

Towards Efficient Neurally-Guided Program Induction for ARC-AGI

Simon Ouellette

Abstract

ARC-AGI is an open-world problem domain in which the ability to generalize out-of-distribution is a crucial quality. Under the program induction paradigm, we present a series of experiments that reveal the efficiency and generalization characteristics of various neurally-guided program induction approaches. The three paradigms we consider are *Learning the grid space*, *Learning the program space*, and *Learning the transformation space*. We implement and experiment thoroughly on the first two, and retain the second one for ARC-AGI submission. After identifying the strengths and weaknesses of both of these approaches, we suggest the third as a potential solution, and run preliminary experiments.

Introduction

Deep Learning is notoriously powerful when addressing a wide variety of problems. Yet, in other areas, it is deficient. One predominant difference between the problem domains where it is highly successful, and the domains where it is less so, is the openness of that domain. In closed-world domains, where the model is able to densely sample and interpolate over the full set of possibilities, deep learning often achieves near-human or sometimes even super-human ability. However, in open-world domains, where the full range of possibilities is so vast that it is practically infeasible to densely sample its entirety, interpolation-based techniques like Deep Learning can be sub-optimal.

This results in Deep Learning’s well-known challenges in generalizing outside the training set distribution. The Abstraction & Reasoning Corpus (Chollet et al. 2024) (ARC-AGI) is one of the foremost problem domains that specifically challenges this closed-world assumption. This is done by using a hidden test set that contains tasks that are qualitatively distinct from any of the publicly available tasks. This encourages research on extending the out-of-distribution generalization capabilities of learning systems. This will be the focus of the work presented here.

The current implementations of the proposed approaches are not yet mature enough to yield an interesting or competitive performance on ARC-AGI. This paper reports proof-of-concept results of experiments made under controlled and limited conditions, exploring ways in which we can efficiently extend the generalization capabilities of a Deep Learning model via search, in a program induction context.

We discuss three paradigms: *Learning the grid space* (LGS), *Learning the program space* (LPS), and *Learning the transformation space* (LTS). The first two have been implemented and thoroughly experimented upon. The second one (LPS) is the best performing method on the ARC evaluation set, so it is retained for ARC-AGI submission. It is the approach that is therefore described in most detail. The third one (LTS) has not been implemented yet, due to lack of time, but it is described in a brief, general way. This third approach is proposed as a way to combine the strengths and weaknesses of the other two approaches. While preliminary experiments are made to support this hypothesis, full implementation and experimentation is left to future work.

Our main contributions We present a few novel results:

1. To the best of our knowledge, we are the first to analyze the out-of-distribution generalization characteristics of program induction approaches based on enumerating the program probability space.
2. We present the first results on ARC-AGI for the LGS approach, and analyze its strengths and weaknesses.
3. We propose a novel probabilistic program enumeration-based search algorithm for program induction, leaning heavily on Transformer-based auto-regressive token sequences, rather than the typical n-gram approach, and analyze its strengths and weaknesses.
4. We sketch the outline of a novel approach (LTS) aimed at addressing the flaws in the previous approaches. We provide preliminary experimentation to support the hypotheses being advanced.

The Problem

The problem setting we use for this paper follows a typical program synthesis setting. We assume that the Domain Specific Language (DSL) includes all the necessary primitives to solve each test task. The goal is to search (as efficiently as possible) for the program that solves each task in the test set, given the provided DSL. Throughout this paper, the tasks experimented upon are mainly taken from the ARC-AGI evaluation set.

A compute and memory time budget is allocated for each task solution, so it must find it within those budget constraints. This can be easy, even trivial for very small programs, but due to combinatorial explosion can become quite

unmanageable for large programs. More formally, given a DSL $\Omega = \{\pi_1, \pi_2, \dots, \pi_N\}$ containing N primitive functions denoted π_i , a search algorithm $F(X, Y)$ given the support input examples X_s and support output examples Y_s must return within the allocated CPU and memory budgets a program P such that $P(X_q) = Y_q$, where X_q and Y_q are the query examples. In other words, the support examples are used to infer the underlying program, and the query examples are used to confirm the correctness of the program.

There is an additional, crucial aspect to our problem setting. While we assume the DSL Ω to fully contain the required primitives to solve both the training and test datasets, we allow the test dataset’s solution programs to follow different structures than the test datasets. In other words, we place the additional constraint to our problem domain that the test set must be out-of-distribution with respect to the training set in terms of the compositional structures of P . This is important because, if it were not the case, a neural network would be sufficient to solve this problem domain.

Learning the Grid Space

The idea of this approach is to learn a model of the space of possible ARC-AGI grids, under a specific DSL. It is important to note that there is no such thing as a pure grid similarity space unconditioned on some DSL, since under two different DSLs the same grid pair can have quite different distances. From there, it is possible to input two different grids and estimate their similarity. This can be used in a search algorithm to gradually build a program in an execution-guided manner. The idea of execution-guided program induction is to leverage some heuristic to guide the choice of primitives during construction. This heuristic requires feedback from partial execution of the program being constructed. As such, program solution construction and execution happen in parallel and impact each other.

At each iteration of the search, the algorithm attempts all possible primitives on the current intermediate grid in the program execution, and selects the one that brings it the closest to the target grid state. Note that this approach assumes that all (or most) steps in program execution result in a grid-to-grid transformation. This is why the experiments on LGS use a DSL that strictly contains grid-to-grid transformations (DSL Version 1). As such, it should be kept in mind that the results reported for this method are therefore the best case scenario where all steps in program execution result in a grid transformation.

The experiments reported here used a Transformer encoder-only model with max pooling that outputs a flattened vector: a grid embedding. The training procedure consists of iteratively feeding a pair of grids into the model, retrieving their respective embeddings, and then penalizing deviations between their dot product and their ground truth distance.

The training data is generated by randomly sampling patches of ARC-AGI training set grids, and then picking a random number of grid transformation primitives to apply to them (zero to eight, in our experiments). This randomly selected sequence of transformations is executed, and the output grid is retained. The input and output grids, along

with the number of transformations applied, constitute a full training sample. The number of transformations is converted to a similarity between zero and one.

It should be noted that the procedurally generated training data often contained over-estimated transformation distance ground truths. This is because it is very easy to accidentally generate redundant transformations. The desirable quality of a good grid distance ground truth is that it represents the shortest distance between two grids, rather than whatever distance was arbitrarily used to generate these two grids. A trivial example of this is that "rotate 90, rotate 90, rotate 90" should yield a grid distance of one (rotate 270), not three steps. Several heuristics were implemented in the training data generation script to attempt to mitigate this issue, but nonetheless the issue remains. This is arguably one of the weaknesses of applying supervised learning to a cost-to-go type of approach.

Learning the Program Space

This is the selected approach, because it performed better than LGS on the ARC-AGI evaluation set. As such, it will be described in more depth than the other two methods. Furthermore, because we are referring to a specific, novel implementation of the overarching paradigm of LPS, we name it *GridCoder*.

At a high-level, the solution consists of training a transformer to output a program, using a pre-determined grammar (DSL) and syntax, that solves the task. Specifically, the transformer outputs a sequence of probabilities over tokens that can be interpreted as a probability distribution over program space. We then use a search algorithm to enumerate that space and test valid programs for correctness. In a sense, the search explores the area of the solution space covering the neural network’s region of uncertainty.

The Search Algorithm

It is a probability-based enumeration of programs, conceptually similar to DreamCoder (Ellis et al. 2021) and Heap Search (Fijalkow et al. 2021). However, unlike these approaches, the probability predictions are not based on n-grams, but instead are based on transformer generated token probability sequences. That is, instead of predicting the probability over children classes conditional on (n-1) parents, we predict the probability over children classes conditional on the entire maximum probability token sequence (as in typical auto-regressive transformer decoding).

Also, unlike the other two paradigms presented in this paper (LGS and LTS), this approach is not execution-guided, in that it does not receive feedback about the intermediate state of the program as it develops a solution. Initially, there was hope that the transformer decoder could learn to do this implicitly, however the empirical results point to this not being the case (at least, not with the current architecture).

A first full decoding loop is executed on the Vision-Language Model (VLM) that is fed one input and output grid pair from a given task’s demonstration set. The full decoding loop stops until either the maximum sequence length (40, in our experiments) is reached, or a token is found where the

only class that has a probability greater than some threshold τ is the End Of Sentence class. See algorithm 2 for the pseudo-code of a full decoding loop.

From there, there is a short *bootstrapping* phase in which a few more inference runs are made. In these, the example pair index and the starting token (or, sometimes, the first two starting tokens) are randomly selected. The full probability distribution is calculated as a mean of the probabilities over tokens for all of these initial decoding loops. Hence, the probability distribution *prob.dist* (see algorithm 1) is a two-dimensional array, where one dimension represents the token positions in the sequence, and the other dimension represents the probability distribution over the possible classes (i.e. primitives in the DSL) at that token position. Therefore, our probability distribution is represented as a flat sequence, rather than a tree in which the probability of children depends on which parents are selected. This is a simplifying assumption of conditional independence, similar to the one made in Bayesian networks, that increases the efficiency of the search in our experiments.

Algorithm 1: GridCoder

Input: X , input for each example
Input: Y , target for each example
Input: DSL , a DSL to search over
Input: T , the timeout parameter
Input: M , the neural network model
Input: K , number of bootstrapping examples
Input: τ , probability threshold
Output: correct program if found

- 1: **// Bootstrapping probabilities**
- 2: $prob_dist \leftarrow \text{getProbSpace}(M, X, Y)$
- 3: $counts = \text{array of ones of length of } prob_dist$
- 4: **for** $k \in \{1..K\}$ **do**
- 5: $idx \leftarrow \text{Pick a random example index}$
- 6: $seq \leftarrow \text{Pick a random first or first two tokens}$
- 7: $tmp_dist \leftarrow \text{getProbSpace}(M, X, Y, idx, seq)$
- 8: $prob_dist[len(seq) :] += tmp_dist$
- 9: $counts[len(seq) :] += 1$
- 10: **end for**
- 11: $prob_dist = prob_dist / counts$
- 12:
- 13: **// Enumerate programs from tokens with $P > \tau$**
- 14: $sorted_progs \leftarrow \text{enumerate and sort programs}$
- 15:
- 16: **// Evaluate the programs and return the correct one**
- 17: **for** $\pi \in sorted_progs$ **do**
- 18: **if** time elapsed $> T$ **then**
- 19: return null
- 20: **end if**
- 21: $success \leftarrow \text{evaluate_program}(\pi, X, Y, DSL)$
- 22: **if** $success$ **then**
- 23: return π
- 24: **end if**
- 25: **end for**
- 26: return null

Algorithm 2: getProbSpace

Input: X , input for each example
Input: Y , target for each example
Input: M , the neural network model
Input: τ , probability threshold
Input: idx , the example pair index
Input: seq , the starting seq, empty by default
Input: L , maximum sequence length
Output: Probability distribution over each token in the sequence.

- 1: $seq_len \leftarrow 0$
- 2: $done \leftarrow \text{False}$
- 3: $shiftedSeq \leftarrow [StartOfSentence] + seq$
- 4: $prob_dist \leftarrow []$
- 5: **while** not done and $seq_len < L$ **do**
- 6: $probs \leftarrow \text{predict}(M, X[idx], Y[idx], shiftedSeq)$
- 7: $prob_dist \leftarrow prob_dist + [probs]$
- 8: $best_token \leftarrow \text{argmax}(probs)$
- 9: **if** $best_token$ is End Of Sentence **then**
- 10: **if** End Of Sentence is the only class with $P > \tau$ **then**
- 11: $done \leftarrow \text{True}$
- 12: **else**
- 13: $best_token \leftarrow \text{second most probable class}$
- 14: **end if**
- 15: **end if**
- 16: $shiftedSeq \leftarrow shiftedSeq + [best_token]$
- 17: $seq_len \leftarrow seq_len + 1$
- 18: **end while**

In the main algorithm (alg. 1), once we have estimated our

probability space over the DSL, we then proceed to enumerating all the possible programs formed by the token classes whose probability is greater than τ , the probability threshold hyperparameter. If we calculate that the total set of possible programs is greater than a pre-determined number, we gradually increase τ until the cut-off level results in a manageable number of programs. This is not shown in the pseudo-code, for the sake of brevity. Then, we proceed to sort the programs in decreasing order of their joint probability: we multiply the probabilities associated with all of the tokens in its sequence, truncating at the first `<End Of Sentence>` token seen.

Conditional Independence

One aspect of the algorithm that is important, yet is not obvious in the pseudo-code, is the underlying data structure. While typical program induction operates on a tree, Grid-Coder operates on a flat sequence of nodes. In other words, our probabilities are represented by a list of lists: the probability distribution over a given token is solely dependent on its position in the sequence, rather than being conditional on the previous choices of tokens. That is, to be more exact, it was conditional on the previous (max probability) tokens at the time of auto-regressively calculating the probability distribution. But, at search time, we do not update the choice of selected tokens and requery the model to find the new, conditionally updated probabilities.

Doing this, and maintaining a tree of probabilities instead, would be theoretically correct, because it would respect the natural conditional dependence between primitives. The probability over the next token is theoretically dependent not just on the token's position in the sequence, but on all of the previous selections that were made. However, we choose to make a simplifying assumption of conditional independence between nodes, much like in Bayesian networks. While this simplifying assumption potentially hurts in some cases, overall we believe it yields a significant gain in efficiency: making this assumption means that we only need to call the neural network a few times while bootstrapping the probabilities.

In fact, in the baseline version of this approach, the neural network need only be called exactly once. From there, all of the necessary probability distributions can be acquired. However, it was found empirically that using a bit of bootstrapping helps. Bootstrapping in this context means that we collect the probability distributions as a mean of K decoding iterations (where K is six in our experiments). Each decoding iteration differs in which example pair is provided to the neural network, and which starting token is selected. This bootstrapping idea helps mitigate the rare cases where conditional independence hurts because the neural network outputs an unusually constrained (i.e., falsely confident) set of probabilities on the first inference. Using a few bootstrapped inference trials seems to open up the space of possibilities a bit, introducing a bit of much needed entropy the process.

The DSLs

A few iterations of the DSL were used during the experiments, in which primitives were gradually added to the pre-

vious version. They are briefly described here. You can refer to the appendix to see a full list of all primitives in each version, along with a brief description of each primitive.

Version 1 The first version of the DSL contains 74 primitives, which exclusively transform one or two input grids into one grid. These grid-to-grid transformations include things such as changing all pixels of one color to another color, applying grid rotations, mirroring, etc. This choice of primitives was initially made to facilitate research on the LGS paradigm.

Version 2 The second version of the DSL contains 89 primitives. It contains all the primitives in version 1, plus some primitives for object detection and manipulation. In particular, six different versions of object detection were implemented, for various concepts of objectness (what defines an object tends to change from task to task). In addition, it contains new primitives that allowing looping through object lists and applying transforms to them, either in-place, or by cropping all the objects and generating a new grid out of them.

Version 3 The third version of the DSL contains 98 primitives. New primitives were added to allow filtering objects based on a few characteristics such as size and symmetry.

VLM Architecture & Training

The neural network used to provide the probabilities that guide the search is based on a Vision Language Model (VLM) paradigm (see figure 1). It is a pair of convolutional neural networks, one receiving the input grid, the other the output grid, eventually bridging their intermediate features into one convolutional network. This serves as an encoder, followed by a transformer decoder that processes the target sequence. In our experiments this architecture converged faster and to a higher validation accuracy on our dataset than a T5 or LongT5 architecture, while allowing us to use a higher parameter count for the same amount of VRAM. Specifically, the model used for the experiments has 452 million parameters, and only one decoder layer with an embedding dimensionality of 512 and 4 attention heads.

Each predicted token sequence forms a program syntax tree in a bottom-up, left-to-right manner. Each token refers to a primitive from the DSL, or to one of the three special tokens: `<End of Sentence>`, `<New Level>`, `<Identity>`. `<Identity>` can be thought of as a primitive that takes as input a grid and returns that grid without any modification. It is often used as a placeholder in order to disambiguate between different possible interpretations of a program. `<New Level>` is a special token that indicates that the tokens that follow define the next level in the program syntax tree. The syntax tree is built bottom-up (as the token sequence is read from left to right), and in each level the function argument ordering is defined from left to right in the same order as the token sequence. Here are a few examples to help better understand the syntax structure.

Suppose the following token sequence: `[lefthalf, righthalf, <New Level>, cellwiseOR, <New Level>, set_fg.color3, <End of Sentence>]`.

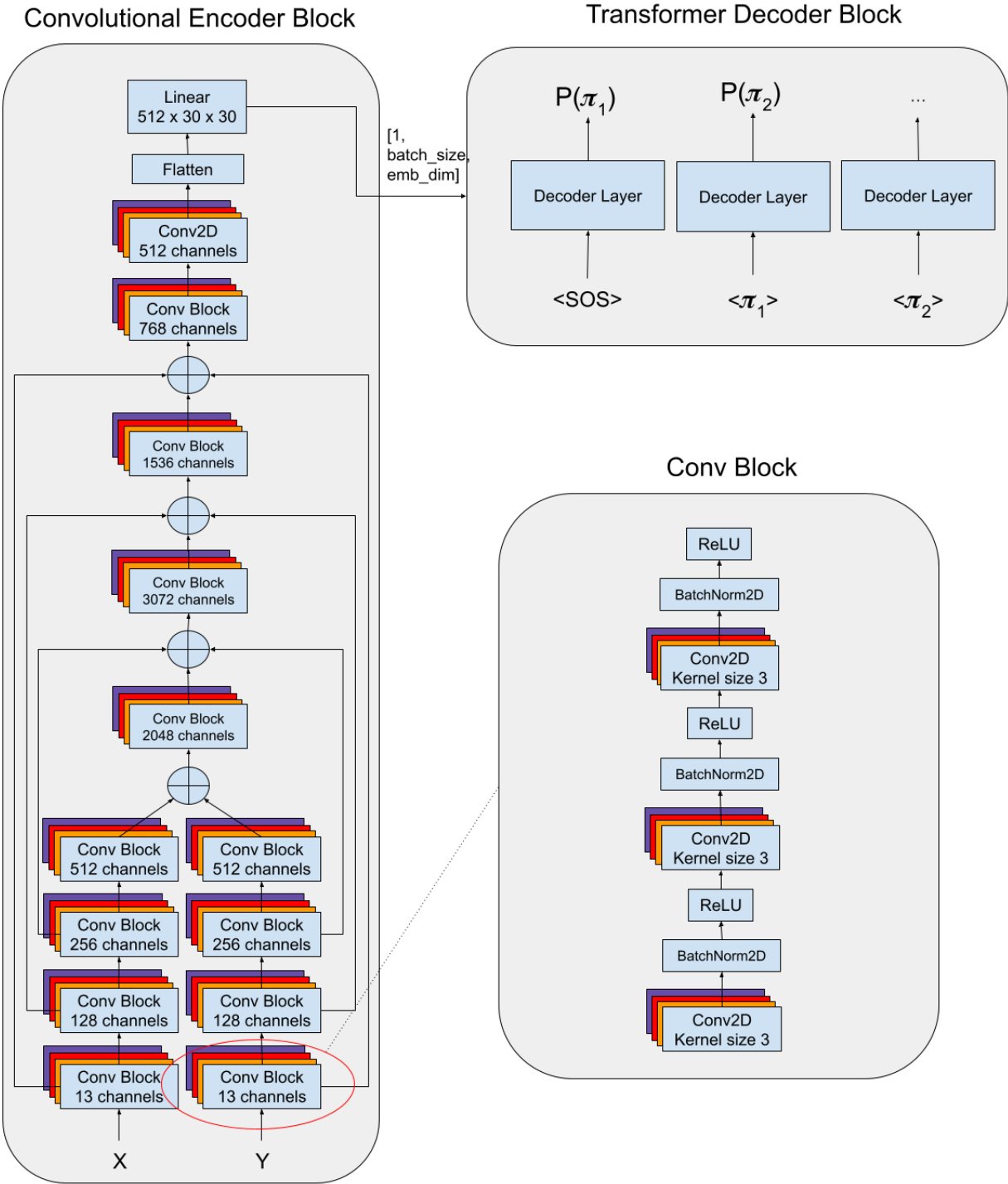


Figure 1: Architecture of VLM

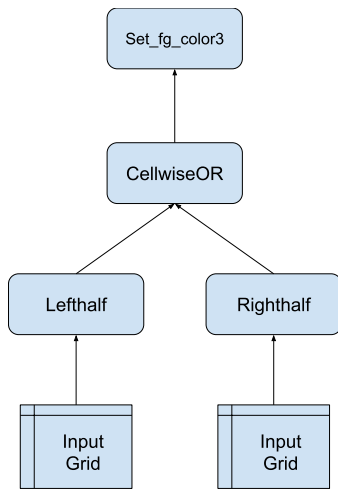


Figure 2: Syntax tree for [lefthalf, righthalf, <New Level>, cellwiseOR, <New Level>, set_fg_color3, <End of Sentence>]

The corresponding syntax tree (see figure 2) indicates that we take the left half and right half of the original grid separately, merge the pixels using OR logic (i.e., if any pixel is a foreground pixel, we place a foreground pixel of the same color on the target grid, prioritizing the color of the first argument if both are foreground pixels), and then change all of the non-zero pixels to color three.

A more advanced example is the following token sequence: [topthird, hcenterthird, bottomthird, <New Level>, <Identity>, cellwiseOR, <New Level>, cellwiseOR, <End of Sentence>]. Here the corresponding syntax tree (figure 3) forms a program that splits the original grid into three sub-grids vertically, and then merges them with OR logic. Specifically, in the second level it passes on the output of `topthird` intact to the upper level, which means that it first merges the output of `hcenterthird` and `bottomthird`. Finally, at the last level, it merges `topthird` with the output from the previous level that merged `hcenterthird` and `bottomthird`.

The training was first done on the DSL version 1 data, until convergence. Then, the DSL version 2 and DSL version 3 models were separately fine-tuned on their respective datasets, starting from the pre-trained DSL version 1 model. The DSL 2 and DSL 3 datasets contain training examples from the previous versions as well (not just exclusively tasks related to the newly added primitives). The learning rate selected was 0.0001, and the weight decay is 0.005. For DSL version 1, a total of around 5M training samples (i.e., input-output grid pairs + ground-truth program sequence for distinct tasks) were generated. For DSL version 2, around 600K training samples were generated, and around 300K for DSL version 3 (because data generation is slower). More data would have been preferable, but time was a significant constraint.

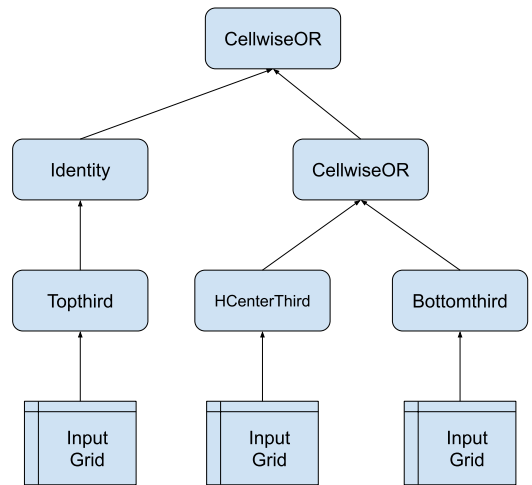


Figure 3: Syntax tree for [topthird, hcenterthird, bottomthird, <New Level>, <Identity>, cellwiseOR, <New Level>, cellwiseOR, <End of Sentence>]

Training Data Generation

The training data for the VLM is the output of procedurally generated tasks. The concept is to identify task "meta-patterns" or categories in the ARC-AGI Training Set, and to implement data generators that reproduce them, while randomly varying several aspects of the tasks. Different DSL versions have different task generators associated with them.

DSL version 1 In this version, there are only three task generators implemented. The input grids used to generate the output grids are randomly sampled from two sources: manually pre-generated grids (mostly drawn from the ARC-AGI Training set), and randomly generated grids. Input grids that match grids found in the ARC-AGI Evaluation set are not allowed, though it is not impossible for the randomly generated grids, by extreme luck, to reproduce grids similar to those in the ARC-AGI Evaluation set.

The first task generator consists of randomly generating trivial tasks that are either made up of one primitive (that takes as input a grid and generates a grid as output), or of two primitives composed together. The two primitive tasks are not entirely random: a script that generates all two-primitive permutations of the DSL was created, but several heuristics were used to prevent non-sensical or redundant compositions. Overall, this task generator has 796 distinct tasks.

The second task generator is based on a "split-merge" pattern that was observed. This task category produces tasks in which the goal is to split the input grid in some way, and then merge them back together into one grid based on *OR*, *XOR*, *AND*, *NOR* or *Difference* logic (see the appendix for details). Finally, one last global grid-to-grid transform can be randomly applied on the resulting grid. The possible split patterns are two-way to four-way horizontally or vertically, or four-way by taking each quadrant separately. These various parameters (split method, logical operator, optional last transform) are randomized and combined to generate a

whole range of possible tasks that follow the same "meta-pattern".

The third task generator is based on a "tiling" pattern. In this task category, the idea is to apply various rotational and mirroring transforms to the input grid, and concatenate its transformed variants together in some way to produce a larger output grid. The possible individual transforms are randomly picked from rotations, the identity function, or flipping the grid horizontally or vertically. The tiling patterns that can be generated are 2x2, 3x3, horizontal concatenation of 2 to 4 grids, and vertical concatenation of 2 to 4 grids.

DSL version 2 This DSL version includes a new general-purpose task generator whose pattern is to randomly pick a number of objects to generate, to randomly pick some transforms to apply to them (including no transform at all as a possibility), and to decide whether to apply those transforms to the objects in-place, or to crop the objects and tile them together into some pattern (which itself is randomly determined among a few possible templates). Finally, there is an optional post-processing transform, applied to the grid as a whole, that can be selected.

DSL version 3 This DSL version includes three new task generators that are a bit more niche than the one added in DSL version 2. The first one, the object selector task generator, consists of generating random objects and deciding randomly whether the task is to keep one, or to filter out some objects based on one of the following characteristics: horizontal symmetry or lack thereof, vertical symmetry or lack thereof, object size (number of foreground pixels it contains), and the number of sub-objects it contains. As for most other task generators, a random post-processing transform can be applied to the output grid.

The windowing task generator produces frame-like objects, i.e. non-filled rectangles of uniform color that can contain randomly generated pixels or objects. The objective of the task is to apply randomly selected transforms to the inside of the frames (with or without the rectangular boundary) and then either crop the object, if there is only one, or apply those transforms in place.

Finally, the object recombiner task generator creates random objects on a grid, decides whether to split the objects horizontally or vertically, and applies some color-based modification to one of the halves. It then recombines the halves into whole objects.

Learning the Transform Space

This approach can be thought of as an evolution of LPS, while bringing back the execution-guided aspect of LGS. The concept is to train a model such that, given an intermediate or starting program state, and a target grid, it predicts the probability distribution over the DSL for the next token. In other words, the main difference with LPS is that we explicitly feed back into each decoding step some notion of the intermediate state of the program. Thus, the algorithm can leverage the efficiency of auto-regressive decoding with the ability to evaluate the transformations in and of them-

selves regardless of whether this program pattern has been seen during training.

The emphasis is on learning what transformations are required to bring the program execution from its current state to the target grid. Initially, this was the intended result of the GridCoder algorithm, however further analysis reveals that the latter does not maintain such a hidden state implicitly, it only superficially relies on the encoding of the input-output grid and the token sequence generated so far. See the Discussion section for a deeper analysis.

LTS borrows the same program syntax and auto-regressive sequence supervision as in GridCoder. In a more basic implementation, the intermediate program state could simply be the last full intermediate grid produced by the program execution. The experiments associated with Table 4 make this basic assumption. In a more sophisticated version, however, there needs to be a hidden latent state maintained throughout the decoding process that gets updated at each step.

The challenging aspect of maintaining this generalized intermediate state is that the output of a primitive can be anything from a Grid to an integer, a Boolean, a list of integers, a list of Grids, etc. A clever mechanism must be designed to allow embedding into a fixed vector space this dynamic and diverse range of outputs, which is left to future work.

Experiments & Results

This section presents the results of three experiments. First, a comparison of the performance of different approaches on the ARC evaluation set is made. Second, the increments of performance of the selected approach (GridCoder) are shown, as new primitives are added to the DSL. Finally, an experiment on structurally out-of-distribution tasks is shown, to illustrate the limitations of the selected approach and the promising capabilities of the suggested new approach (*Learning the transformation space*).

Performance Comparison on ARC-AGI Eval Set

Table 1 reports the success rate (as a percentage of solved tasks) of various approaches on the ARC-AGI evaluation set. We show the numbers in terms of the absolute ARC-AGI evaluation set performance, but also in terms of the subset of the ARC-AGI evaluation set that is theoretically solvable according to the DSL. DSL Version 1 is used for these experiments (see section titled "The DSLs"). This subset refers to the 29 tasks (7.25% of full evaluation set) that can theoretically be solvable in the current DSL. All the other tasks are impossible to solve regardless of how efficient the search is.

For each task, a time budget of up to 15 minutes to find the answer is allowed. For *GridCoder* and *GridCoder cond.*, the 5-minute time budget performance is also shown, in order to better differentiate the impact of the conditional independence assumption. The following algorithms are compared in this experiment:

GridCoder This is the selected approach for the associated ARC-AGI submission, elaborated upon in the *Learning the program space* section.

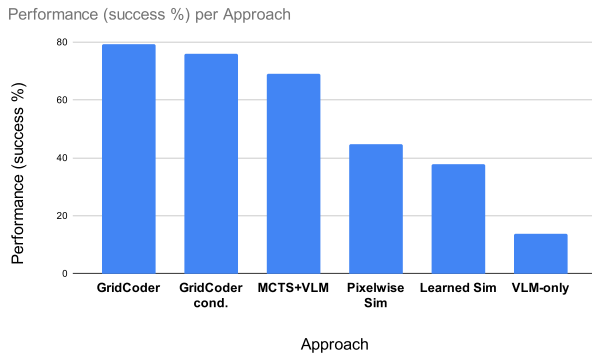


Figure 4: Success rate on the 29-task solvable set from ARC-AGI evaluation set

GridCoder cond. This is a variant of GridCoder where the conditional independence assumption is not made. The objective is to show that the conditional independence assumption is an improvement.

VLM-only This is the pre-trained neural network used in GridCoder, but without the search component. The objective is to show the necessity of search.

MCTS+VLM This uses *Monte-Carlo Tree Search*, instead of GridCoder, for the search algorithm. Specifically, it uses a MuZero-like approach (Schrittwieser et al. 2020) where the probabilities instruct the primitive selection at each node, and a pixelwise similarity metric (same as in *Pixelwise sim*) is used to backpropagate the value.

Pixelwise sim This is an LGS approach. Here, *pixelwise* refers to a simple hand-crafted pixelwise comparison of the grid, and the percentage of identical pixels is the similarity heuristic. There is no learning in this version.

Learned sim Similar to *Pixelwise sim*, but the heuristic is a pre-trained Transformer encoder that outputs a similarity heuristic between 0 and 1, as discussed in the *Learning the Grid Space* section.

It should be noted that the two LGS implementations tend to fail on tasks whose ground truth program has more than 5 or 6 primitives, explaining its relatively low performance. It captures the low-hanging fruit of relatively small programs quite well, but it is not efficient enough to solve the programs that have longer descriptions.

GridCoder on Different DSLs

These results make full use of the three presented versions of DSLs. The objective is to study the scaling characteristics of the LPS approach and determine whether scaling this proof-of-concept to a competitive solution is plausible. This is done by training three separate VLM instances, one on each of the three DSLs and their associated training data generators. This can be thought of as a curriculum progression, where the DSL version 3 model includes all primitives and task generators of DSL version 2, which itself includes all primitives and task generators of DSL version 1.

The performance (success rate) on the ARC-AGI Evaluation set can be found in Table 2. For each, the performance on the DSL version 1 solvable tasks (a subset of 29 tasks from the ARC Evaluation Set), with a 5-minute time budget, is used to show whether there is regression on previous tasks as we grow the DSL. For the sake of transparency, the hidden test set performance is also displayed.

The inference times per task are presented in Table 3, to determine whether the solution time scaling with respect to DSL size is super-linear, linear, or sub-linear. The latter is a more attractive quality, since it suggests that the DSL can be grown to solve more tasks, while the former would suggest a solution that scales unfavorably. The average inference time per task, per DSL, and its standard deviations are reported. Additionally, the normalized time per DSL size, i.e. the number of seconds per task on average divided by the number of primitives in the DSL, is shown in order to directly observe whether we have such a super-linear or sub-linear progression.

Generalization of LTS

The purpose of this experiment is to suggest a promising new method that mitigates the limitations of the selected GridCoder approach (see Discussion section), by indicating preliminary empirical results on its out-of-distribution generalization characteristics. 10 new tasks are hand-crafted, specifically selected to guarantee that there is no structurally similar task ever generated in the training data. One of them is picked from the ARC-AGI evaluation set, task #48131b3c, since it has been observed that GridCoder fails to generalize to this task.

A high-level description of the 10 tasks follows:

1. Task #48131b3c from ARC-AGI evaluation set: tile the original grid 2x2, and invert the colors. The training data generator never presents the post-tiling inversion of colors.
2. Hand-crafted task 1: gravitate the pixels to the left, then gravitate the pixels upward, then change the foreground pixels' color to aquamarine. The training data generator never generates tasks that involves the application of three primitives in a row (or longer). The maximum is two.
3. Hand-crafted task 2: rotate the grid 90 degrees, upscale it horizontally by two, and then upscale it vertically by two.
4. Hand-crafted task 3: same as task 2, but add an extra operation of horizontal mirroring at the end.
5. Hand-crafted task 4: same as task 3, but add an extra operation of color inversion at the end.
6. Hand-crafted task 5: tile the original grid 2x4 while alternating 180-degree rotation with the identity transform on each tile. This is out-of-distribution because a 2x4 tiling is never generated in the training data. The closest are either 2x2 or 1x4.
7. Hand-crafted task 6: tile the smallest object in the grid twice horizontally (i.e. concatenate it with itself horizontally), and use that as output grid. In the training data,

Algorithm	ARC-AGI Eval Set % (15m)	Solvable subset % (15m)	Solvable subset % (5m)
GridCoder	5.75	79.3	79.3
GridCoder cond.	5.5	75.9	72.4
MCTS+LVM	5	69	-
Pixelwise Sim	3.25	44.8	-
Learned Sim	2.75	37.9	-
VLM-only	1	13.8	-

Table 1: Success rate on full ARC-AGI evaluation set at 15-minute budget, on subset of 29 solvable tasks at 15-minute budget, and on subset of 29 solvable tasks at 5-minute budget.

DSL Version	ARC-AGI Eval Set %	DSL 1 subset %	Hidden test set %
Version 1	5.75	79.3	0
Version 2	7.75	79.3	1
Version 3	8.25	79.3	1

Table 2: Success rates of different DSL versions on full ARC-AGI evaluation set at 5-minute budget, and on the 29-task DSL version 1 subset at 5-minute time budget. ARC-AGI hidden test set performance is also presented.

DSL Version	Mean per task	Mean per primitive
Version 1	4.74 (1.07)	0.064
Version 2	5.82 (1.22)	0.064
Version 3	5.66 (1.22)	0.057

Table 3: Mean number of seconds taken to solve a task, by DSL. Standard deviations in parentheses. Mean number of seconds taken to solve a task, normalized by number of primitives in a DSL.

tiling tasks and object selection tasks are never combined.

8. Hand-crafted task 7: crop the the object that contains the largest number of sub-objects, rotate it 90 degrees, and then duplicate the top row and the bottom row. The training data only at most applies 1 post-processing transform on these types of object cropping tasks.
9. Hand-crafted task 8: filter out the largest object in the grid and then rotate the grid 270 degrees. The training data generation never applies rotation primitives to the output of an object filtering task.
10. Hand-crafted task 9: crop the largest object in the grid, split it in half horizontally, and merge the left and right halves with cellwise OR logic. In the training data, object cropping and "split and merge" types of tasks are never combined.

The experiment consists of simulating a model that was trained to decode in an execution-guided way (see the section *Learning the Transformation Space*), receiving execution feedback at each step along the way. While the potential impact on inference time is not simulated, the task attempt is considered a success if it is possible to reach the solution by re-launching the current search algorithm on the intermediate output of one of the program sequences that get evaluated in the previous run. As a result, the program is expected to

Task	GridCoder	Execution-guided
Task #48131b3c	NO	YES
Task 1	YES	YES
Task 2	NO	YES
Task 3	NO	YES
Task 4	NO	YES
Task 5	NO	NO
Task 6	NO	NO
Task 7	NO	YES
Task 8	NO	YES
Task 9	NO	YES
Total (%)	10	80

Table 4: Success on out-of-distribution tasks, comparing the selected GridCoder with the suggested execution-guided search.

provide a probability over the program space from a starting point that is the intermediate output of a previous program – hence it is given the opportunity to compose sub-programs.

This proxy experiment was done, instead of correctly training a model that learns to decode in such a way, due to lack of time. It is intended as an approximation, or suggestion of what is potentially achievable, if we fully train the model to receive an execution output at each decoding step. Table 4 indicates the results on the 10 tasks, comparing the non-execution-guided GridCoder to the execution-guided GridCoder.

Task 5 fails on both approaches because when the neural network sees a large target grid (of 2x4 in this case), instead of attempting a 2x2 solution or a 1x4 solution which can then be scaled by composition to the correct 2x4 solution, it goes directly for 3x3. So the partial solutions that it suggests cannot be used to produce the correct solution.

Task 6 fails on both approaches because it fails to see what

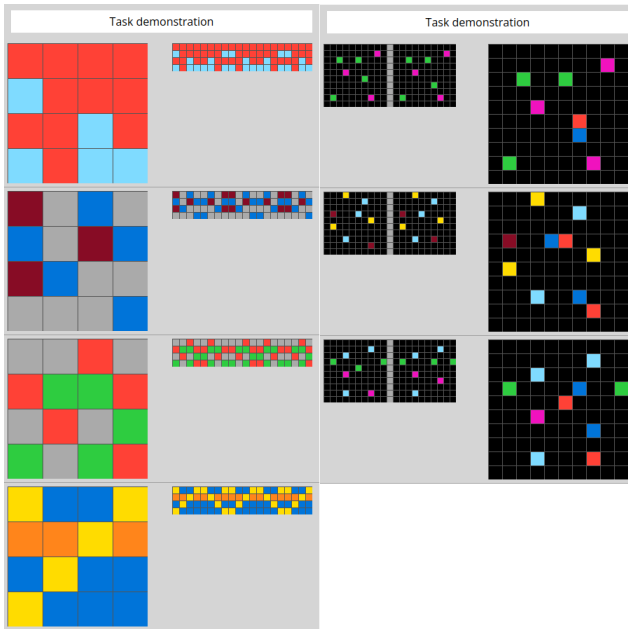


Figure 5: Two examples of GridCoder failures. Left: task bc4146bd. Right: task d47aa2ff.

needs to be done after the initial object cropping. It attributes near-zero probabilities to the primitives that yield the required horizontal tiling. Instead it seems to suggest programs that would mirror the input, horizontally or vertically.

Discussion

Generalization Characteristics

In the first experiment, the gap between *VLM-only* and *GridCoder* is strictly an out-of-distribution generalization gap. The fact that the search-enabled algorithms obtain a significantly better performance than the *VLM-only* approach suggests that searching over the probability space significantly extends the neural network’s reach beyond its training distribution. However, generalization obstacles remain, as indicated by the experiments in out-of-distribution generalization where the GridCoder approach only succeeds in 1 out of the 10 tasks. Furthermore, we note the following additional GridCoder failure cases:

bc4146bd.json This task (see figure 5) consists of “tiling” the original grid five times horizontally, using various grid transformations (e.g. rotations) each time. While the neural network was exposed to conceptually similar tasks during training, it has never seen fivefold horizontal (or vertical) tiling at all. As such, it struggled to generate the program structure that allows this fivefold tiling.

d47aa2ff.json This task (see figure 5) consists of copying over to the output grid the pixels that are common to both halves of the input grid. However, when there is a discrepancy between both halves, if the pixel only appears on the left grid, it must be colored red in the output, and if it is the reverse it must be colored blue. While the DSL allows us

to solve this task in principle, the solution is quite elaborate and there is no conceptually similar task in the training data.

The vast improvement in generalization capabilities from *VLM-only* to *GridCoder* can appear to be in opposition to the poor structural generalization findings from Table 4. We hypothesize that the out-of-distribution uncertainty that gets resolved by the probabilistic search is mainly related to a grid-level distribution shift. The test grids are distinct from the training ones, hence it struggles with properly predicting the types of per-grid transformations. This results in a relatively high entropy when choosing a primitive that produces some kind of grid-to-grid transformation.

However, we observe a very low entropy on tokens that set the overall program structure (for example, when differentiating “tiling” types of tasks from “split and merge” types of tasks). As a result, the search is able to resolve grid-level distribution shift by searching over the relatively high-entropy distribution, while the very low entropy distribution of program structure across task categories means that it cannot generalize out-of-distribution in that sense.

An example will better illustrate this. Consider, for example, the predicted probabilities for task 59341089. json:

Probabilities at position 0:

- <Identity>: 0.56
- vmirror: 0.10
- hmirror: 0.17
- rot90: 0.02
- rot180: 0.11
- rot270: 0.03

Probabilities at position 1:

- <Identity>: 0.43
- vmirror: 0.13
- hmirror: 0.22
- rot90: 0.02
- rot180: 0.15
- rot270: 0.05

Probabilities at position 2:

- <Identity>: 0.24
- vmirror: 0.16
- hmirror: 0.33
- rot90: 0.02
- rot180: 0.20
- rot270: 0.04

Probabilities at position 3:

- <New Level>: 0.02
- <Identity>: 0.31
- vmirror: 0.14
- hmirror: 0.26
- rot90: 0.03
- rot180: 0.19
- rot270: 0.04

Probabilities at position 4:

- <New Level>: 0.99

Probabilities at position 5:

- hconcat: 1.00

Probabilities at position 6:

- hconcat: 0.99

Probabilities at position 7:

- <New Level>: 0.99

Probabilities at position 8:

- hconcat: 0.99

Probabilities at position 9:

- <End of Sentence>: 0.99
- hconcat: 0.01

Only the classes whose probability is 0.01 or greater are shown. This task consists of tiling the input grid four times horizontally, while applying the following transforms to the grid, in order, from left to right: [hmirror, <Identity>, hmirror, <Identity>]. The token positions responsible for the "skeleton", or "structure" of the 4-way tiling program have very high probabilities (low entropy) associated to them. These are the token positions 4 to 9. The token positions responsible to determining what transformations are applied to the grid, and in what order, are, however, high entropy.

Further training the model will not solve these generalization issues. What has been observed is that the solved tasks are solved more rapidly when the neural network is better trained, because the probabilities are more accurate and less noisy. However, a token sub-sequence that has never been observed during training will still retain a probability near zero (obviously, since the probability of that sub-sequence in the training set is exactly zero), and a token sub-sequence that is always observed under a given "type" of input-output grid pairs will have a probability of one, since nothing else has ever been observed during training. This is precisely the behaviour that limits the generalization capabilities of LPS methods such as GridCoder.

Adding a bit of entropy such that no probability is ever zero or one, could help break through the hard probability barriers, at the expense of a vastly increased search space. Furthermore, this is unlikely to allow the search to discover extremely low probability (high complexity) modifications such as the one required to go from the most similar program for task bc4146bd (from figure 5) to its correct solution.

This is where Learning the transformation space (adding an execution-guided feedback and re-evaluating the probabilities on the intermediate grids) can help break through this generalization barrier, as suggested by the experimental results of Table 4.

Scaling the DSL

A potential criticism for the LPS approach is that, as we grow the DSL, the search space might become exponentially

more complex, so the solution time slows down proportionately. This would make this approach unlikely to scale to a competitive solution on ARC-AGI. Another possibility is that, as we add new primitives to the DSL, and train on new tasks, performance degrades on previous tasks.

Tables 2 and 3 paint a positive picture of the scaling properties of the LPS approach. First, we see from Table 2 that as we grow the DSL from version 1 to version 3, the tasks that the approach was able to solve in DSL version 1 are still solved in subsequent sections. There is no forgetting, and no significant loss in search performance that degrades the success rate. Meanwhile, the overall performance on the ARC-AGI evaluation set increases as we add primitives and corresponding training tasks.

Furthermore, we see that from Version 1 to Version 3, the mean solution time per primitive does not increase. In fact, it seems to decrease a bit from Version 2 to Version 3. This may simply be a fluke from having a better trained model in Version 3 and Version 2. As a reminder, DSL Version 1 has 74 primitives, DSL Version 2 has 89 primitives, and DSL Version 3 has 98 primitives. In summary, the solution time scales sub-linearly as we add more primitives to the DSL, which is a desirable property.

Flaws in LGS

The LGS approach reveals itself to be somewhat inefficient, though *a priori* it may be better at structural out-of-distribution generalization (because it is not trained on specific task solutions). Extensive work has been done to understand the causes of this inefficiency. The main hypothesized causes are as follows:

A* versus Q* The idea of a model that takes as input two grids and predicts their similarity suggests an A*-style usage pattern:

1. Each transformation in the DSL is applied to the input grid to obtain some intermediate grid.
2. Each pair of intermediate grid and target grid is fed to the similarity model to obtain their respective similarities.
3. The primitive that yields the biggest increase in similarity is selected.
4. This is done iteratively until a solution is found or some budget limit is reached.

Agostinelli et al. (2024) show that this approach is both temporally and spatially inefficient, and propose Q* search, which can be differentiated as follows:

1. The input grid and output grid pair is fed to the neural network, which in one shot simultaneously predicts all of the similarity (or *cost-to-go*) values for each of the primitives.
2. The best one is selected, and used to generate the next intermediate state.
3. This is done iteratively until a solution is found or some budget limit is reached.

Hence, where the grid-to-grid similarity approach results in N evaluations where N is the size of the DSL, the Q* approach has a constant number of evaluations with respect to

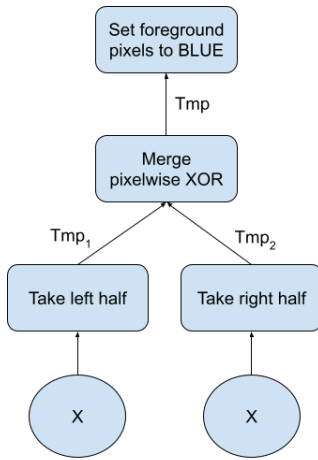


Figure 6: Example of a task solution

the DSL size. Note that this means we are inherently moving from a LGS approach to a LTS approach. Instead of having a model that learns a latent grid space contrastively and outputs similarities between two vectors, we have a model that necessarily has to learn the effect of applying each transformation from the DSL, and evaluate its cost-to-go. This becomes even more evident when we supervise this model to output probabilities over the DSL instead of a cost-to-go.

The choice of *cost-to-go* or probability in this context is debatable, but we hypothesize that learning probabilities is more direct and efficient than the cost-to-go. This is the case both from the standpoint of learning the underlying relationships and from the standpoint of generating training data. It is easy to accidentally produce incorrect cost-to-go labels from procedurally generated programs with redundant primitives.

Incomplete information A program is a syntax tree, not a flat sequence of grid-to-grid operations. Because of that, as we execute a program from the bottom of the tree upwards, there are potentially multiple operations that need to be done in "parallel" (not computationally speaking, but syntactically speaking). Thus, at any given point in execution time, the intermediate state of a program is in fact a list of objects (not necessarily grids, it can be lists of indices, integers, etc.).

To select the next primitive based only on the output of one of the previous transformations means that we rely on incomplete information. As an example, refer to figure 6. In this diagram, *X* represents the input grid. The left half and right half are taken, and merged together using XOR logic. Consider what happens at training time, when the model is learning to predict the cost-to-go for the first primitive (i.e. "Take left half" or "Take right half").

The choice of whether to use XOR, OR, AND, NOR, or any other merge-based logic requires knowing which two input grids are fed to these primitives. There are several different cost-to-go values that correspond to the same pair (output of "Take right half" and the target grid), depending on the left half that is used in the procedurally gen-

erated program. Because the model necessarily involves a similarity between exactly two grids, the model has to learn without knowing the output of the left half operation. Thus, the model is trained on incomplete information, resulting in weak cost-to-go estimates.

ARC-AGI is not a grid-to-grid problem ARC-AGI is fundamentally not a grid-to-grid problem, when we consider the execution trace of a solution. That is, most of the steps involved in solving an ARC-AGI task are not strictly grid-to-grid transformations and thus do not result in an intermediate grid that can be evaluated for its similarity to the target. This is related to the previous argument, in the sense that a solution program for most ARC-AGI tasks is not a simple sequence of grid-to-grid transformations, but instead can be thought of as a sequence of intermediate program states which themselves are lists of various kinds of intermediate objects and values. As such, the grid-to-grid similarity concept only potentially allows guiding a small portion of the overall search for a solution.

In summary, these conceptual flaws, combined with the inferior empirical results, led to the choice of selecting the LPS approach instead. However, it should be noted that the LGS idea also presents some essential characteristics that the selected LPS approach is lacking. This is why we attempt to combine the best of both approaches via the suggested LTS approach.

Limitations and Future work

Limited DSL The work done so far has been a proof-of-concept centered on finding an optimal neurally-guided search algorithm given some pre-determined DSL and a set of tasks that are known to be solvable within this DSL. Very limited work has been done on improving the DSL. This is, therefore, the main limitation of this work.

Scaling the algorithms presented here to a competitive level will certainly involve completing the DSL by introducing all of the necessary primitives required to solve all ARC training and evaluation tasks. This is not as simple as going through the existing tasks and implementing primitives for every high-level functionality seen.

We aspire to making the DSL more flexible and general as well, as we currently deem it too high-level and specific. We aim for a certain compromise between a purely generic programming language and a highly specific DSL such as the current one. We expect this work, on its own, to make our approach competitive on the ARC hidden test set.

Functional DSL Based on the work by Hocquette and Cropper (2024), it appears that a relational type of DSL can lead to more efficient search than the kind of functional paradigm currently being used in these experiments. Intuitively, relational decomposition of tasks means that partial verification of a subset of the rules is possible while searching for a solution. In contrast, our current approaches can only verify the full program in a binary "fully correct" or "fully incorrect" manner. Relational decomposition therefore allows a more gradual progression towards the full solution.

get_objects primitives In retrospect, the current implementation of six different primitives for object detection is very sub-optimal. In spite of having implemented six different notions of what is an object, there are still many cases where object detection primitives fail. Going forward, we will instead implement this primitive as an "object segmentation" neural network that, we hope, will allow us to reduce the number of object detection primitives to one, while simultaneously covering more use cases.

Program syntax The program representation syntax needs to be improved. First, the difference in program structure between four tiles and five tiles is relatively large, when in principle it should be quite small. Here is an example program sequence that tiles an input grid four times horizontally:

```
[hmirror, <Identity>, hmirror, <Identity>,
<New Level>, hconcat, hconcat, <New Level>,
hconcat, <End of Sentence>]
```

Here is the program sequence to tile an input grid five times horizontally:

```
[hmirror, <Identity>, hmirror, <Identity>,
hmirror, <New Level>, hconcat, hconcat,
<Identity>, <New Level>, hconcat,
<Identity>, <New Level>, hconcat, <End
of Sentence>]
```

There is a difference of five tokens between the two, spread throughout all of the levels of the program. Clearly, the conceptual similarity of these two tasks is not reflected in the syntactical similarity of their solutions. This certainly does not help the search.

Another problem is that a lot of duplication is necessary in the lower levels. For example, in the last task used for the results of Table 4, the goal is to crop the largest object, split it in half horizontally and merge it pixelwise using OR logic. In order to take the left half and the right half on the same input, the current program syntax forces us to duplicate the entire sub-tree that generates the cropped object. The corresponding solution program is:

```
[get_objects1, get_objects1, get_objects1,
get_object_size, get_objects1, get_objects1,
get_object_size, <New Level>, <Identity>,
for_each, <Identity>, for_each, <New
Level>, keep_largest, keep_largest, <New
Level>, lefthalf, righthalf, <New Level>,
cellwiseOR, <End of Sentence>]
```

The entire sub-tree [get_objects1, get_objects1, get_object_size, <New Level>, <Identity>, for_each, <New Level>, keep_largest] needs to be duplicated twice, thereby increasing the solution by six tokens. This is obviously not very efficient.

Learning the transform space The empirical results suggest that GridCoder is limited to proposing program structures that have been seen during training. The neural network does not get to observe explicitly or even generate implicitly the intermediate state of a program as it derives its solution. Instead, it currently has to be able to solve all of it in one go based on the provided input-output examples.

In other words, as it generates new tokens in the decoding loop, it does not get to look at the intermediate program state to decide on the next steps. We certainly do not do this explicitly (yet), and we know from how transformer decoders work that it does not naturally maintain a hidden state between decoding steps. It only gets to look at the decoded sequence so far, as well as the output of the encoder. As a result, it tends to lazily rely on positional encoding and pre-determined sequence templates that it learned during training. The consequence of this is that, as the search stumbles into a territory where the model should know what to do to transform the intermediate grids into the target grids, it cannot leverage this knowledge.

It could be argued that a sufficiently flexible transformer decoder module could potentially learn to maintain implicitly the current intermediate state by applying the current token sequence to the encoder output. This is, in fact, what we were hoping to see, and our architectural choice of a particularly lean decoder could be the culprit. However, it's not clear what would encourage it to learn to do this if it can lazily rely on the positional encoding and the token sequence as it does at the moment. Even if this were to be true, however, it is still the case that explicitly feeding back the intermediate state would constitute an inductive bias that can guide the learning process towards faster and better convergence. This makes LTS a promising paradigm regardless.

Related Work

Execution-guided program induction Execution-guided program induction can apply to both the LGS and LTS paradigms, in our current nomenclature. Much work has focused on this approach (Nye et al. 2020; Chen, Song, and Tian 2021; Chen, Liu, and Song 2018). Recently, Kapur, Jenner, and Russell (2024) use an image similarity heuristic to discover programs that transform an input image into a target image. The training procedure follows a denoising paradigm, and the search process is based on a diffusion, in which program edits are gradually made to arrive at the denoised image.

Neurally-guided probabilistic program induction DreamCoder (Ellis et al. 2021) is a program induction approach not unlike the one proposed here, under what we refer to as the LPS paradigm. There is a neural network that learns to output probabilities over the DSL, and a search algorithm suggests programs to evaluate against the example inputs and outputs. It brings as its main innovations a DSL growth mechanism, in which frequently occurring semantically equivalent subroutines are added to the DSL.

GridCoder does not use such a DSL growth mechanism. On one hand, we see the limited amount of natural ARC-AGI examples as an obstacle to making this feature particularly useful. Correctly fine-tuning the granularity of the DSL, and having the relevant reusable primitives, is something we prefer to do manually. Similarly, DreamCoder has a "dreaming" functionality in which new tasks are randomly generated from the DSL in order to train the neural network. We found that a purely random task generation approach tends to generate several incoherent, redundant, unnatural

tasks. Instead, in GridCoder we adopt an intermediate position in which we manually design high-level task categories or concepts, but allow some random variation in their parameters and details.

A variety of other neurally guided program induction algorithms have been proposed (Fijalkow et al. 2021; Kapur, Jenner, and Russell 2024; Devlin et al. 2017). To the best of our knowledge, GridCoder differs from these algorithms mainly in the neural architecture used and in the fact that we lean heavily into the auto-regressive sequence output as the basis for enumerating the probability space.

Most approaches use a neural network that outputs probabilities for each of the n -grams that can be formed by their DSL, relying instead more heavily on search than learning. Learning bigram probabilities, for example, means that the model learns to predict a probability for the children of a primitive p_i conditionally on the choice of p_i (and p_i alone). Our auto-regressive approach, instead, learns the probability over a child node given the entire program sequence (and the encoder embedding) so far.

As an illustration of this significant distinction, Alford (2021) experimented on using DreamCoder on ARC-AGI, and concluded that the main bottleneck was the search efficiency. In our experiments with GridCoder, the main bottleneck was not the search efficiency (since the neural network was very good at identifying program structure), but its challenges at generalizing to unseen program structures (see the Discussion section for more details).

Conclusion

In summary, we compared three different program induction paradigms: LGS, LPS (specifically GridCoder) and LTS. Given the experimentation and analysis on GridCoder, we are confident that scaling up the DSL and training data could make it reach a performance that is competitive with the current top solutions (unless we later discover that these approaches explicitly solve the structural generalization problem instead of training on a dataset that unknowingly includes solutions to the hidden test set).

However, given its fundamental challenges with structural generalization, we are not confident that the LPS approach could solve ARC-AGI in any true sense. Instead, preliminary experimentation on the LTS approach suggests that GridCoder could be extended via execution-guided search to solve these generalization obstacles. This will be the focus of our future work.

References

Agostinelli, F.; Shperberg, S. S.; Shmakov, A.; McAleer, S.; Fox, R.; and Baldi, P. 2024. Q* Search: Heuristic Search with Deep Q-Networks.

Alford, S. 2021. *A Neurosymbolic Approach to Abstraction and Reasoning*. Ph.D. thesis, Massachusetts Institute of Technology.

Chen, X.; Liu, C.; and Song, D. 2018. Execution-guided neural program synthesis. In *International Conference on Learning Representations*.

Chen, X.; Song, D.; and Tian, Y. 2021. Latent execution for neural program synthesis beyond domain-specific languages. *Advances in Neural Information Processing Systems*, 34: 22196–22208.

Chollet, F.; Knoop, M.; Landers, B.; Kamradt, G.; Jud, H.; Reade, W.; and Howard, A. 2024. ARC Prize 2024. <https://kaggle.com/competitions/arc-prize-2024>. Kaggle.

Devlin, J.; Uesato, J.; Bhupatiraju, S.; Singh, R.; Mohamed, A.-r.; and Kohli, P. 2017. Robustfill: Neural program learning under noisy i/o. In *International conference on machine learning*, 990–998. PMLR.

Ellis, K.; Wong, C.; Nye, M.; Sablé-Meyer, M.; Morales, L.; Hewitt, L.; Cary, L.; Solar-Lezama, A.; and Tenenbaum, J. B. 2021. DreamCoder: Bootstrapping Inductive Program Synthesis with Wake-Sleep Library Learning. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, PLDI 2021, 835–850. New York, NY, USA: Association for Computing Machinery. ISBN 9781450383912.

Fijalkow, N.; Lagarde, G.; Matricon, T.; Ellis, K.; Ohlmann, P.; and Potta, A. 2021. Scaling Neural Program Synthesis with Distribution-based Search. *CoRR*, abs/2110.12485.

Hocquette, C.; and Cropper, A. 2024. Relational decomposition for program synthesis. *arXiv preprint arXiv:2408.12212*.

Kapur, S.; Jenner, E.; and Russell, S. 2024. Diffusion On Syntax Trees For Program Synthesis. *arXiv preprint arXiv:2405.20519*.

Nye, M.; Pu, Y.; Bowers, M.; Andreas, J.; Tenenbaum, J. B.; and Solar-Lezama, A. 2020. Representing partial programs with blended abstract semantics. *arXiv preprint arXiv:2012.12964*.

Schrittwieser, J.; Antonoglou, I.; Hubert, T.; Simonyan, K.; Sifre, L.; Schmitt, S.; Guez, A.; Lockhart, E.; Hassabis, D.; Graepel, T.; et al. 2020. Mastering atari, go, chess and shogi by planning with a learned model. *Nature*, 588(7839): 604–609.

Appendix - Description of the DSL

The DSL revolves around the concept of a Grid, which is a custom class that contains the following information about a grid:

- Its pixels, represented as a list of (X, Y, color) tuples.
- A utility method that converts from the pixels list to a 2D grid representation (a tuple of tuples where each value is the color at that position).
- Its width.
- Its height.
- Its upper-left corner position in the outer Grid. Default: (0, 0).

DSL Version 1:

1. `set_fg_color1` to `set_fg_color9`: set all the non-zero pixels to color 1 to 9.
2. `shift_left`: translate pixels once to the left.

3. `shift_right`: translate pixels once to the right.
4. `shift_up`: translate pixels once upward.
5. `shift_down`: translate pixels once downward.
6. `vmirror`: flip the grid vertically.
7. `hmirror`: flip the grid horizontally.
8. `rot90`: rotate the grid 90 degrees.
9. `rot180`: rotate the grid 180 degrees.
10. `rot270`: rotate the grid 270 degrees.
11. `tophalf`: keep only the top half of the grid.
12. `bottomhalf`: keep only the bottom half of the grid.
13. `lefthalf`: keep only the left half of the grid.
14. `righthalf`: keep only the right half of the grid.
15. `symmetrize_left_around_vertical`: copy the left half to the right half, in a mirroring fashion.
16. `symmetrize_right_around_vertical`: copy the right half to the left half, in a mirroring fashion.
17. `symmetrize_top_around_horizontal`: copy the top half to the bottom half, in a mirroring fashion.
18. `symmetrize_bottom_around_horizontal`: copy the bottom half to the top half, in a mirroring fashion.
19. `upscale_horizontal_by_two`: duplicate pixels once, horizontally.
20. `upscale_vertical_by_two`: duplicate pixels once, vertically.
21. `upscale_by_two`: duplicate pixels once, vertically and horizontally.
22. `gravitate_right`: stack the pixels to the right, as if they had fallen rightward.
23. `gravitate_left`: stack the pixels to the left, as if they had fallen leftward.
24. `gravitate_up`: stack the pixels at the top, as if they had fallen upward.
25. `gravitate_down`: stack the pixels at the bottom, as if they had fallen downward.
26. `gravitate_left_right`: stack the leftmost pixels to the left, and the rightmost pixels to the right.
27. `gravitate_top_down`: stack the top-half pixels at the top, and the bottom-half pixels at the bottom.
28. `topthird`: keep only the top third of the grid.
29. `vcenterthird`: keep only the middle third of the grid (vertically).
30. `bottomthird`: keep only the bottom third of the grid.
31. `leftthird`: keep only the leftmost third of the grid.
32. `hcenterthird`: keep only the middle third of the grid (horizontally).
33. `rightthird`: keep only the rightmost third of the grid.
34. `cellwiseOR`: 2-argument primitive that takes in two grids of the same shape, and merges them into one by applying OR logic to the foreground pixels (prioritizing the first argument's pixel color if both grids have a foreground pixel at a certain position).
35. `cellwiseAND`: 2-argument primitive that takes in two grids of the same shape, and merges them into one by applying AND logic to the foreground pixels (using the first argument's pixel color).
36. `cellwiseXOR`: 2-argument primitive that takes in two grids of the same shape, and merges them into one by applying XOR logic to the foreground pixels.
37. `cellwiseDifference`: 2-argument primitive that takes in two grids of the same shape, and merges them into one by setting the target grid cells the same color as the first argument's foreground pixels when the second argument does not contain a foreground pixel at that position.
38. `cellwiseNOR`: 2-argument primitive that takes in two grids of the same shape, and places a foreground pixel of the same color as the dominant color in the first argument if neither grid has a foreground pixel at a given position.
39. `vconcat`: 2-argument primitive that vertically concatenates two grids.
40. `hconcat`: 2-argument primitive that horizontally concatenates two grids.
41. `color_change`: 3-argument primitive that takes as input a grid, and two integers representing colors. Every pixel in the grid that is of the color of the first integer gets converted to the color of the second integer. At search time, the integer values get resolved by an enumerative search built into the program evaluation logic.
42. `invert_colors`: Finds the least dominant color and the most dominant color, and swaps them in the grid.
43. `first_quadrant`: Return the top-left quadrant of the grid.
44. `second_quadrant`: Return the top-right quadrant of the grid.
45. `third_quadrant`: Return the bottom-left quadrant of the grid.
46. `fourth_quadrant`: Return the bottom-right quadrant of the grid.
47. `hfirstfourth`: Returns the leftmost fourth of the grid (horizontal splitting).
48. `hsecondfourth`: Returns the middle-left fourth of the grid (horizontal splitting).
49. `hthirdfourth`: Returns the middle-right fourth of the grid (horizontal splitting).
50. `hlastfourth`: Returns the rightmost fourth of the grid (horizontal splitting).
51. `vfirstfourth`: Returns the topmost fourth of the grid (vertical splitting).
52. `vsecondfourth`: Returns the middle-top fourth of the grid (vertical splitting).
53. `vthirdfourth`: Returns the middle-bottom fourth of the grid (vertical splitting).
54. `vlastfourth`: Returns the bottommost fourth of the grid (vertical splitting).
55. `duplicate_top_row`: duplicates the top row.

56. `duplicate_bottom_row`: duplicates the bottom row.
57. `duplicate_left_column`: duplicates the left column.
58. `duplicate_right_column`: duplicates the right column.
59. `remove_outline`: removes the top row, bottom row, left column, right column.
60. `shear_grid_left`: from the bottom up of the grid, it gradually shifts left the pixels by an increment that increases by 1 each row.
61. `shear_grid_right`: from the bottom up of the grid, it gradually shifts right the pixels by an increment that increases by 1 each row.
62. `shear_grid_zigzag`: from the bottom up of the grid, it shifts the pixels horizontally by an increment that cycles between [-1, 0, +1].
63. `insert_outline`: it adds a black row at the top and bottom, and a black column at the left and the right.
64. `get_major_pixel`: finds the most frequent color in the grid, and outputs a 1x1 grid of that color.
65. `get_minor_pixel`: finds the least frequent color in the grid, and outputs a 1x1 grid of that color.
66. `upscale_by_three`: upsamples the grid by a factor of three.

DSL Version 2 - added primitives:

1. `get_objects1`: returns a list of Grids, each representing an object. This variant is good at isolating rectangular shapes (filled or empty) from noisy backgrounds.
2. `get_objects2`: returns a list of Grids, each representing an object. This variant is good at finding four-way (up, down, left, right) adjacent objects against a black background.
3. `get_objects3`: returns a list of Grids, each representing an object. This variant is good at finding eight-way adjacent objects against a black background.
4. `get_objects4`: returns a list of Grids, each representing an object. This variant is good at finding four-way (up, down, left, right) adjacent objects against of any color, as long as the object's color is uniform.
5. `get_objects5`: returns a list of Grids, each representing an object. This variant is good at finding eight-way adjacent objects against of any color, as long as the object's color is uniform.
6. `get_objects6`: returns a list of Grids, each representing an object. This variant is good at isolating shapes that are potentially disjoint, but clustered in a grid-like pattern.
7. `compress_objects_linear`: this generates an output grid from a list of objects by concatenating them to each other in horizontal or vertically fashion, automatically determined from their relative positions.
8. `compress_objects_quad`: this generates an output grid from a list of objects by concatenating them to each other in a rectangular fashion (mostly 2x2, 2x3, 3x2 and

3x3 patterns), automatically determined from their relative positions.

9. `compress_objects_quad_pad`: similar to `compress_objects_quad`, but introduces one row and column of black padding between the concatenated objects.
10. `foreach`: 2-argument primitive. The first argument is a list of objects to which to apply a function. The second argument is the lambda function to apply to each object. The input argument of that lambda function must be a Grid type.
11. `apply_to_grid`: 2-argument primitive. The first argument is the original grid itself, and the second argument is the list of transformed objects to overlay onto the grid (blanking out the previous objects).
12. `cellwise_OR_list`: like `cellwise_OR`, but takes a list of two or more objects, instead of strictly only taking two arguments.
13. `get_pixels`: returns a list of integers representing all the colors of the pixels in the grid (this is not a set, but a list, so it can be used to count the number of pixels in an object, for example).
14. `get_object_size`: returns an integer, the number of pixels in a Grid instance.
15. `count`: counts the number of elements in a list.

DSL Version 3 - added primitives:

1. `filter_largest`: 2-argument primitive that takes as input a list of objects and a list of integer values. It returns a list of objects in which the object that has the largest associated integer value is removed.
2. `filter_smallest`: 2-argument primitive that takes as input a list of objects and a list of integer values. It returns a list of objects in which the object that has the smallest associated integer value is removed.
3. `keep_largest`: 2-argument primitive that takes as input a list of objects and a list of integer values. It returns the object from the list that has the largest associated integer value.
4. `keep_smallest`: 2-argument primitive that takes as input a list of objects and a list of integer values. It returns the object from the list that has the smallest associated integer value.
5. `is_h_symmetrical`: returns True if the grid is horizontally symmetrical, False otherwise.
6. `is_v_symmetrical`: returns True if the grid is vertically symmetrical, False otherwise.
7. `logical_not`: takes as input a Boolean value, and returns its negation.
8. `keep_boolean`: 2-argument primitive that takes as input a list of objects and a list of Boolean values. It returns the first object from the list that has True as its associated value.
9. `filter_boolean`: 2-argument primitive that takes as input a list of objects and a list of Boolean values. It returns a list of objects in which the objects whose associated value is True are removed.