

# Evolving and Refining Code: A Neuro-Evolutionary Program Synthesis for Compositional Reasoning

Anonymous submission

## Abstract

We propose a neuro-evolutionary framework for program synthesis that combines Genetic Programming (GP) with a frozen Large Language Model (LLM) to tackle problems that require planning and compositional reasoning. Rather than generating programs from scratch, our method evolves candidate solutions using GP and then refines failed programs via LLM-based code edits. The LLM purely operates through a self-feedback prompting loop and serves as a semantic editor for improving GP-generated programs. The system operates in two stages. Stage 1: GP performs evolutionary search over a type-safe Domain-Specific Language (DSL), adopted from prior work, using the Distributed Evolutionary Algorithms in Python (DEAP) framework with strict typing to ensure semantic validity. Stage 2: If GP fails to solve the given task, the best candidate program so far is passed to LLM for repair. LLM receives the final best program, current input and output examples, program fitness score, and iteratively suggests edited programs until fitness improvements or the maximum retry loop is reached. Although recent methods rely heavily on fine-tuned LLMs, evolutionary approaches remain underexplored, especially in the Abstraction and Reasoning Corpus (ARC) benchmark. We show that frozen LLM can enhance genetic programming for code refinement. Experiments on the ARC reveal performance gains of the proposed approach over the GP-only baseline. Our goal is to show that evolutionary learning methods can be enhanced through LLM code editing, forming a neuro-evolutionary pipeline.

## Introduction

The task of program synthesis is to automatically generate a program that fulfills a given specification, often represented by input-output examples (O’Neill and Spector 2020). A robust program synthesizer must compose programs from reusable sub-components and generalize beyond training examples—combining learned subprograms in new ways to synthesize longer or conceptually novel programs that have the ability to reason over multi-step transformations (Shi et al. 2023). While neural based program synthesis approaches have made progress in this area these approaches have important limitations: (a) they are computationally expensive and hard to train, (b) a model has to be trained for each task (program) separately, (c) it is hard to interpret or verify the correctness of the learnt mapping (Parisotto et al. 2016).

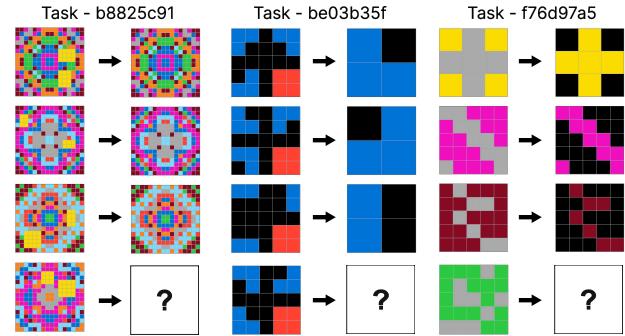


Figure 1: Sample tasks from the ARC benchmark. Each task provides a few input-output grid examples, where the goal is to infer the rule that maps input grids (left) to output grids (right) and apply the discovered rule to the test input (final row). Task b8825c91 is to fill out the yellow patches based on the pattern of the remaining part. Task be03b35f is to guess the red patch based on the characteristics of other objects in the grid. Since all forms of  $N$  times 90-degree clockwise rotation of the left top object appear in the input and output grid, when  $N = 3$ , the rule of the task is to apply rotations on the left top object of the input grid. Task f76d97a5 is to change the color gray to a non-gray color in the grid and the non-gray color to black.

Given these limitations, alternative approaches based on evolutionary search have received renewed interest in program synthesis. In particular, Genetic Programming (GP) has remained a robust and interpretable framework for learning-based program synthesis, due to its ability to evolve structured programs over Domain-Specific Languages (DSLs) using selection, crossover, and mutation (Sobania, Schweim, and Rothlauf 2021). Unlike neural models that require gradient-based training and massive data, GP directly manipulates symbolic program trees, enabling compositional reuse and open-ended exploration of the program space (Koza 1994b). These characteristics make GP especially appealing for synthesis tasks that require multi-step reasoning, hierarchical abstraction, and few-shot generalization, aligning closely with the capabilities expected of a strong synthesizer (Muslea 1997). How-

ever, GP suffers from a well-known limitation: it frequently stalls in local optima, especially in large, rugged program spaces, where small mutations often degrade fitness rather than improve it (Pantridge and Helmuth 2023).

In this work, we explore whether Genetic Programming alone is sufficient for solving compositional and multi-step reasoning tasks in the Abstraction and Reasoning Corpus (ARC) - a benchmark explicitly designed to test few-shot generalization, multi-step reasoning, and goal-directed planning abilities (Chollet 2019). Three example ARC tasks from the training set are shown in Figure 1.

The ARC tasks require abstract transformations over grid-based inputs, often involving compositional rules and reasoning over multiple steps to arrive at the correct output. To investigate this, we adopted a type-safe GP framework to evolve programs over a Domain-Specific Language (DSL), and used it as a base synthesizer for solving ARC tasks. However, consistent with long-standing observations in the GP literature (Koza 1994b; Goldberg 1989), we found that GP often stalls in local optima, particularly when small mutations fail to correct semantically meaningful yet structurally subtle errors in programs.

To overcome this limitation, and drawing on the emerging paradigm where LLMs are used to refine and debug program outputs (Tang et al. 2024; Wei, Xia, and Zhang 2023; Xu and ... 2023), we adopt a hybrid feedback mechanism: when the Genetic Programming baseline gets stuck, a frozen LLM is prompted with the best-evolved program, input-output pairs, and current fitness to perform code repair. We introduce a neuro-evolutionary ARC solver that incorporates LLMs' code repair into the GP pipeline. The LLM acts as a semantic editor, proposing code refinements of the program made by the GP. We observe that this LLM-guided editing process enables GP to escape local optima, producing correct programs that were otherwise unreachable through evolution alone.

Figure 2 outlines the two stages in the pipeline in which Stage 1 runs GP program synthesis with population  $P$  and Stage 2 passes the best individual program (if the GP failed to solve the given task) to the LLM for at most  $R$  repair rounds. The typed DSL enforces safety and keeps all solutions symbolic. The pseudo-code of this framework is shown in Algorithm 1.

**Our contribution.** We propose a neuro-evolutionary ARC program synthesizer that combines GP with an LLM-based code repair stage. When GP reaches near-solutions, the LLM edits the code to fix semantic errors to maximize the fitness. This hybrid flow enables more effective and open-ended exploration.

## Related Work

### Genetic Programming

Genetic Programming (GP) is an evolutionary algorithm for searching combinatorial spaces by evolving typed program trees  $T \in \mathcal{T}$  to maximize a fitness function  $f : \mathcal{T} \rightarrow \mathbb{R}$  (Koza 1994a). Each individual is represented as  $(p, f(p))$ , and a population of  $N$  individuals is denoted

$\mathcal{P} = \{(p_1, f(p_1)), \dots, (p_N, f(p_N))\}$ . Over  $G$  generations, GP applies three core operators:

- **Selection.** SEL :  $\mathcal{T}^N \times \mathbb{R}^N \rightarrow \Delta(\mathcal{T}^N)$ .
- **Crossover.** XO :  $\mathcal{T}^\nu \rightarrow \Delta(\mathcal{T})$ .
- **Mutation.** MUT :  $\mathcal{T} \rightarrow \Delta(\mathcal{T})$ .

However, traditional GP does not enforce that functions are applied to compatible data types, which becomes problematic when programs involve multiple data types and type-specific operations. This lack of constraint can significantly expand the search space and reduce generalization performance. Strongly Typed Genetic Programming (STGP) addresses this by introducing type constraints, ensuring that all functions operate only on appropriately typed inputs. By doing so, STGP reduces invalid program generation and guides the search toward semantically valid solutions (Montaña 1995).

### The ARC Benchmark and Approaches

The **Abstraction and Reasoning Corpus** (ARC) (Chollet 2019) is designed to evaluate fluid intelligence through few-shot grid transformations, where systems must infer latent rules from just a few input-output examples. To aid interpretability and concept-level evaluation, **ConceptARC** (Moskvichev, Odonard, and Mitchell 2023) extends ARC by categorizing tasks into high-level conceptual classes such as symmetry, object movement, and counting. This categorization enables more targeted analysis of model capabilities across distinct reasoning types. Despite recent progress, ARC remains a challenging benchmark for both neural and symbolic systems due to its emphasis on compositional generalization and minimal supervision.

**Evolutionary Learning on ARC.** Some approaches have explored evolutionary synthesis for ARC tasks. DSL with GE (Fischer et al. 2020) used grammatical evolution over DSLs, achieving top 4 % in the Kaggle competition, but solving only 30 tasks. LCSs (Coombe, Howard, and Browne 2024) applied Learning Classifier Systems with limited success (4.8 %), while Claude + EA (Berman 2024; Chollet et al. 2024) and SOAR (Pourcel, Colas, and Oudeyer 2025) introduced hybrid LLM-evolution strategies. SOAR achieves 52 % accuracy, but requires costly model queries and fine-tuning.

**Domain Specific Language (DSL).** A domain-specific language (DSL) is a programming or specification language designed to express solutions within a particular problem domain through concise, high-level abstractions (Van Deursen, Klint, and Visser 2000). For ARC-style reasoning tasks, where solutions involve complex visual and structural transformations over grid inputs, DSLs provide a symbolic and interpretable alternative to raw code generation. One widely adopted DSL in this ARC domain is the DSL introduced by Hodel (Hodel 2024), which defines a set of primitives for object manipulation, color filtering, grid operations, and compositional reasoning (Hodel 2024). Our work builds directly on this DSL, leveraging its typed expression to enable safe and structured search via STGP.

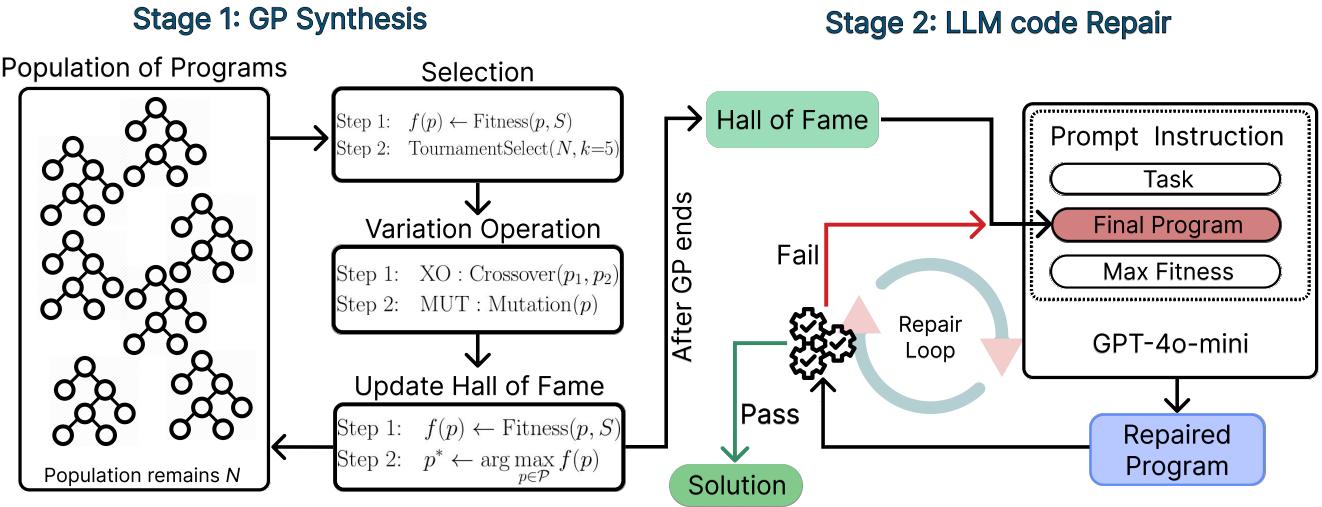


Figure 2: Framework Overview: A population of programs is evolved using Strongly Typed Genetic Programming (STGP) via crossover and mutation within a DSL. After each generation, programs are evaluated for task fitness. If no candidate solves the task, the best program is passed to a frozen LLM along with prompt instructions and fitness feedback. The LLM iteratively attempts to repair the program through semantic edits without fine-tuning.

## System Overview

### Stage 1: Genetic Programming-based Program Synthesis

We begin with STGP that synthesizes programs over a type-safe DSL, adopted from (Hodel 2024), implemented using the DEAP evolutionary computation library (Fortin et al. 2012). Our formulation and operators follow the classical genetic programming paradigm (Koza 1994b; Goldberg 1989). Each candidate program is denoted as  $p \in \mathcal{P}$ , where  $\mathcal{P}$  denotes the program space, and is represented as a typed expression tree  $T = [n_0, n_1, \dots, n_k]$ . Here, each internal node  $n_i$  corresponds to a DSL primitive of the form  $r : (\tau_1, \dots, \tau_m) \rightarrow \tau_{\text{out}}$ , indicating a function with input types  $\tau_1, \dots, \tau_m$  and output type  $\tau_{\text{out}}$ . Note that  $T$  denotes the program’s expression tree, representing the syntactic structure of a candidate program, while  $p$  is the executable function that results from evaluating  $T$  on a given input. The leaf nodes are called terminals, whose values are either constants, booleans, or the input grid. All types  $\tau$  are drawn from a finite set that has been defined in this work (Hodel 2024). Some of the DSL types defined are (`Grid`, `Colour`, `List[Grid]`). The GP system maintains a population  $\mathcal{P} = \{p_1, p_2, \dots, p_N\}$ , where each individual  $p$  is evaluated on a training set  $S = \{(x_i, y_i)\}_{i=1}^N$ . The fitness of a program is computed using a pixel-wise similarity function  $\text{sim}(\cdot, \cdot)$  as used in most prior works (Xu, Khalil, and Sanner 2023).

$$f(p) = \sum_{i=1}^N \text{sim}(p(x_i), y_i)$$

The goal of evolution is to identify the highest-fitness program:

$$p^* = \arg \max_{p \in \mathcal{P}} f(p)$$

The initial population is created using `genFull` that creates a full program tree by first picking the target depth and then recursively expanding nodes until every root-to-leaf path has  $h$  depth. At each generation, the population undergoes three genetic operations adopted from (Fortin et al. 2012). The use of strong typing ensures that all generated programs are executable and meaningful. This not only avoids runtime errors but also reduces the search space significantly, enabling GP to explore functional program structures. Note that throughout this paper, we will use the terms STGP and GP interchangeably.

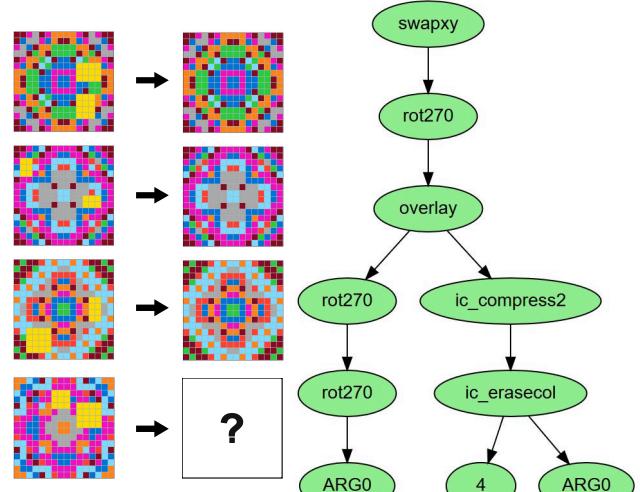


Figure 3: Task `b8825c91` and the symbolic expression tree generated by our GP pipeline.

---

**Algorithm 1:** Neuro-evolutionary Program Synthesizer via STGP and LLM Repair

---

```

Input: Training examples  $\mathcal{S} = \{(x_i, y_i)\}_{i=1}^N$ 
Output: function solve(inp) mapping  $x \mapsto y$ 
1: Stage 1: GP-based Program Synthesis
2: Initialize population  $\mathcal{P}^{(0)} = \{p_1, \dots, p_N\} \sim \mathcal{U}(\mathcal{P})$ 
3: for generation  $g = 1$  to  $G$  do
4:   foreach  $p \in \mathcal{P}^{(g)}$  do
5:     Compute fitness:
6:      $f(p) = \sum_{i=1}^N \text{sim}(p(x_i), y_i)$ 
7:     if  $\exists p^* \in \mathcal{P}^{(g)}$  such that  $f(p^*) = f_{\max}$  then
8:       return  $p^*$ 
9:   Select parents via tournament selection
10:  Apply type-safe crossover and mutation
11:  Form next generation:  $\mathcal{P}^{(g+1)}$ 
12:  Update Hall-of-Fame with top-performing individual
13: Let  $p^* = \arg \max_{p \in \mathcal{P}^{(g)}} f(p)$ 
14: Stage 2: LLM-based Program Repair
15: for repair round  $r = 1$  to  $R$  do
16:   Construct prompt  $\pi_r$  from I/O examples, fitness and DSL code for  $p^*$ 
17:   Query LLM with  $\pi_r$  to obtain candidate:
18:    $p'_r \sim \text{LLM}(\pi_r)$ 
19:   Evaluate fitness:  $f(p'_r) = \sum_{i=1}^N \text{sim}(p'_r(x_i), y_i)$ 
20:   if  $f(p'_r) > f(p^*)$  then
21:     Update  $p^* \leftarrow p'_r$ 
22: return Final program  $p^*$ 

```

---

211 **Program Representation.** Our system represents candidate programs as symbolic expression trees. Figure 3 shows an example generated by our GP pipeline. The tree is built bottom-up from terminal nodes such as constants and input grids, which are composed through typed DSL primitives. In Task b8825c91, the program rotates the input, applies `ic_erasercol` and `ic_compress2`, combines results with `overlay`, and finishes with `swapxy`. These trees reflect the pixel-wise optimization objective, with type correctness enforced throughout by the DSL.

**221 Stage 2: LLM-based Code Repair**

222 While STGP enables type-safe exploration of the DSL space, we found that it often struggles to escape local optima. In this case, GP may converge to syntactically correct but semantically incorrect programs that partially match the outputs without solving the task completely, as shown in Figure 4. To address this and inspired by (Xu and ... 2023), we incorporate a second refinement phase using a frozen LLM, specifically GPT-4o-mini, to perform code repair. After GP completes 500 generations, if the best evolved program does not achieve the perfect fitness score, a repair loop is instantiated for up to  $R$  rounds. In each round, we constructed a prompt that includes the final program, the input

234 and output training examples, the program’s fitness score, 235 and the max score it has to achieve for the given task. This 236 repair continues until either perfect fitness is reached or the 237 repair rounds are exhausted. Let  $p^* \in \mathcal{P}$  be the best-evolved 238 candidate after  $G$  generations. If  $f(p^*) < f_{\max}$ , the system 239 enters a repair loop for up to  $R$  refinement rounds. At each 240 round, a natural language prompt is constructed using three 241 as follows.

- 242 • The full set of input-output training pairs
  - 243 • The symbolic DSL source code of  $p^*$
  - 244 • The current program and max fitness score of the task
- 245 This prompt ( $\pi_r$ ) is submitted to the LLM, which returns a revised program  $p'_r \sim \text{LLM}(\pi_r)$  expressed in Python. 246 The revised program is executed and evaluated. If it yields a 247 higher fitness score than the previous best, it is accepted.

$$f(p'_r) > f(p^*) \Rightarrow p^* \leftarrow p'_r$$

249 This procedure stops either when a perfect solution is found 250 or when the maximum number of rounds  $R$  is reached. For 251 the experiment result shown in Figure 5 and 6, the repair 252 loop used is set to five ( $R = 5$ ). Code edit round ( $R = 0$ , 253  $R = 5$ ) refers to the number of times that the LLM can edit 254 its own code. As shown in Figure 2, there is a loop letting the 255 LLM edit its code, which is  $R$ . For our experiment setting, 256 we compare the result when  $R = 0$  and  $R = 5$ .

**257 Case Study: Success and Failure Analysis**

258 This case study serves as a concrete example illustrating 259 both the overall functioning of our two-stage framework 260 and the complementary roles of GP and LLM-based repair. 261 The ARC task shown in Figure 4 is from the concept class 262 with task number `Center7`. As shown in the Figure 4, 263 the first stage processes the task using the GP synthesis 264 process, resulting in the synthesized program shown in step 265 of Figure 4. This program represents the final output of 266 the GP search process and uses the following primitives: 267 `creament_tuple`, `sign_tuple`, `mostcolor_objs`, 268 `interval_int`, `mostcolor_objs`, `subtract_int`, 269 `add_int` with terminals 4, 2, and 5. While the program is 270 syntactically well-formed and type-safe, it fails to capture 271 the key semantic operation: detecting and reasoning over 272 object-level color content in the input grid. As a result, 273 the GP solution is semantically misaligned with the task 274 requirements scoring fitness score of 0. Step 3 LLM Code 275 Repair shows the Stage 2 of our framework, where the 276 LLM-based code repair module is invoked. This module 277 receives the best GP program, along with the input-output 278 grid pairs, the fitness score, and a prompt requesting code 279 edits. The program shown in stage 3 of Figure 4 is the edited 280 code which raises the fitness score from 0 (fitness score 281 by GP) to 0.7(fitness score of the repaired program). This 282 shows how the LLM fixes high-level semantic errors that 283 GP alone could not resolve and help it achieve high fitness 284 score than GP alone.

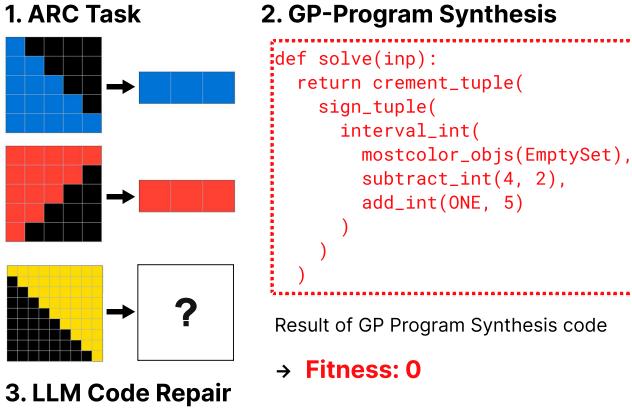


Figure 4: Case study showing GP program failure and LLM repair on ConceptARC task id: Center7 (Moskvichev, Odouard, and Mitchell 2023). The flow of the process starts with (1) illustrating the input-output training examples (2) Stage 1 applies GP to evolve a program over a symbolic DSL. The resulting expression is syntactically valid but semantically incorrect, producing an invalid output grid. (3) Stage 2 invokes a frozen LLM to repair the GP-generated program by editing the function directly based on example pairs and fitness. The new function the LLM edited takes a grid as input and returns the grid as an output, which matches the DSL type system defined. This example demonstrates how LLM-based repair can resolve semantic failures that GP alone cannot.

## Experiments

### Experiment Setting

This section starts with the experimental settings and the benchmark problem, and visualizes the evaluation result of our framework on the 400 ARC training data set and 160 ConceptARC dataset (Chollet 2019; Moskvichev, Odouard, and Mitchell 2023), which include tasks requiring abstraction, pattern matching, and symbolic reasoning. Each task, as shown in Figure 1, consists of input-output grid pairs, and the objective is to synthesize a function `solve(inp)` that generalizes to unseen grids.

For evaluation, we adopt two widely used metrics from the program synthesis literature. First, pixel-wise accuracy which measures the fraction of correctly predicted pixels across all test grids (Bober-Irizar and Banerjee 2024). Second, we report Pass@ $k$ , a probabilistic metric introduced by (Chen et al. 2021), which estimates the likelihood that at least one of the top- $k$  synthesized programs produces a fully correct output on all test cases. We report Pass@ $k$  for evaluating both the GP baseline and the full framework that

Table 1: Performance comparison between our approach and previous systems. Results for prior work are reported as presented in (Bober-Irizar and Banerjee 2024). All systems were assessed using the official ARC evaluation procedure, allowing up to three attempts per task, and a task is considered solved if it is completed correctly.

System	Solved Tasks
ARGA (Xu, Khalil, and Sanner 2023)	49 / 400
Dreamcoder (Alford 2021)	23 / 400
Dreamcoder-DC (Bober-Irizar and Banerjee 2024)	70 / 400
GP Program Synthesis (Stage 1 of our work)	65 / 400

305 incorporates LLM-based repair.

306 The terminal set for the GP includes integer color values  
307 (ranging from 0 to 9), boolean values (True, False),  
308 and an Input grid (ARG0). The primitive set includes DSLs  
309 from (Hodel 2024).

310 Each synthesis attempt runs the GP engine with a popula-  
311 tion of 1,000 individuals (population size remains the same  
312 throughout the generation) for up to 500 generations with  
313 crossover and mutation probability of 0.5 and max depth  
314 of tree set to 80. All results are averaged over five random  
315 seeds. If no perfect solution is found via GP alone, the LLM-  
316 based repair stage is invoked.

### Result Analysis

317 As shown in Table 1 our approach solves 65 out of 400  
318 tasks, outperforming ARGA solving 49 tasks and Dream-  
319 Coder (Xu, Khalil, and Sanner 2023; Alford 2021). Prior  
320 results are taken directly taken as reported in (Xu, Khalil, and  
321 Sanner 2023). ARGA uses a hand-engineered DSL based  
322 on graph abstractions, which represents solutions as a se-  
323 quence of symbolic graph transformations to solve 49 tasks.  
324 DreamCoder (Alford 2021) solves 23 tasks by representing  
325 programs using functional language and synthesized using  
326 enumeration. Its improved variant (Bober-Irizar and Baner-  
327 jee 2024) represents programs as symbolic functional pro-  
328 grams by incorporating learned priors and crafting DSLs  
329 for their system, achieving 70 tasks. Our system uses ge-  
330 netic programming over the DSL designed by Hodel, with-  
331 out introducing any new or custom DSLs. It solves 65 tasks,  
332 demonstrating that general-purpose evolutionary search can  
333 be competitive with more domain-specialized symbolic ap-  
334 proaches. We do not compare our method to recent LLM-  
335 based techniques that rely heavily on fine-tuning with large  
336 synthetic datasets, as our focus is on zero-shot or low-  
337 resource synthesis without synthetic data augmentation.

338 The results in Figure 5 show the change in program fitness  
339 after applying LLM code edits on top of the GP-only base-  
340 line. Figure 5a shows the improvement across 50 tasks from  
341 the ConceptARC benchmark, with an average fitness gain  
342 of +6.54%. Figure 5b presents results on 112 tasks sampled  
343 from the 400 ARC training tasks, yielding an average gain  
344 of +6.54% in fitness. In both Figure 5a) and 5b), each bar  
345 represents a task and shows the fitness difference (improve-  
346 ment) called ( $\Delta$ ) after the LLM edits, with tasks sorted by

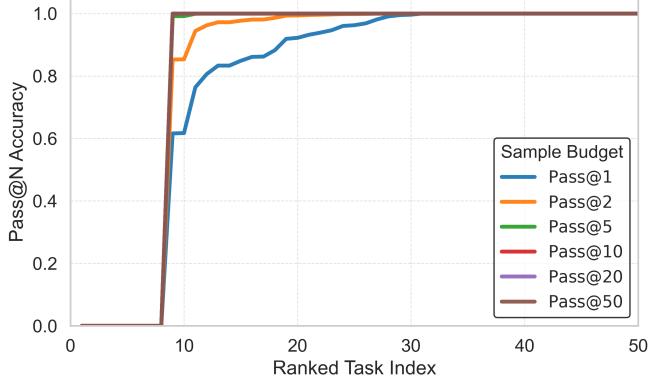


(a) 50 ConceptARC Tasks

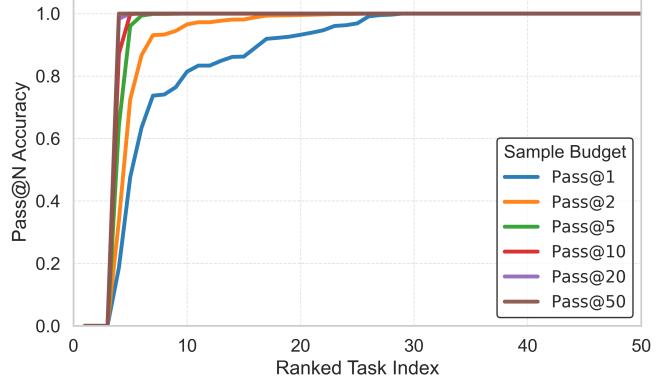


(b) 112 Tasks from ARC Training Set

Figure 5: Results show the improvement in fitness score ( $\Delta$ ) per task after LLM-based code edits for two different datasets. On average, the fitness score improved by +6.76% and +6.54% in each dataset. Tasks are sorted by fitness improvement in both.



(a) GP-only Program Synthesis



(b) GP with LLM Code Editing

Figure 6: Simulated pixel-level Pass@ $k$  accuracy curves across 50 tasks. Each line estimates the probability that at least one of  $k$  samples achieves  $\geq 90\%$  pixel-level similarity. Compared to GP-only, LLM code edits yield higher early-stage success probabilities (Pass@1, Pass@2), showing faster convergence and improved initial solution quality.

348 ( $\Delta$ ) value. These plots demonstrate that LLM code editing  
349 increases the fitness score across the tasks.

350 Figure 6 shows pixel-level Pass@ $k$  accuracy for 50 ARC  
351 tasks using (a) GP-only synthesis and (b) GP with LLM  
352 code editing. For each task, the fitness score achieved by  
353 the method is normalized by the maximum possible fitness  
354 for that task, resulting in a ratio between 0 and 1. This  
355 ratio is then treated as an estimator of how likely a gener-  
356 ated program is to reach 90% pixel-level accuracy. Using  
357 this estimated per-sample probability, the Pass@ $k$  accuracy  
358 is computed as the probability that at least one of the  $K$   
359 independent samples would reach that threshold. The key take-  
360 away is that LLM code editing with GP improves estimated  
361 success probabilities in the early samples, such as Pass@1  
362 and Pass@2, indicating that LLM code editing improves the  
363 quality of candidate programs and accelerates convergence  
364 toward high-accuracy outputs with limited samples.

365 Figure 7 presents a comparative analysis of the per-  
366 formance of the baseline GP, using LLM ( $R = 0$ ), LLM  
367 ( $R = 5$ ), and GP with LLM repair loop ( $R = 5$ ). The eval-  
368 uation covers 160 tasks from the ConceptARC (Moskvichev,  
369 Odouard, and Mitchell 2023), where each task belongs

370 to one of 16 high-level concepts listed. Each high-level  
371 concept on the x-axis represents a group of 10 tasks  
372 that share a common abstract high-level reasoning. The  
373 y-axis reflects the average pixel-wise performance within  
374 each concept, averaged over its corresponding tasks.  
375 The results show that the GP with LLM code edits  
376 outperforms the baseline GP model in the majority of  
377 high-level concepts. Specifically, substantial improvements  
378 are observed in concepts such as ExtractObjects,  
379 Count, FilledNotFilled, MoveToBoundary,  
380 ExtendToBoundary, and TopBottom. The LLM with  
381 repair rounds consistently outperforms both the GP-only  
382 and LLM zero-shot baselines. In contrast, the differences  
383 are less pronounced in concepts like CleanUp and  
384 CompleteShape, where both models perform similarly,  
385 showing that the baseline GP model is already sufficiently  
386 effective at handling these tasks, and LLM-based edits  
387 offer limited benefit. The LLM with code repair round  
388 of 0 performs worse in most of the categories, showing  
389 the importance of the initial GP search in constraining  
390 the solution space before invoking the LLM refinement.  
391 These findings support the view that LLMs can complement

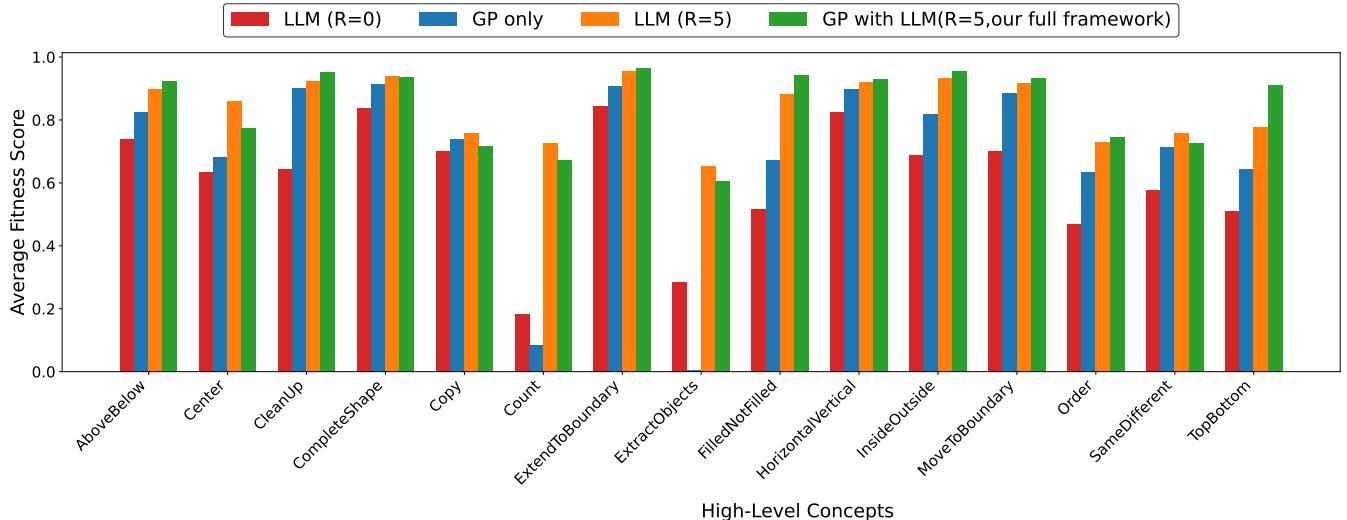


Figure 7: Performance comparison of LLM-only without repair round ( $R = 0$ ), GP-only, LLM-only with five repair rounds ( $R = 5$ ), and GP with LLM code edit with five repair rounds ( $R = 5$ , our full framework) on the ConceptARC over the high-level concept categorization.

evolutionary program synthesis by injecting semantic refinements into existing code. Moreover, as shown in the results, increasing the number of LLM refinement rounds from 1 to 5 improves performance. Integrating LLM-driven edits into evolutionary loops appears to be a promising strategy to enhance evolutionary program synthesis. Out of 160 tasks of ConceptARC dataset our full framework solve 25 tasks.

stay strictly within the space. In some cases, such as shown in Figure 4, the LLM synthesis solution uses Python operations, discarding the original DSL expressions. This case will be interpreted as the LLM dynamically generating a new DSL-like function, one that still maps input grids into output grids in a consistent and composable way, since we have other DSLs in our DSL pool that map input grids to output grids—this function made by the LLM also has the same functionality. From this perspective, the LLM is expanding the search space by injecting novel program representations that can be treated as valid individuals. This opens future works in which LLM-generated programs are not only fall-backs but can be reintegrated into the GP population for further evolution and recombination.

## Discussion

### Why is GP used as the base rather than starting from LLM-generated candidates or neural guidance?

We chose not to start from LLM-generated programs because our primary goal was to explore how GP, well-suited for multi-step transformation and planning tasks, can be enhanced by integrating with foundation models. Our aim was not to optimize ARC benchmark performance but to use ARC as a test bed for studying how evolutionary methods can be made more flexible and effective by leveraging LLMs. We are inspired by GP’s strength in structured search spaces and sought to show that these classic methods can benefit from LLMs without being replaced. Although research applying GP to ARC is limited, we intend to show that evolutionary processes can complement foundation models. In this way, our work builds on and reinforces prior directions in evolutionary learning by showing how such methods can remain open-ended and adaptive in the LLM era.

### What happens when the LLM produces a solution that no longer uses the existing DSL—can such code still be considered part of the symbolic search space, and how does this fit into the overall evolutionary process?

While our framework prompts the LLM to refine the symbolic programs produced by the genetic programming using the predefined DSL, we do not constrain the model to

## Conclusion

We propose a neuro-evolutionary framework that combines STGP with LLMs to synthesize programs for ARC tasks. Our approach evolves symbolic programs over a DSL, and when GP alone is insufficient, it invokes an LLM to refine them through semantic code edits. Experiments on the ARC and ConceptARC benchmarks show that LLM-based repair improves upon GP-only baselines, enabling synthesis on tasks otherwise unreachable by GP. While our framework currently treats LLM-edited programs as a terminal repair step, our findings suggest that these programs, which satisfy the type constraints of the DSL primitives, could also be reintegrated into the evolutionary population. This opens a promising direction where LLM-generated code not only patches failures but expands the search space, enabling open-ended and adaptive synthesis.

## References

- 456
- 457 Alford, S. 2021. *A neurosymbolic approach to abstraction and reasoning*. Ph.D. thesis, Massachusetts Institute of  
458 Technology.
- 460 Berman, J. 2024. How I Got a Record 53.6%  
461 on ARC-AGI-Pub Using Claude 3.5 and Evolutionary  
462 Test-Time Compute. <https://jeremyberman.substack.com/p/how-i-got-a-record-536-on-arc-agi>. Blog post, ac-  
464 cessed 17 Jul 2025.
- 465 Bober-Irizar, M.; and Banerjee, S. 2024. Neural networks  
466 for abstraction and reasoning: Towards broad generalization  
467 in machines. *arXiv preprint arXiv:2402.03507*.
- 468 Chen, M.; Tworek, J.; Jun, H.; and et al. 2021. Evaluating  
469 Large Language Models Trained on Code. *arXiv preprint*  
470 *arXiv:2107.03374*.
- 471 Chollet, F. 2019. On the measure of intelligence. *arXiv*  
472 *preprint arXiv:1911.01547*.
- 473 Chollet, F.; Knoop, M.; Kamradt, G.; and Landers, B. 2024.  
474 ARC Prize 2024: Technical Report. Technical Report  
475 arXiv:2412.04604, ARC Prize Foundation.
- 476 Coombe, C.; Howard, D. G.; and Browne, W. N. 2024.  
477 Learning Classifier Systems as a Solver for the Abstrac-  
478 tion and Reasoning Corpus. In *Companion Proceedings*  
479 of the Genetic and Evolutionary Computation Conference  
480 (GECCO '24).
- 481 Fischer, R.; Jakobs, M.; Mücke, S.; and Morik, K. 2020.  
482 Solving Abstract Reasoning Tasks with Grammatical Evo-  
483 lution. In *LWDA 2020 – Proceedings of the Workshop Day*  
484 “Lernen, Wissen, Daten, Analysen”, volume 2738 of *CEUR*  
485 *Workshop Proceedings*, 65–76.
- 486 Fortin, F.-A.; De Rainville, F.-M.; Gardner, M.-A.; Parizeau,  
487 M.; and Gagné, C. 2012. DEAP: Evolutionary Algorithms  
488 Made Easy. *Journal of Machine Learning Research*, 13:  
489 2171–2175.
- 490 Goldberg, D. E. 1989. *Genetic Algorithms in Search, Opti-  
491 mization and Machine Learning*. Reading, MA: Addison-  
492 Wesley.
- 493 Hodel, M. 2024. arc-dsl: A Domain-Specific Language for  
494 ARC. <https://github.com/michaelhodel/arc-dsl>. Accessed:  
495 2025-07-25.
- 496 Koza, J. R. 1994a. Genetic Programming as a Means for  
497 Programming Computers by Natural Selection. *Statistics*  
498 and Computing
- 499 Koza, J. R. 1994b. *Genetic Programming: On the Program-  
500 ming of Computers by Means of Natural Selection*. Cam-  
501 bridge, MA: MIT Press.
- 502 Montana, D. J. 1995. Strongly typed genetic programming.  
503 *Evolutionary computation*, 3(2): 199–230.
- 504 Moskvichev, A.; Odouard, V. V.; and Mitchell, M.  
505 2023. The conceptarc benchmark: Evaluating understand-  
506 ing and generalization in the arc domain. *arXiv preprint*  
507 *arXiv:2305.07141*.
- 508 Muslea, I. 1997. SINERGY: A linear planner based on ge-  
509 netic programming. In *European Conference on Planning*,  
510 312–324. Springer.
- 511 O’Neill, M.; and Spector, L. 2020. Automatic programming:  
512 The open issue? *Genetic Programming and Evolvable Ma-  
513 chines*, 21(1): 251–262.
- 514 Pantridge, E.; and Helmuth, T. 2023. Solving novel program  
515 synthesis problems with genetic programming using para-  
516 metric polymorphism. In *Proceedings of the Genetic and*  
517 *Evolutionary Computation Conference*, 1175–1183.
- 518 Parisotto, E.; Mohamed, A.-r.; Singh, R.; Li, L.; Zhou, D.;  
519 and Kohli, P. 2016. Neuro-symbolic program synthesis.  
520 *arXiv preprint arXiv:1611.01855*.
- 521 Pourcel, J.; Colas, C.; and Oudeyer, P. 2025. Self-Improving  
522 Language Models for Evolutionary Program Synthesis: A  
523 Case Study on ARC-AGI. In *Proceedings of the 42nd Inter-  
524 national Conference on Machine Learning (ICML)*. Poster.
- 525 Shi, K.; Hong, J.; Deng, Y.; Yin, P.; Zaheer, M.; and Sut-  
526 ton, C. 2023. Exedec: Execution decomposition for compo-  
527 sitional generalization in neural program synthesis. *arXiv*  
528 *preprint arXiv:2307.13883*.
- 529 Sobania, D.; Schweim, D.; and Rothlauf, F. 2021. Recent  
530 developments in program synthesis with evolutionary algo-  
531 rithms. *arXiv preprint arXiv:2108.12227*.
- 532 Tang, H.; Hu, K.; Zhou, J.; Zhong, S. C.; Zheng, W.-L.;  
533 Si, X.; and Ellis, K. 2024. Code repair with llms gives an  
534 exploration-exploitation tradeoff. *Advances in Neural Infor-  
535 mation Processing Systems*, 37: 117954–117996.
- 536 Van Deursen, A.; Klint, P.; and Visser, J. 2000. Domain-  
537 specific languages: An annotated bibliography. *ACM Sig-  
538 plan Notices*, 35(6): 26–36.
- 539 Wei, Y.; Xia, C. S.; and Zhang, L. 2023. Copiloting the  
540 copilots: Fusing large language models with completion en-  
541 gines for automated program repair. In *Proceedings of the*  
542 *31st ACM Joint European Software Engineering Conference*  
543 *and Symposium on the Foundations of Software Engineer-  
544 ing*, 172–184.
- 545 Xu, W.; and ... 2023. Self-Refine: Iterative Refinement with  
546 Self-Feedback. In *NeurIPS*.
- 547 Xu, Y.; Khalil, E. B.; and Sanner, S. 2023. Graphs, con-  
548 straints, and search for the abstraction and reasoning corpus.  
549 In *Proceedings of the AAAI Conference on Artificial Intelli-  
550 gence*, volume 37, 4115–4122.

551

## Reproducibility Checklist

### 552 Instructions for Authors:

553 This document outlines key aspects for assessing repro-  
 554 ductibility. Please provide your input by editing this `.tex`  
 555 file directly.

556 For each question (that applies), replace the “Type your  
 557 response here” text with your answer.

### 558 Example: If a question appears as

```
559 \question{Proofs of all novel claims  

  560   are included} { (yes/partial/no) }  

  Type your response here
```

561 you would change it to:

```
562 \question{Proofs of all novel claims  

  563   are included} { (yes/partial/no) }  

  yes
```

564 Please make sure to:

- 565 • Replace ONLY the “Type your response here” text and  
 566 nothing else.
- 567 • Use one of the options listed for that question (e.g., **yes**,  
 568 **no**, **partial**, or **NA**).
- 569 • **Not** modify any other part of the `\question` com-  
 mand or any other lines in this document.

570 You can `\input` this `.tex` file right before  
 571 `\end{document}` of your main file or compile it as  
 572 a stand-alone document. Check the instructions on your  
 573 conference’s website to see if you will be asked to provide  
 574 this checklist with your paper or separately.

### 575 1. General Paper Structure

- 576 1.1. Includes a conceptual outline and/or pseudocode de-  
 577 scription of AI methods introduced (yes/partial/no/NA)  
 578 **yes**
- 579 1.2. Clearly delineates statements that are opinions, hypoth-  
 580 esis, and speculation from objective facts and results  
 581 (yes/no) **yes**
- 582 1.3. Provides well-marked pedagogical references for less-  
 583 familiar readers to gain background necessary to repli-  
 584 cate the paper (yes/no) **yes**

### 585 2. Theoretical Contributions

- 586 2.1. Does this paper make theoretical contributions?  
 587 (yes/no) **no**
- 588 If yes, please address the following points:
- 589 2.2. All assumptions and restrictions are stated clearly  
 590 and formally (yes/partial/no)
- 591 2.3. All novel claims are stated formally (e.g., in theorem  
 592 statements) (yes/partial/no)

593 2.4. Proofs of all novel claims are included  
 594 (yes/partial/no)

595 2.5. Proof sketches or intuitions are given for complex  
 596 and/or novel results (yes/partial/no)

597 2.6. Appropriate citations to theoretical tools used are  
 598 given (yes/partial/no)

599 2.7. All theoretical claims are demonstrated empirically  
 600 to hold (yes/partial/no/NA)

601 2.8. All experimental code used to eliminate or disprove  
 602 claims is included (yes/no/NA)

### 603 3. Dataset Usage

- 604 3.1. Does this paper rely on one or more datasets? (yes/no)  
 605 **yes**
- 606 If yes, please address the following points:
  - 607 3.2. A motivation is given for why the experi-  
 608 ments are conducted on the selected datasets  
 609 (yes/partial/no/NA) **yes**
  - 610 3.3. All novel datasets introduced in this paper are in-  
 611 cluded in a data appendix (yes/partial/no/NA) **NA**
  - 612 3.4. All novel datasets introduced in this paper will be  
 613 made publicly available upon publication of the pa-  
 614 per with a license that allows free usage for research  
 615 purposes (yes/partial/no/NA) **NA**
  - 616 3.5. All datasets drawn from the existing literature (po-  
 617 tentially including authors’ own previously pub-  
 618 lished work) are accompanied by appropriate cita-  
 619 tions (yes/no/NA) **yes**
  - 620 3.6. All datasets drawn from the existing litera-  
 621 ture (potentially including authors’ own pre-  
 622 viously published work) are publicly available  
 623 (yes/partial/no/NA) **yes**
  - 624 3.7. All datasets that are not publicly available are de-  
 625 scribed in detail, with explanation why publicly  
 626 available alternatives are not scientifically satisfying  
 627 (yes/partial/no/NA) **NA**

### 628 4. Computational Experiments

- 629 4.1. Does this paper include computational experiments?  
 630 (yes/no) **yes**
- 631 If yes, please address the following points:
  - 632 4.2. This paper states the number and range of values  
 633 tried per (hyper-) parameter during development of  
 634 the paper, along with the criterion used for selecting  
 635 the final parameter setting (yes/partial/no/NA) **yes**
  - 636 4.3. Any code required for pre-processing data is in-  
 637 cluded in the appendix (yes/partial/no) **partial**

- 638 4.4. All source code required for conducting and analyzing  
639 the experiments is included in a code appendix  
640 (yes/partial/no) **yes**
- 641 4.5. All source code required for conducting and analyzing  
642 the experiments will be made publicly  
643 available upon publication of the paper with a license  
644 that allows free usage for research purposes  
645 (yes/partial/no) **yes**
- 646 4.6. All source code implementing new methods have  
647 comments detailing the implementation, with references  
648 to the paper where each step comes from  
649 (yes/partial/no) **yes**
- 650 4.7. If an algorithm depends on randomness, then the  
651 method used for setting seeds is described in  
652 a way sufficient to allow replication of results  
653 (yes/partial/no/NA) **no**
- 654 4.8. This paper specifies the computing infrastructure used for running experiments (hardware and  
655 software), including GPU/CPU models; amount  
656 of memory; operating system; names and versions of relevant software libraries and frameworks  
657 (yes/partial/no) **No**
- 660 4.9. This paper formally describes evaluation metrics  
661 used and explains the motivation for choosing these  
662 metrics (yes/partial/no) **yes**
- 663 4.10. This paper states the number of algorithm runs used  
664 to compute each reported result (yes/no) **yes**
- 665 4.11. Analysis of experiments goes beyond single-dimensional summaries of performance (e.g., average;  
666 median) to include measures of variation, confidence, or other distributional information (yes/no)  
667 **yes**
- 670 4.12. The significance of any improvement or decrease in  
671 performance is judged using appropriate statistical  
672 tests (e.g., Wilcoxon signed-rank) (yes/partial/no) **no**
- 673 4.13. This paper lists all final (hyper-)parameters used  
674 for each model/algorithm in the paper's experiments  
675 (yes/partial/no/NA) **yes**