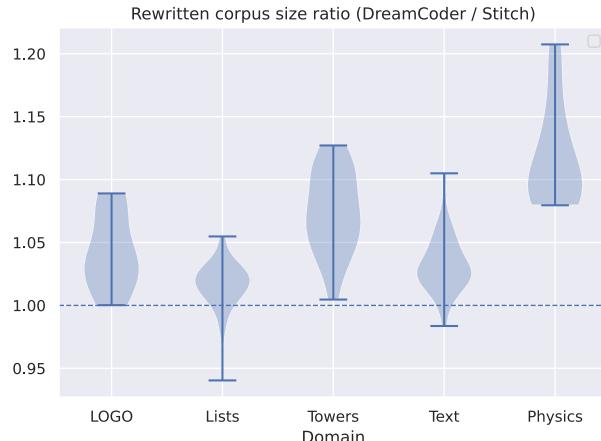


Fig. 6. Compression rates obtained when running STITCH on the five domains considered in 6.1 relative to those of DreamCoder. Higher is better for STITCH: a ratio above 1.0 indicates that STITCH achieves greater compression than DreamCoder. The specific compression metric used in the ratio is given in Eq. (17).

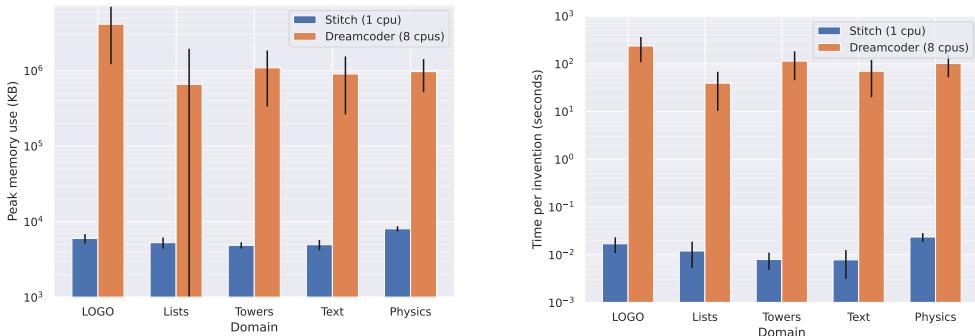


Fig. 7. Peak memory usage of STITCH and DreamCoder while running on the five domains considered in 6.1, averaged over all benchmarks. Lower is better; black lines indicate \pm one standard deviation. Note the logarithmic y-axis.

Fig. 8. Wall-clock time required to find (and rewrite under) one abstraction in each of the five domains from 6.1, averaged over all benchmarks. Lower is better; black lines indicate \pm one standard deviation. Note the logarithmic y-axis.

It is trivial to implement this best-of-task metric in STITCH, so we use this for the comparisons with DreamCoder.

Results. We compare DreamCoder and STITCH for library quality and resource efficiency.

Library quality. We first examine the quality of the libraries learned by both STITCH and DreamCoder using the compression metric in Eq. (17): Fig. 6 shows the ratio between them across all of the benchmarks for each domain. A ratio of 1.0 indicates that programs are the *exact same length* under the libraries learned by DreamCoder and STITCH, a ratio greater than 1.0 indicates that STITCH learns more compressive libraries, and a ratio less than 1.0 indicates that STITCH learns

libraries which are less compressive. For example, a ratio of 1.1 indicates that DreamCoder rewrote to produce a corpus 10% larger than that of STITCH, so it achieved less compression.

These results show that STITCH generally learns libraries of comparable and often greater quality than DreamCoder when matching the number of abstractions learned by the latter. In the *logo*, *towers*, and *physics* domains, STITCH always finds abstractions that are of equal—and often considerably greater—compressive quality than DreamCoder does; in the *list* and *text* domains, STITCH more often than not still obtains better compression, but sometimes loses out to DreamCoder. This is a result of the fact that STITCH cannot learn higher-order abstractions, which are useful in these domains; although we emphasize our focus is on scalability, we will later present an extension of STITCH capable of handling most higher-order functions in Section 6.5, based on the formalism developed in Section 5. Nonetheless, we conclude that **STITCH learns libraries whose quality is comparable to and often better than those found by DreamCoder.**

Resource efficiency. In addition to the quality of the libraries found, we are interested in how the two methods compare in terms of time and space requirements. Since we a priori believe STITCH to be significantly faster, for this evaluation we allow DreamCoder to use 8 CPUs but limit STITCH to single threading⁵. The results are shown in Fig. 7 and Fig. 8; in summary, we find that STITCH takes *tens of milliseconds* to discover abstractions across all five domains—achieving a *3-4 order of magnitude* speed-up over DreamCoder—while also requiring more than *2 orders of magnitude* less memory. We thus conclude that **STITCH is dramatically more efficient than the state-of-the-art deductive baseline** when replicating the iterative library learning experiments of Ellis et al. 2021.

6.2 Large-scale corpus library learning

Experimental setup. While the previous experiment allowed us to benchmark STITCH directly against a state-of-the-art deductive baseline, the iterated learning setting considered by Ellis et al. 2021 only evaluates library learning on relatively small corpora of short programs discovered by the synthesizer. Our second experiment is instead designed to evaluate STITCH in a more traditional learning setting, in which we aim to learn libraries of abstractions from a large corpus of existing programs all at once.

We source our larger-scale program datasets from Wong et al. 2022, which present a series of datasets designed as a benchmark for comparing human-level abstraction learning and graphics program writing against automated synthesis and library learning models. These datasets are divided into two distinct high-level domains (*technical drawings* and *block-tower planning tasks*), each consisting of four distinct subdomains containing 250 programs:

- **Technical drawing domains: nuts and bolts; vehicles; gadgets; furniture:** CAD-like graphics programs that render technical drawings of common objects, written in an initial DSL consisting of looped transformations (scaling, translation, rotation) over simple geometric curves (lines and arcs).
- **Tower construction domains: bridges; cities; houses; castles:** planning programs that construct complex architectures by placing blocks, written in an initial DSL that moves a virtual hand over a canvas and places horizontal and vertical bricks.

We choose these datasets not only for their size and scale (full dataset statistics in Table 1), but also for the complexity of their potential abstractions: Wong et al. 2022 explicitly design their corpora

⁵While this may seem unfair to STITCH, it is worth noting that it would be unlikely to benefit from multithreading when running on the order of milliseconds anyway; DreamCoder, on the other hand, would struggle greatly in these domains without the aid of parallelism.

883 to contain complex hierarchical structures throughout the programs, making them an interesting
 884 setting for library learning.

| 886 Domain | 886 #Programs | 886 Average program length | 886 Average program depth |
|------------------|---------------|----------------------------|---------------------------|
| 888 nuts & bolts | 888 250 | 888 76.03 ± 24.22 | 888 15.18 ± 2.13 |
| 889 gadgets | 889 250 | 889 142.85 ± 87.32 | 889 20.88 ± 2.37 |
| 890 furniture | 890 250 | 890 171.74 ± 48.41 | 890 31.83 ± 5.33 |
| 891 vehicles | 891 250 | 891 141.70 ± 40.23 | 891 21.22 ± 1.35 |
| 892 bridges | 892 250 | 892 137.03 ± 59.71 | 892 92.35 ± 39.80 |
| 893 cities | 893 250 | 893 161.70 ± 55.56 | 893 109.80 ± 37.66 |
| 894 castles | 894 250 | 894 189.09 ± 60.18 | 894 128.27 ± 40.77 |
| 895 houses | 895 250 | 895 168.13 ± 55.75 | 895 114.85 ± 37.60 |

900 Table 1. Summary statistics about the domains from Wong et al. 2022. Program length is the number of
 901 terminal symbols in the program; program depth is the length of the longest path from root to leaf in the
 902 program tree. Both are reported as the mean over the entire dataset \pm one standard deviation.

903
 904 When performing library learning in a synthesis setting by compressing a corpus of solutions, it's
 905 desirable to find abstractions that would be useful for solving *new* tasks, as opposed to abstractions
 906 that overfit to the existing solutions. To evaluate how well the abstractions we learn apply to
 907 heldout programs in the domain, we split the corpora into train and test sets, running STITCH on
 908 the train set and evaluating its compression on the test set. Since Wong et al. 2022 do not present a
 909 train/test split of their datasets, we use stochastic cross validation to evaluate the generalization of
 910 the libraries found by STITCH. For each domain, we randomly sample 80% of the dataset to train
 911 on and reserve the last 20% as a held out test set; we repeat this procedure 50 times. We then ask
 912 STITCH to learn a library consisting of ≤ 10 abstractions with a maximum arity of 3, and average
 913 the results across the different random seeds.
 914

915 To obtain a baseline to compare its performance against, we once again turned to DreamCoder
 916 [Ellis et al. 2021]. However, we found that DreamCoder was unable to discover even a single
 917 abstraction when run directly on any of the datasets from Wong et al. 2022, despite being given
 918 hours of runtime and 256GB of RAM. We also experimented with heavily sub-sampling the training
 919 dataset before passing it to DreamCoder, but failed to find a configuration under which DreamCoder
 920 finds any interesting abstractions at all due to the fact that it immediately blows up on programs
 921 as long as these. As a result, we resort to presenting STITCH's performance metrics without any
 922 baseline to compare against; we stress that this is a direct result of the fact that STITCH is the first
 923 library learning tool that scales to such a challenging setting.

924 **Results.** The results are summarized in Table 2. We find that STITCH scales up to even the most
 925 complex sub-domains, running in ≈ 2 minutes with a peak memory usage way below 1GB on
 926 *castles*. On four out of eight of the domains, STITCH finishes in single-digit seconds and consumes
 927 only tens of megabytes. This stands in stark contrast to DreamCoder, which we were unable to run
 928 on the very simple nuts & bolts domain even with 256GB RAM and several hours worth of compute
 929 budget. These results thus support our claim that **STITCH scales to corpora of programs that**
 930 **would be intractable with prior library learning approaches.**

| 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 | Domain | Compression Ratio | | Runtime (s) | Peak mem. usage (MB) |
|---|------------------|-------------------|-------------------|--------------------|----------------------|
| | | Training set | Test set | | |
| nuts & bolts | 12.00 ± 0.25 | 11.57 ± 0.49 | 0.28 ± 0.07 | 11.22 ± 0.25 | |
| gadgets | 4.03 ± 0.15 | 3.91 ± 0.31 | 2.92 ± 0.44 | 30.48 ± 0.90 | |
| furniture | 4.95 ± 0.07 | 4.85 ± 0.26 | 5.52 ± 0.57 | 31.59 ± 0.60 | |
| vehicles | 4.28 ± 0.14 | 4.14 ± 0.25 | 2.32 ± 0.25 | 27.31 ± 0.84 | |
| bridges | 4.36 ± 0.06 | 3.78 ± 0.14 | 28.50 ± 1.38 | 199.46 ± 10.57 | |
| cities | 3.15 ± 0.05 | 3.06 ± 0.14 | 97.24 ± 3.50 | 431.99 ± 9.90 | |
| castles | 2.57 ± 0.07 | 2.55 ± 0.08 | 134.44 ± 7.76 | 715.42 ± 33.89 | |
| houses | 8.92 ± 0.21 | 8.85 ± 0.57 | 27.53 ± 1.72 | 243.62 ± 2.89 | |

947 Table 2. Results for the large-scale library learning experiment in Sec. 6.2. The compression ratio refers to
 948 how many times smaller the corpus is after rewriting under the learned library compared to the original
 949 corpus; higher is thus better. All results are given as the mean \pm one standard deviation over 50 runs with
 950 different random seeds for the dataset splitting.

951
 952
 953 We hope that by providing our results on these datasets in full, future work in this field will
 954 benefit from having a directly comparable baseline.
 955

956 6.3 Robustness to early search termination

957 *Experimental setup.* This experiment is designed to evaluate how early into the search procedure
 958 STITCH finds what will eventually prove to be the optimal set of abstractions. This is highly relevant
 959 in settings where the set of training programs is too large to run the search to completion. The
 960 experiment showcases one of STITCH’s more subtle strengths: corpus-guided top-down abstraction
 961 search is an *anytime* algorithm, and thus does not need to be run to completion to give useful
 962 results.
 963

964 We re-use the domains from Wong et al. 2022 and once again evaluate the quality of the library
 965 learned (measured in program compression), similarly to what was done in the previous experiment.
 966 However, since we are interested in how quickly STITCH finds *what it perceives to be* the optimal
 967 abstractions, we measure compressivity of the training dataset itself (rather than a held-out test
 968 set) and capture the compression ratio obtained by each *candidate* abstraction found *during* search
 969 (rather than just the compression ratio obtained when search has been run to completion). Thus,
 970 we are able to investigate how early on during the search procedure STITCH converges on a chosen
 971 library. We restrict STITCH to learning a single abstraction with a maximum arity of 3.

972 *Results.* The results are shown in Fig. 9. These results validate our hypothesis that **STITCH is**
 973 **empirically robust to terminating the search procedure early:** in every sub-domain except for
 974 *nuts & bolts* STITCH converges to the optimal library very early on, having only completed a tiny
 975 fraction of the total search.⁶ We believe that this has great importance for STITCH’s applicability in
 976 data-rich settings since it suggests that a nearly-optimal library can be found even if the search
 977

978 ⁶Given that significant attention has already been given to wall-clock run-times of STITCH on similar workloads in 6.1, we
 979 here use the number of nodes explored instead of wall-clock time to ensure deterministic and easily reproducible results.
 980

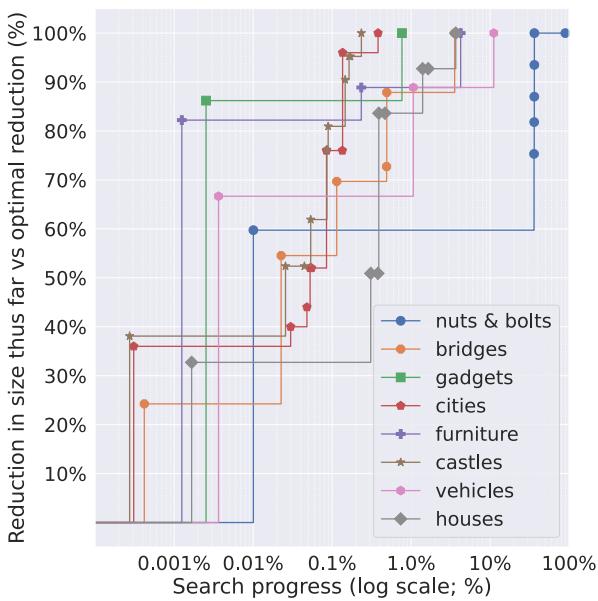


Fig. 9. Reduction in size obtained on the training set when rewritten under the best abstraction found thus far vs. the number of nodes expanded during search. The y-axis is normalized with respect to the optimal abstraction; the x-axis is normalized with respect to the total number of nodes explored by STITCH. Lines thus end earlier the quicker STITCH finds the optimal abstraction; however in all runs STITCH *continues to search* until reaching 100% on the x-axis, exploring the rest of the abstraction space without finding a new best abstraction. Note the logarithmic x-axis.

must be terminated early (e.g. after a fixed amount of time has passed), making early stopping an empirically useful way of speeding up the library learning process.

6.4 Ablation study

Experimental setup. STITCH implements several different optimizations, which we have argued hasten the search for abstractions. To verify that this claim holds in practice, we now carry out a brief ablation study.

Since the space of every possible combination of optimizations is too large to present succinctly, we focus our attention on four ablations:

- **no-arg-capture**, which disables the pruning of abstractions which are only ever used with the exact same set of arguments (and these arguments could therefore just be in-lined for greater compression)
- **no-upper-bound**, which disables the upper-bound based pruning
- **no-redundant-args**, which disables the pruning of multi-argument abstractions that have a redundant argument that could be removed because it is always the same as another argument
- **no-opt**, which disables all of STITCH’s optimizations

To isolate the impact of disabling optimizations, we run a single iteration of abstraction on each of the 8 domains from 6.2 and collect the number of nodes explored during the search. Running a single iteration reduces the impact of compound effects on the results, and focusing on the number

| System \ Domain | bridges | castles | cities | gadgets | furniture | houses | nuts & bolts | vehicles |
|-------------------|---------|---------|--------|---------|-----------|--------|--------------|----------|
| System | | | | | | | | |
| no-arg-capture | ⌚ | ⌚ | ⌚ | 1375.05 | ⌚ | ⌚ | 33.74 | 7.18 |
| no-upper-bound | 12.04 | 23.27 | 27.96 | 241.82 | 63.43 | 35.44 | 159.33 | 302.25 |
| no-redundant-args | 1.92 | 1.45 | 1.39 | 1.01 | 1.00 | 1.01 | 1.04 | 1.00 |
| no-opt | RAM | ⌚ | ⌚ | RAM | ⌚ | ⌚ | RAM | ⌚ |

Table 3. Results from the ablation study. Each cell contains the ratio between the number of nodes explored during search by that particular system and the number of nodes explored by the baseline on the same domain. Lower is better; 1.00 means performance is identical to the baseline. Cells labeled  crashed due to reaching the 50GB virtual address space limit, while those labeled  reached the 90 minute time restriction.

of nodes explored (rather than for example latency) allows for deterministic results. We also fix the maximum arity of abstractions to 3 in all runs, aiming to strike a good balance between compute requirements and how much the optimization will be exposed. We limit each run to 50GB of virtual address space, as well as 90 minutes of compute.

Results. The results are shown in Table 3. We first note as a sanity check that each ablation does indeed lead to reduced performance in general (i.e. explores a search space larger than the baseline does). Furthermore, the results suggest that upper-bound based pruning is the most important in the *nuts & bolts* and *vehicles* domains, while pruning out argument capture abstractions is the most important in the other six domains; it is noteworthy that this latter ablation by itself causes STITCH to reach the memory limit or timeout on more than half of the domains, which suggests that function applications in these domains contain very little variation in the arguments supplied. On the other hand, disabling the pruning of redundant arguments has a relatively modest impact on the size of the search space, but still leads to an almost 2x improvement in the bridge domain.

Perhaps the most important takeaway from this ablation study is that when *all* optimizations are disabled, STITCH fails to find an abstraction within the resource budget on any of the domains. This verifies our hypothesis that corpus-guided pruning of the search space is the key factor involved in making top-down synthesis of abstractions tractable.

6.5 Learning Libraries of Higher-Order Functions

Experimental setup. Our experiments up till now have focused on performant and scalable library learning. This comes at the expense of some expressivity: deductive rewrite systems can, in principle, express broader spaces of refactorings. For example, a rewrite based on β -reduction allows inventing auxiliary λ abstractions, which helps learn higher-order functions: in Ellis et al. 2021, DreamCoder is shown to recover higher-order functions such as `map`, `fold`, `unfold`, `filter`, and `zip_with`, starting from the Y-combinator. This works by constructing a version space which encodes every refactoring which is equivalent up to β -reduction rewrites. But DreamCoder’s coverage comes at steep cost: inverting β -reduction is expensive.

In this experiment, we seek to give evidence that it is possible to make STITCH recover almost all of these higher-order functions by layering it on top of version space obtained after a single step of

1079 DreamCoder's β -reduction-inverting refactoring, following the formalism outlined in Section 5⁷.
 1080 We then compare this modified version of STITCH against DreamCoder on its ability to learn these
 1081 higher-order functions from programs generated by intermediate DreamCoder iterations, and
 1082 measure the runtime of each approach.

1083 *Probabilistic Re-ranking.* While the approach outlined in Section 5 should suffice to layer STITCH
 1084 on top of general version space-based deductive rewrite systems, some extra care needs to be taken
 1085 to combine it with DreamCoder. This is because DreamCoder implements a probabilistic Bayesian
 1086 objective for judging candidate abstractions (exploiting the connection between compression and
 1087 probability [Shannon 1948]), seeking the library L which maximizes $P(L) \prod_{p \in \mathcal{P}} P(p|L)$ for a given
 1088 or learned prior prior $P(L)$ and likelihood $P(p|L)$. STITCH, on the other hand, effectively judges
 1089 compression quality by counting the syntactic size of the programs, as has been discussed at great
 1090 length in Section 4.

1091 To heuristically hybridize these methods, we simply run STITCH on the version spaces as-is
 1092 but then re-score each complete abstraction popped off the queue under the Bayesian objective,
 1093 using DreamCoder's models of the prior and likelihood. To make this integration easier, we re-
 1094 implement STITCH in Python, giving a prototype version called PYSTITCH which is only used for
 1095 this experiment.

| | System | Fold | Unfold | Map | Filter | ZipWith | Time (s) |
|------|------------|-----------|--------|-----|--------|---------|----------|
| 1100 | PYSTITCH | base | ✗ | ✗ | ✗ | ✓ | 254 |
| 1101 | | +Bayes | ✗ | ✗ | ✗ | ✓ | 475 |
| 1102 | | +VS | ✓ | ✓ | ✗ | ✓ | 1103 |
| 1103 | | +Bayes+VS | ✓ | ✓ | ✓ | ✗ | 1616 |
| 1104 | DreamCoder | Step 1 | ✗ | ✗ | ✓ | ✗ | 67 |
| 1105 | | Step 2 | ✗ | ✗ | ✓ | ✓ | 116 |
| 1106 | | Step 3 | ✗ | ✗ | ✓ | ✓ | 2254 |
| 1107 | | Step 4 | ✓ | ✓ | ✓ | ✓ | 228048 |

1112 Table 4. Comparing library learners on functional programming exercises. DreamCoder, n -steps: deductive
 1113 baseline rewriting n steps of β -reduction. PYSTITCH: Python reimplementation of STITCH, which enables
 1114 better interoperability with DreamCoder's version space algebra (+VS) and probabilistic models (+Bayes).
 1115 PYSTITCH+VS+Bayes learns essentially all of the same higher-order functions, using < 1% of the compute. (✓:
 1116 half credit. For filter, PYSTITCH learns a filter-esque function with type $\alpha \rightarrow (\alpha \rightarrow \beta \rightarrow \text{bool}) \rightarrow \text{list}(\beta) \rightarrow$
 1117 $\text{list}(\beta)$.)

1118
 1119
 1120 *Results.* The results are shown in Table 4. We note first the large discrepancy in runtimes; running
 1121 DreamCoder with 4 steps of rewriting takes roughly 2.5 days of compute on this domain, while
 1122 even the slowest version of PYSTITCH still finishes in less than half an hour. In terms of the number
 1123 of higher-order abstractions found, DreamCoder only finds all five when run in its most expensive

1124 ⁷Note that DreamCoder was only used as a convenient base upon which to graft the STITCH; we believe the same hybridization
 1125 recipe should apply to other rewrite frameworks as well, such as e-graphs [Willsey et al. 2021], but leave it to future work
 1126 to investigate this in detail.

1128 configuration; reducing the computation cost quickly decreases its expressivity. For `PYSTITCH`, the
 1129 Bayesian re-ranking alone (`PYSTITCH+Bayes`) does not yield any improvements, while running it
 1130 on the version spaces (`PYSTITCH+VS`) yields 3 out of 5 functions. However, it is only when these
 1131 two adaptations are used in conjunction (`PYSTITCH+Bayes+VS`) that the `STITCH`-based method is
 1132 able to find almost all of the higher-order functions, failing only in that it learns an unconventional
 1133 sibling of the `filter` function.

1134 In summary, we find that running CTS after a single step of version space rewriting and then
 1135 probabilistically re-ranking the results suffices to recover almost all of the core higher-order
 1136 functions that DreamCoder learns, while using <1% of its compute. Since our methodology may
 1137 easily be extended to other version space-based deductive systems, and even e-graph based ones
 1138 with some modifications to the underlying formalism, we thus conclude that **running STITCH on**
 1139 **top of a deductive rewrite system closes the expressivity gap, while retaining superior**
 1140 **performance.**

1141

1142 7 RELATED WORK

1143 `STITCH` is related to two core ideas from prior work: **deductive refactoring and library learning**
 1144 **systems**, which introduce the idea of learning abstractions that capture common structure across
 1145 a set of programs, but have largely been driven by deductive algorithms; and **guided top-down**
 1146 **program synthesis systems**, which use cost functions to guide top-down enumerative search
 1147 over a space of programs, but have largely been used to synthesize whole programs for individual
 1148 tasks in prior work.

1149

1150 7.1 Deductive refactoring and library learning.

1151 Recent work shares `STITCH`'s goal of learning *libraries* of program abstractions which capture
 1152 reusable structure across a corpus of programs [Allamanis and Sutton 2014; Cropper 2019; Dechter
 1153 et al. 2013; Ellis et al. 2021, 2018; Iyer et al. 2019; Jones et al. 2021; Shin et al. 2019; Wong et al. 2021].
 1154 Several of these prior approaches also introduce a utility metric based on *program compression* in
 1155 order to determine the most useful candidate abstractions to retain [Dechter et al. 2013; Ellis et al.
 1156 2021; Iyer et al. 2019; Lázaro-Gredilla et al. 2019; Wong et al. 2021].

1157 Prior work generally follows a bottom-up approach to abstraction learning, combining a bottom-
 1158 up traversal across individual training programs with a second stage to extract shared abstractions
 1159 from across the training. This approach includes systems that work through direct *memoization* of
 1160 subtrees across a corpus of programs [Dechter et al. 2013; Lázaro-Gredilla et al. 2019; Lin et al. 2014];
 1161 *antiunification* (caching tree templates that can be unified with training program syntax trees)
 1162 [Ellis et al. 2018; Henderson 2013; Hwang et al. 2011; Iyer et al. 2019]; or by more sophisticated
 1163 *refactoring* using one or more rewrite rules to expose additional shared structure across training
 1164 programs [Chlipala et al. 2017; Ellis et al. 2021, 2018; Liang et al. 2010].

1165 These latter library learning algorithms draw more generally on *deductive synthesis* approaches
 1166 that apply local rewrite rules in a bottom-up fashion to program trees in order to refactor them –
 1167 historically, to synthesize programs from a declarative specification of desired function [Burstall and
 1168 Darlington 1977; Manna and Waldinger 1980]. Deductive approaches to library learning, however,
 1169 confront fundamental memory and search-time scaling challenges as the corpus size and depth of
 1170 the training programs increases; prior deductive approaches such as [Chlipala et al. 2017; Ellis et al.
 1171 2021, 2018] use version spaces [Lau et al. 2003; Mitchell 1977] to mitigate the memory usage during
 1172 bottom-up abstraction proposal. Still, deductive approaches are challenging to bound and prune
 1173 (unlike the top-down approach we take in `STITCH`), as they generally traverse individual program
 1174 trees locally and must store possible abstraction candidates in memory before the extraction step.

1175

1177 7.2 Guided top-down program synthesis.

1178 STITCH uses a corpus-guided top-down approach to learning library abstractions that is closely
1179 related to recent *guided enumerative synthesis techniques*. This includes methods that leverage
1180 type-based constraints on holes [Feser et al. 2015; Polikarpova et al. 2016], over- and under-
1181 approximations of the behaviors of holes [Chen et al. 2020; Lee et al. 2016], and probabilistic
1182 techniques to heuristically guide the search [Balog et al. 2016; Ellis et al. 2020, 2021; Nye et al.
1183 2021; Shah et al. 2020]. These approaches have largely been applied in to synthesize entire
1184 programs based on input/output examples or another form of specification, in contrast to the
1185 abstraction-learning goal in our work [Allamanis et al. 2018; Balog et al. 2016; Chen et al. 2018;
1186 Ellis et al. 2018; Ganin et al. 2018; Koukoutos et al. 2017]. Like STITCH however, these approaches
1187 sometimes use cost functions (such as the likelihood of a partially enumerated program under a
1188 hand-crafted or learned probabilistic generative model over programs) in order to direct search
1189 towards more desirable program trees. However, STITCH’s cost function leverages a more direct
1190 relationship between partially-enumerated candidate functions and the existing training corpus,
1191 unlike the cost functions typically applied in inductive synthesis, which must be estimated from
1192 input/output examples.

1193 8 CONCLUSION AND FUTURE WORK

1194 We have presented *corpus-guided top-down synthesis* (CTS)—an efficient new algorithm for syn-
1195 thesizing libraries of functional abstractions capturing common functionality within a corpus of
1196 programs. CTS directly synthesizes the abstractions, rather than exposing them through a series of
1197 rewrites as is done by deductive systems. Key to its performance is the usage of a guiding utility
1198 function, which allows CTS to effectively search (and prune out large portions of) the space of
1199 possible abstractions.

1200 We implement this algorithm in STITCH, an open-source library learning tool which exposes both
1201 Rust and Python bindings. We evaluate STITCH across five experimental settings, demonstrating
1202 that it learns comparably compressive libraries with 2 orders of magnitude less memory and 3-4
1203 orders of magnitude less time compared to the state-of-the-art deductive algorithm of Ellis et al.
1204 2021. We also find that STITCH scales to learning libraries of abstractions from much larger datasets
1205 of deeper program trees than is possible with prior work, and that the *anytime* property of top-down
1206 corpus-guided search — abstractions discovered via top-down search are already compressive early
1207 in search and improve as it continues — offers the opportunity for high-quality library learning
1208 even in complex domains through early search termination.

1209 For future work, we aim to build upon the hybrid approach introduced in Section 5 by more
1210 closely integrating STITCH into a deductive rewrite system like egg [Willsey et al. 2021], further
1211 extending the sorts of equivalences it can reason over. Beyond this, we also aim to explore applying
1212 CTS to corpora of human-written code, such as Haskell and OCaml codebases. We are also exploring
1213 the application of CTS to finding abstractions in dataflow DAGs and more general graph structures
1214 like molecules, since the basic ideas of matching and bounding are still applicable. Finally, we are
1215 interested in the new ways that abstraction learning can be applied when it is made this fast; for
1216 example, by directly integrating abstraction learning into a synthesis search procedure.

1217 REFERENCES

- 1218 Miltiadis Allamanis, Earl T Barr, Premkumar Devanbu, and Charles Sutton. 2018. A survey of machine learning for big code
1219 and naturalness. *ACM Computing Surveys (CSUR)* 51, 4 (2018), 1–37.
1220 Miltiadis Allamanis and Charles Sutton. 2014. Mining idioms from source code. In *Proceedings of the 22nd ACM SIGSOFT
1221 international symposium on foundations of software engineering*. 472–483.

- 1226 Matej Balog, Alexander L Gaunt, Marc Brockschmidt, Sebastian Nowozin, and Daniel Tarlow. 2016. Deepcoder: Learning to
 1227 write programs. *arXiv preprint arXiv:1611.01989* (2016).
- 1228 Rod M Burstall and John Darlington. 1977. A transformation system for developing recursive programs. *Journal of the ACM (JACM)* 24, 1 (1977), 44–67.
- 1229 Qiaochu Chen, Xinyu Wang, Xi Ye, Greg Durrett, and Isil Dillig. 2020. Multi-modal synthesis of regular expressions. In
 1230 *Proceedings of the 41st ACM SIGPLAN International Conference on Programming Language Design and Implementation, PLDI 2020, London, UK, June 15-20, 2020*, Alastair F. Donaldson and Emina Torlak (Eds.). ACM, 487–502. <https://doi.org/10.1145/3385412.3385988>
- 1231 Xinyun Chen, Chang Liu, and Dawn Song. 2018. Execution-guided neural program synthesis. In *International Conference on Learning Representations*.
- 1232 Adam Chlipala, Benjamin Delaware, Samuel Duchovni, Jason Gross, Clément Pit-Claudel, Sorawit Suriyakarn, Peng Wang,
 1233 and Katherine Ye. 2017. The end of history? Using a proof assistant to replace language design with library design. In
 1234 *2nd Summit on Advances in Programming Languages (SNAPL 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- 1235 Andrew Cropper. 2019. Playgol: Learning programs through play. *arXiv preprint arXiv:1904.08993* (2019).
- 1236 Nicolaas Govert de Bruijn. 1972. Lambda calculus notation with nameless dummies, a tool for automatic formula manip-
 1237 ulation, with application to the Church-Rosser theorem. In *Indagationes Mathematicae (Proceedings)*, Vol. 75. Elsevier,
 1238 381–392.
- 1239 Eyal Dechter, Jonathan Malmaud, Ryan Prescott Adams, and Joshua B Tenenbaum. 2013. Bootstrap learning via modular
 1240 concept discovery. In *Proceedings of the International Joint Conference on Artificial Intelligence*. AAAI Press/International
 1241 Joint Conferences on Artificial Intelligence.
- 1242 Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Luc Cary, Lucas Morales, Luke Hewitt, Armando Solar-
 1243 Lezama, and Joshua B Tenenbaum. 2020. Dreamcoder: Growing generalizable, interpretable knowledge with wake-sleep
 1244 bayesian program learning. *arXiv preprint arXiv:2006.08381* (2020).
- 1245 Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Lucas Morales, Luke Hewitt, Luc Cary, Armando Solar-
 1246 Lezama, and Joshua B Tenenbaum. 2021. Dreamcoder: Bootstrapping inductive program synthesis with wake-sleep
 1247 library learning. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*. 835–850.
- 1248 Kevin M Ellis, Lucas E Morales, Mathias Sablé-Meyer, Armando Solar Lezama, and Joshua B Tenenbaum. 2018. Library
 1249 learning for neurally-guided bayesian program induction. (2018).
- 1250 Matthias Felleisen and Robert Hieb. 1992. The Revised Report on the Syntactic Theories of Sequential Control and State.
 1251 *Theor. Comput. Sci.* 103, 2 (sep 1992), 235–271. [https://doi.org/10.1016/0304-3975\(92\)90014-7](https://doi.org/10.1016/0304-3975(92)90014-7)
- 1252 John K Feser, Swarat Chaudhuri, and Isil Dillig. 2015. Synthesizing data structure transformations from input-output
 1253 examples. *ACM SIGPLAN Notices* 50, 6 (2015), 229–239.
- 1254 Yaroslav Ganin, Tejas Kulkarni, Igor Babuschkin, SM Ali Eslami, and Oriol Vinyals. 2018. Synthesizing programs for images
 1255 using reinforced adversarial learning. In *International Conference on Machine Learning*. PMLR, 1666–1675.
- 1256 Sumit Gulwani, José Hernández-Orallo, Emanuel Kitzelmann, Stephen H Muggleton, Ute Schmid, and Benjamin Zorn. 2015.
 1257 Inductive programming meets the real world. *Commun. ACM* 58, 11 (2015), 90–99.
- 1258 Robert John Henderson. 2013. Cumulative learning in the lambda calculus. (2013).
- 1259 Irvin Hwang, Andreas Stuhlmüller, and Noah D Goodman. 2011. Inducing probabilistic programs by Bayesian program
 1260 merging. *arXiv preprint arXiv:1110.5667* (2011).
- 1261 Srinivasan Iyer, Alvin Cheung, and Luke Zettlemoyer. 2019. Learning programmatic idioms for scalable semantic parsing.
 1262 *arXiv preprint arXiv:1904.09086* (2019).
- 1263 R Kenny Jones, David Charatan, Paul Guerrero, Niloy J Mitra, and Daniel Ritchie. 2021. ShapeMOD: macro operation
 1264 discovery for 3D shape programs. *ACM Transactions on Graphics (TOG)* 40, 4 (2021), 1–16.
- 1265 Manos Koukoutos, Mukund Raghothaman, Etienne Kneuss, and Viktor Kuncak. 2017. On repair with probabilistic attribute
 1266 grammars. *arXiv preprint arXiv:1707.04148* (2017).
- 1267 Tessa Lau, Steven A Wolfman, Pedro Domingos, and Daniel S Weld. 2003. Programming by demonstration using version
 1268 space algebra. *Machine Learning* 53, 1 (2003), 111–156.
- 1269 Miguel Lázaro-Gredilla, Dianhuan Lin, J Swaroop Guntupalli, and Dileep George. 2019. Beyond imitation: Zero-shot task
 1270 transfer on robots by learning concepts as cognitive programs. *Science Robotics* 4, 26 (2019), eaav3150.
- 1271 Mina Lee, Sunbeam So, and Hakjoo Oh. 2016. Synthesizing regular expressions from examples for introductory automata
 1272 assignments. In *Proceedings of the 2016 ACM SIGPLAN International Conference on Generative Programming: Concepts
 1273 and Experiences, GPCE 2016, Amsterdam, The Netherlands, October 31 - November 1, 2016*, Bernd Fischer and Ina Schaefer
 1274 (Eds.). ACM, 70–80. <https://doi.org/10.1145/2993236.2993244>
- 1275 Percy Liang, Michael I Jordan, and Dan Klein. 2010. Learning programs: A hierarchical Bayesian approach. In *Proceedings of
 1276 the 27th International Conference on Machine Learning (ICML-10)*. 639–646.
- 1277

- 1275 Dianhuan Lin, Eyal Dechter, Kevin Ellis, Joshua B Tenenbaum, and Stephen H Muggleton. 2014. Bias reformulation for
1276 one-shot function induction. (2014).
- 1277 Zohar Manna and Richard Waldinger. 1980. A deductive approach to program synthesis. *ACM Transactions on Programming
Languages and Systems (TOPLAS)* 2, 1 (1980), 90–121.
- 1278 Tom M Mitchell. 1977. Version spaces: A candidate elimination approach to rule learning. In *Proceedings of the 5th
international joint conference on Artificial intelligence-Volume 1*. 305–310.
- 1280 Maxwell Nye, Yewen Pu, Matthew Bowers, Jacob Andreas, Joshua B Tenenbaum, and Armando Solar-Lezama. 2021.
1281 Representing Partial Programs with Blended Abstract Semantics. In *International Conference on Learning Representations*.
- 1282 Nadia Polikarpova, Ivan Kuraj, and Armando Solar-Lezama. 2016. Program synthesis from polymorphic refinement types.
1283 *ACM SIGPLAN Notices* 51, 6 (2016), 522–538.
- 1284 Oleksandr Polozov and Sumit Gulwani. 2015. Flashmeta: A framework for inductive program synthesis. In *Proceedings of
the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*.
1285 107–126.
- 1286 Ameesh Shah, Eric Zhan, Jennifer Sun, Abhinav Verma, Yisong Yue, and Swarat Chaudhuri. 2020. Learning Differentiable
1287 Programs with Admissible Neural Heuristics. In *Advances in Neural Information Processing Systems*, H. Larochelle,
1288 M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 4940–4952. <https://proceedings.neurips.cc/paper/2020/file/342285bb2a8cadef22f667eeb6a63732-Paper.pdf>
- 1289 Claude Elwood Shannon. 1948. A mathematical theory of communication. *The Bell system technical journal* 27, 3 (1948),
1290 379–423.
- 1291 Eui Chul Shin, Miltiadis Allamanis, Marc Brockschmidt, and Alex Polozov. 2019. Program synthesis and semantic parsing
1292 with learned code idioms. *Advances in Neural Information Processing Systems* 32 (2019).
- 1293 Max Willsey, Chandrakana Nandi, Yisu Remy Wang, Oliver Flatt, Zachary Tatlock, and Pavel Panchekha. 2021. Egg:
1294 Fast and Extensible Equality Saturation. *Proc. ACM Program. Lang.* 5, POPL, Article 23 (jan 2021), 29 pages. <https://doi.org/10.1145/3434304>
- 1295 Catherine Wong, Kevin M Ellis, Joshua Tenenbaum, and Jacob Andreas. 2021. Leveraging language to learn program
1296 abstractions and search heuristics. In *International Conference on Machine Learning*. PMLR, 11193–11204.
- 1297 Catherine Wong, William McCarthy, Gabriel Grand, Jacob Andreas, Joshua B Tenenbaum, Robert Hawkins, and Judy Fan.
1298 2022. Identifying concept libraries from language about object structure. In *CogSci*. To appear.
- 1299
- 1300
- 1301
- 1302
- 1303
- 1304
- 1305
- 1306
- 1307
- 1308
- 1309
- 1310
- 1311
- 1312
- 1313
- 1314
- 1315
- 1316
- 1317
- 1318
- 1319
- 1320
- 1321
- 1322
- 1323