# STARK: STRATEGIC TEAM OF AGENTS FOR REFINING KERNELS

#### **Anonymous authors**

000

002 003 004

010 011

012

013

014

015

016

017

018

019

021

023

025

026027028

029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

#### **ABSTRACT**

The efficiency of GPU kernels is central to the progress of modern AI, yet optimizing them remains a difficult and labor-intensive task due to complex interactions between memory hierarchies, thread scheduling, and hardware-specific characteristics. While recent advances in large language models (LLMs) provide new opportunities for automated code generation, existing approaches largely treat LLMs as single-shot generators or naive refinement tools, limiting their effectiveness in navigating the irregular kernel optimization landscape. We introduce an LLM agentic framework for GPU kernel optimization that systematically explores the design space through multi-agent collaboration, grounded instruction, dynamic context management, and strategic search. This framework mimics the workflow of expert engineers, enabling LLMs to reason about hardware trade-offs, incorporate profiling feedback, and refine kernels iteratively. We evaluate our approach on KernelBench, a benchmark for LLM-based kernel optimization, and demonstrate substantial improvements over baseline agents: our system produces correct solutions where baselines often fail, and achieves kernels with up to 16× faster runtime performance. These results highlight the potential of agentic LLM frameworks to advance fully automated, scalable GPU kernel optimization.

# 1 Introduction

Artificial intelligence (AI) has advanced at an unprecedented pace, transforming both research and real-world applications. While innovations in model architectures and training algorithms have been central to this progress, the efficiency of the computational infrastructure that executes them is equally critical. At the core of modern AI systems are *GPU kernels*, which implement fundamental operations such as matrix multiplication and convolution. Even modest improvements in GPU kernel efficiency can translate into significant reductions in training time, inference latency, and deployment cost, making kernel optimization a cornerstone for sustaining AI's rapid growth.

Despite their importance, designing and optimizing GPU kernels remains a major challenge. The performance of a kernel depends on sub-

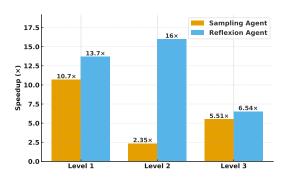


Figure 1: Speedup of STARK over baseline agents on KernelBench (L1–L3) with same number of attempts. Bars report GPU wall-clock speedups ( $\times$ ) relative to the *Sampling* and *Reflexion* agents; higher is better. STARK reaches up to  $16\times$  over Reflexion (L2) and  $10.7\times$  over Sampling (L1).

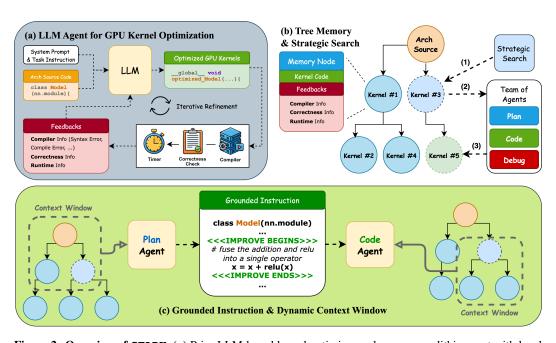
tle interactions between thread scheduling, memory hierarchy utilization, synchronization, and hardware-specific characteristics. Small changes in tiling strategies, loop unrolling, or memory alignment can yield disproportionate effects on runtime. As a result, the kernel optimization land-scape is highly irregular, architecture-dependent, and difficult to navigate. Existing approaches largely fall into two categories: *manual optimization* by expert engineers, which is effective but labor-intensive and difficult to scale; and *automated compilers and domain-specific languages* (DSLs) such as TVM and Triton (Chen et al., 2018; Tillet et al., 2019), which apply heuristics or search but often struggle with irregular operators and hardware variability (Zheng et al., 2020a;b).

The rapid progress of large language models (LLMs) opens a new opportunity for kernel optimization. Beyond their ability to generate correct code, LLMs can be guided to reason about hardware trade-offs, adapt to profiling feedback, and iteratively refine implementations. However, prior work has mostly treated LLMs as single-shot code generators or simple refinement tools (Ouyang et al., 2025), which underutilizes their potential for structured exploration of the kernel design space. To build a more powerful agent, we identify and address three critical limitations in existing methods:

- 1. **Naive exploration strategy.** Current agents typically refine code linearly, learning only from the immediately preceding attempt. This simplistic process neglects the rich history of prior attempts and fails to effectively balance the exploration-exploitation trade-off, often getting trapped in local optima.
- 2. **Monolithic agent design.** Kernel optimization is a multifaceted task requiring distinct capabilities for planning, implementation, and reflection. By assigning all these responsibilities to a single, generalist LLM, current agents operate inefficiently.
- 3. **Planning-implementation gap.** We observe a failure mode particularly acute in this domain: LLMs frequently devise a correct high-level optimization plan (e.g., "apply memory tiling") but fail to translate it into valid low-level CUDA code. This gap stems from the relative scarcity of expert-level kernel code in the models' training data.

To address these limitations, we introduce STARK (*Strategic Team of Agents for Refining Kernels*), a novel framework for automated GPU-kernel optimization. Our contributions are threefold:

- Collaborative multi-agent workflow. We design a workflow with specialized agents for planning, coding, and reflection, mirroring an expert development cycle and overcoming the inefficiencies of monolithic designs.
- Bridging the planning-implementation gap. We propose two mechanisms—grounded instruction and dynamic context windows—that translate high-level strategies into precise, actionable code edits, ensuring robust coordination across agents.
- Strategic search for refinement. We incorporate a search policy that balances exploration and exploitation over prior attempts, enabling systematic discovery of strong kernels.



**Figure 2: Overview of STARK.** (a) Prior LLM-based kernel optimizers rely on a monolithic agent with local iterative refinement. (b) STARK replaces this with a collaborative multi-agent workflow (plan/code/debug) coupled with strategic search over a tree memory. (c) The plan agent issues *grounded instructions* that anchor edits to code spans; *dynamic context windows* surface role-specific history; and the debug agent repairs failures. See Section 4 for details.

We evaluate our framework on **KernelBench** Ouyang et al. (2025), a benchmark designed to assess LLM-based GPU kernel optimization. Experiments show that combining these improvements leads to an agent system significantly more competitive than the baseline agents in both runtime performance and success rate across diverse kernel problems, authoring competitive kernels for the challenging problems in KernelBench where the baseline agents struggle to even find a working solution. Notably, STARK achieves more than  $10\times$  speedup over kernels produced by the baseline agents (i.e., the optimized kernels run in under one-tenth the time of the baseline.). Overall, our work suggests that LLM-driven agents represent a promising step toward fully automated GPU kernel optimization.

# 2 RELATED WORK

Due to space constrain, we only review the most relevant prior work here and defer the complete discussion to Appendix E.

The optimization of GPU kernels has progressed from empirical auto-tuning frameworks that perform black-box parameter searches (van Werkhoven, 2019; Nugteren & Codreanu, 2015) and compiler-based approaches with static heuristics (Yang et al., 2010), to the use of machine learning (ML). ML-based techniques have been used to replace hand-tuned heuristics in production compilers (Trofin et al., 2021), learn cost models to guide optimization (Chen et al., 2018), and even learn directly from raw source code without manual feature engineering (Cummins et al., 2017). A significant leap was the use of deep reinforcement learning to discover fundamentally new algorithms, as demonstrated by AlphaTensor's success in finding faster matrix multiplication methods (Fawzi et al., 2022). While powerful, these prior works either optimize within a fixed search space or operate in purely formal domains. Our work addresses these limitations by operating directly on source code to implement novel, structural changes.

The emergence of powerful Large Language Models (LLMs) has revolutionized programmatic interaction with source code, demonstrating a remarkable proficiency in generating code for diverse applications from competitive programming to compiler testing (Gu, 2023; Zhong & Wang, 2024; Jain et al., 2025). This capability has catalyzed a paradigm shift away from single-shot code generation and toward the development of autonomous LLM agents. An agent enhances a base LLM with planning, memory, and tool-use capabilities to direct its own workflow (Weng, 2023). The success of frameworks like SWE-agent in independently resolving complex GitHub issues has validated the power of this approach for software engineering (SWE) (Yang et al., 2024). While the application of LLM agents to SWE is a burgeoning field of research (Yang et al., 2024; Antoniades et al., 2024; Yang et al., 2025), their potential in the specialized domain of GPU kernel optimization remains largely unexplored. To fill this gap, we designed STARK, an agent framework with capabilities tailored to the unique challenges of this domain.

#### 3 PRELIMINARY

## 3.1 LLMs and Autoregressive Generation

Given an input sequence  $x=(x_1,x_2,\ldots,x_n)$  (e.g., the task instruction) as the context, an LLM  $p_\theta$  with parameters  $\theta$  generates an output sequence  $y=(y_1,y_2,\ldots,y_m)$  where  $y_t\in\mathcal{Y},t\in\{1,\ldots,m\}$  are tokens. Pretrained on a massive corpus of text, LLMs autoregressively generate the next token  $y_t$  conditioning on x and all the previously generated token  $y_{< t}=(y_1,\ldots,y_{t-1})$ . Specifically, at each time t, the LLM first computes the logits  $z_\theta(y|y_{< t},x)$  for each token y in the vocabulary  $\mathcal Y$  and generate  $y_t$  following the conditional distribution

$$p_{\theta}(y_t|y_{< t}, x) = \frac{\exp(z_{\theta}(y_t|y_{< t}, x)/\tau)}{\sum_{y' \in \mathcal{Y}} \exp(z_{\theta}(y'|y_{< t}, x)/\tau)}.$$
 (1)

The temperature parameter  $\tau>0$  modulates the randomness of an LLM's output. Higher values of  $\tau$  flatten the next token distribution in Equation 1, encouraging creative and diverse responses. Conversely, lower values sharpen the distribution, promoting deterministic and high-fidelity outputs.

This trade-off is critical in complex tasks, as different sub-problems demand different behaviors. For instance, **planning** and **exploration** benefit from a high temperature to generate novel strategies,

whereas tasks requiring precision and factual correctness, such as **code implementation**, necessitate a low temperature to ensure reliability. A single agent with a fixed temperature is ill-equipped to handle this dichotomy. This observation is a core motivation for STARK's multi-agent design, which allows specialized agents to operate at distinct temperatures tailored to their roles, i.e., a high  $\tau$  for the creative plan agent and a low  $\tau$  for the precise code agent.

#### 3.2 Kernelbench

**KernelBench** Ouyang et al. (2025) is a recently proposed benchmark specifically designed for assessing LLM-based GPU kernel optimization. Unlike prior evaluations that focus only on code correctness or small-scale operator tests, KernelBench provides a principled and reproducible testbed that measures both correctness and runtime efficiency across a broad spectrum of GPU workloads. KernelBench comprises a suite of optimization tasks, categorized into three difficulty levels. For each task, the objective is to create a custom GPU kernel that is functionally equivalent to a provided PyTorch reference implementation while minimizing its wall-clock execution time. See an example of the KernelBench task in Appendix D.

Specifically, **Level 1** tasks focus on single, common operators such as matrix multiplication and convolution, serving as a baseline for fundamental optimization capabilities; **Level 2** tasks comprise tasks with multiple operators fused into a single kernel, testing the ability to manage more complex dataflows and scheduling; **Level 3** tasks represent the highest difficulty, featuring popular full ML architectures such as the ResNet (He et al., 2016) and LSTM (Hochreiter & Schmidhuber, 1997), which involve highly irregular computations and intricate memory access patterns that are challenging for both human experts and automated systems to optimize effectively.

# 4 STARK: STRATEGIC TEAM OF AGENTS FOR REFINING KERNELS

**Framework Overview.** We now present STARK, an agentic framework for GPU-kernel optimization. STARK organizes kernel refinement into three layers: (i) a *multi-agent workflow* that separates planning, coding, and debugging, (ii) *coordination mechanisms* with *grounded instruction* to anchor planned edits to concrete code spans and *dynamic context windows* that surface role-specific history (e.g., prior attempts, failures, profiler feedback) to each agent, and (iii) a *strategic search* policy that balances exploration and exploitation across iterative attempts. Notably, multi-agent workflow and grounded instruction improve reliability even under a single-attempt budget, whereas dynamic context windows and strategic search deliver most of their gains when multiple attempts are allowed. Figure 2 provides an overview; the following subsections detail each component in turn.

#### 4.1 Multi-Agent Collaboration

Optimizing GPU kernels is inherently multifaceted and mirrors expert team workflows. A single agent typically fails to balance correctness, performance, and exploration across a vast, irregular design space. In particular, *strategy discovery* (e.g., fusion, vectorization, shared-memory tiling) benefits from higher-temperature generation that encourages diversity whereas *strategy realization*, i.e., committing those ideas to code, requires low-temperature precision to avoid errors. We therefore adopt a multi-agent framework that enables role specialization through LLMs.

**Multi-Agent Design (MAD).** Specifically, STARK decomposes kernel optimization into three roles – *plan*, *code*, and *debug*. Using a role-specific context window (Section 4.4) with selected prior attempts and execution outcomes, the **plan** agent proposes targeted transformations to either the source kernel or a candidate chosen by the strategic search policy (Section 4.2), emitting *grounded instructions* (Section 4.3) that anchor edits to explicit code spans. The **code** agent consumes grounded instructions and translates them into executable GPU-kernel code, conditioning on its own context window to improve adherence and code quality. The **debug** agent repairs promising but failing candidates by consulting the plan agent's instructions and compiler/runtime diagnostics, producing a working kernel that realizes the intended transformation.

Benefits of MAD. Role specialization lets each agent use prompts and base LLMs matched to its objective. In our instantiation, we choose Claude Sonnet 4 with temperature  $\tau$ =0.8 for the plan agent to encourage strategy diversity, and the same model with  $\tau$ =0.1 for the code and debug agents

to enforce precision. Despite this simple setup, MAD already performs strongly (see Section 5). We underscore that because the design is modular, we can swap in planners with richer kernel-optimization priors or code-specialist reasoning models to further improve results. In addition, modularity also exposes bottlenecks. We observe that the dominant bottleneck is code-synthesis fidelity: LLMs often need multiple attempts to faithfully implement a given instruction. Finally, MAD makes targeted post-training straightforward: we can fine-tune the base LLM for a specific agent (e.g., the code agent) without affecting the others, improving stability and predictability. However,, a systematic study of agent-specific post-training is orthogonal to our core contributions and is left to future work.

# 4.2 STRATEGIC SEARCH WITH TREE MEMORY

Prior LLM-driven kernel optimizers typically use either best-of-K sampling that generates multiple candidates independently and select the fastest correct one or  $iterative\ refinement$  which repeatedly edits the latest kernel (Ouyang et al., 2025). However, best-of-K is unguided and wasteful: all the new attempts ignore feedback from earlier attempts and repeatedly probe redundant regions of the design space. On the other hand, iterative refinement is feedback-aware but myopic: by building only on the most recent candidate, it is prone to getting trapped in narrow, suboptimal basins.

To address these limitations, STARK reframes kernel optimization as **strategic search** over a persistent **tree memory**. We maintain a search tree T whose nodes store candidates and their observations (runtime, correctness, and compiler diagnostics). The root represents the source architecture; each edge corresponds to applying a grounded instruction from the plan agent and realizing it via the code agent (or repairing via the debug agent). Each node n is assigned a score s(n) reflecting competitiveness; in our implementation we use the straightforward kernel runtime as s(n) and treat *lower is better*. For kernels that are incorrect or failing to compile, we give them scores of  $+\infty$ . At each step, we (1) *select* a node to expand using a strategic policy, (2) *expand* by invoking the plan/code (or debug) agents to produce a child candidate, (3) *evaluate* for correctness and runtime, and (4) *record* results in T to inform subsequent selections. This converts ad-hoc trial-and-error into a directed, feedback-driven process.

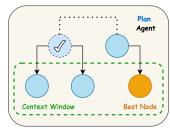
**Policy choice and an adapted**  $\epsilon$ **-greedy rule.** We compared representative search policies including Monte-Carlo Tree-Search (MCTS), evolutionary, greedy, and  $\epsilon$ -greedy policies and found that  $\epsilon$ -greedy consistently performs best under the same budget constraint. Importantly, we observe that kernel optimization poses domain-specific challenges that are root dominance (it is very challenging to even outperform the source architecture in the root node) and frequent compilation/runtime failures. To address these challenges, we adapt the canonical rule as follows: (1) **Root throttling:** cap the number of direct children of the root at  $n_{\text{root}}$  to avoid redundant first-hop edits; once the cap is reached, the root is ineligible for selection; (2) **Dead-branch pruning:** if a node has more than  $n_{\text{child}}$  children and all current children fail, mark the node ineligible to prevent wasting trials; (3) **High exploration rate:** use a relatively large  $\epsilon$  (empirically 0.3–0.4) to counteract local traps; (4) **Leaf-biased exploration:** with probability  $\epsilon$ , sample uniformly from expandable leaves (not only failing nodes), encouraging discovery beyond the immediate failure set.

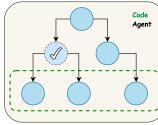
#### 4.3 GROUNDED INSTRUCTION

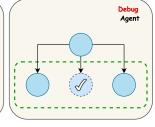
We introduce grounded instruction for kernel enhancement. The plan agent must not only propose an optimization, but also insert **explicit span anchors** in the kernel source that mark exactly where the change should occur. Each anchor is a short, machine-checkable tag (i.e, <<<IMPROVE BEGINS>>> ... <<<IMPROVE ENDS>>>) wrapped around the target site, such as a load/store, loop body, or the launch configuration. The code agent consumes this annotated scaffold and resolves each anchor by emitting concrete CUDA that realizes the instruction. Grounded instruction tightens plan—code alignment, curbs hallucinated guidance, and narrows the coder's search space. It also improves traceability: every proposal leaves a visible, verifiable footprint in the final code. In practice, we observe fewer misinterpretations and markedly fewer faulty kernels. Despite its simplicity, the mechanism is especially effective on Level 3 KernelBench tasks with deeper architectures (e.g., VGG).

#### 4.4 DYNAMIC CONTEXT WINDOW

Past attempts provide rich, actionable signals for subsequent decisions, but different agents benefit from different *views* of this history. We therefore maintain a *dynamic, agent-specific context window* that is rebuilt at each selection step for different agents. See Figure 3 for a visual demonstration. Throughout this section, let node i be the node selected by the search policy defined in Section 4.2. We use  $\mathcal{W}(i)$  to denote the context window containing a subset of historical attempts and their evaluation outcomes (e.g., compiler information and runtime). As we always include the source architecture as part of the prompt for agents,  $\mathcal{W}(i)$  always includes the root node  $n_{\text{root}}$ . For a naive search algorithm without dynamic context window,  $\mathcal{W}(i) = \{i, n_{\text{root}}\}$  only includes node i in addition to the root.







**Figure 3:** Dynamic Context Window. Nodes with  $\checkmark$ 's represent selected nodes.

**Tree relations.** We use tree relations to build agent-specific context windows. Let p(i) be the parent of node i. Define the *siblings* of i as  $S(i) = \{j : p(j) = p(i)\}$ . Moreover, define the set of child nodes of a node i as D(i). We also maintain a small global leaderboard C of top-performing nodes.

Plan agent (local & contrastive global context). For a selected node i, the plan agent conditions on a context window  $\mathcal{W}_{\text{plan}}(i)$  that aggregates node i's children and a small set of global leaders from the leaderboard  $\mathcal{C}$ . Formally,

$$\mathcal{W}_{\text{plan}}(i) = \{i, n_{\text{root}}\} \cup \mathcal{D}(i) \cup \text{Top} - r(\mathcal{C}),$$

where  $\mathcal{D}(i)$  contains all evaluated children of i with their observations, and  $\operatorname{Top}-r(\mathcal{C})$  returns the r highest-scoring distinct kernels from the global leaderboard (excluding i's subtree) to discourage duplication.

This design serves three purposes. (i) **reflection**: the plan agent can revise or stack its prior instructions rather than rediscovering them; (ii) **ambition calibration**: top competitors prevent redundant exploration and provide transferable motifs such as warp-shuffle reductions, vectorized LD/ST, and shared-memory tiling; (ii) **capability estimation**: by inspecting how past instructions were realized or failed by the code agent, the next instruction is adapted to what the code agent can reliably execute, *improving first-pass success and avoiding instructions beyond current ability*. To achieve this, we explicitly require the plan agent to adapt its instruction to the code agent's demonstrated capabilities observed in  $\mathcal{D}(i)$ .

Code agent (extended context). For kernel code emission at node i, the code agent conditions on

$$\mathcal{W}_{\text{code}}(i) = \{i, n_{\text{root}}\} \cup \mathcal{D}(i) \cup \{j : p(j) \in S(i)\}.$$

The nodes in  $\{j: p(j) \in S(i)\}$  are essentially the children of node i's siblings. Our insight is that these nodes typically share near-identical scaffolds with node i from a common planning lineage, so successful patches and micro-optimizations transfer with high probability; conversely, seeing failures in closely related contexts helps the coder avoid repeating the same mistakes. Hence, this extended window serves two aims: (i) **reduce implementation errors** by letting the coder imitate successful patches from closely related scaffolds and avoid previously observed failure modes; (ii) **surface stronger implementations** by transferring micro-optimizations (e.g., warp-shuffle motifs, vectorized LD/ST, shared-memory tiling) that have already worked on cousin nodes.

**Debug agent (local context).** For fault repair, we construct the context window for the debug code as

$$\mathcal{W}_{\text{debug}}(i) = \{i, n_{\text{root}}\} \cup \mathcal{S}(i),$$

We choose this design mainly for two reasons. Most fixes are structural and local, e.g., off-by-one guards, stride/indexing alignment, launch-parameter tweaks, or shared-memory sizing often transfer directly among siblings that share the same scaffold. Moreover, restricting the window to  $\mathcal S$  avoids distracting the debug agent with globally unrelated kernels, improving precision and reducing hallucinated edits.

#### 4.5 Framework Overview

324

325

326

327

328

329 330

331 332

333 334

335

336

337

338

339

340 341

342

364

366

367 368

369

370

371 372

373

374

375

376

377

Here we provide an overview of our framework STARK and describe its execution process. Algorithm 1 presents its pseudocode.

At a high level, STARK repeatedly (i) selects a promising node (a prior attempt) from a search tree, (ii) builds agent-specific context windows from local history and global leaders, (iii) asks the *planning agent* to propose a concrete optimization along with *grounded instruction* anchors inserted into a scaffold, (iv) asks the *code agent* to realize those anchors into an executable kernel, (v) if the selected node has a problematic kernel, we build the debugger's dynamic context window and request a minimal fix. The new attempt is evaluated, appended as a child node, and the leaderboard  $\mathcal C$  is updated. We repeat this process until we reach a pre-specified max attempts  $\mathcal B$ .

## Algorithm 1 STARK: Strategic Team of Agents for Refining Kernels

```
343
          Require: Budget B (max attempts), selection policy \pi_{\text{select}} (adapted \varepsilon-greedy), leaderboard size r
344
           1: Initialize search tree T with root n_{\text{root}} (PyTorch reference)
           2: Initialize leaderboard \mathcal{C} \leftarrow \{n_{\text{root}}\}
345
           3: for t = 1, 2, ..., B do
346
           4:
                    i \leftarrow \pi_{\text{select}}(T, \mathcal{C})
                                                                                                      ⊳ pick a node to refine
347
           5:
                    if HASBUG(i) then
                                                                              \triangleright compile fail or unit-test fail recorded at i
348
           6:
                        W_{\text{debug}}(i) \leftarrow \text{BUILDCONTEXTDEBUG}(i, T)
349
           7:
                        kernel' \leftarrow DebugAgent(W_{dbg}, i.kernel, i.logs)
350
           8:
                         (ok, correct, runtime, logs) ← EVALUATE(kernel') > compile, correctness check,
351
               timing
352
           9:
                         (plan, anchors) \leftarrow (i.plans, i.anchors)
353
          10:
                    else
354
                         \mathcal{W}_{\mathsf{plan}}(i) \leftarrow \mathsf{BUILDCONTEXTPLAN}(i, T, \mathcal{C})
          11:
355
                         (plan, anchors) \leftarrow PLANAGENT(W_{plan})
          12:
356
           13:
                         \mathcal{W}_{\text{code}}(i) \leftarrow \text{BuildContextCode}(i, T)
357
                         kernel' \leftarrow CODEAGENT(W_{code}, plan, anchors)
          14:
          15:
                         (ok, correct, runtime, logs) \leftarrow EVALUATE(kernel')
359
          16:
360
          17.
                    j \leftarrow ADDCHILD(T, i, kernel', plan, anchors, ok, correct, runtime, logs)
                    C \leftarrow \text{UPDATELEADERS}(C, j, r)
361
          18:
          19: end for
362
          20: return BEST(\mathcal{C})
                                                                                         363
```

#### 5 EXPERIMENTS

We use KernelBench (Ouyang et al., 2025), a recently proposed benchmark consisting of comprehensive and challenging GPU kernel tasks, to validate the effectiveness of our proposed approaches.

Baselines and Metrics. We compare our framework STARK with the following list of approaches:

- Torch Eager: the out-of-box PyTorch modules without any compilation or optimization.
- Torch Compile: We use torch.compile to produce optimized versions of the given PyTorch modules. While torch.compile offers different compilation modes, we compare to two of the most representative and competitive ones default and max-autotune.
- Sampling Agent: the single agent framework originally proposed and used by Kernel-Bench to evaluate the difficulty of the tasks in KernelBench and the ability of LLMs to

write efficient kernels. This agent repeatedly samples responses when given the source model to optimize and chooses the best generated custom kernel as the solution.

381 382

384

386

• **Reflexion Agent**: this agent follows the Reflexion paradigm (Shinn et al., 2023), where at each optimization step, it tries to update its last attempt using its corresponding observations such as the compiler and runtime information.

387

389 390 391

392 394 395

397 398 399

400

407

427

431

We report the following metrics to comprehensively understand the agents' performances: (i)  $\mathbf{Fast}_1$ rate is the percentage of the problems for which the agent can generate kernels that are at least as fast as the torch baselines; (ii) Success rate represents the percentage of the problems for which the agent can generate compiled and correct kernels; (iii) Speed: To better understand how good the generated kernels are, we also report the average speed across all tasks.

Comparison with Torch Baselines. In Table 1, we present the results about success rate, Fast<sub>1</sub> rate and speed over all 3 levels of KernelBench challenges. For each task, we let all agents to have a maximum of B=30 attempts. Due to limited computation resource, we evaluate on the representative subset of KernelBench (Ouyang et al., 2025). We use Claude Sonnet 4 as the base LLMs for all the LLM-based baselines and our agents. Due to space constraint, we defer implementation and evaluation details to Appendix B.

The results in Table 1 demonstrate that our proposed framework, STARK, consistently outperforms both the Sampling and Reflexion baselines across all KernelBench difficulty levels. At Level 1, STARK not only achieves a perfect 100% success rate but also delivers up to a  $3.0\times$  speedup over Torch Eager baselines, while Sampling and Reflexion agents frequently generate kernels that are slower than the baselines. This advantage becomes even more pronounced at Level 2, where the complexity of the kernels increases. Here, STARK maintains a perfect success rate and achieves speedups of 2.7×, whereas the Reflexion agent, despite attaining 100% correctness, produces kernels that run slower than the baseline. At Level 3, which involves the most irregular and challenging tasks, both Sampling and Reflexion degrade significantly, with success rates falling and runtimes dropping below baseline. In contrast, STARK continues to maintain full success while producing kernels that outperform the Torch implementations by up to  $1.6\times$ . These results highlight that STARK not only generates correct kernels but also delivers substantial performance improvements, even as task difficulty increases.

		Torch Eager		Default		Max-autotune	
Level 1	Success ↑	$ $ Fast <sub>1</sub> $\uparrow$	Speed ↑	$\mathbf{Fast}_1 \uparrow$	<b>Speed</b> ↑	$ $ Fast $_1 \uparrow$	<b>Speed</b> ↑
Sampling Agent	57.1%	14.3%	0.81×	7.1%	0.46×	7.1%	0.81×
Reflexion Agent	92.6%	28.6%	$1.24 \times$	14.3%	$0.57 \times$	35.7%	$0.92 \times$
STARK	100%	71.4%	$3.03 \times$	78.6%	$2.37 \times$	78.6%	$2.76 \times$
Level 2	Success	$\mathbf{Fast}_1$	Speed	$\mathbf{Fast}_1$	Speed	Fast <sub>1</sub>	Speed
Sampling Agent	87.5%	50%	$1.06 \times$	37.5%	$0.91 \times$	37.5%	$0.91 \times$
Reflexion Agent	100%	75%	$0.88 \times$	62.5%	$0.78 \times$	62.5%	$0.78 \times$
STARK	100%	100%	$2.69 \times$	87.5%	$2.51 \times$	87.5%	$2.52 \times$
Level 3	Success	$\mathbf{Fast}_1$	Speed	$\mathbf{Fast}_1$	Speed	Fast <sub>1</sub>	Speed
Sampling Agent	100%	50%	$0.87 \times$	12.5%	$0.67 \times$	12.5%	$0.66 \times$
Reflexion Agent	67.5%	25%	$0.79 \times$	12.5%	$0.62 \times$	12.5%	$0.61 \times$
STARK	100%	87.5%	$1.58 \times$	87.5%	$1.27 \times$	87.5%	$1.26 \times$

Table 1: Performance of LLM Agents on the KernelBench Tasks. Fast 1 represents the percentage of problems for which the agent can generate custom kernels that are correct and as fast as the Torch baselines (higher is better). Speed is computed as the ratio of the kernel runtime of the baseline to that of the generated kernel.

**Comparison between Agents.** We investigate deeper into the behavior of our agent STARK with the two baseline agents to better understand their optimization behaviors. A deeper analysis of compile and correctness rates, shown in Table 2, provides further insight into why STARK succeeds where baselines struggle. While all agents achieve relatively high compile rates (mostly above 80%),

the fraction of kernels that are both compilable and correct varies widely. The Sampling agent, for example, compiles over 90% of its outputs on Level 1 but only 43% of these are functionally correct. Reflexion improves correctness slightly through iterative refinement, but its correctness rate remains below 55% at all levels. In contrast, STARK achieves the highest correctness rates across the board, reaching 61.2% on Level 2 tasks. This suggests that STARK 's structured planning and feedback-driven refinement not only increase the chance of generating efficient kernels but also reduce wasted attempts on invalid or incorrect code. Finally, Figure 1 highlights the dramatic runtime improvements of STARK relative to baseline agents. On Level 1 tasks, STARK achieves over a  $10\times$  speedup compared to Sampling and a  $13.7\times$  speedup over Reflexion. On Level 2, these gains rise as high as  $16\times$ , and even at the most challenging Level 3 tasks STARK maintains  $5-6\times$  improvements. These relative gains indicate that while baselines occasionally achieve correctness, they rarely deliver true runtime efficiency. By contrast, STARK 's ability to jointly optimize for correctness and speed allows it to close both gaps simultaneously. Taken together, these findings confirm that multi-agent collaboration and strategic search are key enablers for scaling LLMs to the demands of GPU kernel optimization.

	<b>Compile Rate</b> ↑			<b>Correct Rate</b> ↑			
KernelBench Level	1	2	3	1	2	3	
Sampling Agent	90.8%	97.0%	84.9%	43%	44.0%	15.1%	
Reflexion Agent	86.0%	86.2%	78.9%	48.3%	53.4%	28.4%	
STARK	84.5%	90.7%	83.4%	50.6%	61.2%	35.5%	

**Table 2:** Percentages of Successfully Compiled and Correct Kernels.

**Ablations.** We ablate the agentic components of our system. We compare (i) **Search Agent**, which is a single-agent model equipped with our strategic search, and (ii) **MA-only**, which employs the multi-agent workflow (plan/code/debug with grounded instruction and dynamic context) using best-of-K sampling instead of search. As shown in Table 3, both variants outperform the *Sampling* baseline, confirming that each component helps. When combined in STARK, the effects compound: strategic search exploits the structured proposals produced by the multi-agent workflow, yielding the largest gains.

	Torch Eager		Defa	ault	Max-autotune	
	$\mathbf{Fast}_1 \uparrow$	<b>Speed</b> ↑	$\mathbf{Fast}_1 \uparrow$	<b>Speed</b> ↑	$\mathbf{Fast}_1 \uparrow$	<b>Speed</b> ↑
Sampling Agent	50%	0.87×	12.5%	$0.67 \times$	12.5%	0.66×
Search Agent	67.5%	$0.89 \times$	25%	$0.71 \times$	25%	$0.70 \times$
MA-Only	67.5%	$1.11 \times$	25%	$0.92 \times$	25%	$0.91 \times$
STARK	87.5%	$1.58 \times$	87.5%	$1.27\times$	87.5%	$1.26 \times$

**Table 3:** Ablation on the Proposed Agentic Features.

#### 6 Conclusion

In this work, we introduced an agentic framework for GPU kernel optimization that combines multiagent role play, dynamic context management, and strategic search. Our evaluation on KernelBench demonstrated that the proposed framework consistently outperforms baseline methods in both success rate and runtime efficiency, across tasks of varying complexity. These results highlight the value of moving beyond single-agent or unguided sampling approaches, and point to the promise of collaborative, feedback-driven optimization. Looking forward, we envision that agentic LLM frameworks will play an increasingly important role in automated system optimization. Extending our approach to broader classes of operators, diverse hardware architectures, and cross-kernel scheduling decisions are natural directions for future research. More broadly, our work suggests that multi-agent LLMs can meaningfully accelerate the co-design of AI algorithms and infrastructure, pushing the boundaries of what is possible in efficient large-scale computation.

# ETHICS STATEMENT

This study complies with the ICLR Code of Ethics. All datasets employed are publicly available and open-source under licenses that permit research use. No private or personally identifiable information was accessed, and no new data were collected from human subjects. The research does not pose privacy, security, or fairness concerns. The authors declare no conflicts of interest and no external sponsorship.

#### REPRODUCIBILITY STATEMENT

We document all algorithmic and implementation details in the paper and appendix, including the exact prompts for every agent, full hyperparameters, and ablation settings. Upon acceptance, we will release an open-source repository with configuration files, an environment file and Dockerfile, and step-by-step commands to recreate every table and figure.

#### REFERENCES

- Miltiadis Allamanis, Earl T Barr, Prem Devanbu, and Charles Sutton. A survey of machine learning for big code and naturalness. *ACM Computing Surveys (CSUR)*, 51(4):1–37, 2018.
- Jason Ansel et al. Baco: A bayesian optimization framework for compiler autotuning. In Proceedings of the 32nd International Conference on Compiler Construction, 2021.
- Antonia Antoniades, Albert Örwall, Kexun Zhang, Yuxi Xie, Anirudh Goyal, and William Wang. Swe-search: Enhancing software agents with monte carlo tree search and iterative refinement. arXiv preprint arXiv:2410.20285, 2024.
- Simran Arora, Alex L Zhang, William Hu, and Christopher Ré. Kernelbench: Can Ilms write efficient gpu kernels? *arXiv preprint arXiv:2502.10517*, 2025.
- Ian Buck. Gpu program optimization. In Matt Pharr (ed.), GPU Gems 2. NVIDIA, 2008.
- Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Hangu Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Tvm: An automated end-to-end optimizing compiler for deep learning. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, 2018.
- Chris Cummins, Pavlos Petoumenos, Zheng Wang, and Hugh Leather. End-to-end deep learning of optimization heuristics. In 2017 26th International Conference on Parallel Architectures and Compilation Techniques (PACT), pp. 219–232. IEEE, 2017.
- DeepMind. Discovering faster matrix multiplication algorithms with alphatensor. https://www.deepmind.com/blog/discovering-faster-matrix-multiplication-algorithms-with-alphatensor, 2022. Accessed: 2024.
- Alhussein Fawzi, Matej Balog, Aja Huang, Thomas Hubert, Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Francisco JR Ruiz, Julian Schrittwieser, Grzegorz Swirszcz, et al. Discovering faster matrix multiplication algorithms with reinforcement learning. *Nature*, 610(7930):47–53, 2022.
- Martin Fowler. Some initial thoughts on ai in software development. 2025. Accessed: 2024.
- Flavius Gruian. Estimating performance gains of compiler optimizations for gpu applications. In *MSc Thesis, Lund University*, 2021.
- Qiuhan Gu. Llm-based code generation method for golang compiler testing. In *Proceedings of the 31st ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp. 2201–2203, 2023.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
  - Stijn Heldens and Ben van Werkhoven. Kernel launcher: C++ library for optimal-performance portable cuda applications. In 2023 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW), pp. 556–565. IEEE, 2023.
  - Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
  - Shu Ishida, Gianluca Corrado, George Fedoseev, Hudson Yeo, Lloyd Russell, Jamie Shotton, João F Henriques, and Anthony Hu. Langprop: A code optimization framework using large language models applied to driving. *arXiv preprint arXiv:2401.10314*, 2024.
  - Naman Jain, King Han, Alex Gu, Wen-Ding Li, Fanjia Yan, Tianjun Zhang, Sida Wang, Armando Solar-Lezama, Koushik Sen, and Ion Stoica. Livecodebench: Holistic and contamination free evaluation of large language models for code. In *The Thirteenth International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=chfJJYC3iL.
  - Carlos E. Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. SWE-bench: Can language models resolve real-world github issues? In *International Conference on Learning Representations*, 2024.
  - Jakub Kurzak, Stanimire Tomov, and Jack Dongarra. Autotuning gemm kernels for the fermi gpu. *IEEE Transactions on Parallel and Distributed Systems*, 23(11):2045–2054, 2012.
  - LangChain Inc. Langgraph: Build resilient language agents as graphs, 2025. URL https://github.com/langchain-ai/langgraph. Version 0.6.7; accessed 2025-09-24.
  - Milo Lurati. Bringing concurrency to kernel tuner: Gpu auto-tuning in multi-node systems. *MSc Thesis, University of Amsterdam*, 2024.
  - Teodor V Marinov and Alekh Agarwal. Offline imitation learning from multiple baselines with applications to compiler optimization. In *NeurIPS 2024 Workshop on Machine Learning for Systems*, 2024.
  - Cedric Nugteren and Valeriu Codreanu. Cltune: A generic auto-tuner for opencl kernels. In 2015 IEEE International Parallel and Distributed Processing Symposium Workshop, pp. 846–851. IEEE, 2015.
  - NVIDIA Developer. Automating gpu kernel generation with deepseek-rl and inference-time scaling. https://developer.nvidia.com/blog/automating-gpu-kernel-generation-with-deepseek-rl-and-inference-time-scaling/, 2024a. Accessed: 2024.
- NVIDIA Developer. Improving gemm kernel auto-tuning efficiency on nvidia gpus
  with heuristics and cutlass 4.2. https://developer.nvidia.com/blog/
  improving-gemm-kernel-auto-tuning-efficiency-on-nvidia-gpus-with-heuristics-and-cu
  2024b. Accessed: 2024.
  - Anne Ouyang, Simon Guo, Simran Arora, Alex L Zhang, William Hu, Christopher Re, and Azalia Mirhoseini. Kernelbench: Can LLMs write efficient GPU kernels? In *Forty-second International Conference on Machine Learning*, 2025. URL https://openreview.net/forum?id=yeoNliQTlx.
  - Zhaoran Qian et al. Message-passing graph neural networks can simulate interior-point methods for linear optimization. *arXiv preprint arXiv:2402.13825*, 2024.
    - Ari Rasch, Markus Haidl, and Sergei Gorlatch. Atf: A generic auto-tuning framework. In 2017 IEEE 19th International Conference on High Performance Computing and Communications (HPCC), pp. 64–71. IEEE, 2017.

- Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik Narasimhan, and Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning. *Advances in Neural Information Processing Systems*, 36:8634–8652, 2023.
- Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pp. 10–19, 2019.
- Mircea Trofin, Yundi Qian, Eugene Brevdo, Zinan Lin, Krzysztof Choromanski, and David Li. Mlgo: a machine learning guided compiler optimizations framework. *arXiv preprint arXiv:2101.04808*, 2021.
- Ben van Werkhoven. Kernel tuner: A search-optimizing gpu code auto-tuner. In *Future Generation Computer Systems*, volume 90, pp. 347–358. Elsevier, 2019.
- Lilian Weng. Llm-powered autonomous agents. https://lilianweng.github.io/posts/2023-06-23-agent/, 2023. Accessed: 2024.
- John Yang, Carlos E Jimenez, Alexander Wettig, Kilian Lieret, Shunyu Yao, Karthik R Narasimhan, and Ofir Press. SWE-agent: Agent-computer interfaces enable automated software engineering. In *Thirty-eighth Conference on Neural Information Processing Systems*, 2024.
- John Yang, Kilian Lieret, Carlos E Jimenez, Alexander Wettig, Kabir Khandpur, Yanzhe Zhang, Binyuan Hui, Ofir Press, Ludwig Schmidt, and Diyi Yang. Swe-smith: Scaling data for software engineering agents. *arXiv preprint arXiv:2504.21798*, 2025.
- Yi Yang, Ping Xiang, Jingfei Kong, and Huiyang Zhou. A gpgpu compiler for memory optimization and parallelism management. In *Proceedings of the 31st ACM SIGPLAN Conference on Programming Language Design and Implementation*, pp. 86–97. ACM, 2010.
- Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, et al. Ansor: Generating {High-Performance} tensor programs for deep learning. In *14th USENIX symposium on operating systems design and implementation (OSDI 20)*, pp. 863–879, 2020a.
- Size Zheng, Yun Liang, Shuo Wang, Renze Chen, and Kaiwen Sheng. Flextensor: An automatic schedule exploration and optimization framework for tensor computation on heterogeneous system. In *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 859–873, 2020b.
- Li Zhong and Zilong Wang. Can llm replace stack overflow? a study on robustness and reliability of large language model code generation. In *Proceedings of the AAAI conference on artificial intelligence*, volume 38, pp. 21841–21849, 2024.

# A USE OF LARGE LANGUAGE MODELS (LLMS)

Large Language Models (LLMs) were used solely as general-purpose assistive tools to improve the clarity and readability of the manuscript. Specifically, we used an LLM to help rephrase and polish text that we had already drafted.

# B IMPLEMENTATION AND EVALUATION

Agent Implementation. We use Claude Sonnet 4 as the base LLMs for all agent baselines and agents of STARK. For both sampling and reflexion agents, we follow KernelBench to set temperature  $\tau=0.7$  during generating, with other generation parameters such as top-p set to the default value. For STARK, we use Claude Sonnet 4 with temperature  $\tau=0.8$  for the plan agent, and  $\tau=0.1$  for the code and debug agents. For all tasks, all agent baselines and STARK have a maximum of B=30 attempts to optimize each task. Regarding the hyperparameters of STARK, we choose the root throtting number to be 5, dead-branch pruning number to be 3,  $\epsilon=0.3$  for the search, r=2 for the leaderboard  $\mathcal C$ . To prevent exploding context to the LLMs, we set an upper limit for the number of nodes in the dynamic context window: whenever the dynamic context window has more than 5 nodes, we randomly sample 5 from all the nodes in the window. We implement STARK with LangGraph (LangChain Inc., 2025).

**Runtime Evaluation.** We evaluate all the Pytorch baselines and LLM-generated kernels on the same NVIDIA A100 40GB GPU. We use the source code provided by KernelBench at its official repo<sup>1</sup> to benchmark the kernels' runtime. In particular, to ensure stable measurement, runtime is measured with CUDA events after warm-up runs using fixed input shapes; we choose a large number of 100 warm-up runs to ensure accurate measurement.

#### C PROMPTS

 Our prompts follow the templates of KernelBench (Ouyang et al., 2025), which has four components: system message, in-context example, architecture source code, and instruction.

As we have multiple agents in STARK with different tasks, they require different prompts to fulfill their tasks. Specifically, we reuse the system message and in-context example from KernelBench for all agents and include the architecture source code regardless of which node is selected for optimization. To motivate the agents to use the already optimized modules such as cuBLAS, we include an additional instruction in the system message to consider using existing highly optimized kernels. We show the system prompt in Figure 5, the in-context example in Figures 9 and 10. In addition, we include the information within the *dynamic context window* and *role-specific instructions* for different agents. See Figure 4 for the prompt template of STARK. We show the role-specific instruction in Figures 6, 7, and 8.

# D EXAMPLE KERNELBENCH TASKS

Here we show three examples of KernelBench tasks, one for each level. See Figures 11, 12, and 13 for example tasks in Level 1, 2, and 3. We refer interested readers to Ouyang et al. (2025) for the complete list.

#### E RELATED WORK

Optimizing GPU kernels to extract maximum performance from underlying hardware is a long-standing and formidable challenge. The vast, non-convex, and hardware-specific search space of possible kernel implementations necessitates sophisticated optimization strategies. The evolution of these strategies can be broadly categorized into three paradigms: empirical auto-tuning, compiler-and model-driven optimization, and most recently, generative approaches using Large Language

https://github.com/ScalingIntelligence/KernelBench

```
702
           {System Message}
         2
704
          Here's an example to show you the syntax of inline embedding
          custom CUDA operators in torch: The example given architecture is:
705
          {Example Architecture Source Code}
706
          The example new arch with custom CUDA kernels looks like this:
707
          {Example New Architecture Source Code}
708
709
          You are given the following architecture:
         8
710
           {Architecture Source Code}
         9
        10
711
          Here is your latest attempt:
        11
712
          {Source Code of the Selected Node}
        12
713
        13
714
           [Dynamic Context Window] You should use the following observations
715
           regarding your historical attempts to provide better
           implementations:
716
          - Learn from the failed examples to avoid bugs and write
717
           successful kernels.
718
          - Learn from the successful examples to design improved kernels.
719
        17
720
           **Kernel Source Code #1**
        18
          {Source Code of Historical Attempt}
        19
721
        20
722
          **Compiler Observation**
723
          {Compiler Log}
        22
724
        23
725
        24
           **Kernel Execuation Result**
           {Runtime or Correctness Error}
        25
726
        26
727
           **Kernel Source Code #2**
        27
728
          [...skipped]
        28
729
        29
730
          {Role-specific Instruction}
```

**Figure 4:** Prompt Template for Agents.

Models (LLMs). Our work builds upon this trajectory by introducing a fully autonomous agent that manages the entire optimization lifecycle.

#### E.1 EMPIRICAL AND COMPILER-BASED OPTIMIZATION

The foundational approach to GPU performance tuning is empirical auto-tuning, which treats the problem as a black-box search over a set of tunable parameters, such as thread block dimensions, memory tiling factors, and loop unrolling factors Lurati (2024); van Werkhoven (2019). Traditional methods often rely on an exhaustive or brute-force search, where thousands of potential kernel configurations are generated, compiled, and benchmarked to identify the top performer Kurzak et al. (2012). While effective, this process is prohibitively time-consuming; for instance, an exhaustive search for a single GEMM kernel can take over 700 minutes to complete NVIDIA Developer (2024b).

To mitigate this cost, heuristic-driven methods prune the search space. NVIDIA's nvMatmulHeuristics, for example, uses a predictive model to recommend a small subset of high-potential configurations, achieving near-optimal performance in a fraction of the time NVIDIA Developer (2024b). Frameworks like Kernel Tuner van Werkhoven (2019), ATF Rasch et al. (2017), and CLTune Nugteren & Codreanu (2015) provide robust environments for orchestrating these searches and support more advanced strategies like Bayesian Optimization, which builds a probabilistic performance model to guide the search more intelligently Ansel et al. (2021); Heldens & van Werkhoven (2023).

```
## System Message

You are an expert in writing efficient code.
You write custom CUDA kernels to replace the pytorch operators in the given architecture to get speedups.

You have complete freedom to choose the set of operators you want to replace. You may make the decision to replace some operators with custom CUDA kernels and leave others unchanged. You may replace multiple operators with custom implementations, consider operator fusion opportunities (combining multiple operators into a single kernel, for example, combining matmul+relu), or algorithmic changes (such as online softmax). You are only limited by your imagination.

You should consider using CUDA's existing highly optimized kernels and operations whenever appropriate. Try building on these optimized blocks and further improve it with your custom kernels.
```

Figure 5: System Message for All Agents.

Concurrently, compiler-based approaches aim to automate optimization through a series of program transformations applied to an intermediate representation Yang et al. (2010). GPU compilers employ passes for memory coalescing, data prefetching, vectorization, and loop optimizations to adapt naive code to the hardware architecture Buck (2008). However, these two paradigms exist in a state of tension. Compilers must employ conservative heuristics that are generally applicable but rarely optimal for a specific kernel, as they lack runtime context Gruian (2021). Auto-tuners achieve high, specialized performance but are constrained to a fixed, human-defined parameter space. While these approaches excel at finding optimal configurations within a predefined search space, they cannot fundamentally alter the kernel's algorithm. Our work introduces an agent that reasons about performance bottlenecks to implement novel, structural code changes, moving beyond simple parameter tuning.

#### E.2 MACHINE LEARNING FOR CODE OPTIMIZATION

Machine learning (ML) has emerged as a powerful tool to transcend the limitations of hand-crafted heuristics. Early work focused on using ML to make better decisions within existing compiler and tuning frameworks. Systems like TVM employ a learned cost model to predict the performance of kernel variants, guiding the search process and avoiding exhaustive empirical evaluation Chen et al. (2018). More recent efforts have integrated ML directly into production compilers. Google's MLGO framework uses reinforcement learning (RL) to train policies for classic compiler optimizations like function inlining and register allocation, demonstrating significant improvements in code size and performance over decades-old, manually-tuned heuristics in LLVM Trofin et al. (2021); Marinov & Agarwal (2024). These models can learn from massive code corpora and discover complex feature interactions that are opaque to human experts Cummins et al. (2017). Further, Graph Neural Networks (GNNs) have shown promise in learning from the inherent graph structure of programs to inform decisions on complex combinatorial problems like parallelism mapping and scheduling Allamanis et al. (2018); Qian et al. (2024).

A more profound application of ML has been in algorithmic discovery. DeepMind's AlphaTensor framed the search for faster matrix multiplication algorithms as a single-player game, using a deep RL agent based on AlphaZero to navigate the enormous search space of tensor decompositions Fawzi et al. (2022). This approach successfully discovered novel, provably correct algorithms that outperform human-derived state-of-the-art methods, including improving upon Strassen's algorithm for  $4 \times 4$  matrices for the first time in over 50 years Fawzi et al. (2022); DeepMind (2022). This work marked a critical shift from using ML to *configure* existing optimization strategies to using it to *invent* new ones from first principles. However, AlphaTensor operated in a clean, formal mathematical domain. Translating this power to the messy, syntactic, and hardware-constrained domain of GPU kernel programming presents a distinct challenge. Our work addresses this by employing

```
810
          ## Instruction
811
        2
812
        3
          - Optimize the architecture named Model with custom CUDA operators
813
          - Give explicit and actonable advice to improve the efficiency, in
814
           terms of the GPU wall-clock time, of the architecture named Model
815
816
        5
817
          - Give ONE advice of the top priority! Don't over-request.
        6
818
          - Include necessary details such as how to change pointers and
          indices or how to achieve shared memory tiling so that your advice
819
           can be correctly implemented.
820
821
          - After your advice, modify and return the given source code in
822
          the following way:
823
              - Identify the code block whose efficiency can be improved (
          that is where your advice should be implemented)
824
                and mark it with comments '<<<IMPROVE BEGIN>>>' at the
        11
825
          beginning and '<<<IMPROVE END>>>' at the end
826
              - The markers '<<<IMPROVE BEGIN>>>' and '<<<IMPROVE END>>>'
          should be valid comments for the marked coding language. For
          example, when marking source code of custom kernels, you need to
828
          use comments for the C++ language as ^{\prime}// <<<IMPROVE BEGIN>>>^{\prime} and
829
          '// <<<IMPROVE END>>>'; when marking source code of Python, you
830
          should use '## <<<IMPROVE BEGIN>>>' and '## <<<IMPROVE END>>>'
831
              - Add your advice as comments at the identified code block to
        13
832
          help the following agent's implementation
833
               - There will be another agent focusing on improving the
        14
          efficiency of the identified code block.
834
        15
             - Return the complete code block with the identified code
835
          block as its subpart.
836
        16
837
          - You should consider using CUDA's existing highly optimized
838
          kernels and operations whenever appropriate. Try building on these
           optimized blocks and further improve it with your custom kernels.
839
          - When presented with multiple prior attempts, you should consider
840
           exploration of more diverse optimization strategies.
841
          - Pay careful attention to the implementation agent's capability
          demonstrated from the historical implementations.
843
          - Adjust your advice accordingly to ensure that it can
          successfully implement.
844
845
```

Figure 6: Instruction for the Plan Agent.

an agent that operates directly on source code, navigating the complexities of syntax, compilation, and hardware-specific performance characteristics.

#### E.3 LLM-POWERED AUTONOMOUS AGENTS

The capabilities of LLMs have given rise to a new paradigm of autonomous agents. An LLM agent uses a core LLM as its "brain" or controller, augmented with capabilities for planning, memory, and tool use to perform complex tasks autonomously Weng (2023); Fowler (2025). The key distinction from simple LLM prompting is the agent's ability to decompose a high-level goal into a sequence of manageable subtasks, execute them iteratively, and use reflection to gauge progress and self-correct. This agentic workflow involves the LLM interacting with an external environment through a set of tools, such as a code interpreter or a web search API, to gather information and perform actions. While this paradigm is powerful for general problem-solving, its application to specialized domains like software engineering requires tailored tools and reasoning processes. Our work specializes this agentic concept for the domain of performance optimization, which presents unique challenges not found in general-purpose agent tasks.

```
864
          ## Instruction
865
        2
866
          Optimize the architecture named Model with custom CUDA operators!
867
          - Think about the given advice from human experts and implement
868
          the ones that you believe are correct and you are confident
869
          implementing.
870
          - Only focus on the code block marked with '<<<IMPROVE BEGIN>>>'
871
          and '<<<IMPROVE END>>>'
           - Write custom cuda kernel to replace the pytorch operators within
872
           the marked code block to improve its efficiency
873
          - Name your optimized output architecture ModelNew.
874
          - Output the new code in codeblocks.
875
          - Explain your implementation and how you follow the advice.
876
          - Using the given tool to return your final structured answer.
877
        12
        13 Please generate real code, NOT pseudocode, make sure the code
878
          compiles and is fully functional.
879
        14
          NO testing code!
```

Figure 7: Instruction for the Code Agent.

```
## Instruction
1
2
3
 Fix the issues of your implementation named ModelNew, which should
  improve efficiency of the source model named Model.
 - ModelNew and Model should have the same functionality, that is,
5
 the same input-output mapping.
 - The given architecture ModelNew ieither does not compile, or has
  run-time error, or has different functionality to the source
 Model.
 - Use the given observations to infer bugs and then fix them.
 - Explain how the bugs happen and how you fix it.
 - Return the fixed bug-free code and name your optimized output
 architecture ModelNew.
```

Figure 8: Instruction for the Debug Agent.

#### E.4 LLMs for Code Optimization and Generation

The advent of powerful LLMs has opened a new frontier in performance engineering. Instead of being confined to specific heuristic decisions or mathematical domains, LLMs can directly generate and reason about source code. Initial work has demonstrated that while one-shot generation is often suboptimal, LLMs excel within an iterative refinement loop where they receive feedback from the execution environment. Frameworks like LangProp Ishida et al. (2024) and conceptual systems like the "GPU Kernel Scientist"? formalize this process: an LLM generates code, the code is compiled and benchmarked, and the resulting performance data and errors are fed back to the LLM to guide the next iteration. To systematically evaluate these capabilities, the KernelBench benchmark was introduced, which tasks an LLM with converting a PyTorch reference implementation into a correct and performant CUDA kernel Arora et al. (2025); NVIDIA Developer (2024a). However, these systems represent a form of orchestrated refinement, where the high-level optimization strategy is still encoded in a human-written script.

To grant LLMs greater autonomy, the agentic paradigm has been adapted specifically for software engineering. The success of systems like SWE-agent, which autonomously resolves complex bugs in large GitHub repositories, has demonstrated the viability of this approach Yang et al. (2024). SWE-agent equips an LLM with a specialized Agent-Computer Interface (ACI) containing tools for

```
918
           ## Example Architecture Source Code
919
           import torch
920
           import torch.nn as nn
921
           import torch.nn.functional as F
922
923
         7
           class Model(nn.Module):
924
         8
               def __init__(self) -> None:
925
         9
                    super().__init__()
926
        10
        11
               def forward(self, a, b):
927
        12
                    return a + b
928
        13
929
        14
930
           def get_inputs():
        15
931
               # randomly generate input tensors based on the model
           architecture
932
               a = torch.randn(1, 128).cuda()
        17
933
               b = torch.randn(1, 128).cuda()
        18
934
               return [a, b]
        19
935
        20
936
        21
           def get_init_inputs():
937
        22
               # randomly generate tensors required for initialization based
938
           on the model architecture
939
               return []
        24
940
```

**Figure 9:** In-context Example Architecture for All Agents.

file navigation, editing, and test execution, enabling it to perform long-horizon tasks far beyond the scope of simple code generation Yang et al. (2024); Jimenez et al. (2024). While these agents are a significant step towards autonomous software engineering, their focus has primarily been on functional correctness, such as bug fixing. Our work extends this agentic software engineering paradigm to the non-functional, performance-oriented domain of GPU kernel optimization. We introduce an agent that not only interacts with a codebase but also with hardware profiling tools, allowing it to autonomously diagnose performance bottlenecks, form hypotheses, and conduct experiments to iteratively improve kernel efficiency, thus acting as a true autonomous performance engineer.

```
972
973
           ## Example New Architecture Source Code
974
         2
975
         3
           import torch
976
         4
           import torch.nn as nn
           import torch.nn.functional as F
977
         5
           from torch.utils.cpp_extension import load_inline
         6
978
979
         8
           # Define the custom CUDA kernel for element-wise addition
980
           elementwise_add_source = """
         9
981
         10 #include <torch/extension.h>
982
        11 #include <cuda_runtime.h>
983
           __global__ void elementwise_add_kernel(const float* a, const float
         13
984
           * b, float* out, int size) {
985
               int idx = blockIdx.x * blockDim.x + threadIdx.x;
         14
986
         15
               if (idx < size) {
                    out[idx] = a[idx] + b[idx];
         16
987
         17
988
        18 }
989
         19
990
           torch::Tensor elementwise_add_cuda(torch::Tensor a, torch::Tensor
991
992
               auto size = a.numel();
        21
               auto out = torch::zeros_like(a);
        22
993
        23
994
               const int block_size = 256;
        24
995
               const int num_blocks = (size + block_size - 1) / block_size;
        25
996
        26
               elementwise_add_kernel << num_blocks, block_size >>> (a.data_ptr <
997
        27
           float>(), b.data_ptr<float>(), out.data_ptr<float>(), size);
998
        28
999
               return out;
        29
1000
        30
           11 11 11
1001
        31
1002
        32
        33 elementwise_add_cpp_source = (
1003
               "torch::Tensor elementwise_add_cuda(torch::Tensor a, torch::
        34
           Tensor b);"
1005
        35
        36
           # Compile the inline CUDA code for element-wise addition
1007
        37
           elementwise_add = load_inline(
        38
1008
               name="elementwise_add",
        39
1009
        40
               cpp_sources=elementwise_add_cpp_source,
1010
        41
               cuda_sources=elementwise_add_source,
               functions=["elementwise_add_cuda"],
1011
        42
               verbose=True,
         43
1012
               extra_cflags=[""],
        44
1013
               extra_ldflags=[""],
        45
1014
           )
        46
1015
        47
1016
        48
1017
        49
           class ModelNew(nn.Module):
               def __init__(self) -> None:
        50
1018
                    super().__init__()
        51
1019
                    self.elementwise_add = elementwise_add
        52
1020
        53
1021
               def forward(self, a, b):
        54
                    return self.elementwise_add.elementwise_add_cuda(a, b)
         55
1022
```

Figure 10: In-context Optimized Example Architecture for All Agents.

```
1030
1031
1032
1033
1034
1035
1036
1037
1038
           import torch
           import torch.nn as nn
        2
1039
        3
1040
        4
           class Model(nn.Module):
1041
        5
1042
               Simple model that performs a LogSoftmax activation.
        6
1043
                11 11 11
        7
               def __init__(self, dim: int = 1):
1044
                    super(Model, self).__init__()
        9
1045
        10
                    self.dim = dim
1046
        11
1047
               def forward(self, x: torch.Tensor) -> torch.Tensor:
        12
1048
        13
                    Applies LogSoftmax activation to the input tensor.
1049
        14
        15
1050
                    Args:
       16
1051
                        x (torch.Tensor): Input tensor of shape (batch_size,
1052
                             dim).
1053
       18
                    Returns:
       19
1054
                        torch. Tensor: Output tensor with LogSoftmax applied,
       20
1055
                             same shape as input.
1056
       21
1057
                    return torch.log_softmax(x, dim=self.dim)
       22
1058
       23
           batch\_size = 4096
1059
           dim = 393216
       25
1060
       26
1061
           def get_inputs():
       27
1062
               x = torch.rand(batch_size, dim)
       28
1063
               return [x]
       29
1064
       30
           def get_init_inputs():
       31
1065
               return [] # No special initialization inputs needed
1066
1067
```

Figure 11: Example KernelBench Level 1 Task.

```
1081
1082
1083
1084
1085
1086
1087
1088
1089
           import torch
1090
           import torch.nn as nn
1091
1092
           class Model(nn.Module):
        5
1093
               Model that performs a matrix multiplication, division,
        6
1094
                  summation, and scaling.
1095
1096
               def __init__(self, input_size, hidden_size, scaling_factor):
        8
                   super(Model, self).__init__()
1097
                   self.weight = nn.Parameter(torch.randn(hidden_size,
1098
                       input_size))
1099
                   self.scaling_factor = scaling_factor
       11
1100
       12
1101
               def forward(self, x):
       13
1102
       14
1103
        15
                   Args:
                        x (torch.Tensor): Input tensor of shape (batch_size,
       16
1104
                            input_size).
1105
                   Returns:
1106
                        torch. Tensor: Output tensor of shape (batch_size,
1107
                            hidden_size).
                    " " "
       19
1108
                   x = torch.matmul(x, self.weight.T) # Gemm
       20
1109
                   x = x / 2 \# Divide
1110
                   x = torch.sum(x, dim=1, keepdim=True) # Sum
       22
1111
                   x = x * self.scaling_factor # Scaling
1112
       24
                   return x
1113
1114
          batch_size
                        = 1024
       27
1115
                        = 8192
           input_size
       28
          hidden_size = 8192
1116
       29
          scaling_factor = 1.5
1117
       31
1118
          def get_inputs():
       32
1119
               return [torch.rand(batch_size, input_size)]
       33
1120
1121
          def get_init_inputs():
       35
               return [input_size, hidden_size, scaling_factor]
1122
1123
```

Figure 12: Example KernelBench Level 2 Task.

```
1134
           import torch
1135
        2
           import torch.nn as nn
1136
           import torch.nn.functional as F
           import math
1137
1138
           class NewGELU(nn.Module):
        6
1139
1140
               Implementation of the GELU activation function currently in
1141
                   Google BERT repo (identical to OpenAI GPT).
               Reference: Gaussian Error Linear Units (GELU) paper: https://
1142
        9
                   arxiv.org/abs/1606.08415
1143
        10
1144
               def
                     _init__(self):
       11
1145
                   super(NewGELU, self).__init__()
       12
1146
       13
1147
               def forward(self, x):
                   return 0.5 * x * (1.0 + torch.tanh(math.sqrt(2.0 / math.pi
       15
1148
                       ) * (x + 0.044715 * torch.pow(x, 3.0)))
1149
       16
1150
           class CausalSelfAttention(nn.Module):
       17
1151
       18
1152
               A vanilla multi-head masked self-attention layer with a
       19
                   projection at the end.
1153
               It is possible to use torch.nn.MultiheadAttention here but I
       20
1154
                   am including an
1155
               explicit implementation here to show that there is nothing too
       21
1156
                    scary here.
1157
       22
       23
1158
               def __init__(self, n_embd, n_head, attn_pdrop, resid_pdrop,
1159
                   max_seqlen):
1160
       25
                   super().__init__()
1161
                   assert n_embd % n_head == 0
1162
                   # key, query, value projections for all heads, but in a
                       batch
1163
                   self.c_attn = nn.Linear(n_embd, 3 * n_embd)
       28
1164
       29
                   # output projection
1165
                   self.c_proj = nn.Linear(n_embd, n_embd)
1166
       31
                   # regularization
                   self.attn_dropout = nn.Dropout(attn_pdrop)
1167
       32
                   self.resid_dropout = nn.Dropout(resid_pdrop)
1168
       33
                   # causal mask to ensure that attention is only applied to
1169
                       the left in the input sequence
1170
                   self.register_buffer("bias", torch.tril(torch.ones(
       35
1171
                       max_seqlen, max_seqlen))
1172
                                                   .view(1, 1, max_seqlen,
       36
                                                      max_seqlen))
1173
       37
                   self.n_head = n_head
1174
                   self.n\_embd = n\_embd
       38
1175
                   [...skipped]
       39
1176
          class Model(nn.Module):
1177
       41
               """ an unassuming Transformer block """
       42
1178
       43
1179
       44
               [...skipped]
1180
       45
1181
               def forward(self, x):
       46
                   x = x + self.attn(self.ln_1(x))
1182
       47
                   x = x + self.mlpf(self.ln_2(x))
       48
1183
                   return x
       49
1184
       50
1185
           [...skipped]
1186
```

**Figure 13:** Example KernelBench Level 3 Task.