# EVOLUTIONARY PROMPT OPTIMIZATION ENABLES EMERGENT MULTIMODAL REASONING STRATEGIES IN VISION-LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

#### ABSTRACT

We present a framework for optimizing prompts in vision-language models to elicit multimodal reasoning without model retraining. Using an evolutionary algorithm to guide prompt updates downstream of visual tasks, our approach improves upon baseline prompt-updating algorithms, which lack evolution-esque "survival of the fittest" iteration. Crucially, we find this approach enables the language model to independently discover progressive problem-solving techniques across several evolution generations. For example, the model reasons that to "break down" visually complex spatial tasks, making a tool call to a Python interpreter to perform tasks such as cropping, image segmentation, or saturation changes would improve performance significantly. Our experimentation shows that explicitly evoking this "tool calling" call, via system-level XML ... <tool>... </tool>... tags, can effectively flag Python interpreter access for the same language model to generate relevant programs, generating advanced multimodal functionality. This functionality can be crystallized into a system-level prompt that induces improved performance at inference time, and our experimentation suggests up to  $\approx 50\%$  relative improvement across select visual tasks. Downstream performance is trained and evaluated across subtasks from MathVista, M3CoT, and GeoBench-VLM datasets. Importantly, our approach shows that evolutionary prompt optimization guides language models towards self-reasoning discoveries, which results in improved zero-shot generalization across tasks.

031 032 033

034

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026

027

028

029

#### 1 INTRODUCTION

Vision-language models (VLMs) have advanced rapidly in their ability to jointly process images and text (Chameleon, 2024; Zhou et al., 2024; Zhang et al., 2024). Despite notable progress, many current approaches treat textual and visual inputs as loosely connected components (Chen et al., 2024b), underutilizing the potential for more integrated multimodal reasoning. Meanwhile, chain-of-thought prompting in text-only domains has shown that generating intermediate reasoning steps can substantially improve large language model (LLM) performance (Wei et al., 2023; Jin et al., 2024). Extending these ideas to the multimodal setting, however, often requires expensive retraining or specialized data (Zhang et al., 2024; Hao et al., 2024).

In this paper, we propose an *evolutionary prompt optimization* framework for vision-language models that operates purely at inference time, without finetuning model weights. Our method selfreferentially evolves a population of *system prompts* through iterative mutations and fitness selection. The best-performing prompts adaptively incorporate strategies for combining visual and textual information, including explicit calls to external code or processing modules. This multimodal, multi-model approach enables the VLM to leverage tool use at inference (e.g., Python-based image manipulation or segmentation (Kirillov et al., 2023)) for more detailed reasoning about images. Critically, the entire process unfolds at test time, permitting sophisticated behaviors to emerge from search over prompt space.

We demonstrate our approach on three benchmarks: MathVista (Lu et al., 2024), M3CoT (Chen et al., 2024a), and GeoBench-VLM (Danish et al., 2024), covering tasks ranging from spatial reasoning to complex counting. Results show that evolutionary prompt optimization can uncover nontrivial

problem decomposition, including subdividing an image into smaller regions or iteratively applying specialized code snippets. Such prompts substantially improve baseline performance, achieving up to  $\approx 50\%$  relative gains on certain subtasks. Moreover, these improvements require no supervised updates to the core model and can requires as few as 20 labeled examples per subtask to achieve generalizable improvements, making the approach widely applicable to any large vision-language model.

- Our contributions are the following:
  - We show that *evolutionary search* in natural language prompts can uncover multimodal reasoning strategies in VLMs without retraining.
  - We propose a purely *inference-time* framework that integrates an auxiliary interpreter model, enabling tool usage and dynamic problem decomposition to emerge naturally.
  - We provide extensive evaluations across multiple benchmarks, demonstrating substantial improvements over baseline approaches and illustrating the power of prompt evolution in driving advanced multimodal reasoning.
  - 2 RELATED WORKS
  - 2.1 CHAIN-OF-THOUGHT PROMPTING IN LANGUAGE AND VISION-LANGUAGE MODELS

075 Chain-of-Thought (CoT) prompting has emerged as a powerful paradigm for enabling Large Lan-076 guage Models (LLMs) to solve complex tasks by generating intermediate reasoning steps before 077 arriving at a final answer (Zhang et al., 2022; Wei et al., 2023). The success of CoT prompting has extended to Vision-Language Models (VLMs) in recent work. Zhang et al. (2024) highlight that 079 relying solely on brief annotations constrains the depth of multimodal reasoning. By distilling more comprehensive rationales from GPT-4 and incorporating reinforcement learning signals, they signif-081 icantly enhance the interpretability and robustness of VLM outputs. Similarly, Chen et al. (2024b) examine the consistency of VLM reasoning and propose methods for systematically quantifying and improving step-by-step visual grounding. While these advances have led to more transparent VLM 083 behaviors, they frequently rely on large-scale datasets or specific fine-tuning stages for reliable CoT 084 generation. 085

086 087

880

062

063

064

065

066 067

068

069

071

072 073

074

2.2 UNIFIED MULTIMODAL REASONING IN VISION-LANGUAGE MODELS

Recent studies have sought to broaden the scope of modern foundation models from reasoning 089 purely over text tokens to fully multimodal reasoning. Chameleon (Chameleon, 2024) and Transfu-090 sion (Zhou et al., 2024) enable multimodal reasoning by allowing transformers to natively generate 091 mixed-modal tokens. However, no works have combined such omnimodal models with advanced 092 methods in inference-time prompt optimization for multi-step reasoning. With the recent successes 093 in inference-time approaches for large models, such as self-consistency decoding (Wang et al., 2023), 094 there is growing interest in purely at-inference strategies that can guide VLMs to deeper analysis. 095 Additionally, Hao et al. (2024) shift from utilizing purely textual chains of thought to continuous 096 latent spaces, and Zhang et al. (2024); Chen et al. (2024b) leverage teacher-distilled rationales for 097 better VLM interpretability. These methods typically require specialized training or data collection, 098 while we focus on prompt optimization without model updates.

099 100

101

2.3 TOOL USAGE AND MULTI-MODEL INTERACTIONS

Our work also connects to a line of research that emphasizes *tool usage* and *multi-model interactions* for expanding a model's capabilities at inference time. For instance, Kirillov et al. (2023) introduce segment-anything modules that can be integrated with text-based pipelines but rely on carefully orchestrated external calls. In our approach, the evolutionary prompt optimization naturally yields prompts containing structured "tool calls," which are then parsed and executed by an auxiliary interpreter model. This multi-model synergy has been relatively underexplored for vision-language tasks, especially in the context of purely inference-time methods.

## 108 2.4 EVOLUTIONARY ALGORITHMS FOR PROMPT OPTIMIZATION

Parallel to these developments in vision-language reasoning, evolutionary algorithms have been increasingly employed to optimize prompts in LLMs. Guo et al. (2024) introduce EvoPrompt, demonstrating that discrete natural language prompts can be systematically evolved to enhance task accuracy. Jin et al. (2024) similarly leverage evolutionary strategies to refine zero-shot chain-of-thought
prompts, highlighting that diverse mutations can mitigate blind spots in static prompts. Further,
Fernando et al. (2023) propose Promptbreeder, where self-referential prompt mutation outperforms
standard CoT on arithmetic and commonsense benchmarks.

In addition to evolutionary algorithms, there has been extensive recent work on automated prompt 117 optimization that does not rely on EAs. For instance, RLPrompt (Deng et al., 2022) employs re-118 inforcement learning to optimize discrete text prompts; InstructZero (Chen et al., 2023) and Ad-119 versarial In-Context Learning (Do et al., 2024) adapt prompt instructions using black-box feedback 120 signals; INSTINCT (Lin et al., 2024) introduces a neural bandit for prompt refinement; and Teach 121 Better or Show Smarter? (Wan et al., 2024) explores how best to optimize instructions versus ex-122 emplars in prompting. These methods share the common goal of systematically refining prompting 123 strategies with minimal overhead. However, they remain mostly text-centric and do not directly ad-124 dress complex multimodal tasks, where visual grounding and incremental self-correction might be 125 necessary. 126

127 2.5 Self-Improving and Self-Referential Frameworks

A key theme in our work is *self-improvement* of prompts, where the system iteratively refines its 129 own instructions in order to better solve the task at hand. Approaches to self-referential optimization 130 date back at least to Schmidhuber (1993) and others (Schmidhuber, 1990; 1992), who introduced 131 early methods for neural networks that can adapt their own weights or configurations using internal 132 feedback loops. A central insight from these lines of work is that networks may exploit their own 133 internal representations to search for better control policies or parameter updates. Later, Gödel 134 Machines (Schmidhuber, 2003) formalized a self-referential agent that can rewrite aspects of its 135 own code upon proving the rewrite yields better expected utility. 136

Though these self-referential approaches were primarily concerned with model-level self-137 modification, the idea of iterative improvement without relying on external human supervision is 138 highly relevant to evolving system prompts for modern foundation models. In our problem setting, 139 we pursue an inference-time evolutionary scheme that effectively performs a lightweight form of 140 self-improvement in the prompt space. This resonates with recent interest in building foundation 141 models that continue to learn or refine themselves beyond their initial training data-sometimes us-142 ing imperfect internal or external verifiers. While our technique neither modifies the model weights 143 nor requires proofs of correctness, the broad notion of a system rewriting its own instructions to 144 enhance performance aligns with the self-referential tradition (Schmidhuber, 1993). As foundation 145 models scale, the role of self-improving methods—especially ones that can expand or refine behavior at test time-grows increasingly important for tasks that are not easily solved by static prompts 146 or purely supervised approaches. 147

- 149 3 METHODOLOGY
- 148 149 150 151

128

### 3.1 PROBLEM FORMULATION

We address the challenge of evolving task-specific prompts for vision-language models (VLMs) applied to downstream multimodal reasoning tasks through an evolutionary prompt optimization framework. Given a high-level multimodal cognitive task such as counting or classification, we seek to discover optimal system prompts that enhance model performance across diverse problem instances within that given task.

Let Q represent the complete set of question instances for a specific task. For any question  $q \in Q$ , we define a system prompt p from the prompt space  $\mathcal{P}$  that is prepended to the question. The concatenation operation  $\oplus : \mathcal{P} \times Q \to S$  maps a prompt-question pair to the final input string space S. Our objective function is then:  $p^* = \arg \max_{p \in \mathcal{P}} \mathbb{E}_{q \in Q}[\text{Score}(p \oplus q)]$  where Score :  $S \to [0, 100]$  evaluates the quality of the LLM's response. The score normalization to [0, 100] enables consistent comparison across different task types, with task-specific metrics (e.g., accuracy for classification, precision for counting) mapped to this standardized range.

We partition Q into training and test sets A and B, respectively. Our method relies on a fundamental assumption about the objective function's behavior across these sets:

$$\underset{p \in \mathcal{P}}{\arg\max} \mathbb{E}_{q \in A}[\operatorname{Score}(p \oplus q)] \approx \underset{p \in \mathcal{P}}{\arg\max} \mathbb{E}_{q \in B}[\operatorname{Score}(p \oplus q)]$$
(1)

This assumption allows us to optimize prompts on the training set with the expectation that they will generalize effectively to unseen test instances. For practical utility of our technique on downstream tasks, we operate in the regime where |B| >> |A|, with the assumption that the evolutionary framework can generalize past the train set and find the globally optimal task-specific prompt  $p^*$ . The practical validity of this assumption is demonstrated empirically in Section 5. The small training set serves as a "few-shot training set" for evolution, so that the large test set is indeed the real target.

#### 3.2 EVOLUTIONARY ALGORITHM DESIGN

Our evolutionary framework operates across three hierarchical spaces: the task prompt space  $\mathcal{P}$ , mutation prompt space  $\mathcal{M}$ , and meta-mutation prompt space  $\mathcal{H}$ . This hierarchical structure enables both direct optimization of task prompts and meta-learning of effective mutation strategies.

#### 3.2.1 POPULATION EVOLUTION

167 168

175

176 177

178

179

181

182 183

208

209

213

214

Algo	rithm 1 Binary Tournament Evolution	
1: <b>f</b>	or generation $g = 1$ to G do	
2:	Sample prompts $p_1, p_2 \sim \mathcal{P}_q$ without replacement	
3:	$p_w \leftarrow \arg \max_{p \in \{p_1, p_2\}} \operatorname{Fitness}(p)$	
4:	$p_l \leftarrow \arg\min_{p \in \{p_1, p_2\}} Fitness(p)$	
5:	$m \sim \mathcal{M}$	Sample mutation prompt
6:	$p'_w \leftarrow LLM(m \oplus p_w)$	▷ Mutate winner
7:	$\mathcal{P}_{g+1} \leftarrow (\mathcal{P}_g \setminus \{p_l\}) \cup \{p'_w\}$	
8: <b>e</b>	nd for	

Our framework uses a binary tournament genetic algorithm to evolve task prompts. The population evolves by replacing  $p_l$  with  $p'_w$ , maintaining size while improving fitness over G generations. The binary tournament balances exploration and exploitation, reduces computational overhead, and creates selection pressure towards better solutions. Mutation combines structured guidance from  $\mathcal{M}$ with the LLM's flexibility, enabling discovery of effective prompts that manual or fully automated approaches might miss.

## 200 3.2.2 MUTATION OPERATORS

Our framework employs a hierarchical system of mutation operators that combines both zero-order and first-order optimization strategies. The fundamental mutation process occurs in the prompt space  $\mathcal{P}$ , where each prompt  $p \in \mathcal{P}$  represents a strategy for solving a given task. These mutations are guided by prompts from the mutation space  $\mathcal{M}$  and hyper-mutation space  $\mathcal{H}$ .

We employ both first-order and zero-order prompt optimization techniques. For first-order optimization, we generate a new task prompt by applying the mutation prompt to the current prompt:

$$p' = \mu_1(p,m) = \mathcal{L}(m \oplus p) \tag{2}$$

where  $\mathcal{L}$  represents the language model's text generation function and  $\oplus$  denotes concatenation.

1

For zero-order optimization, we generate a new task prompt independently by concatenating the problem description D with a hint-generation template:

- $p'' = \mu_0(D) = \mathcal{L}(\text{"A list of 100 hints:"} \oplus D)$ (3)
- 215 This allows for the generation of novel task prompts that are closely tied to the original problem description, providing diversity in the evolutionary process.

The mutation prompt m itself evolves through both first-order and zero-order hyper-mutation operators. The first-order hyper-mutation operator is defined as:

$$m' = \nu_1(m, h) = \mathcal{L}(h \oplus m) \tag{4}$$

220 where  $h \in \mathcal{H}$  is a hyper-mutation prompt.

The zero-order hyper-mutation operator generates new mutation prompts by combining the problem description with a hint-generation template, similar to the zero-order mutation operator.

$$m'' = \nu_0(D, t) = \text{"A list of 100 hints:"} \oplus D \tag{5}$$

We adapt this paradigm of zero and first-order prompt optimization from Promptbreeder, and find that it generalizes well across vision-language tasks when initial prompt populations are visionspecific. This hierarchical system allows for both direct optimization of task prompts and adaptation of mutation strategies, while maintaining simplicity and interpretability in the evolutionary process. The combination of zero-order and first-order operators ensures both exploration of new ideas and refinement of existing solutions.

231 232

233

219

221

222

223

224 225

#### 3.3 FITNESS EVALUATION

The fitness function  $F : \mathcal{P} \to \mathbb{R}$  evaluates task prompts through a weighted combination of task performance and prompt quality metrics, defined as  $F(p) = (1 - \lambda)F_{\text{task}}(p) + \lambda F_{\text{aux}}(p)$ . The task fitness component  $F_{\text{task}}(p) = \frac{1}{k} \sum_{q \in C} \text{Score}(p \oplus q)$  measures empirical performance on a stochastically determined minibatch  $C \subset A$  of size k, where Score :  $S \to [0, 100]$  quantifies the quality of the LLM's response to task instance q when using prompt p. The score function is taskspecific-for example, in the case of a counting task, the score function is the percentage of correct answers.

241 The auxiliary fitness component  $F_{aux}(p) = \mathcal{L}_{critic}(critique \oplus p)$  employs an LLM-based critique 242 system that evaluates the sensibility and adherence of task prompts to their intended goals. This critique system acts as a regularizer for the evolutionary search process, steering the optimization 243 towards prompts that are not only effective but also semantically meaningful and aligned with the 244 task objectives. This aligns with previous literature that leverage LLMs' expansive knowledge base 245 for optimization tasks, even in settings of sparse reward (Yang et al., 2024). We find that critique 246 prompts that emphasize the coherence, explicitness, and adherence of mutated task prompts to stan-247 dard formatting perform the best, effectively guiding the evolutionary search process. By incorpo-248 rating this linguistic prior, we significantly improve sample efficiency, as demonstrated empirically 249 in Section 5, while maintaining the discovery of high-performing prompts. 250

The weighting coefficient  $\lambda$  balances the trade-off between empirical performance and prompt quality, with this value determined through ablation studies. We find that across subtasks, when we enforce  $F_{aux}(p) \rightarrow [0, 100]$ , that  $\lambda = 0.25$  performs well empirically. This balanced approach ensures that the evolutionary process discovers prompts that are both effective and interpretable, while the LLM-based critique system provides a computationally efficient mechanism for maintaining semantic coherence throughout the optimization process. This novel modification reduces the need for more complicated mutation mechanisms adopted by other works, such as Promptbreeder.

257

## 258 3.4 EVOLUTIONARILY EMERGENT TOOL SYNTHESIS259

A key discovery in our evolutionary framework is the emergence of self-referential tool generation capabilities. Through our robust and performant evolutionary search procedure, as well as highquality initial universes of mutation prompts and task prompts, we find that evolutionary search procedures for certain visual tasks yield task prompts that attempt to modify and re-ingest the input image(s) for multiple passes of reasoning. A successful example of this evolutionary reasoning is shown in Figure A.4, in contrast to an unsuccessful naive prompting example in Figure 2.

Rather than predefining a fixed tool universe, we observe that evolved system prompts naturally
 develop the ability to decompose problems into tool-like operations. We then leverage the natural
 capacity of LLMs to generate performant code from natural-language instructions by converting the
 natural language tool description into Python code with an auxillary LLM and executing it on the
 input image(s).

The evolutionary process operates on system prompts  $s \in S$  that guide the primary language model  $\mathcal{L}_1$  in processing inputs. Through mutation and selection pressure, these prompts evolve to incorporate structured tool suggestions enclosed in XML tags:

$$\mathcal{L}_1(s,x) \to (\dots < \texttt{tool} > \tau_i < /\texttt{tool} > \dots)_{i=1}^k \tag{6}$$

where each  $\tau_i$  represents a natural language description of a proposed tool operation. These tool suggestions emerge organically as the system discovers that breaking down complex tasks into composable operations improves performance. A secondary language model  $\mathcal{L}_2$  acts as an interpreter, translating each tool suggestion into executable Python code:

284

285

286

287 288

289

290

274

275

 $\mathcal{L}_2(\tau_i) \to c_i \in \mathcal{C} \tag{7}$ 

where C is the space of valid Python programs. This creates a flexible tool synthesis pipeline where  $\mathcal{L}_2$  leverages its code generation capabilities to implement operations like image manipulation, mathematical computations, or data processing based on natural language descriptions.

The composed transformation on input *x* becomes:

$$T(x) = \operatorname{eval}(c_k \circ \dots \circ c_1)(x) \tag{8}$$

where the composition emerges from the sequential application of synthesized tools. Critically, this 291 approach allows for open-ended tool discovery, because the system isn't constrained by predefined 292 tools. Additionally, this approach allows for recursive refinement, as the tool outputs can be fed back 293 into  $\mathcal{L}_1$  for iterative processing. This represents a novel reasoning paradigm for traditional vision-294 language models, as they can have multiple iterative reasoning passes at the same image, allowing 295 for more complex reasoning patterns such as examining different patches of the image multiple 296 times, increasing the brightness/contrast of patches, and applying external models such as Meta's 297 Segment Anything (SAM) tool Kirillov et al. (2023). Due to the expressivity of natural language 298 and the efficacy of LLMs in converting natural language instructions to executable code, tool usage 299 patterns become increasingly effective on downstream tasks with respect to generation count.

The emergence of structured tool suggestions in evolved prompts indicates that the system has discovered a fundamental principle: complex tasks often benefit from decomposition into smaller, well-defined operations. Crucially, this discovery happens naturally through the evolutionary process, as prompts that effectively utilize this pattern tend to produce better results across diverse inputs, resulting in a iteratively-optimized final prompt.

4 Results

308 Our experimental results demonstrate significant improvements across multiple vision-language 309 reasoning benchmarks through evolutionary prompt optimization. In all cases we benchmark results 310 on OpenAI's model 40 mini (OpenAI, 2023). Table 1 presents a comprehensive comparison of our 311 approach against several baselines, including the base model with no Chain-of-Thought prompting 312 (40 mini), standard Chain-of-Thought prompting (+CoT) using "Let's think step by step" as in 313 Wei et al. (2023), and PromptBreeder (+PB). Our evolutionary tool synthesis approach (+Tools) 314 achieves substantial gains across nearly all tasks, with particularly notable improvements in tasks 315 requiring complex spatial and physical reasoning.

316

306

307

The experimental results clearly demonstrate that our evolutionary prompt optimization framework markedly enhances multimodal reasoning in vision-language models. As evidenced in Table 1, while the standard chain-of-thought (CoT) prompting yield only incremental improvements, our evolved prompts (denoted as "+Ours") already push performance higher across tasks, and the addition of tool interpreter access ("++Tools") consistently achieves the best outcomes. For example, in the MathVista benchmark, performance on Visual QA improves from 49.5 with 40 mini to 53.3 with our approach, and further to an impressive 60.5 when tool usage is enabled. Similar trends are observed across other tasks—including Figure QA and Math Word Problems—indicating that the

Benchmark	Task	40 mini	+CoT	+PB	+Ours	++Tools
MathVista	Visual QA	49.5	51.0	49.6	53.3	60.5
	Figure QA	58.6	60.1	58.7	61.5	64.1
	Math Word Problem	61.8	63.2	61.9	64.5	68.0
МЗСоТ	Geometry	37.8	39.1	37.9	35.2	42.1
	Theory	6.1	9.0	6.2	6.3	-
	Physical Commonsense	42.6	43.9	42.7	47.2	61.7
GeoBench-VLM	Damaged Building Count	21.5	22.2	21.6	21.0	32.1
	Crop Type Classification	9.8	10.1	9.9	9.8	10.0
	Farm Pond Change Detection	12.3	12.7	12.4	14.1	20.2

Table 1: Performance across benchmarks showing the impact of different reasoning approaches. Best results for each task are bolded.

CoT = Chain of Thought, PB = PromptBreeder, Ours = Our method (with vision initial population) without tool interpreter access, +Tools = our method with tool interpreter access. A dash (–) in the +Tools column indicates that Tools were not elicited due to the nature of the subtask, so the performance is identical to the Ours column for that subtask.

evolutionary process not only refines prompt instructions for better task alignment but also facilitates
the spontaneous emergence of sophisticated strategies such as hierarchical problem decomposition
and dynamic tool synthesis. These findings underscore the potential of inference-time prompt evolution to unlock latent reasoning capabilities in vision-language models, thereby offering a scalable
and efficient pathway toward more robust multimodal AI systems.

#### 4.1 ANALYSIS OF EMERGENT BEHAVIORS

A particularly interesting finding is the emergence of sophisticated tool-use patterns through evolution. The evolved prompts frequently develop structured approaches to problem decomposition, often breaking complex tasks into sequences of simpler operations. For instance, in the Math Word Problem task, we observe prompts that systematically partition large images into manageable sections.



These behaviors emerged naturally through the evolutionary process, without explicit programming
or human demonstration. Similarly, in physical reasoning tasks, the evolved prompts often exhibit a
form of "mental simulation," breaking down complex physical scenarios into sequences of simpler
state transitions.

The results demonstrate that our evolutionary framework not only improves raw performance metrics but also discovers interpretable and generalizable reasoning strategies. The emergent behaviors often mirror human problem-solving approaches, suggesting that the framework is finding natural and effective solutions to complex reasoning tasks.

- 386 387
- 388
- 389

#### 4.2 GENERALIZATION EXPERIMENTS

We also measure the generalization capabilities of our 390 evolution framework compared to the base-line method 391 in multimodal reasoning domains. Due to cost con-392 straints induced by evaluating on larger datasets, we 393 run these generalization experiments on MolmoE-1B-394 0924 (Deitke et al., 2024) rather than 40 mini. Evolu-395 tionary prompt optimization only poses significant util-396 ity on downstream tasks if our assumption holds-that 397 the optimal system prompt generated through prompt optimization and evaluation on a train set generalizes 398 to perform near-optimally on a withheld test set. We 399 compare the baseline PromptBreeder (Fernando et al., 400 2023) method to our vision-language-specific approach 401 (both with tool usage enabled and without), and find 402 that both our approaches generate high-fitness system 403 prompts with just 20-30% of the total dataset, which is 404 often under 20 individual samples. We attribute this to 405 our improved LLM-augmented fitness function, which 406 serves as an effective regularizer to the search process 407 and yields higher sample efficiency. 408



Figure 1: Generalization performance of various prompt optimization techniques on Damaged Building Count vision-heavy reasoning task

#### 5 DISCUSSION

Our results highlight the strong potential of evolutionary prompt optimization for enhancing multi modal reasoning in vision-language models. This approach sheds light on artificial reasoning, the
 role of guided search in prompt space, and the future of multimodal AI.

414

409

410

415 416

#### 5.1 Self-referential Evolutionary Search as a Path to Advanced Reasoning

Notably, using the LLM as the mutation operator enables a "cognitive bootstrap," where its lin-417 guistic understanding refines prompts through mutation operators and its hyper mutation operators 418 in tandem evolve those mutation prompts in a self-referential way. That prompts exhibit human-419 like strategies (e.g., hierarchical decomposition) purely through the search process suggests that this 420 guided exploration can uncover interpretable, effective reasoning. The mixture of fitness scores and 421 auxiliary critic scores in the evolution process increases the sample efficiency of the mutation process 422 by ensuring that evolved prompts maintain coherence and align with the task objectives. This iter-423 ative loop itself fosters self-improvement at inference time, providing a mechanism for uncovering 424 emergent reasoning abilities that will only improve in performance as the auxiliary critic, mutation 425 and hypermutation models themselves improve. This yields a form of self-referential evolution: the 426 same family of models that ultimately needs improved instructions is generating mutations to those 427 instructions.

428 429

430

#### 5.2 SYSTEM PROMPTS AS LIGHTWEIGHT NEURAL PROGRAMS

431 System prompts can encode advanced computational strategies without altering model weights. This has three advantages: (1) it only requires inference-time computation, thus avoiding retraining costs;

(2) it can augment existing architectures without major modifications; and (3) it acts like a program, supporting explicit control flow, decomposition, and error handling. Natural language itself thus becomes a powerful medium to steer neural systems.

435 436

437

#### 5.3 THE CASE FOR NATIVE MULTIMODAL REASONING

438 The emergence and strong performance of image-based tool-calling behavior within the evolved 439 prompts motivates native multimodal reasoning as a future reasoning paradigm. Our method of 440 decomposing complex multimodal tasks into smaller image-based patches, performing text-based 441 reasoning on the subproblems, and then aggregating back up outperforms purely text-based rea-442 soning methods across several subdomains. Through evolutionary tool usage, the system prompts learn to decompose complex tasks into smaller patches and reason over these patches. This can be 443 seen as a primitive of reasoning natively over both text and image modalities flexibly, and hence 444 motivates the development of more advanced multimodal reasoning models. Additionally, allowing 445 the vision-language model to conduct multiple passes over the visual patches enables more robust 446 feature extraction and relationship understanding. The emergence of this multi-pass behavior in our 447 evolved prompts suggests that effective multimodal reasoning requires not just the ability to process 448 different modalities, but also the capability to dynamically revisit and reinterpret information as the 449 reasoning process unfolds. This finding has implications for architectural design choices in future 450 multimodal systems, particularly in how attention mechanisms and information flow are structured 451 across multiple reasoning steps. By incorporating separate tools and interpreters, our approach 452 shows how multi-model interactions can enable further self-improvement at test time, aligning with 453 broader frameworks that explore multi-agent or multi-module synergy.

454 455

456

- 5.4 SCALING GUIDED SEARCH TOWARD OMNIMODAL AI
- <sup>457</sup> Our findings suggest broad implications for future vision-language systems:

Guided Search as a Development Paradigm: Success here indicates guided search in prompt space can be a powerful strategy for emerging AI capabilities, especially where human intuition is limited (DeepSeek-AI, 2025).

Towards Omnimodal AI: Evolving multimodal reasoning strategies could ultimately yield systems
 integrating diverse modalities. Future work might explore multi-model or specialized modules for
 self-improving feedback loops.

Novel Test-Time Scaling Laws: Our approach achieves notable emergence through continual system prompt evolution, with emergent paradigms such as multimodal tool usage and dynamic programming emerging as we increase the amount of compute expended on evolution. This suggests a potential future inference-time scaling law in prompt / system program space, which may act as complementary to the recent advances in inference-time compute.

470

#### 5.5 LIMITATIONS AND FUTURE WORK

471 472

Though the evolutionary search effectively discovers sophisticated reasoning, it fails to improve on 473 some highly abstract tasks, suggesting limitations in handling theoretical domains. We tested only a 474 handful of benchmarks (MathVista, M3CoT, GeoBench-VLM) and one base model, so more diverse 475 tasks and architectures are needed to confirm generalization. Moreover, reliance on the same LLM 476 for mutations and auxiliary fitness introduces variability, as the whole evolution process can con-477 verge to prompts that please the model's own notion of correctness, rather than universal correctness. 478 We address this echo chamber risk by incorporating both objective fitness scores rather than only 479 self-referential prompt evaluation, but future work could explore alternative methods to mitigate this 480 potential issue. While our method is more efficient than retraining, the added inference overhead 481 may remain problematic in resource-constrained environments. Exploring reinforcement learning 482 or latent-space search strategies could further enhance multi-modal prompt optimization. Finally, 483 mechanisms behind emergent behaviors such as hierarchical decomposition remain underexplored, and combining few-shot optimization with continuous learning stands as a promising direction. 484 These areas may extend evolutionary prompt optimization to richer multi-model frameworks that 485 enable flexible, self-improving VLMs without human supervision or retraining.

## 486 REFERENCES

527

528

529

- Chameleon. Chameleon: Mixed-modal early-fusion foundation models, 2024. URL https://arxiv.org/abs/2405.09818.
- Lichang Chen, Jiuhai Chen, Tom Goldstein, Heng Huang, and Tianyi Zhou. Instructzero: Efficient instruction optimization for black-box large language models, 2023. URL https://arxiv. org/abs/2306.03082.
- 493
   494
   494
   495
   495
   496
   496
   497
   498
   498
   499
   499
   499
   490
   490
   490
   491
   492
   493
   493
   494
   495
   496
   496
   497
   498
   498
   498
   498
   499
   499
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
   490
- Yangyi Chen, Karan Sikka, Michael Cogswell, Heng Ji, and Ajay Divakaran. Measuring and improving chain-of-thought reasoning in vision-language models, 2024b. URL https://arxiv.org/abs/2309.04461.
- Muhammad Sohail Danish, Muhammad Akhtar Munir, Syed Roshaan Ali Shah, Kartik Kuckreja, Fahad Shahbaz Khan, Paolo Fraccaro, Alexandre Lacoste, and Salman Khan. Geobench-vlm: Benchmarking vision-language models for geospatial tasks, 2024. URL https://arxiv. org/abs/2411.19325.
- 505 DeepSeek-AI. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 506 2025. URL https://arxiv.org/abs/2501.12948.
- 507 Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Moham-508 madreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin 509 Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-510 Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne 511 Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Bor-512 chardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar 513 Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali 514 Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for state-of-515 the-art vision-language models, 2024. URL https://arxiv.org/abs/2409.17146. 516
- Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng
  Song, Eric Xing, and Zhiting Hu. RLPrompt: Optimizing discrete text prompts with reinforcement learning. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3369–3391,
  Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.222. URL https://aclanthology.org/2022.
  emnlp-main.222/.
- Xuan Long Do, Yiran Zhao, Hannah Brown, Yuxi Xie, James Xu Zhao, Nancy F. Chen, Kenji Kawaguchi, Michael Shieh, and Junxian He. Prompt optimization via adversarial in-context learning, 2024. URL https://arxiv.org/abs/2312.02614.
  - Chrisantha Fernando, Dylan Banarse, Henryk Michalewski, Simon Osindero, and Tim Rocktäschel. Promptbreeder: Self-referential self-improvement via prompt evolution, 2023. URL https: //arxiv.org/abs/2309.16797.
- Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview.net/forum?id=ZG3RaNIsO8.
- Shibo Hao, Sainbayar Sukhbaatar, DiJia Su, Xian Li, Zhiting Hu, Jason Weston, and Yuandong Tian. Training large language models to reason in a continuous latent space, 2024. URL https: //arxiv.org/abs/2412.06769.
- 539 Feihu Jin, Yifan Liu, and Ying Tan. Zero-shot chain-of-thought reasoning guided by evolutionary algorithms in large language models, 2024. URL https://arxiv.org/abs/2402.05376.

- Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Spencer Whitehead Tete Xiao, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything, 2023. URL https://arxiv.org/abs/2304.02643.
- Xiaoqiang Lin, Zhaoxuan Wu, Zhongxiang Dai, Wenyang Hu, Yao Shu, See-Kiong Ng, Patrick
   Jaillet, and Bryan Kian Hsiang Low. Use your instinct: Instruction optimization for llms using
   neural bandits coupled with transformers, 2024. URL https://arxiv.org/abs/2310.
   02905.
- Pan Lu, Hritik Bansal, Tony Xia, Jiacheng Liu, Chunyuan Li, Hannaneh Hajishirzi, Hao Cheng, Kai-Wei Chang, Michel Galley, and Jianfeng Gao. Mathvista: Evaluating mathematical reasoning of foundation models in visual contexts, 2024. URL https://arxiv.org/abs/2310.02255.
- OpenAI. Gpt-4vision system card. 2023. URL https://api.semanticscholar.org/ CorpusID:263218031.
  - J. Schmidhuber. Making the world differentiable: On using fully recurrent self-supervised neural networks for dynamic reinforcement learning and planning in non-stationary environments. Technical Report FKI-126-90, Technische Universität München, 1990.
  - J. Schmidhuber. Learning to control fast-weight memories: An alternative to dynamic recurrent networks. *Neural Computation*, 4(1):131–139, 1992. ISSN 0899-7667. doi: 10.1162/neco.1992. 4.1.131.
  - J. Schmidhuber. A 'self-referential' weight matrix. In *Proceedings of ICANN '93*, pp. 446–450, London, 1993. Springer. ISBN 978-1-4471-2063-6. doi: 10.1007/978-1-4471-2063-6\_107.
  - J. Schmidhuber. Gödel machines: Self-referential universal problem solvers making provably optimal self-improvements. https://arxiv.org/abs/cs/0309048, 2003. arXiv preprint cs/0309048.
  - Xingchen Wan, Ruoxi Sun, Hootan Nakhost, and Sercan O. Arik. Teach better or show smarter? on instructions and exemplars in automatic prompt optimization, 2024. URL https://arxiv.org/abs/2406.15708.
- 571 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh 572 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models,
   573 2023. URL https://arxiv.org/abs/2203.11171.
  - Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun
   Chen. Large language models as optimizers, 2024. URL https://arxiv.org/abs/2309.
   03409.
- Ruohong Zhang, Bowen Zhang, Yanghao Li, Haotian Zhang, Zhiqing Sun, Zhe Gan, Yinfei Yang, Ruoming Pang, and Yiming Yang. Improve vision language model chain-of-thought reasoning, 2024. URL https://arxiv.org/abs/2410.16198.
- Zhuosheng Zhang, Aston Zhang, Mu Li, and Alex Smola. Automatic chain of thought prompting in large language models, 2022. URL https://arxiv.org/abs/2210.03493.
- Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model, 2024. URL https://arxiv.org/abs/2408. 11039.
- 591

555

556

558

559

561

562

563 564

565

566

567

568

569

570

574

575

576

577

A	Appendix
A.1	Full Evolutionary Algorithm Outline
Algo	orithm 2 Evolutionary Prompt Optimization
1:	Input:
	• Training set A
	• Initial prompt population $\mathcal{P}_0$ of size N
	• Mutation prompt space $\mathcal{M}$ and hyper-mutation space $\mathcal{H}$
	• Fitness function $F(p) = (1 - \lambda)F_{tack}(p) + \lambda F_{my}(p)$
	• Maximum number of generations $C_{r}$
2.	<b>Output:</b> Optimized prompt $n^*$
3:	Initialization:
4:	Set generation counter $a \leftarrow 0$
5:	Initialize prompt population $\mathcal{P}_0$ with N candidate prompts
6:	while $g < G$ do
7:	Selection and Evaluation:
8:	Randomly sample two distinct prompts $p_1, p_2 \in \mathcal{P}_g$
9:	Compute fitness scores $F(p_1)$ and $F(p_2)$
10:	if $F(p_1) \ge F(p_2)$ then
11:	Set winner $p_w \leftarrow p_1$ and loser $p_l \leftarrow p_2$
12:	else
13:	Set winner $p_w \leftarrow p_2$ and loser $p_l \leftarrow p_1$
14:	ena n Mutation
15. 16.	Sample a mutation prompt $m \in M$
10. 17·	Generate mutated prompt: $n' \leftarrow f(m \oplus n_m)$
18:	Population Update:
19:	Update population: $\mathcal{P}_{a+1} \leftarrow (\mathcal{P}_a \setminus \{p_l\}) \cup \{p'_u\}$ Apply hyper-mutation operators using
	$h \in \mathcal{H}.$
20:	Increment generation: $q \leftarrow q + 1$
21:	end while
22:	<b>Return:</b> $p^* \leftarrow \arg \max_{p \in \mathcal{P}_G} F(p)$

#### A.2 NAIVE PROMPTING EXAMPLE - FAILURE

#### Direct Prompting



Figure 2: A naive example of directly prompting an input image using the original dataset prompt. Note that the LLM misses one of the metallic shapes, leading to an incorrect conclusion. Given the visual complexity of the input image, Vision Language Models may struggle to accurately analyze subcomponents without further guidance.

## 648 A.3 EVOLUTIONARY PROMPTING EXAMPLE - SUCCESS



Figure 3: Walk through of an example where an initial prompt fails to elicit a correct answer, while a successful evolutionarily optimized prompt including a tool call (cropping) succeeds. Via the evolved prompt, the model elicits a tool call that crops the original image, allowing the LLM to better ingest the image's contents. With the improved quadrant division of visual analysis, the model is able to correctly answer the question.

#### 702 A.4 PROMPT FITNESS AS A FUNCTION OF EVOLUTION TIME



Figure 4: Baseline evolutionary prompt optimization method (Promptbreeder, Fernando et. al. 2023)
fails to generalize to vision-language reasoning domains. We find this is because their instructionfollowing prompting for LLM mutation, hypermutation, and their initial universes are not suited for
vision-language reasoning tasks.





Figure 5: Our naive method outperforms baselines in the evolution process, due to significant improvements in mutation methods, our auxillary loss preventing significant and nonsensical deviations from current task prompts, and our initial universes of task prompts, mutation prompts, and hypermutation prompts, that are tuned specifically for image tasks.



777

756

778

779 Figure 6: With tool use enabled, and an improved set of mutation and hypermutation prompts that 780 encourages emergence of tool use, population fitness scales positively with time. The notable drops in average performance (red curve) are the critical windows in which tool usage emerges through 781 evolution. Initially, fitness falls, because the tool usage and reasoning paradigms are nascent, but 782 as they are evolved more, they become high-performing. Towards the end, the evolution process 783 guides the highest-performing system prompts towards another layer of tool calling. This emergent 784 strategy results in an immediate drop in performance, due to its nascence and incompleteness. We 785 hypothesize that continuing the evolution process further would lead to even more advanced reason-786 ing paradigms like these results indicate, but under our fixed computation budget, this remains an 787 avenue for future research.

788 789 790

791

798

800

#### A.5 AUXILIARY LLM CRITIC IMPLEMENTATION DETAILS

The auxiliary critic component ensures evolved prompts maintain coherence and stay aligned with task objectives. Implemented using GPT-4o-mini, the critic evaluates prompts through a multi-dimensional rubric designed to prioritize clarity, logical structure, and task relevance. We assign weights to different components of the given task prompt's quality, which are determined empirically. The critic operates via a structured evaluation template that emphasizes task fidelity:

```
797 Evaluate this prompt for:
```

```
799 Relevance to the stated visual reasoning task
```

801 Logical flow between instructions

```
802
803 Clarity of language and specificity
```

```
Score each dimension 1-5, then compute weighted total (0-100).
```

```
Flag any instructions that deviate from the task's core requirements.
```

806

```
807
```

808

## A.6 MUTATION PROMPTS AND TASK PROMPTS INITIAL UNIVERSES

#### 812 Table 2: Sample Vision-Language Starting Prompts 813 814 Initially Evolved Prompt (Vision-Task Aligned) **Reasoning Strategy** 815 816 "Generate a diagram highlighting the fundamental shapes and Focuses on style-invariant structures to key objects in the image. Use these as anchors to guide your final capture essential spatial and content 817 answer (such as a numeric value)." information. 818 819 "Translate the relevant visual features into symbolic or textual Balances interpretability and precision, 820 notations, aiming for both clarity and accuracy. Then refine this handling trade-offs between simplicity 821 representation to produce a final answer." and completeness. 822 "Segment the image according to the task requirements, focusing Applies task-specific segmentation to 823 on regions most relevant to the question. Prioritize these segments isolate key areas, reducing distraction 824 for deeper analysis." from less important regions. 825 826 "Iteratively refine your approach by generating bounding boxes Uses a population-based or iterative 827 or region proposals for the image. Retain only the proposals that mechanism to refine localized views of 828 significantly improve the clarity or correctness of your final rethe image. 829 sult." 830 "Simulate common edge cases or distortions (like occlusion and Incorporates robustness testing and itunusual lighting) to see how they affect your answer. Refine the erative prompt fixes for improved fault 831 prompt steps that cause ambiguous or incorrect responses." tolerance. 832 833 "Construct a hierarchical representation of objects in the image, Organizes local and global features in a 834 multi-scale structure for more holistic capturing relationships at multiple scales. Merge the partial find-835 ings into one cohesive conclusion." reasoning. 836 "Generate several possible answers for the question by varying Uses contrastive evaluation to identify 837 the approach. Compare how well each aligns with the visual dethe answer best supported by the evi-838 tails, and select the most fitting explanation." dence in the image. 839 840 "Apply a mix of symbolic and sub-symbolic steps to interpret the Combines rule-based reasoning with 841 image. For instance, if the question involves counting objects, learned representations for inter-842 express it in conditional form (IF more than X, THEN...). Evaluate pretable and flexible analysis. 843 which approach yields the clearest final result." 844 845

864	Table 3: Sample Mutator Prompts					
865						
866	Prompt					
867	"Downite the instance of the it for any instance of the instan					
868	"Rewrite the instruction so that it focuses on <b>breaking down</b> any complex parts into simpler steps. Include a helpful tip for someone struggling"					
869	into simpler steps. Include a netpjul up jor someone straggling.					
870	"Note there is likely a <b>critical error</b> in the last response. A corrected version would					
871	be:"					
872	"Parkrass the instruction as if you are quiding someone who does not have visual					
074	stimulus making sure every detail is crystal clear"					
074	sumans, making sure every delan is erystat creat.					
070	"Imagine you must <b>teach this instruction to a peer</b> who could be easily confused.					
070	Simplify it, but offer one surprising or creative example."					
070	"Imagine a shortcut for this task if you had infinite resources and canabilities					
870	SHORTCL/T="					
880						
881	"Flip the point of view: rewrite the instruction as if the user is already an expert,					
882	and you are simply double-checking their approach."					
883	"Break the instruction into <b>two different methods</b> —one for someone who learns					
884	best by doing, and another for someone who prefers planning."					
885						
886	"Encourage <b>outside-the-box thinking</b> : rephrase the instruction so it allows for an					
887	unconventional or imaginative angle, but keep it workable."					
888	"Create a more <b>visual-oriented version</b> of the instruction by prompting the user to					
889	sketch out key steps or components before proceeding."					
890						
891	"Rewrite the instruction in a <b>step-by-step checklist</b> format, then add a final insight					
892	or reminder that ensures the goal is met.					
893						
894						
895						
896						
897						
898						
899						
900						
901						
902						
903						
904						
905						
906						
907						
908						
909						
910						
311						
JIZ 012						
313						
914 015						
916						
917						