# SELF-DIRECTED DISCOVERY: HOW LLMS EXPLORE AND OPTIMIZE CONFIGURABLE LLM SOLUTIONS

**Anonymous authors** 

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

024

025 026 027

028 029

031

033

034

037

038

040

041

042

043

044

045

046

047

048

049

051

052

Paper under double-blind review

#### **ABSTRACT**

LLM solutions are increasingly replacing specialized machine learning models across various industry domains. While offering simplicity and maintainability advantages, optimizing these workflows remains heavily dependent on expertdriven experimentation. In this paper we test off-the-shelf powerful LLMs for the task of automatically building and iteratively improving LLM-based solutions. We propose a configuration-driven framework which defines such workflows by specifying model parameters, prompt templates, and data transformations. We show that this standardized representation enables automatic iterative improvement loops via optimization agents which are themselves defined within the same framework. When evaluated on challenging datasets, it discovers improved solutions while maintaining interpretability and human verifiability. The improvement loop we instantiate is generic and minimal (using prompts that have not been engineered) and self-referential (improved solutions are discovered with improved documentation and examples). The proposal is a self-improving system that bridges the gap between generic code-generating automatic optimization and more narrowly-focused techniques such as prompt engineering.

#### 1 Introduction

Applications across domains such as financial services, healthcare, customer support, content moderation, and many others, increasingly rely on LLM solutions where the same inference pattern is applied to thousands or millions of data points. Such unified LLM-based solutions are increasingly preferred over specialized ML solutions, offering the appeal of a single, versatile approach that can handle diverse tasks without requiring domain-specific model training and maintenance (Chkirbene et al., 2024; Raza et al., 2025; Saleh et al., 2025). While these LLM inference workflows may initially lack the efficiency of highly optimized, custom machine learning (ML) solutions developed over years, maintainable systems often outweigh the performance gains of complex, specialized alternatives. However, unlike conversational AI systems where each interaction is unique, these LLM inference scenarios involve repeatedly executing identical prompt templates with only the input data varying, creating additional opportunities for optimization, for both quality and other performance goals. While it is becoming commonplace to utilize techniques such as prompt caching (Xiao et al., 2023; Zhu et al., 2023), automatic prompt engineering ((Ramnath et al., 2025; Zhou et al., 2022; Li et al., 2025a; Debnath et al., 2025)) or model selection (Wang et al., 2023; Tanaka et al., 2023; Tang et al., 2024), optimizing repeated inference patterns holistically still heavily relies on expert-driven experimentation (Li et al., 2024).

On the other hand, generative models are also transforming ML *research* by acting as powerful assistants that automate routine tasks, augment human creativity, and accelerate the pace of discovery. While LLMs are not yet capable of fully replacing human ML scientists, they significantly enhance efficiency by automating coding, data analysis, hypothesis generation, and manuscript writing (Tang et al., 2025; Wang & et al., 2025; Yang et al., 2024b; Schramowski et al., 2025). In this paper we take a step further in the direction of using generative models to aid in the development of ML solutions by showing that LLMs can successfully drive the optimization of LLM inference workflows and help discover improved solutions. This proposal addresses the gap between generic code-generating automatic optimization and more narrowly-focused techniques such as prompt engineering.

The paper makes the following contributions:

- We address a specific class of ML solutions based on LLM inference and propose a configuration-driven framework that supports *interpretability*. We demonstrate that this constrained yet expressive solution space enables powerful off-the-shelf LLMs to autonomously create and improve workflows across various tasks with minimal task-specific guidance. The generated workflows execute without errors over 95% of the time, achieve task performance that matches standardized prompts on average, and consistently discover significantly better solutions.
- We further implement a simple automatic iterative improvement loop, where other LLMs improve
  the solutions for custom goals such efficiency or quality. We propose a very simple linear iterative
  improvement loop using an Analyzer and an Improver agent, and show it can find better solutions
  to various challenging data sets. Since the Analyzer and Improver are in themselves LLM-based,
  we define them within the same framework and can thus be the object of optimization as well.
- We prioritize documentation-driven optimization over prompt engineering, which requires extensive manual tuning for each task and optimization objective. Instead of engineered meta-prompts, our LLMs receive minimal generic prompts alongside rich contextual artifacts: framework documentation and workflow examples. This approach enables robust cross-task generalization and sets the stage for a self-evolving system where accumulated knowledge artifacts improve performance across diverse domains without task-specific prompt crafting.

#### 2 CONFIGURATION-DRIVEN LLM AGENT DEFINITION

In this section we test whether off-the-shelf LLMs can generate high-performing agents for new tasks given a well-defined agent specification framework and examples of agents implemented using it. To achieve this we propose a configuration-driven LLM agent framework which promotes interpretability (a human user can easily understand and verify the generated agent) and enables systematic optimization (another LLM can iteratively improve on these agents).

**LLM Agent** We define an LLM agent as a workflow that makes a *single* LLM inference call, optionally preceded by input pre-processing and followed by output processing. <sup>1</sup>. More precisely, let  $\mathcal{D} = \{(x_1, y_1), \ldots, (x_n, y_n)\}$  denote a dataset where each data point consists of input  $x_i$  and optional ground truth  $y_i$ . An LLM agent A is defined by:

- ullet A prompt template T containing fixed text and placeholder variables
- A set of input functions  $\{f_1, f_2, \dots, f_k\}$  that transform input data  $x_i$
- An LLM model  $(\theta_{\rm LLM})$  alongside inference hyper-parameters  $\rho$  (temperature, max tokens, etc.)

For input  $x_i$ , the agent generates a prompt  $p_i$  and then output  $o_i$  as:

$$p_i = T(f_1(x_i, c), \dots, f_k(x_i, c), c)$$
 and  $o_i = A(x_i) = \theta_{\text{LLM}}(p_i; \rho)$ 

where c stands for additional context (examples, external data sources, etc.). Dataset-level performance is computed as  $\bar{E} = \frac{1}{n} \sum_{i=1}^{n} E(o_i, x_i, y_i)$  where E is a task-appropriate evaluation function.

**Agent Implementation** At the level of implementation, we describe an LLM agent via three files: Configuration (json), Agent Class Implementation (Python) and Input Schema File (json). See Appendix B for more details. The *configuration* file specifies the agent's behavior: model settings, prompt template, input processing functions, and output format. The agent class implements the pre-processing functions referenced in the configuration. These functions handle data transformations, feature extraction, example retrieval, and other per-datapoint operations. The input schema file defines the agent's interface with other components, specifying what data the agent can access and how it communicates with upstream and downstream tasks. This three-file structure enforces clear separation between declarative configuration, implementation logic, and interface contracts, enabling systematic optimization while maintaining interpretability.

At run-time, agents process data points sequentially through a standardized pipeline: load and validate configuration files, apply input transformations, render the prompt template, call the LLM,

<sup>&</sup>lt;sup>1</sup>This definition differs from the broader usage of "agent" in the literature, which typically encompasses multi-step reasoning, tool usage, memory systems, and iterative planning. Our constrained definition focuses on the fundamental building block of LLM-based solutions.

and parse the response according to the specified output format. Each agent is constrained to exactly one LLM inference call—multi-step reasoning requiring multiple LLM interactions must be orchestrated at the workflow level using separate agents.

**Modifiability Control** To guide other LLMs in improving task agents, we introduce modifiability control. Each configuration section includes a flag  $modifiable \in true$ , false that helps meta-agents understand which sections can be modified. For example, model.modifiable = true allows modification of model parameters, while output.modifiable = false preserves the response schema. This enables systematic exploration while maintaining configurable constraints based on the specific improvement loop requirements—for example, preserving output schemas for downstream compatibility allowing flexible optimization across all components.

Appendix C shows a complete math problem solving agent example and Appendix D shows how this design supports diverse implementations including ML integration, in-context learning, and dynamic input transformation.

#### 2.1 LLM-BASED AGENT GENERATION

To evaluate the configuration-driven framework, we first examine whether LLMs can generate valid task agents  $A_t$  given only framework documentation and minimal task context. This evaluation aims to validate the framework's clarity and usability, as successful agent generation from documentation and minimal context indicates that the abstractions are sufficiently well-designed for automated solution synthesis.

#### 2.1.1 EXPERIMENTAL SETUP

We evaluate agent generation capabilities across the following problems:

**AIME2024 and AIME 2025** contain 30 competition-level mathematics problems each, requiring integer answers between 000-999 (MAA Committees). These challenging problems span algebra, geometry, number theory, and combinatorics, requiring complex mathematical reasoning that remains difficult for current (non-reasoning) models. We evaluate exact match accuracy and use AIME2024 Part I for development, with remaining datasets for testing.

**Math500:** 500 competition mathematics problems with detailed solutions across 7 subjects (Algebra, Counting & Probability, Geometry, Intermediate Algebra, Number Theory, Precalculus, Prealgebra) (Hendrycks et al., 2021; Lightman et al., 2024). For Math500, evaluation uses mathematical expression normalization to account for equivalent mathematical representations<sup>2</sup>. We sample 100 data points for development and leave the remaining 400 for test.

**MMLU-Pro** is a benchmark designed to challenge LLMs beyond the original MMLU (Wang et al., 2024). It features reasoning-focused multiple-choice questions across 14 domains with ten answer options instead of four, creating a more demanding evaluation. We use validation data for development and sample 200 test problems.

We test if LLMs can generate working agents for these tasks using agent framework documentation and minimal task information. Specifically we implement the agent generation using a single meta-agent defined in the same framework. Table 1 shows the prompt template used by this meta-agent, where the placeholders stand for:

- 1. (Framework Documentation) and (Agent Example): The complete LLM Agent README (in supplementary material, totaling 5k tokens) and a single complete agent example. For all tasks we use the basic AIME2024 agent shown in Appendix C. Note that this example is out-of-domain for MMLU-Pro.
- 2. (Task Data Point) and (Task description): Task data points are randomly drawn from a different data split, with the purpose of exemplifying the task to be solved and the input/output format. Task descriptions are 1-2 sentences describing the task objective (see Table 1).

 $<sup>^2</sup>We$  use the grader at https://github.com/openai/prm800k/blob/main/prm800k/grading/grader.py

Table 1: Prompt template used by the agent generating (meta-)agent, which generates task agents  $A_t$  using the task information in this table, alongside the  $\langle Framework \ Documentation \rangle$  and  $\langle Agent \ Example \rangle$  described above.

$\langle  ext{Task Data Point}  angle$	⟨Task Description⟩	Generator Prompt Template
Every morning Aya goes for	AIME (American Invitational	Instructions: You are an expert AI agent developer.
č	,	Task: (Task Description)
	1 2	Requirements: Generate exactly 3 files: (1) Agent
		configuration JSON, (2) Agent implementation
		Python, (3) Input schema JSON Schema. The
		agent must accept problem data as input, generate
	final numerical answer.	detailed reasoning, output correct format, handle
		parsing errors gracefully, and follow framework
		patterns.
E .		Framework Doc: (Framework Documentation)
		Agent Example: \( \text{Agent Example} \) Sample Data: \( \text{Task Data Point} \)
· ·		Output Format: JSON only with agent_config,
	MATIL muchlama magyina an	agent_code, input_schema fields.
		Generate the agent now:
		Generate the agent now:
Subject. Frecalculus, Level. 2		
<u> </u>		
Which of the following ren-		-
0 1		
	<b>j</b>	
	` ,	
, , ,	and select the correct answer	
Answer: B		
	Every morning Aya goes for a 9-kilometer-long walk and stops at a coffee shop afterwards. When [] Find the number of minutes the walk takes her, including the $t$ minutes spent in the coffee shop. Answer: 204 Solution: detailed solution Find the sum of all integer bases $b > 9$ for which $17_b$ is a divisor of $97_b$ . Answer: $07_b$ Simplify $0.00000000000000000000000000000000000$	Every morning Aya goes for a 9-kilometer-long walk and stops at a coffee shop afterwards. When [] Find the number of minutes the walk takes her, including the $t$ minutes spent in the coffee shop. Answer: 204 Solution: detailed solution Find the sum of all integer bases $b>9$ for which $17_b$ is a divisor of $97_b$ . Answer: $070$ Simplify $0.00$ tanger: $0.00$ Answer: $0.00$ and $0.00$ the final numerical answer.  Which of the following represents an accurate statement concerning arthropods? A. They possess an exoskeleton composed primarily of peptidoglycan. B. They possess an open circulatory []  B. They possess an open circulatory []  C. []   AIME (American Invitational Mathematics Examination) problems require integer answers between 000 and 999. The agent must show detailed reasoning and provide the final numerical answer.  MATH problems require integer answers in LaTeX format within \boxed tags. The agent must show detailed reasoning and provide the final answer in the correct format.  MMLU-Pro problems are multiple-choice questions across academic subjects (biology, math, physics, etc.) with 10 answer choices (A-J). The agent must provide detailed step-by-step reasoning and select the correct answer choice.

The LLM must infer an appropriate agent structure, input processing logic, prompt design, and output schema solely from the framework documentation, complete agent example and data sample; no additional examples of agent configurations, code samples, or task-specific guidance are provided. However in the future we see this as a self-evolving framework, where the amount of artifacts increases with every problem solved, creating better and richer context for the LLM or code assistant.

**Evaluation and baselines** In order to facilitate comparisons, we standardize Claude 3.5 Sonnet v2 for all task agents by removing other model references from the framework documentation. The meta-agent uses the same model with temperature 0.7 for exploration.

We run the meta-agent 10 times per task and perform static validation on generated configurations (structure, model names, required fields). Runtime issues include undefined placeholders (non-fatal) and syntax/execution errors (fatal). Success rate measures the percentage of agents executing without fatal errors on development data. For successful agents, we further evaluate task performance using multiple runs due to output variability: 5 runs for AIME datasets and 2 runs for others. We report accuracy averaged across all runs (Pass@1).

We compare against standard prompts for each task: basic CoT for AIME datasets, CoT with one example for Math500, and 5-shot CoT for MMLU-Pro. Exact prompts and evaluation methodology follow https://www.vals.ai/benchmarks/.

While the framework supports building sophisticated agents beyond prompt variations, we expect limited use of its full expressive power in these experiments given the minimal framework documentation and single agent example provided.

#### 2.1.2 RESULTS

Results are shown in Table 2.

The meta-agent successfully generated executable agents across all datasets with 100% success rate, even for MMLU-Pro—a previously unseen task provided with only a single example and brief description. The generated agents demonstrated strong performance with median performance approaching that of the standard prompts for the tasks. On all tasks, the best agents significantly out-

220 221 222

229

230

231

232

233

234 235

236

237

238

239

240

241

242

243

244 245

246 247

248

249 250

251

252

253

254 255

256

257

258

259

260

261

262

264

265

266

267

268

269

Table 2: Performance when generating 10 agents for each task: Success rates (percentage of generated agents that execute without errors) and task performance metrics (median and max). Base accuracy reports performance of standard prompts used for these tasks (see https://www.vals. ai/)

Dataset	Success	Base Acc.	1 Task dat	a point	5 Task data points		
	Rate		Median Acc.	Max Acc.	Median Acc.	Max Acc.	
AIME24 I	100%	22.7	21.3	25.3	20.7	25.3	
Math500	100%	70.0	66.0	74.0	69.0	73.0	
MMLU-Pro	100%	82.0	82.1	87.1	82.9	90.0	

performed the baselines, reaching maximum accuracies of 25% on AIME2024, 74% on Math500, and 90% on MMLU-Pro (an 8% absolute improvement). The primary challenge for agent generation was ensuring consistency between the output format specification in the configuration, the formatting instructions in the prompt template, and the evaluation requirements (which expect the same output fields as shown in the task examples). Misalignment across these components resulted in task agents that produced un-parseable outputs, leading to evaluation failures and zero accuracy scores.

All agents implemented dataset-specific preprocessing and incorporated detailed chain-of-thought reasoning through structured step-by-step problem breakdown. The AIME2024 best agent wrote a 6-step processes: "1. Read and understand the problem carefully, 2. Break down key information, 3. Plan your solution [...] 6. Verify answer is integer between 000-999". In term of temperature. 80% of agents used temperature 0.7, while the rest used used 0.2 or 0.3. Input/output token patterns shows MMLU-Pro's best agent used more tokens than the average (+13% input, +17% output), while the other tasks showed minimal token differences. No agent incorporated worked examples or few-shot demonstrations. Finally, the agent generation performance remains similar across datasets regardless of whether 1 or 5 task examples are provided, demonstrating that models can understand new tasks from a single example.

#### 3 DATA-DRIVEN AGENT IMPROVEMENT

We next explore whether LLMs can not only create agents but also improve them over time through iterative optimization, using meta-agents to analyze performance and generate better configurations.

We implement a simple self-improvement orchestration pattern involving four agents: a main task agent  $A_{\text{task}}$  that is the target of improvement, evaluator  $A_{\text{eval}}$ , analyzer  $A_{\text{analyze}}$ , and improver  $A_{\text{improve}}$ . We employ simple scorers as  $A_{\text{eval}}$  rather than LLM-based ones to ensure reliable performance measurements. Although these experiments focus on optimizing  $A_{task}$ , the analyzer and improver agents themselves can serve as optimization targets in the same iterative process.

More precisely, given a development dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$  and performance metric  $\rho$ , we seek to find agents  $A'_{\text{task}}$  that improve  $\Sigma_i \rho(A'_{\text{task}}(x_i), y_i)$ .

**Improvement Loop:** At iteration t, the system executes:

$$Out_t = A_{\text{task}}^{t-1}(\mathcal{D}, Ctx) \tag{1}$$

$$Eval_t = A_{eval}(Out_t, \mathcal{D})$$
 (2)

$$A_t = A_{\text{analyze}}(\text{Eval}_t, \text{Out}_t, \mathcal{D})$$
(3)

$$A_{\text{task}}^{t} = A_{\text{improve}}(A_{\text{task}}^{t-1}, A_{\text{task}}^{best}, A_{t}, \mathcal{H}_{t-1}, \text{Ctx})$$

$$\mathcal{H}_{t} = \mathcal{H}_{t-1}; (\text{Out}_{t}, \text{Eval}_{t}, A_{t})$$

$$(4)$$

$$\mathcal{H}_t = \mathcal{H}_{t-1}; (\text{Out}_t, \text{Eval}_t, A_t)$$
 (5)

where  $Out_t$  are task outputs,  $Eval_t$  contains evaluation metrics (aggregate and individual results),  $A_t$ provides analytical insights, and  $\mathcal{H}_t$  is the improvement history.

The analyzer's role it to identify patterns in performance and provide qualitative insights; we define an analyzer that uses aggregate results and randomly samples correct and incorrect examples from the data. The improver uses current and best task agents, current performance analysis and historical data, and is instructed to write a new agent while respecting the *modifiable* flag. The process starts with an initial agent  $A^0_{\rm task}$ , and terminates when a target performance is met or maximum number of iterations is reached. With respect to the optimization objective, the loop implements metricagnostic improvement by embedding the target metric within the evaluator, allowing the meta-agents to implicitly understand and optimize toward the desired goal.

The Analyzer and Improver agents are given in Appendix E and F respectively, while their prompt templates are summarized in Appendix I. Importantly, as it can be observed, these agents have minimal prompts and rely on the documentation artifacts that are provided to them (Ctx): the README associated to the LLM framework (as in the previous section) and an additional README file describing the iterative improvement loop.

The Analyzer and Improver agents are detailed in Appendices E and F, with prompt templates in Appendix I. These meta-agents use minimal prompts and instead rely on documentation (Ctx): the LLM framework README and an additional README describing the iterative improvement process. This is an artifact-driven approach where rich contextual information replaces engineered prompts.

#### 3.1 EXPERIMENTS

We included math tasks in our experiments because they remain challenging for pre-reasoning state-of-the-art models and exhibit test-time scaling where longer outputs improve accuracy (Snell et al., 2024; Muennighoff et al., 2025). This property creates a natural optimization space: agents can achieve better quality through longer reasoning chains at higher computational cost, or maintain quality while reducing token usage, making them ideal testbeds for iterative improvement loops.

This section introduces additional pseudo-random number generator (PRNG) datasets specifically designed to test test-time scaling properties. Unlike other tasks where models perform pattern recognition, PRNG sequence prediction offers no shortcuts and requires precise algorithmic execution. We use three PRNG algorithms (BBS/LCG/MWC) with deliberately small constants to be more approachable for LLMs, which do not excel at large number arithmetic. We set the task of predicting the k-th number in the sequence with k=2, based on the observation that even strong models perform poorly on this task. In contrast, reasoning models (GPT-OSS-120b) generate long reasoning traces and achieve 100% accuracy even for larger k values. We implement an initial PRNG solver agent as a generic CoT code execution agent. See Appendix G for details on data set creation and PRNG agent.

#### 3.1.1 ACCURACY OPTIMIZATION

We conduct 10-iteration improvement cycles using the same Analyzer/Improver across all tasks. We test Sonnet 3.5 v2 and GPT-OSS-120b as meta-agent LLMs and restrict task models to Sonnet 3.5 v2. As in previous sections, we use an Evaluator that computes accuracy and provides aggregate and per-problem scores.

Table 3 shows results on development and test splits. The "Original agent" uses the standard prompt from the previous section and serves as the starting point for all improvement loops. The "Best agent" column shows the highest-performing configuration from the previous section's agent generation experiments (performed on the development set).

AIME and MMLU-Pro All agents achieve improved performance on development sets. The challenging AIME datasets show modest improvements ( $0.22 \rightarrow 0.27$  with Sonnet 3.5 v2). MMLU-Pro demonstrates strong overall performance, with substantial development improvements ( $0.82 \rightarrow 0.89/0.90$ ) translating to more modest but consistent test gains. Overall, the results indicate that only larger development improvements transfer to test gains, suggesting insufficient development data or that baseline configurations had already undergone significant optimization in prior work. In terms of solutions found, the best AIME agent identifies problem types using word matching and provides targeted guidance such as "Count systematically by cases" for grid problems and "Factor completely" for number theory problems.

Appendix H shows a typical analyzer output from Sonnet analyzers. On AIME for example, analyzers commonly: identify where agents make arithmetic errors or use incorrect formulas, observe performance across problem types (e.g. low on geometry), highlight cases where multiple solu-

Table 3: Performance of initial configurations and after improvement loops. Original agent is the standard task prompt and serves as the  $A_{\text{task}}^0$ . Best agent reports on the best configuration detected in the previous section.

	Original Agent	Best Agent	Improvement Loop		
Dataset			Sonnet 3.5 v2	GPT-OSS-120b	
AIME2024 I (dev)	0.22	0.25	0.27	0.25	
AIME2024 II	0.07	0.05	0.11	0.03	
AIME2025	0.05	0.03	0.03	0.03	
MMLU-Pro-dev	0.82	0.90	0.89	0.84	
MMLU-Pro	0.79	0.83	0.80	0.77	
PRNG-bbs	0.21	_	0.57	1.0	
PRNG-lcg	0.0	-	0.02	1.0	
PRNG-mwc	0.0	-	1.0	1.0	

tions yield different answers (and recommend lower temperature), reason about the optimization trajectory, or highlight the need for more exploration by increasing temperature. Improvement loops (plotted in Appendix K) are not monotonic, reflecting that the optimization process is mostly exploratory, with agents often showing temporary performance degradation before discovering more effective approaches.

**PRNG tasks** For PRNG tasks, the two meta-agents discovered very different solutions. By iterations 3 or 4, GPT-OSS-120b invariantly discovered agents that pre-compute correct answers and instruct the LLM to output them. Specifically, the Analyzer, which sees correct/incorrect examples, inferred that all inputs list the same algorithm and instructed the Improver to create agents that execute this algorithm. The exact implementation varied: for example inserting a prompt section "verification\_note" which instructs the LLM to check the answer against the pre-computed answer or a section "generate\_example\_trace" which despite the name contains the actual computed execution trace for that data point. While not technically cheating, these solution assume that the observed examples are representative of the true data distribution. In contrast the Sonnet loops only discovered this solution for the MWC algorithm and failed to find a working solution for the LCG algorithm. For BBS, Sonnet 3.5 finds improved solutions which list explicit modulo arithmetic rules (which the Analyzer identifies as difficult), add examples and clear validation guidance, but do not compute the function deterministically (reaching 57% performance). Regarding test-time scaling, the system did not identify it as a reliable performance improvement strategy. While one improvement loop produced agents that generated 50% more output tokens (with 10% longer input), this increased verbosity did not consistently translate to better performance. Appendix J shows an extreme case where detailed reasoning instructions resulted in outputs twice the length of the starting agent.

**Variation across loops** We observe that results vary when running identical improvement loops multiple times. To explore this, we run each Sonnet loop 4 additional times and observe max scores in the range [0.84-0.89] for MMLU-Pro, [0.24-0.28] for AIME, and [0.57-0.93] for PRNG-bbs. Combined with the observation that larger development improvements generalize better, this suggests the improvement loop in its current form serves as a discovery tool rather than monotonic optimization. Despite this variability, a *single* 10-iteration loop discovers superior solutions within an effectively infinite search space of pre-processing functions, parameters, and prompt templates, showing very good optimization efficiency.

#### 3.1.2 ACCURACY+COST OPTIMIZATION

We next explore optimization under different objectives using more complex initial configurations. To balance accuracy and efficiency, we introduce a cost-accuracy evaluator where cost equals input\_tokens + 3 × output\_tokens (reflecting real-world pricing). The scoring function assigns 0 points for incorrect answers and 1/cost for correct ones, prioritizing accuracy while minimizing cost as a secondary objective. We modify optimization goals by changing the evaluator's scoring computation and description. Figure 1 shows fragments of an evaluator output.

```
"overall_accuracy": 0.267,
"overall_score": 2.1e-05,
380 3 "avg_input_tokens": 13006.1,
381 4 ...
382 5 "optimization_goal": "Minimize cost while maintaining accuracy",
"scoring_explanation": "Score: 0 if wrong answer, 1/cost if [...]"
```

Figure 1: Example evaluator output for cost+accuracy optimization. The explanation fields are intended to guide the optimization goals of the meta-agents

We experiment on AIME and MMLU-Pro datasets. MMLU-Pro uses the CoT+ICL agent from the previous section as  $A_{\rm task}^0$ . For AIME experiments, we start with an ICL agent that samples 2 examples per problem from the s1k dataset (Muennighoff et al., 2025), which contains similar problems with reasoning traces from powerful models. This dataset has been shown to improve AIME performance when used to fine-tune models.

Table 4: Performance of initial configurations and after improvement loops for cost+accuracy optimization: Accuracy and average input/output tokens.

	Ori	ginal Co	nfig	Improvement Loop Sonnet 3.5 v2			GPT-OSS-120b		
Dataset	Acc	In	Out	Acc	In	Out	Acc	In	Out
AIME2024 I (dev)	0.24	12.8k	340	0.33	6.1k	347	0.29	478	197
AIME2024 II	0.09	13.9k	368	0.08	5.9k	350	0.07	483	209
AIME2025	0.05	14.2k	365	0.05	6.0k	334	0.03	488	216
MMLU-Pro (dev)	0.82	760	190	0.84	301	176	0.84	230	196
MMLU-Pro	0.79	805	192	0.78	256	172	0.79	275	196

Results are shown in Table 4. Both Sonnet and GPT models explored much more diverse solutions than in previous experiments, likely due to the more complex starting agent (for AIME) and the ambiguous cost function allowing both accuracy and cost optimizations. For AIME, agents implemented strategies including: 1) automatic problem categorization into geometric, algebraic, and combinatorial types, 2) targeted prompts for each problem type (e.g., for geometry: "1. Draw and label key elements, 2. List given measurements and relationships..."), and 3) type-based ICL example selection. This approach yielded the winning Sonnet solution. However, the winning GPT-OSS-120b agent (Appendix L) used none of these techniques, instead replacing random ICL sampling with a fixed geometry example. While cost-effective, this did not generalize to good test performance. The variation in discovered solutions suggests the development set may be too small for reliable AIME optimization. On MMLU-Pro, all agents were able to find much more efficient solutions in terms of cost with roughly 70% reduction in input tokens and no impact on test performance.

Overall the cost+accuracy optimization showed the potential for creative solution discovery, with agents autonomously developing problem classification systems, targeted prompting strategies, and token compression techniques. More sophisticated improvement patterns could incorporate larger development sets, extended exploration phases, or multi-objective optimization strategies in order to further expand the range of discoverable solutions.

#### 4 RELATED WORK

**Large Language Models as Optimizers** LLMs have emerged as powerful optimization tools across diverse domains. Automatic prompt engineering (Ramnath et al., 2025; Zhou et al., 2022; Li et al., 2025a; Debnath et al., 2025) has demonstrated LLMs' ability to generate and refine prompts, with APE Zhou et al. (2022) introducing instructions as *programs* that can be systematically optimized. Similar work addresses automatic model selection (Wang et al., 2023; Tanaka et al., 2023; Tang et al., 2024), ICL example selection (Liu et al., 2022; Zhang et al., 2022; Do et al., 2024),

or prompt compression (Jiang et al., 2023; Li et al., 2025b) among many others. Recent advances explore LLMs as general-purpose optimizers (Jiang et al., 2025a; Yang et al., 2024a; Fernando et al., 2023), leveraging natural language understanding to interpret problems, generate solutions, and analyze feedback iteratively.

Building on these, our work introduces workflow optimization that extends beyond individual components to optimize entire agent configurations simultaneously: the system shows solution discovery capabilities that go beyond prompt variations. The approach enables metric-agnostic optimization through modular evaluator components, supporting flexible natural language-expressed optimization goals. Additionally, the self-referential optimization where meta-agents are defined within the same configurable framework as the task agents they optimize, enables the improvement process to enhance itself. Finally we designed this as an artifact-driven framework which leverages rich contextual documentation and examples to create a self-evolving system that accumulates knowledge rather than requiring extensive meta-prompt engineering.

**Autonomous Agents and LLM Workflow Optimization Frameworks** Research in LLM-based systems has advanced along two complementary but very related directions: autonomous agents for dynamic problem-solving and optimization frameworks for systematic workflow improvement.

Autonomous Agent Systems have demonstrated remarkable capabilities in dynamic environments. ReAct Yao et al. (2023) integrates reasoning with external tools through "thought-action-observation" loops, enabling LLMs to reason, perform actions, and learn from outcomes. PAL Gao et al. (2022) focuses on programmatic reasoning by generating executable code as intermediate steps, e.g. offloading complex calculations to interpreters for improved accuracy on mathematical tasks. AutoGPT and similar frameworks create agents that independently set goals, break down tasks, and execute multi-step plans with minimal human intervention. AIDE Jiang et al. (2025b) introduces tree-search strategies for systematic code and ML pipeline optimization, structuring solutions hierarchically for targeted improvements based on incremental changes.

Many optimization frameworks have emerged to systematically improve LLM workflows. Some examples are DSPy (Khattab et al., 2023), which enables "programming, not prompting" through modular pipeline definition and automatic prompt optimization. TextGrad (Huang et al., 2023) introduces "automatic differentiation via text," using LLM-based feedback for iterative refinement within computational graphs. GPTSwarm (Zhuge et al., 2024) represents agents as optimizable computational graphs, unifying prompt engineering techniques through node-level optimization and edgelevel orchestration. Teola Zheng et al. (2024) provides end-to-end optimization through dataflow graph representations, enabling rule-based optimizations across workflow modules.

Our approach bridges these paradigms through simple primitives that enable structured, interpretable optimization that combines the systematic exploration of optimization frameworks with the adaptive improvement mechanisms of autonomous agents. Unlike open-ended autonomous exploration that can produce black-box improvements, the configuration-driven approach ensures all generated solutions remain human-verifiable. Instead of replacing humans, the system can serve as a solution discovery tool that aids human decision-making through transparent improvement trajectories. This is possible through the constrained but expressive solution space, resulting in a middle ground between rigid optimization frameworks and unconstrained exploration systems.

#### 5 Future Work

Current improvement cycles analyze development data point-by-point. While necessary for runtime execution, development stages could benefit from global dataset processing to enable aggregate analysis and dataset-wide optimization strategies not visible in individual examples. Additionally, more sophisticated improvement patterns could incorporate multi-objective optimization or Monte Carlo solution search, and beyond fixed meta-agents, the configuration-driven design provides a foundation for joint optimization of both task and meta-agents.

Finally, the framework can be extended to support cross-domain knowledge transfer by accumulating improvement artifacts across domains. This would build a repository of successful patterns and strategies, with accumulated knowledge informing improvements on new tasks and enabling transfer learning at the agent architecture level.

#### REFERENCES

- Zina Chkirbene, Ridha Hamila, Ala Gouissem, and Unal Devrim. Large language models (llm) in industry: A survey of applications, challenges, and trends. In 2024 IEEE 21st International Conference on Smart Communities: Improving Quality of Life using AI, Robotics and IoT (HONET), pp. 229–234, 2024. doi: 10.1109/HONET63146.2024.10822885.
- Tonmoy Debnath, Nurul Absar Siddiky, Muhammad Enayetur Rahman, Prosenjit Das, and Antu Kumar Guha. A comprehensive survey of prompt engineering techniques in large language models. *arXiv preprint arXiv:2503.04403*, 2025.
- 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. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7308–7327, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.395. URL https://aclanthology.org/2024.acl-long.395/.
- 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.
- Luyu Gao, Aniruddha Madaan, Seyed Shakeri, Diyi Yu, Ling Wang, Jun Han, Han Chen, and Jian Lin. PAL: Program-aided language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 577–588, 2022.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the MATH dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks 1*, NeurIPS Datasets and Benchmarks 2021, virtual, December 2021. URL https://datasets-benchmarks-proceedings.neurips.cc/paper/2021/hash/be83ab3ecd0db773eb2dc1b0a17836a1-Abstract-round2.html.
- Tianyi Huang, Huimin Yao, Jiacheng He, Zhiyi Zheng, and Zhou Lin. TextGrad: Autodifferentiating through arbitrary text with large language models. *arXiv preprint arXiv:2308.02450*, 2023.
- Caigao Jiang, Xiang Shu, Hong Qian, Xingyu Lu, Jun Zhou, Aimin Zhou, and Yang Yu. Llmopt: Learning to define and solve general optimization problems from scratch, 2025a. URL https://arxiv.org/abs/2410.13213.
- Huiqiang Jiang, Qianhui Wu, Chin-Yew Lin, Yuqing Yang, and Lili Qiu. Llmlingua: Compressing prompts for accelerated inference of large language models, 2023. URL https://arxiv.org/abs/2310.05736.
- Zhengyao Jiang, Dominik Schmidt, Dhruv Srikanth, Dixing Xu, Ian Kaplan, Deniss Jacenko, and Yuxiang Wu. Aide: Ai-driven exploration in the space of code. 2025b. URL https://arxiv.org/abs/2502.13138.
- Omar Khattab, Arnav Singh, Zhiqing Li, Hamza Khalid, Yi Zhang, Mirian Sanjabi, Amir Gholami, Matei Zaharia, and Matei Zaharia. DSPy: Compiling declarative language model calls into state-of-the-art pipelines. *arXiv preprint arXiv:2311.00062*, 2023.
- Jiaheng Li, Junhao Hu, Jionghao Xu, Yuxin Liu, Yilun Zhang, Li Huang, Ruichen Zhao, Yunpu Wang, Jihong Liu, and Meng Wu. Large language models for constructing and optimizing machine learning workflows: A survey. *arXiv preprint arXiv:2411.10478*, 2024.
- Wenwu Li, Xiangfeng Wang, Wenhao Li, and Bo Jin. A survey of automatic prompt engineering: An optimization perspective. *arXiv preprint arXiv:2502.11560*, 2025a.
  - Zongqian Li, Yinhong Liu, Yixuan Su, and Nigel Collier. Prompt compression for large language models: A survey. In Luis Chiruzzo, Alan Ritter, and Lu Wang (eds.), *Proceedings of the 2025*

Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 7182–7195, Albuquerque, New Mexico, April 2025b. Association for Computational Linguistics. ISBN 979-8-89176-189-6. doi: 10.18653/v1/2025.naacl-long.368. URL https://aclanthology.org/2025.naacl-long.368/.

- Hunter Lightman, Vineet Kosaraju, Yuri Burda, Harrison Edwards, Bowen Baker, Teddy Lee, Jan Leike, John Schulman, Ilya Sutskever, and Karl Cobbe. Let's verify step by step. In *The Twelfth International Conference on Learning Representations*, Vienna, Austria, May 7-11 2024. Open-Review.net. URL https://openreview.net/forum?id=v8L0pN6EOi.
- Jiachang Liu, Dinghan Shen, Yizhe Zhang, Bill Dolan, Lawrence Carin, and Weizhu Chen. What makes good in-context examples for GPT-3? In Eneko Agirre, Marianna Apidianaki, and Ivan Vulić (eds.), *Proceedings of Deep Learning Inside Out (DeeLIO 2022): The 3rd Workshop on Knowledge Extraction and Integration for Deep Learning Architectures*, pp. 100–114, Dublin, Ireland and Online, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.deelio-1.10. URL https://aclanthology.org/2022.deelio-1.10/.
- MAA Committees. AIME Problems and Solutions. Art of Problem Solving Wiki. URL https://artofproblemsolving.com/wiki/index.php/AIME\_Problems\_and\_Solutions. Accessed: September 25, 2025.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Li Fei-Fei, Hannaneh Hajishirzi, Luke Zettlemoyer, Percy Liang, Emmanuel Candès, and Tatsunori Hashimoto. s1: Simple test-time scaling, 2025. URL https://arxiv.org/abs/2501.19393.
- Kiran Ramnath, Kang Zhou, Sheng Guan, Soumya Smruti Mishra, Xuan Qi, Zhengyuan Shen, Shuai Wang, Sangmin Woo, Sullam Jeoung, Yawei Wang, Haozhu Wang, Han Ding, Yuzhe Lu, Zhichao Xu, Yun Zhou, Balasubramaniam Srinivasan, Qiaojing Yan, Yueyan Chen, Haibo Ding, Panpan Xu, and Lin Lee Cheong. A systematic survey of automatic prompt optimization techniques, 2025. URL https://arxiv.org/abs/2502.16923.
- Mubashar Raza, Zarmina Jahangir, Muhammad Riaz, and Muhammad Sattar. Industrial applications of large language models. *Scientific Reports*, 15, 04 2025. doi: 10.1038/s41598-025-98483-1.
- Yasmeen Saleh, Manar Abu Talib, Qassim Nasir, and Fatima Dakalbab. Evaluating large language models: A systematic review of efficiency, applications, and future directions. *Frontiers in Computer Science*, 7:1523699, 2025. doi: 10.3389/fcomp.2025.1523699.
- Patrick Schramowski, Max Eichenseer, Konstantina Stamati, Christian Klein, David Roth, and Christian Bizer. Exploring the role of large language models in the scientific process. *Nature Machine Intelligence*, 2025. URL https://www.nature.com/articles/s44387-025-00019-5.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling Ilm test-time compute optimally can be more effective than scaling model parameters, 2024. URL https://arxiv.org/abs/2408.03314.
- Shun Tanaka, Takehiko Kadowaki, Yuya Nagai, Yuki Nishiyama, and Masaharu Kakiuchi. Autonlp: A framework for automated model selection in natural language processing. In 2023 IEEE 47th Annual Computers, Software, and Applications Conference (COMPSAC), pp. 2011–2016. IEEE, 2023.
- Jiabin Tang, Lianghao Xia, Zhonghang Li, and Chao Huang. Ai-researcher: Autonomous scientific innovation, 2025. URL https://arxiv.org/abs/2505.18705.
- Shicheng Tang, Yongxu Xu, Xiang Li, Zhixu Shi, and Chengxi Zhai. Automatic prompt selection for large language models. *arXiv preprint arXiv:2404.02717*, 2024.
- Shizhe Wang, Xiaoyu Hu, Yu He, Guoliang Xu, Hongguang Liu, and Zhenhua Lin. Automatic model selection with large language models for mathematical reasoning. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 835–845, 2023.

Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, Tianle Li, Max Ku, Kai Wang, Alex Zhuang, Rongqi Fan, Xiang Yue, and Wenhu Chen. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark, 2024. URL https://arxiv.org/abs/2406.01574.

Ziyue Wang and et al. Ai-enabled scientific revolution in the age of generative ai. *Nature*, 1(1): 1–10, 2025. doi: 10.1038/s44387-025-00018-6.

- Yuanzhi Xiao, Guoxin Li, Ziyue Wang, Xinyu Zhao, Zhenhua Lin, Chenxiao Liu, Fan Chen, Chen Zhang, Guoping Liu, and Qi Li. Prompt cache: Modular attention reuse for low-latency inference. *arXiv preprint arXiv:2311.04934*, 2023.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *The Twelfth International Conference on Learning Representations*, 2024a. URL https://openreview.net/forum?id=Bb4VGOWELI.
- Tianxing Yang, Yilun Zhang, Jiaheng Li, Junhao Hu, and Jionghao Xu. Hypothesis generation with large language models. *arXiv preprint arXiv:2404.04326*, 2024b.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Luo, Kai Xu, Jian Li, Jian Li, George Liang, Tianyu Huang, and Bo Hu. ReAct: Synergizing reasoning and acting in language models. *arXiv* preprint *arXiv*:2302.05380, 2023.
- Yiming Zhang, Shi Feng, and Chenhao Tan. Active example selection for in-context learning. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9134–9148, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022. emnlp-main.622. URL https://aclanthology.org/2022.emnlp-main.622/.
- Junyi Zheng, Ziyue Wang, Xinyu Zhao, Guoxin Li, Zekun Wang, Zhenhua Lin, Chenxiao Liu, Fan Chen, Chen Zhang, Guoping Liu, and Qi Li. Teola: Towards end-to-end optimization of llm-based applications. *arXiv preprint arXiv:2407.00326*, 2024.
- Yongchao Zhou, Andrei Song, Shizoz Ma, Jiao Zhang, Yuhan Han, Tianci Liu, Ning Ma, Quanquan Liu, Guang Sun, Ming Xu, Hong Yang, Ruibin Li, and Zhijie Zeng. Large language models are human-level prompt engineers. *arXiv* preprint arXiv:2211.01910, 2022.
- Banghua Zhu, Ying Sheng, Lianmin Zheng, Clark Barrett, Michael I. Jordan, and Jiantao Jiao. On optimal caching and model multiplexing for large model inference, 2023. URL https://arxiv.org/abs/2306.02003.
- Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jürgen Schmidhuber. Gptswarm: language agents as optimizable graphs. In *Proceedings of the 41st International Conference on Machine Learning*, ICML'24. JMLR.org, 2024.

#### A GENERATIVE AI USE

AI tools were used throughout this research to enhance productivity and improve the quality of the work.

**Framework Development Process** The framework and the framework documentation was developed iteratively with assistance from a code assistant powered by Sonnet 3.5 and Sonnet 4.0, depending on availability. The assistant was given specific instructions to improve documentation clarity and update the framework based on feedback. This iterative process was validated by testing the framework on toy prediction tasks (separate from the experimental tasks) to ensure the generated agents functioned correctly.

Paper writing LLMs assisted in refining the clarity and coherence of the writing, suggesting more concise phrasings and identifying areas where technical concepts could be explained more clearly for broader accessibility. Additionally, AI tools facilitated data analysis and presentation by generating LaTeX tables and figures from raw experimental results. All AI-generated content was carefully reviewed and validated by the authors to ensure accuracy and appropriateness for the scientific context.

### B LLM AGENT IMPLEMENTATION

**LLM Agent Implementation** At the level of implementation, we describe an LLM agent via three components: Configuration (json), Agent Class Implementation (Python) and Input Schema File.

The **configuration** file defines the agent's behavior, model settings, input processing, prompt structure, and output format. A configuration tuple C consists of:

 $C = \langle \text{desc}, \text{schema}, \text{model}, \text{inputs}, \text{prompt}, \text{output}, \text{settings} \rangle$ 

where:

- desc: Human-readable description of agent purpose
- input schema: JSON Schema path defining expected input data structure (including additional context data).
- model =  $\langle \text{name}, \rho \rangle$ : LLM specification ( $\theta_{\text{LLM}}$  and inference parameters  $\rho$ )
- inputs =  $\{f_1, \dots, f_k\}$ : Data transformation functions operating on  $x_i$  and context c
- prompt =  $\langle T, \text{sections} \rangle$ : Template T with reusable sections and placeholders
- output = \( \)format, schema, errors \( \): Expected output format which drives the LLM response parsing
- settings: Execution parameters such as timeouts

Notably, inputs can be of two types: "data" or "computed" (with function name and arguments), in which case the function needs to be defined in the agent implementation class. For example the following configuration fragment shows the input section and references <code>raw\_problem\_data</code> which needs to exist in the input data processed, and <code>problem\_text</code> which is computed via a function and takes the raw data as input.

Listing 1: Configuration example showing computed inputs

```
"inputs": {
    "raw_problem_data": {
        "source": "data"
    }
    "problem_text": {
        "source": "computed",
        "function": "extract_problem_text",
        "args": {"problem_data": "raw_problem_data"}
}
}
```

The role of the **agent class** is to implement the functions referenced in the configuration's computed inputs. These functions operate at data-point level and can apply custom pre-processing (such as normalization, feature extraction, problem type classification, etc), retrieve examples from static resources provided as context c, etc. In the example above, the python agent class needs to implement the *extract\_problem\_text* function, for example as:

```
def extract_problem_text(self, problem_data):
    """Function referenced in config's computed inputs"""
    return problem_data.get('problem', '')
```

The **input schema file** and the output section of the configuration define the communication with upstream and downstream tasks. The agent can chose to ignore parts of the input data but can't access any resources outside what is defined in the input schema file.

**Modifiability Control** In order to simplify the process of other LLMs improving task agents, we introduce modifiability control. Each configuration section includes a flag  $modifiable \in \{\text{true}, \text{false}\}$  that guides other agents wrt which sections can be modified. For example, model.modifiable = true indicates that model name or model hyper-parameters can be modified, while output.modifiable = false indicates that the response schema needs to be preserved. This enables systematic exploration

of the configuration space while maintaining structural invariants essential for downstream processing, as well as the option to configure preferences (for example to prevent changes in the underlying LLM model if so desired).

**Agent Execution** The agents operate synchronously, processing one input data point at a time in sequential order. The framework reads the agent definition from three files (JSON configuration, Python implementation, and input schema), performs static validation of the configuration and dynamic validation of runtime input data. At runtime, the framework processes input data through defined transformations, generates the final prompt by resolving placeholders, and calls the LLM with the rendered prompt. The LLM response is then parsed by an output parser according to the output format specified in the configuration.

### MATH PROBLEM SOLVING AGENT

#### Listing 2: AIME 2024 Task Agent Configuration

```
813
814
          "description": "AIME 2024 task agent - generates detailed reasoning followed by 3-digit
815
               integer answer",
          "input_schema": "agents/task/task_input_schema.json",
816
817
            "_modifiable": true,
            "name": "anthropic.claude-3-5-sonnet-20241022-v2:0",
818
            "parameters": {
819
              "temperature": 0.7,
              "max_tokens": 8192
820
821
     11
    12
           "inputs": {
822
              _modifiable": true,
     13
823
     14
             "problem_text": {
    15
               "source": "computed",
824
              "function": "extract_problem_text",
     16
825
     17
               "args": {"problem_data": "raw_problem_data"},
              "description": "Extracted problem text from AIME2024 format"
    18
826
     19
827
    20
             "raw_problem_data": {
    "source": "data",
    21
828
              "path": "$",
829
    23
              "required": true,
               "description": "Complete problem data from AIME2024 dataset"
830
    24
    25
831
    26
    27
           "prompt": {
832
             "_modifiable": true,
    28
833
    29
            "sections": {
               'instructions": "You are a mathematical problem solver specializing in AIME problems.",
    30
834
              "task_description": "AIME problems require integer answers between 000 and 999. Show
    31
835
                   detailed step-by-step reasoning, then provide the final numerical answer.",
836 32
              "output_format": "You must respond with a JSON object containing exactly two fields:\n-
                   837
                    string containing your 3-digit answer (pad with leading zeros if needed)"
    33
838
    34
            "template": "### Instructions:\n{prompt.sections.instructions}\n\n### Task:\n{prompt.
839
                 sections.task_description}\n\n### Problem:\n{inputs.problem_text}\n\n### Output
                 840
                 the JSON object:'
841
    35
    36
           "output": {
842
     37
            "_modifiable": false,
843
    38
            "format": "json",
    39
            "schema": {
844
    40
              "type": "object",
845
    41
              "required": ["reasoning", "answer"],
    42
846
                "reasoning": {"type": "string", "description": "Detailed step-by-step solution"},
"answer": {"type": "string", "pattern": "^[0-9]{3}$", "description": "3-digit integer
     43
847
    44
848
    45
              }
849
    46
    47
             "error_values": {
850
              "reasoning": "Error occurred during reasoning.",
851
    49
              "answer": "PARSE_ERROR"
852
    51
853
    52
           "settings": {
            "_modifiable": false,
"class_name": "AIME2024TaskAgent",
    53
854
855
    55
            "timeout_seconds": 300
    56
856
    57
857
```

#### Listing 3: AIME 2024 Task Agent Implementation

```
859
         from llm_agent import LLMAgent
860
        from typing import Dict, Any
861
     3
862
     4
         class AIME2024TaskAgent(LLMAgent):
863
     6
             AIME2024-specific reasoning + answer agent.
```

```
864
            Processes AIME mathematical competition problems and generates detailed step-by-step
            reasoning followed by answers.
865
866
867 11
            def __init__(self, name="AIME2024_Task", config_path=None, dry_run=False):
                super().__init__(name, config_path, dry_run)
868 12
    13
869 14
            def extract_problem_text(self, problem_data: Dict[str, Any]) -> str:
    15
870
                Extract problem text from AIME2024 data format.
    16
871 17
                AIME2024 format:
872 18
    19
873 20
                     "id": int,
                     "problem": str, # Main problem statement
    21
874
                     "url": str
    22
875 23
    24
876
    25
                return str(problem_data.get('problem', ''))
877
```

#### Listing 4: AIME 2024 Task Agent Input Schema

```
879
880
             "type": "object",
"required": ["id", "problem"],
"properties": {
881
      3
882
               "id": {
883
                 "type": "integer",
                 "description": "Problem identifier"
884
885
               "problem": {
                 "type": "string",
"description": "The AIME problem statement"
     10
886
      11
887
     12
     13
                "url": {
888
                 "type": "string",
"description": "URL to problem source (optional)"
     14
889
      15
     16
890
      17
            }
891
     18
892
```

#### D FRAMEWORK CAPABILITIES

The LLM Framework provides extensive flexibility for implementing diverse agent architectures and optimization strategies through its configuration-driven design. The framework's expressive power enables sophisticated agent implementations across multiple dimensions:

**Feature Extraction and ML Integration:** Agents can incorporate external machine learning models and feature extraction pipelines through initialization methods. Complex preprocessing can be implemented in the agent's \_\_init\_\_ method, loading pre-trained models, statistical analyzers, or domain-specific feature extractors. External datasets for feature computation can be specified in the configuration's input section and loaded during agent initialization:

```
"inputs": {
    "feature_dataset": {
        "source": "data",
        "path": "external_features.json",
        "description": "External dataset for feature extraction"
    },
    "extracted_features": {
        "source": "computed",
        "function": "extract_ml_features",
        "args": {"raw_data": "problem_data", "feature_db": "feature_dataset"}
    }
}
```

**In-Context Learning Implementation:** ICL can be most flexibly implemented as computed inputs that dynamically generate examples using various example selection strategies, including similarity-based retrieval, difficulty-matched sampling, or domain-specific filtering. ICL examples can be sourced from multiple datasets with automatic data leakage prevention:

```
"icl_examples": {
    "source": "computed",
    "function": "generate_icl_examples",
    "args": {
        "icl_data_path": "data/training_examples.json",
        "num_examples": 5,
        "similarity_metric": "cosine",
        "exclude_current": true
    }
}
```

**Dynamic Input Transformation:** The computed input system enables arbitrary feature engineering and input pre-processing. Agents can implement domain-specific transformations, adaptive input formatting, etc. Each transformation function can access multiple input sources and apply complex processing logic:

```
"processed_input": {
    "source": "computed",
    "function": "transform_input",
    "args": {
        "text_data": "raw_text",
        "numerical_data": "raw_numbers",
        "context_data": "external_context",
        "transformation_type": "hierarchical_fusion"
    }
}
```

Model and Parameter Optimization: The framework supports dynamic model selection and parameter tuning through the modifiable model configuration. Agents can be optimized across dif-

ferent LLM models with varying capabilities and costs. Model parameters including temperature, max\_tokens, and model-specific settings can be systematically explored:

```
"model": {
   "_modifiable": true,
   "name": "anthropic.claude-3-5-sonnet-20241022-v2:0",
   "parameters": {
      "temperature": 0.7,
      "max_tokens": 4096,
      "top_p": 0.9
   }
}
```

**Prompt Engineering and Structure:** The modular prompt system enables sophisticated prompt architectures with independently optimizable sections. Agents can implement hierarchical prompting, conditional prompt sections, dynamic example integration, and task-specific reasoning templates. The template system supports complex placeholder substitution and structured prompt composition.

**Orchestration Constraints:** While the framework provides extensive flexibility, certain constraints ensure orchestration compatibility. The output schema should remain fixed (\_modifiable: false) to maintain consistency across the evaluation-analysis-improvement loops. This constraint ensures that orchestration agents can reliably process results while allowing freedom in input processing, model selection, and reasoning strategies.

### E ANALYZER AGENT

1026

1027 1028

Listing 5: Analyzer Agent Configuration

```
1030 1
              "description": "Generic analyzer agent - analyzes task performance and identifies
1031
             improvement opportunities",
"input_schema": "analyzer_input_schema.json",
"model": {
1032 3
1033 5
                "_modifiable": true,
                "name": "anthropic.claude-3-5-sonnet-20241022-v2:0",
1034 6
                "parameters": {
1035
                   "temperature": 0.3,
1036 9
                   "max_tokens": 8192
1037 \frac{10}{11}
             },
"inputs": {
1038 12
1039 \frac{13}{14}
                "_modifiable": true,
                "current_task_config": {
    "source": "data",
    "path": "current_task_config",
1040 15
1041 \frac{16}{17}
                   "required": true,
1042 18
                   "description": "Current task agent configuration as JSON string"
1043 \frac{19}{20}
                "evaluation_summary": {
                   "source": "data",
"path": "evaluation_summary",
1044 21
1045_{23}^{22}
                   "required": true,
1046 24
                   "description": "Aggregated evaluation metrics from evaluator"
1047 <sup>25</sup> <sub>26</sub>
                "individual_results": {
                   "source": "data",
"path": "individual_results",
1048 27
1049 28
                   "required": true,
1050 30
                   "description": "Per-problem results with model outputs, ground truth, and scores"
1051 31
                 "correct_examples": {
1052 33
                   "source": "computed",
1053 34
35
                   "function": "extract_correct_examples",
1054 36
                     "individual_results": "individual_results",
1055 <sup>37</sup><sub>38</sub>
                     "num_examples": 3,
                      "score_field": "score_field",
1056 39
                      "perfect_score": "perfect_score"
1057 40
1057 41
                   "description": "Examples where the task agent performed correctly"
1058 42
1058 <sup>12</sup>
1059 <sup>43</sup>
44
                 "incorrect_examples": {
                   "source": "computed",
                   "function": "extract_incorrect_examples",
1060 45
1061 <sup>46</sup> <sub>47</sub>
                   "args": {
                     "individual_results": "individual_results",
1062 48
                     "num_examples": 5,
1063 <sup>49</sup> <sub>50</sub>
                      "score_field": "score_field",
                      "perfect_score": "perfect_score"
1064 51
1065 52
1065 53
                   "description": "Examples where the task agent made errors"
                }
1066 54
1067 55
1067 56
                 "template": "### Instructions:\n{prompt.sections.instructions}\n\n### Current Task Agent
                      \label{local_config} $$\operatorname{Configuration:} \pi_{\operatorname{inputs.current\_task\_config}} \n\theta_{\#\#} $$ Evaluation Summary: $$ \operatorname{Summary:} \theta_{\operatorname{inputs.current\_task\_config}} $$
1068
                       evaluation_summary}\n\n### Examples Where Task Agent Was CORRECT:\n{inputs.
1069
                       correct_examples}\n\n### Examples Where Task Agent Made ERRORS:\n{inputs.
1070
                      {\tt incorrect\_examples} \\ {\tt nnAnalyze \ the \ task \ agent \ performance \ and \ identify \ improvement} \\
                      opportunities:"
1071 57
              "output": {
1072 58
                "_modifiable": false,
"format": "json",
1073 <sup>59</sup> <sub>60</sub>
                "schema": {
  "type": "object",
1074 61
1075 <sup>62</sup> <sub>63</sub>
                   "required": ["analysis"],
"properties": {
1076 64
                      "analysis": {
    "type": "string",
1077 <sup>65</sup> <sub>66</sub>
1078 67
                        "description": "Complete analysis of task agent performance"
1079 \frac{68}{69}
      70
```

```
1080 71
               },
"settings": {
1081 72
1082 <sup>73</sup> <sub>74</sub>
                   "_modifiable": false,
                   "class_name": "AnalyzerAgent",
1083 75
                   "timeout_seconds": 600
1084 \begin{array}{c} 76 \\ 77 \end{array}
1085
```

### Listing 6: Analyzer Agent Implementation

```
1087
1088 1
         from llm_agent import LLMAgent
1089 2
         from typing import Dict, List, Any
         import json
1090 4
         import random
1091 6
         class AnalyzerAgent(LLMAgent):
1092 7
             Generic analyzer agent for performance analysis.
1093 9
             Analyzes task agent performance and identifies improvement opportunities
109410
1095 11
             for any domain that follows the standard evaluator contract.
1096 13
             def __init__(self, name="Analyzer", config_path=None, dry_run=False):
1097 14
                  super().__init__(name, config_path, dry_run)
109816
1099 17
             def extract_correct_examples(self, individual_results: List[Dict[str, Any]],
                                          num_examples: int,
score_field: str = "score"
110019
1101 20
                                           perfect_score: float = 1.0) -> str:
110222
                 Extract examples where the task agent performed correctly.
1103 23
1103 24
                 correct_examples = [result for result in individual_results
                                     if result.get(score_field, 0) >= perfect_score)
110425
1105 <sup>26</sup> <sub>27</sub>
                 selected = random.sample(correct examples, min(num examples, len(correct examples)))
110628
                 return json.dumps(selected, indent=2)
1107 29 30
             def extract_incorrect_examples(self, individual_results: List[Dict[str, Any]],
110831
                                             num_examples: int,
score field: str = "score",
1109 32
                                             perfect_score: float = 1.0) -> str:
                  ,,,,,
111034
                 Extract examples where the task agent made errors.
1111 <sup>35</sup> <sub>36</sub>
                  incorrect_examples = [result for result in individual_results
111237
1113 38
                                       if result.get(score_field, 0) < perfect_score]</pre>
111440
                 selected = random.sample(incorrect_examples, min(num_examples, len(incorrect_examples))
111541
                  return json.dumps(selected, indent=2)
```

### F IMPROVER AGENT

1134

1135 1136

1183

#### Listing 7: Improver Agent Configuration

```
1137
1138 1
           "description": "Generic improver agent - generates improved task agent based on analysis",
1139 \frac{1}{3}
           "input_schema": "improver_input_schema.json",
1140 4
             "_modifiable": true,
1141
             "name": "anthropic.claude-3-5-sonnet-20241022-v2:0",
1142 7
             "parameters": {
               "temperature": 0.3,
1143 9
               "max_tokens": 25000
1144 10
1145 12
           "inputs": {
1146 13
             "_modifiable": true,
             "analysis_results": {
1147 15
               "source": "data",
               "path": "analysis_results",
1148 16
               "required": true,
     17
1149 18
               "description": "Analysis results from the analyzer agent"
1150 19
     20
              "current task agent": {
1151 21
               "source": "data",
1152 <sup>22</sup> 23
               "path": "current_task_agent",
               "required": true,
1153 24
               "description": "Current task agent configuration"
1154 <sup>25</sup>
26
1155 <sup>27</sup>
             "best_task_agent": {
   "source": "data",
1156 <sup>28</sup>
               "path": "best_task_agent",
               "required": false,
     29
1157 30
               "description": "Best performing task agent configuration and implementation"
1158 31
32
1159 33
           "prompt": {
              "template": "### Instructions:\n{prompt.sections.instructions}\n\n### Current Task Agent:\
1160 34
                  n{inputs.current_task_agent}\n\n### Best Performing Agent:\n{inputs.best_task_agent}\
1161
                  1162
                  configuration based on the analysis:"
1163 36
           "output": {
1164 37
             "_modifiable": false,
             "format": "json",
     38
1165 39
             "schema": {
  "type": "object",
1166 40
1167 41
1167 42
               "required": ["explanation", "new_task_implementation", "new_task_config"],
               "properties": {
1168 43
                  "explanation": {
                   "type": "string"
1169 45
                   "description": "Detailed explanation of improvements"
1170 46
                 "new_task_implementation": {
   "type": "string",
   "description": "Complete Python file content"
1171 48
1172 49
     50
1173 51
                  "new_task_config": {
1174 52
                    "type": "string",
                    "description": "Complete JSON config content"
1175 54
1176 55
     56
             }
1177 57
1178 58
           "settings": {
             "_modifiable": false,
1179\overset{37}{60}
             "class_name": "ImproverAgent",
1180 61
             "timeout_seconds": 1200
     62
1181 63
1182
```

#### Listing 8: Improver Agent Implementation

```
1184
1185 \frac{1}{2}
         from llm_agent import LLMAgent
1186 3
         class ImproverAgent(LLMAgent):
1187
     5
             Generic improver agent - generates improved task agents based on analysis.
```

```
1188 7
1189 8
            def __init__(self, name="Improver", config_path=None, dry_run=False):
1190 9
                super().__init__(name, config_path, dry_run)
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
```

#### G PRNG EXPERIMENTS

 We use the three algorithms BBS/LCG/MWC, and deliberately set small constants to be more approachable for LLMs which do not excel at large number arithmetic:

- LCG:  $X_{n+1} = (25173 \cdot X_n + 13849) \mod 65536$  smaller modulus than typical implementations
- MWC:  $X_{n+1} = (18782 \cdot X_n + \text{carry}) \& 0 \text{xFFF} 12 \text{-bit output with carry propagation}$
- **BBS**:  $X_{n+1} = X_n^2 \mod 209$  where  $209 = 11 \times 19$  small primes

We create train and test data sets (each containing 70 sequences) where models must predict the k'th number in the sequence given: (1) complete algorithm implementation in Python, (2) hyperparameters (e.g. the prime numbers used) and (3) initial seed value (distinct in train/test splits). We set the task as that of predicting the second number in the sequence (k = 2) based on the initial observation that even state-of-the-art models perform poorly on the task. In contrast, we confirm that a reasoning model (OpenAI 120b) generates long reasoning traces and achieves 100% accuracy.

Table 5: PRNG agents for BBS/LCG/MWC algorithms. The names of the algorithms are obfuscated in case the LLMs are aware of these algorithms and have memorized sequences.  $\langle Function \rangle$  is the Python code,  $\langle N \rangle$  is the number to be predicted (in this case at position 2),  $\langle Initial\ Value \rangle$  the seed,  $\langle Hyperparams \rangle$  are constants used in the algorithms. Note that this is a generic code execution agent, as all the placeholders are only resolved at run-time when individual data points are processed.

Instructions: You are a function evaluator. You will be given a function implementation, its hyperparameters, and an initial value. Your task is to predict what the function will return at a specific position. Task: Given the function implementation and parameters, determine what value the function would return at position  $\langle N \rangle$  when called with the initial value. Note: The first number generated is at position 1, the first iteration is position 1. Show your detailed step-by-step reasoning, then provide the final prediction.

Function Implementation: 〈Function〉 Hyperparameters: 〈Hyperparams〉

Initial Value: (Initial Value)

Approach: 1. Carefully analyze the function implementation 2. Understand what the hyperparameters control 3. Trace through the function execution step by step 4. For modulo operations (%), decompose them: compute the division first, then find the remainder - do NOT compute modulo in a single step 5. For ANY arithmetic expression with more than 2 terms (e.g., a \* b + c, x + y - z, etc.), break it down step by step and compute intermediate results - do NOT evaluate complex expressions in one shot 6. Calculate the value at position  $\langle N \rangle$  7. Present your reasoning followed by the final answer Output Format: JSON object with "reasoning" (detailed step-by-step analysis) and "answer" (predicted

number as string) fields.

#### 1350 ANALYZER OUTPUT EXAMPLES 1351 1352 Analysis Output Example 1353 1354 1. Performance Overview: 1355 • Current score: 0.18666 (28% improvement from initial 0.17333) 1356 Inconsistent performance across iterations (scores: 1357 $0.17 \rightarrow 0.16 \rightarrow 0.0 \rightarrow 0.22 \rightarrow 0.0 \rightarrow 0.22$ 1358 • Best performance achieved in iterations 3 and 5 with score 0.22666 1359 2. Key Error Patterns: 1360 A. Mathematical Reasoning Failures: 1361 • Complex geometric problem (#70) shows systematic failure to: 1362 - Properly handle cyclic quadrilateral properties 1363 Apply collinearity constraints 1364 - Use similar triangle relationships effectively 1365 • Function intersection problem (#63) reveals: - Inadequate periodic function analysis 1367 - Poor handling of composite functions 1368 - Failure to systematically count intersection points 1369 **B. Framework Implementation Issues:** 1370 • Overcomplicated input processing attempts (iterations 1-2) 1371 · Inconsistent JSON output formatting 1372 3. Successful Elements: 1373 · Basic mathematical structure identification 1374 • Geometric component recognition 1375 Clear prompt organization with validation steps 1376 · Strong JSON output examples 1377 1378 4. Improvement Pattern Analysis: Most Effective Changes (Iterations 3/5): 1379 1380 · Balanced mathematical structure detection 1381 · Clear geometric reasoning guidance 1382 • Simple but effective input processing 1383 • Strong output format examples 1384 **Counterproductive Changes:** 1385 • Overly complex pattern analysis (iteration 1) 1386 • Excessive preprocessing (iteration 2) 1387 • Oversimplified approach (iteration 4) 1388 5. Recommended Focus Areas: 1389 A. Input Processing: 1390 • Keep clean\_problem and math\_structure 1391 • Focus geometric\_analysis on fundamental properties 1392 · Avoid complex pattern detection 1393 1394 **B. Prompt Engineering:** 1395 · Strengthen systematic proof structure 1396 · Enhance case analysis guidance 1397 • Maintain clear validation steps 1398 C. Mathematical Framework: 1399 • Add explicit proof technique selection 1400 · Improve periodic function analysis 1401 Strengthen geometric relationship handling 1402 6. Framework-Compatible Implementation: 1403

• Use computed inputs for basic math structure [...]

### I ANALYZER AND IMPROVER AGENTS: PROMPT TEMPLATES

Table 6: Analyzer and Improver agents. Temperature is set to 0.7 to encourage exploration.

#### Analyzer Agent

Instructions: You are the Analyzer agent in the Iterative Task Improvement Pattern that evaluates LLM agent performance on tasks by examining inputs, outputs, and ground truth. Analyze error patterns and suggest framework-compatible improvements.

LLM Agent Framework: (Framework Documentation)

You are an analyzer agent that works as part of this iterative improvement workflow: (Iterative Improvement Documentation)

Current Task Agent Configuration: (Current Task Config)

Current Task Agent Implementation: (Current Task Implementation)

1416 Evaluation Summary: (Evaluation Summary)

Examples Where Task Agent Was CORRECT: (Correct Examples)
Examples Where Task Agent Made ERRORS: (Incorrect Examples)

Iteration History: (Iteration History)

Analysis Focus: Error patterns across multiple runs, Framework constraints (single LLM call, config-driven), General improvements (avoid overfitting to sample data), Past iteration results and performance trends, Which changes had positive/negative effects on performance, Correlation between specific improvements and score changes

Output Format: (Output Format Requirements)

#### Improver Agent

Instructions: You are the Improver agent in the Iterative Task Improvement Pattern. Generate improved task agents that follow LLM Agent Framework constraints (single LLM call, config-driven architecture). IMPORTANT: The goal is to create agents that generalize well to UNSEEN data from the same distribution. The current development data is only for evaluation - focus on improvements that will work on new, unseen problems rather than overfitting to the specific development examples.

LLM Agent Framework: (Framework Documentation)

Iterative Improvement Pattern: (Iterative Improvement Documentation)

Iteration History: (Iteration History)

Current Task Agent: (Current Task Agent)

Best Performing Agent (Score: \langle Best Agent Score \rangle from iteration \langle Best Agent Iteration \rangle): \langle Best Task Agent \rangle

Analysis Results: (Analysis Results)

Improvement Guidelines: Framework Implementation Rules: All improvements must work within single LLM call architecture, Add computed inputs with processing functions for data transformation, Reference new inputs in prompt templates using {inputs.input\_name}, Only implement class methods referenced by computed inputs. Explain changes like: 'Added computed input X that processes Y using function Z, referenced in prompt template as {inputs.X} to achieve improvement W'. IMPORTANT: Remember you can only modify config fields with \_modifiable set to True.

Format Requirements: (Output Format Requirements)

#### J AGENT GENERATED BY SONNET IN A PRNG-BBS IMPROVEMENT LOOP

1458

```
1460
1461 2
             "description": "Enhanced Func1 solver with comprehensive arithmetic validation",
1462 3
             "input_schema": "func1_solver_input_schema.json",
             "model": {
1463
                "_modifiable": true,
                "name": "anthropic.claude-3-5-sonnet-20241022-v2:0",
1464 6
                "parameters": {
1465
                  "temperature": 0.0,
"max_tokens": 50000
1466 9
1467 <sup>10</sup> <sub>11</sub>
1468 12
             "inputs": {
                "_modifiable": true,
1469 <sup>13</sup><sub>14</sub>
                "function_implementation": {
   "source": "computed",
1470 15
                  "function": "extract_implementation",
1471 \frac{16}{17}
                  "args": {
1472 18
                     "problem_data": "raw_problem_data"
1473 \frac{19}{20}
1474 21
                "max_value_validator": {
1475 22
                  "source": "computed",
                  "function": "max_value_validator",
1476 24
                  "args": {
1477_{26}^{25}
                     "problem_data": "raw_problem_data"
1478 27
1479 <sup>28</sup> <sub>29</sub>
                "value_tracker": {
                  "source": "computed",
"function": "value_tracker",
1480 30
1481 <sup>31</sup>
<sub>32</sub>
                  "args": {
                     "problem_data": "raw_problem_data"
1482 33
                  }
1483 34
35
               },
"checksum_generator": {
    "computed",
1484 36
                  "source": "computed",
1485 <sup>37</sup><sub>38</sub>
                  "function": "checksum_generator",
1486 39
                     "problem_data": "raw_problem_data"
1487 \frac{40}{41}
1488 42
                "division_validator": {
1489 43
1489 44
                  "source": "computed",
                  "function": "division_validator",
1490 45
                  "args": {
1490 13
1491 46
1491 47
                     "problem_data": "raw_problem_data"
1492 48
1493 <sup>49</sup> <sub>50</sub>
                "quotient_checker": {
                  "source": "computed",
1494 51
                  "function": "quotient_checker",
1495 {52 \atop 53}
                  "args": {
                     "problem_data": "raw_problem_data"
1496 54
1496 54
1497 55
                "range_validator": {
                   "source": "computed",
1498 57
                  "function": "range_validator",
1499 58
1499 59
                  "args": {
1500 60
                     "problem_data": "raw_problem_data"
1501 <sup>61</sup> <sub>62</sub>
                  }
1502 63
                "function_steps": {
    "source": "computed",
1503_{\ 65}^{\ 64}
                  "function": "extract_function_steps",
                  "args": {
1504 66
1505 <sup>67</sup> <sub>68</sub>
                     "problem_data": "raw_problem_data"
1506 69
1507_{71}^{70}
                "hyperparams": {
                  "source": "computed",
                  "function": "extract_hyperparams",
1508 72
1509 <sup>73</sup><sub>74</sub>
                  "args": {
                     "problem_data": "raw_problem_data"
1510 75
1511 <sup>76</sup>
                "initial_value": {
      78
                  "source": "computed",
```

```
"function": "extract_initial_value",
1513 80
                "args": {
1514 81
                  "problem_data": "raw_problem_data"
1515 <sup>82</sup> <sub>83</sub>
1516 84
1517 85
1517 86
                "source": "computed",
                "function": "get_n_parameter",
1518 87
                "args": {
                  "n": 2
1519 89
               }
1520 90
1521 91
1521 92
              'raw_problem_data": {
                "source": "data",
                "path": "$",
1522 93
     94
                "required": true
1523 95
1524 96
     97
            "prompt": {
1525 98
              " modifiable": true.
1526 99
              "sections": {
1527<sup>100</sup>
                "instructions": "You are a precise arithmetic calculator focused on careful validation
                     of each calculation step. Follow ALL validation requirements exactly and show every
                      intermediate value and check.",
1528
1529
                \hbox{\tt "task\_description": "Calculate the value at position \{inputs.n\} using rigorous}
                     validation. Start with initial_value={inputs.initial_value}. Record and verify ALL
                intermediate values.",
"validation_requirements": "Large Number Validation:\n{inputs.max_value_validator}\n\
1530
1531 102
                     nValue Chain Tracking:\n{inputs.value_tracker}\n\nChecksum Validation:\n{inputs.
1532
                     \verb"nQuotient Checking:\\ \verb"n{inputs.quotient\_checker}\\ \verb"n|nRange Validation:\\ \verb"n{inputs.quotient\_checker}\\ \verb"n".
1533
                     range_validator}",
1534103
                "function_steps": "Calculation Steps:\n{inputs.function_steps}",
"output_format": "Respond with JSON:\n{\n \"reasoning\": \"Show each arithmetic step
1535
                     with ALL intermediate values, validation checks, checksums, and cross-validations
                     \",\n \"answer\": \"final numeric result as string\"\n\"
1536
    105
1537106
              "template": "### Instructions:\n{prompt.sections.instructions}\n\n### Task:\n{prompt.
                   \n### Hyperparameters:\n{inputs.hyperparams}\n\n### Required Validation Steps:\n{
1539
                   prompt.sections.validation_requirements}\n\n### Calculation Process:\n{prompt.
                   sections.function_steps}\n\n### Output Format:\n{prompt.sections.output_format}\n\
1540
                   nCalculate step by step, showing all validation:"
1541107
1542108
            "output": {
1543<sub>110</sub>
              "_modifiable": false,
             "format": "json",
1544111
              "schema": {
112
1545<sub>113</sub>
                "type": "object",
                "required": [
1546114
                  "reasoning",
1547<sub>116</sub>
                  "answer"
1548117
                "properties": {
1549<sup>118</sup><sub>119</sub>
                  "reasoning": {
                    "type": "string"
1550120
1550<sup>121</sup>
1551<sub>122</sub>
                  "answer": {
                    "type": "string"
1552123
    124
1553125
1554126
              "error_values": {
   "reasoning": "Error occurred during reasoning.",
1555<sub>128</sub>
                "answer": "PARSE_ERROR"
1556129
1557<sub>131</sub>
            "settings": {
              "_modifiable": true,
1558132
133
1559<sub>134</sub>
              "class_name": "Func1SolverAgent",
              "timeout seconds": 300
1560135
1561<sup>136</sup>
1562
1563 <sup>1</sup><sub>2</sub>
         #!/usr/bin/env python3
         """Enhanced Func1 solver with comprehensive arithmetic validation"""
1564 3
         import sys
1565
     5
         import os
     6
```

```
sys.path.append(os.path.join(os.path.dirname(__file__), '..', '..', '..', '..', 'llm_agent'))
1567 8
1568 9
         from llm_agent import LLMAgent
1569 10
11
         class Func1SolverAgent(LLMAgent):
1570 12
              """Enhanced Func1 solver with arithmetic validation."""
1571 13
14
              def extract_implementation(self, problem_data):
1572 15
                   """Extract function implementation from problem data."""
                  return problem_data.get('implementation', '')
1573 17
              def extract_hyperparams(self, problem_data):
1574 18
                   """Extract and format hyperparameters."""
1575_{20}^{19}
                  hyperparams = problem_data.get('hyperparams', {})
1576 21
                  if isinstance(hyperparams, dict):
                       formatted = []
1577 22 23
                       for key, value in hyperparams.items():
                           formatted.append(f"{key}: {value}")
1578 24
                       return "\n".join(formatted)
1579 25
26
                  return str(hyperparams)
1580 27
              def extract_initial_value(self, problem_data):
1581 29
                    ""Extract initial value from problem data."""
                  return str(problem_data.get('initial_value', ''))
158230
1583 31
32
              def get_n_parameter(self, n):
                    ""Return the position parameter."""
158433
1585 34
35
                  return str(n)
              def max_value_validator(self, problem_data):
1586 36
1587 37
38
                    ""Generate validation steps for large number arithmetic."""
                  steps = [
                       s = [
"Large Number Validation:",
"'-lightion A B:",
1588 39
                       "1. For multiplication A
     40
1589 41
                           a. Split into single digits: A =
                                                                  a a ... a , B = b b ... b ",
1590 42
                           b. Calculate partial products:",
                              - For each digit pair (a , b ): p
- Track carries explicitly",
1591 43
1591 44
159245
                           c. Sum partial products with position tracking",
1593 46
47
                           d. Verify: result = sum of all partial products",
                       "2. Cross-validate using reverse calculation:",
159448
                           - Split result into chunks",
1595 <sup>49</sup><sub>50</sub>
                           - Verify chunks match partial products"
1596 51
                  return "\n".join(steps)
1597 <sup>52</sup> <sub>53</sub>
              def value_tracker(self, problem_data):
1598 54
                   """Generate value tracking steps for calculation chain integrity."""
1599 <sub>56</sub>
                  steps = [
                       "Value Chain Tracking:",
                       "1. For each position i in calculation:",
1600 57
                           a. Record value[i] = current value",
b. Record squared[i] = value[i] ",
1601 58
1601 59
                           c. Record quotient[i] = squared[i]
160260
                                                                      modulus",
                           d. Record remainder[i] = value[i+1]",
1603 <sup>61</sup><sub>62</sub>
                       "2. Verify chain integrity:",
160463
                          - remainder[i] = squared[i] - (quotient[i] modulus)",
1605 <sup>64</sup> <sub>65</sub>
                                    remainder[i] < modulus",
                           - value[i+1] = remainder[i]"
160666
                  return "\n".join(steps)
1607 <sup>67</sup> <sub>68</sub>
1608 69
              def checksum generator(self, problem data):
1609 <sup>70</sup><sub>71</sub>
                   """Generate checksum validation steps between calculations."""
                  steps = [
161072
                       "Checksum Validation:",
                       "1. For each step i:",
1611 <sup>73</sup><sub>74</sub>
                           a. Calculate checksum[i] = (value[i]
                                                                        31 + squared[i]) mod 997",
                           b. Verify: checksum[i] matches independent calculation",
161275
                       "2. Cross-step validation:",
1613 76
                           - Verify: checksum[i+1] = (remainder[i]
                                                                            31) mod 997",
                           - Track checksum chain for each position"
161478
1615 %
                  return "\n".join(steps)
161681
1617<sup>82</sup><sub>83</sub>
              def division_validator(self, problem_data):
                   """Generate enhanced division validation steps."""
161884
                  modulus = problem_data.get('hyperparams', {}).get('modulus', 0)
1619 85
86
                       f"Division Validation Steps:",
    87
                       f"1. For dividing number A by {modulus}:",
```

```
1620 88
                               a. Calculate quotient Q = A {modulus}",
b. Verify: Q {modulus} A",
1621 89
                               b. Verify: Q {modulus} A", c. Verify: (Q + 1) {modulus} > A",
1622<sup>90</sup>
                               d. Record Q for remainder calculation",
1623 92
                               e. Track sub-steps:",
                                  - Record partial quotients",
1624 93
94
                                   - Verify digit-by-digit multiplication",
1625 95
                                   - Cross-validate remainder"
1626 <sup>96</sup>
                    return "\n".join(steps)
1627 98
               def quotient_checker(self, problem_data):
1628 99
100
                      ""Generate enhanced quotient checking steps."""
1629101
                    steps = [
1630<sup>102</sup>
                          "Quotient Verification:",
    103
                         "1. Record squared value (S)",
1631104
                         "2. Calculate and verify quotient (Q):",

" -Q = S modulus (integer division)",

" - Verify: Q modulus S",
1632 105
    106
1633107
                             - Verify: (Q + 1) modulus > S",
                              - Perform digit-by-digit multiplication check",
1634^{108}
                          "3. Calculate remainder:",
    109
1635110
                             R = S - (Q modulus)",
1636<sup>111</sup>
                          "4. Verify remainder is correct:",
                          " - 0 R < modulus",
" - Cross-validate with value chain"</pre>
    112
1637113
1638<sup>114</sup>
                    return "\n".join(steps)
    115
163916
1640<sup>117</sup>
               def range_validator(self, problem_data):
                      ""Generate enhanced range validation steps."""
    118
1641119
                    modulus = problem_data.get('hyperparams', {}).get('modulus', 0)
1642<sup>120</sup>
                    steps = [
                         f"Range Validation:",
1643122
                         f"1. For each position i:",
                         f" - Verify value[i] 0",
f" - Verify value[i] < {modulus}",</pre>
1644123
    124
1645125
                         f"2. For squared values at position i:",
                         f" - Record exact value squared[i]",
f" - Verify squared[i] = value[i] value[i]",
f" - Verify using digit-by-digit multiplication",
1646^{126}
     127
1647128
1648<sup>129</sup>
                         f"3. Position-specific bounds:",
                         f" - Track maximum possible value at each position",
f" - Verify values stay within position bounds"
    130
1649131
1650<sup>132</sup>
                    return "\n".join(steps)
1651134
               def extract_function_steps(self, problem_data):
1652135
                     """Extract and format key arithmetic operations."""
1653137
1654<sup>138</sup>
                         "1. Initialize calculation chain:",
                         **
                            - Set value[0] = initial_value",
1655140
                              - Calculate checksum[0]",
1656<sup>141</sup>
                         "2. For each position i:",
                            a. Square Operation:",
1657143
                                  - Calculate value[i]
                                                              using digit-by-digit multiplication",
1658<sup>144</sup>
145
                                  - Record squared[i] and partial products",
                                  - Verify checksum[i]",
1659146
                              b. Division and Modulo:",
1660<sup>147</sup>
                                 - Calculate quotient[i] with validation",
                                  - Calculate remainder[i] = value[i+1]",
1661149
                                  - Verify chain integrity",
                              c. Cross-validate:",
1662150
     151
                                  - Verify all checksums match",
1663152
                                  - Confirm value chain consistency",
1664^{153}
                                  - Validate position-specific bounds"
    154
1665155
                    return "\n".join(steps)
1666
```

## K ACCURACY-DRIVEN IMPROVEMENT LOOPS

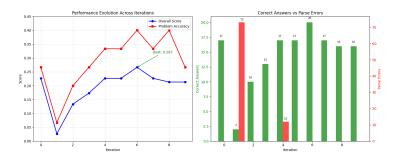


Figure 2: Aime 2024 improvement loop using Sonnet 3.5 v2 as meta agents.

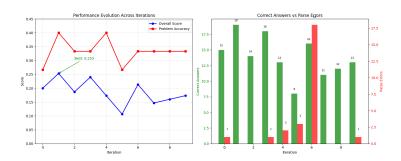


Figure 3: Aime 2024 improvement loop using OpenAI 120b reasoning model as meta agents.

# L AGENT GENERATED BY ACCURACY+COST OPTIMIZATION LOOP BY GPT-OSS-120B

1728

1729

```
1732 1 "description": "AIME 2024 ICL task agent \u2013 focused rhombus example, explicit bound,
                             streamlined prompt",
1733 <sub>2</sub>
                        "input_schema": "aime2024_icl_solver_input_schema.json",
1734 3
                        "model": {
                             "_modifiable": true,
1735 <sup>-</sup><sub>5</sub>
                             "name": "anthropic.claude-3-5-sonnet-20241022-v2:0",
                             "parameters": {
  "temperature": 0.2,
1736 6
1737 8
                                 "max_tokens": 500
1738 9
1739 \frac{10}{11}
                       },
"inputs": {
    difial
                             "_modifiable": true,
1740 12
                            "problem_text": {
    "source": "computed",
1741 \frac{13}{14}
                                 "function": "extract_problem_text",
1742 15
                                 "args": {
1743 \frac{16}{17}
                                     "problem_data": "raw_problem_data"
1744 18
                                }
1745 \frac{19}{20}
                             "rhombus_property": {
 1746 21
                                "source": "computed",
"function": "rhombus_property",
1747_{23}^{22}
                                 "args": {}
1748 24
1749 <sup>25</sup> <sub>26</sub>
                            "solve_perpendicular_condition": {
                                "source": "computed",
"function": "solve_perpendicular_condition",
1750 27
1751 <sup>28</sup> <sub>29</sub>
                                 "args": {
                                     "problem_data": "raw_problem_data"
1752 30
                                }
1753 31
32
                             "bd2_bound": {
1754 33
                                "source": "computed",
1755 34 35
                                 "function": "bd2_bound",
                                 "args": {}
1756 36
1757 <sup>37</sup> <sub>38</sub>
                             "rhombus_icl_example": {
                                "source": "computed",
"function": "rhombus_icl_example",
1758 39
1759 40
41
                                 "args": {}
1760 42
                             "guard_text": {
1760 12
1761 43
44
                                 "source": "computed",
                                 "function": "guard_text",
1762 45
                                 "args": {}
1763 46
                             "raw_problem_data": {
1764 48
                                 "source": "data",
1765 <sup>49</sup> <sub>50</sub>
                                 "path": "$",
                                 "required": true
 1766 51
1760 51
1767 52
53
                        "prompt": {
 1768 54
                               _modifiable": true,
1769 55
56
                             "sections": {
                                 "instructions": "You are a mathematical problem solver. Solve the given AIME problem and
1770
                                           output the answer as a three\u2011digit integer.",
1771 57
                                  "output_format": "You must respond with a JSON object containing exactly two fields:\n-
                                            \"reasoning\": string with your detailed solution\n- \"answer\": string with a
1772
                                            three\u2011digit answer (pad with leading zeros if needed)\n\nFormat: {\"reasoning
                                           \": \"Your detailed solution here\", \"answer\": \"XXX\"}\n\nExample:\n{\"reasoning \": \"To solve..., I ...\", \"answer\": \"042\"}"
1773
1774 58
1775 59
                             "template": "### Instructions:\n{prompt.sections.instructions}\n\n### Geometry hint:\n{
                                       inputs.rhombus_property}\n\### Algebraic bound hint:\n\{inputs.
                                       1776
                                       u00b2 is {inputs.bd2_bound}.\n\### Example to follow:\n{inputs.rhombus_icl_example}\
1777
                                       \verb|n|n### Problem: \\ | n = \frac{1}{n} . \\ | n = \frac{
1778
                                       output_format}\n\n{inputs.guard_text}"
1779 60
61
                        "output": {
 1780 62
                             "_modifiable": false,
                            "format": "json",
1781 <sup>63</sup> <sub>64</sub>
                            "schema": {
           65
                               "type": "object",
```

```
1782 66
                   "required": [
1783 67
                      "reasoning",
1784 68
                     "answer"
1785 70
                   "properties": {
                      "reasoning": {
    "type": "string"
1786 \, {71 \atop 72}
1787 73
1788 <sup>74</sup> <sub>75</sub>
                     "answer": {
    "type": "string",
1789 76
                        "pattern": "^[0-9]{3}$"
1790 <sup>77</sup> <sub>78</sub>
                  }
1791 79
                },
"error_values": {
1792 <sup>80</sup> <sub>81</sub>
                   "reasoning": "Error occurred during reasoning.",
"answer": "PARSE_ERROR"
1793 82
1794 <sup>83</sup> <sub>84</sub>
1795 85
              "settings": {
                "_modifiable": false,
"class_name": "AIME2024ICLSolverAgent",
1796 \frac{86}{87}
1797 88
                "timeout_seconds": 300,
"exit_on_parse_failure": false
1798 \frac{89}{90}
1799 91
1800
1801
           from llm_agent import LLMAgent
1802 <sup>2</sup>
           from typing import Dict, Any
1803 4
           class AIME2024TaskAgent(LLMAgent):
1804 5
                AIME2024-specific reasoning + answer agent.
1805 7
1806 <sup>8</sup>
                Processes AIME mathematical competition problems and generates detailed step-by-step
                reasoning followed by answers.
1807 g
1808^{\,10}
1809 11
1809 12
                def __init__(self, name="AIME2024_Task", config_path=None, dry_run=False):
                     super().__init__(name, config_path, dry_run)
1810^{\,13}
                def extract_problem_text(self, problem_data: Dict[str, Any]) -> str:
1811 15
1812 <sup>16</sup> 17
                     Extract problem text from AIME2024 data format.
```

 $1814_{20}^{19}$ 

1815 21

1816<sup>22</sup>

 $1817_{24}^{23}$ 

AIME2024 format:

"id": int,

"url": str

"problem": str, # Main problem statement

return str(problem\_data.get('problem', ''))