# BUILDING LEARNING CONTEXT FOR AUTONOMOUS AGENTS THROUGH GENERATIVE OPTIMIZATION

**Anonymous authors**Paper under double-blind review

000

001

002003004

010 011

012

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

031

033

035

037

038

040

041 042

043

044

046

047

048

049

051

052

#### **ABSTRACT**

Building intelligent agents that learn involves designing systems that can evolve their behavior based on experiences. While early approaches to large language models (LLMs) agent learning relied mostly on structured memory and in-context learning, they often led to behavioral instability, poor interpretability, and difficulty in control. Recent success in generative optimization, where an LLM is used as an optimizer, has shown the possibility of creating autonomous software agents. By separating behavior logic (workflow) and how that logic is updated (optimizer), the agent designer can exhibit more control over the agent. In this work, we show the surprising fact that the agent learning problem is *under-specified* with the generative optimization framework. If we want an agent to learn the right behavior, we must set up the right context that will induce such behavior. We investigate three types of software engineering problems that span data science, computer security, game playing, and question answering. We show that the original generative optimization framework can only learn robustly under one of the three settings. To address the issue, we propose to construct a *meta*-graph through templates to introduce the right learning context to an LLM optimizer. With this addition, we demonstrate that defining the right learning context enables agents to discover behaviors aligned with the designer's objectives. In particular, we show the first known result of using generative optimizers to learn executable programs that play Atari games, where the resulting agents achieve performance comparable to deep reinforcement learning while requiring 50%-90% less training time.

#### 1 Introduction

Artificial intelligent *agents* — computer programs that "operate autonomously, perceive their environment, persist over time, adapt to changes, create and pursue goals" (Russell & Norvig, 2016) — have regained significant attention recently, due to the maturation of large language models (LLMs) (Achiam et al., 2023). The ease of accessing LLMs gave rise to new programming paradigms, such as language model programs (Khattab et al., 2024) and multi-agent orchestration frameworks (Wu et al., 2024), all of which leverage calls to LLMs to handle a wide range of tasks in human society, from computer use (Fourney et al., 2024), and software engineering (Jimenez et al., 2024), to scientific discovery (Yang et al., 2024b).

Despite these advances, building agents still requires a substantial amount of human engineering. Often agent developers need to design complex decision rules to orchestrate an agent's behaviors, build pipelines to parse information from the environment into the agent's percepts, and engineer prompts to control LLMs. However, results in competitive programming (El-Kishky et al., 2025) have shown that it is paramount to find a general-purpose method that can scale with compute, rather than engineering domain-specific solutions (Sutton, 2019). These human-engineered components should be automatically *learned* through agent's experiences, enabling agents that can program themselves through learning.

Traditional techniques such as reinforcement learning (RL) (Sutton et al., 1998) are general-purpose algorithms that theoretically can be applied to agent learning. However, since state-of-the-art agents are made of large language models (LLMs), gradient updates must be performed on LLM weights. Initial explorations showed promising results (Bai et al., 2025; Guo et al., 2025), although it is unclear how well the learned LLM weights generalize to out-of-domain tasks, and how many tasks are needed

Figure 1: A Diagram for An LLM-based System. represents a repeated workflow execution graph (we denote as *workflow graph*). We use and to represent trainable parameter nodes (string, code, etc.). The optimizer updates parameter nodes during learning. We show that as agent designer, we can choose to optimize the parameters under different learning contexts (interactive, batch, and episodic). We show how to leverage repeatable *workflow graphs* (c) can be concatenated to construct an *agent learning graph* (c).

to reach a *generalist* agent, as multi-task learning has always been a challenge in RL (Kirkpatrick et al., 2017; Taiga et al., 2023).

Recently, there is an emerging end-to-end perspective that instead treats a given software program as a computational graph and optimizes its parameters (i.e., the human-made decisions discussed above) via a generalized form of "back-propagation gradient descent" (Cheng et al., 2024b; Yuksekgonul et al., 2025; Wang et al., 2024b). By using generative models (like LLMs) as optimizers (Pryzant et al., 2023; Nie et al., 2023; Yang et al., 2024a), this perspective has led to frameworks that can automatically tune parameters and generate new code in non-differentiable programs, which achieved promising results in producing distributed systems programs (Wei et al., 2024).

This approach offers several advantages. It enables an agent to adapt to solve tasks through automatic optimization, directly using rich feedback (such as compilation errors, system reports, and end-user feedback), a process that traditionally requires manual trial and error. Moreover, by learning, agents can discover solutions that exceed what human experts can design. We have seen this trend with deep learning (Silver et al., 2016; Brown & Sandholm, 2019; Mirhoseini et al., 2020; Bellemare et al., 2020). This perspective turns agent learning into a closed-loop optimization problem by calling LLMs over multiple iterations to refine the agent. This is in contrast to the predominant use of LLM-as-agents when the agents change their behaviors based on each query and hope to produce the correct sequence of actions in one-shot.

Prior works that incorporate environment feedback to change an agent's behavior have met with mixed success. While there are very successful applications (Cheng et al., 2024b), limitations have also been observed: the optimization process can be unstable (Huang et al., 2024) and the self-improvement phenomenon only persists for a few rounds (Shinn et al., 2023; Madaan et al., 2023). We argue that this issue arises in part because the problem of agent learning is *under-specified* with generative optimization. An agent needs to learn solutions that can generalize different contexts, while generative optimization defines an optimization problem under a single problem context.

In this paper, we analyze sources of this under-specification issue and propose constructive remedies. We show that agent learning should not be specified only as an execution (workflow) graph of its own internal operations, but a *meta*-graph on a stream of experiences to capture the learning context. We introduce operators ( $\oplus$  and  $\Rightarrow$ ) to insert workflow graphs into a template to construct an agent learning graph, which correctly specifies the agent's learning objective and enables generative optimization to learn parameters effective for the agent designer's goal. In addition, we discuss how to structure the agent's internal workflow to improve optimization results (similar to how architectural choices in neural networks facilitate better learning outcomes). We note that this factor has been overlooked in previous attempts to design self-improvement loops (Chen et al., 2023; Huang et al., 2024; Snell et al., 2024). Finally, we discuss a few choices of enhancing and amplifying feedback for different stages of the learning process, analogous to reward shaping.

We show that these insights allow us to apply generative optimization to solve a wide range of tasks. Automatic software engineering, such as creating an agent to write machine learning programs, can be seen as an interactive learning task. We show that on the MLAgentBench (Huang et al., 2023), we can learn an agent that can output high-quality models that surpassed 86.6% submissions on Kaggle leaderboard than the baseline agent, which only surpassed 70.8% submissions. We can also optimize an LLM based workflow to improve its performance by as much as 14.5% on GSM8K and 65.1% on BBEH. Finally, we show that we can learn a static Python program that can play Atari games, nearly

matching the performance of Deep RL baselines but with 50%-90% less compute time. All of these show the versatility of this paradigm and the power of the automatic agent learning process.

#### 2 BACKGROUND AND RELATED WORK

**History of Learning Agents** The term AI agent has a long history (Genesereth & Nilsson, 1987). In this paper, we follow Russell & Norvig (2016) and define a learning agent as a program that can sense percepts, take actions, and adapt with experiences in a digital/physical environment. For an agent to learn, it implies that the software has components that are modifiable and can influence its behaviors; these components are called *parameters*. For example, a tabular agent has a lookup table as its parameter, and a learning algorithm such as policy iteration or value iteration would be suitable (Bertsekas, 1987). A deep RL agent has a neural network as its parameter, and learn from rewards using algorithms like proximal policy optimization (PPO) (Schulman et al., 2017).

**Adaptive Workflow** Increasingly, intelligent systems are being built with LLMs. For the state-of-the-art LLM systems (Wang et al., 2024a;c; Fourney et al., 2024), their parameters can be model weights, or more generally, system prompts, code that pre-processes input, and code that modifies the returned results from the LLM. Besides the obvious approach of fine-tuning the LLM's weights (Scheurer et al., 2023), there isn't a dominant approach on how to change the system's behavior on the fly. Some inference-time learning methods have been introduced, with the prevailing strategy utilizing databases, referred to as "memories" in RAG (Lewis et al., 2020). Recently, a new perspective of building intelligent agent emerges, which leverages an LLM's ability to write coherent programs to accomplish a purpose (Cheng et al., 2024b; Zhang et al., 2024). This view separates an LLM agent into two parts: the workflow that represents the behavioral logic of the agent, and an optimizer that updates such behavioral logic.

Generative Optimization Generative optimization algorithms have been proposed to update an LLM workflow. They typically use a generative model (like an LLM) as part of its optimizer to analyze problems and propose updates. A generative optimizer takes as inputs (1) a problem context, (2) parameters, (3) a computational graph involving the parameters, (4) a feedback signal, and outputs a parameter value. Several generative optimizer implementations have been proposed, such as DSPy (Khattab et al., 2024), OptoPrime (Cheng et al., 2024b), TextGrad (Yuksekgonul et al., 2025), and GASO (Wang et al., 2024b). They differ in how they represent and reason about the graph and kinds of feedback they can process. For instance, optimizers in DSPy work with scalar feedback, while OptoPrime/TextGrad/GASO uses any feedback that an LLM can interpret. OptoPrime formats the entire graph into a single LLM prompt, while TextGrad/GASO processes the graph iteratively.

The Framework of OPTO Recently OPTO (Optimization with Trace Oracle) (Cheng et al., 2024b) was proposed as a unified math setup for describing iterative generative optimization problems. An OPTO problem (a generalization of numerical optimization) is described by a tuple  $(\Theta, \omega, \mathcal{T})$ , where  $\Theta$  is the parameter space,  $\omega$  is the problem context and  $\mathcal{T}$  is a Trace Oracle. For a parameter  $\theta \in \Theta$ , the Trace Oracle  $\mathcal{T}$  returns a tuple (f,g) where g is a computational graph involving  $\theta$  and f is a feedback signal provided to exactly one node of g (the output node). An autonomous agent that learns through experience in this setup corresponds to a workflow design and an optimizer that can update the workflow. We emphasize that a workflow itself is not an autonomous agent, but a workflow combined with an optimizer that can rewrite its own behavior according to feedback is an agent.

## 3 BUILDING LEARNING AGENTS WITH GENERATIVE OPTIMIZATION

In agent learning, we wish to optimize an agent's parameters in a stream of experience. Our approach takes inspiration from deep learning, which accomplishes machine learning via optimization on differentiable computational graphs. In deep learning, we specify a neural network architecture (a computational graph which is parameterized by tensors) and a numerical oracle (e.g., a loss function to minimize) to provide feedback at the output of the computational graph. Following the OPTO framework, our approach is built similarly with these two components but using generative optimization. The main differences are that here differentiability is not required and that the agent is not limited to learning from numerical feedback only. In the following, we show constructive

templates to define the parameterized computational graph and discuss design principles for the feedback oracle and the agent's computational graph.

We demonstrate how computational graphs can naturally describe different agent learning problems. We suppose that a workflow is given and is represented as a computational graph  $W_{\theta}$ , where  $\theta$  denotes the parameters. Without loss of generality, we suppose the workflow takes a single input x and returns a single output  $y = W_{\theta}(x)$  (x, y can be arrays for modeling multi-input-multi-output cases). We call  $W_{\theta}$  the workflow graph. Recall that the workflow here means the full software program, which internally may be composed of multiple calls to LLMs and decision rules (in other words, with abuse of notation, the workflow here can represent also an entire multi-agent orchestration (Wu et al., 2024)). As a result, the workflow graph does not need to be static and it can vary with input or due to the internal randomness of the workflow.

In contrast to the workflow, we denote the full computational graph that will be presented to the generative optimizer (i.e., g in OPTO) as  $A_{\theta}$ , which we call the agent learning graph. The agent learning graph is composed of the workflow graph and other nodes derived from the learning problem, such that the learning structure can be captured correctly. Lastly, a feedback oracle maps (x,y) into feedback f, which can be numerics, texts, images, or structured objects (e.g. a dictionary). We assume the feedback is not adversarial and contains some information of the agent's performance. With these assumptions, an OPTO problem instance can be created where the problem context  $\omega$  can be fixed to a string such as "Update the parameters to incorporate the feedback."

Now we discuss how to construct the agent graph for common agent learning problems using templates to build *agent learning graphs*, as shown in Fig. 1.

Interactive Learning Template. Here the agent learns on the fly as it interacts with the world (Shalev-Shwartz et al., 2012). At each time step, it sees an input x, outputs y, receives feedback f, and then updates its parameter  $\theta$ . These problems encompass online learning and bandit variants that are well-studied in the literature. The meta graph here simply shows the input x is transformed by an operator (i.e., the agent) and then the feedback is provided to y. When the agent graph  $A_{\theta}$  is inserted into this template, it yields the workflow graph  $W_{\theta}$ .

Batch Learning Template. Different from interactive learning, a batch-learning agent learns from a given dataset  $D = \left[ (x_i, z_i) \right]_{i=1}^N$  of size N, where  $x_i$  and  $z_i$  denote the input and the information to learn from for the ith data point (Hastie et al., 2009). For example,  $z_i$  can be the desired agent output when seeing  $x_i$  (supervised learning), or it can be a positive-negative pair (preference-based learning), etc. In batch learning, the agent learns from an oracle that takes  $(x_i, y_i, z_i)$  as input (where  $y_i = W_{\theta}(x_i)$ ) and provides feedback such as a loss. To handle such a batch problem with the iterative setup of OPTO, we appeal to the idea of online-to-batch conversion (Shalev-Shwartz et al., 2012) and mini-batching. As shown in Fig. 1, in each iteration of OPTO, we sample a minibatch and construct the graph for the sampled batch. We introduce a batchify operator  $\oplus$  that concatenates different inputs. For a minibatch  $\left[ (x_i, z_i) \right]_{i=1}^B$  of size B, we first obtain  $\left[ (f_i, g_i) \right]_{i=1}^B$ , where  $g_i$  and  $f_i$  result from input  $x_i$ . Suppose  $o_i$  is the output node of  $g_i$ . Then we concatenate the outputs from  $\{g_i\}_{i=1}^B$  to create a new node  $\hat{o} = \bigoplus_{i=1}^B o_i$  and give the concatenated feedback  $\bigoplus_{i=1}^B f_i$  to  $\hat{o}$ .

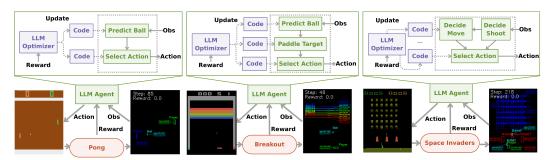
**Episodic Learning Template.** Here we adopt a broader definition of RL, which describes the agent learning in a sequential decision process with feedback of reward signals (Sutton et al., 1998) or richer signals like natural language (Cheng et al., 2024a; Chen et al., 2024). We consider an episodic setting. In each iteration, the agent interacts with the environment for multiple steps, receives feedback for each step (the feedback can be empty), and then updates its parameters at the end of the episode. To represent this structure as a computational graph, first we describe the interaction process of how observations and actions are generated. In Fig. 1, this is shown as a chain similar to a Markov decision process; notice there is an arrow going from action to the next observation via an operator denoted as  $\Rightarrow$ , which captures the causality. Then we apply the batchify operator  $\oplus$  on the observations generated  $\hat{o} = \bigoplus_{i=1}^T o_i$ , where T is the episode length, and similarly concatenates the feedback.

**Remark.** Using the right meta graph for a learning setup is important as it provides the learning context to the generative optimizer; otherwise objective misalignment can happen. For instance, if we desire a batch learning solution (i.e. a parameter that works well across a dataset of examples) but use the meta-graph for online learning in OPTO, we can get sub-par optimization results (for instance, an unstable parameter that is sensitive to the order in which individual examples are presented).

Similarly, if an agent's behavior has long-term consequences, then we should only change its behavior logic after an episode terminates. The separation between behavior logic and when/how to change them allows us to specify the right learning objective and allow the agent to learn the right behavior.

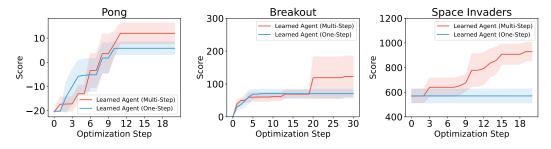
## 4 ATARI GAME PLAYING AGENTS THROUGH EPISODIC LEARNING

Game playing has been a central focus in reinforcement learning (Mnih et al., 2013; Silver et al., 2016; Brown & Sandholm, 2019). Recently, LLMs have demonstrated abilities to play long-horizon games such as Pokémon Blue (Karten et al., 2025; Anthropic, 2025). However, all of these successes utilize direct weight updates for the neural network, through training on collected in-game experience or massive pre-training on tutorials and forum posts. In this section, we want to demonstrate that, shockingly, with LLM as the optimizer and a correct learning template, we can learn a python program (not weights) that can play games that were typically mastered by neural networks.



**Figure 2:** We show the workflow design of different decision-making program components for each Atari game agent. The LLM agent receives an object-centric dictionary of information of the game state and uses Python code to process and output an action.

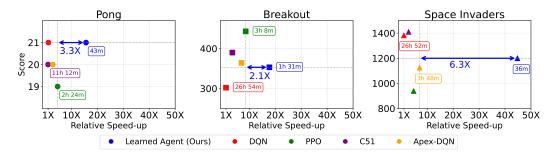
The Arcade Learning Environment (ALE) of Atari games has remained an important benchmark for evaluating RL algorithms for training neural network-based policies (Mnih et al., 2013). ALE can be used to evaluate an RL algorithm in several ways: 1) The algorithm's learning efficiency both in terms of number of interactions with the environment and the overall wallclock time (Hessel et al., 2018); 2) The diverse set of environments allow the evaluation of generalization of learning (Lee et al., 2022).



**Figure 3:** Performance of agent under different learning graphs (one-step vs multi-steps) across 5 trials. Episodic learning template concatenates workflow graph at each step to build an episodic learning graph that correctly specifies the temporal dependency in the agent's learning objective.

We use object-centric Atari Environments (OCAtari) (Delfosse et al., 2024) to parse the pixel-based observation from ALE to object-based representation. OCAtari provides the coordinates, size and velocity of the object on screen, game termination condition ("lives"), and current reward (see Figure A.9). We do not perform additional transformations to make the observation more readable.

**Workflow Design.** The agent is designed slightly differently for each game. The design decision is driven by a high-level modularization of the decision-making process. Both Pong and Breakout agents have select\_action as the final component. They use predict\_ball\_trajectory as an intermediate step, where the prediction is provided to select\_action to decide how the paddle can



**Figure 4:** We show the relative speedup running the generative optimization process compared to traditional RL methods. Learned Agent result is the highest score it achieved in 5 trials. We report RL results from an open-source implementation of RL algorithms (Huang et al., 2022b) and publicly available experiment logs (Huang et al., 2022a). RL algorithms are trained with 8 parallel environment instances. Note that most recent Deep RL with 32 environment instances can earn score of 450 on Breakout in 33m, see Appendix F.5. The difference of scores (19-21) in Pong is not a meaningful difference and caused mostly by rounding.

be moved. For Breakout, we introduced a goal prediction component (generate\_paddle\_target) to strategically determine where the paddle should go to maximize the reward. For Space Invaders Agent, we simply have decide\_shoot and decide\_movement to decompose the decision space of choosing when to fire and when to move the game avatar.

**Learning Graph Design.** Due to the length of the context window for the LLM we use, we are only able to trace a fixed number of temporal steps. The number of steps is determined by the token spent representing the observation, the complexity of the agent design for the game, and the length of the solution. The training rollout is 300 steps for Breakout, 400 steps for Pong, and 10 steps for Space Invaders.

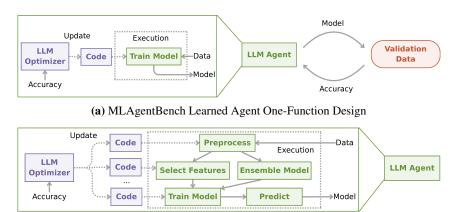
**Feedback Design.** We notice that only providing feedback based on the reward in the training rollout leads to performance plateaus, particularly in games where the game mechanism changes based on player progress. For example, in Breakout, the higher-value bricks in the upper rows deflect the ball at greater speeds, creating a distribution shift between the training context (primarily lower bricks) and the evaluation context (including higher bricks). This observation inspires two feedback design choices: 1) we provide staged feedback to instruct the model to pay attention to different game mechanisms or share high-level winning strategies; 2) we evaluate the performance of the agent with longer rollouts (up to 4000 steps) and use that reward as feedback to the generative optimizer.

**Results.** We find that even with sparse representation of game states and rewards in the form of trajectories, LLM optimizer (OptoPrime) demonstrate remarkable ability to infer game mechanics and environmental constraints from traced trajectories. While our approach provides docstrings that describe high-level game objectives and mechanics, we experiment with deliberately omitting specific implementation details like exact boundary coordinates or collision physics. Despite this, OptoPrime consistently infers these crucial details through analysis of the trajectory data. For example, in Breakout (see an example observation in Figure A.9), OptoPrime identifies the exact positions of the left wall (x=9) and right wall (x=152) by observing ball position and velocity changes across multiple steps. It correctly implements ball physics calculations including bounce mechanics without being explicitly told these details. This emergent understanding of game physics and boundaries demonstrates the LLM's ability to perform causal inference from sequential observations.

## 5 DATA SCIENCE AGENT WITH INTERACTIVE LEARNING GRAPH

The interactive learning setup can describe the learning objective for the majority of LLM agent benchmarks. The hallmark of these benchmarks is that even though an LLM agent needs to take multiple intermediate steps to complete a task, such step *does not* cause state transition in the environment that changes the reward the agent would receive. Even for benchmark that requires the agent to carry out multiple actions to successfully complete a task, the intermediate actions do not cause an internal, stateful change in the environment – the reward is often only associated with the final output of the agent (such as a customer response or an executable code).

MLAgentBench is a benchmark specifically designed to measure the effectiveness of machine learning agents in automating ML experimentation processes (Huang et al., 2023). There are different designs of the ML agents. The majority of agents created to solve this task have a primitive self-improvement loop, where the agent simply looks at its previous output and self-refine (Huang et al., 2023; Wang et al., 2024c; Chan et al., 2024). All tasks in MLAgentBench involves training a machine learning model, figuring out preprocessing data, feature selection, choosing hyperparameters of the model, and deciding training details. We chose two tabular tasks for the purpose of easy experimentation, as they consume the least amount of compute resources compared to large datasets.



(b) MLAgentBench Learned Agent Many-Function Design

**Figure 5: Agent Design for MLAgentBench**. We can design different agent workflow graphs for the learned agent to solve the task of training a machine learning model given a dataset.

**Workflow Design.** We design our agent to have specific components that are changable by the optimizer. That is, to set up a generative optimization process that is automatic, only with human engineering focused on the initial configurations. To design the agent's internal operation, we experiment with two different kinds of workflow design to highlight the influence of workflow design on the optimization outcome. One design asks the optimizer to program *one* code block that does everything, labeled as "One" in Table 1. The other design properly decomposes the model training tasks into five steps: preprocess, select\_features, create\_ensemble\_model, train\_model, and predict, labeled as "Many" in Table 1. We illustrate these two choices in Figure 5.

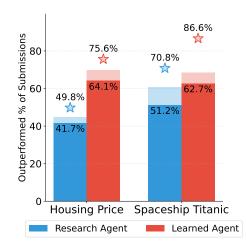
**Learning Design.** We use the OptoPrime as the generative optimizer (Cheng et al., 2024b). The entire graph is represented in the LLM context window and the LLM perform parameter update for all parameters all at once. We perform a train-validation split on the dataset to create a validation partition and use the task-specific metric on the validation dataset as the optimization objective (i.e., *maximize accuracy* or *minimize error*). We use the final learned agent's machine learning model to produce predictions on the hidden test set and submit to Kaggle website to compare against hidden ground truth. We *do not* use the Kaggle test score as the reward signal in the optimization loop.

**Feedback Design.** We apply fine-grained style feedback to the generative optimizer at different stages of validation accuracy (see Figure A.1). We additionally experimented with improvement style feedback where the model fails to train a machine learning model that has a higher validation accuracy than the previous step, we append an improvement suggestion to the feedback string.

**Results.** To make the comparison fair, we pre-downloaded the datasets for the Research Agent and made sure it could produce a machine learning model with valid test submission files for Kaggle (Huang et al., 2023). We track the average performance of the model produced by the both agents as well as the best result. After 20 optimization steps, we submit the model with the highest validation accuracy to the Kaggle competition to get the test score and leaderboard ranking. On both tasks, the gap between the Research Agent (Huang et al., 2023) and our learned agent is around 11.5%-22.4% on average, and the best machine learning model produced in the learned agent surpasses 86.6% of human submissions. Surprisingly, letting the optimizer continuously updating one function (block of code) is better for Housing Price, but not for Spaceship Titanic. This result highlights the importance of experimenting with different workflow designs.

		Housing Price RMSE (↓)	Spaceship Titanic Accuracy (†)
	Research	Agent (Huang et	al., 2023)
Average	_	0.149	78.17
Best	_	0.145	79.84
	Le	earned Agent (Ou	rs)
Average	One	0.135	79.65
	Many	0.147	79.69
Best	One	0.129	80.00
	Many	0.141	80.43

**Table 1:** MLAgentBench Result. We run both agents 5 times and compute the average and best test score. Single and Many refer to a one-function vs many-functions workflow design for the agent (Figure 5). Both agents use the same underlying LLM (Claude Sonnet-3.5-v2).



**Figure 6:** Leaderboard ranking for the best agent across design choices.

#### 6 LANGUAGE UNDERSTANDING AGENT THROUGH BATCH LEARNING

From document processing to logical deductions, LLM agents are used for general language understanding tasks where the agent designer needs to write a pre/post-processing program as well as instructions/prompts to the LLM API call. The crucial learning context here is to allow the agent to write *one* prompt and a *fixed* program that generalizes to different kinds of questions and tasks. We explore the effect of setting the right learning context in BigBench Extra Hard (BBEH) (Kazemi et al., 2025).

**Agent Design.** We show the workflow design graph in Figure 7. This agent is made of two components, one is a *call llm* function that takes in the task query and a optimizable prompt. The other is an *answer extraction* function that parses the return from the LLM call.

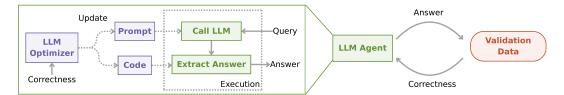


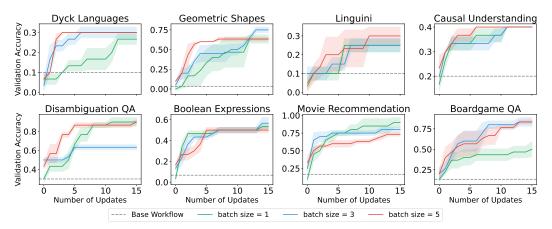
Figure 7: Agent Design for BigBench Extra Hard (BBEH). Note that in here, we only show an interactive learning setup, where the agent graph only contains one query and receives one correctness feedback. In our experiment, we insert the execution graph into a template through batchify  $\oplus$  operator to construct a batch learning graph over multiple queries, answers, and feedbacks.

**Learning Design** We apply the batch learning template to construct the learning graph for the agent. In order to apply this template, we sample a batch of inputs from the dataset  $\mathcal{D}$ . We roll out the workflow graph on each of the input, and by the end, we concatenate all the workflow graphs together to form the batch learning graph. We only use 20 example inputs for learning and the rest are holdout test set for evaluation. Batch learning graph helps agent understand how one shared program and prompt need to adapt to different kinds of inputs.

**Feedback Design** We provide feedback as a list of strings that contain whether the workflow's response for each input is correct or incorrect and revealing the solution to the optimizer when the answer is incorrect.

**Results.** We systematically change the batch learning graph's size (batch size) during the learning phase over 20 examples and measure how well the learned prompt and postprocessing code generalize to the unseen examples in the BBEH set. Surprisingly, even though generally presenting more than one example at a time (Batch Size > 1) gives better results (see Table 2), it seems that different batch sizes lead to different performances for different tasks. We also see different patterns of learning convergence on a 5-example validation set that we selected from the 20 examples (Figure 8). Larger

batch sizes allow the agent to learn faster but also plateau more quickly (Geometric Shapes). This highlights the benefit of constructing learning graphs – it exposes the hyperparameters of learning that agent designers must decide. We additionally report results on GSM-8K with the learned LangGraph agent in Section E and Table A.3.



**Figure 8:** Performance of the learned workflow for each task across 3 trials. Each task starts with the same prompt and answer extraction code. Shaded area shows standard error. The training dataset size is fixed to 15 examples. Validation set has 10 examples.

MiniBatch Size	Dyck Languages	Geometric Shapes	Linguini	Causal Understanding
Un-Optimized	$\textbf{0.114} \pm 0.007$	$\textbf{0.074} \pm 0.003$	$0.183 \pm \textbf{0.010}$	$0.114 \pm 0.005$
Batch Size=1 Batch Size=3 Batch Size=5	$0.183 \pm 0.049$ $0.063 \pm 0.010$ $0.190 \pm 0.031$	$0.343 \pm 0.039$ $0.389 \pm 0.040$ $0.200 \pm 0.099$	$0.149 \pm 0.024$ $0.234 \pm 0.012$ $0.170 \pm 0.030$	$0.375 \pm 0.146$ $0.408 \pm 0.097$ $0.531 \pm 0.018$

**Table 2:** Holdout Test Set Performance for BBEH tasks. Bold indicates best accuracy per column; standard error is shown in smaller gray text. The test dataset excludes examples used for train and val, and usually includes 175 examples. The full table for 8 tasks are in Appendix A.2.

#### 7 CONCLUSION AND LIMITATION

We demonstrate generative optimization on computational graph is a powerful new paradigm for agent learning. We identify common misalignment issues in practice and provide constructive guidelines to address them. With these insights, we demonstrate successful agent learning results on a wide range of problems (GSM8K, BBEH, MLAgentBench, and Atari games) across interactive, batch, and reinforcement learning scenarios. These experimental results push the boundary of problems where generative optimization has been applied in the literature and provide strong evidence that generative optimization can be the key to the next breakthrough of agent learning and an effective method to leverage inference time compute to find optimal solution automatically.

However, we should also highlight that our current results have limitations. Although we solved the objective misalignment with a principled reductionist approach, our current recommendations on agent and feedback design are still heuristic-driven. We also notice that the optimization process can be unstable. The guideline can largely mitigate the issue, but efforts are still required to configure the initial condition and optimization procedure correctly. All of these warrant future research to explore paths to create a fully automated, goal-driven, *generalist* agent.

## REFERENCES

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. GPT-4 technical report. arXiv preprint arXiv:2303.08774, 2023.

- Anthropic. Claude's extended thinking, February 2025. URL https://www.anthropic.com/news/visible-extended-thinking. Accessed: 2025-05-31.
  - Hao Bai, Yifei Zhou, Li Erran Li, Sergey Levine, and Aviral Kumar. Digi-q: Learning q-value functions for training device-control agents. *arXiv preprint arXiv:2502.15760*, 2025.
    - Marc G Bellemare, Salvatore Candido, Pablo Samuel Castro, Jun Gong, Marlos C Machado, Subhodeep Moitra, Sameera S Ponda, and Ziyu Wang. Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*, 588(7836):77–82, 2020.
    - Dimitri P Bertsekas. *Dynamic programming: deterministic and stochastic models*. Prentice-Hall, Inc., 1987.
  - Noam Brown and Tuomas Sandholm. Superhuman ai for multiplayer poker. *Science*, 365(6456): 885–890, 2019.
  - Jun Shern Chan, Neil Chowdhury, Oliver Jaffe, James Aung, Dane Sherburn, Evan Mays, Giulio Starace, Kevin Liu, Leon Maksin, Tejal Patwardhan, et al. Mle-bench: Evaluating machine learning agents on machine learning engineering. *arXiv preprint arXiv:2410.07095*, 2024.
  - Angelica Chen, Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Samuel R Bowman, Kyunghyun Cho, and Ethan Perez. Learning from natural language feedback. *Transactions on Machine Learning Research*, 2024.
  - Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. *arXiv preprint arXiv:2304.05128*, 2023.
  - Ching-An Cheng, Andrey Kolobov, Dipendra Misra, Allen Nie, and Adith Swaminathan. LLF-bench: Benchmark for interactive learning from language feedback. In *ICLR Workshop on Large Language Model (LLM) Agents*, 2024a.
  - Ching-An Cheng, Allen Nie, and Adith Swaminathan. Trace is the next autodiff: Generative optimization with rich feedback, execution traces, and LLMs. In *NeurIPS*, 2024b.
  - Quentin Delfosse, Jannis Blüml, Bjarne Gregori, Sebastian Sztwiertnia, and Kristian Kersting. OCAtari: Object-centric Atari 2600 reinforcement learning environments. *Reinforcement Learning Journal*, 1:400–449, 2024.
  - Ahmed El-Kishky, Alexander Wei, Andre Saraiva, Borys Minaev, Daniel Selsam, David Dohan, Francis Song, Hunter Lightman, Ignasi Clavera, Jakub Pachocki, et al. Competitive programming with large reasoning models. *arXiv preprint arXiv:2502.06807*, 2025.
  - Adam Fourney, Gagan Bansal, Hussein Mozannar, Cheng Tan, Eduardo Salinas, Friederike Niedtner, Grace Proebsting, Griffin Bassman, Jack Gerrits, Jacob Alber, et al. Magentic-one: A generalist multi-agent system for solving complex tasks. *arXiv preprint arXiv:2411.04468*, 2024.
  - Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. Pal: Program-aided language models. In *International Conference on Machine Learning*, pp. 10764–10799. PMLR, 2023a.
  - Luyu Gao, Aman Madaan, Shuyan Zhou, Uri Alon, Pengfei Liu, Yiming Yang, Jamie Callan, and Graham Neubig. PAL: Program-aided Language Models, January 2023b.
  - Michael R Genesereth and Nils J Nilsson. *Logical foundations of artificial intelligence*. Morgan Kaufmann Publishers Inc., 1987.
  - Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*, pp. 249–256, 2010.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
  - Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, 2 edition, 2009.

- Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI*, 2018.
- Dan Horgan, John Quan, David Budden, Gabriel Barth-Maron, Matteo Hessel, Hado van Hasselt, and David Silver. Distributed prioritized experience replay. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=H1Dy---0Z.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. Large language models cannot self-correct reasoning yet. In *ICLR*, 2024.
- Qian Huang, Jian Vora, Percy Liang, and Jure Leskovec. Mlagentbench: Evaluating language agents on machine learning experimentation. *arXiv preprint arXiv:2310.03302*, 2023.
- Shengyi Huang, Rousslan Fernand Julien Dossa, Antonin Raffin, Anssi Kanervisto, and Weixun Wang. The 37 implementation details of proximal policy optimization. In *ICLR Blog Track*, 2022a. URL https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/. https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/.
- Shengyi Huang, Rousslan Fernand Julien Dossa, Chang Ye, Jeff Braga, Dipam Chakraborty, Kinal Mehta, and JoÃGo GM AraÚjo. Cleanrl: High-quality single-file implementations of deep reinforcement learning algorithms. *Journal of Machine Learning Research*, 23(274):1–18, 2022b.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *ICLR*, 2024.
- Seth Karten, Andy Luu Nguyen, and Chi Jin. Pok\'echamp: an expert-level minimax language agent. *arXiv preprint arXiv:2503.04094*, 2025.
- Mehran Kazemi, Bahare Fatemi, Hritik Bansal, John Palowitch, Chrysovalantis Anastasiou, Sanket Vaibhav Mehta, Lalit K Jain, Virginia Aglietti, Disha Jindal, Peter Chen, et al. Big-bench extra hard. *arXiv preprint arXiv:2502.19187*, 2025.
- Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan A, Saiful Haq, Ashutosh Sharma, Thomas T. Joshi, Hanna Moazam, Heather Miller, Matei Zaharia, and Christopher Potts. DSPy: Compiling declarative language model calls into state-of-the-art pipelines. In *ICLR*, 2024.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114 (13):3521–3526, 2017.
- LangChain. Langgraph: A low-level orchestration framework for building controllable agents, 2024. URL https://langchain-ai.github.io/langgraph/.
- Kuang-Huei Lee, Ofir Nachum, Mengjiao Sherry Yang, Lisa Lee, Daniel Freeman, Sergio Guadarrama, Ian Fischer, Winnie Xu, Eric Jang, Henryk Michalewski, et al. Multi-game decision transformers. *Advances in Neural Information Processing Systems*, 35:27921–27936, 2022.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *NeurIPS*, pp. 9459–9474, 2020.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Self-refine: Iterative refinement with self-feedback. In *NeurIPS*, 2023.
- Azalia Mirhoseini, Anna Goldie, Mustafa Yazgan, Joe Jiang, Ebrahim Songhori, Shen Wang, Young-Joon Lee, Eric Johnson, Omkar Pathak, Sungmin Bae, et al. Chip placement with deep reinforcement learning. *arXiv preprint arXiv:2004.10746*, 2020.

- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
   Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. arXiv preprint
   arXiv:1312.5602, 2013.
  - Allen Nie, Ching-An Cheng, Andrey Kolobov, and Adith Swaminathan. Importance of directional feedback for LLM-based optimizers. In *NeurIPS 2023 Foundation Models for Decision Making Workshop*, 2023.
  - Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with "gradient descent" and beam search. In *EMNLP*, 2023.
  - Stuart J Russell and Peter Norvig. Artificial intelligence: a modern approach. Pearson, 2016.
  - Jérémy Scheurer, Jon Ander Campos, Tomasz Korbak, Jun Shern Chan, Angelica Chen, Kyunghyun Cho, and Ethan Perez. Training language models with language feedback at scale. *arXiv preprint arXiv:2303.16755*, 2023.
  - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
  - Shai Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and Trends*® *in Machine Learning*, 4(2):107–194, 2012.
  - Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: language agents with verbal reinforcement learning. In *NeurIPS*, 2023.
  - David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
  - Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling LLM test-time compute optimally can be more effective than scaling model parameters. *arXiv* preprint arXiv:2408.03314, 2024.
  - Richard S. Sutton. The bitter lesson, 2019. URL http://www.incompleteideas.net/IncIdeas/BitterLesson.html.
  - Richard S Sutton, Andrew G Barto, et al. *Reinforcement learning: An introduction*, volume 1. MIT Press Cambridge, 1998.
  - Adrien Ali Taiga, Rishabh Agarwal, Jesse Farebrother, Aaron Courville, and Marc G Bellemare. Investigating multi-task pretraining and generalization in reinforcement learning. In *The eleventh international conference on learning representations*, 2023.
  - Mark Towers, Ariel Kwiatkowski, Jordan Terry, John U Balis, Gianluca De Cola, Tristan Deleu, Manuel Goulão, Andreas Kallinteris, Markus Krimmel, Arjun KG, et al. Gymnasium: A standard interface for reinforcement learning environments. *arXiv preprint arXiv:2407.17032*, 2024.
  - Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *Transactions on Machine Learning Research*, 2024a.
  - Wenyi Wang, Hisham A Alyahya, Dylan R Ashley, Oleg Serikov, Dmitrii Khizbullin, Francesco Faccio, and Jürgen Schmidhuber. How to correctly do semantic backpropagation on language-based agentic systems. *arXiv preprint arXiv:2412.03624*, 2024b.
  - Xingyao Wang, Boxuan Li, Yufan Song, Frank F Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, et al. Opendevin: An open platform for ai software developers as generalist agents. *arXiv preprint arXiv:2407.16741*, 2024c.
  - Anjiang Wei, Allen Nie, Thiago SFX Teixeira, Rohan Yadav, Wonchan Lee, Ke Wang, and Alex Aiken. Improving parallel program performance through dsl-driven code generation with llm optimizers. *arXiv preprint arXiv:2410.15625*, 2024.

- Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, Ahmed Hassan Awadallah, Ryen W White, Doug Burger, and Chi Wang. AutoGen: Enabling next-gen LLM applications via multi-agent conversations. In *COLM*, 2024.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *ICLR*, 2024a.
- Zonglin Yang, Xinya Du, Junxian Li, Jie Zheng, Soujanya Poria, and Erik Cambria. Large language models for automated open-domain scientific hypotheses discovery. In *ACL*, pp. 13545–13565, 2024b.
- Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Pan Lu, Zhi Huang, Carlos Guestrin, and James Zou. Optimizing generative ai by backpropagating language model feedback. *Nature*, 639(8055):609–616, 2025.
- Shaokun Zhang, Jieyu Zhang, Jiale Liu, Linxin Song, Chi Wang, Ranjay Krishna, and Qingyun Wu. Offline training of language model agents with functions as learnable weights. In *Forty-first International Conference on Machine Learning*, 2024.
- Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Claire Cui, Olivier Bousquet, Quoc V Le, and Ed H. Chi. Least-to-most prompting enables complex reasoning in large language models. In *ICLR*, 2023.

#### A LARGE LANGUAGE MODEL ACCESS CARD

The experiments were conducted during the period of February 2025 to April 2025. The model used was Anthropic's Claude Sonnet 3.5v2, and the exact model name was "anthropic.claude-3-5-sonnet-20241022-v2:0". All of the agents, including the baseline agents from other papers (ResearchAgent (Huang et al., 2023), PAL agent (Gao et al., 2023a), Self-Refine agent (Madaan et al., 2023)) were all rerun using the same model endpoint as our learned agent during the same period of time.

#### B DISCUSSION ON OTHER IMPORTANT FACTORS FOR AGENT LEARNING

#### B.1 Specifying Agent Behavior Through Workflows

The same agent can be described by different workflow graphs, but these graphs may have different optimization properties. We discuss factors that affect optimization difficulty.

**Modularization.** Breaking down the reasoning process from a monolithic block into multiple smaller blocks has proven useful for complex reasoning tasks (Zhou et al., 2023). It remains an empirical question to what extent one should modularize, but the guiding principle here is to decompose a monolithic process into a computational graph with smaller operators.

**Parametrization.** There are some parts of the workflow that do not need much design exploration. For example, loading in a text file given a file name would always result in similar code (even with error handling). Designing an optimizable workflow requires engineers to think carefully and parametrize the part of their workflow where the optimal solution is not known a priori to them, where exploring new parameters is valuable.

**Initialization.** Similar to the research on neural network initialization (e.g. (Glorot & Bengio, 2010)), the optimizable functions contain initial code and docstrings designated by engineers. We found that a workflow that involves operators with ambiguous docstrings is difficult to optimize via generative optimization. We advise using initialization that conveys a clear (desired) behaviors of operators used in the graph.

#### B.2 GUIDING AGENT LEARNING PROCESS THROUGH EFFECTIVE FEEDBACK

Prior works have studied the importance of providing feedback to the optimization process (Chen et al., 2023; Nie et al., 2023; Wei et al., 2024). We highlight and summarize a few useful types of implementations below.

**Fine-grained Feedback.** The simplest form of feedback is to include the correct/incorrect information or the numerical reward in text. However, other feedback designs have also proven to be useful. In Section 4, we show that we can have *staged* feedback, where a different feedback is used when the agent reaches different reward regions (Table A.6). This allows flexibility on how to guide the generative optimizer to search the solution space. Another way is to identify trigger keywords from environment information (such as a system profiler) and retrieve corresponding pre-written feedback (Wei et al., 2024).

**Suggestive Feedback.** The best type of feedback tells the generative optimizer exactly how to change the parameter or output – it should be actionable. If it is not possible to know this information, a suggestion should still be made with proper degree of suggestions in the phrasing. Earlier this was referred to as "directional" feedback (Nie et al., 2023) and later Wei et al. (2024) showed suggestive feedback allowed the optimizer to find better solutions than explanation-based feedback.

#### C MLAGENTBENCH DETAILS

#### C.1 AGENT DESIGN DETAILS

The ML agent shares a similar design for both tasks with modular components for different steps of the machine learning pipeline.

#### C.2 FEEDBACK DESIGN DETAILS

We provided task-specific feedback instructions when the agent reaches different performance level. We show the feedback template in Figure A.1. The {SUGGESTION} block if filled by the suggestive feedback in Table A.1.

```
1 Epoch {epoch}/20:
2
3 Accuracy: {val_accuracy:.4f}
4 F1: {val_f1:.4f}
5 Precision: {val_precision:.4f}
6 Recall: {val_recall:.4f}.
7
8 {SUGGESTION}
```

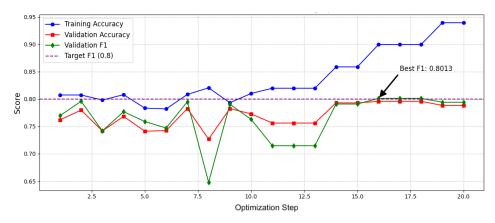
Figure A.1: The feedback template used for the ML Agent for Spaceship-Titanic.

Accuracy	Suggestive Feedback
Val F1 < 0.5	"Model performance is poor. Try better feature engineering and preprocessing."
$0.5 \le \text{Val F1} < 0.7$	"Model is showing promise but needs improvement. Consider class balancing techniques."
$0.7 \le \text{Val F1} < 0.8$	"Model is performing well. Fine-tune hyperparameters for further improvements."
Val F1 $\geq 0.8$	"Excellent performance! Focus on preventing overfitting."

Table A.1: Staged suggestive feedback for the ML agent at different accuracy levels for the Spaceship-Titanic task.

#### C.3 LLM AGENT LEARNING RESULTS

We perform a training and validation split outside of the agent and only pass the training set as input to the agent. This is due to the fact that generative optimization requires an optimization signal. Kaggle does not permit more than 5 submissions on the test set per day, therefore, we do not use the test set as our optimization signal. We randomly split the training data, providing 80% to the agent –



**Figure A.2:** We show the training/validation accuracy and F1 score from machine learning model outputted by the ML Agent after each optimization step. Note that the machine learning model outputted could be trained internally for hundreds of epochs. The x-axis describes the number of optimization steps in the generative optimizer to update the parameter of the ML Agent.

the agent is allowed to further split that data into train and validation. We use the 20% as our test set to evaluate the agent's machine learning model's performance. We show the learning progress of one trial run in Figure A.2. The x-axis of this figure shows the optimization steps. Although on a cursory glance, this graph seems to be depicting typical model overfitting behavior as training accuracy goes up and validation accuracy goes down as optimization continues, it is however not the case. At each optimization step, the agent is producing a fully trained machine learning model using the training dataset, with however many numbers of training iterations it chooses. This figure shows the phenomenon of meta-overfitting, where the generative optimizer updates the agent to choose hyperparameters, training procedures of the model that overfits the training set, even though the feedback reward comes purely from the validation performance.

#### C.4 EFFECT OF WORKFLOW DESIGN

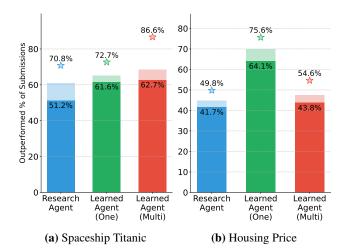


Figure A.3: Leaderboard performance metrics for learned model submissions.

#### C.5 Examples of the Learned ML Agent

We provide initial code (with docstrings) for each game and the final learned code. For Spaceship-Titanic Agent, initial code is in Figure A.10, A.11, A.12, A.13, A.14. Some of the functions are not heavily updated, but we showcase the learned final functions in Figure. Essentially, the generative

optimizer chose to tune the numbers (hyperparameters) of machine learning models and ensemble method.

#### D BATCH LEARNING AGENT DETAILS

MiniBatch	Dyck Languages	Geometric Shapes	Linguini	Causal Understanding
Un-Optimized	$0.114 \pm 0.007$	$\textbf{0.074} \pm 0.003$	$0.183 \pm \textbf{0.010}$	$0.114 \pm 0.005$
Batch=1 Batch=3 Batch=5	$0.183 \pm 0.049$ $0.063 \pm 0.010$ $0.190 \pm 0.031$	$0.343 \pm 0.039$ $0.389 \pm 0.040$ $0.200 \pm 0.099$	$0.149 \pm 0.024$ $0.234 \pm 0.012$ $0.170 \pm 0.030$	$0.375 \pm 0.146$ $0.408 \pm 0.097$ $0.531 \pm 0.018$

MiniBatch	Disambiguation QA	Boolean Expressions	Movie Recommendation	Boardgame QA
Un-Optimized	$0.358 \pm 0.013$	$0.076 \pm 0.009$	$0.238 \pm 0.007$	<b>0.371</b> ± 0.003
Batch=1 Batch=3	<b>0.537</b> ± 0.036 <b>0.295</b> ± 0.091	$0.177 \pm 0.005$ $0.238 \pm 0.006$	$0.889 \pm 0.038$ $0.683 \pm 0.119$	$0.341 \pm 0.032 \\ 0.278 \pm 0.009$
Batch=5	$0.293 \pm 0.091$ $0.526 \pm 0.035$	$0.154 \pm 0.004$	$0.810 \pm 0.016$	$0.276 \pm 0.009$ $0.276 \pm 0.007$

**Table A.2:** Performance across tasks and batching strategies. Best test accuracy per task is bolded; standard error is shown in smaller gray text. For Boardgame QA, we observe meta-overfitting: the learned workflow had strong validation scores but failed to generalize to test examples. Base model is Claude Sonnet-3.5-v2.

#### E BATCH LEARNING LANGGRAPH AGENT DETAILS

The emergence of LLM multi-agent frameworks (Wu et al., 2024; LangChain, 2024) allow static programs to have dynamic behaviors enabled by LLMs. Here we apply generative optimization to such frameworks to enable an agentic workflow to improve itself.

**Agent Design.** We used LangGraph (LangChain, 2024) to implement two popular LLM agent designs. The first one is program-aided language model (PAL) (Gao et al., 2023a). This agent design consists of two components: first, it tries to produce a Python program conditioned on a prompt and the input. Then it executes the program to get the final answer. We learn both of these components. The second agent is a self-refine agent (Madaan et al., 2023), where the agent would use a function to solve the problem, verify its solution, and if the solution is wrong, it will try to refine the solution until it passes the verification step.

**Learning Design** We use LangGraph to build the workflow but use OptoPrime (Cheng et al., 2024b) as the optimizer to learn all the mentioned modules. At each iteration, execution traces from a minibatch of examples are captured, evaluated in terms of feedback, and aggregated. Feedback is concatenated into an aggregated feedback, which is then processed by the optimizer. Before implementing batch learning, optimizing by example would overfit, leading to over-specialized improvements that failed to generalize.

**Feedback Design** We provide feedback as string with templates for both correct and incorrect responses, revealing the solution to the optimizer when the answer is incorrect.

**Results.** Empirical results in Table A.3 confirm the efficacy of the generative optimization framework. We evaluate on GSM8K (Kazemi et al., 2025) and BBEH (Kazemi et al., 2025). For GSM8K, we use the same train/validation/test split as used in DSPy (Khattab et al., 2024). For BBEH, we chose two tasks as representative examples to verify our pipeline. The baselines are both PAL and self-refine agent implemented with good initial starter code and working prompts. However, the learned agent (both code and prompts are learned) performs much better on GSM8K, increasing the performance from 78.9% to 93.4%. For BBEH, the initial agent was not able to output answers with a valid format without few-shot examples. In this zero-shot setup, the optimizer is able to find good prompts and valid code to ensure the produced answer is correct.

		BBEH	
	GSM8K	Causal Understanding	BoardgameQA
PAL Agent (Gao et al., 2023b)	78.9	5.0	5.0
Learned PAL Agent (Ours)	<b>93.4</b>	<b>42.5</b>	<b>33.0</b>
Self-Refine Agent (Madaan et al., 2023)	78.2	0.0	0.0
Learned Self-Refine Agent (Ours)	<b>86.8</b>	<b>44.0</b>	<b>32.5</b>

Table A.3: Comparison of baseline LLM agent design and their optimized design on GSM8K and BBEH.

#### E.1 AGENT DESIGN DETAILS

The PAL (Program-Aided Language Model) agent (Gao et al., 2023a) is designed to have two functions: parse\_problem and execute\_code. parse\_problem makes a call to the LLM with the following prompt "Read the problem and output a Python expression to compute the answer and store it into 'result' variable. Problem: {}". The execute\_code uses Python's "exec" function to execute the program written by parse\_problem. Both functions are updated by the generative optimizer.

The self-refine agent (Madaan et al., 2023) is designed to have three functions: solve\_problem, verify\_solution, and refine\_solution. All three are LLM calling functions with prompt strings defined within the function.

- solve\_problem: "Solve the following problem step by step and give the final answer: {question}.
   Solution:"
- verify\_solution: "You are a math expert. Verify the solution below for correctness. Problem: {question}. Solution and Answer: {solution}. Is the answer correct? If not, explain the error."
- refine\_solution: "The previous answer was found to be incorrect. {verification\_feedback}. Please solve the problem again correctly: {question}. Correct Solution:"

Each function takes the LangGraph's state dictionary as input and was added as nodes to LangGraph's StateGraph for execution. The self-refinement loop is controlled by LangGraph's conditional routing strategy.

#### E.2 FEEDBACK DESIGN DETAILS

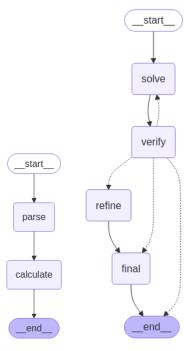
The feedback given to the optimizer when the answer is correct is "ANSWER IS CORRECT/SUCCESS" and when the answer is wrong, we reveal the reference solution since there is a training phase for the agent: "WRONG ANSWER / FAILED - your answer: {answer} vs. good answer: {solution}".

#### F ATARI GAME DETAILS

#### F.1 GAME SETUP

The training configuration is reported in Table A.4 and the environment setup is reported in Table A.5. The Atari Gym offers many wrappers to help with learning. Atari environments by default uses frameskip (repeat actions) to reduce the horizon length and use sticky action probability to randomly repeat the previous action with given probability. Both were designed to enable better training for the deep neural network. In our experiment, we found that not using sticky action results in better optimization of the model.

We generate data on-the-fly for Atari games using object-centric Atari Environments (OCAtari) (Delfosse et al., 2024), a wrapper for the Gymnasium API (Towers et al., 2024) that provides object-centric representation of the game screen at each



(a) PAL (b) Self-refine agent

Parameter	Value	
Environment name	{env}-NoFrameskip-v4	
Action repeat (frameskip)	4	
Sticky action probability	0.0	
Optimization iterations	30	
Rollout length	15/300/400 steps	
Memory size (optimizer context)	5	
Evaluation episode length	$\sim$ 4000 steps	
LLM optimizer	OptoPrime	
LLM Backend	Claude-3.5 Sonnet-20241022-v2:0	
Access Date	3/20/2025	

**Table A.4:** Atari Gym environment and training configurations

Parameter	Breakout	Pong	Space Invaders
Rollout horizon	300 steps	400 steps	15 steps
Action space	LEFT/RIGHT/ NOOP	UP/DOWN/ NOOP	LEFT/RIGHT FIRE/NOOP
Env special mechanics	Auto-fire on life loss	None	Fire cooldown

Table A.5: Atari game-specific experiment configurations

timestep. For instance, for the game Pong, OCAtari returns the position (x,y), size (width, height), and velocity (dx,dy) of the player paddle, ball, and enemy paddle. This representation abstracts away from raw pixel inputs, providing the LLM optimizer and our agent with structured state information that facilitates targeted improvements to the agent's prediction and action selection. The actual input observation to the agent is shown in Figure A.9, and an annotated screen through OCAtari can be seen in Figure A.5, A.6, and A.7.

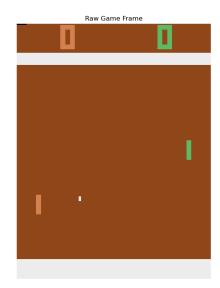
**Pong** In Pong, the player controls a paddle on the right side of the screen to deflect the ball into the enemy's goal. The player scores a point if the enemy misses the ball. The game ends when one side scores 21 points.

**Breakout** In Breakout, the player moves a bottom paddle horizontally to deflect a ball that scores against brick walls upon contact. The brick wall consists of six rows of different colored bricks, with higher bricks worth more points. Hitting higher bricks would deflect the ball faster, increasing the difficulty in catching the ball. The player wins after scoring 864 points. The player loses one life when failing to catch the ball and the ball moves out of range. The player has five lives in total.

**Space Invaders** In Breakout, the player moves a bottom paddle horizontally to deflect a ball that scores against brick walls upon contact. The brick wall consists of six rows of different colored bricks, with higher bricks worth more points. Hitting higher bricks would deflect the ball faster, increasing the difficulty in catching the ball. The player wins after scoring 864 points. The player loses one life when failing to catch the ball and the ball moves out of range. The player has five lives in total.

#### F.2 FEEDBACK DESIGN DETAILS

We provided game-specific feedback instructions when the agent reaches different reward regions.





**Figure A.5: Pong**: An annotated screenshot to show how OCAtari (Delfosse et al., 2024) translates objects from pixels to obejcts with annotations.



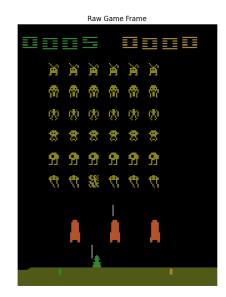


**Figure A.6: Breakout**: An annotated screenshot to show how OCAtari (Delfosse et al., 2024) translates objects from pixels to obejcts with annotations.

## F.3 AGENT DESIGN DETAILS

**Pong** In order to succeed at Pong, the agent should accurately predict where the ball will intersect with the player's paddle plane, accounting for bounces off the top and bottom walls. Thus, we adapt our base agent architecture to focus on ball trajectory prediction and paddle positioning (predict\_ball\_trajectory() and select\_action()). We initialize predict\_ball\_trajectory() to return the current y coordinate of the ball and select\_action() to return a random action of UP or DOWN. In the docstring, we provide detailed description of the game screen, including screen dimensions and paddle positions. We show the initialized agent in Figure A.16 and the optimized agent in Figure A.17.

**Breakout** Breakout has a similar emphasis of considering wall boucing but with a focus on brick targeting. Like Pong, we adapt our base agent architecture to focus on predicting the trajectory of the ball (predict\_ball\_trajectory()), but also prioritizing hitting bricks with higher scores (generate\_paddle\_target())





**Figure A.7: Space Invaders**: An annotated screenshot to show how OCAtari (Delfosse et al., 2024) translates objects from pixels to obejcts with annotations.

Performance Level	Feedback
High (Reward $\geq$ 19)	"Good job! You're close to winning the game! You're scoring 20 points against the opponent, only 1 points short of winning."
Medium (0 < Reward < 19)	"Keep it up! You're scoring 12 points against the opponent but you are still 9 points from winning the game. Try improving paddle positioning to prevent opponent scoring."
Low (Reward $\leq 0$ )	"Your score is $-5$ points. Try to improve paddle positioning to prevent opponent scoring."

Table A.6: Staged feedback for the Pong agent at different performance levels

and selecting paddle action based on the analysis (<code>select\_paddle\_action()</code>). We initialize both <code>predict\_ball\_trajectory()</code> and <code>generate\_paddle\_target()</code> to return None, and <code>select\_paddle\_action()</code> to move the paddle <code>LEFT</code> or <code>RIGHT</code> by comparing the paddle location to the target position generated by <code>generate\_paddle\_target()</code>, which is None upon initialization. In the docstring, we describe the game screen, such as locations of left and right wall, but we leave the exact location out for the LLM to infer based on the traced trajectory. We also describe the point system of the brick wall and some generic strategic considerations (without telling the agent how to implement these strategies). We show the initialized agent in Figure A.18, A.19 and the optimized agent in Figure A.20, A.21.

**Space Invaders** For Space Invaders, we adapt our base agent architecture to into two tasks of deciding whether to shoot (decide\_shoot()) and deciding where to move (decide\_move()), and finally combining the two decisions in (combine\_actions()). We initialize decide\_shoot() and decide\_movement() to return random actions, and combine\_actions() map the outputs of the previous two functions to the correct action. In the docstring, we describe the game setup and the presence of shield objects. We show the initialized agent in Figure A.22, A.23 and the learned agent in Figure A.24.

#### F.4 LLM AGENT LEARNING RESULT

Performance Level	Example Feedback
High (Reward ≥ 300)	"Good job! You're close to winning the game! You're scoring 320 points against the opponent, try ensuring you return the ball, only 30 points short of winning."
Medium (0 < Reward < 300)	"Keep it up! You're scoring 50 points against the opponent but you are still 300 points from winning the game. Try improving paddle positioning to return the ball and avoid losing lives."
Low (Reward $\leq 0$ )	"Your score is -5 points. Try to improve paddle positioning to return the ball and avoid losing lives."

**Table A.7:** Staged feedback for the Breakout agent at different performance levels

Performance Level	Feedback
High (Reward ≥ 1000)	"Great job! You're performing well with an average score of 1005. Try to score more even more points"
Medium (500 < Reward < 1000)	"Good progress! Your average score is 570. Focus on better timing for shooting and avoiding enemy projectiles."
Low (Reward ≤ 500)	"Your average score is 270. Try to improve your strategy for shooting aliens and dodging projectiles."

Table A.8: Staged feedback for the Space Invaders agent at different performance levels

## F.5 DEEP RL RESULT

Due to a large variation in how people report Atari game results and the fact that many state-of-the-art deep RL models are not released as open-source, the numbers we reported in Table A.9 are from CleanRL report (Huang et al., 2022b), the published ICLR blog post (Huang et al., 2022a) and the experiment log<sup>1</sup>. In terms of runtime, we directly compute the time from the Weights & Biases log. For Breakout and Space Invaders, the agent performances were continuously improving, so we reported the duration of the full experiment run. For Pong, the RL policy plateaued before the experiment finished, so we found the time step where the policy achieved the highest performance reliably and computed training time starting from the launch of the experiment to that time step.

It is worth noting that we reporeted the Deep RL results with 8 parallel environment instances in Table A.9. However, there are faster implementations of Deep RL training on Atari games. For example, Apex-DQN (Horgan et al., 2018) would train the actor and critic model separately in a truly asynchronous fashion, resulting in massive reduction of training time. EnvPool is a C++-based batched environment pool that enabled fast sampling and interaction with the game environment. All of these changes enabled faster learning. For example, on Breakout, with 32 to 64 parallel environments, Advantage Actor-Critic (A2C) can learn a high performing policy in 33m 19s<sup>2</sup>. However, Trace only uses 1 environment instance and has not gone through any special speed-related algorithm/hardware optimization.

#### F.6 Examples of the Learned Atari Agent

We provide initial code (with docstrings) for each game and the final learned code. For Pong Agent, initial code is in Figure A.16, and final agent in Figure A.17. For Breakout agent, initial code is in

https://wandb.ai/cleanrl/cleanrl.benchmark/reports/Atari--VmlldzoxMTExNTI https://wandb.ai/costa-huang/cleanRL/reports/Breakout-v5--VmlldzoxNDI1MTIx

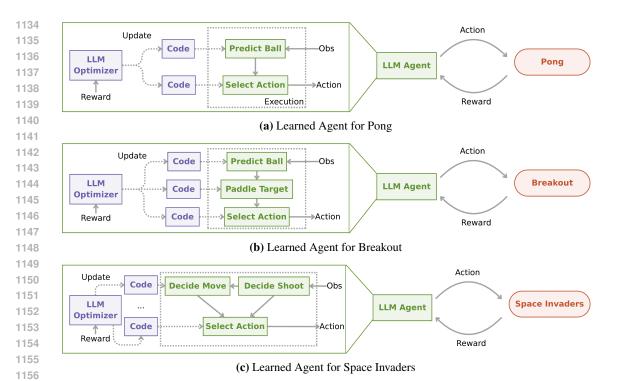


Figure A.8: Agent Design for Atari Games. We can design different agent workflow graphs for the learned agent to achieve high scores in three Atari games.

```
TracedEnv.step.step16 = {
    'Player': {'x': 99, 'y': 189, 'w': 16, 'h': 4, 'dx': 0, 'dy': 0},
    'Ball': {'x': 7, 'y': 193, 'w': 2, 'h': 4, 'dx': -4, 'dy': 4},
    'RB': [{'x': 8, 'y': 57, 'w': 144, 'h': 6}],
    'OB': [{'x': 8, 'y': 63, 'w': 144, 'h': 6}],
    'YB': [{'x': 8, 'y': 69, 'w': 144, 'h': 6}],
    'GB': [{'x': 8, 'y': 75, 'w': 144, 'h': 6}],
    'AB': [{'x': 8, 'y': 81, 'w': 144, 'h': 6}],
    'BB': [{'x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 8, 'y': 87, 'w': 144, 'h': 6}],
    'YB': ['x': 87, 'w':
```

Figure A.9: Example of a single traced step from Breakout

Figure A.18, A.19 and final agent in Figure A.20, A.21. For Space Invaders agent, initial code is in Figure A.22, A.23 and the final agent in Figure A.24.

Game	Learned Agent	DQN (Time)	PPO (Time)	Human
Pong	21 (43m)	20 (10h 6m)	19 (2h 24m)	14.59
Breakout	353 (1h 31m)	302 (26h 54m)	443 (3h 8m)	30.47
Space Invaders	1200 (36m)	1383 (26h 52m)	939 (5h 39m)	1668.67

**Table A.9:** Comparison of Algorithm Performance on Atari Games with Time. Due to high variations of numbers reported by different papers, we report results from an open-source implementation of RL algorithms (Huang et al., 2022b) and publicly available experiment logs (Huang et al., 2022a). RL algorithms are trained with 8 parallel environment instances. Our agent is trained on 1 environment instance. Note that highly optimized Deep RL with 32 environment instances can reach ~450 on Breakout in 33m, see Appendix F.5.

# G LARGE LANGUAGE MODEL USE FOR WRITING

A small amount of paragraphs have been polished by GPT-5. The process is – the author wrote the sentence or paragraph first and then send into the LLM with the prompt of "Polish the following writing and correct the grammar mistakes." LLM was never used to directly produce a paragraph without an original human-written input. The LLM assistance was only used to enhance the calrity and readability of the paragraph only.

# H EXAMPLES OF LEARNED AGENT

#### H.1 ML AGENT

## H.2 ATARI GAME AGENTS

```
1 import trace
1243
        3 @trace.model
        4 class SpaceshipTitanicPipeline(Module):
1245
1246
              def __call__(self, x, y=None, test_data=None):
    processed_data = self.preprocess(x)
1247
                  selected_features = self.select_features(processed_data)
1248
                  if y is not None:
1249
        10
                      ensemble = self.ensemble_model(selected_features, processed_data)
        11
1250
                      model = self.train_model(ensemble, selected_features, processed_data, y)
                      if test data is not None:
1251
                          processed_test_data = self.preprocess(test_data)
1252
                           filtered_test_data = self.filter_features(selected_features,
        15
1253
               processed_test_data)
        16
                          return self.predict(model, filtered_test_data)
1254
                      filtered_data = self.filter_features(selected_features, processed_data)
                      return self.predict(model, filtered_data)
1255
        18
        19
                  else:
1256
        20
                      ensemble = self.ensemble_model(selected_features, processed_data)
        21
                       processed_test_data = self.preprocess(x)
1257
                       filtered_test_data = self.filter_features(selected_features, processed_test_data)
1258
                      model = self.train_model(ensemble, selected_features, processed_data,
              pd.Series([False] * len(processed_data)))
1259
                      return self.predict(model, filtered_test_data)
1260
1261
        26
              def filter_features(self, selected_features, data):
                  return data[selected_features]
1262
        29
              @trace.bundle(trainable=True)
1263
        30
              def preprocess(self, data):
1264
        31
                  Preprocessing Steps (some examples on how you could do this, however you can use
1265
        32
               your own method if it works better):
1266
        33
                  1. Missing Value Handling:
        34
                      - Numerical features: Intelligent imputation (median, mean, or 0)
1267
                      - Categorical features: Mode filling or meaningful defaults
1268
                      - Outlier detection and treatment
1269
                  2. Feature Engineering:
1270
                       - Passenger ID parsing:
                          * Extract group and individual identifiers
1271
                           * Create group-related features
1272
                      - Cabin information extraction:
                           * Deck identification
1273
                           * Cabin number parsing
                          * Side (port/starboard) classification
                      - Name feature parsing:
1275
                          * Title extraction
1276
                           * Potential family relationship inference
1277
                  3. Advanced Feature Creation:
1278
                      - Family size computation
                      - Total and relative spending calculations
1279
                      - Amenity usage patterns
1280
                      - Spatial features (cabin location metrics)
1281
                  4. Categorical Variable Handling:
1282
                      - One-hot encoding
                       - Label encoding
1283
                      - Embedding techniques for high-cardinality features
1284
        60
        61
                  5. Numerical Feature Transformation:
1285
                      - Scaling (StandardScaler, MinMaxScaler)
        62
1286
        63
                       - Skewness correction (log, square root, Box-Cox)
                      - Normalization techniques
1287
        64
        65
1288
        66
                  Args:
                      data (pd.DataFrame): Raw input dataset
1289
        68
1290
        69
                  Returns:
                  pd.DataFrame: Preprocessed dataset with engineered features
1291
        70
1292
                  return data
1293
```

**Figure A.10:** Initial code for Spaceship-Titanic ML Agent (Part 1). Docstrings are generated by ChatGPT and then edited by humans.

```
1299
1300
1301
1302
        1@trace.model
1303
        2 class SpaceshipTitanicPipeline(Module):
1304
              # (continued from above)
              @trace.bundle(trainable=True)
1305
              def select_features(self, processed_data):
1306
                  Select the most relevant features for predicting whether passengers were transported.
1307
                  Selection Methodology (some examples on how you could do this, however you can use
1308
              your own method if it works better):
1309
        10
                  1. Statistical Feature Importance:
                      - Correlation analysis
1310
                      - Mutual information
1311
                      - Chi-squared tests
                      - Model-based feature importance
1312
                  2. Feature Weighting Criteria:
                      - Predictive power for transportation status
1314
                      - Domain-specific relevance
1315
                      - Minimal multicollinearity
        19
        20
                      - Computational efficiency
1316
        21
1317
                  3. Key Feature Categories:
                      - Demographic Signals
        23
1318
                      - Travel Characteristics
1319
                      - Economic Indicators
1320
        27
                  4. Selection Mechanism:
1321
                      - Probabilistic feature selection
                      - Dynamic weight adjustment
1322
                      - Prevent overfitting through selective inclusion
        30
1323
        31
1324
        32
                  Args:
                      processed_data (pd.DataFrame): Preprocessed dataset
        33
1325
        34
                  Returns:
1326
        35
                  list: Optimally selected feature names with associated weights """
1327
        37
                  all_features_with_weights = {col: 0.5 for col in processed_data.columns}
        38
        39
1329
        40
                  available_features = {k: v for k, v in all_features_with_weights.items() if k in
               processed_data.columns}
1330
        41
1331
                  feature_names = list(available_features.keys())
        42
        43
                  feature_weights = list(available_features.values())
1332
        44
1333
        45
                  num_features = min(len(feature_names), int(len(feature_names) * 0.8))
        46
1334
        47
                  selected_features = np.random.choice(
1335
        48
                      feature_names,
        49
                      size=num_features,
1336
        50
                      replace=False
1337
        51
                      p=[w/sum(feature_weights) for w in feature_weights]
        52
                  ).tolist()
1338
        53
1339
                  selected_features = [f for f in selected_features if f in processed_data.columns]
1340
                  return selected_features
1341
```

Figure A.11: Initial code for Spaceship-Titanic ML Agent (Part 2). Docstrings are generated by ChatGPT and then edited by humans.

```
1352
1353
1354
1355
         1@trace.model
1356
         2 class SpaceshipTitanicPipeline(Module):
              # (continued from above)
1357
              @trace.bundle(trainable=True)
1358
              def ensemble_model(self, features, data):
1359
                   Create an ensemble model for predicting passenger transport status.
1360
                   Ensemble Strategy (some examples on how you could do this, however you can use your
1361
               own method if it works better):
1362
        10
                  1. Model Diversity:
                       - Tree-based models (Random Forest, Gradient Boosting)
1363
        11
                       - Linear models (Logistic Regression variants)
1364
                       - Support Vector Machines
                       - Probabilistic classifiers
        14
1365
        15
1366
                   2. Ensemble Techniques:
        16
                       - Voting, boosting, bagging, stacking
1367
                       - Stacking with meta-learners
        18
1368
        19
                       - Weighted model combination
                       - Regularization-aware model selection
1369
        20
1370
                  3. Hyperparameter Optimization:
                       - Cross-validated parameter tuning
- Regularization strength balancing
        23
1371
1372
        25
                       - Learning rate and depth control
1373
                       - Subsample and feature sampling strategies
1374
                   4. Computational Considerations:
                       - Computational complexity management
1375
                       - Memory-efficient model design
1376
                       - Scalable ensemble construction
1377
        33
1378
                       features (list): Selected feature names
                       data (pd.DataFrame): Processed dataset
1379
1380
                   Returns:
                   sklearn Classifier: Configured ensemble model ready for training
1381
1382
                   models = [
                      ('rf', RandomForestClassifier(n_estimators=150, max_depth=10,
        41
1383
               min_samples_split=5, min_samples_leaf=4, max_features='sqrt',random_state=42)),
1384
                       ('gbr', GradientBoostingClassifier(n_estimators=200, learning_rate=0.03,
               max_depth=3, subsample=0.8, min_samples_split=5, random_state=42)),
1385
                       ('xgb', XGBClassifier(n_estimators=200, learning_rate=0.03, max_depth=3,
1386
               colsample_bytree=0.6, subsample=0.8, reg_alpha=0.1, reg_lambda=1.0, gamma=1,
               random_state=42)),
1387
                       ('lasso', LogisticRegression(penalty='11', C=0.1, random_state=42)), ('ridge', LogisticRegression(penalty='12', C=20.0, random_state=42))
1388
        45
                       ('elastic', LogisticRegression(penalty='elasticnet', C=0.1, l1_ratio=0.8,
        46
1389
               random_state=42))
1390
        47
        48
1391
                   ensemble = VotingClassifier(
        49
1392
                      estimators=models,
        50
                       voting='soft'
        51
1393
                       weights=[2, 3, 3, 2, 1, 1]
        52
1394
        53
                  )
        54
1395
        55
                   return
1396
```

**Figure A.12:** Initial code for Spaceship-Titanic ML Agent (Part 3). Docstrings are generated by ChatGPT and then edited by humans.

```
1409
1410
1411
1412
1413
1414
1415
         1@trace.model
1416
         2 class SpaceshipTitanicPipeline(Module):
1417
              # (continued from above)
              @trace.bundle(trainable=True)
1418
              def train_model(self, ensemble_model, features, data, results):
1419
                   Train machine learning models to predict whether passengers were transported.
1420
                  Training Methodology (some examples on how you could do this, however you can use
1421
        9
               your own method if it works better):
1422
        10

    Data Preparation:

                       - Feature subset preparation
        12
                       - Cross-validation splitting
1424
        13
                       - Stratified sampling
1425
        15
                   2. Class Imbalance Handling:
1426
                       Weighted loss functionsSMOTE oversampling
        16
1427
                       - Synthetic data generation
1428
                       - Class-aware regularization
        20
1429
                  3. Regularization Techniques:
1430
                       - L1/L2 penalty integration
- Dropout-like regularization
1431
                       - Early stopping mechanisms
1432
                       - Gradient clipping
1433
                   4. Training Optimization:
1434
                       - Adaptive learning rates
                       - Ensemble member performance tracking
1435
                       - Dynamic weight adjustment
                       - Prediction confidence calibration
1437
                  Args:
1438
                       ensemble_model: Configured ensemble model
                       features (list): Selected feature names
1439
                       data (pd.DataFrame): Processed training dataset
1440
                       results (pd.Series): Training labels
1441
                   Trained ensemble model optimized for passenger transportation prediction """
1442
        40
1443
1444
1445
```

**Figure A.13:** Initial code for Spaceship-Titanic ML Agent (Part 4). Docstrings are generated by ChatGPT and then edited by humans.

```
1462
1463
1464
1465
1466
1467
1468
1469
1470
                                    1@trace.model
1471
                                    2 class SpaceshipTitanicPipeline(Module):
                                                         # (continued from above)
1472
                                                         @trace.bundle(trainable=True)
1473
                                                         def predict(self, model, data):
1474
                                                                         Make predictions on whether passengers were transported.
1475
                                                                        Prediction Workflow (some examples on how you could do this, however you can use
1476
                                                           your own method if it works better):
1477
                                                                        1. Probabilistic Prediction:
                                                                                       - Soft classification probabilities
                                 11
1478
                                                                                         - Confidence-based thresholding
1479
                                                                                         - Ensemble prediction aggregation
                                 13
1480
                                 14
                                                                         2. Post-processing Techniques:
1481
                                                                                         - Calibration curves
                                 16
                                                                                         - Probability scaling
1482
                                 17
                                                                                         - Ensemble diversity preservation
                                 18
1483
                                 19
                                                                        3. Output Formatting:
                                 20
1484
                                                                                       - Binary classification output
1485
                                                                                          - Kaggle submission compatibility
                                                                                         - Interpretable prediction format
1486
                                 23
                                 24
1487
                                                                        4. Prediction Quality Assessment:
                                 25
                                                                                          - Uncertainty quantification
1488
                                 26
                                                                                          - Prediction reliability scoring
1489
                                                                                          - Anomaly detection
1490
                                 30
                                                                        Args:
1491
                                                                                          model (VotingClassifier): Trained ensemble model
                                 31
1492
                                 32
                                                                                          data (pd.DataFrame): Processed test dataset
                                 33
1493
                                                                         Returns:
                                 34
                                                                         np.ndarray: Binary predictions for passenger transportation status """ % \left( \frac{1}{2}\right) =\frac{1}{2}\left( \frac{1}{2}\right) \left( \frac{1}{2}\right) 
1494
                                 35
                                 36
1495
                                 37
                                                                         predictions = model.predict(data)
                                                                         predictions = np.array(predictions, dtype=bool)
1496
                                                                         return predictions
1497
1498
```

**Figure A.14:** Initial code for Spaceship-Titanic ML Agent (Part 5). Docstrings are generated by ChatGPT and then edited by humans.

```
1516
1517
1518
1519
1520
1521
1522
         1 import trace
1523
         3 @trace.model
1524
         4 class SpaceshipTitanicPipeline(Module):
1525
               @trace.bundle(trainable=True)
1526
               def preprocess(self, data):
                     """(same as before, skipped to save space)"""
1527
                    \ensuremath{\text{\#}} Create a copy to avoid modifying original data
1528
                    # Handle missing values in numeric columns
         10
                    numeric_columns = ["Age", "RoomService", "FoodCourt", "ShoppingMall",
                        "Spa", "VRDeck"]
1530
                    for col in numeric_columns:
         13
1531
         14
                        processed_data[col] = processed_data[col].fillna(processed_data[col].median())
         15
1532
                    # Handle boolean/categorical columns
                    processed_data["VIP"] = processed_data["VIP"].fillna(False)
processed_data["CryoSleep"] = processed_data["CryoSleep"].fillna(False)
1533
1534
         19
1535
                    # Convert HomePlanet to numeric using label encoding
                    if "HomePlanet" in processed_data.columns:
    processed_data["HomePlanet"] = processed_data["HomePlanet"].fillna("Unknown")
    planet_map = {"Earth": 0, "Europa": 1, "Mars": 2, "Unknown": 3}
1536
         23
1537
                        processed_data["HomePlanet"] = processed_data["HomePlanet"].map(planet_map)
1538
                    # (skipped some code)
1539
1540
                    # Age-related features
                    processed_data["Age"] = processed_data["Age"].fillna(processed_data["Age"].median())
1541
                    processed_data["AgeGroup"] = pd.qcut(
1542
                        processed_data["Age"], q=6, labels=[0, 1, 2, 3, 4, 5]
1543
                    ).astype(int)
1544
                    # Interaction features
                   processed_data["CryoSleepVIP"] = processed_data["CryoSleep"].astype(int) *
1545
                processed_data["VIP"].astype(int)
1546
                   processed_data["SpendingPerAge"] = processed_data["TotalSpending"] /
                processed_data["Age"].clip(lower=1)
1547
                    processed_data["HasSpent"] = (processed_data["TotalSpending"] > 0).astype(int)
1548
         38
                    processed_data["SpendingVariety"] = (processed_data[spending_columns] >
                0).sum(axis=1)
1549
1550
                    # ... standard scaling, dropping columns, etc.
         40
         41
1551
                    # Final check for NaN values
1552
         43
                    processed_data = processed_data.fillna(0)
                    return processed_data
1553
1554
```

Figure A.15: Final learned code for Spaceship-Titanic ML Agent (Part 1). Docstrings are generated by ChatGPT and then edited by humans.

1616

```
1567
1568
1569
         1 import trace
1570
1571
        3 @trace.model
         4 class Policy(Module):
1572
              def __call__(self, obs):
1573
                  predicted_ball_y = self.predict_ball_trajectory(obs)
                  action = self.select_action(predicted_ball_y, obs)
1574
                  return action
1575
        10
              @trace.bundle(trainable=True)
1576
              def predict_ball_trajectory(self, obs):
        11
1577
                  Predict the y-coordinate where the ball will intersect with the player's paddle by
1578
               calculating its trajectory,
1579
        14
                  using ball's (x, y) and (dx, dy) and accounting for bounces off the top and bottom
               walls.
1580
        15
                  Game Setup:
1581
                  - Screen dimensions: The game screen has boundaries where the ball bounces
1582
                    - Top boundary: approximately y=30
1583
                    - Bottom boundary: approximately y=190
                  - Paddle positions:
1584
                    - Player paddle: right side of screen (x = 140)
1585
                    - Enemy paddle: left side of screen (x = 16)
1586
                  Args:
                      obs (dict): Dictionary containing object states for "Player", "Ball", and
1587
               "Enemy".
1588
                                  Each object has position (x,y), size (w,h), and velocity (dx,dy).
1589
                  Returns:
1590
                      float: Predicted y-coordinate where the ball will intersect the player's paddle
1591
               plane.
        30
                             Returns None if ball position cannot be determined.
1592
        31
1593
        32
                  if 'Ball' in obs:
1594
                      return obs['Ball'].get("y", None)
1595
                  return None
        35
        36
1596
              @trace.bundle(trainable=True)
        37
              def select_action(self, predicted_ball_y, obs):
1597
        38
        39
        40
                  Select the optimal action to move player paddle by comparing current player position
               and predicted_ball_y.
1599
        41
1600
                  IMPORTANT! Movement Logic:
        42
                   - If the player paddle's y position is GREATER than predicted_ball_y: Move DOWN
1601
        43
               (action 2)
1602
        44
                    (because the paddle needs to move downward to meet the ball)
                   - If the player paddle's y position is LESS than predicted_ball_y: Move UP (action 3)
1603
        45
        46
                    (because the paddle needs to move upward to meet the ball)
1604
        47
                  - If the player paddle is already aligned with predicted_ball_y: NOOP (action 0)
                    (to stabilize the paddle when it's in position)
1605
        48
        49
                  Ensure stable movement to avoid missing the ball when close by.
1606
        50
1607
        51
                  Args:
                      \verb|predicted_ball_y (float)|: predicted y coordinate of the ball or None|\\
        52
1608
                      obs(dict): Dictionary of current game state, mapping keys ("Player", "Ball",
        53
               "Enemy") to values (dictionary of keys ('x', 'y', 'w', 'h', 'dx', 'dy') to integer
1609
               values)
1610
        54
                  Returns:
                  int: 0 for NOOP, 2 for DOWN, 3 for UP
1611
        56
1612
1613
        58
                  if predicted_ball_y is not None and 'Player' in obs:
        50
                      return random.choice([2, 3])
1614
                  return 0
1615
```

Figure A.16: Initial code for Pong Agent.

```
1620
1621
         1 import trace
1622
         3 @trace.model
1623
         4 class Policy(Module):
1624
              def __call__(self, obs):
                   predicted_ball_y = self.predict_ball_trajectory(obs)
1625
                   action = self.select_action(predicted_ball_y, obs)
1626
                   return action
1627
               @trace.bundle(trainable=True)
        10
               def predict_ball_trajectory(self, obs):
1628
        11
                    """(same as before, skipped to save space)"""
        12
1629
                   if 'Ball' in obs:
        13
                       ball = obs['Ball']
1630
        14
                       # If ball moving away from player, return None
if ball.get('dx', 0) < 0:</pre>
        15
1631
        16
1632
                            return None
        18
1633
                       # Calculate time to reach paddle
        19
                       paddle_x = 140
1634
        20
                       ball_x = ball.get('x', 0)
        21
1635
                       ball_dx = ball.get('dx', 0)
                       if ball_dx == 0:
1636
        24
                            return ball.get('y', None)
1637
        25
                       time_to_paddle = (paddle_x - ball_x) / ball_dx
1638
        26
1639
                       # Calculate predicted y position with improved accuracy
        28
                       ball_y = ball.get('y', 0)
ball_dy = ball.get('dy', 0)
1640
        29
        30
1641
                       predicted_y = ball_y + ball_dy * time_to_paddle
        31
1642
        32
        33
                       # Account for bounces with improved accuracy
1643
        34
                       num bounces = 0
1644
        35
                        while predicted_y < 30 or predicted_y > 190:
                           if predicted_y < 30:</pre>
        36
1645
        37
                                predicted_y = 30 + (30 - predicted_y)
1646
        38
                            if predicted_y > 190:
        39
                                predicted_y = 190 - (predicted_y - 190)
1647
        40
                            num_bounces += 1
1648
        41
                            if num_bounces > 4: # Limit bounce calculations
1649
        43
                       return predicted_y
1650
                 return None
        45
1651
               @trace.bundle(trainable=True)
        47
               def select_action(self, predicted_ball_y, obs):
1653
                   '''(same as before, skipped to save space)'''
                   if predicted_ball_y is not None and 'Player' in obs:
1654
                        # Calculate center of paddle
1655
                       paddle_center = obs['Player']['y'] + obs['Player']['h']/2
1656
                       # Increase margin and add dynamic adjustment based on ball distance
1657
                       base_margin = 4
                       if 'Ball' in obs:
1658
                            ball_x = obs['Ball'].get('x', 0)
dist_factor = (140 - ball_x) / 140 # Normalized distance factor
1659
                            margin = base_margin * (1 + dist_factor) # Larger margin when ball is far
1660
        60
1661
                            # Add momentum-based adjustment
        61
                            if obs['Ball'].get('dx', 0) > 0:
    ball_dy = obs['Ball'].get('dy', 0)
1662
        62
        63
1663
        64
                                # Scale adjustment based on distance
1664
                                predicted_ball_y += ball_dy * dist_factor
        65
                       else:
        66
1665
                            margin = base_margin
        67
1666
        68
                       # More aggressive movement thresholds
        69
1667
                       if paddle_center > predicted_ball_y + margin:
        70
                            return 2 # Move down
        71
1668
                        elif paddle_center < predicted_ball_y - margin:</pre>
1669
        73
                           return 3 # Move up
                       return 0 # Stay in position
1670
                   return 0
1671
```

**Figure A.17:** Final learned code for Pong Agent.

```
1674
        1@trace.model
1675
        2 class Policy(Module):
1676
              def __call__(self, obs):
                  pre_ball_x = self.predict_ball_trajectory(obs)
1677
                  target_paddle_pos = self.generate_paddle_target(pre_ball_x, obs)
                  action = self.select_paddle_action(target_paddle_pos, obs)
1678
                  return action
1679
              @trace.bundle(trainable=True)
1680
        10
              def predict_ball_trajectory(self, obs):
1681
        12
                  Predict the x-coordinate where the ball will intersect with the player's paddle by
1682
              calculating its trajectory,
1683
        13
                 using ball's (x, y) and (dx, dy) and accounting for bounces off the right and left
              walls.
1684
        14
1685
                  Game setup:
                  - Screen dimensions: The game screen has left and right walls and brick wall where
1686
        16
               the ball bounces
1687
                    - Left wall: x=9
                    - Right wall: x=152
1688
                  - Paddle positions:
1689
                    - Player paddle: bottom of screen (y=189)
        20
                  - Ball speed:
1690
                    - Ball deflects from higher-scoring bricks would have a higher speed and is harder
1691
               to catch.
                   - The paddle would deflect the ball at different angles depending on where the ball
        23
1692
               lands on the paddle
1693
        24
                  Args:
1694
                     obs (dict): Dictionary containing object states for "Player", "Ball", and blocks
1695
               "{color}B" (color in [R/O/Y/G/A/B]).
                                 Each object has position (x,y), size (w,h), and velocity (dx,dy).
1696
1697
                     float: Predicted x-coordinate where the ball will intersect the player's paddle
              plane.
1698
        30
                            Returns None if ball position cannot be determined.
1699
        32
                  if 'Ball' not in obs:
1700
        33
                     return None
1701
              @trace.bundle(trainable=True)
        35
1702
              def generate_paddle_target(self, pre_ball_x, obs):
        36
1703
        37
                  Calculate the optimal x coordinate to move the paddle to catch the ball (at
1704
        38
              predicted ball x)
1705
                 and deflect the ball to hit bricks with higher scores in the brick wall.
        39
        40
1706
        41
                  Logic:
1707
                  - Prioritize returning the ball when the ball is coming down (positive dy)
        42
                  - The brick wall consists of 6 vertically stacked rows from top to bottom:
        43
1708
                   - Row 1 (top): Red bricks (7 pts)
        44
1709
                   - Row 2: Orange (7 pts)
        45
                   - Row 3: Yellow (4 pts)
        46
1710
                   - Row 4: Green (4 pts)
        47
1711
        48
                   - Row 5: Aqua (1 pt)
                   - Row 6 (bottom): Blue (1 pt)
        49
1712
        50
                   - Strategic considerations:
1713
        51
                    - Breaking lower bricks can create paths to reach higher-value bricks above
        52
                   - Creating vertical tunnels through the brick wall is valuable as it allows
1714
        53
                      the ball to reach and bounce between high-scoring bricks at the top
1715
        54
                    - Balance between safely returning the ball and creating/utilizing tunnels
        55
                      to access high-value bricks
1716
                  - Ball speed increases when hitting higher bricks, making it harder to catch
        56
1717
        57
        58
1718
                  Args:
        59
                      pre_ball_x (float): predicted x coordinate of the ball intersecting with the
1719
               paddle or None
        60
                      obs (dict): Dictionary containing object states for "Player", "Ball", and blocks
1720
               "{color}B" (color in [R/O/Y/G/A/B]).
1721
                                 Each object has position (x,y), size (w,h), and velocity (dx,dy).
                  Returns:
1722
                      float: Predicted x-coordinate to move the paddle to.
1723
                          Returns None if ball position cannot be determined.
1724
                  if pre_ball_x is None or 'Ball' not in obs:
1725
                      return None
                  return None
1726
```

**Figure A.18:** Initial code for Breakout Agent (Part 1).

```
1731
1732
1733
1734
1735
1736
1737
1738
         1 import trace
1739
        3 @trace.model
1740
         4 class Policy(Module):
1741
              # (continued from above)
1742
1743
              @trace.bundle(trainable=True)
              def select_paddle_action(self, target_paddle_pos, obs):
1744
1745
        11
                  Select the optimal action to move player paddle by comparing current player position
              and target_paddle_pos.
1746
1747
                  Movement Logic:
                   If the player paddle's center position is GREATER than target_paddle_pos: Move
        14
1748
               LEFT (action 3)
                  - If the player paddle's center position is LESS than target_paddle_pos: Move RIGHT
1749
        15
               (action 2)
1750
                  - If the player paddle is already aligned with target_paddle_pos: NOOP (action 0)
1751
                    (to stabilize the paddle when it's in position)
                  Ensure stable movement to avoid missing the ball when close by.
        18
1752
1753
        20
                  Args:
                      target_paddle_pos (float): predicted x coordinate of the position to best
        21
1754
               position the paddle to catch the ball,
1755
                          and hit the ball to break brick wall.
                      obs (dict): Dictionary containing object states for "Player", "Ball", and blocks
        23
1756
               "{color}B" (color in [R/O/Y/G/A/B]).
1757
                          Each object has position (x,y), size (w,h), and velocity (dx,dy).
                  Returns:
                  int: 0 for NOOP, 2 for RIGHT, 3 for LEFT """
1758
        26
1759
                  if target_paddle_pos is None or 'Player' not in obs:
1760
        29
                      return 0
1761
        30
                  paddle = obs['Player']
        31
1762
                  paddle_x = paddle['x']
        32
                  paddle_w = paddle['w']
1763
                  paddle_center = paddle_x + (paddle_w / 2)
        34
1764
        35
1765
                  # Add deadzone to avoid oscillation
        36
        37
                  deadzone = 2
1766
        38
                  if abs(paddle_center - target_paddle_pos) < deadzone:</pre>
                      return 0 # NOOP if close enough
1767
        39
        40
                  elif paddle_center > target_paddle_pos:
1768
                     return 3 # LEFT
        41
1769
        42
                  else:
                      return 2 # RIGHT
1770
```

**Figure A.19:** Initial code for Breakout Agent (Part 2).

```
1785
1786
1787
1788
1789
1790
1791
1792
1793
         1@trace.model
         2 class Policy(Module):
1794
1795
               def __call__(self, obs):
                    pre_ball_x = self.predict_ball_trajectory(obs)
1796
                    target_paddle_pos = self.generate_paddle_target(pre_ball_x, obs)
1797
                    action = self.select_paddle_action(target_paddle_pos, obs)
                    return action
1798
               @trace.bundle(trainable=True)
               def predict_ball_trajectory(self, obs):
1800
                    """(same as before, skipped to save space)"""
if pre_ball_x is None or 'Ball' not in obs or 'Player' not in obs:
1801
                        return None
1802
                    ball = obs['Ball']
1803
                    paddle = obs['Player']
1804
1805
                    # Default to centering paddle on predicted ball position
                    target_x = pre_ball_x
         20
1806
1807
                    # Adjust paddle position based on current ball direction and brick locations
                    if ball['dy'] > 0: # Ball moving down
1808
                        if ball['y'] < 120: # Ball in upper half - aim for tunnels to high bricks
                             # Look for gaps in brick rows to target
1809
                             high_brick_x = None
1810
                             for color in ['RB', 'OB']: # Check red and orange rows
                                 if color in obs:
1811
                                     bricks = obs[color]
1812
                                      if len(bricks) > 0:
         30
                                          brick = bricks[0]
1813
         31
                                          high_brick_x = brick['x'] + (brick['w'] / 2)
         33
                                          break
1815
         34
                             if high_brick_x is not None:
         35
1816
                                 # Adjust paddle to deflect ball toward high-value bricks
if ball['x'] < high_brick_x:
   target_x = pre_ball_x - 4 # Hit ball on right side</pre>
1817
         38
1818
                                 else:
         39
                                     target_x = pre_ball_x + 4 # Hit ball on left side
1819
         40
         41
1820
         42
                    # Ensure target is within screen bounds
                    paddle_half_width = paddle['w'] / 2
1821
         43
         44
                    target_x = max(9 + paddle_half_width, min(152 - paddle_half_width, target_x))
1822
         45
1823
         46
                    return target_x
1824
```

Figure A.20: Learned code for Breakout Agent (Part 1).

```
1838
1839
1840
1841
1842
         1@trace.model
1843
         2 class Policy(Module):
1844
               # (continued from above)
1845
               @trace.bundle(trainable=True)
1846
               def generate_paddle_target(self, pre_ball_x, obs):
    """(same as before, skipped to save space)"""
1847
                    if 'Ball' not in obs:
1848
                        return None
         10
1849
                    ball = obs['Ball']
1850
         12
                    ball_x = ball['x']
1851
                    ball_y = ball['y']
         14
                    ball_dx = ball['dx']
ball_dy = ball['dy']
1852
         15
         16
                    # If ball is not moving, can't predict
1854
         18
                    if ball_dx == 0 and ball_dy == 0:
         19
1855
                        return None
         20
1856
         21
                    # Calculate time to reach paddle
1857
                    paddle_y = 189  # Player paddle y-position
if ball_dy != 0:  # Avoid division by zero
         23
1858
         25
                         time\_to\_paddle = (paddle\_y - ball\_y) / ball\_dy
1859
                        if time_to_paddle > 0: # Only predict if ball is moving toward paddle
    # Account for wall bounces
1860
                             x = ball_x + (ball_dx * time_to_paddle)
1861
                             num_bounces = 0
         29
1862
         30
                             while x < 9 or x > 152:
         31
                                 if x < 9:
1863
                                      x = 9 + (9 - x) # Reflect off left wall
         32
1864
         33
                                  if x > 152:
         34
                                      x = 152 - (x - 152) # Reflect off right wall
1865
         35
                                  num_bounces += 1
1866
                                  if num_bounces > 10: # Prevent infinite bounces
                                      break
1867
         39
                    return None
1869
               @trace.bundle(trainable=True)
               def select_paddle_action(self, target_paddle_pos, obs):
1870
                    """(same as before, skipped to save space)""
1871
                    if target_paddle_pos is None or 'Player' not in obs:
1872
1873
                    paddle = obs['Player']
                    paddle_x = paddle['x']
1874
                    paddle_w = paddle['w']
1875
                    paddle_center = paddle_x + (paddle_w / 2)
1876
                    # Add deadzone to avoid oscillation
1877
                    deadzone = 2
1878
                    if abs(paddle_center - target_paddle_pos) < deadzone:</pre>
                        return 0 # NOOP if close enough
1879
                    elif paddle_center > target_paddle_pos:
1880
         57
                        return 3 # LEFT
                    else:
1881
                        return 2 # RIGHT
1882
```

Figure A.21: Learned code for Breakout Agent (Part 2).

```
1893
1894
1895
         1 @trace.model
         2 class Policy(Module):
1896
1897
              def __call__(self, obs):
1898
                   shoot_decision = self.decide_shoot(obs)
                   move_decision = self.decide_movement(obs)
1899
                   return self.combine_actions(shoot_decision, move_decision)
1900
              @trace.bundle(trainable=True)
1901
              def decide_shoot(self, obs):
        10
1902
                   Decide whether to shoot based on enemy positions and existing projectiles.
1903
1904
        14
                   Args:
               obs (dict): Game state observation containing object states for "Player", "Shield0", "Shield1", "Alien0", "Alien1", etc.
Each object has position (x,y), size (w,h), and velocity (dx,dy).
        15
1905
1906
        16
                       Player bullets have negative dy velocity and alien bullets have positive dy
1907
               velocity
1908
        18
        19
                   Strategy tips:
                   - You can only have one missile at a time
        20
                   - Try to shoot when aliens are aligned with your ship
1910
        21
                   - Prioritize shooting at lower aliens as they're closer to you
1911
                   - Consider the movement of aliens when deciding to shoot
        23
1912
        24
        25
                   Returns:
1913
                   bool: True if should shoot, False otherwise
,,,
1914
1915
        29
                   # There can only be one player bullet on the field at a time
1916
        30
                   # Check for player bullets (which have negative dy velocity)
        31
                   for key, obj in obs.items():
1917
        32
                      if key.startswith('Bullet') and obj.get('dy', 0) < 0:
1918
        33
                           return False
        34
1919
        35
                   return random.choice([True, False])
1920
              @trace.bundle(trainable=True)
        37
1921
               def decide_movement(self, obs):
1922
        39
        40
                   Decide movement direction based on enemy positions and projectiles.
1923
        41
1924
                   Args:
                       obs (dict): Game state observation containing object states for "Player",
        43
1925
               "Shield0", "Shield1", "Alien0", "Alien1", etc.
        44
                       Each object has position (x,y), size (w,h), and velocity (dx,dy).
                       Player bullets have negative dy velocity and alien bullets have positive dy
        45
1927
               velocity
1928
                   Strategy tips:
1929
                   - Move to dodge enemy projectiles
                   - Position yourself under aliens to shoot them
1930
                   - Stay away from the edges of the screen
        50
1931
        51
                   - Consider moving toward areas with more aliens to increase score
1932
                   int: -1 for left, 1 for right, 0 for no movement
1933
1934
        55
1935
        57
                   player = obs['Player']
1936
                   return random.choice([-1.0.1])
1937
```

**Figure A.22:** Initial code for Space Invaders Agent (Part 1).

```
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
         1@trace.model
1958
         2 class Policy(Module):
1959
               # (continued from above)
1960
1961
               @trace.bundle(trainable=True)
               def combine_actions(self, shoot, movement):
1962
1963
                   Combine shooting and movement decisions into final action.
1964
                   Args:
                       shoot (bool): Whether to shoot
1965
                       movement (int): Movement direction
1966
1967
                   Action mapping:
                   - 0: NOOP (no operation)
- 1: FIRE (shoot without moving)
        16
1968
                   - 2: RIGHT (move right without shooting)
1969
                   - 3: LEFT (move left without shooting)
1970
                   - 4: RIGHT+FIRE (move right while shooting)
        20
                   - 5: LEFT+FIRE (move left while shooting)
1971
        21
1972
                   Returns:
                       int: Final action (0: NOOP, 1: FIRE, 2: RIGHT, 3: LEFT, 4: RIGHT+FIRE, 5:
1973
        24
               LEFT+FIRE)
1974
        25
1975
        26
                   if shoot and movement > 0:
                       return 4 # RIGHT+FIRE
                   elif shoot and movement < 0:</pre>
1977
        29
                       return 5 # LEFT+FIRE
        30
1978
                   elif shoot:
        31
                       return 1 # FIRE
1979
        32
                   elif movement > 0:
        33
1980
                       return 2 # RIGHT
        34
                   elif movement < 0:
    return 3 # LEFT</pre>
1981
        35
1982
                   return 0 # NOOP
1983
1984
```

Figure A.23: Initial code for Space Invaders Agent (Part 2).

```
1998
         1@trace.model
1999
         2 class Policy(Module):
2000
               def __call__(self, obs):
                    shoot_decision = self.decide_shoot(obs)
                    move_decision = self.decide_movement(obs)
2002
                    return self.combine_actions(shoot_decision, move_decision)
2003
2004
               @trace.bundle(trainable=True)
               def decide_shoot(self, obs):
    """(same as before , skipped to save space)"""
         10
         11
2006
                    # There can only be one player bullet on the field at a time
2007
                    # Check for player bullets (which have negative dy velocity)
         14
                    for key, obj in obs.items():
2008
                        if key.startswith('Bullet') and obj.get('dy', 0) < 0:</pre>
         16
2009
                             return False
2010
         18
                    player = obs['Player']
         19
2011
                    for key, obj in obs.items():
         20
2012
                        if key.startswith('Alien'):
                             # Check if alien is aligned with player (within 5 pixels)
if abs(obj['x'] - player['x']) < 5:
2013
         23
                                 # Prioritize lower aliens (higher y value)
if obj['y'] > 60: # Adjust this threshold as needed
2014
         25
2015
                                      return True
2016
                    return False
2017
               @trace.bundle(trainable=True)
         29
               def decide_movement(self, obs):
    """(same as before , skipped to save space)"""
    player = obs['Player']
2018
         30
         31
2019
         32
                    move = 0
2020
         33
         34
                    threat_left = 0
2021
         35
                    threat_right = 0
2022
         36
                    aliens_left = 0
         37
                    aliens\_right = 0
2023
         38
2024
         30
                    for key, obj in obs.items():
         40
                        if key.startswith('Alien'):
2025
         41
                            if obj['x'] < player['x']:</pre>
2026
         42
                                 aliens_left += 1
         43
                             else:
2027
         44
                                 aliens_right += 1
2028
         45
                        elif key.startswith('Bullet') and obj['dy'] > 0: # Enemy bullet
                            if obj['x'] < player['x']:</pre>
         46
2029
                                 threat_left += 1
                             else:
2030
                                  threat_right += 1
2031
2032
                    # Move away from threats
                   if threat_left > threat_right:
2033
                        move = 1
                    elif threat_right > threat_left:
2034
                        move = -1
2035
                    # If no immediate threat, move towards more aliens
2036
                    elif aliens_left > aliens_right:
                        move = -1
                    elif aliens_right > aliens_left:
2038
         60
                        move = 1
         61
2039
                    return move
         62
2040
         63
               @trace.bundle(trainable=True)
         64
2041
               def combine_actions(self, shoot, movement):
         65
2042
                    """(same as before , skipped to save space)"""
         66
                    if shoot and movement > 0:
         67
2043
                        return 4 # RIGHT+FIRE
         68
                    elif shoot and movement < 0:</pre>
2044
         69
                        return 5 # LEFT+FIRE
         70
2045
                    elif shoot:
         71
         72
                        return 1 # FIRE
2046
                    elif movement > 0:
         73
2047
         74
                        return 2 # RIGHT
2048
         75
                    elif movement < 0:
                       return 3 # LEFT
2049
                    return 0 # NOOP
2050
```

Figure A.24: Learned code for Space Invaders Agent.