

---

# Trace is the Next AutoDiff: Generative Optimization with Rich Feedback, Execution Traces, and LLMs

---

**Ching-An Cheng\***  
Microsoft Research  
chinganc@microsoft.com

**Allen Nie\***  
Stanford  
anie@cs.stanford.edu

**Adith Swaminathan\***  
Netflix  
aswaminathan@netflix.com

## Abstract

We study a class of optimization problems motivated by automating the design and update of AI systems like coding assistants, robots, and copilots. AutoDiff frameworks, like PyTorch, enable efficient end-to-end optimization of differentiable systems. However, general computational workflows can be non-differentiable and involve rich feedback (e.g. console output or user’s responses), heterogeneous parameters (e.g. prompts, codes), and intricate objectives (beyond maximizing a score). We investigate *end-to-end generative optimization* – using generative models such as LLMs within the optimizer for automatic updating of general computational workflows. We discover that workflow execution traces are akin to back-propagated gradients in AutoDiff and can provide key information to interpret feedback for efficient optimization. Formally, we frame a new mathematical setup, Optimization with Trace Oracle (OPTO). In OPTO, an optimizer receives an execution trace along with feedback on the computed output and updates parameters iteratively. We provide a Python library, Trace, that efficiently converts a workflow optimization problem into an OPTO instance using PyTorch-like syntax. Using Trace, we develop a general LLM-based generative optimizer called OptoPrime. In empirical studies, we find that OptoPrime is capable of first-order numerical optimization, prompt optimization, hyper-parameter tuning, robot controller design, code debugging, etc., and is often competitive with specialized optimizers for each domain. We envision Trace as an open research platform for devising novel generative optimizers and developing the next generation of interactive learning agents. Website: <https://microsoft.github.io/Trace/>.

## 1 Introduction

Computational workflows that integrate large language models (LLMs), machine learning (ML) models, orchestration, retrievers, tools, etc., power many state-of-the-art AI applications [1]: from chatbots [2], coding assistants [3], robots [4], to multi-agent systems [5]. However designing a computational workflow requires laborious engineering because many heterogeneous parameters (e.g. prompts, orchestration code, and ML hyper-parameters) are involved. Moreover, after deployment any erroneous behaviors of the workflow persist unless a developer manually updates it.

We study a class of optimization problems motivated by automating the design and update of computational workflows. Computational workflows produce optimization problems with heterogeneous parameters, rich feedback (e.g. console output and user’s verbal responses), and intricate objectives (beyond maximizing a score). Moreover, a workflow can have interdependent steps (e.g. adaptive orchestration, feedback control loops) and/or involve non-differentiable, semi-black-box, stochastic operations (e.g. ML models, simulations) whose behavior cannot be succinctly captured. As a result, the structure of the computation may change as the parameters and the inputs of the workflow vary.

---

\*Equal contribution

Due to its complexity, computational workflow optimization is usually framed as a black-box [6] or algorithm configuration [7] problem, and is tackled by general techniques like Bayesian Optimization [8], Evolutionary Algorithms [9], Reinforcement Learning (RL) [10] using scalar scores as feedback. But one observation of scalar feedback alone does not provide an improvement signal, so these algorithms are very inefficient when the parameter space is large (e.g. codes or natural language prompts). Recently LLM-based optimizers [11–16] have been proposed as *generative optimizers* to improve efficiency, leveraging the prior of generative models learned from large pre-training corpora to optimize complex prompts and codes. In this paper, we investigate how *generative optimization* can be applied more broadly and systematically to optimize a general computational workflow end-to-end. Appendix B discusses related works in generative optimization.

### 1.1 Toward Efficient End-to-End Optimization of Computational Workflows

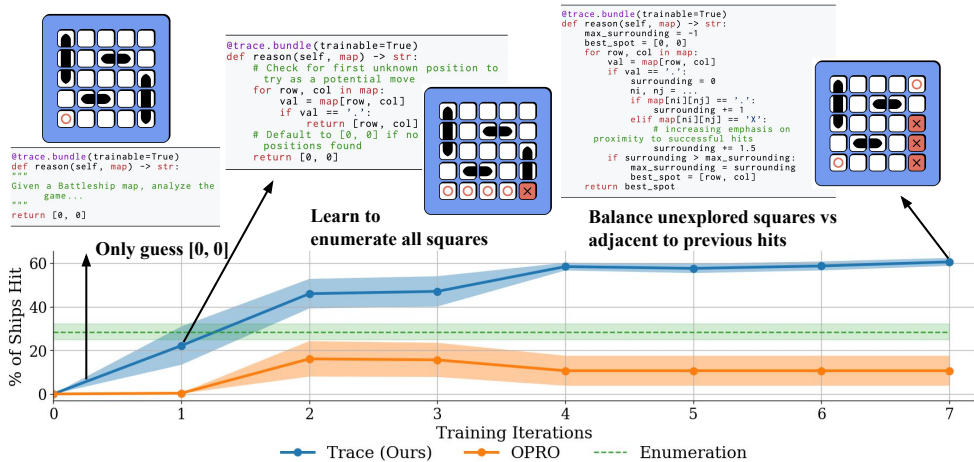
Inspired by back-propagation [17], we take an end-to-end approach to computational workflow optimization. AutoDiff frameworks [18, 19] have scaled the back-propagation algorithm to optimize differentiable workflows (e.g., neural networks) with billions of parameters. We extend the idea of AutoDiff and design a new framework, Trace, for jointly optimizing all of the heterogeneous parameters in general computational workflows, which may not be differentiable.

Trace treats a computational workflow as a graph like a neural network, where nodes are either inputs, parameters or the results of computation steps, and directional edges denote how nodes are created from others. But, instead of gradients, Trace propagates the *execution trace* of a workflow (which records the intermediate computed results and how they are used), via the notion of minimal subgraph (see Section 3.3). We show that propagating the execution trace subsumes back-propagation for differentiable workflows, and remains applicable even for non-differentiable workflows. Viewing a workflow as a graph and using its execution trace is standard for software engineering; for instance, human developers use such traces to debug distributed systems [20]. Our novel insight is that execution traces also enable end-to-end generative optimization, because they provide the information needed to relate rich feedback to the parameters in general workflows.

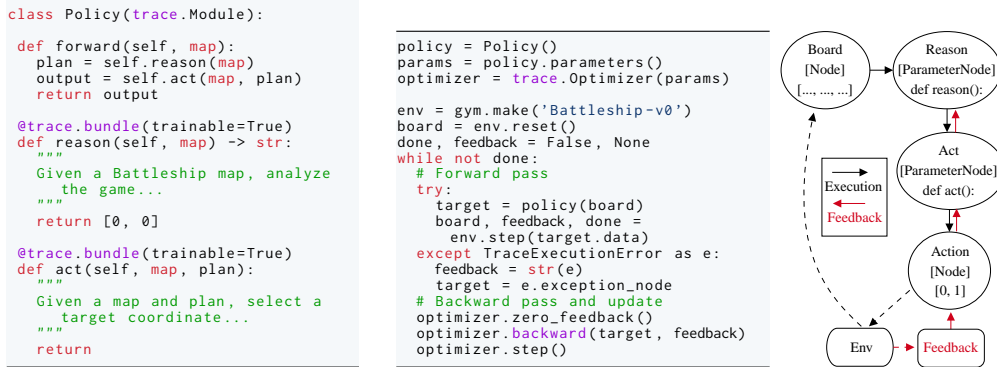
### 1.2 Example of Trace in Action

Trace is available as a Python library with an API inspired by PyTorch [19]. A user declares the parameters to be optimized using a `trainable` flag, decorates the workflow with node and bundle wrappers, and runs a Trace optimizer – just like how they would declare and train neural networks.

Consider building an AI agent for the Battleship game (Fig. 1). The agent’s policy (Fig. 2a) has two components (`reason` and `act`), which are chained together to react to different board configurations. The Battleship environment provides feedback (binary reward in texts) to tell if the agent’s action hit the hidden ships, and the goal is to hit all hidden ships as fast as possible. Consider how a human programmer might approach the problem. They may run the policy and change the code based on the



**Figure 1: Learning Example in Battleship.** Trace automatically optimizes heterogeneous parameters (multiple codes) to train an agent to Battleship. Means and standard errors are computed over 10 random seeds.



**Figure 2: Python Code of the Battleship Example.** To build a self-adapting agent with Trace, we only need to annotate some empty functions (reason, act) and set up an optimizer following PyTorch semantics. For space, we trim the docstrings of the empty functions with “...” and list them in Appendix J. Trace then builds a DAG as the workflow executes and updates the parameters (see Fig. 1 for the result).

observed feedback. They may rewrite the code a few times to try different heuristics to solve this problem. They will fix any execution errors (e.g., out-of-bounds exceptions) by using stack traces.

Our Trace framework accomplishes the programmer’s goal automatically without adding complexity to the Python code. The user declares reason and act as trainable (Fig. 2a) and then runs the agent in a PyTorch-like training loop (Fig. 2b). During the execution, Trace records a directed acyclic graph (DAG) (Fig. 2c) and uses it to compute the execution trace for optimization. Trace also automatically catches errors (e.g., syntax/semantic errors) and can use them as feedback. In Fig. 1, we show what the agent learns as Trace optimizes<sup>2</sup> its policy, where the learned policy is evaluated on new randomly generated games for generalization. The agent can quickly improve its performance and learn strategies that are increasingly complex. Remarkably, there is no mention of the specific Battleship environment API, nor details on how the functions reason and act should behave or adapt in Fig. 2a. The Trace optimizer figures out all the details through interactions as the computational graph unfolds and the feedback on the output is observed. Beyond code as parameters in this example, we also have experiments in Section 5 where prompts and other heterogenous parameters are optimized.

### 1.3 A New World of Optimization

The design of Trace is based on a new mathematical setup of iterative optimization, which we call Optimization with Trace Oracle (OPTO). In OPTO, an optimizer selects parameters and receives a computational graph as well as feedback on the computed output. Trace is a tool to efficiently convert the optimization of computational workflows into OPTO problems in practice.

We argue that framing computational workflow optimization as OPTO can lead to faster convergence than a black-box approach. We present a constructive proof: We design a general-purpose efficient generative optimizer called OptoPrime, for OPTO. OptoPrime turns OPTO to a sequence of pseudo-algorithm problems. In each iteration of OPTO, we format the execution trace and output feedback as a pseudo-algorithm question and present it to an LLM for solution (GPT-4 using a ReAct-CoT prompt listed in Appendix G). In experiments, we apply OptoPrime unchanged to many disparate applications like prompt optimization, first-order numerical optimization, hyper-parameter tuning, and robot controller design. We find that the general purpose OptoPrime is competitive with specialized optimizers for each domain, e.g., achieving 10% higher accuracy on BigBenchHard [21] when optimizing a DSPy [22] program compared to their hand-designed optimizer.

Trace, OPTO, and OptoPrime together provide the first tractable algorithm for optimizing general computational workflows end-to-end. The Trace framework *a)* leverages a workflow’s structure and *b)* can incorporate rich feedback beyond scores, extending AutoDiff to complicated, non-differentiable computational workflows. With Trace, we conjecture that “training deep agent networks” (which fluidly mix computation of tensors, LLMs, and other programmable tools) will soon be possible.

<sup>2</sup>We use a new general-purpose LLM-based optimizer OptoPrime that we detail in Section 4.

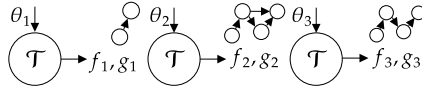
## 2 Optimization with Trace Oracle

OPTO is the foundation of Trace. In this section, we define this graph-based abstraction of iterative optimization and discuss how OPTO covers various computational workflow optimization problems.

**Preliminary** We review the definition of a computational graph (see Fig. 2c). A computational graph  $g$  is a DAG, where a node represents an object (such as tensors, strings, etc.) and an edge denotes an input-output relationship. We call a node without parents a root and a node without children a leaf, which are the inputs and outputs of the computational graph. In the context of optimization, some inputs are marked as trainable *parameters*, which are denoted as  $\{X_\theta\}$ . For a node  $X$ , its parents are the inputs to an operator that creates  $X$ . The descendants of node  $X$  are those that can be reached from  $X$  following the directed edges; the ancestors are defined conversely. Without loss of generality, we suppose that all computational operators have a unitary output<sup>3</sup>. In this way, we can associate the operator that creates the child node with the child node, and the full computation can be represented compactly as a DAG without explicitly representing the operators. The execution trace of a computational workflow is defined as the sequence of operations and their execution results invoked when computing the output from a set of inputs; execution traces can be expressed as a computational graph as defined above.

### 2.1 Problem Definition of OPTO

OPTO is an *abstract* setup of iterative computational workflow optimization. An OPTO problem instance is defined by a tuple  $(\Theta, \omega, \mathcal{T})$ , where  $\Theta$  is the parameter space,  $\omega$  is the context of the problem, and  $\mathcal{T}$  is a Trace Oracle. In each iteration, the optimizer selects a parameter  $\theta \in \Theta$ , which can be heterogeneous. Then the Trace Oracle  $\mathcal{T}$  returns a *trace feedback*, denoted as  $\tau = (f, g)$ , where  $g$  is the execution trace represented as a DAG (the parameter is contained in the root nodes of  $g$ ), and  $f$  is the feedback updated to exactly one of the output nodes of  $g$ . Finally, the optimizer uses the trace feedback  $\tau$  to update the parameter according to the context  $\omega$  and proceeds to the next iteration, as shown in Fig. 3.



**Figure 3: Iterations of OPTO.** When  $\theta \in \Theta$  is selected, the Trace Oracle  $\mathcal{T}$  returns trace feedback  $\tau = (f, g)$ , where  $g$  is a computational graph using  $\theta$  as input and  $f$  is the feedback given to the output of  $g$ .

In OPTO, the output feedback  $f$  is generic, such as scores, gradients, hints/explanations expressed in natural language, and console messages. The context  $\omega$  provides invariant information to interpret the output feedback  $f$  as well as any known side-information, e.g., desired properties of the parameters. The context  $\omega$  is fixed for an OPTO problem instance (similar to an instruction, or a problem definition), whereas the output feedback  $f$  can change with the parameter  $\theta \in \Theta$  and the resulting computation  $g$ . For example,  $\omega$  may be “Minimize a loss function” and  $f$  is a loss. Alternatively,  $\omega$  can be open-ended, like “Follow the feedback” and  $f$  describes how an output should be changed. In Section 3.2, we discuss how to define the context and output feedback when constructing OPTO problems in practice. In this paper, we focus on OPTO problems where  $f$  and  $\omega$  can be expressed compactly in text. This covers a wide range of problems [23], including those with scalar feedback.

OPTO differs from a black-box setup in that the execution trace  $g$  shows the computational path toward the output, which provides information to construct a parameter update direction from  $f$  and  $\omega$ . In the minimization example above, when the execution trace  $g$  is missing, it is unclear how the parameter can be improved given only a point evaluation of  $f$ . On the other hand, with  $g$  documenting the functions used to create the output, an update direction (e.g., a gradient) can be derived. We highlight that the structure of the computational graph  $g$  returned by the Trace Oracle  $\mathcal{T}$  can be different in each iteration in the general case (as in Fig. 3) because some workflows can change with different inputs and parameters.

To ground the OPTO setup, we show how OPTO is related to some existing problems with examples. We discuss other examples like hyperparameter tuning and multi-agent systems in Appendix C.

**Example 1** (Neural network with back-propagation). The parameters are the weights.  $g$  is the neural computational graph and  $f$  is the loss. An example context  $\omega$  can be “Minimize loss”. The back-propagation algorithm can be embedded in the OPTO optimizer, e.g., an OPTO optimizer can use  $\tau$  to compute the propagated gradient at each parameter, and apply a gradient descent update.

<sup>3</sup>A multi-output operator can always be modeled by a single-output operator and single-output indexers.

**Example 2 (RL).** The parameters are the policy.  $g$  is the trajectory (of states, actions, rewards) resulting from running the policy in a Markov decision process; i.e.,  $g$  documents the graphical model of how a generated action is applied to the transition, which then returns the observation and reward.  $f$  can be a success flag.  $\omega$  can be “Maximize return (cumulative rewards)” or “Maximize success”.

**Example 3 (Prompt Optimization of an LLM Agent).** The parameters are the prompt of an LLM workflow.  $g$  is the computational graph of the agent and  $f$  is the feedback about the agent’s behavior (which can be scores or natural language).  $\omega$  can be “Maximize score” or “Follow the feedback”.

### 3 Trace: The Next AutoDiff

We design a framework, Trace, to bring OPTO from an abstract concept to reality. Trace provides a lightweight Python<sup>4</sup> tool to implement the Trace Oracle of OPTO for optimizing computational workflows. Through the OPTO framing, Trace separates the design of optimizers and domain-specific components so that general-purpose optimizers can be built that work across diverse domains.

#### 3.1 Design of Trace

Trace is designed based on two primitives:

- `node` is the wrapper of Python objects. When wrapped, a Python object is registered as a unique node in the computational graph of Trace. A node can be set `trainable`, which would make the node a parameter in OPTO. In addition, when using `node` to declare a parameter, one can also describe (in natural language) constraints that the parameter should obey.
- `bundle` is the decorator to turn Python methods into operators. When a function is decorated, its docstring and source code are recorded as the definition of the operator, which infer how the output feedback should change the parameters. Moreover, functions decorated by `bundle` can be set `trainable` as well, which means that the code of the decorated method becomes a parameter.<sup>5</sup>

For any workflow, using Trace involves the following steps (see Fig. 2). First, the user declares the workflow’s parameters using `node` and `bundle`, and defines the workflow’s conceptual blocks as operators using `bundle`. Then the user creates an OPTO optimizer (such as `OptoPrime` in Section 4), and optionally provides the context  $\omega$  for the problem. (A default context  $\omega$  of `OptoPrime` is “Follow the feedback”). In addition, the user defines a mechanism to provide feedback to the computed result (e.g., scores, natural language suggestions, etc.), in analogy to defining a loss function in neural network training. After the declaration, optimization via Trace repeats the following: 1) Execute the decorated workflow. As it runs, a DAG is built in the backend, logging the computed results and their connections. 2) Initiate the propagation of the output feedback to the parameters by calling `backward`. (Any execution error is also treated as feedback; see Appendix D.) Internally, Trace extracts the minimal subgraph  $g$  connecting the parameters and the output and sends the OPTO optimizer the trace feedback  $\tau = (f, g)$ . 3) Call the OPTO optimizer’s `step` method to update the parameters.

#### 3.2 Using Trace Primitives for Effective Execution Tracing

There are many ways to represent a computational workflow as a computational graph, from abstracting the entire process as one big operator to listing all low-level steps as operators in the graph. In Trace, the level of abstraction is decided by how `bundle` is applied, as all steps underneath `bundle` are abstracted as one operator. The design of `bundle` allows tracing most Python codes, except for those modifying the content of an object reference in place. However, such a case can be avoided by first duplicating the object and then applying the modification to the copied object, similar to how a recurrent neural network is implemented.

Different abstraction choices trade off the graph complexity and the description needed for each operator. Abstracting everything into a single operator makes a trivial graph but requires more descriptions to faithfully capture the workflow. On the other hand, not all details matter in optimization, so exposing every low-level operator can make the graph unnecessarily cluttered. The best representation depends on the application and OPTO optimizer at hand. This problem, we believe, is similar to the design of neural network architectures. Here, we suggest defining the operators to

<sup>4</sup>The current implementation is in Python. The DAG-based design of Trace can be adapted to other programming languages also, just like AutoDiff libraries that are available in C++.

<sup>5</sup>This would add an extra parent (i.e., the trainable code) to the computed child node.



mimic the whiteboard system diagram of the computational workflow. This level of abstraction in our experiments strikes a good balance between the ease of documentation and the graph size.

Apart from architecture design, another design question is what information goes into the context  $\omega$  versus the description of each operator? For a *single* problem, there is no difference in principle; one can choose to provide details of all operators in  $g$  through the context  $\omega$ . However, this will require manually crafting a context for every workflow. We suggest instead providing a description of the operators when they are defined using `bundle`. Then Trace will automatically generate the workflow-specific information while the same context  $\omega$  is shared across *many* workflows.

### 3.3 Backward Feedback Propagation: Realizing the Trace Oracle of OPTO

Trace uses a recursive graph traversal algorithm (Algorithm 1) to propagate feedback in the reversed topological ordering. With different propagators, Algorithm 1 can implement various forward-backward schemes including back-propagation.<sup>6</sup> We propose a general propagator, Minimal Subgraph Propagator (MSP), in Algorithm 2. MSP propagates the trace feedback  $\tau = (f, g)$ , where  $g$  is implemented as a priority queue. Running Algorithm 1 with MSP (Algorithm 2) together implements the Trace Oracle of OPTO, which extracts the *minimal subgraph* between parameter nodes and output.<sup>7</sup> connecting the parameters and an output. Appendix E proves the following:

**Theorem 1.** *For a graph with  $N$  nodes and maximum degree  $W$ , Algorithms 1 and 2 have time complexity  $O(WN^2 \log N)$  and space complexity<sup>8</sup>  $O(WN)$ .*

By contrast, back-propagation has time and space complexities of  $O(WNd^2)$  and  $O(d)$  respectively, where  $d$  is the maximal dimension of tensors. The difference is because in the most general setting of computational graphs and feedback, the propagated feedback (no matter how it is represented) does not have a constant size and needs the full subgraph.

**Theorem 2.** *For generic computational graphs of  $N$  nodes, in the worst case, the propagated feedback needs a description length  $\Omega(N)$  to construct an improvement direction.*

Despite the worst case complexity of MSP, in practice the difference is negligible. Since MSP only involves merging priority queues of references, most actual computation happens in the forward pass (and also the optimizer’s step method). However for very large problems with millions of nodes in the minimal subgraph, we anticipate that computational issues of MSP could arise.

## 4 Design of the First OPTO Optimizer

We introduce an LLM-based generative optimization algorithm OptoPrime for any text-based OPTO problem. We believe that this is one of many possible optimization algorithms for these problems and there is a large space to be explored for identifying efficient optimization methods for OPTO.

**Subgraph Representation** One core challenge of designing an LLM-based OPTO optimizer is how to represent the execution trace subgraph  $g$  (which can involve various graph structures and heterogeneous data) to LLMs, in a way that LLMs can understand and reason about the downstream effects of parameter updates. We leverage the LLMs’ remarkable coding and debugging ability [3].

<sup>6</sup>In back-propagation, the message is the gradient  $\nabla_i$  and the propagate function returns  $J_i^\top \sum_j \nabla_j$  to its  $i$ th parent, where  $J_i$  is the Jacobian to the  $i$ th parent and the  $\nabla_j$  gradient received from the  $j$ th child.

<sup>7</sup>The *minimal subgraph*  $g_{\mathcal{X} \rightarrow Y}$  between nodes  $\mathcal{X}$  and a node  $Y$  is defined as  $g_{\mathcal{X}, Y} := \mathcal{X} \cup \{Y\} \cup \{Z, P \mid Z \in \text{ancestors}(Y), Z \in \text{descendants}(X), X \in \mathcal{X}, P \in \text{parents}(Z), P \in \text{parents}(Y)\}$ .

<sup>8</sup>The space complexity refers to the extra space needed for the backward pass, not including the forward pass.

---

#### Algorithm 1 Backward Message Passing

---

**Input:** Node  $output$ , feedback  $f$ , propagator  $P$

```

1:  $\tau \leftarrow P.\text{init}(f)$ 
2:  $output.\text{add\_feedback}(\text{"User"}, \tau)$ 
3:  $queue \leftarrow \text{MinHeap}([output])$ 
4: while  $queue$  is not empty do
5:    $node \leftarrow queue.\text{pop}()$ 
6:    $feedback \leftarrow P.\text{propagate}(node)$ 
7:   for  $parent$  in  $node.\text{parents}$  do
8:      $\tau \leftarrow feedback[parent]$ 
9:      $parent.\text{add\_feedback}(node, \tau)$ 
10:    if  $parent \notin queue$  then
11:       $queue.\text{push}(parent)$ 

```

---



---

#### Algorithm 2 Minimal Subgraph Propagator

---

**Input:** A child node  $node$

*// The pseudo code implements propagate.*  
*//  $\text{init}(f)$  returns  $(f, \{\})$ .*

```

1:  $g \leftarrow \{node\} \cup \{parent \text{ in } node.\text{parents}\}$ 
2: for  $(f_i, g_i)$  in  $node.\text{feedback}$  do
3:    $g \leftarrow g \cup g_i$ 
4:    $f \leftarrow f_i$  // all  $f_i$  are the same.
5: return  $\{p : (f, g) \text{ for } p \text{ in } node.\text{parents}\}$ 

```

---

We present the trace feedback  $(f, g)$  computed by Trace as a pseudo-algorithm problem: the subgraph  $g$  is expressed as a report of code execution with information about the computed values and descriptions of operators in  $g$ . Then we prompt the LLM to update the parameters in  $g$  based on feedback  $f$  given to the output. Fig. 4 shows an example. It is crucial to note that even though the lines look like an actual program, it is not the real program but the computational graph defined by `bundle` (see Section 3.2).

**Parameter Update** We prompt the LLM with a ReAct-CoT style prompt (Appendix G.2) in one query, asking it to generate reasoning based on the graph, and a suggestion on the parameter changes. If the suggestion can be extracted from the response, we update the parameters.

**Optimization Memory** OptoPrime optimizes most workflows reasonably well using just the traced graph and feedback, but it can run into issues when single feedback alone is not informative (e.g., the output feedback is rewards, but there is no description of how the rewards are generated). For robustness, we have a basic memory module in OptoPrime, which tracks the past parameter-feedback pairs as in-context examples. See Appendix G for details.

```
#Code:
a = bar(x)
y = add(b, a)
z = mul(a, y)
#Definitions:
[mul] This is a multiply operator.
[add] This is an add operator.
[bar] This is a method that does
      negative scaling.
#Inputs:
b=1.0
#Others:
a=2.0
y=3.0
#Output
z=6.0
#Variable
x=-1.0
#Feedback:
Output should be larger.
```

**Figure 4:** An example pseudo-code report generated by Trace for a program of  $x = \text{Node}(-1.0)$ ;  $z = \text{bar}(x) * (\text{bar}(x)+1)$  and the objective of  $\max_x z$ .

## 5 Experiments

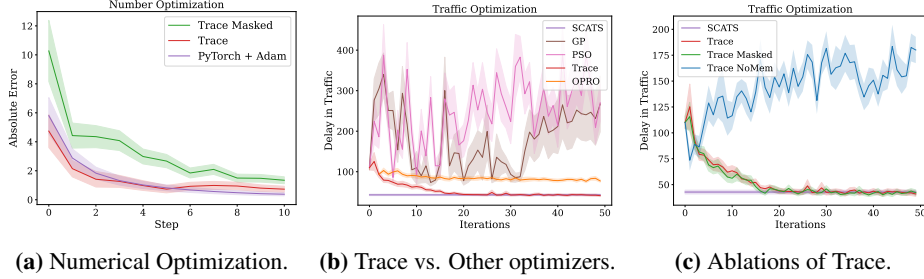
We evaluate the Trace framework with OptoPrime. We implement the state-of-the-art LLM optimizer OPRO [13] as a baseline; in comparison with OptoPrime, OPRO does not use the execution trace but relies on the memory of parameter and feedback pairs. For these experiments, we use GPT-4-0125-Preview. We run the experiments on a standard PC with 16 GB RAM, and Trace introduces no measurable overhead on executing the workflow. We also conduct experiments to compare Trace and OptoPrime with a concurrent AutoDiff-like framework, TextGrad [24], which was released after Trace was submitted to NeurIPS. We show that TextGrad can be easily implemented as an optimizer in Trace, and OptoPrime achieves similar or better performance than TextGrad while using much less computation time. In the rest of this section, we will denote Trace+OptoPrime simply as Trace. We report the token usages of all approaches in all experiments in Appendix I.1.

### 5.1 Validating with Numerical Optimization

First, we want to validate if OptoPrime can solve classical differentiable optimization problems, since they are a special case of OPTO. Consider the problem of  $\min_x |h(x) - y^*|$  for a target  $y^*$ . We construct a synthetic task environment that randomly creates  $y^*$  and the computational graph of  $h$  with arbitrarily complex connections between numerical variables (see Appendix I.3 for details). We evaluate OptoPrime (denoted as Trace) and a variant that does not see the graph (Trace Masked); both the optimizers do not use memory. The output feedback is “The output should be <larger/smaller>” (this feedback has the same information as the gradient w.r.t.  $h$ ). We compare also the performance of Adam optimizer [25]. We run 30 trials over different randomly generated problems. All methods see the same randomness. Trace is able to match the best-in-class Adam; on the other hand, without access to the full computational graph, the feedback-alone optimizer struggles to find  $x^*$  (Figure 5a).

### 5.2 Tuning Hyperparameters to Orchestrate Complex Systems

We tested Trace in a traffic control problem, which is an instance of hyper-parameter tuning. We used UXSim [26] to simulate traffic at a four-way intersection, where the trainable parameters are 2 integers in  $[15, 90]$ , which are the green light duration for each direction of traffic flow. The feedback is the estimated delay experienced by all vehicles due to intersections, and the goal of an optimizer is to minimize the delay using the fewest number of traffic simulations. To this end, this optimizer must find the right trade-off for temporally distributed and variable demands. In Fig. 5 we report the performance of a SOTA heuristic from the traffic control literature, SCATS [27] as well as two black-box optimization techniques: Gaussian Process Minimization (GP) [8] and Particle Swarm Optimization (PSO) [28]. All methods use the same starting parameters. Trace denotes OptoPrime using memory, and Trace NoMem denotes OptoPrime without memory. We report further details in



**Figure 5:** Numerical Optimization and Traffic Optimization Results.

Appendix I.4. GP and PSO appear bad because 50 iterations are insufficient for their convergence; given enough iterations, both will eventually perform well. Trace is quickly competitive with the SCATS heuristic, whereas OPRO is not. Moreover, we find that memory is crucial for Trace to perform well for this task. But we note that Trace consumes extra overhead compared to other methods, since Trace has to materialize the resulting computation graph and query an LLM with effectively a longer prompt than that of OPRO.

### 5.3 Unifying Prompts and Functions Optimization

Many LLM agents today, e.g., specified by LangChain [29] and DSPy [22], have many components. These libraries provide optimization tools to optimize a small portion of their workflows, predominantly the prompt that goes into an LLM call. However, for building self-adapting agents that can modify their own behavior, only allowing the change to one part of a workflow but not others can be limiting. In this experiment, we test Trace’s ability in joint prompt optimization and code generation. Specifically, we optimize a given DSPy-based LLM agent and tune its three components: the meta-prompt `prompt_template`, a function `create_prompt` that modifies the prompt with the current question, and a function `extract_answer` that post-processes the output of an LLM call.

We set up an end-to-end prompt-and-code optimization pipeline. We use an automatic evaluation function to compare the LLM’s output with the ground truth, which requires the LLM agent to generate outputs not only with the correct answer but also in the correct format (following the guidelines of [30]). We use the Big-Bench Hard dataset [21] (15 examples for training, 5 for validation, and the rest for testing). We compare Trace with DSPy’s COPRO module (which optimizes the meta-prompt). In Table 1, we show that Trace is able to optimize a DSPy program beyond what DSPy’s COPRO optimizer can, especially on algorithmic tasks. This result shows how Trace can concretely improve upon existing LLM prompt optimization libraries. We show learned codes in Appendix J.

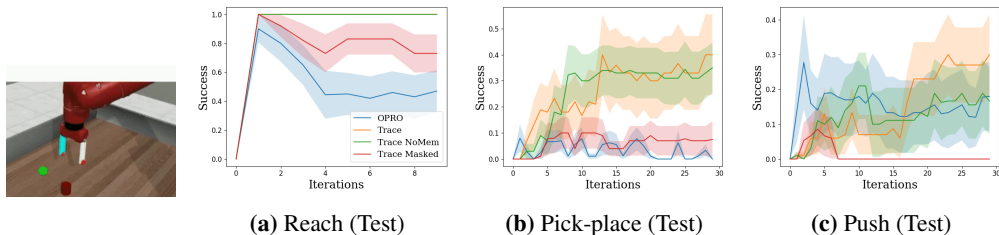
	BBH all (23 tasks)	NLP (12 tasks)	Algorithmic (11 tasks)		BBH all (23 tasks)	NLP (12 tasks)	Algorithmic (11 tasks)
0-shot				0-shot			
DSPy	41.6	53.8	32.6	DSPy + CoT	70.4	73.7	68.0
DSPy-PO	55.3	69.0	45.2	DSPy-PO + CoT	71.6	73.9	70.0
Trace	<b>59.5</b>	<b>70.9</b>	<b>51.1</b>	Trace + CoT	<b>78.6</b>	<b>75.8</b>	<b>80.6</b>

**Table 1:** End-to-end workflow optimization for an LLM benchmark (Big-Bench Hard) in 0-shot setup. CoT refers to Chain-of-Thought prompting and PO refers to DSPy’s own prompt optimizer (COPRO). We use Trace to optimize a DSPy program, starting from the same program and prompt template specified by DSPy.

### 5.4 Long-Horizon Robot Manipulator Control

We test the ability of Trace to optimize long-horizon workflows with complex dependencies and to “back-propagate through time”. We use Trace to train the controller code (in Python) for a simulated Sawyer robot manipulator. We use the Meta-World environment from LLF-Bench [23] as the simulator and consider three tasks: Reach, Pick-place, and Push. For each task, LLF-Bench provides a task instruction and meaning of the action space, which we use as the context  $\omega$  of the OPTO problem. The observation is a Python dictionary of vectors, indicating the end-effector position, the goal position, the gripper status, etc. The action space is a 4-dimensional vector to control the relative position of the end-effector and the gripper state. In each time step, the LLF-Bench Meta-World simulator returns the observation and natural language feedback to guide the robot. An





**Figure 6:** Learning the feedback control policy (code) for a simulated Sawyer manipulator in LLL-Bench Meta-World. In each iteration (x-axis), one episode of rollout (10 steps) is performed, and then the policy is updated. The mean and standard error of the success rate over 10 seeds are shown.

episode ends if the robot successfully solves the problem or because of time-out. We consider an episodic training setting. The initial condition for all iterations in training is the same. We evaluate the learned policy in terms of success, starting from 10 held-out initial conditions. The task horizon is 10 steps, which is sufficient for task completion, and each training iteration has one rollout. The output feedback in OPRO is a textual representation of task success. In addition to the controller code, we also decorated the reset and step functions of the gym environment so that the entire rollout can be traced end-to-end. We compare Trace with OPRO; because of the streaming OPRO setting, our OPRO implementation only proposes one candidate in each iteration, which is then evaluated and provided with the output feedback.

The experimental results are summarized in Fig. 6. Trace denotes OptoPrime using memory, and Trace NoMem denotes OptoPrime without memory. We show learned code in Appendix J. OptoPrime is clearly the top-performing optimizer, especially the version with memory. OPRO is able to solve Reach at the start, but its performance degraded over iterations (this instability was mentioned in [13]) and gets a similar performance as OptoPrime (without memory) in Push. To validate that the performance of OptoPrime is indeed due to using the execution trace, we include an ablation where we mask out the execution trace, which leads to a significant decline in performance and stability. This experiment features the most complex graph structures among all the experiments. The experimental results here are quite impressive, showing that Trace is able to learn a sophisticated control logic in dozens of interactions, not only working on the training initial conditions but also on held-out testing ones too. We discuss some limitations in Appendix I.6.

## 5.5 Comparison with TextGrad

After the submission of our work to NeurIPS, another AutoDiff-like framework, TextGrad [24], was released, which shares the same goal of end-to-end optimizing AI workflows as Trace. In comparison, TextGrad propagates text feedback, whereas Trace propagates minimal subgraphs (see Section 3.3). The graph-based design of Trace, which separates the tracing infrastructure and optimization algorithms, makes it more flexible. In fact, we easily implemented TextGrad as an optimizer in the Trace framework, but the reverse is not possible (because TextGrad couples the infrastructure and the optimization algorithm together). In addition, unlike TextGrad, Trace supports jointly optimizing heterogeneous parameters and can be applied to directly trace a given computational workflow without the need to rewrite the workflow using pre-defined templates. Please see Appendix H for more discussion comparing the two frameworks.

In this experiment, we apply Trace to directly decorate the evaluation code released with the TextGrad library and optimize the parameters following their training/evaluation pipeline line-by-line. This experimental design makes the comparison fair by allowing each optimizer to access the same LLM APIs around the same time, and showcases the flexibility of Trace framework to optimize any computational workflow. We pick the Solution optimization [24, Table 2] and Prompt optimization [24, Table 3] for the reasoning tasks experiments. Please see [24] for details on the exact setup. We compare OptoPrime, TextGrad<sup>9</sup>, and a reimplementation<sup>10</sup> of TextGrad as an optimizer in Trace. We find that all these algorithms achieve similar success rates in these experiments. One noticeable difference is that OptoPrime is about 3x faster wall-clock time than TextGrad since OptoPrime makes a single call to LLM in each optimization step, whereas TextGrad calls linear to the graph’s size.

<sup>9</sup>The numbers in the original paper cannot be reproduced exactly despite using the released TextGrad code.

<sup>10</sup>All algorithms in Trace see the same graph, which is slightly different from the graph of TextGrad.

	OptoPrime (Trace)	Time	TextGrad (24-10-30)	Time	TextGrad (Trace)	TextGrad (Reported)
MMLU-Machine Learning	<b>86.6</b> (0.2)	1.7 (0.6)	86.1 (0.5)	3.5 (1.1)	86.3 (0.2)	88.4
MMLU-College Physics	<b>94.1</b> (0.8)	1.2 (0.3)	93.1 (0.7)	2.3 (0.4)	93.3 (0.6)	95.1
Google-proof QA	<b>59.6</b> (1.3)	12.2 (1.4)	53.2 (0.6)	19.5 (1.9)	54.0 (0.7)	55.0
BBH Counting	<b>89.4</b> (0.1)	55.9 (4.5)	89.2 (1.2)	142.9 (9.3)	87.6 (1.7)	91.9
BBH Word Sorting	71.6 (3.1)	82.5 (10.1)	<b>72.0</b> (0.4)	211.1 (16.8)	71.4 (2.5)	79.8
GSM8K	<b>82.5</b> (0.1)	—	82.4 (0.6)	—	82.0 (0.2)	81.1

**Table 2: Comparison between Trace and TextGrad.** The optimizer is GPT-4o-2024-08-06, and the student model is GPT-35-turbo-1106. The results show the mean and the standard error of success rate of the last iterate computed by 5 seeds. The experiment time reported is in minutes (the time involves not just training but also validation and testing by running TextGrad’s original pipeline); the time of GSM8K experiment is omitted as the experiment time (>8hrs) is determined primarily by the evaluation not optimization.

## 6 Limitations

We highlight that Trace, OPTO and OptoPrime are a first step towards end-to-end generative optimization and building self-adapting workflows. They have limitations in their current form. OPTO captures rich feedback, but it is important to specify a solution concept as well as the feedback source. We provide guidance for feedback design in Section 3.2 and discuss notions of optimality in Appendix F. We believe designing feedback will be as important as designing loss function in deep learning, both of which are open research questions. Also, Trace cannot convert all computational workflows into OPTO problems, e.g., stateful functions that modify their state in place cannot be represented as a DAG without modification, and distributed/parallel computing workflows are incompatible with the current implementation (though in theory Trace can run in an asynchronous way so long as the overall graph does not end up with cycles). Finally, while Trace is designed to be generic and future-proof, the OptoPrime optimizer is preliminary. Although we demonstrated that OptoPrime could work well with moderate-size graphs, it is not a provably optimal algorithm and uses more tokens than OPRO, though, in our experiments, OPRO’s performance does not improve even when given a large token budget. The debugging ability and context limits of the LLM used in OptoPrime crucially determine the scale of problems that we can practically address today. Consequently, more research is needed for designing token-efficient generative optimization algorithms.

## 7 Conclusion and Future Work

We created Trace that can convert a computational workflow optimization problem into an OPTO problem, and we demonstrated a tractable OPTO optimizer, OptoPrime. This is just a first step towards a new paradigm of end-to-end *generative optimization*, with exciting avenues for future work. We discuss a few selected ones below. Please see Appendix A for a longer discussion.

In OptoPrime, we connect optimization to an LLM’s reasoning capability. Techniques that have been proposed to improve LLM reasoning, e.g. Chain-of-Thought [31], Few-Shot Prompting [32], Tool Use [33], and Multi-Agent Workflows [5] could also help improve OptoPrime or design new OPTO optimizers. We conjecture that a hybrid workflow of LLM and search algorithms, can enable a truly general-purpose OPTO optimizer. Along the way, we must settle on how to delineate the agent vs. the optimizer. How to trade off the generality of optimizer vs. crafting side-information in the context  $\omega$  to achieve task-specific performance is an open question.

In Trace, we chose a specific propagator (MSP), which maximally preserves information for a general computation graph. We can instead specialize it for specific computations, e.g. to accommodate very large graphs with a hierarchical graph representation. Going a step beyond the basic memory module we experimented with in OptoPrime, we anticipate that an optimizer that can reason about how a workflow will behave under counterfactual parameter settings can be more efficient than OptoPrime and can enable a divide-and-conquer approach to OPTO. More research is needed to study the theoretical properties of OPTO (such as optimization landscape and complexity). We hope our preliminary effort in Appendix F can provide some guidance. Finally, in this paper, we focused on output feedback and context that can be compactly textualized. We anticipate that computational workflows with rich non-textual contexts and output feedback will also benefit from automatic generative optimization through appropriate applications of Trace (e.g., with VLMs).

## Acknowledgments and Disclosure of Funding

We would like to thank John Langford, Ahmed Awadallah, Jennifer Neville, Andrey Kolobov, Ricky Loynd and Paul Mineiro for thought-provoking discussions. We would also like to thank Tobias Schnabel, Ruijie Zheng, Wanqiao Xu, and Kaiwen Wang for their valuable feedback on an early draft of this manuscript. Additionally, we thank Anqi Li, Omar Khattab, David Hall, Yifan Mai, Bryan He, Yash Chandak, Emma Brunksill, and Dawen Liang for their suggestions and feedback. The work was partially done while Adith Swaminathan was at Microsoft Research.

## References

- [1] Matei Zaharia, Omar Khattab, Lingjiao Chen, Jared Quincy Davis, Heather Miller, Chris Potts, James Zou, Michael Carbin, Jonathan Frankle, Naveen Rao, and Ali Ghodsi. The shift from models to compound AI systems. <https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/>, 2024.
- [2] 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.
- [3] Albert Ziegler, Eirini Kalliamvakou, X. Alice Li, Andrew Rice, Devon Rifkin, Shawn Simister, Ganesh Sittampalam, and Edward Aftandilian. Productivity assessment of neural code completion. In *SIGPLAN International Symposium on Machine Programming*, page 21–29, 2022.
- [4] Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as policies: Language model programs for embodied control. In *IEEE International Conference on Robotics and Automation*, pages 9493–9500, 2023.
- [5] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. Autogen: Enabling next-gen LLM applications via multi-agent conversation framework. *arXiv preprint arXiv:2308.08155*, 2023.
- [6] Andrew R Conn, Katya Scheinberg, and Luis N Vicente. *Introduction to derivative-free optimization*. SIAM, 2009.
- [7] Frank Hutter, Holger H Hoos, and Kevin Leyton-Brown. Sequential model-based optimization for general algorithm configuration. In *Learning and Intelligent Optimization*, pages 507–523. Springer, 2011.
- [8] Peter I Frazier. Bayesian optimization. *Recent Advances in Optimization and Modeling of Contemporary Problems*, pages 255–278, 2018.
- [9] Thomas Bäck and Hans-Paul Schwefel. An overview of evolutionary algorithms for parameter optimization. *Evolutionary computation*, 1(1):1–23, 1993.
- [10] Josep Ginebra and Murray K Clayton. Response surface bandits. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 57(4):771–784, 1995.
- [11] Yongchao Zhou, Andrei Ioan Muresanu, Ziwon Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. Large language models are human-level prompt engineers. In *ICLR*, 2023.
- [12] Reid Pryzant, Dan Iter, Jerry Li, Yin Lee, Chenguang Zhu, and Michael Zeng. Automatic prompt optimization with “gradient descent” and beam search. In *EMNLP*, pages 7957–7968, 2023.
- [13] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. Large language models as optimizers. In *ICLR*, 2024.
- [14] Tobias Schnabel and Jennifer Neville. Prompts as programs: A structure-aware approach to efficient compile-time prompt optimization. *arXiv preprint arXiv:2404.02319*, 2024.

- [15] Eric Zelikman, Eliana Lorch, Lester Mackey, and Adam Tauman Kalai. Self-taught optimizer (STOP): Recursively self-improving code generation. *arXiv preprint arXiv:2310.02304*, 2023.
- [16] Allen Nie, Ching-An Cheng, Andrey Kolobov, and Adith Swaminathan. The importance of directional feedback for LLM-based optimizers. *arXiv preprint arXiv:2405.16434*, 2024.
- [17] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986.
- [18] Atilim Gunes Baydin, Barak A. Pearlmutter, Alexey Andreyevich Radul, and Jeffrey Mark Siskind. Automatic differentiation in machine learning: a survey. *Journal of Machine Learning Research*, 18(153):1–43, 2018.
- [19] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. PyTorch: an imperative style, high-performance deep learning library. In *NeurIPS*, pages 8026–8037, 2019.
- [20] Karthik Nagaraj, Charles Killian, and Jennifer Neville. Structured comparative analysis of systems logs to diagnose performance problems. In *USENIX Symposium on Networked Systems Design and Implementation*, pages 353–366, 2012.
- [21] Mirac Suzgun, Nathan Scales, Nathanael Schärli, Sebastian Gehrmann, Yi Tay, Hyung Won Chung, Aakanksha Chowdhery, Quoc Le, Ed Chi, Denny Zhou, et al. Challenging BIG-Bench tasks and whether chain-of-thought can solve them. In *Findings of the ACL*, pages 13003–13051, 2023.
- [22] Omar Khattab, Arnav Singhvi, Paridhi Maheshwari, Zhiyuan Zhang, Keshav Santhanam, Sri Vardhamanan, Saiful Haq, Ashutosh Sharma, Thomas T Joshi, Hanna Moazam, et al. DSPy: Compiling declarative language model calls into self-improving pipelines. *arXiv preprint arXiv:2310.03714*, 2023.
- [23] Ching-An Cheng, Andrey Kolobov, Dipendra Misra, Allen Nie, and Adith Swaminathan. LLF-Bench: Benchmark for interactive learning from language feedback. *arXiv preprint arXiv:2312.06853*, 2023.
- [24] Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and James Zou. TextGrad: Automatic “differentiation” via text. *arXiv preprint arXiv:2406.07496*, 2024.
- [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [26] Toru Seo. UXsim: An open source macroscopic and mesoscopic traffic simulator in python—a technical overview. *arXiv preprint arXiv:2309.17114*, 2023.
- [27] Courtney Slavin, Wei Feng, Miguel Figliozzi, and Peter Koonce. Statistical study of the impact of adaptive traffic signal control on traffic and transit performance. *Transportation Research Record*, 2356(1):117–126, 2013.
- [28] James Kennedy and Russell Eberhart. Particle swarm optimization. In *ICNN*, volume 4, pages 1942–1948, 1995.
- [29] LangChain Team. Langchain tracing. <https://blog.langchain.dev/tracing/>, 2023.
- [30] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. AgentBench: Evaluating LLMs as agents. In *ICLR*, 2024.
- [31] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H Chi, Quoc V Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models. In *NeurIPS*, pages 24824–24837, 2022.

- [32] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020.
- [33] Timo Schick, Jane Dwivedi-Yu, Roberto Dessi, Roberta Raileanu, Maria Lomeli, Eric Hambro, Luke Zettlemoyer, Nicola Cancedda, and Thomas Scialom. Toolformer: Language models can teach themselves to use tools. In *NeurIPS*, 2023.
- [34] James Wexler. Artificial Intelligence in Games: A look at the smarts behind Lionhead Studio’s “Black and White” and where it can and will go in the future. <https://www.cs.rochester.edu/~brown/242/assts/termprojs/games.pdf>, 2002.
- [35] Weiran Yao, Shelby Heinecke, Juan Carlos Niebles, Zhiwei Liu, Yihao Feng, Le Xue, Rithesh R N, Zeyuan Chen, Jianguo Zhang, Devansh Arpit, Ran Xu, Phil L Mui, Huan Wang, Caiming Xiong, and Silvio Savarese. Retroformer: Retrospective large language agents with policy gradient optimization. In *ICLR*, 2024.
- [36] Alessandro Sordoni, Xingdi Yuan, Marc-Alexandre Côté, Matheus Pereira, Adam Trischler, Ziang Xiao, Arian Hosseini, Friederike Niedtner, and Nicolas Le Roux. Joint prompt optimization of stacked LLMs using variational inference. In *NeurIPS*, 2023.
- [37] Mingchen Zhuge, Wenyi Wang, Louis Kirsch, Francesco Faccio, Dmitrii Khizbullin, and Jurgen Schmidhuber. Language agents as optimizable graphs. *arXiv preprint arXiv:2402.16823*, 2024.
- [38] Ruotian Ma, Xiaolei Wang, Xin Zhou, Jian Li, Nan Du, Tao Gui, Qi Zhang, and Xuanjing Huang. Are large language models good prompt optimizers? *arXiv preprint arXiv:2402.02101*, 2024.
- [39] Xinyuan Wang, Chenxi Li, Zhen Wang, Fan Bai, Haotian Luo, Jiayou Zhang, Nebojsa Jojic, Eric Xing, and Zhiting Hu. PromptAgent: Strategic planning with language models enables expert-level prompt optimization. In *ICLR*, 2024.
- [40] Xinyu Tang, Xiaolei Wang, Wayne Xin Zhao, Siyuan Lu, Yaliang Li, and Ji-Rong Wen. Unleashing the potential of large language models as prompt optimizers: An analogical analysis with gradient-based model optimizers. *arXiv preprint arXiv:2402.17564*, 2024.
- [41] Xinyun Chen, Maxwell Lin, Nathanael Schärli, and Denny Zhou. Teaching large language models to self-debug. In *ICLR*, 2024.
- [42] Zeyuan Ma, Hongshu Guo, Jiacheng Chen, Guojun Peng, Zhiguang Cao, Yining Ma, and Yue-Jiao Gong. LLaMoCo: Instruction tuning of large language models for optimization code generation. *arXiv preprint arXiv:2403.01131*, 2024.
- [43] Tennison Liu, Nicolás Astorga, Nabeel Seedat, and Mihaela van der Schaar. Large language models to enhance bayesian optimization. In *ICLR*, 2024.
- [44] Michael R. Zhang, Nishkrit Desai, Juhan Bae, Jonathan Lorraine, and Jimmy Ba. Using large language models for hyperparameter optimization. *arXiv preprint arXiv:2312.04528*, 2023.
- [45] Gábor Bartók, Dean P Foster, Dávid Pál, Alexander Rakhlin, and Csaba Szepesvári. Partial monitoring—classification, regret bounds, and algorithms. *Mathematics of Operations Research*, pages 967–997, 2014.
- [46] Laurent Hascoet and Mauricio Araya-Polo. Enabling user-driven checkpointing strategies in reverse-mode automatic differentiation. *arXiv preprint cs/0606042*, 2006.
- [47] Amirreza Shaban, Ching-An Cheng, Nathan Hatch, and Byron Boots. Truncated back-propagation for bilevel optimization. In *AISTATS*, pages 1723–1732, 2019.
- [48] Ching-An Cheng, Jonathan Lee, Ken Goldberg, and Byron Boots. Online learning with continuous variations: Dynamic regret and reductions. In *AISTATS*, pages 2218–2228, 2020.
- [49] Sham Machandranath Kakade. *On the Sample Complexity of Reinforcement Learning*. PhD thesis, University College London, 2003.



## A Perspective: Deep Agent Workflows

We posit that the current practice of manually engineering computational workflows to build AI systems is analogous to programmers in the early 2000s hand-coding neural network weights to create engaging AI characters in video games [34]. Just like AutoDiff enabled the automatic and scalable optimization of deep neural networks with billions of parameters, we believe that Trace is the first step towards automatic and scalable optimization of “Deep Agent Workflows” to power even more capable AI systems. However, there are several limitations of the current implementation of Trace that need to be addressed to build Deep Agent Workflows.

When designing interactive AI systems that learn from their interactions, we need to define the *parameters* and *feedback* of the system. Parameters are the internal attributes that can be updated by the learning algorithm employed by the system. Feedback are the things observed and recorded by the system as a product of its interactions, and that provide signal for learning. Trace enables the development of new learning algorithms (e.g. through OptoPrime) that incorporate rich feedback to update heterogeneous parameters. In contrast, AutoDiff for deep neural networks uses numerical feedback (e.g. rewards or loss functions) to optimize numerical parameters (e.g. tensors). Black-box optimization techniques (e.g. Reinforcement Learning) can use numerical feedback to optimize heterogeneous parameters (e.g. codes, hyper-parameters as well as tensors), though they are inefficient. In the experiments, we saw that Trace was more efficient than black-box methods by using a generalization of back-propagation. Finally, Trace can use rich feedback (e.g. language) to extract more signal for learning.

**What can be traced?** Trace cannot convert all computational workflows into OPTO problems. Workflows with recursive bundle operators or those requiring distributed/parallel computing are not compatible with the current implementation. A future work would be to expand the Trace implementation to support these scenarios. In addition, the current implementation of Trace does not trace the execution within an operator defined by `bundle`, though in principle this is possible. There are also ambiguities in how an existing workflow can be traced and represented as a DAG. One example is when there is some sub-workflow following an `if` condition and another one following `else`. One can choose to wrap the entire code, including `if` and `else`, by `bundle` as a single operator. On the other hand, one can also just wrap the sub-workflows and not trace the `if` condition nor represent it as part of the DAG. (That is, suppose the `if` condition is true; from the DAG, one cannot see the alternate path under `else`). The latter choice has a flavor of applying “stop-gradient” on the boolean condition, whereas the first choice enables back-propagation through also the logical condition. We summarized some design considerations in Section 3.2. We foresee the choice of what to trace and how to trace in building Deep Agent Workflows will be an on-going research problem, similar to neural network architecture design.

**Where do we get rich feedback?** The OPTO framework captures an abstraction of rich feedback, called trace feedback (the execution trace and the output feedback), but it requires specification of operator descriptions and output feedback source to guide the optimization effectively. OPTO can be more efficiently solvable than black-box problems only when the trace feedback provides information beyond reward signals; otherwise, information-wise, OPTO is no easier than black-box problems. The Trace framework automates the generation of the execution trace in OPTO, when users of Trace decorate the workflow end-to-end. Currently, OptoPrime uses the docstring of operators in the execution trace to understand the operators; nonetheless, Trace logs also the source code of the decorated methods, which can also be used in the future to design optimizers that uses more details of the operators. In our experiments, we focus on output feedback that is automatically generated. In other contexts, e.g. users interacting with chatbots, we can natively gather natural language feedback or synthesize feedback [16] based on raw observations. We also anticipate that feedback in the form of images (e.g. users’ gestures) or videos (e.g. videogame player showing a desired correction to agent’s behavior) will be readily available and can be used as the output feedback in the OPTO framework. While Trace can handle this feedback, the current design of OptoPrime does not handle non-textual content (either in parameters, inputs or feedbacks). We anticipate future work along the lines of [16] on enhancing feedback design and developing guidelines for designing more informative and directive optimizer and feedback mechanisms.

**How to design adaptive optimizer for general OPTO problems?** The proposed OptoPrime optimizer shows the possibility of designing a single optimization algorithm to solve a range OPTO problems from diverse domains. However, we remark that OptoPrime is just the first step; it is akin to the vanilla gradient descent algorithm, which shows a proof of concept but is not scalable for large problems. The current design of OptoPrime has several scalability limitations due to its summary approach that updates parameters through one call to an LLM. While we show in Section 5.1, that OptoPrime can compete with ADAM in small problems, OptoPrime is not as computational efficient and cannot scale as well as ADAM; therefore, for large-scale numerical problems, Trace and OptoPrime with the current design does not replace classical AutoDiff. OptoPrime is also limited by the ability of LLMs. It has difficulty in handling parameters or nodes that cannot be compactly represented in text, which prevents it from optimizing neural network weights, or reasoning with large, stateful objects like a database. Similarly, it likely cannot handle large graphs (with thousands of nodes) at the moment, as such a large graph would result in a huge context which may be beyond what LLMs can understand and reason about reliably. This limitation is aggravated when dealing with noisy feedback or systems, as we need to present multiple graphs and feedback in the context at once. We also do not know how to rigorously define the concept of step size, which however we expect is important to handle noisy or local feedback. We need further research on graph simplification and representation, to reduce complexity and improve the efficiency of feedback propagation. Lastly, right now OptoPrime represents the constraints on parameters as part of the text description, but we have observed that LLMs do not always follow it. An effective workaround is to implement constraint checking in the workflow to throw exceptions (which are then handled as feedback). More specialized constraint handling techniques is an interesting research direction (e.g. “projecting” OPTO solver proposals for parameters onto their feasible sets), but they are not implemented in OptoPrime yet.

## B Related Work

**Framework for Computational Workflows** Frameworks such as LangChain [29], Semantic Kernels, AutoGen [5], DSPy [22] allow composing computational workflows and provide hand-engineered optimizers to tune an LLM’s context (i.e. prompt templates, few shot examples, or tool libraries) using scalar feedback with black-box search techniques. They support tracing of the workflow to aid in profiling, debugging and visualization. In contrast, Trace uses tracing for *automatic optimization*, and constructs a different representation of the computational graph which is suited for that purpose. Moreover, Trace is designed to be general-purpose and agnostic to the underlying frameworks of computational workflows users choose. In principle, one can apply Trace to decorate and tune a workflow based on a mix of AutoGen, LangChain, DSPy codes. In fact, our experiments in Appendix I use workflows declared using both AutoGen and DSPy.

**Optimization of Graphs of LLM Workflows** There are multiple efforts to optimize the computational graph of LLM workflows, which is a special case of the OPTO problem. These algorithms focus on optimizing prompts. SAMMO [14] is an example for prompts that uses additional graph structure to make the optimizer efficient. SAMMO represents the prompt parameter itself as a program so as to enable more efficient black-box search through the space of programs. DSPy [22] can optimize directly the prompts or the few-shot examples to include using scalar reward feedback. Retroformer [35] uses another small language model (LM) to provide suggestions/feedback (i.e. changing prompts) to improve the behavior of an actor LLM, where the small LM is tuned by offline RL. Deep Language Networks [36] view all of the prompts in an LLM workflow as tunable parameters and jointly optimizes them. They discovered that optimizing each parameter in isolation instead produces subpar results. [37] frames LLM systems as graph where nodes are operations and edges are messages/connections. (Note that this is different from the DAG used in Trace; here nodes are messages and edges are input-output of operators) and optimizes for the connection on edges (binary variables) by REINFORCE using scalar reward feedback and prompts by LLMs. They optimize each component separately without considering each other; for example, the prompts are optimized individually without considering the graph topology or how they are used down the road. We suspect this approach can be less stable. Their prompt optimization part also does not take output feedback, but simply use an LLM to self-check whether the prompt meets the need of generating desired functions the user specified. In contrast to these works, through the OPTO framing, Trace supports *joint* optimization of all parameters (prompts, hyperparameters, codes) with *rich* feedback, and is agnostic to graph structures (e.g., changing these parameters can dynamically change the graph

structure and connections between nodes). Users of Trace are free to specify which parameters they want to automatically optimize via online interactions.

**LLM-Optimizers for Prompts and Codes** There is a huge and fast growing literature on using LLMs as optimizers to improve prompts [12, 11, 38–40] or codes [41, 15, 42, 38]. Different from the works mentioned above, here the focus has been on an isolated problem (e.g., changing the behavior of a single LLM or improving the code generation in the question-answering format) rather than considering a non-trivial workflow or agent with multiple components like above. They do not consider optimizing prompts or codes as one component of a bigger workflow (e.g. implementing an autonomous agent), which is harder and requires the right credit assignment. In addition, these LLM-based optimizers, including OPRO [13], often propose only principles of how prompts should be designed and requires crafting problem specific prompts (as opposed to a single optimization prompt that can be applied to different problems). For adapting them to new problems, users need to design new prompts. Trace can also be applied to optimize trivial OPTO problems where the returned graph has just a single node of the parameter (which are the scenarios considered by these works). Nonetheless, the main focus of this paper is to study how optimization can be done efficiently as the graph becomes nontrivial and for diverse applications. Trace achieves this by using the abstract OPTO problem framing. Since OPTO encapsulates domain specific info in the graph, it enables designing fully instantiated optimizers that can be applied to multiple problems, rather than just principles which then requires hand crafting prompts for individual problems like previous works.

**LLM-Optimizers for Hyperparameters** Recent works like [43, 44] use LLMs to optimize numerical hyperparameters, as an alternate to Bayesian optimization. Here in the experiments we show that Trace + OptoPrime also can effectively learn hyperparameters, faster than Bayesian optimization. The main difference between Trace and the aforementioned work is the representation of the problem. In Trace, we provide the graph to the LLM-based optimization (through the pseudo-algorithm representation), and we consume rich language feedbacks on the output, both of which accelerates hyper-parameter optimization.

**OPTO Related Setups** OPTO is a generalization of partial monitoring games [45]. If there exists a latent loss function that the feedback  $f$  adheres to (e.g. as in [23]), those OPTO instances can be written as partial monitoring game. However OPTO admits a more general notion of feedback  $f$ , and we discuss solution concepts for them in Appendix F. On the other hand, OPTO can be also viewed as a special case of Learning from Language Feedback (LLF) setup defined in [23] with observations as the trace feedback. This is a framing of a *meta* LLF problem. In the LLFBench Meta-World experiments of this paper (Section 5), we show Trace can be used to learn policy for LLF problems grounded to an application too.

**AutoDiff and Back-propagation** Back-propagation has been shown to be a very effective tool in optimizing differential computational workflows. Our design of Trace is inspired by back-propagation and the ease of use of the AutoDiff framework PyTorch [19]. Nonetheless, we highlight that back-propagation (Backward Mode Differentiation) is not the only AutoDiff algorithm. For example, the gradient can be computed in a forward mode (Forward Mode Differentiation) as well, and there are also techniques of Checkpointing [46] and Truncated Back-Propagation approximation [47] for efficiency. What are the equivalent ideas of these methods for general computational workflows? We think this is an interesting future research direction.

## C Examples of OPTO

To ground the OPTO setup, we show how OPTO is related to some existing problems with examples.

**Example 4** (Neural network with back-propagation). The parameters are the weights.  $g$  is the neural computational graph and  $f$  is the loss. An example context  $\omega$  can be “Minimize loss”. The back-propagation algorithm, in view of the OPTO formulation, is embedded in the OPTO optimizer. For example, an OPTO optimizer here is a composition of back-propagation and gradient descent, where back-propagation takes  $\tau$  to compute the propagated gradient at the parameter.

**Example 5** (Code Debugging). The parameters are the codes.  $g$  denotes the stacked trace and  $f$  is the error message returned by a compiler.  $\omega$  can be “Make no error”.

**Example 6 (RL).** The parameters are the policy.  $g$  is the trajectory (of states, actions, rewards) resulting from running the policy in a Markov decision process; that is,  $g$  documents the graphical model of how an action generated by the policy, applied to the transition dynamics which then returns the observation and reward, etc.  $f$  can be the termination signal or a success flag.  $\omega$  can be “Maximize return” or “Maximize success”.

**Example 7 (Hyperparameter Tuning of ML Pipeline).** The parameters are e.g. learning rates and architectures.  $g$  describes the stages of the ML pipeline and the evaluation on the validation set, and  $f$  is the validation loss.  $\omega$  can be “Minimize validation error”.

**Example 8 (Prompt Optimization of an LLM Agent).** The parameters are the prompt of an LLM workflow.  $g$  is the computational graph of the agent and  $f$  is the feedback about the agent’s behavior (which can be scores or natural language).  $\omega$  can be “Maximize score” or “Follow the feedback”.

**Example 9 (Multi-Agent Collaboration).** The parameters are each agent’s prompts.  $g$  describes the entire conversation flow between agents, and  $f$  is the feedback about whether the task is successful after each agent performs their action.  $\omega$  can be “A group of agents coordinate to finish a task”.

As mentioned, the computational graph  $g$  returned by the Trace Oracle  $\mathcal{T}$  may have different graph structures. The length of the execution trace, e.g., in the debugging example above depends on how far the code executes. Similarly, the rollout length of in the RL problem can be randomly determined. The formulation of the Trace Oracle abstracts the details of a computational workflow, so problems from different domains can be framed in the unified framework. This abstraction allows us to design the computational tool Trace for various applications.

## D Trace Handles Error in Execution as Feedback

It is worth mentioning that execution error can be directly used as feedback to optimize parameters in Trace. When execution error happens within a method decorated by `bundle`, Trace would add a special exception node to the global computational graph and throw an `TraceExecutionError` to stop the computation. The computational graph ends at where the execution error happens. This exception node becomes the new output of the inputs to the decorated method (since the original method raises an error) and is the output of the truncated computational graph. Messages in `TraceExecutionError` can then be used as the feedback  $f$  in OPTO and propagated from the exception node to the parameters. By calling an OPTO optimizer, the parameters can be updated to avoid causing the same execution error. See the exception handling code in Fig. 2.

We find that this error handling mechanism has two convenient usages. First, this allows using Trace to automatically debug issues in the workflow due to incorrect parameter settings. Such errors can happen frequently especially when codes are parameters, as during optimization codes not satisfying syntax or downstream API requirements can happen. The second usage is to enforce constraints the workflow has to satisfy at different stages of computation. With Trace, if an intermediate computed result does not satisfy the constraint, we can simply throw an exception which states the desired constraint. This error signal would be caught by Trace and can then provide early feedback to efficiently improve the parameters, since the graph is truncated at the error.

## E Analysis of Trace

### E.1 Proof of Complexity

Algorithm 2 propagates the subgraph, represented by a priority queue (implemented as a min-heap). At a time, it needs to maintain the subgraphs coming from  $W$  children separately. This leads to the space complexity of  $O(WN)$ . This  $O(WN)$  space complexity leads to the extra  $WN \log N$  factor in the time complexity of MSP compared with back-propagation, which is the time needed for merging  $W$  subgraphs of size  $O(N)$ .

### E.2 Proof of Lower bounds

Consider an OPTO problem whose goal is to find a parameter matching a  $k$ -digit binary number. The computation checks each digit against a reference number in an arbitrary order. The feedback is either “ $N^{th}$  check failed” or “All checks succeeded”. Propagated feedback must communicate  $k$  bits

of information to interpret the feedback correctly; and the minimal subgraph conveys exactly that information. Updating the parameter using the minimal subgraph is trivial, whereas without it there are  $2^k$  possibilities to check.

## F When is OPTO Efficiently Solvable?

We show that OPTO covers a wide range of complicated optimization problems. This shows that if OPTO can be efficiently solved, then many complex workflows can be efficiently optimized. However, the generality of OPTO also raises some fundamental questions, such as if OPTO is well defined and when OPTO can be efficiently solved. These questions stem from its generality of the context  $\omega$  and the output feedback  $f$  in OPTO, since e.g. they can be anything describable texts. This flexibility makes the scope of OPTO go beyond standard mathematical optimization problems, where a setup has a fixed context  $\omega$  (e.g., “First-order optimization”) and a fixed type of output feedback  $f$  (a descent direction). Fully characterize the properties of OPTO, due to its generality, is beyond the scope of this paper and would require years of future research to come. Nonetheless, here we attempt to provide some preliminary answers and point out some research questions.

### F.1 What is a solution?

Classical mathematical optimization problems have a problem definition which itself is the solution concept. For example, in a minimization problem, it is clear we want to find the minimum of an objective function; even for problems as abstract and general as an equilibrium problem, the problem setup clearly states the solution concept of finding a point/set satisfying an equilibrium inequality [48]. One common pattern of these problems is that the solution concept is something that can be described as conditions on feedback that the parameter should satisfy.

By contrast, in a OPTO problem  $(\Theta, \omega, \mathcal{T})$ , by varying the context  $\omega$ , the desired parameter can change from one extreme to another. For example  $\omega$  may state “Follow the feedback” or “The feedback is adversarial.”. Therefore, we need define the solution concept of OPTO differently, rather than just using the feedback. We need to also consider the context  $\omega$  appropriately. Below we make an attempt to give an axiom of OPTO for its solution to be well defined.

**Axiom 1** (Verifiability). There is an verification oracle (a human, a machine learning model, or a polynomial-time algorithm) when given  $(\theta, \omega, f)$  can verify whether  $\theta$  is a solution or not.

Notice the verification oracle in Axiom 1 is not limited to just algorithms. This is intentional because we currently do not have algorithms that are intelligent enough to process the wide range of contexts and feedback that OPTO allows. Therefore, we include human judgement or the use of LLMs or other AI systems as part of the definition, while acknowledging the impreciseness of the statement due to OPTO’s soft computing nature. Lastly we note the verifiability is only defined with respect to the context  $\omega$  and the output feedback  $f$ , not the execution trace  $g$ . That is, the verification of a solution depends only on the output of computation.

### F.2 Does a solution exist?

Under Axiom 1, we can start to ask the basic question of whether a solution to an OPTO problem exists or not. There are clearly problems where no solution exists (that is, no parameter in  $\Theta$  can be verified by the verification oracle). For example if the feedback  $f$  is contradicting and yet the context  $\omega$  is “Follow the feedback.”, then there would be no solution that is satisfactory. On the other hand, if  $\omega$  is “Ignore the feedback”, all parameters can be solutions. In the following, we assume solutions of OPTO under consider exist. This assumption would rule out problems, e.g., where the feedback is adversary to the context, and makes solving OPTO is a well-defined search problem.

**Assumption 1.** For an OPTO problem  $(\theta, \omega, \mathcal{T})$ , we assume there is at least a parameter  $\theta \in \Theta$  such that it can be verified as a solution by the verification oracle.

### F.3 Can OPTO be efficiently solved?

So far our discussion establishes OPTO as a well-defined search problem, based on qualification on the context  $\omega$  and the output feedback  $f$ . However search problems can be NP-hard. In other words,



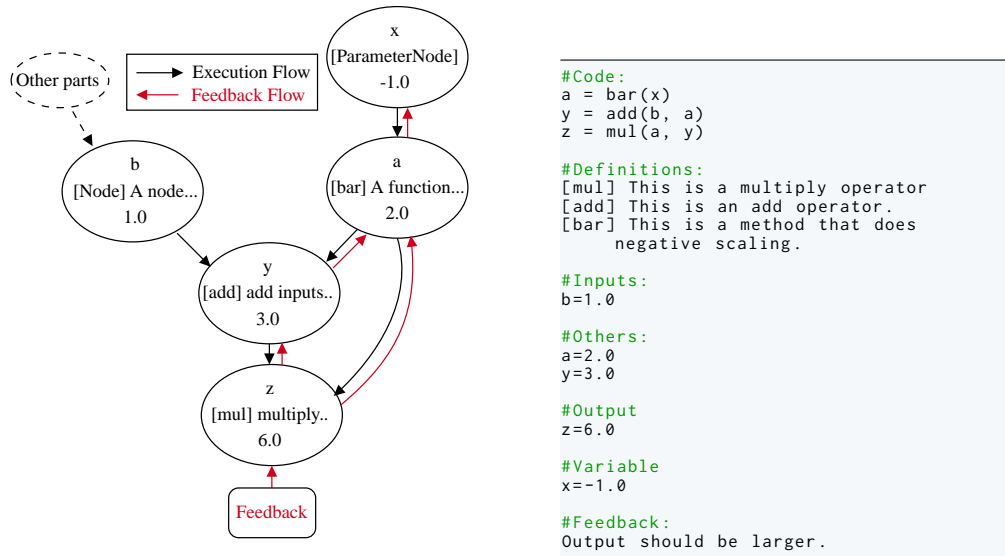
we know that, without the execution trace, there are search problem instances modeled by some  $\omega$  and oracle giving  $f$  that cannot be efficiently solved. Take RL for a tabular MDP as an example of OPTO problem. Without the execution trace (i.e., not seeing the Markovian structure and trajectories), the problem has an exponential complexity (due to the size of the policy space) and we know by using the execution trace here, tabular RL can be solved approximately in polynomial time [49]. Another example is training of neural networks. Without the execution trace, we have a complex black-box optimization with a loss value, without gradients, whereas an execution trace allows implementation of back-propagation to compute the gradients at the parameters.

More broadly speaking, if we consider a “human” as an optimizer for OPTO, we see that (expert) engineers/researchers, when equipped with additional computational tools, can efficiently solve a broad range of OPTO problems (such as by using the execution trace. From these observations, we conjecture using information in the execution trace is the key to unlock efficient OPTO. More precisely, we conjecture that OPTO is efficiently solvable when the context and the trace feedback need to provide information to construct a corrective search direction. For example, when the output feedback back is just a scalar loss, and yet the context + execution trace feedback does not provide enough information to compute a descending direction then OPTO reduces back to a black box problem. (See the problem instance in Appendix E.2). Nonetheless, identifying which subsets of OPTO are efficiently solvable is a big open research question.

## G Additional Details of Trace and OptoPrime

### G.1 Backward Step of Trace

The MSP extracts the minimal subgraph of the full computational graph of the workflow. Here we show a visualization using the example in Fig. 4.



(a) This is an illustrative example of the graph constructed by Trace and how feedback is backpropagated to the parameter  $x$ .

(b) We create a succinct summary of the computation graph using a language that mimics a program.

**Figure A.1: Optimization Representation.** For a program of  $x = \text{node}(-1.0)$ ;  $a = \text{bar}(x)$ ;  $y = a + 1$ ;  $z = a * y$  and the optimization objective of  $\max_x z$ , Trace automatically constructs a computation graph and represent the optimization problem as a debugging report. Note that the real program and the traced execution graph are different.

### G.2 Prompts used in OptoPrime

OptoPrime is an LLM-based optimizer. Its prompt is composed of the following parts.

1. System Prompt: Representation Prompt (Fig. A.2) + ReAct+CoT Output Prompt (Fig. A.3)

## 2. User Prompt (Fig. A.4 or Fig. A.5)

where + denotes concatenation. We list the prompt templates of different components below.

```
1 You're tasked to solve a coding/algorithm problem. You will see the instruction, the code,
  the documentation of each function used in the code, and the feedback about the execution
  result.
2
3 Specifically, a problem will be composed of the following parts:
4 - #Instruction: the instruction which describes the things you need to do or the question
  you should answer.
5 - #Code: the code defined in the problem.
6 - #Documentation: the documentation of each function used in #Code. The explanation might
  be incomplete and just contain high-level description. You can use the values in
  #Others to help infer how those functions work.
7 - #Variables: the input variables that you can change.
8 - #Constraints: the constraints or descriptions of the variables in #Variables.
9 - #Inputs: the values of other inputs to the code, which are not changeable.
10 - #Others: the intermediate values created through the code execution.
11 - #Outputs: the result of the code output.
12 - #Feedback: the feedback about the code's execution result.
13
14 In #Variables, #Inputs, #Outputs, and #Others, the format is:
15
16 <data_type> <variable_name> = <value>
17
18 If <type> is (code), it means <value> is the source code of a python code, which may
  include docstring and definitions.
19
```

**Figure A.2:** Representation Prompt that phrases the OPTO update as a pseudo-algorithm question.

```
1 Output_format: Your output should be in the following json format, satisfying the json
  syntax:
2
3 {{
4 "reasoning": <Your reasoning>,
5 "answer": <Your answer>,
6 "suggestion": {{
7 <variable_1>: <suggested_value_1>,
8 <variable_2>: <suggested_value_2>,
9 }}
10 }}
11
12 In "reasoning", explain the problem: 1. what the #Instruction means 2. what the
  #Feedback on #Output means to #Variables considering how #Variables are used in #Code
  and other values in #Documentation, #Inputs, #Others. 3. Reasoning about the suggested
  changes in #Variables (if needed) and the expected result.
13
14 If #Instruction asks for an answer, write it down in "answer".
15
16 If you need to suggest a change in the values of #Variables, write down the suggested
  values in "suggestion". Remember you can change only the values in #Variables, not
  others. When <type> of a variable is (code), you should write the new definition in the
  format of python code without syntax errors, and you should not change the function
  name or the function signature.
17
18 If no changes or answer are needed, just output TERMINATE.
19
```

**Figure A.3:** ReAct+CoT Output Prompt that instructs LLMs should respond in the format of (reasoning, answer, suggestion) and explains the output format.

```

1 Now you see problem instance:
2
3 =====
4 {actual_problem_instance}
5 =====
6
7 Your response:
8

```

**Figure A.4:** User Prompt for OptoPrime without Memory

```

1 Now you see problem instance:
2
3 =====
4 {actual_problem_instance}
5 =====
6
7 Below are some variables and their feedbacks you received in the past.
8
9 {
10   "variables": {
11     {variable1_name}: {variable1_value1}
12     {variable2_name}: {variable2_value1}
13     ...
14   },
15   "feedback": {feedback_1}
16 }
17
18 {
19   "variables": {
20     {variable1_name}: {variable1_value2}
21     {variable2_name}: {variable2_value2}
22     ...
23   },
24   "feedback": {feedback_2}
25 }
26 ...
27
28
29 Your response:
30

```

**Figure A.5:** User Prompt for OptoPrime with Memory

```

1 #Instruction
2 {instruction}
3
4 #Code
5 {code}
6
7 #Documentation
8 {documentation}
9
10 #Variables
11 {variables}
12
13 #Constraints
14 {constraints}
15
16 #Inputs
17 {inputs}
18
19 #Others
20 {others}
21
22 #Outputs
23 {outputs}
24
25 #Feedback:
26 {feedback}

```

**Figure A.6:** Problem Template used to fill the User Prompt. By default the Instruction (which is the context  $\omega$  of OPTO) is “You need to change the <value> of the variables in #Variables to improve the output in accordance to #Feedback.”

## H Comparison between Trace and TextGrad

Beyond the empirical results described in Section 5.5, there are several first principles reasons to prefer the Trace design to build future LLM-powered generative optimizers.

**Joint vs. Individual Parameter Optimization:** When there are multiple parameters to be optimized in a workflow, Trace takes a “joint optimization” view (constructing the minimal subgraph involving all of the parameters described in the initialization of OptoPrime). This flexible design also allows an alternative approach of “individual optimization” akin to co-ordinate descent if desired, i.e., fix all but one parameter and individually optimize each parameter. Such an optimization heuristic can be accomplished, e.g., by instantiating an OptoPrime instance per parameter and iterating through their updates; however, for many problems (like the toy example below), this strategy, which TextGrad employs, can be sub-optimal because a consistent optimizer must additionally maintain state about any suggested updates to all other parameters when reasoning about an individual parameter update.

```
1 @bundle()
2 def function1(x):
3     return x > 0
4 @bundle()
5 def function2(y):
6     return y % 2 == 0
7 def xor_test(x, y):
8     return function1(x).neq(function2(y))
9 input1 = node(3, trainable=True); input2 = node(4, trainable=True)
10 xor_test(input1, input2).backward(feedback="Find a set of inputs to make the return True.")
```

**Listing 1:** Individual Parameter Optimization is suboptimal in this example compared to Trace.

**Heterogenous Parameters:** Trace encapsulates all Python datatypes as node, including numbers, floats, strings, etc. In contrast, many generative optimization libraries like TextGrad restrict their “variables” (i.e. parameters and intermediate values) to only be strings; consequently running numerical optimization experiments, or optimizing code alongwith prompts and hyper-parameters can be challenging to set up.

**Modularity:** Trace provides a clean separation between the infrastructure (a platform akin to PyTorch) and a generative optimizer (OptoPrime akin to the Adam optimizer). Trace provides an object “TraceGraph” that describes the computational graph. This graph is independent of optimizers and LLMs. Such independent computation graph objects do not exist in TextGrad. We hope the Trace-provided interface can spur the development of many other generative optimizers, like how the OpenAI Gym API allowed the development of several Deep RL algorithms targeting a unified representation of RL problems.

**Flexible Representation:** The bundle functionality of Trace is more flexible in representing computations to an LLM than pre-defined templates used in other optimization frameworks. TextGrad relies on special functions to chain variables together, such as using “FormattedLLMCall”, while Trace supports any user-defined functions through decorators like bundle.

**Scalability:** The strategy taken by OptoPrime can scale better to larger graphs (both in terms of error accumulation across multiple LLM calls, and the costs of multiple LLM calls) compared to requiring one LLM call per computation node (and we see with bundle that not all computation nodes present equal difficulty for an optimizer). The difference in terms of computation costs between the two strategies can be large, especially when not all operators in the graph are as complex as querying LLM. However there may be even better optimizers that blend OptoPrime and TextGrad so as to divide-and-conquer huge computation graphs, we are excited about this avenue for future research.

# I Experiment Details

## I.1 Token Counts

We include token counts for the OPRO and OptoPrime prompts used across all our experiments, at the first iteration of optimization (note that OPRO’s token usages grows *linearly* with iterations). We can see that indeed OptoPrime consumes significantly more tokens than OPRO. However, we observe consistently that even allowing 7-10x more iterations of OPRO so as to equalize token costs, the OPRO performance plateaus to a worse level than OPTOPrime (e.g. Figure 1: OPRO at Iter 7 vs. OptoPrime at Iter 2; Figure 5b: OPRO at Iter 50 vs. OptoPrime at Iter 5; Figure 6b: OPRO at Iter 30 vs. OptoPrime at Iter 10, etc.). OPRO is suboptimal not due to a token limit but instead a lack of information, which is captured and represented using Trace.

Domain	OPRO	OptoPrime
Numerical Opt	175	918
BigBench-Hard	400	1883
Traffic Opt	198	1679
MetaWorld	470	7101
Battleship	437	1305

Table A.1: Token counts.

## I.2 Battleship

We implement a simple battleship game board in Python. The exact code is in the supplement. The game offers a string-based visualization of the board. It randomly places different types of ships on a 2-dimensional board with pre-specified width and height when it initializes. The agent does not see the ship location and has to select a coordinate on the board to hit next. One additional rule of this game is that the agent can go again if their previous coordinate selection (fire) is a hit, not counting as the finish of a turn. In Figure 1, we ran 10 trials, where in each trial, we ran 20 iterations of training. We measure the reward as % of ship squares hit (over all squares occupied by ships). The reward plateaued at 60% because the game has a chance element (heuristics and strategies can only go so far – strategy is only in effect if a hit happens. Otherwise, there is no information about where ships might be).

## I.3 Numerical Optimization

Any classical numerical optimization problem can be framed as an OPTO problem. Consider  $h(x)$  and a target  $y^*$ , in a context  $\omega$  finding the  $y^*$  by changing  $x$ ; we know the most useful corrective  $f$  feedback to change  $x$  is the gradient  $\nabla_h x$ . Similar to Trace, AutoDiff packages like PyTorch’s AutoGrad have implemented dynamic graph construction with special classes like `torch.Tensor`. We want to validate whether it is possible to rely on binary text feedback, a graph automatically constructed by Trace, and OptoPrime to update  $x$  in the context of minimizing  $|y - y^*|$ .

We constructed a synthetic task environment where we can create a complex computation graph with arbitrarily complex connections between numerical variables. The focus of this environment is on the complexity of the graph, not on the complexity of the numerical operators. Therefore, we only use one-dimensional input and basic arithmetic operators to create a numerical optimization problem solvable by a first-order optimizer. This environment constructs a computational graph by sampling a number of times. At each time, it will either use a previously computed variable or sample a new variable, and an operation will be sampled to combine them. The optimization task is, for a fixed number of steps, an optimizer needs to output  $x$  that minimizes  $y$ .

We evaluate the following baseline methods. Basic Agent: a basic LLM agent that simply stores past information of  $(x_{t-1}, y_{t-1})$  in context before choosing the next  $x_t$ . OPRO Agent: a basic LLM agent but we implement the state-of-the-art LLM optimizer OPRO [13], which updates the meta-prompt of the basic LLM agent. Torch + Adam: the problem we construct is end-to-end differentiable. Therefore, we simply pass in `torch.Tensor(x)` as input and use Adam optimizer to update. We tune the learning rate slightly and found  $1e-1$  to work well. We compare two kinds of Trace-based



optimizers: Trace, where we allow OptoPrime to read in the entire computation graph before updating  $x$ , or Trace Masked, where we hide the computation graph.

We run 30 trials over different computation graphs and start all methods with the same initial  $x, y^*$ . We compute the absolute error, which is  $|y - y^*|$ . On average, Trace is able to match the best-in-class first-order gradient optimizer Adam [25]. It is not entirely surprising that all the other baselines are performing worse due to a lack of access to the computation graph. To our surprise, OPRO, by only accessing the history of input and output, as well as changing the meta-prompt, is able to eventually discover the correct solution. This confirms why there were early signs of success using LLMs for black-box optimization in a simple plug-and-play style. However, OPRO is not an efficient optimizer because it lacks access to the Trace oracle. We show OPRO struggles even more when the computation graph gets more complex.

#### I.4 Traffic Control

We tested OptoPrime in a traffic control problem which is an instance of hyper-parameter tuning. We used UXSim [26] to simulate traffic at a four-way intersection, where the tunable parameters are the duration of the green lights for each direction of traffic flow. The feedback is a scalar loss calculated by monitoring the flow of a pseudo-random sequence of vehicles arriving at the intersection over a period of 30 minutes. The loss computes an estimate of the delay experienced by all vehicles due to the intersection, as well as variability in this estimate for every link in the network; lower values are better. The goal of an optimizer is to identify values for all of the green light duration so as to minimize the loss using the fewest number of traffic simulations. If the green light duration for a given traffic flow direction is set too low, then vehicles will queue up over time and experience delays, thereby lowering the score for the intersection. However, if the green light duration for a given direction is set too high, vehicles in other directions will queue up and experience delays, thereby lowering the score for the intersection. Hence an optimizer must find the right trade-off for temporally distributed and variable demands.

In Figure 5 we report the performance of a SOTA heuristic from the traffic control literature, SCATS [27] (adapted to this toy setting) as well as two black-box optimization techniques: Gaussian Process Minimization (GP) [8] and Particle Swarm Optimization (PSO) [28]. All methods are initialized to evaluate the same starting parameter. GP and PSO further evaluate 5 random parameters; moreover, if they query a previously evaluated point, that query is replaced by a randomly sampled parameter. GP constructs a surrogate model to mimic the black-box traffic simulation function which maps from parameters to observed score. Then it minimizes a utility function (e.g. the lower confidence bound) using the surrogate model to pick the next parameter to evaluate. PSO on the other hand maintains 5 particles in parameter space, each with a position and velocity. At each iteration of PSO, particles update their positions according to their previous positions and velocity, evaluate the function at the updated positions, and update the velocities of all particles using the observed values. Although GP and PSO are both black-box methods, GP can be thought to replace Trace oracle with instead a smooth differentiable surrogate function; whereas PSO is very different and maintains a candidate set of parameters (can be thought of as conceptually related to OptoPrime with memory).

GP appears to be bad because even when it converged, the exploration heuristic randomly samples parameters rather than pick the converged parameter. PSO appears bad because 10 iterations is insufficient for its convergence. Note that given enough number of iterations, black-box approaches will eventually perform well. Trace is quickly competitive with the SCATS heuristic, whereas OPRO is not. Moreover, we find that memory is crucial for OptoPrime to perform well at this task. Finally, Trace consumes additional overhead compared to black-box methods; beyond the space and time complexity for running the traffic simulation, Trace additionally materializes the computation graph per iteration. Thus it can also be more expensive per LLM call compared to OPRO.

#### I.5 BigBench-Hard

Perhaps more surprisingly, there are many components that a workflow needs to learn. Some of these components can be the prompt to generate output from an LLM, while other components can be code that needs to further process these outputs. In many workflows today, enabled by LangChain [29] and DSPy [22], only a small part of this workflow, predominantly, the input to an LLM API call, is optimized. These libraries optimize input to an LLM, and human engineers process that input

```

1 import trace
2
3 class Predict(trace.Module):
4     def __init__(self):
5         self.prompt_template = trace.node("""
6             Given 'question', produce the 'answer'.
7             question: {}
8             answer:
9             """,
10            trainable=True)
11
12     def forward(self, question):
13         user_prompt = self.create_prompt(question)
14         response = self.call_llm(user_prompt)
15         answer = self.extract_answer(question,
16            response)
17         return answer
18
19     @trace.bundle(trainable=True)
20     def create_prompt(self, question):
21         """formulate the prompt with the question"""
22         return self.prompt_template.format(question)
23
24     @trace.bundle(trainable=True)
25     def extract_answer(self, question, response):
26         """Extract the answer out of LLM response"""
27         answer = response.split("answer:")[1].strip()
28         return answer

```

(a) We write a workflow that prompts an LLM for a question and extracts the answer.

```

1 import trace
2
3 dataset = load_dataset(link, task)
4 predict = Predict()
5 optimizer =
6     trace.Optimizer(predict.parameters())
7
8 for ex in dataset:
9     # Forward pass
10    response = dp.forward(ex['question'])
11    if response == ex['answer']:
12        feedback = "The answer is correct"
13    else:
14        feedback = "The answer is incorrect"
15
16    # Backward pass and update
17    optimizer.zero_feedback()
18    optimizer.backward(response, feedback)
19    optimizer.step()

```

(b) The optimizer class takes in any parameter regardless of whether it is code or text. Although the actual optimization implementation can provide different treatments to many input types, the user interface stays consistent.

**Figure A.7: LLM-based Workflow Optimization Example.**

and integrate it into other systems. Indeed, both libraries can enable robust and swift large-scale engineering efforts to build LLM-based software. However, if our goal is to develop self-adapting agents that can modify their own behavior, we should not ignore one of LLM’s greatest strengths: code generation. Trace allows us to unify prompt optimization and code generation, which enables the creation of agents capable of fast learning.

In this example of an LLM-based workflow (Figure A.7), there are three parameters that are flagged as trainable for the optimizer: `prompt_template`, `create_prompt`, and `extract_answer`. Note that two of them require the LLM to generate Python code, and one of them requires the LLM to modify a text. Trace abstracted away the different data types and enabled direct update and optimization of them. Furthermore, a human engineer is often tasked with writing an error-free `extract_answer`. The output of an LLM can be highly stochastic and can often change over time; the code that is used to extract the response of an LLM has to be extremely robust and, therefore, arduous to create. Whenever a major distribution shift happens in the LLM output, this code needs to be rewritten by a human engineer, and it is hard for humans to anticipate all of LLM’s output patterns.

We set up the task of end-to-end workflow optimization. Unlike a typical LLM benchmark evaluation, where a lot of effort went into creating the perfect `evaluate(answer, target)` method so that all kinds of LLM outputs were post-processed, cleaned, and formatted to match the ground truth, we choose a simple evaluation function (that extracts a segment or does exact string matching) and place the burden on the workflow itself to figure out how to create the right answer to satisfy the evaluation metric. We choose Big-Bench Hard [21] as our task because it has 23 subtasks and contains both language and algorithmic tasks.

We split each task dataset into training, validation, and test. For Trace and Trace-CoT, we use the first 15 examples for training, 5 examples for validation (picking the best learned workflow), and then evaluate the performance on test examples. We use template-based positive and negative feedback during training. The positive feedback is “The answer is correct! No need to change anything.” The negative feedback is “The answer is wrong. We expect the output of your answer to be {ANSWER}. Please modify the prompt and relevant parts of the program to help LLM produce the right answer.” DSPy’s prompt optimization method does not explicitly require a validation set, therefore, we just used all 20 examples for training. For both, we only optimize for 1 epoch. We either start with the vanilla boilerplate prompt template used by DSPy or we use the slightly sophisticated template used by DSPy’s CoT module. Trace optimizes both DSPy’s original design and outperforms their own optimizer COPRO by 10% on algorithmic tasks.

Big-Bench Hard requires different answer outputs. Out of 23 tasks, 14 tasks require a multiple-choice answer with options provided in the question. 4 of them require yes/no. 1 task requires True/False, while 1 task requires valid/invalid. And the 3 remaining tasks require answers that contain words or

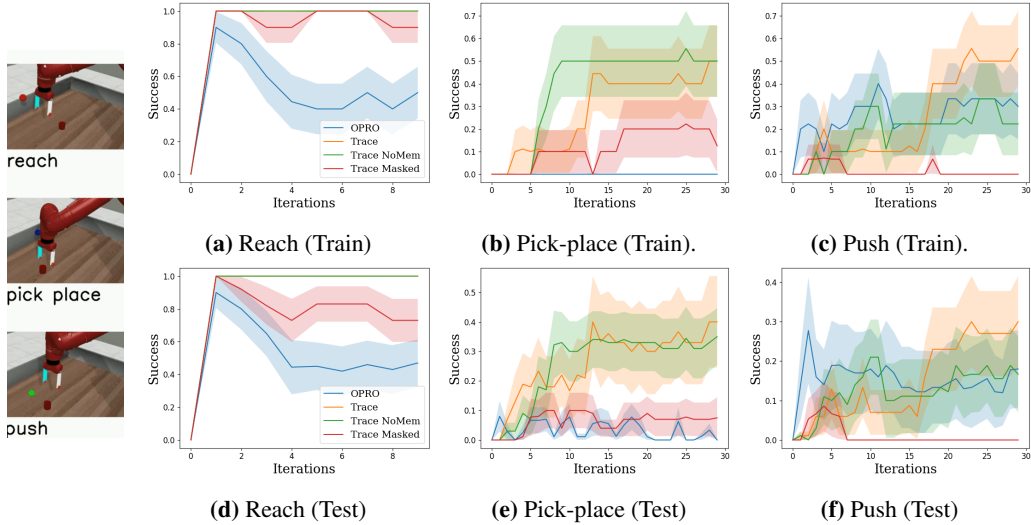
Task Name	DSPy	+CoT	DSPy-PO	+CoT	Trace	+CoT
tracking shuffled objects 7 objects	37.39	90.0	90.43	90.43	37.8	87.8
salient translation error detection	51.3	70.87	51.3	69.57	63.0	70.0
tracking shuffled objects 3 objects	39.13	94.35	97.39	93.91	38.7	96.5
geometric shapes	50.43	62.17	59.13	60.43	49.6	62.2
object counting	0.0	40.0	0.0	74.35	42.2	80.4
word sorting	0.0	0.0	0.0	0.43	84.3	74.3
logical deduction five objects	70.0	80.43	70.0	74.78	48.7	75.7
hyperbaton	74.78	86.52	74.78	88.26	78.3	91.7
sports understanding	0.0	0.0	0.0	0.0	79.6	45.7
logical deduction seven objects	68.7	64.78	68.7	64.78	45.7	69.6
multistep arithmetic two	0.0	93.04	0.0	93.04	94.8	88.7
ruin names	84.35	87.83	84.35	87.83	87.8	90.0
causal judgement	7.78	70.66	74.25	70.66	70.1	54.5
logical deduction three objects	85.22	97.39	85.22	97.83	91.7	97.0
formal fallacies	1.74	81.3	62.17	81.3	73.5	67.8
snarks	86.08	87.34	86.08	87.97	81.6	87.3
boolean expressions	0.0	98.26	64.35	98.26	88.7	96.5
reasoning about colored objects	53.04	91.3	89.13	91.3	91.3	95.7
dyck languages	0.0	8.7	7.83	8.7	26.5	9.6
navigate	0.0	95.65	0.0	97.39	59.6	92.6
disambiguation qa	67.83	66.09	73.91	66.09	75.7	59.1
temporal sequences	99.57	99.13	97.39	99.13	97.8	98.3
web of lies	0.0	0.0	0.0	0.0	49.6	90.4
tracking shuffled 5 objects	37.83	96.09	37.83	96.09	58.3	88.7
penguins in a table	69.84	92.86	97.62	92.86	81.7	91.3
movie recommendation	83.48	76.09	83.48	76.09	81.3	75.7
date understanding	69.13	85.65	69.13	85.65	70.4	85.7

**Table A.2:** Big Bench-Hard Per-Task Result. 0-shot performance. Some 0.0 here shown is because DSPy cannot find the clean/stripped output that matches what the automatic evaluation method expects. With additional human engineering, these numbers can improve.

numbers. Even though DSPy’s meta-prompt optimization is trained on each task individually, the output of LLM to the evaluation method is still not post-processed, resulting in low performances of these tasks. However, Trace can optimize code and LLM prompt jointly to successfully deliver the response expected by an automatic evaluation method.

## I.6 LLFBench Meta-World

We test the ability of Trace to optimize long-horizon workflows with complex dependencies. We experiment with using Trace to train controller (python code) for a simulated Sawyer robot manipulator. We use the Meta-World environment of LLF-Bench [23] as the simulator and consider three tasks reach, pick-place and push. LLF-Bench is a simulated benchmark with gym interface for testing an agent’s ability to learn from language feedback. In these LLF-Bench Meta-World tasks, the observation is a dictionary where each field denotes a feature of the state and has a vector value (e.g., the end-effector position, the goal position, the gripper status, etc.). The keys of the observation dictionary can differ for each task. The action space is 4-dimensional, which controls the relative position of the end-effector and the state of the gripper. In each time step, the LLF-Bench Meta-World simulator returns the observation dictionary and natural language feedback to guide the robot (we use the ‘a’ mode of LLF-Bench, with which the language feedback would contain information about the current performance, explanation of past successes and failures, and suggestions for the next step). An episode ends if the robot successfully solve the problem or because of time-out. For each task, LLF-Bench also provides a task instruction explaining that the task is about controlling a Sawyer robot arm and the meaning of the action space (see [23]). We use that as the context  $\omega$  of the OPTO problem. We consider an episodic setting. For each experiment (a random seed), we randomly sample an initial configuration. Then for each iteration of optimization, we reset the simulator to that sampled



**Figure A.8:** Learning the feedback control policy (code) for a simulated Sawyer manipulator in LLF-Bench Metaworld. In each iteration (x-axis), one episode of rollout is performed and then the policy is updated. Mean and standard error of success rate over 10 seeds are shown.

initial configuration and run the robot policy for  $10^{11}$  steps or until the episode termination due to success. We compute the sum of rewards and gives the output feedback  $f$  in texts in the format of “Success: <true/false> Return: <score>”. Note that the initial condition for all iterations within an experiment is the same so that the optimization problem is deterministic. To evaluate the learned policy’s performance, for each experiment, we additionally run the learned policy starting from 10 held-out initial conditions, different from the fixed training initial condition. For each training algorithm discussed, we run it with 30 iterations, where each iteration consists of one episode rollout and one update.

To optimize the controller with Trace, we declare the control code as the parameter using the bundle decorator with `trainable` set to `True`; the initial control code simply outputs a zero vector  $[0, 0, 0, 0]$ . We decorate also the reset and the step function of the gym environment, so that the entire rollout of an episode can be traced end-to-end. In our implementation, a prototypical rollout would create a graph with around 30 operations where the controller code parameter is used multiple times. This graph structure is similar to that of running a recurrent neural network. For Trace, we experiment with OptoPrime with and without a memory of size 10. In addition to Trace, we implement the state-of-the-art LLM optimizer OPRO [13] as a baseline. Compared with Trace, OPRO does not use the execution trace information but rely on just memory of parameter and feedback pairs<sup>12</sup> To make the comparison with Trace fair, we append the feedback observation from LLF-Bench given at each time step to the final feedback received by OPRO; on the other hand, Trace uses the simple final feedback of success and return and has to read the per-step feedback from the execution trace. To run OPRO in the OPTO setting, our implementation only proposes a single candidate in each iteration, which is then evaluated and provided with the output feedback. Since in [13] OPRO generates about 10 samples per iteration, so one iteration in [13] is roughly equivalent to 10 iterations here.

The experimental results are summarized in Fig. A.8, where we show the success rates at both the training initial condition as well as the held-out testing initial conditions over 10 seeds. OptoPrime is clearly the top-performing optimizer, especially the version with memory. OPRO is able to solve Reach at the start but its performance degraded over iterations (this instability was observed in [13]) and gets similar performance as OptoPrime (without memory) in Push. To validate that the performance of OptoPrime is indeed due to using the execution trace, we include an ablation where we mask out information in `#Inputs`, `#Others`, `#Code`, `#Definition` in the LLM context (see Fig. A.1b),

<sup>11</sup>We set the problem horizon to be 10 steps, as we find the expert policies implemented in LLF-Bench can solve these problems within 10 steps.

<sup>12</sup>The original version of OPRO uses parameter-score pairs. Since we’re interested in the more general setup of OPTO, we extend it to use parameter-feedback pairs.

which lead to significant degrade in performance and stability. This ablation shows that additionally using the execution trace provides more informed search direction compared with just using just the output feedback, which agrees with our hypothesis.

This experiment features the most complex graph structures, and using Trace for optimization here is similar to back-propagation over time. The experimental results here are quite impressive, showing that Trace is able to learn a complex control logic in a dozens of interactions, not only working on the training initial conditions but also on the held-out testing ones too. Nonetheless, we want to point out some limitations in the current experimental results. We find that the success rate of the learned policy varies largely across random seeds. Except for Reach (the simplest task), in a seed, often either it finds a policy close to 1.0 success rate or 0.0 success rate. Therefore, the plots can roughly be interpreted as how long it takes to find a working policy. In addition, in these experiments, we find that providing task-related context is necessary. We find the context needs to be informative enough for humans to understand the problem<sup>13</sup>; otherwise, the optimization can be solved efficiently with the time scale considered here. Nonetheless, this requirement is reasonable, as there is no free lunch.

---

<sup>13</sup>The original instructions in the v2 environments of LLF-Bench does not contain task specific background, but only the task name. We find this task name alone, except for Reach, is too vague even for humans to understand how to read the observation dictionary is related to the problem, so we added additionally one line to explain the task. Pick-place: “The goal of the task is to pick up a puck and put it to a goal position.”. Push: “The goal of the task is to push a puck to a goal position.”. Without this extra information, none of the LLM-based optimizers works in the experiments.

## J Examples of the Optimized Parameters in the Experiments

```
1 class Policy(trace.Module):
2     def init(self, width, height):
3         pass
4
5     def __call__(self, map):
6         return self.select_coordinate(map).data
7
8     def select_coordinate(self, map):
9         plan = self.reason(map)
10
11         act = self.act
12         output = act(map, plan)
13         return output
14
15     @trace.bundle(trainable=True)
16     def act(self, map, plan):
17         """
18         Given a map, select a target coordinate in a Battleship game. In map, 0 denotes
19         misses, X denotes successes, and . denotes unknown positions.
20         """
21         return
22
23     @trace.bundle(trainable=True)
24     def reason(self, map) -> str:
25         """
26         Given a map, analyze the board in a Battleship game. In map, 0 denotes misses, X
27         denotes successes, and . denotes unknown positions.
28         """
29         return [0, 0]
30
31 def train(epochs=10):
32     policy = Policy()
33     optimizer = FunctionOptimizerV2Memory(policy.parameters())
34     board = BattleshipBoard()
35
36     feedback = ""
37     rewards = []
38     board = board.reset()
39     obs = node(board.get_shots()) # init observation
40     i = 0
41     while i < epochs:
42         try:
43             output = policy.select_coordinate(obs)
44             # not traced
45             obs, reward, terminal, feedback = user_fb_for_placing_shot(board, output.data)
46         except TraceExecutionError as e: # this is a retry
47             output = e.exception_node
48             feedback = output.data
49             reward, terminal = 0, False
50
51         if terminal:
52             board = reset_board()
53             obs = node(board.get_shots()) # init observation
54
55         # Update
56         optimizer.zero_feedback()
57         optimizer.backward(output, feedback)
58         optimizer.step(verbose=True)
```

Figure A.9: Battleship Agent and optimization code.



```

1 ## Iteration 1 ; Success : False
2 def reason(self, map) -> str:
3     """
4     Given a map, analyze the board in a game. On map, 0 denotes misses, X denotes successes,
5     and . denotes unknown positions.
6     """
7     # Check for first unknown position to try as a potential move
8     for row_index, row in enumerate(map):
9         for col_index, val in enumerate(row):
10            if val == '.':
11                return [row_index, col_index]
12 # Default to [0, 0] if no unknown positions found
13 return [0, 0]
14 ## Iteration 2 ; Success : True
15 def reason(self, map) -> str:
16     """
17     Given a map, analyze the board in a game. On map, 0 denotes misses, X denotes successes,
18     and . denotes unknown positions.
19     """
20 # Improved strategic approach by checking . surrounded by more . which might indicate
21 undetected area
22 max_surrounding = -1
23 best_spot = [0, 0]
24 for row_index, row in enumerate(map):
25     for col_index, val in enumerate(row):
26         if val == '.':
27             surrounding = 0
28             for di in [-1, 0, 1]:
29                 for dj in [-1, 0, 1]:
30                     ni, nj = row_index + di, col_index + dj
31                     if 0 <= ni < len(map) and 0 <= nj < len(map[0]) and map[ni][nj] ==
32                         '.':
33                         surrounding += 1
34                     if surrounding > max_surrounding:
35                         max_surrounding = surrounding
36                         best_spot = [row_index, col_index]
37 return best_spot
38 ## Iteration 8 ; Success : True
39 def reason(self, map) -> str:
40     """
41     Given a map, analyze the board in a game. On map, 0 denotes misses, X denotes successes,
42     and . denotes unknown positions.
43     """
44 max_surrounding = -1
45 best_spot = [0, 0]
46 for row_index, row in enumerate(map):
47     for col_index, val in enumerate(row):
48         if val == '.':
49             surrounding = 0
50             for di in [-1, 0, 1]:
51                 for dj in [-1, 0, 1]:
52                     ni, nj = row_index + di, col_index + dj
53                     if 0 <= ni < len(map) and 0 <= nj < len(map[0]):
54                         if map[ni][nj] == '.':
55                             surrounding += 1
56                         elif map[ni][nj] == 'X':
57                             surrounding += 1.5 # increasing emphasis on proximity to
58 successful hits
59                     if surrounding > max_surrounding:
60                         max_surrounding = surrounding
61                         best_spot = [row_index, col_index]
62 return best_spot

```

**Figure A.10:** Learned Battleship Code (reason). Functions with the same name are learned during different iterations or trials.

```

1 ## Iteration 0 (Initialization); Success : False
2
3
4 ## Iteration 1 ; Success : False
5 def act(self, map, plan):
6     """
7     Given a map, select a target coordinate in a game. On map, 0 denotes misses, X denotes
8     successes, and . denotes unknown positions.
9     """
10    # Example implementation, more sophisticated logic may be required based on intended
11    # gameplay mechanics
12    return [0, 1]
13
14 ## Iteration 2 ; Success : True
15 def act(self, map, plan):
16    """
17    Given a map, select a target coordinate in a game. On map, 0 denotes misses, X denotes
18    successes, and . denotes unknown positions.
19    """
20    # Use information from plan, which is output from __code2, checking strategic position,
21    # not just first available
22    best_score = -1
23    best_move = [0, 0] # initially no move is considered
24    for row_index, row in enumerate(map):
25        for col_index, val in enumerate(row):
26            if val == '.': # consider all unknown spots
27                hopes_surrounding = 0 # will hold number of '.'s around the spot
28                for di in [-1, 0, 1]:
29                    for dj in [-1, 0, 1]:
30                        ni, nj = row_index + di, col_index + dj
31                        if 0 <= ni < len(map) and 0 <= nj < len(map[0]):
32                            if map[ni][nj] == '.':
33                                hopes_surrounding += 1
34                    score = hopes_surrounding # more '.' around, higher the chance of hit
35                    if score > best_score:
36                        best_score = score
37                        best_move = [row_index, col_index]
38    if map[plan[0]][plan[1]] == '.' and best_score == -1:
39        return plan # if no better move found, and original place is still unknown
40    return best_move if best_score > -1 else plan # return the best move found, or stick
41    with the plan if still valid
42
43 ## Iteration 8 ; Success : True
44 def act(self, map, plan):
45    """
46    Given a map, select a target coordinate in a game. On map, 0 denotes misses, X denotes
47    successes, and . denotes unknown positions.
48    """
49    best_score = -1
50    best_move = [0, 0]
51    for row_index, row in enumerate(map):
52        for col_index, val in enumerate(row):
53            if val == '.' and [row_index, col_index] != plan:
54                hopes_surrounding = 0
55                for di in [-1, 0, 1]:
56                    for dj in [-1, 0, 1]:
57                        ni, nj = row_index + di, col_index + dj
58                        if 0 <= ni < len(map) and 0 <= nj < len(map[0]):
59                            if map[ni][nj] == '.':
60                                hopes_surrounding += 1 # favoring positions with more
61                                unknowns surrounding
62                            elif map[ni][nj] == 'X':
63                                hopes_surrounding += 2 # increased incentive for moves near
64                                successful spots
65                    score = hopes_surrounding
66                    if score > best_score:
67                        best_score = score
68                        best_move = [row_index, col_index]
69    if best_score > -1:
70        return best_move
71    return plan

```

**Figure A.11:** Learned Battleship Code (act). Functions with the same name are learned during different iterations or trials.

```

1 @trace_class
2 class Predict(LLMCallable):
3     def __init__(self):
4         super().__init__()
5
6         self.demos = []
7         self.prompt_template = dedent("""
8         Given the fields 'question', produce the fields 'answer'.
9
10        ---
11
12        Follow the following format.
13
14        Question:
15        Answer:
16
17        ---
18        Question: {}
19        Answer:
20        """)
21
22        self.prompt_template = trace.node(self.prompt_template, trainable=True,
23                                         description="[ParameterNode] This is the Prompt
24                                         Template to the LLM...")
25
26 @trace.bundle(trainable=True)
27 def extract_answer(self, prompt_template, question, response):
28     """
29     Need to read in the response, which can contain additional thought, deliberation and
30     an answer.
31     Use code to process the response and find where the answer is.
32     Can use self.call_llm("Return the answer from this text: " + response) again to
33     refine the answer if necessary.
34
35     Args:
36     prompt_template: The prompt that was used to query LLM to get the response
37     question: Question has a text describing the question but also "Options"
38     response: LLM returned a string response
39     Process it and return the answer in the exact format that the
40     evaluator wants to see.
41     Be mindful of the type of answer you need to produce.
42     It can be (A)/(B), a number like 8, or a string, or Yes/No.
43     """
44     answer = response.split("Answer:")[1].strip()
45     return answer
46
47 @trace.bundle(trainable=True)
48 def create_prompt(self, prompt_template, question):
49     """
50     The function takes in a question and then add to the prompt for LLM to answer.
51     Args:
52     prompt_template: some guidance/hints/suggestions for LLM
53     question: the question for the LLM to answer
54     """
55     return prompt_template.format(question)
56
57 def forward(self, question):
58     """
59     question: text
60
61     We read in a question and produces a response
62     """
63     user_prompt = self.create_prompt(self.prompt_template, question)
64     response = self.call_llm(user_prompt)
65     answer = self.extract_answer(self.prompt_template, question, response)
66     return answer

```

**Figure A.12:** Starting Code for BigBench. We write it in a similar style to DSPy's Predict module.

```

1 @trace_class
2 class PredictCoT(LLMCallable):
3     def __init__(self):
4         super().__init__()
5
6         self.demos = []
7         self.prompt_template = dedent("""
8         Given the fields 'question', produce the fields 'answer'.
9
10        ---
11
12        Follow the following format.
13
14        Question: question
15        Reasoning: Let's think step by step in order to produce the answer. We ...
16        Answer: answer
17
18        ---
19        Question: {}
20        """)
21
22         self.prompt_template = trace.node(self.prompt_template, trainable=True,
23                                         description="[ParameterNode] This is the Prompt
24                                         Template to the LLM...")
25
26     @trace.bundle(trainable=True)
27     def extract_answer(self, prompt_template, question, response):
28         """
29         Need to read in the response, which can contain additional thought, deliberation and
30         an answer.
31         Use code to process the response and find where the answer is.
32         Can use self.call_llm("Return the answer from this text: " + response) again to
33         refine the answer if necessary.
34
35         Args:
36             response: LLM returned a string response
37             Process it and return the answer in the exact format that the
38             evaluator wants to see.
39             Be mindful of the type of answer you need to produce.
40             It can be (A)/(B), a number like 8, or a string, or Yes/No.
41             question: Question has a text describing the question but also "Options"
42         """
43         answer = response.split("Answer:")[1].strip()
44         return answer
45
46     @trace.bundle(trainable=True)
47     def create_prompt(self, prompt_template, question):
48         """
49         The function takes in a question and then add to the prompt for LLM to answer.
50         The prompt should instruct the LLM to reason, think.
51         Args:
52             prompt_template: some guidance/hints/suggestions for LLM
53             question: the question for the LLM to answer
54         """
55         return prompt_template.format(question)
56
57     def forward(self, question):
58         """
59         question: text
60
61         We read in a question and produces a resposne
62         """
63         user_prompt = self.create_prompt(self.prompt_template, question)
64         response = self.call_llm(user_prompt)
65         answer = self.extract_answer(self.prompt_template, question, response)
66         return answer

```

**Figure A.13:** Starting Code for BigBench. We write it in a similar style to DSPy’s Predict CoT (0-shot Chain-of-Thought) module.

```

1 ## Iteration 0 ( initialization )
2 def create_prompt(self, prompt_template, question):
3     """
4     The function takes in a question and then add to the prompt for LLM to answer.
5     Args:
6     prompt_template: some guidance/hints/suggestions for LLM
7     question: the question for the LLM to answer
8     """
9     return prompt_template.format(question)
10
11 ## Iteration > 0
12 def create_prompt(self, prompt_template, question):
13     """
14     The function takes in a question and then add to the prompt for LLM to answer.
15     The prompt should now further instruct the LLM to carefully track the ball swaps
16     occurring step-by-step.
17     Args:
18     prompt_template: some guidance/hints/suggestions for LLM
19     question: the question for the LLM to answer
20     """
21     prompt_template = 'Process this carefully: Step-by-step.' + prompt_template
22     return prompt_template.format(question)

```

**Figure A.14:** Learned Predict module for BigBench. Functions with the same name are learned during different iterations or trials.

```

1 ## Iteration 0 ( initialization )
2 def extract_answer(self, prompt_template, question, response):
3     """
4     Need to read in the response, which can contain additional thought, deliberation and an
5     answer.
6     Use code to process the response and find where the answer is.
7     Can use self.call_llm("Return the answer from this text: " + response) again to refine
8     the answer if necessary.
9
10    Args:
11        prompt_template: The prompt that was used to query LLM to get the response
12        question: Question has a text describing the question but also "Options"
13        response: LLM returned a string response
14        Process it and return the answer in the exact format that the evaluator
15        wants to see.
16        Be mindful of the type of answer you need to produce.
17        It can be (A)/(B), a number like 8, or a string, or Yes/No.
18    """
19    answer = response.split("Answer:")[1].strip()
20    return answer
21
22 ## Iteration > 0
23 def extract_answer(self, response):
24     """
25     Need to read in the response, which can contain additional thought, deliberation and an
26     answer.
27     Use code to process the response and find where the answer is.
28     Can use self.call_llm("Return the answer from this text: " + response) again to refine
29     the answer if necessary.
30
31    Args:
32        response: LLM returned a string response
33        Process it and return the answer in the exact format that the evaluator
34        wants to see.
35        Be mindful of the type of answer you need to produce.
36        It can be (A)/(B), a number like 8, or a string, or Yes/No.
37        question: Question has a text describing the question but also "Options"
38    """
39    answer = ''
40    segments = response.split('\n')
41    for segment in segments:
42        if 'Answer:' in segment:
43            answer = segment.split('Answer:')[1].strip()
44    refined_answer = self.call_llm('Return the refined answer from this text: ' + answer)
45    return refined_answer
46
47 def extract_answer(self, prompt_template, question, response):
48     """
49     Processes the LLM response and extracts the final answer in the required format.
50     """
51     # Assuming that the relevant part of the response is after 'Answer:' and before any
52     # further commentary
53     extracted_part = response.split('Answer: ')[1].split(' ')[0].strip()
54     # Find the section of the answer and return it directly
55     result = re.search('\([A-E]\)', extracted_part)
56     if result:
57         return result.group()
58     else:
59         return 'No valid answer found'
60
61 def extract_answer(self, prompt_template, question, response):
62     """
63     Processes the LLM response, extracting and formatting the final answer.
64     Uses code to meticulously parse the response to locate the answer section.
65     Optionally refines the answer by querying the LLM again if necessary.
66
67    Args:
68        response: string from LLM, expected format contains 'Answer:' followed by the answer.
69        question: Description of the question being addressed, may include 'Options'
70    """
71    answer = response.split('Answer:')[1].split('.')[0].strip().lower().split(',')[0].strip()
72    return answer

```

**Figure A.15:** Learned Predict module for BigBench. Functions with the same name are learned during different iterations or trials.



```
1 ## Iteration 0 ( initialization )
2 prompt_template = """
3 Given the fields 'question', produce the fields 'answer'.
4
5 ---
6
7 Follow the following format.
8
9 Question:
10 Answer:
11
12 ---
13 Question: {}
14 Answer:
15 """
16
17 ## Iteration > 0
18 prompt_template = """
19 Given the fields 'question', produce the fields 'answer'.
20
21 ---
22
23 Follow the following format.
24
25 Question:
26 Answer: Please respond precisely in the format required.
27
28 ---
29 Question: {}
30 """
31
```

**Figure A.16:** Learned BigBench Prompt Template.

```

1
2 ## Iteration 0 (initialization); Success: False
3 def controller(obs):
4     """
5     A feedback controller that computes the action based on the observation.
6
7     Args:
8         obs: (dict) The observation from the environment. Each key is a string (indicating a
9             type of observation) and the value is a list of floats.
10            Output:
11            action: (list or nd.array) A 4-dimensional vector.
12            """
13            return [0, 0, 0, 0]
14 ## Iteration 1; Success: False
15 def controller(obs):
16     """
17     A feedback controller that computes the action based on the observation.
18
19     Args:
20         obs: (dict) The observation from the environment. Each key is a string (indicating a
21             type of observation) and the value is a list of floats.
22            Output:
23            action: (list or nd.array) A 4-dimensional vector.
24            """
25            hand_pos = obs['hand_pos']
26            goal_pos = obs['goal_pos']
27            action = [goal_pos[0] - hand_pos[0], goal_pos[1] - hand_pos[1], goal_pos[2] -
28                    hand_pos[2], 0]
29            # Adjusting the last element of the action based on whether the hand is close enough to
30            # trigger grasping
31            if abs(goal_pos[0] - hand_pos[0]) < 0.02 and abs(goal_pos[1] - hand_pos[1]) < 0.02 and
32                abs(goal_pos[2] - hand_pos[2]) < 0.02:
33                action[3] = 1 # Assuming action[3] = 1 triggers closing the gripper
34            return action
35 ## Iteration 5; Success: False
36 def controller(obs):
37     """
38     A feedback controller that computes the action based on the observation.
39
40     Args:
41         obs: (dict) The observation from the environment.
42            Output:
43            action: (list or nd.array) A 4-dimensional vector.
44            """
45            hand_pos = obs['hand_pos']
46            puck_pos = obs['puck_pos']
47            goal_pos = obs['goal_pos']
48            gripper_open = 1 if obs['gripper_distance_apart'][0] > 0.5 else 0
49
50            # Adjusting strategy to close the gripper when in proximity of the puck
51            if gripper_open and ((abs(hand_pos[0] - puck_pos[0]) < 0.05) and (abs(hand_pos[1] -
52                puck_pos[1]) < 0.05) and (abs(hand_pos[2] - puck_pos[2]) < 0.05)):
53                action = [0, 0, 0, 1] # Close the gripper
54            elif not gripper_open:
55                direction_to_goal = [0.09 - hand_pos[0], 0.95 - hand_pos[1], 0.12 - hand_pos[2]]
56                action = [direction_to_goal[0], direction_to_goal[1], direction_to_goal[2], 0] #
57                Move towards the suggested pose once puck is grasped
58            else:
59                # Move towards the puck first if not carrying it
60                direction_to_puck = [puck_pos[0] - hand_pos[0], puck_pos[1] - hand_pos[1],
61                puck_pos[2] - hand_pos[2]]
62                action = [direction_to_puck[0], direction_to_puck[1], direction_to_puck[2], 0]
63            return action

```

**Figure A.17:** Learned Code for LLFBench Meta-World Pick-Place (Part 1). Functions with the same name are learned during different iterations or trials.

```

1
2
3 ## Iteration 10; Success: False
4 def controller(obs):
5     """
6     A feedback controller that computes the action based on the observation.
7
8     Args:
9         obs: (dict) The observation from the environment.
10    Output:
11        action: (list or nd.array) A 4-dimensional vector.
12    """
13    hand_pos = obs['hand_pos']
14    puck_pos = obs['puck_pos']
15    goal_pos = obs['goal_pos']
16    gripper_open = 1 if obs['gripper_distance_apart'][0] > 0.5 else 0
17
18    normalize = lambda x: [i / max(abs(max(x, key=abs)), 1) for i in x]
19
20    # Close the gripper when close to the puck and the gripper is open
21    if gripper_open and ((abs(hand_pos[0] - puck_pos[0]) < 0.05) and (abs(hand_pos[1] -
22    puck_pos[1]) < 0.05) and (abs(hand_pos[2] - puck_pos[2]) < 0.05)):
23        return [0, 0, 0, 1]
24    # When the puck is grasped, move towards the goal position with a normalized direction
25    elif not gripper_open and ((abs(hand_pos[0] - puck_pos[0]) < 0.1) and (abs(hand_pos[1] -
26    puck_pos[1]) < 0.1) and (abs(hand_pos[2] - puck_pos[2]) < 0.1)):
27        direction_to_goal = [goal_pos[0] - hand_pos[0], goal_pos[1] - hand_pos[1],
28    goal_pos[2] - hand_pos[2]]
29        return normalize(direction_to_goal) + [0]
30    # When the gripper is open and not close enough to the puck, move towards the puck
31    else:
32        direction_to_puck = [puck_pos[0] - hand_pos[0], puck_pos[1] - hand_pos[1],
33    puck_pos[2] - hand_pos[2]]
34        return normalize(direction_to_puck) + [0]
35
36 ## Iteration 13; Success: True
37 def controller(obs):
38     """
39     A feedback controller that computes the action based on the observation.
40
41     Args:
42         obs: (dict) The observation from the environment.
43    Output:
44        action: (list or nd.array) A 4-dimensional vector.
45    """
46    hand_pos = obs['hand_pos']
47    puck_pos = obs['puck_pos']
48    goal_pos = obs['goal_pos']
49    gripper_open = 1 if obs['gripper_distance_apart'][0] > 0.5 else 0
50
51    normalize = lambda x: [i / max(abs(max(x, key=abs)), 1) for i in x]
52    # Close the gripper when close to the puck and the gripper is open
53    if gripper_open and ((abs(hand_pos[0] - puck_pos[0]) < 0.05) and (abs(hand_pos[1] -
54    puck_pos[1]) < 0.05) and (abs(hand_pos[2] - puck_pos[2]) < 0.05)):
55        action = [0, 0, 0, 1] # Close the gripper
56    elif not gripper_open and ((abs(hand_pos[0] - puck_pos[0]) < 0.1) and (abs(hand_pos[1] -
57    puck_pos[1]) < 0.1) and (abs(hand_pos[2] - puck_pos[2]) < 0.1)):
58        # When the puck is grasped, adjust direction towards the goal with improved
59        # precision and ensure successful pickup
60        direction_to_goal = normalize([goal_pos[0] - hand_pos[0], goal_pos[1] - hand_pos[1],
61    goal_pos[2] - hand_pos[2]])
62        action = direction_to_goal + [1] # Keep the gripper closed
63    else:
64        # Move towards the puck first if not carrying it
65        direction_to_puck = normalize([puck_pos[0] - hand_pos[0], puck_pos[1] - hand_pos[1],
66    puck_pos[2] - hand_pos[2]])
67        action = direction_to_puck + [0] # Open the gripper to prepare for grasping
68
69    return action

```

**Figure A.18:** Learned Code for LLFBench Meta-World Pick-Place (Part 2). Functions with the same name are learned during different iterations or trials.

```

1 # Iteration 1
2 {
3   "reasoning": "The task involves controlling a robotic arm to perform a pick-place
4     operation by interpreting observations and generating appropriate control actions (4D
5     vectors) to move the arm towards the goal. The objective is to adjust the controller's
6     code (__code0) based on feedback received at each step of the process. The feedback within
7     #Others suggests whether the movements made by the arm are leading it closer to or
8     farther from the goal, along with advice on what coordinates to aim for next (e.g.,
9     [-0.07 0.68 0.12 0. ]). The controller's code currently returns a static action [0,
10    0, 0, 0], which does not adjust based on observations and hence does not efficiently
11    move the robotic arm towards the goal. Given the repetitive feedback suggesting to go to
12    [-0.07 0.68 0.12 0. ], it appears that the controller needs to dynamically compute
13    actions based on the difference between the current positions (hand_pos, puck_pos) and
14    the goal_pos, along with possibly manipulating the gripper state for pick and place
15    operation.",
16   "answer": "",
17   "suggestion": {...}
18 }
19
20 # Iteration 5
21 {
22   "reasoning": "The instruction involves creating a controller for a Sawyer robot arm to
23     solve a pick-and-place task by providing sequential control actions based on
24     observations of the environment and state of the robot. The feedback indicates that
25     despite the actions taken, the success criteria have not been met, with the latest
26     attempt achieving a return of 0.9574427394250022. The code provided for __code0 defines
27     a controller function using the observed state to compute control actions aimed at
28     grasping a puck and moving it to a goal position while adjusting the gripper based on
29     its distance from the puck. Feedback throughout the execution indicates that the
30     controller adapts its strategy based on proximity to the puck and whether it has been
31     picked up. Feedback also suggests that the gripper's opening and closing actions are
32     central to the task's success, with corrections being made to close the gripper when
33     necessary. The controller's logic for determining movements toward the puck or goal and
34     adjusting the gripper state appears reasonable, but the feedback implies that
35     adjustments might be required to ensure task completion. This could include fine-tuning
36     the conditions under which the gripper is closed or adjusting the thresholds for puck
37     pickup and delivery to the goal.",
38   "answer": "",
39   "suggestion": {}
40 }
41
42 # Iteration 10
43 {
44   "reasoning": "The task involves programming a robot arm to pick up a puck and place it
45     at a goal position by providing sequential 4-dimensional action vectors. The action
46     vectors represent target positions and gripper states (open or close) based on the
47     robot's and puck's current positions. The feedback indicates that despite optimizing the
48     controller logic to close the gripper when near the puck, move towards the puck when
49     it's not grasped, and head towards the goal once the puck is grasped, the actions have
50     not resulted in task completion. Specifically, repeated suggestions to close the
51     gripper (indicating it remained open when it should have been closed) and guidance on
52     where to move next suggest that the action selection logic might not be effectively
53     responding to the environment's state or might not be properly engaging the gripper at
54     the right times. The latest code includes a condition to normalize the direction of
55     movement and checks for gripper position, but there seems to be a misunderstanding or
56     oversight in adjusting the robot's behavior based on the scenario to successfully
57     complete the pick-place task.",
58   "answer": "",
59   "suggestion": {...}
60 }
61 }

```

**Figure A.19:** Example Responses of LLM for LLFBench Meta-World Pick-Place (Part 1); Content in suggestion is omitted for space (represented as ...); please see the code examples.

```

1 # Iteration 12
2 {
3   "reasoning": "In this task, the controller function for the Sawyer robot arm manipulates
the robot's actions based on the observed environment to achieve a pick-and-place task.
The function correctly defines actions based on the robot's hand position, the puck's
position, the goal position, and the gripper's state. However, the feedback after each
step in the documentation indicates that although the robot is making progress towards
the goal, there is an issue with the gripper's action state not being switched to
'close' (1) at the crucial moment when the robot is close enough to the puck. Most
feedback comments suggest that the robot is moving in the right direction but fails to
grasp the puck due to the gripper not closing. This inefficiency appears to arise from
the conditions set for closing the gripper and possibly the normalized direction
towards the goal after the puck has been grasped. The condition for closing the gripper
(when close enough to the puck) seems to be met, but the action selected afterwards
does not effectively lead to puck manipulation or goal achievement in the most
optimized manner, as indicated by the 'Feedback' messages suggesting movements that
imply needing adjustments in action selection.",
4   "answer": "",
5   "suggestion": {...}
6 }
7
8 # Iteration 13 (Success)
9 LLM response:
10 {
11   "reasoning": "The instruction requests improvements based on feedback for a Sawyer robot
arm's controller function, aiming to accomplish a pick-and-place task. The controller
function accepts observations about the robot and the environment and computes an
action to move the robot accordingly. \n\nThe feedback provided indicates the task was
successfully completed with a suggestion that the final command should have been
'[0.04, 0.88, 0.24, 1]', which probably indicates the precise position and gripper
state (closed) needed to achieve the task goal. The feedback suggests the need for
fine-tuning actions especially towards the end of the task to ensure the robot's
gripper is appropriately positioned and the gripper state is correct for completing the
task. Since the feedback for the current code positively notes 'Success: True' with a
high return value, no modifications are recommended for '__code0'. \n\nPrevious feedback
instances with lower success metrics and recommendations for change highlight the
iterative improvements to the controller logic, specifically related to grip control
and movement precision. The aim was to adjust the controller's response to observation
inputs to accurately position the gripper and handle the puck. The current code
evidently manages these aspects well by normalizing directions and managing the
gripper's open-close state effectively under varying circumstances (proximity to puck
or goal and whether the puck is already grasped). These adjustments have likely led to
the successful output in the described problem.",
12   "answer": "No changes are recommended as the feedback indicates successful task
completion and suggests the controller's current logic adequately addresses the task's
requirements.",
13   "suggestion": {}
14 }

```

**Figure A.20:** Example Responses of LLMs for LLFBench Meta-World Pick-Place (Part 2); Content in suggestion is omitted for space (represented as ...); please see the code examples.

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [Yes]

Justification: The contributions (Trace, OPTO, OptoPrime) are clearly stated in the introduction and detailed in Sections 2,3,4, and claimed improvements are validated in Section 5. Meanwhile the aspirational goal of self-improving interactive agents is elaborated separately in the concluding Section 7.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [Yes]

Justification: The limitations are discussed in Section 6.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?



Answer: [Yes]

Justification: There are theorems about the computational complexity of Trace and an information theoretic lower bound on propagated feedback. Both theorems are proved in Appendix E.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Experimental setups are outlined in Section 5 and detailed in Appendix I. Supplementary material additionally provides the exact code that was run to produce all results, and only requires user to supply an OpenAI API key.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: All of the code to implement Trace and OptoPrime are provided in the supplementary material. All of the experiments described in the paper are reproducible with the supplied code.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: Experimental setups are outlined in Section 5 and detailed in Appendix I.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: All experiments include standard error bars from  $\geq 10$  replications.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: All the experiments use the same compute resources (LLM API for OptoPrime and machine to run Trace), which are listed at the start of Section 5.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The contributions of this paper (a framework analogous to PyTorch to generalize the backpropagation algorithm) do not have direct safety or security implications.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: In Sections 6 and 7 we discuss both the aspirational goals and their broader impacts.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

#### 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: There are no models, data or APIs released along with the paper.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

#### 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: Datasets used in the experiments (e.g. BigBench-Hard) are properly attributed.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. **New Assets**

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: There is code implementing Trace and OptoPrime that are released along with the paper. The code repository contains licenses and several tutorial notebooks documenting each functionality. There are no datasets or models in the release, hence datasheets and modelsheets are not applicable for this release.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. **Crowdsourcing and Research with Human Subjects**

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: There are no crowdsourcing or human subject studies conducted in this paper.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects**

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: The paper does not involve crowdsourcing nor research with human subjects.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.