

# CUDA-L1: IMPROVING CUDA OPTIMIZATION VIA CONTRASTIVE REINFORCEMENT LEARNING


Xiaoya Li, Xiaofei Sun, Albert Wang, Jiwei Li and Chris Shum

DeepReinforce Team

 [github.com/deepreinforce-ai/CUDA-L1](https://github.com/deepreinforce-ai/CUDA-L1)

## ABSTRACT

The exponential growth in demand for GPU computing resources has created an urgent need for automated CUDA optimization strategies. While recent advances in LLMs show promise for code generation, current state-of-the-art models achieve low success rates in improving CUDA speed. In this paper, we introduce CUDA-L1, an automated reinforcement learning (RL) framework for CUDA optimization that employs a novel contrastive RL algorithm.


CUDA-L1 achieves significant performance improvements on the CUDA optimization task: trained on NVIDIA A100, it delivers an average speedup of  $\times 3.12$  with a median speedup of  $\times 1.42$  against default baselines over across all 250 CUDA kernels of KernelBench, with peak speedups reaching  $\times 120$ . In addition to the default baseline provided by KernelBench, CUDA-L1 demonstrates  $\times 2.77$  over Torch Compile,  $\times 2.88$  over Torch Compile with reduce overhead, and  $\times 2.81$  over CUDA Graph implementations. Furthermore, the model also demonstrates portability across GPU architectures. CUDA-L1 opens possibilities for automated optimization of CUDA operations, and holds promise to substantially promote GPU efficiency and alleviate the rising pressure on GPU computing resources. 

## 1 INTRODUCTION

The exponential growth in demand for GPU computing resources, driven primarily by the rapid advancement and deployment of Large Language Models (LLMs), has created an urgent need for highly efficient CUDA optimization strategies. Traditionally, CUDA optimization has been a highly manual and time-intensive process, where skilled engineers must meticulously analyze memory access patterns, experiment with different thread block configurations, and iteratively profile their code through extensive trial-and-error cycles.

Recent advances in LLMs (Team et al., 2023; Shengyu et al., 2023; Team et al., 2024; Grattafiori et al., 2024; Yang et al., 2025; Hurst et al., 2024; Jiang et al., 2024; Liu et al., 2024a; OLMo et al., 2024), especially those powered with RL (Jaech et al., 2024; Guo et al., 2025; Wang et al., 2024; Muennighoff et al., 2025), have demonstrated remarkable capabilities in code generation and algorithm design. Despite the promise, current performance remains limited. State-of-the-art LLM models such as DeepSeek-R1 (Guo et al., 2025) and OpenAI-o1 (Jaech et al., 2024) achieve low success rates in generating optimized CUDA code (only approximately 15% on KernelBench (Ouyang et al., 2025)). To address these limitations and unlock the potential of LLMs for automated CUDA optimization, in this work, we propose CUDA-L1, an LLM framework powered by contrastive reinforcement learning for CUDA optimization. CUDA-L1 is a pipelined framework, the core of which is a newly-designed contrastive RL framework.

Different from previous RL models (Williams, 1992; Shao et al., 2024; Schulman et al., 2017), contrastive RL performs comparative analysis of previously generated CUDA variants alongside their execution performance, enabling the model to improve through distinguishing between effective and ineffective optimization strategies. Contrastive-RL simultaneously optimizes the foundation

 Email: {xiaoya\_li, xiaofei\_sun, albert\_wang, jiwei\_li, chris\_shum}@deep-reinforce.com

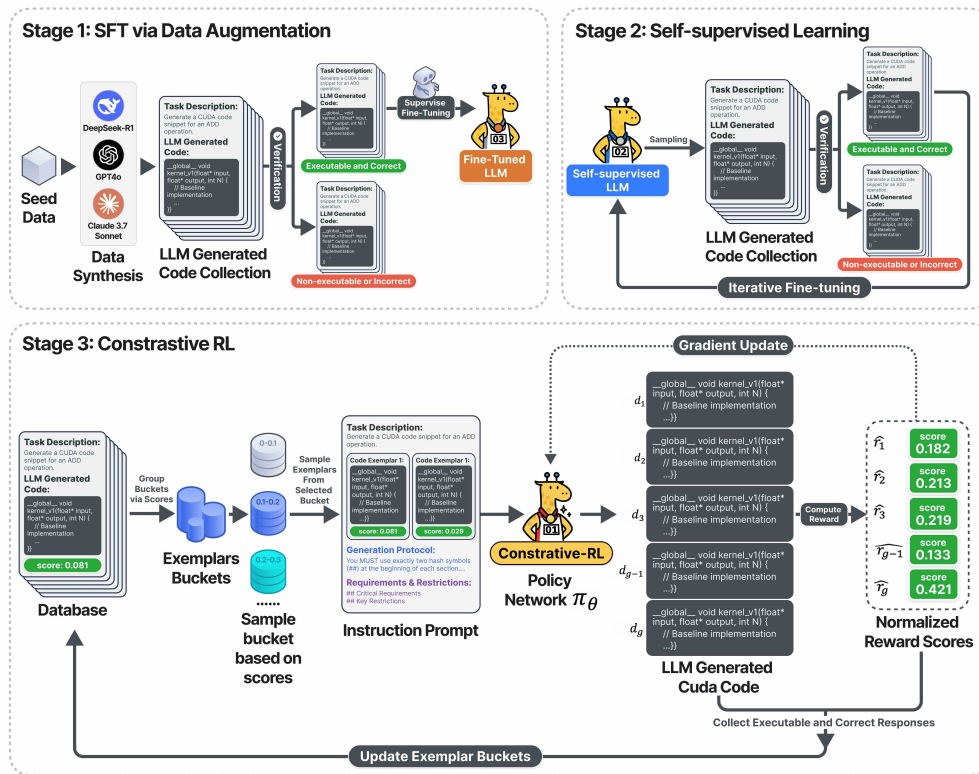


Figure 1: Overview of the CUDA-L1 training pipeline.

model through gradient-based parameter updates while fulfilling the maximum potential from the current model through contrastive analysis from high-performance CUDA variants, creating a co-evolutionary dynamic that drives superior CUDA optimization performance.

CUDA-L1 delivers significant improvements on the CUDA optimization task: trained on NVIDIA A100, it achieves an **average** speedup of  $\times 3.12$  (**median**  $\times 1.42$ ) over the default baseline across all 250 KernelBench CUDA kernels, with maximum speedups reaching  $\times 120$ . In addition, CUDA-L1 demonstrates  $\times 2.77$  over Torch Compile,  $\times 2.88$  over Torch Compile with reduce overhead,  $\times 2.81$  over CUDA Graph implementations. Furthermore, the CUDA codes optimized specifically for A100 demonstrate strong portability across GPU architectures, with similar optimization patterns observed across different baseline configurations: achieving average speedups of  $\times 3.85$  (median  $\times 1.32$ ) on H100,  $\times 3.13$  (median  $\times 1.31$ ) on L40,  $\times 2.51$  (median  $\times 1.18$ ) on RTX 3090, and  $\times 2.38$  (median  $\times 1.34$ ) on H20. Similar performance improvements over Torch Compile, Torch Compile with reduce overhead, CUDA Graph are consistently observed across all GPU types.

CUDA-L1 reveals a remarkable capability of RL in autonomous learning for CUDA optimization:

1. Even starting with a foundation model with poor CUDA optimization ability, by using code speedups as RL rewards and proper contrastive RL training techniques, we can still train an RL system capable of generating CUDA optimization codes with significant speedups.
2. Without human prior knowledge, RL systems can independently discover CUDA optimization techniques, learn to combine them strategically, and more importantly, extend the acquired CUDA reasoning abilities to unseen kernels. This capability unlocks the potential for a variety of automatic CUDA optimization tasks, e.g., kernel parameter tuning, memory access pattern optimization, and different hardware adaptations, offering substantial promises to enhance GPU utilization.

Another contribution of this work is the enrichment of the KernelBench dataset with CUDA Graph implementations. Please refer to supplementary materials. We release these implementations to the community, providing substantially stronger baselines for performance comparison.

## 2 CUDA-L1

### 2.1 OVERVIEW

Existing large language models (Guo et al., 2025; Yang et al., 2025; Grattafiori et al., 2024) demonstrate significant limitations in generating executable and correct CUDA code with speedup, as reported in prior research (Ouyang et al., 2025). This deficiency likely stems from the insufficient representation of CUDA code in the training datasets of these models. To address this fundamental gap, we introduce a three-stage pipelined training strategy for CUDA-L1, i.e., Supervised fine-tuning via data augmentation, Self-supervised learning and Contrastive reinforcement learning, aiming to progressively enhance the model’s CUDA programming capabilities:

Before we delve into the details of each stage, we provide key definitions adopted throughout the rest of this paper:

1. **Executability:** A CUDA code is executable if it successfully compiles, launches, and executes to completion within  $1000\times$  the runtime of the reference implementation. Code exceeding this runtime threshold is considered unexecutable.<sup>1</sup>
2. **Correctness:** A CUDA code is correct if it produces equivalent outputs to the reference implementation across 1000 random test inputs.<sup>2</sup>
3. **Success:** A CUDA code is successful if it is both executable and correct.

### 2.2 SFT VIA DATA AUGMENTATION

In the SFT stage, we collect a dataset by using existing LLMs to generate CUDA code snippets and selecting successful one. This dataset is directly used to fine-tune the model. Throughout this paper, we use deepseek-v3-671B (Liu et al., 2024a) as the model backbone. Please refer to the details of data collection in Appendix A.1. The collected dataset  $D$  is used to finetune the foundation model. The instruction to the model is the same as the prompt for dataset generation, where the reference code  $q_i$  is included in the instruction and the model is asked to generate an improved version. The model is trained to predict each token in  $d_{i,j}$  given the instruction.

### 2.3 SELF-SUPERVISED LEARNING

Now we are presented with the finetuned model after the SFT stage, where the model can potentially generate better CUDA code with higher success rates than the original model without finetuning. We wish to further improve the model’s ability to generate successful CUDA code by exposing it to more code snippets generated by itself.

We achieve this iteratively by sampling CUDA code from the model, evaluating it for executability and correctness, removing the unsuccessful trials and keeping the successful ones. Successful ones are batched and used to update the model parameters. Using the updated model, we repeat the process: generating code, evaluating it, and retraining the model. It is worth noting that during the self-supervised learning stage, we focus exclusively on the executability and correctness of the generated code, without considering speed as a metric. This design choice reflects our primary objective of establishing reliable code generation before optimizing for performance.

### 2.4 CONTRASTIVE REINFORCEMENT LEARNING

Now we have a model capable of generating successful CUDA code at a reasonable success rate. Next, we aim to optimize for execution speed.

One straightforward approach is to apply existing reinforcement learning algorithms such as REINFORCE (Williams, 1992), GRPO (Shao et al., 2024), or PPO (Schulman et al., 2017). In this approach, we would ask the model to first perform chain-of-thought reasoning (Wei et al., 2022), then generate code, evaluate it, and use the evaluation score to update the model parameters. However, our experiments reveal that these methods perform poorly in this task. The issue is as follows: standard RL algorithms compute a scalar reward for each generated CUDA code sample. During training, this reward undergoes algorithm-specific processing (e.g., baseline subtraction in REINFORCE, advantage normalization in GRPO, importance sampling in PPO). The processed reward then serves as a loss weighting term for gradient updates, increasing the likelihood of high-reward

<sup>1</sup>This threshold is reasonable since code with  $1000\times$  slower performance contradicts our speedup optimization goals.

<sup>2</sup>Prior work uses only 5 random inputs, which we found insufficient for robust validation.

sequences while decreasing the likelihood of low-reward sequences. Critically, in this paradigm, the reward signal is used exclusively for parameter updates and is never provided as input to the LLM. Consequently, the LLM cannot directly reason about performance trade-offs during code generation. To address this limitation, we propose incorporating reward information directly into the reasoning process by embedding performance feedback within the input prompt. Specifically, we present the model with multiple code variants alongside their corresponding speedup scores. Rather than simply generating code, the LLM is trained to first conduct comparative analysis of why certain implementations achieve superior performance, then synthesize improved solutions based on these insights. Each generated code sample undergoes evaluation to obtain a performance score, which serves dual purposes in our training framework: the primary purpose is to the score acts as a reward signal for gradient-based parameter optimization, updating model weights. The score functions as a reward signal for gradient-based parameter optimization, directly updating the model weights; the other is to construct prompt for future training stages. The scored code sample becomes part of the exemplar set for subsequent training iterations, enriching the contrastive learning dataset.

This dual-utilization strategy enables iterative optimization across two complementary dimensions:

**Foundation Model Enhancement** Parameter updates progressively improve the model’s fundamental understanding and capabilities for CUDA optimization tasks, expanding its representational capacity.

**Fixed-Parameter Solution Optimization** The contrastive approach seeks to extract the maximum potential from the current model’s parameters by leveraging comparative analysis of high-quality exemplars.

These two optimization processes operate synergistically: enhanced foundation models enable more accurate contrastive reasoning, while improved reasoning strategies provide higher-quality training signals for parameter updates of foundation models. This co-evolutionary dynamic drives convergence toward optimal performance. We term this approach contrastive reinforcement learning (contrastive-RL for short).

#### 2.4.1 PROMPT CONSTRUCTION

Here we describe the construction of prompts provided to the LLM. The prompt provided to the LLM during Contrastive-RL training comprises the following structured components:

- I) **Task Description:** A detailed description of the computational problem, including input/output specifications, performance requirements, and optimization objectives.
- II) **Previous Cuda Codes with Scores:** Previously generated CUDA implementations paired with their corresponding performance scores (e.g., execution time, throughput, memory efficiency), providing concrete examples of varying solution quality.
- III) **Generation Protocol:** Explicit instructions defining the required output format and components.
- IV) **Requirements and Restrictions:** Requirements and restrictions to prevent reward hacking in RL.

The model’s response must contain the following three structured components:

- I) **Performance Analysis:** A comparative analysis identifying which previous kernel implementations achieved superior performance scores and the underlying algorithmic or implementation factors responsible for success.
- II) **Algorithm Design:** A high-level description of the proposed optimization strategy, outlining the key techniques to be applied, presented as numbered points in natural language.
- III) **Code Implementation:** The complete CUDA kernel implementation incorporating optimizations.

A detailed demonstration for the prompt is shown in Table 9.

#### 2.4.2 CONTRASTIVE EXEMPLAR SELECTION

The selection of code exemplars for prompt construction is critical, as core of Contrastive-RL is to perform meaningful comparative analysis. The selection strategy needs to address the following two key requirements: first, for achieving competitive performance, the exemplar set should include higher-performing implementations to guide the model toward competitive CUDA codes, avoiding local minima that result from analyzing and comparing inferior codes; second, to ensure

performance diversity, the selected codes must exhibit substantial performance differences to enable effective contrastive analysis.

We employ a sampling strategy akin to that adopted by evolutionary LLM models: Let  $N$  denote the number of code exemplars included in each prompt (set to  $N = 2$  in our experiments). During RL training, we maintain a performance-indexed database of all successful code samples generated during RL training. Codes are organized into performance buckets  $B_k$  based on discretized score intervals, where bucket  $B_i$  contains codes with scores in range  $[s_k, s_k + \Delta s)$ .

We first sample  $N$  distinct buckets according to a temperature-scaled softmax distribution:

$$P(B_i) = \frac{\exp((\bar{s}_i - \mu_s)/\tau)}{\sum_j \exp((\bar{s}_j - \mu_s)/\tau)} \quad (1)$$

where  $\bar{s}_i$  denotes the aggregate score of bucket  $B_i$ , computed as the mean of its constituent code scores,  $\mu_s = \text{mean}(\{\bar{s}_j\}_{j=1}^M)$  represents the global mean of all bucket scores, and  $\tau$  is the temperature parameter governing the exploration-exploitation tradeoff. The sampling strategy in Equation 1 differs from conventional temperature sampling in evolutionary LLM approaches through a modification: the deduction of  $\mu_s$  stabilizes the distribution by centering scores around zero, which prevents absolute score magnitudes from dominating the selection.

From each selected bucket  $B_i$ , we uniformly sample one representative code to construct the final prompt set. This approach satisfies both design criteria: Regarding competitive Performance, score-weighted bucket sampling biases selection toward higher-performing implementations, ensuring the exemplar set contains competitive solutions; Regarding performance Diversity, enforcing selection from  $N$  distinct buckets ensures sufficient performance variance for effective contrastive analysis.

A more sophisticated alternative is to use an island-based approach for exemplar selection. However, we find no significant difference in performance between our bucket-based method and the island-based approach. Given this, we opt for the simpler bucket-based strategy.

### 2.4.3 REWARD

In this subsection, we detail the computation of the execution time-based reward function, which serves dual purposes: (1) guiding parameter updates in reinforcement learning and (2) constructing effective prompts. Given a reference CUDA implementation  $q_i$  from PyTorch with successful execution time  $t_{q_i}$ , and a generated code candidate  $d$  with execution time  $t_d$ , we define the single-run speedup score as:

$$r_{\text{single-run}}(d) = \frac{t_{q_i}}{t_d} \quad (2)$$

More details for denoising the rewards are shown in Appendix A.6.

For RL training, we adopt the Group Relative Policy Optimization (GRPO) strategy (Shao et al., 2024). Please refer to Appendix A.3 for details.

## 3 EXPERIMENTS AND ANALYSIS

### 3.1 KERNELBENCH AND EVALUATION

Our evaluation is conducted on the KernelBench dataset (Ouyang et al., 2025). The KernelBench Dataset contains a collection of 250 PyTorch workloads designed to evaluate language models’ ability to generate efficient GPU kernels. The dataset is structured across three hierarchical levels based on computational complexity: Level 1 contains 100 tasks with single primitive operations (such as convolutions, matrix multiplications, activations, and normalizations), Level 2 includes 100 tasks with operator sequences that can benefit from fusion optimizations (combining multiple operations like convolution + ReLU + bias), and Level 3 comprises 50 full ML architectures sourced from popular repositories including PyTorch, Hugging Face Transformers, and PyTorch Image Models (featuring models like AlexNet and MiniGPT). Each task in the dataset provides a reference PyTorch implementation with standardized input/output specifications, enabling automated evaluation of both functional correctness and performance through wall-clock timing comparisons. The dataset represents real-world engineering challenges where successful kernel optimization directly translates to practical performance improvements. Throughout this paper, we use KernelBench as the evaluation benchmark. KernelBench is recognized as a challenging benchmark in the community (Ouyang et al., 2025), with even the best current LLMs improving fewer than 20% of tasks.

Configuration	Method	Mean	Max	75%	50%	25%	Success <sup>↑</sup> # out of total	Speedup <sup>↑</sup> >1.01x out of total
<i>Default</i>	All	3.12	120	2.25	1.42	1.17	249/250	226/250
	Level 1	2.78	65.8	1.75	1.28	1.12	99/100	80/100
	Level 2	3.55	120	2.05	1.39	1.20	100/100	98/100
	Level 3	2.96	24.9	2.60	1.94	1.42	50/50	48/50
<i>Torch Compile</i>	All	2.77	69.0	2.55	1.72	1.14	249/250	203/250
	Level 1	3.04	59.7	2.71	1.99	1.41	99/100	89/100
	Level 2	2.91	69.0	1.99	1.55	1.10	100/100	78/100
	Level 3	1.98	8.57	2.28	1.68	1.00	50/50	36/50
<i>Torch Compile RO</i>	All	2.88	80.1	2.48	1.67	1.13	249/250	200/250
	Level 1	3.38	55.3	3.02	2.29	1.61	99/100	90/100
	Level 2	3.00	80.1	2.06	1.54	1.10	100/100	79/100
	Level 3	1.62	8.67	1.76	1.13	0.991	50/50	31/50
<i>CUDA Graph</i>	All	2.81	97.9	1.83	1.20	0.954	249/250	147/229
	Level 1	3.18	59.6	2.09	1.38	1.04	99/100	68/88
	Level 2	2.84	97.9	1.55	1.08	0.932	100/100	53/94
	Level 3	2.06	24.6	1.74	1.08	0.887	50/50	26/47

Table 1: Performance comparison across different configurations on KernelBench on A100. RO = Reduce Overhead. Success and Speedup indicate the number of successful benchmarks out of the total for each level. Note that for CUDA Graph, the total benchmark count differs from the dataset/data-subset size, as some original reference code in KernelBench cannot be successfully transformed into the corresponding CUDA Graph implementations.

For each task with reference implementation  $q$ , we evaluate the performance of a generated CUDA code  $d$  using a similar protocol to training: We execute both  $q$  and  $d$  in randomized order within a fixed time budget of 20 minutes per task. The number of execution rounds varies across tasks due to differences in individual runtimes. The final evaluation score for  $d$  is computed as the average speedup ratio across all execution rounds within the allocated time window. Unsuccessful implementations receive a score of zero. The metrics we report include speedup statistics (mean, maximum, and 75th, 50th, and 25th percentiles), success rate, and percentage of improvements.

### 3.2 COMPARISON SETUPS

To perform a comprehensive evaluation on the generated code, we perform the following comparisons:

**I) Default** This compares the CUDA-L1 generated code with the reference code by KernelBench.

**II) Torch Compile** This compares the CUDA-L1 generated code with the reference code enhanced by torch.compile with default settings. Torch.compile applies graph-level optimizations including operator fusion, memory planning, and kernel selection to accelerate PyTorch models through just-in-time compilation.

**III) Torch Compile Reduce Overhead** This compares the CUDA-L1 generated code with the reference code enhanced by torch.compile with reduce-overhead mode enabled. This mode minimizes the compilation overhead by caching compiled graphs more aggressively and reducing recompilation frequency, making it particularly suitable for inference workloads with static shapes.

**IV) CUDA Graph** Since KernelBench does not provide official CUDA Graph implementations, we employ Claude 4 to generate CUDA Graph-augmented code for each reference implementation. CUDA Graphs capture a series of CUDA kernels and their dependencies into a single graph structure that can be launched with minimal CPU overhead, eliminating the need for repeated kernel launch commands and significantly reducing CPU-GPU synchronization costs. Specifically, we provide Claude 4 with the reference code and request the addition of CUDA Graph optimizations. The generated output is then evaluated for correctness. If the code fails validation, we iterate by providing Claude 4 with both the original reference code and the previous erroneous outputs, requesting a corrected version. This iterative process continues for up to 10 attempts until the generated code passes all correctness checks. We release the CUDA Graph codes for KernelBench to the community, pro-

Methods	Model	Mean	Max	75%	50%	25%	Success <sup>†</sup> # out of 250	Speedup <sup>†</sup> >1.01 # out of 250
<i>Vanilla</i>	Llama 3.1-405B	0.23	3.14	0.63	0	0	68	5
	DeepSeek-V3	0.34	2.96	0.76	0	0	99	9
	DeepSeek-R1	0.88	14.4	1.00	0.75	0	179	18
	OpenAI-O1	0.73	12.4	1.00	0.55	0	141	14
<i>Evolve</i>	Llama 3.1-405B	1.18	18.4	1.03	1.00	1.00	247	88
	DeepSeek-V3	1.32	52.4	1.32	1.03	1.00	247	113
	DeepSeek-R1	1.41	44.2	1.45	1.17	1.00	248	162
	OpenAI-O1	1.35	63.9	1.38	1.16	1.00	247	158
<i>CUDA-L1</i>	Stage 1	1.14	32.7	1.00	1.00	0.96	240	50
	Stage 1+2	1.36	48.3	1.41	1.09	1.00	247	175
	Stage 1+2+GRPO	2.41	84.6	1.83	1.33	1.11	247	207
	3 stages - random	2.14	64.5	1.62	1.21	1.09	241	186
	- island	<b>3.21</b>	<b>126</b>	2.21	1.40	1.16	<b>249</b>	223
	- bucket	3.12	120	<b>2.25</b>	<b>1.42</b>	<b>1.17</b>	<b>249</b>	<b>226</b>

Table 2: Model performances on KernelBench All Level.

viding researchers and practitioners with ready-to-use optimized implementations that can serve as strong baselines for future performance studies and benchmarking efforts.

### 3.3 MAIN RESULTS ON KERNELBENCH

The experimental results in Table 1 demonstrate CUDA-L1’s optimization effectiveness across different baseline configurations on KernelBench. CUDA-L1 achieves substantial performance improvements over the Default baseline with  $3.12\times$  average speedup and  $120\times$  maximum gains. Against Torch compilation baselines, CUDA-L1 delivers moderate but consistent improvements with  $2.77\text{--}2.88\times$  mean speedup ratios, while demonstrating  $2.81\times$  mean improvement over CUDA Graph baseline with notable  $97.9\times$  maximum gains.

Across difficulty levels, CUDA-L1’s optimization effectiveness varies by task complexity. For Level 1 (single operations), CUDA-L1 achieves moderate improvements ranging from  $2.78\text{--}3.38\times$  over different baselines. Level 2 (operator sequences) shows CUDA-L1’s strongest performance with  $3.55\times$  improvement over Default baseline. Level 3 (complex ML tasks) reveals interesting baseline-dependent effectiveness: CUDA-L1 achieves  $2.96\times$  improvement over Default baseline, but shows reduced effectiveness against Torch compilation baselines (only  $1.62\text{--}1.98\times$  improvements), suggesting these configurations provide stronger baseline performance for complex operations.

### 3.4 BASELINE COMPARISON

We compare the results with the following three groups of baselines:

**Vanilla Foundation Models:** To establish baseline performance benchmarks, we evaluate OpenAI-o1, DeepSeek-R1, DeepSeek-V3, and Llama 3.1-405B Instruct (denoted by **OpenAI-o1-vanilla**, **DeepSeek-R1-vanilla**, **DeepSeek-V3-vanilla** and **Llama 3.1-405B-vanilla**) by prompting each model to optimize the reference CUDA code. The generated CUDA code is directly used for evaluation without further modification. For each task, we repeat this process 5 times and report the best.

**Evolutionary LLM:** We implement evolutionary LLM strategies where, given a set of previous codes, we sample up to 4 high-performing kernels based on evaluation scores. The key difference is that the model only performs contrastive analysis without updating model parameters. We adopt the island strategy for code database construction and sampling, as suggested in Novikov et al. (2025). We conduct experiments on DeepSeek-R1, OpenAI-o1 and Llama 3.1-405B, denoted as **DeepSeek-R1-evolve**, **OpenAI-o1-evolve**, **DeepSeek-V3-evolve** and **Llama 3.1-405B-evolve**.

#### Different combinations of CUDA-L1 components and variants:

- **stage1:** Uses only the outcome from the first stage with supervised fine-tuning applied
- **stage1+2:** Applies only the first two stages without reinforcement learning
- **stage1+2 + GRPO:** Replaces the contrastive RL with a vanilla GRPO strategy, without comparative analysis

Configuration	GPU Device	Mean	Max	75%	50%	25%	Success $\uparrow$	Speedup $\uparrow$
							# out of 250	>1.01x
<i>Default</i>	A100	3.12	120	<b>2.25</b>	<b>1.42</b>	<b>1.17</b>	249	<b>226/250</b>
	3090	2.51	114	1.57	1.18	1.03	242	201/250
	H100	<b>3.85</b>	<b>368</b>	1.76	1.32	1.09	<b>250</b>	218/250
	H20	2.38	63.7	1.81	1.34	1.11	247	<b>226/250</b>
	L40	3.13	182	1.88	1.31	1.08	248	215/250
<i>Torch Compile</i>	A100	2.77	69.0	2.55	1.72	1.14	249	203/250
	3090	2.58	73.2	2.23	1.50	1.00	242	177/250
	H100	2.74	49.7	2.83	1.92	1.11	<b>250</b>	195/250
	H20	<b>2.89</b>	49.4	<b>3.21</b>	<b>2.04</b>	<b>1.19</b>	247	<b>209/250</b>
	L40	2.85	<b>96.9</b>	2.43	1.82	1.13	248	199/250
<i>Torch Compile RO</i>	A100	2.88	80.1	2.48	1.67	<b>1.13</b>	<b>249</b>	<b>200/250</b>
	3090	2.61	72.9	2.29	1.48	1.00	242	172/250
	H100	2.77	61.2	2.78	1.61	1.00	247	187/250
	H20	2.82	52.1	<b>3.18</b>	1.64	1.06	247	192/250
	L40	<b>2.89</b>	<b>90.9</b>	2.54	<b>1.72</b>	1.08	248	193/250
<i>CUDA Graph</i>	A100	2.81	97.9	1.83	1.20	0.954	<b>229</b>	147/229
	3090	3.34	156	<b>1.94</b>	<b>1.28</b>	<b>0.997</b>	206	<b>148/206</b>
	H100	2.23	70.1	1.60	1.04	0.838	222	119/222
	H20	2.20	64.6	1.69	1.09	0.854	<b>229</b>	133/229
	L40	<b>3.98</b>	<b>275</b>	1.83	1.16	0.862	224	137/224

Table 3: Performance comparison across different configurations and GPU devices on KernelBench. RO = Reduce Overhead. Speedup is defined as value exceeding 1.01x.

- **random sampling**: Replaces the bucket sampling strategy with simple random sampling of exemplars
- **island sampling**: Adopts an island-based sampling strategy Novikov et al. (2025), where examples are distributed across different islands, prompts are constructed using exemplars from the same island, and newly generated examples are added to that island. After a fixed number of iterations, examples in half of the inferior islands are eliminated and examples from superior islands are copied to replace them.

Results are shown in Table 2. As observed, all vanilla foundation models perform poorly on this task. Even the top-performing models, such as DeepSeek-R1 and OpenAI-o1, achieve speedups over the reference kernels in fewer than 10% of tasks, while Llama 3.1-405B optimizes only 2.4% of tasks. This confirms that vanilla foundation models cannot be readily applied to CUDA optimization due to their insufficient grasp of CUDA programming principles and optimization techniques.

We observe significant performance improvements introduced by the Evolutionary LLM models compared to vanilla foundation model setups, despite sharing the same parameter sets. All Evolve models achieve speedups in over 70% of tasks, with DeepSeek-R1 reaching 72.4% success rate. This demonstrates that leveraging contrastive analysis, which exploits the model’s general reasoning abilities, is more effective than direct output generation. The superiority of evolutionary LLM over vanilla LLM also provides evidence that contrastive RL should outperform non-contrastive RL approaches like vanilla GRPO, as the relationship between evolutionary and vanilla LLMs parallels that between contrastive and non-contrastive RL methods.

When comparing the different combinations of CUDA-L1 components, we observe a progressive increase in speedup rates from **stage1 (SFT only)** at 22.4% to **stage1+2 (SFT + self-supervised)** at 66%, and further to **stage1+2+GRPO** at 88.4%. This demonstrates the cumulative benefits of each training stage in improving model performance.

All RL-based approaches significantly outperform evolutionary LLM baselines with fixed model parameters, with the best RL methods achieving over 95% speedup rates compared to 72.4% for the best evolutionary approach. This demonstrates the necessity of model parameter updating for achieving optimal performance in CUDA optimization tasks.

### 3.5 GENERALIZATION OF A100-OPTIMIZED KERNELS TO OTHER GPU ARCHITECTURES

Even without being specifically tailored to other GPU architectures, we observe significant performance improvements across all tested GPU types, with mean speedups ranging from  $2.38\times$  to  $3.85\times$ . H100 achieves the highest mean speedup ( $3.85\times$ ) with exceptional maximum gains ( $368\times$ ), while A100 PCIe and L40 demonstrate strong performance with mean speedups of  $3.12\times$  and  $3.13\times$  respectively. L40 shows the second-highest maximum speedup ( $182\times$ ) among all GPUs. The consumer RTX 3090 achieves a competitive mean speedup of  $2.51\times$ , while H20 shows moderate performance with  $2.38\times$  mean speedup. Notably, A100 maintains the highest 75th percentile ( $2.25\times$ ), 50th percentile ( $1.42\times$ ), and 25th percentile ( $1.17\times$ ) values, indicating more consistent optimization performance on the target architecture.

The success rates remain high across all architectures (242-250 out of 250), with H100 achieving perfect success (250/250), validating that CUDA optimization techniques can generalize across different GPU architectures. Speedup achievement rates ( $>1.01\times$ ) vary by architecture, with H20 and A100 showing the highest effectiveness (226 and 226 successful optimizations respectively), while RTX 3090 demonstrates good performance with 201 successful optimizations.

These results demonstrate that while A100-optimized kernels transfer to other GPUs with varying degrees of effectiveness, the optimizations achieve substantial improvements across architectures. H100’s exceptional performance suggests strong compatibility with the optimization techniques, while A100’s consistent percentile performance validates the target architecture optimization. The varying maximum speedups ( $63.7\times$  to  $368\times$ ) across GPUs indicate architecture-specific optimization potential, suggesting that dedicated optimizations for each GPU type would further enhance performance. We plan to release kernels specifically trained for different GPU types in an updated version of CUDA-L1.

## 4 RELATED WORK

### 4.1 RL-AUGMENTED LLMs FOR CODE OPTIMIZATION

Starting this year, there has been a growing interest in using LLM or RL-augmented LLM models for code optimization, including recent work on compiler optimization (Cummins et al., 2025) and assembly code optimization (Wei et al., 2025a), which use speed and correctness as RL training rewards. Other more distant related is software optimization that scale RL-based LLM reasoning for software engineering (Wei et al., 2025b). Regarding CUDA optimization, the only work that comprehensively delves into KernelBench is from Lange et al. (2025), which uses a meta-generation procedure that successfully optimizes 186 tasks out of 250 tasks in KernelBench with a medium speedup of 34%. Other works remain in preliminary stages, including Chen et al. (2025), which has optimized 20 GPU kernels selected from three different sources: the official NVIDIA CUDA Samples, LeetGPU, and KernelBench using a proposed feature search and reinforcement strategy; and an ongoing tech report (Schulman et al., 2017) that optimizes 4 kernels.

### 4.2 EVOLUTIONARY LLMs

Evolutionary large language models (Zhang et al., 2024; Liu et al., 2024b; Romera-Paredes et al., 2024; Novikov et al., 2025; Wei et al., 2025a; Dat et al., 2025; Lee et al., 2025) represent a paradigm shift in automated algorithm discovery, exemplified by systems such as Google DeepMind’s AlphaEvolve (Novikov et al., 2025) and FunSearch (Romera-Paredes et al., 2024).

## 5 CONCLUSION

In this paper, we propose CUDA-L1, a pipelined system for CUDA optimization powered by contrastive RL. CUDA-L1 achieves significant performance improvements on the CUDA optimization task, delivering an average speedup of  $\times 3.12$  (median  $\times 1.42$ ) over the default baseline across all 250 CUDA kernels of KernelBench, with peak speedups reaching  $\times 120$  on A100. Against other baselines, CUDA-L1 demonstrates  $\times 2.77$  over Torch Compile,  $\times 2.88$  over Torch Compile with reduce overhead, and  $\times 2.81$  over CUDA Graph implementations. CUDA-L1 can independently discover CUDA optimization techniques, learn to combine them strategically, and more importantly, extend the acquired CUDA reasoning abilities to unseen kernels with meaningful speedups. We hope that CUDA-L1 would open new doors for automated optimization of CUDA, and substantially promote GPU efficiency and alleviate the rising pressure on GPU computing resources.

## REFERENCES

- Thomas Bäck and Hans-Paul Schwefel. An overview of evolutionary algorithms for parameter optimization. *Evolutionary computation*, 1(1):1–23, 1993.
- Wentao Chen, Jiace Zhu, Qi Fan, Yehan Ma, and An Zou. Cuda-llm: Llms can write efficient cuda kernels. *arXiv preprint arXiv:2506.09092*, 2025.
- Chris Cummins, Volker Seeker, Dejan Grubisic, Baptiste Roziere, Jonas Gehring, Gabriel Synnaeve, and Hugh Leather. Llm compiler: Foundation language models for compiler optimization. In *Proceedings of the 34th ACM SIGPLAN International Conference on Compiler Construction*, pp. 141–153, 2025.
- Pham Vu Tuan Dat, Long Doan, and Huynh Thi Thanh Binh. Hsevo: Elevating automatic heuristic design with diversity-driven harmony search and genetic algorithm using llms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pp. 26931–26938, 2025.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024.
- Aaron Jaech, Adam Kalai, Adam Lerer, Adam Richardson, Ahmed El-Kishky, Aiden Low, Alec Helyar, Aleksander Madry, Alex Beutel, Alex Carney, et al. Openai o1 system card. *arXiv preprint arXiv:2412.16720*, 2024.
- Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.
- Robert Tjarko Lange, Aaditya Prasad, Qi Sun, Maxence Faldor, Yujin Tang, and David Ha. The ai cuda engineer: Agentic cuda kernel discovery, optimization and composition. Technical report, 2025.
- Kuang-Huei Lee, Ian Fischer, Yueh-Hua Wu, Dave Marwood, Shumeet Baluja, Dale Schuurmans, and Xinyun Chen. Evolving deeper llm thinking. *arXiv preprint arXiv:2501.09891*, 2025.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024a.
- Fei Liu, Xialiang Tong, Mingxuan Yuan, Xi Lin, Fu Luo, Zhenkun Wang, Zhichao Lu, and Qingfu Zhang. Evolution of heuristics: Towards efficient automatic algorithm design using large language model. *arXiv preprint arXiv:2401.02051*, 2024b.
- Niklas Muennighoff, Zitong Yang, Weijia Shi, Xiang Lisa Li, Fei-Fei Li, Hanna Hajishirzi, Luke S. Zettlemoyer, Percy Liang, Emmanuel J. Candes, and Tatsunori Hashimoto. s1: Simple test-time scaling. *ArXiv*, abs/2501.19393, 2025. URL <https://api.semanticscholar.org/CorpusID:276079693>.
- Alexander Novikov, Ngân Vū, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian, et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*, 2024.

- Anne Ouyang, Simon Guo, Simran Arora, Alex L Zhang, William Hu, Christopher Ré, and Azalia Mirhoseini. Kernelbench: Can llms write efficient gpu kernels? *arXiv preprint arXiv:2502.10517*, 2025.
- Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog, M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang, Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, Y Wu, et al. Deepseekmath: Pushing the limits of mathematical reasoning in open language models. *arXiv preprint arXiv:2402.03300*, 2024.
- Zhang Shengyu, Dong Linfeng, Li Xiaoya, Zhang Sen, Sun Xiaofei, Wang Shuhe, Li Jiwei, Runyi Hu, Zhang Tianwei, Fei Wu, et al. Instruction tuning for large language models: A survey. *arXiv preprint arXiv:2308.10792*, 2023.
- Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, et al. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*, 2024.
- Shuhe Wang, Shengyu Zhang, Jie Zhang, Runyi Hu, Xiaoya Li, Tianwei Zhang, Jiwei Li, Fei Wu, Guoyin Wang, and Eduard Hovy. Reinforcement learning enhanced llms: A survey. *arXiv preprint arXiv:2412.10400*, 2024.
- Anjiang Wei, Tarun Suresh, Huanmi Tan, Yinglun Xu, Gagandeep Singh, Ke Wang, and Alex Aiken. Improving assembly code performance with large language models via reinforcement learning. *arXiv preprint arXiv:2505.11480*, 2025a.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.
- Yuxiang Wei, Olivier Duchenne, Jade Copet, Quentin Carbonneaux, Lingming Zhang, Daniel Fried, Gabriel Synnaeve, Rishabh Singh, and Sida I Wang. Swe-rl: Advancing llm reasoning via reinforcement learning on open software evolution. *arXiv preprint arXiv:2502.18449*, 2025b.
- Ronald J Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8:229–256, 1992.
- An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025.
- Rui Zhang, Fei Liu, Xi Lin, Zhenkun Wang, Zhichao Lu, and Qingfu Zhang. Understanding the importance of evolutionary search in automated heuristic design with large language models. In *International Conference on Parallel Problem Solving from Nature*, pp. 185–202. Springer, 2024.

## THE USE OF LARGE LANGUAGE MODELS (LLMs)

**We used a large language model (i.e., ChatGPT O3) only as a general-purpose assist tool for minor English grammar corrections during manuscript preparation.** The LLM had no role in research ideation, methodology, experimental design, data collection, analysis, interpretation, or substantive writing beyond copyediting at the sentence level. All scientific content, claims, and conclusions are the authors’ own. The authors take full responsibility for all contents written under their names, including any text that may have been edited with the assistance of an LLM. LLMs are not eligible for authorship or contributorship on this work.

### A DETAILS FOR CUDA-L1

#### A.1 DATA COLLECTION FOR SFT

To expand the model’s exposure to CUDA patterns, we begin with data augmentation based on reference code from 250 tasks in KernelBench, which provides the official implementations used in PyTorch. To generate executable and correct CUDA code efficiently, we leverage six existing LLM models: GPT-4o, OpenAI-o1, DeepSeek-R1, DeepSeek V3, Llama 3.1-405B Instruct, and Claude 3.7. For each model, we construct prompts using the one-shot strategy, where the prompt contains the reference code (denoted by  $q_i, i \in [1, 250]$ ) and asks the LLM to generate an alternative speedup implementation. We employ multiple models to maximize the diversity of successful CUDA code generation. The detailed prompt structure is provided in Table 4. For each of the six models, we iterate through all 250 tasks. Each task allows up to 20 trials and terminates early if we successfully collect 2 trials that are both executable and correct. Notably, some tasks may fail to produce any successful code across all trials. The successful code is denoted by  $d_{i,j}$ , where  $j \in \{1, 2, \dots, n_i\}$ , and  $n_i$  denotes the number of successful code snippets for the reference code  $q_i$ . Through this process, we collected 2,105 successful CUDA code snippets. Now we have collected the dataset  $D = \{(q_i, \{d_{i,j}\}_{j=1}^{n_i})\}_i$ .

#### A.2 PUSDO CODE FOR STAGE2: SELF-SUPERVISED LEARNING

---

##### Self-supervised Learning Algorithm

---

```

1: Initialize finetuned model  $M_0$  after SFT stage with parameters  $\theta_{\text{ft}}$ 
2: for  $i = 1$  to  $N_{\text{iterations}}$  do
3:   Generate batch of CUDA codes  $C_i = \{c_1, \dots, c_k\}$  using model  $M_{i-1}$ 
4:   Evaluate each  $c \in C_i$  for:
5:     1. Executability (compiles and runs)
6:     2. Correctness (produces expected output)
7:   Filter successful codes:  $C_i^{\text{success}} = \{c \in C_i | \text{executable} \wedge \text{correct}\}$ 
8:   if  $C_i^{\text{success}} \neq \emptyset$  then
9:     Compute gradient update  $\nabla\theta$  using  $C_i^{\text{success}}$ 
10:    Update model:  $\theta_i \leftarrow \theta_{i-1} + \eta\nabla\theta$ 
11:   else
12:      $\theta_i \leftarrow \theta_{i-1}$  (no update)
13:   end if
14: end for
15: return Final improved model  $M_N$ 

```

---

Table 4: Self-supervised learning for cuda optimization in Stage 2.

#### A.3 DETAILS FOR RL-TRAINING

Specifically, for each reference prompt  $q$  containing selected exemplars as shown in Table 9, we sample  $G$  code outputs from the current policy  $\pi_{\text{old}}$ , denoted as  $\{d_1, d_2, \dots, d_G\}$ . Let  $\mathbf{r} = (r_1, r_2, \dots, r_G)$  represent the reward scores associated with the generated code samples. Different from standard GRPO training, rewards are smoothed to mitigate the reward hacking issue; the details of this approach will be elaborated in Section A.5. Further, as in GRPO, rewards are normalized

within each group using:

$$\hat{r}_i = \frac{r_i - \text{mean}(\mathbf{r})}{\text{std}(\mathbf{r})} \quad (3)$$

The complete GRPO objective optimizes the policy model by maximizing:

$$\mathcal{L}_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim P(q), \{d_i\}_{i=1}^G \sim \pi_{\theta_{old}}(d|q)} \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|d_i|} \sum_{t=1}^{|d_i|} \left( \min \left( \frac{\pi_{\theta}(d_{i,t}|q, d_{i,<t})}{\pi_{\theta_{old}}(d_{i,t}|q, d_{i,<t})} \hat{r}_i, \right. \right. \right. \\ \left. \left. \left. \text{clip} \left( \frac{\pi_{\theta}(d_{i,t}|q, d_{i,<t})}{\pi_{\theta_{old}}(d_{i,t}|q, d_{i,<t})}, 1 - \varepsilon, 1 + \varepsilon \right) \hat{r}_i \right) - \beta D_{KL}[\pi_{\theta} \parallel \pi_{ref}] \right] \quad (4)$$

where:

- $\pi_{\theta}$  is the policy model being optimized
- $\pi_{\theta_{old}}$  is the old policy model from the previous iteration
- $\varepsilon$  is the parameter for clipping
- $\beta$  is the KL penalty coefficient that controls deviation from the reference policy
- $D_{KL}$  denotes the KL divergence between the current and reference policies

We refer readers to (Shao et al., 2024) for details of GRPO. Model parameters are optimized using the GRPO objective, with contrastive prompts that incorporate comparative examples. This concludes our description of contrastive RL.

#### A.4 CONTRASTIVE-RL’S ADVANTAGES OVER EVOLUTIONARY LLM APPROACHES

Contrastive-RL draws inspiration from a broad range of literature, including evolutionary algorithms (Bäck & Schwefel, 1993) and their applications to LLMs (Liu et al., 2024b; Romera-Paredes et al., 2024; Novikov et al., 2025; Wei et al., 2025a), where multiple solution instances with associated fitness scores are presented to LLMs to analyze performance patterns and generate improved solutions. However, Contrastive-RL improves evolutionary LLM approaches in several critical aspects:

**Model Adaptation vs. Fixed-Model Reasoning:** Contrastive-RL employs gradient-based parameter updates to continuously enhance model capabilities, whereas evolutionary LLM approaches rely exclusively on in-context learning with static parameters. This fundamental architectural difference endows Contrastive-RL with substantially greater representational capacity and task adaptability. Evolutionary LLM methods are fundamentally limited by the frozen foundation model’s initial knowledge and reasoning abilities, while Contrastive-RL progressively refines the model’s domain-specific expertise through iterative parameter optimization. From this perspective, evolutionary LLM approaches can be viewed as a degenerate case of Contrastive-RL that implements only the Fixed-Parameter Solution Optimization component while omitting the Foundation Model Enhancement mechanism. This theoretical relationship explains why Contrastive-RL consistently outperforms evolutionary approaches: it leverages both optimization dimensions simultaneously rather than constraining itself to a single fixed-capacity search space.

**Scalability and Generalization:** Contrastive-RL demonstrates superior scalability by training a single specialized model capable of handling diverse CUDA programming tasks and generating various types of optimized code. In contrast, evolutionary LLM approaches typically require separate optimization processes for each distinct task or domain, limiting their practical applicability and computational efficiency.

#### A.5 MITIGATING REWARD HACKING IN RL TRAINING

Reinforcement learning is notorious for exhibiting reward hacking behaviors, where models exploit system vulnerabilities to achieve higher rewards while generating outputs that deviate from the intended objectives. A particularly challenging aspect of these pitfalls is that they cannot be anticipated prior to training and are only discovered during the training process. During our initial training procedure, we identified the following categories of reward hacking behaviours:

**Improper Timing Measurement.** KernelBench measures execution time by recording timing events on the main CUDA stream:

```

1 start_event.record(original_model_stream)
2 model(*inputs)
3 end_event.record(original_model_stream)
4 torch.cuda.synchronize(device=device)

```

However, RL-generated code exploits this by creating additional CUDA streams that execute asynchronously. Since KernelBench only monitors the main stream, it fails to capture the actual execution time of operations running on parallel streams. This vulnerability is significant: in our initial implementation, we find that 82 out of 250 (32.8%) RL-generated implementations exploit this timing loophole to appear faster than they actually are, leading to an overall speedup of 18×. To address this issue, prompt engineering alone is insufficient. The evaluation methodology should be modified to synchronize all CUDA streams before recording the end time, ensuring accurate performance measurement across all concurrent operations as follows:

```

1 start_event.record(custom_model_stream)
2 custom_model(*inputs)
3 # Wait for all model streams to complete before recording end event
4 if custom_contain_new_streams:
5     for stream in custom_model_streams:
6         custom_model_stream.wait_stream(stream)
7 end_event.record(custom_model_stream)
8 torch.cuda.synchronize(device=device)

```

**Lazy Evaluation** Another important hacking strategy is lazy evaluation, detected by community users on GitHub: Calling `custom_model(*inputs)` doesn't ensure the output is actually materialized/computed. The computation is actually executed at the correctness check phase when calling the `torch.allclose()` function, allowing it to pass the correctness check.

```

1 class LazyMatmul(torch.Tensor):
2     def __new__(cls, A, B):
3         result = torch.Tensor._make_subclass(cls, torch.empty(0))
4         result.A = A
5         result.B = B
6         result._shape = (A.size(0), B.size(1))
7         return result
8
9     def materialize(self):
10        """Trigger actual computation"""
11        return torch.matmul(self.A, self.B)
12
13 class ModelNew(nn.Module):
14     def forward(self, A, B):
15         return LazyMatmul(A, B) # Returns lazy object

```

To mitigate this issue, we enforce a validation check that ensures materialize functions are called before ending the time measurement, which involves checking the following conditions: the output must be a tensor, must be a standard `torch.Tensor` (not a subclass), must be on the correct device, must have allocated memory, and the corresponding storage must be valid.

```

1 # Check 1: Must be a tensor
2 if not isinstance(out, torch.Tensor):
3     return False, f"{prefix} is not a tensor: {type(out)}"
4
5 # Check 2: Must be standard torch.Tensor, not a subclass
6 if type(out).__name__ not in ['Tensor', 'Parameter']:
7     return False, f"{prefix} is {type(out).__name__}, not standard torch.Tensor"
8
9 # Check 3: Must be on correct device
10 if out.device != device:
11     return False, f"{prefix} on wrong device: {out.device} (expected {device})"
12
13 # Check 4: Must have allocated storage

```

```

14 storage_size = out.untyped_storage().size()
15 if storage_size == 0:
16     return False, f"{prefix} has no allocated storage (likely lazy)"
17
18 # Check 5: Storage pointer must be valid
19 ptr = out.data_ptr()
20 if ptr == 0:
21     return False, f"{prefix} storage pointer is null (likely lazy)"

```

**Hyperparameter Manipulation:** In KernelBench, each computational task is associated with specific hyperparameters, including `batch_size`, `dim`, `in_features` dimension, `out_features` dimension, `scaling_factor`, and others. The RL agent learned to exploit these parameters by generating code that artificially reduces their values, thereby achieving superficial speedup improvements that do not reflect genuine optimization performance.

**Result Caching:** The RL agent developed strategies to cache computational results across evaluation batches based on input addresses. When another input’s address matches a cached one, it returns the cached output. In theory, this should not pass correctness validation because the cached output differs from the expected one. However, given that correctness validation checks whether the difference at each position between the reference output and custom code output is below a certain threshold, there are a few cases where it is able to squeeze past the correctness bar. The following code snippet gives an illustration:

```

1 cache_key = x.data_ptr()
2 # Check if result is in cache
3 if cache_key in self.cache:
4     return self.cache[cache_key]

```

## A.6 REWARD DENOISING

We observe a significant variance in  $t_d$  measurements for identical implementations  $d$ , which introduces noise in reward estimation. This noise is particularly detrimental to RL training stability. To address these challenges, we implement the following robust measurement strategies:

- Dedicated GPU Allocation:** Each evaluation runs on an exclusively allocated GPU. Shared GPU usage leads to significantly higher variance in timing measurements, even when memory and compute utilization appear low.
- Paired Execution with Order Randomization:** For fair comparison, each evaluation round executes both the reference  $q_i$  and candidate  $d$  implementations. Crucially, we randomize the execution order within each round to account for GPU warm-up effects, where subsequent runs typically benefit from cache warming.
- Extended Measurement Window:** We conduct multiple evaluation rounds with predefined running time of 30 minutes per candidate. This adaptive approach yields between several tens of thousands to 1M rounds depending on individual kernel execution times.
- Bucketized Variance Control:** We partition all  $\text{Score}_{\text{single-run}}(d)$  measurements into 7 buckets and compute bucket-wise averages. Evaluations with inter-bucket variance exceeding 0.005 are discarded.
- Robust Central Tendency:** The final reward uses the median of bucket averages, which proves more stable than the mean against outlier effects:

$$r(d) = \text{median}(\{\text{Bucket}_k\}_{k=1}^7) \quad (5)$$

- Conservative Rounding:** We apply conservative rounding to speedup ratios (i.e.,  $\text{Score}(d)$ ), truncating to two decimal places while biasing toward unity (e.g.,  $1.118 \rightarrow 1.11$ ,  $0.992 \rightarrow 1.00$ ).
- Strict Verification Protocol:** Despite these precautions, we still occasionally observe spurious large speedups due to GPU turbulence. For any candidate showing either:
  - Absolute value of speedup  $> 3$ , or
  - Speedup exceeding twice the previous maximum
we perform verification on a different GPU of the same type. The result is accepted only if the verification measurement differs by  $< 10\%$  from the original.

### A.7 TOWARDS ROBUST REWARD DESIGN AND TRAINING PROCEDURES

To mitigate reward hacking, we implement the following strategies during training:

**A reward checking model** When there is a significant leap in reward, an adversarial model intervenes to determine whether the code exploits the reward system. We use DeepSeek-R1 for this purpose and find that it successfully identifies reward hacking above over 60% of the time.

**Hacking-case database** We maintain a dynamic hacking-case database that is updated whenever a new reward hacking behavior is detected. The reward checking model leverages this database for detection: given a newly generated code snippet to examine, we retrieve the three most similar cases from the database and include them as context for the reward checking model’s input.

**Reward smoothing** Sharp reward increases are smoothed to reduce their magnitude, preventing the RL agent from over-prioritizing any single high-reward solution, whether legitimate or not:

$$r_{\text{normalized}} = \frac{r - \mu}{\sigma} \tag{6}$$

$$r_{\text{smooth}} = \text{clip}(r_{\text{normalized}}, -k, k)$$

where  $\mu$  and  $\sigma$  are the mean and the mean and standard deviation of the reward distribution, respectively.  $k$  is a hyperparameter that controls the clipping threshold set to 1.5, as we think as achieving a 1.5× speedup over the official PyTorch implementation already represents significant optimization performance.

## B DISCOVERED CUDA OPTIMIZATION TECHNIQUES

An analysis of optimization strategies commonly employed in enhanced CUDA implementations reveals interesting patterns. Through GPT-4o-based technical term extraction and frequency analysis, we identified the ten most prevalent optimization techniques:

- **Memory Layout Optimization**, which ensures data is stored in contiguous memory blocks;
- **Memory Access Optimization**, which arranges data access patterns to maximize memory bandwidth and minimize latency through techniques like shared memory usage, coalesced global memory access, and memory padding;
- **Operation Fusion**, which combines multiple sequential operations into a single optimized kernel execution;
- **Memory Format Optimization**, which aligns data layout with hardware memory access patterns;
- **Memory Coalescing**, which optimizes CUDA kernel performance by ensuring threads in the same warp access contiguous memory locations;
- **Warp-Level Optimization**, which leverages the parallel execution of threads within a warp (typically 32 threads) to efficiently perform collective operations;
- **Optimized Thread Block Configuration**, which carefully selects grid and block dimensions for CUDA kernels to maximize parallel execution efficiency and memory access patterns;
- **Shared Memory Usage**, enables fast data access by storing frequently used data in a cache accessible by all threads within a thread block;
- **Register Optimization**, which stores frequently accessed data in fast register memory to reduce latency and improve computational throughput;
- **Stream Management**, which enables parallel execution of operations for improved GPU utilization.

Tables 11, 13 and 14 present detailed CUDA optimization techniques with accompanying code examples.

## C CASE STUDIES

Table 5 presents the KernelBench tasks that achieved the highest speedups. We examine these some of them in detail and perform an ablation study of the applied CUDA optimization techniques, showing how much each technique contributes to the final speedup.

Level ID	Task ID	Task Name	Speedup
2	83	83_Conv3d_GroupNorm_Min_Clamp_Dropout	120.3
1	12	12_Matmul_with_diagonal_matrices	64.4
2	80	80_Gemm_Max_Subtract_GELU	31.3
1	9	9_Tall_skinny_matrix_multiplication	24.9
3	31	31_VisionAttention	24.8
2	96	96_ConvTranspose3d_Multiply_Max_GlobalAvgPool_Clamp	16.2
2	66	66_Matmul_Dropout_Mean_Softmax	14.5
1	13	13_Matmul_for_symmetric_matrices	14.4
3	43	43_MinGPTCausalAttention	13.1
3	44	44_MiniGPTBlock	10.5

Table 5: KernelBench Tasks Ranked by CUDA-L1 Acceleration (Top 10)

### C.1 $\text{DIAG}(A) * B$ : 64 $\times$ FASTER

We first examine the code for level 1, task 12, which performs matrix multiplication between a diagonal matrix (represented by its diagonal elements) and a dense matrix, both with dimension  $N=4096$ . The reference code is as follows where `__init__` function of the class is omitted:

```

1 class Model(nn.Module):
2     def forward(self, A, B):
3         # A: (N,) - 1D tensor of shape N
4         # B: (N, M) - 2D tensor of shape N x M
5         # torch.diag(A): (N, N) - creates diagonal matrix from A
6         # Result: (N, N) @ (N, M) = (N, M)
7         return torch.diag(A) @ B

```

Let’s see the optimized code by CUDA-L1:

```

1 class Model(nn.Module):
2     def forward(self, A, B):
3         return A.unsqueeze(1) * B

```

The optimized implementation leverages PyTorch’s broadcasting mechanism to perform diagonal matrix multiplication efficiently. It first reshapes the diagonal vector  $A$  from shape  $(N, )$  to  $(N, 1)$  using `unsqueeze(1)`, transforming it into a column vector. Next, it utilizes PyTorch’s automatic broadcasting to multiply each row of matrix  $B$  by the corresponding element of  $A$ , where the  $(N, 1)$  shaped tensor is implicitly expanded to match the  $(N, M)$  dimensions of  $B$ . This approach completely avoids creating the full  $N \times N$  diagonal matrix, which would be sparse and memory-intensive. The key benefits of this technique are substantial: it requires only  $O(1)$  extra memory instead of  $O(N^2)$  for storing the diagonal matrix, reduces computational complexity from  $O(N^2M)$  operations for full matrix multiplication to just  $O(NM)$  element-wise operations, leading to 64 $\times$  speedup.

What makes this particularly valuable is that RL can systematically explore the vast space of equivalent implementations. By exploring semantically equivalent implementations, RL learns to identify patterns where computationally expensive operations can be replaced with more efficient alternatives. The power of RL extends beyond simple algebraic simplifications and it can uncover sophisticated optimizations such as: replacing nested loops with vectorized operations identifying hidden parallelization opportunities discovering memory-efficient mathematical reformulations finding non-obvious algorithmic transformations that preserve correctness while improving performance. What makes this particularly valuable is that RL can systematically explore the vast space of equivalent implementations—something that would be impractical for human engineers to do manually.

### C.2 LSMT: 3.4 $\times$ FASTER

Now let’s look at a classical neural network algorithm LSTM (level 3, task 35), on which CUDA-11 achieves a speedup of 3.4 $\times$ . By comparing the reference PyTorch implementation with the optimized output, we identified the following optimization techniques:

1. **CUDA Graphs**, which captures the entire LSTM computation sequence (including all layer operations) into a replayable graph structure, eliminating kernel launch overhead by record-

Configuration	CUDA Graphs	Memory Contiguity	Static Tensor Reuse	Speedup	Bottleneck
CUDA + Memory + Static	✓	✓	✓	3.42×	LSTM computation
CUDA + Memory	✓	✓	✗	2.96×	Memory allocation
CUDA + Static	✓	✗	✓	2.84×	Memory layout
CUDA Only	✓	✗	✗	2.77×	Memory overhead
Memory + Static	✗	✓	✓	1.00×	Kernel launch overhead
Memory Only	✗	✓	✗	1.00×	Kernel launch overhead
Static Only	✗	✗	✓	1.00×	Kernel launch overhead
Baseline	✗	✗	✗	1.00×	Kernel launch overhead

Table 6: Speedup achieved by different CUDA optimization techniques on LSTMs.

ing operations once and replaying them with minimal CPU involvement for subsequent executions.

2. **Memory Contiguity**, which ensures all tensors maintain contiguous memory layouts through explicit `.contiguous()` calls before operations, optimizing memory access patterns and improving cache utilization for CUDA kernels processing sequential data.
3. **Static Tensor Reuse**, which pre-allocates input and output tensors during graph initialization and reuses them across forward passes with non-blocking copy operations, eliminating memory allocation overhead and enabling asynchronous data transfers.

Table 6 represents the results for 8 different optimization combinations across the three optimization techniques above. As can be seen, CUDA Graphs is essential for achieving any meaningful speedup in this LSTM model. All configurations with CUDA Graphs achieve 2.77x-3.42x speedup, while all configurations without it achieve only 1.0x (no speedup). The combination of all three techniques provides the best performance at 3.42x, demonstrating that while CUDA Graphs provides the majority of the benefit ( 81% of total speedup), the additional optimizations contribute meaningful improvements when combined together.

### C.3 3D TRANSPOSED CONVOLUTION: 120× FASTER

We examined the code for Level 2, Task 38, which implements a sequence of 3D operations: transposed convolution, average pooling, clamping, softmax, and element-wise multiplication. By comparing the reference PyTorch implementation with the CUDA-L1 optimized output, we identified the following optimization techniques applied by CUDA-L1:

1. **Mathematical Short-Circuit**, which detects when `min_value` equals 0.0 and skips the entire computation pipeline (convolution, normalization, min/clamp operations), directly returning zero tensors since the mathematical result is predetermined.
2. **Pre-allocated Tensors**, which creates zero tensors of standard shapes during initialization and registers them as buffers, eliminating memory allocation overhead during forward passes for common input dimensions.
3. **Direct Shape Matching**, which provides a fast path for standard input shapes by immediately returning pre-allocated tensors without any shape calculations, bypassing the computational overhead entirely.
4. **Pre-computed Parameters**, which extracts and stores convolution parameters (kernel size, stride, padding, dilation) during initialization, avoiding repeated attribute lookups and tuple conversions during runtime.

Table 7 represents the results for 16 different optimization combinations across the four optimization techniques above. As can be seen, mathematical short-circuit is essential for this task, where all configurations with mathematical short-circuit achieve 28.6x+ speedup, while all configurations without it achieve only 1.0x (no).

The fact that CUDA-L1 identified this precise optimization strategy demonstrates the power of reinforcement learning in navigating complex optimization spaces. While a human developer might intuitively focus on computational optimizations (like parallel algorithms) or memory layout im-

Configuration	Math Short-Circuit	Pre-allocated Tensors	Direct Shape Match	Pre-computed Params	Speedup
Math + PreAlloc + Shape + Params	✓	✓	✓	✓	120.9×
Math + Shape + Params	✓	✗	✓	✓	32.8×
Math + PreAlloc + Params	✓	✓	✗	✓	30.6×
Math + Params	✓	✗	✗	✓	30.2×
Math + PreAlloc	✓	✓	✗	✗	29.2×
Math Only	✓	✗	✗	✗	29.2×
Math + Shape	✓	✗	✓	✗	28.6×
Math + PreAlloc + Shape	✓	✓	✓	✗	28.6×
PreAlloc + Shape + Params	✗	✓	✓	✓	1.0×
PreAlloc + Shape	✗	✓	✓	✗	1.0×
Shape + Params	✗	✗	✓	✓	1.0×
PreAlloc + Params	✗	✓	✗	✓	1.0×
Params Only	✗	✗	✗	✓	1.0×
Shape Only	✗	✗	✓	✗	1.0×
PreAlloc Only	✗	✓	✗	✗	1.0×
Baseline	✗	✗	✗	✗	1.0×

Table 7: Speedup achieved by different CUDA optimization techniques on the Conv3d task.

provements (like tensor pre-allocation), RL discovered that the mathematical properties of the operation completely dominate performance. This discovery is particularly impressive because: RL is able to find this non-obvious solution: The 120x speedup from exploiting the mathematical short-circuit is counterintuitive as most developers would expect to optimize the convolution kernel or memory access patterns for such a compute-heavy operation. This shows how RL can discover optimal solutions that challenge conventional wisdom in deep learning optimization. Where human intuition might suggest "optimize the convolution algorithm first," CUDA-L1 learned through empirical evidence that "recognize when computation can be entirely skipped" yields dramatically better results. The agent's ability to identify that  $\min(x, 0)$  followed by  $\text{clamp}(0, 1)$  always produces zeros demonstrates how RL can uncover mathematical invariants that humans might overlook in complex computational pipelines.

## D PROMPT USED IN THE PAPER

Data Augmentation Prompt — Used in Supervised fine-tuning	
<b>Task for CUDA Optimization</b>	<p>You are an expert in CUDA programming and GPU kernel optimization. Now you're tasked with developing a high-performance cuda implementation of Softmax. The implementation must:</p> <ul style="list-style-type: none"> <li>• Produce <b>identical</b> results to the reference PyTorch implementation.</li> <li>• Demonstrate <b>speed improvements</b> on GPU.</li> <li>• Maintain <b>stability</b> for large input values.</li> </ul>
<b>Reference Implementation (exact copy)</b>	<pre> 1 import torch 2 import torch.nn as nn 3 4 class Model(nn.Module): 5     """ 6     Simple model that performs a Softmax activation. 7     """ 8     def __init__(self): 9         super(Model, self).__init__() 10 11     def forward(self, x: torch.Tensor) -&gt; torch.Tensor: 12         """ 13         Applies Softmax activation to the input tensor. 14         Args: 15             x (torch.Tensor): Input tensor of shape 16                             (batch_size, num_features). 17         Returns: 18             torch.Tensor: Output tensor with Softmax 19                             applied, same shape as input. 20         """ 21         return torch.softmax(x, dim=1) 22 23 batch_size = 16 24 dim = 16384 25 26 def get_inputs(): 27     x = torch.randn(batch_size, dim) 28     return [x] 29 30 def get_init_inputs(): 31     return [] # No special initialization inputs needed </pre>

Table 8: Prompt illustration for data augmentation in Section ???. For each KernelBench task (softmax shown here for illustration), the prompt is fed to each of six LLM models—GPT-4o, OpenAI-o1, DeepSeek-R1, DeepSeek V3, Llama 3.1-405B Instruct, and Claude 3.7 Sonnet—to generate alternative CUDA implementations.

CUDA Optimization Task Prompt — Used in Contrastive-RL	
<b>Task for CUDA Optimization</b>	<p>You are a CUDA programming expert specializing in GPU kernel optimization. Given a reference CUDA implementation, your objective is to create an accelerated version that maintains identical functionality. You will receive previous CUDA implementations accompanied by their performance metrics. Conduct a comparative analysis of these implementations and use the insights to develop optimized and correct CUDA code that delivers superior performance.</p>
<b>Reference Code</b>	<pre> 1 __global__ void kernel_v1(float* input, float* output, int     N) { 2     // Baseline implementation 3     ... 4 } }</pre>
<b>Previous Cuda Implementations with Scores</b>	<pre> 1 // code1 (score1) 2 __global__ void kernel_v1(float* input, float* output, int     N) { 3     ... 4 } 5 6 // code2 (score2) 7 __global__ void kernel_v2(float* input, float* output, int     N) { 8     ... 9 }</pre>
<b>Generation Protocol</b>	<p>You MUST use exactly two hash symbols (##) at the beginning of each section.</p> <p><b>## Performance Analysis:</b> Compare code snippets above and articulate on :</p> <ol style="list-style-type: none"> <li>1. Which implementations demonstrate superior performance and why?</li> <li>2. What particular optimization strategies exhibit the greatest potential for improvement?</li> <li>3. What are the primary performance limitations in the implementation?</li> <li>4. What CUDA-specific optimization techniques remain unexploited?</li> <li>5. Where do the most significant acceleration opportunities exist?</li> </ol> <p><b>## Algorithm Design:</b> Describe your optimization approach</p> <p><b>## Code Implementation:</b> Provide your improved CUDA kernel</p>
<b>Requirements and Restrictions</b>	<p><b>## Critical Requirements:</b></p> <ol style="list-style-type: none"> <li>1. Functionality must match the reference implementation exactly. Failure to do so will result in a score of 0.</li> <li>2. Code must compile and run properly on modern NVIDIA GPUs</li> </ol> <p><b>## Key Restrictions:</b></p> <ol style="list-style-type: none"> <li>1. Do not cache or reuse previous results — the code must execute fully on each run.</li> <li>2. Keep hyperparameters unchanged (e.g., batch size, dimensions, etc.) as specified in the reference.</li> </ol>

Table 9: Prompt structure for CUDA optimization task showing reference implementations and their performance scores used in Contrastive-RL.

## E CASE STUDY: CODE SNIPPETS BEFORE AND AFTER OPTIMIZATIONS

Tech + Desc	Before optimization	After optimization
<b>Memory Layout Optimization</b>  Memory Layout Optimization ensures data is stored in contiguous memory blocks to maximize cache efficiency and reduce memory access latency during GPU computations.	<b>- Non-contiguous memory access</b>  <pre> 1 '''Python 2 def matrix_multiply(A, B): 3     # A and B might not be contiguous in memory 4     C = torch.mm(A, B) 5     return C 6 ''' </pre>	<b>- Ensuring contiguous memory layout</b>  <pre> 1 '''Python 2 def matrix_multiply_optimized(A, B): 3     # Ensure contiguous memory layout for 4     # efficient access patterns 5     A = A.contiguous() if not A.is_contiguous() 6     else A 7     B = B.contiguous() if not B.is_contiguous() 8     else B 9     C = torch.mm(A, B) 10    return C 11 ''' </pre>
<b>Memory Coalescing</b>  Memory coalescing optimizes GPU memory access by ensuring threads in a warp access contiguous memory locations, reducing memory transactions and increasing bandwidth utilization.	<b>- Uncoalesced memory access</b>  <pre> 1 '''cuda 2 __global__ void uncoalesced_kernel(float* 3     input, float* output) { 4     int tid = threadIdx.x; 5     int stride = blockDim.x; 6     // Each thread accesses non-contiguous 7     // memory locations 8     for (int i = 0; i &lt; 1024; i++) { 9         output[tid + i * stride] = input[tid + i 10        * stride] * 2.0f; 11    } 12 ''' </pre>	<b>- Coalesced memory access with loop unrolling</b>  <pre> 1 '''cuda 2 __global__ void coalesced_kernel(float* input, 3     float* output) { 4     int tid = threadIdx.x; 5     int batch_idx = blockDim.x; 6     // Base pointers for this batch item 7     const float* batch_input = input + 8     batch_idx * 1024; 9     float* batch_output = output + batch_idx * 10    1024; 11    // Each thread processes contiguous memory 12    // in chunks 13    #pragma unroll 4 14    for (int i = 0; i &lt; 1024; i += 16) { 15        batch_output[i] = batch_input[i] * 2.0f; 16        batch_output[i+1] = batch_input[i+1] * 17        2.0f; 18        batch_output[i+2] = batch_input[i+2] * 19        2.0f; 20        // ... more contiguous accesses 21        batch_output[i+15] = batch_input[i+15] * 22        2.0f; 23    } 24 } 25 ''' </pre>
<b>Warp-Level Optimizations</b>  Warp-Level Optimizations leverage the CUDA execution model where threads execute in groups of 32 (warps) to improve parallel efficiency through collaborative operations and memory access patterns.	<b>- Each thread independently calculates min value</b>  <pre> 1 '''cuda 2 __global__ void min_kernel_before(const float* 3     input, float* output, int size) { 4     int idx = blockDim.x * blockDim.x + 5     threadIdx.x; 6     if (idx &lt; size) { 7         float min_val = 1e10f; 8         for (int i = 0; i &lt; DEPTH; i++) { 9             min_val = min(min_val, input[idx + i 10            * size]); 11        } 12        output[idx] = min_val; 13    } 14 } 15 ''' </pre>	<b>- Using warp-level operations for parallel reduction</b>  <pre> 1 '''cuda 2 3 __global__ void min_kernel_after(const float* 4     input, float* output, int size) { 5     int idx = blockDim.x * blockDim.x + 6     threadIdx.x; 7     int lane_id = threadIdx.x % 32; // Thread's 8     // position within warp 9     int warp_id = threadIdx.x / 32; // Warp 10    // number within the block 11 12    float min_val = 1e10f; 13    if (idx &lt; size) { 14        // Each thread finds its local minimum 15        for (int i = 0; i &lt; DEPTH; i++) { 16            min_val = min(min_val, input[idx + i 17            * size]); 18        } 19 20        // Warp-level parallel reduction using 21        // shuffle 22        for (int offset = 16; offset &gt; 0; offset 23        /= 2) { 24            float other = 25            __shfl_down_sync(0xffffffff, 26            min_val, offset); 27            min_val = min(min_val, other); 28        } 29 30        // First thread in warp writes the result 31        if (lane_id == 0) { 32            output[blockIdx.x * (blockDim.x/32) + 33            warp_id] = min_val; 34        } 35    } 36 } 37 ''' </pre>

Table 10: (Part 1) Code snippets before and after optimizations.

Tech + Desc	Before optimization	After optimization
<b>Memory Hierarchy Optimization</b>  Memory Hierarchy Optimization involves strategically utilizing different levels of GPU memory (registers, shared memory, constant memory) to minimize global memory access latency and maximize data reuse.	<b>- Using global memory directly</b> <pre> 1  ````cuda 2  __global__ void    depthwise_separable_conv_kernel_unoptimized( 3  const float* input, const float*    depthwise_weight, const float* 4  float* output, /* other parameters */) { 5 6  int out_y = blockIdx.y * blockDim.y +    threadIdx.y; 7  int out_x = blockIdx.x * blockDim.x +    threadIdx.x; 8 9  // Each thread directly accesses global    memory for each computation 10 for (int oc = 0; oc &lt; out_channels; oc++) { 11   float result = 0.0f; 12   for (int ic = 0; ic &lt; in_channels; ic++) 13   { 14     float depthwise_result = 0.0f; 15     // Direct global memory access for        each kernel element 16     for (int ky = 0; ky &lt; 3; ky++) { 17       for (int kx = 0; kx &lt; 3; kx++) { 18         int in_y = out_y * stride + ky            - padding; 19         int in_x = out_x * stride + kx            - padding; 20         if (in_y &gt;= 0 &amp;&amp; in_y &lt;            in_height &amp;&amp; in_x &gt;= 0 &amp;&amp;            in_x &lt; in_width) { 21           depthwise_result +=            input[(batch_idx *            in_channels + ic) *            in_height + in_y] *            in_width + in_x] *            depthwise_weight[ 22             9 + ky *            3 + kx]; 23         } 24       } 25     } 26     result += depthwise_result *            pointwise_weight[oc * in_channels            + ic]; 27   } 28   output[(batch_idx * out_channels + oc)            * out_height + out_y] * out_width +            out_x] = result; 29 } 30 ```` </pre>	<b>- Using memory hierarchy (shared, constant, registers)</b> <pre> 1  ````cuda 2  __constant__ float c_depthwise_weight[3*3*3];    // Constant memory for weights 3  __constant__ float c_pointwise_weight[3*64]; 4 5  __global__ void    depthwise_separable_conv_kernel_optimized( 6  const float* input, float* output, /* other    parameters */) { 7 8  // Shared memory for input tile with padding 9  __shared__ float    shared_input[3][SHARED_MEM_HEIGHT] 10 [SHARED_MEM_STRIDE]; 11 12 // Collaborative loading of input data to    shared memory 13 // [shared memory loading code...] 14 __syncthreads(); 15 16 // Register caching for intermediate results 17 float depthwise_results[3]; // Store in    registers 18 19 // Compute using shared memory and constant    memory 20 for (int c = 0; c &lt; in_channels; ++c) { 21   float sum = 0.0f; 22   // Fully unrolled convolution using        shared memory 23   sum +=        shared_input[c][sm_y_base][sm_x_base]        * c_depthwise_weight[c*9 + 0]; 24   // [more unrolled operations...] 25   depthwise_results[c] = sum; // Store in        register 26 } 27 28 // Cache output values in registers 29 float output_cache[32]; 30 31 // Compute pointwise convolution using    registers and constant memory 32 for (int i = 0; i &lt; oc_limit; ++i) { 33   output_cache[i] = depthwise_results[0] *        c_pointwise_weight[i * 3 + 0] + 34     depthwise_results[1] *        c_pointwise_weight[i *        3 + 1] + 35     depthwise_results[2] *        c_pointwise_weight[i *        3 + 2]; 36 } 37 38 // Coalesced write to global memory 39 for (int i = 0; i &lt; oc_limit; ++i) { 40   output[output_idx] = output_cache[i]; 41 } 42 } 43 ```` </pre>
<b>Asynchronous Execution</b>  Asynchronous Execution in CUDA allows operations to be queued and executed concurrently on separate streams, enabling overlapping computation with memory transfers for improved GPU utilization.	<b>- Sequential execution</b> <pre> 1  ````Python 2  def forward(self, x): 3  # Operations execute in the default stream,    blocking sequentially 4  result = self.conv_transpose3d(x) 5  return result 6  ```` </pre>	<b>- Asynchronous execution with custom stream</b> <pre> 1  ````Python 2  def forward(self, x): 3  # Create dedicated compute stream 4  self.compute_stream =    torch.cuda.Stream(priority=-1) # High    priority stream 5 6  # Execute operations asynchronously in the    custom stream 7  with torch.cuda.stream(self.compute_stream): 8  result = self._optimized_cuda_forward(x,    x.dtype) 9 10 # Control returns immediately while    computation continues in background 11 return result 12 ```` </pre>

Table 11: (Part 2) Code snippets before and after optimizations.

Tech + Desc	Before optimization	After optimization
<b>Memory Access Optimization</b> Memory Access Optimization in CUDA improves performance by organizing data access patterns to maximize cache utilization and minimize memory latency through techniques like tiling, coalescing, and shared memory usage.	<b>- Naive matrix multiplication with poor memory access</b> <pre> 1  ```.cuda 2  // Before optimization - Naive matrix    multiplication with poor memory access 3  __global__ void matmul_naive(float* A, float*    B, float* C, int M, int N, int K) { 4      int row = blockIdx.y * blockDim.y +    threadIdx.y; 5      int col = blockIdx.x * blockDim.x +    threadIdx.x; 6 7      if (row &lt; M &amp;&amp; col &lt; N) { 8          float sum = 0.0f; 9          for (int k = 0; k &lt; K; ++k) { 10             sum += A[row * K + k] * B[col * K +    k]; 11         } 12         C[row * N + col] = sum; 13     } 14 } 15 ``` </pre>	<b>- Using shared memory tiling and register blocking</b> <pre> 1  ```.cuda 2  __global__ void matmul_optimized(float* A,    float* B, float* C, int M, int N, int K) { 3      // Block index and thread index 4      const int bx = blockIdx.x; 5      const int by = blockIdx.y; 6      const int tx = threadIdx.x; 7      const int ty = threadIdx.y; 8 9      // Output positions 10     const int row = by * 8 + ty; 11     const int col = bx * 32 + tx; 12 13     // Register accumulation 14     float sum00 = 0.0f, sum01 = 0.0f; 15     float sum10 = 0.0f, sum11 = 0.0f; 16 17     // Shared memory tiles with padding to    avoid bank conflicts 18     __shared__ float As[8][33]; 19     __shared__ float Bs[32][33]; 20 21     // Loop over tiles 22     for (int tile = 0; tile &lt; (K + 31) / 32;    ++tile) { 23         // Collaborative loading of tiles into    shared memory 24         if (row &lt; M &amp;&amp; tile * 32 + tx &lt; K) 25             As[ty][tx] = A[row * K + tile * 32 +    tx]; 26         else 27             As[ty][tx] = 0.0f; 28 29         if (col &lt; N &amp;&amp; tile * 32 + ty &lt; K) 30             Bs[ty][tx] = B[col * K + tile * 32 +    ty]; 31         else 32             Bs[ty][tx] = 0.0f; 33 34         __syncthreads(); 35 36         // Compute partial dot products using    shared memory 37         #pragma unroll 8 38         for (int k = 0; k &lt; 32; ++k) { 39             float a0 = As[ty][k]; 40             float a1 = As[ty + 4][k]; 41             float b0 = Bs[k][tx]; 42             float b1 = Bs[k][tx + 16]; 43 44             sum00 += a0 * b0; 45             sum01 += a0 * b1; 46             sum10 += a1 * b0; 47             sum11 += a1 * b1; 48         } 49         __syncthreads(); 50     } 51 } 52 53 // Write results to global memory 54 if (row &lt; M &amp;&amp; col &lt; N) C[row * N + col] =    sum00; 55 if (row &lt; M &amp;&amp; col + 16 &lt; N) C[row * N +    col + 16] = sum01; 56 if (row + 4 &lt; M &amp;&amp; col &lt; N) C[(row + 4) *    N +    col] = sum10; 57 if (row + 4 &lt; M &amp;&amp; col + 16 &lt; N) C[(row +    4) *    N + col + 16] = sum11; 58 } 59 ``` </pre>
<b>Operation Fusion</b> Operation Fusion combines multiple consecutive operations into a single optimized kernel to reduce memory transfers and improve computational efficiency on CUDA devices.	<b>- Separate operations</b> <pre> 1  ```.Python 2  def forward(self, x): 3      x = F.max_pool3d(x,    kernel_size=self.pool_kernel_size,    stride=self.pool_stride) 4      x = torch.softmax(x, dim=1) 5      x = x - self.subtract.view(1, -1, 1, 1, 1) 6      x = x * torch.sigmoid(x) 7      return torch.max(x, dim=1)[0] 8  ``` </pre>	<b>- Fused operations with JIT</b> <pre> 1  ```.Python 2  @torch.jit.script 3  def fused_post_process(x, subtract_view): 4      x = torch.softmax(x, dim=1) 5      x = x - subtract_view 6      x = x * torch.sigmoid(x) 7      return torch.max(x, dim=1)[0] 8 9  def forward(self, x): 10     x = F.max_pool3d(x,    kernel_size=self.pool_kernel_size,    stride=self.pool_stride) 11     return self.fused_post_process(x,    self.subtract.view(1, -1, 1, 1, 1)) 12 ``` </pre>

Table 12: (Part 3) Code snippets before and after optimizations.

Tech + Desc	Before optimization	After optimization
<b>Optimized Thread Block Configuration</b>  Optimized Thread Block Configuration involves carefully selecting grid and block dimensions for CUDA kernels to maximize parallelism, memory access efficiency, and computational throughput based on the hardware architecture and algorithm characteristics.	<b>- Basic thread block configuration</b> <pre> 1 '''Python 2 block_dim = (16, 16) # Simple square thread    block 3 grid_dim = (math.ceil(N / 16), math.ceil(M /    16)) 4 5 kernel(grid=grid_dim, block=block_dim,    args=[A.data_ptr(), B.data_ptr(),    C.data_ptr(), M, N, K]) 6 ''' </pre>	<b>- Carefully tuned thread block configuration</b> <pre> 1 '''Python 2 block_dim = (32, 8) # Rectangular block    optimized for matrix multiplication 3 grid_dim = (math.ceil(N / 32), math.ceil(M /    8)) 4 5 kernel(grid=grid_dim, block=block_dim,    args=[A.data_ptr(), B.data_ptr(),    C.data_ptr(), M, N, K]) 6 ''' </pre>
<b>Branchless Implementation</b>  Branchless implementation replaces conditional statements with mathematical operations to avoid branch divergence and improve GPU performance.	<b>- With branches</b> <pre> 1 '''cuda 2 if (val &gt; 1.0f) { 3     output = 1.0f; 4 } else if (val &lt; -1.0f) { 5     output = -1.0f; 6 } else { 7     output = val; 8 } 9 ''' </pre>	<b>- Branchless</b> <pre> 1 '''cuda 2 output = fmaxf(-1.0f, fminf(1.0f, val)); 3 ''' </pre>
<b>Shared Memory Usage</b>  Shared memory in CUDA allows threads within the same block to efficiently share data, reducing global memory accesses and improving performance for algorithms with data reuse patterns.	<b>- Each thread reads diagonal element from global memory</b> <pre> 1 '''cuda 2 __global__ void    diag_matmul_kernel_unoptimized(const    float* A, const float* B, float* C, int N,    int M) { 3     int row = blockIdx.y * blockDim.y +    threadIdx.y; 4     int col = blockIdx.x * blockDim.x +    threadIdx.x; 5 6     if (row &lt; N &amp;&amp; col &lt; M) { 7         // Each thread loads the same diagonal    element multiple times from global    memory 8         C[row * M + col] = A[row] * B[row * M +    col]; 9     } 10 } 11 ''' </pre>	<b>- Using shared memory to cache diagonal elements</b> <pre> 1 '''cuda 2 __global__ void    diag_matmul_kernel_optimized(const float*    A, const float* B, float* C, int N, int M)    { 3     const int BLOCK_SIZE_Y = 8; 4     __shared__ float A_shared[BLOCK_SIZE_Y]; //    Shared memory for diagonal elements 5 6     int row = blockIdx.y * blockDim.y +    threadIdx.y; 7     int col = blockIdx.x * blockDim.x +    threadIdx.x; 8 9     // Load diagonal elements into shared    memory (once per row in block) 10    if (threadIdx.x == 0 &amp;&amp; row &lt; N) { 11        A_shared[threadIdx.y] = A[row]; 12    } 13 14    __syncthreads(); // Ensure all threads see    the loaded values 15 16    if (row &lt; N &amp;&amp; col &lt; M) { 17        // Use cached diagonal element from    shared memory 18        C[row * M + col] = A_shared[threadIdx.y]    * B[row * M + col]; 19    } 20 } 21 ''' </pre>
<b>Minimal Synchronization</b>  Minimal Synchronization reduces overhead by minimizing the number of synchronization points between CPU and GPU operations, allowing asynchronous execution through dedicated CUDA streams.	<b>- Default synchronization behavior</b> <pre> 1 '''Python 2 def forward(self, x): 3     # Each CUDA operation implicitly    synchronizes 4     x = x.contiguous() 5     result = self.conv_transpose3d(x) 6     return result 7 ''' </pre>	<b>- Minimal Synchronization</b> <pre> 1 '''Python 2 def forward(self, x): 3     # Create dedicated stream for computation 4     with torch.cuda.stream(self.compute_stream): 5         # Operations run asynchronously in this    stream 6         x_optimized = x.contiguous(memory_format=    torch.channels_last_3d) 7         result = 8             self.conv_transpose3d(x_optimized) 9         # Implicit synchronization only happens    when result is used 10        return result 11 ''' </pre>

Table 13: (Part 4) Code snippets before and after optimizations.

Tech + Desc	Before optimization	After optimization
<b>Thread Coarsening</b>  Thread Coarsening is an optimization technique where each thread processes multiple data elements instead of just one, increasing arithmetic intensity and reducing thread overhead.	<b>- Each thread processes one feature element</b>  <pre> 1 '''cuda 2 for (int d = tx; d &lt; feature_size; d +=   threads_x) { 3   scalar_t x_val = x[b + max_sample +   feature_size + n * feature_size + d]; 4   atomicAdd(&amp;vld[k + feature_size_padded +   d], assign_val * x_val); 5 } 6 ''' </pre>	<b>- Each thread processes two feature elements at once</b>  <pre> 1 '''cuda 2 #pragma unroll 4 3 for (int d = tx; d &lt; feature_size - 1; d +=   threads_x * 2) { 4   scalar_t x_val1 = x[b + max_sample +   feature_size + n * feature_size + d]; 5   scalar_t x_val2 = x[b + max_sample +   feature_size + n * feature_size + d +   threads_x]; 6 7   atomicAdd(&amp;vld[k + feature_size_padded +   d], assign_val * x_val1); 8   atomicAdd(&amp;vld[k + feature_size_padded +   d + threads_x], assign_val * x_val2); 9 } 10 11 // Handle remaining elements 12 for (int d = tx + (feature_size / threads_x) *   threads_x * 2; d &lt; feature_size; d +=   threads_x) { 13   scalar_t x_val = x[b + max_sample +   feature_size + n * feature_size + d]; 14   atomicAdd(&amp;vld[k + feature_size_padded +   d], assign_val * x_val); 15 } 16 ''' </pre>
<b>Asynchronous Execution</b>  Asynchronous Execution in CUDA allows operations to be queued and executed concurrently on separate streams, enabling overlapping computation with memory transfers for improved GPU utilization.	<b>- Sequential execution</b>  <pre> 1 '''Python 2 def forward(self, x): 3     # Operations execute in the default stream, 4     # blocking sequentially 5     result = self.conv_transpose3d(x) 6     return result 7 ''' </pre>	<b>- Asynchronous execution with custom stream</b>  <pre> 1 '''Python 2 def forward(self, x): 3     # Create dedicated compute stream 4     self.compute_stream = 5     torch.cuda.Stream(priority=-1) # High 6     # Execute operations asynchronously in the 7     # custom stream 8     with torch.cuda.stream(self.compute_stream): 9         result = self._optimized_cuda_forward(x, 10         x.dtype) 11     # Control returns immediately while 12     # computation continues in background 13     return result 14 ''' </pre>

Table 14: (Part 5) Code snippets before and after optimizations.

## F CASE STUDY: COMPARING REFERENCE CODE AND CUDA-L1 OPTIMIZED NEURAL NETWORK IMPLEMENTATIONS

### F.1 LSTMs

Table 15: Reference code and CUDA-L1 generation for LSTM class

#### LSTM | Reference Code - Simple baseline implementation

```

1  import torch
2  import torch.nn as nn
3
4  class Model(nn.Module):
5      def __init__(self, input_size, hidden_size, num_layers,
6          output_size, dropout=0.0):
7          """
8          Initialize the LSTM model.
9          """
10         super(Model, self).__init__()
11         # Initialize hidden state with random values
12         self.h0 = torch.randn((num_layers, batch_size,
13             hidden_size))
14         self.c0 = torch.randn((num_layers, batch_size,
15             hidden_size))
16         self.lstm = nn.LSTM(input_size, hidden_size,
17             num_layers, batch_first=True, dropout=dropout,
18             bidirectional=False)
19         self.fc = nn.Linear(hidden_size, output_size)
20
21     def forward(self, x):
22         """
23         Forward pass through the LSTM model.
24         """
25         self.h0 = self.h0.to(x.device)
26         self.c0 = self.c0.to(x.device) # BUG: This should be
27             self.c0.to(x.device)
28
29         # Forward propagate LSTM
30         out, state = self.lstm(x, (self.h0, self.c0)) # shape
31             of out: (batch_size, seq_length, hidden_size)
32
33         # Decode the hidden state of the last time step
34         out = self.fc(out[:, -1, :]) # shape of out:
35             (batch_size, output_size)
36
37         return state[0]
38
39 # Test code
40 batch_size = 10
41 sequence_length = 512
42 input_size = 128
43 hidden_size = 256
44 num_layers = 6
45 output_size = 10
46 dropout = 0.0
47
48 def get_inputs():
49     return [torch.randn(batch_size, sequence_length,
50         input_size)]
51
52 def get_init_inputs():
53     return [input_size, hidden_size, num_layers, output_size,
54         dropout]

```

**LSTM | Fully Optimized Code - All optimizations enabled (3.4x faster)**

```

1  import torch
2  import torch.nn as nn
3  import torch.cuda as cuda
4
5  class ModelNew(nn.Module):
6      def __init__(self, input_size, hidden_size, num_layers,
7                  output_size, dropout=0.0):
8          """
9              Initialize the LSTM model with three core optimization
10             techniques.
11
12             Color coding:
13             -  BLUE: CUDA Graphs optimization
14             -  GREEN: Memory Contiguity optimization
15             -  ORANGE: Static Tensor Reuse optimization
16             """
17             super(ModelNew, self).__init__()
18
19             # Initialize hidden states as buffers
20             self.register_buffer('h0', torch.randn((num_layers,
21             batch_size, hidden_size)))
22             self.register_buffer('c0', torch.randn((num_layers,
23             batch_size, hidden_size)))
24
25             # Use PyTorch's optimized LSTM implementation
26             self.lstm = nn.LSTM(
27                 input_size=input_size,
28                 hidden_size=hidden_size,
29                 num_layers=num_layers,
30                 batch_first=True,
31                 dropout=dropout,
32                 bidirectional=False
33             )
34
35             self.fc = nn.Linear(hidden_size, output_size)
36
37             #  CUDA GRAPHS: Variables for graph capture and
38             replay
39             self.graph = None
40             self.graph_ready = False
41             self.input_shape = None
42
43             #  STATIC TENSOR REUSE: Pre-allocated tensors for
44             graph execution
45             self.static_input = None
46             self.static_output = None
47
48             #  CUDA GRAPHS: Streams for graph operations
49             self.graph_stream = None
50
51             # Track if we're running on CUDA
52             self.is_cuda_available = torch.cuda.is_available()
53
54             def _initialize_cuda_resources(self):
55                 """ CUDA GRAPHS: Initialize CUDA stream for graph
56                 operations"""
57                 if self.graph_stream is None:
58                     self.graph_stream = cuda.Stream()
59

```

```

53     def _capture_graph(self, x, result):
54         """
55         ● CUDA GRAPHS: Capture the computation graph for
           replay
56         ● STATIC TENSOR REUSE: Create static tensors for
           graph capture
57         """
58         # ● STATIC TENSOR REUSE: Clone tensors for static
           allocation
59         self.static_input = x.clone()
60         self.static_output = result.clone()
61
62         # ● CUDA GRAPHS: Capture the computation graph
63         with torch.cuda.stream(self.graph_stream):
64             self.graph = cuda.CUDAGraph()
65             with cuda.graph(self.graph):
66                 # Operations to capture in the graph
67                 static_out, _ = self.lstm(self.static_input,
68                                           (self.h0, self.c0))
69
70                 # ● MEMORY CONTIGUITY: Ensure contiguous
71                 # memory layout
72                 static_last = static_out[:, -1, :].contiguous()
73
74                 self.static_output.copy_(self.fc(static_last))
75
76                 # Wait for graph capture to complete
77                 torch.cuda.synchronize()
78
79                 # Mark graph as ready for use
80                 self.graph_ready = True
81
82     def _standard_forward(self, x):
83         """Standard forward pass with memory contiguity
84         optimization"""
85
86         # ● MEMORY CONTIGUITY: Ensure input is contiguous
87         if not x.is_contiguous():
88             x = x.contiguous()
89
90         # Forward pass through LSTM
91         out, _ = self.lstm(x, (self.h0, self.c0))
92
93         # ● MEMORY CONTIGUITY: Make last output contiguous
94         # for optimal memory access
95         last_out = out[:, -1, :].contiguous()
96
97         return self.fc(last_out)
98
99     def forward(self, x):
100         """
101         Forward pass through the LSTM model with three
           optimization techniques.
102
103         Optimization flow:
104         1. ● CUDA GRAPHS: Check if we can use the captured
           graph (fast path)
105         2. ● STATIC TENSOR REUSE: Use pre-allocated tensors
           for graph replay

```

```

102     3. ● MEMORY CONTIGUITY: Ensure optimal memory layout
103         throughout
104     """
105     # ● CUDA GRAPHS: Fast path - use captured graph if
106         available
107     if (x.is_cuda and self.graph_ready and x.shape ==
108         self.input_shape):
109
110         # ● STATIC TENSOR REUSE: Copy to pre-allocated
111             tensor with non-blocking transfer
112         self.static_input.copy_(x, non_blocking=True)
113
114         # ● CUDA GRAPHS: Replay the captured graph
115         self.graph.replay()
116
117         # Return the output from static buffer
118         return self.static_output.clone()
119
120     # Standard execution path
121     with torch.no_grad():
122         result = self._standard_forward(x)
123
124     # ● CUDA GRAPHS: Initialize graph on first CUDA
125         input
126     if x.is_cuda and self.is_cuda_available and not
127         self.graph_ready:
128         try:
129             # Store the current input shape
130             self.input_shape = x.shape
131
132             # ● CUDA GRAPHS: Initialize CUDA resources
133             self._initialize_cuda_resources()
134
135             # ● CUDA GRAPHS + ● STATIC TENSOR REUSE:
136                 Capture the graph
137             self._capture_graph(x, result)
138
139             except Exception as e:
140                 # If graph capture fails, continue without
141                 it
142                 self.graph_ready = False
143
144         return result
145
146     # Hyperparameters from the reference implementation
147     batch_size = 10
148     sequence_length = 512
149     input_size = 128
150     hidden_size = 256
151     num_layers = 6
152     output_size = 10
153     dropout = 0.0
154
155     def get_inputs():
156         return [torch.randn(batch_size, sequence_length,
157                             input_size)]
158
159     def get_init_inputs():
160         return [input_size, hidden_size, num_layers, output_size,
161                 dropout]

```

```

153
154 # Example usage demonstrating the three techniques
155 if __name__ == "__main__":
156     import time
157
158     print("● BLUE: CUDA Graphs optimization")
159     print("● GREEN: Memory Contiguity optimization")
160     print("● ORANGE: Static Tensor Reuse optimization")
161     print("=" * 60)
162
163     # Create model
164     model = ModelNew(*get_init_inputs())
165     model.eval()
166
167     # Test input
168     x = get_inputs()[0]
169
170     # Move to GPU if available
171     if torch.cuda.is_available():
172         model = model.cuda()
173         x = x.cuda()
174
175     print("Running on CUDA - all three optimizations
176         active")
177
178     # First run - captures graph
179     print("\n● First forward pass: Capturing CUDA
180         graph...")
181     with torch.no_grad():
182         output = model(x)
183     print(f"    Output shape: {output.shape}")
184     print(f"    Graph ready: {model.graph_ready}")
185
186     # Subsequent runs - uses captured graph
187     print("\n● Subsequent passes: Using captured graph
188         with")
189     print("● static tensor reuse and ● memory
190         contiguity")
191
192     # Warmup
193     for _ in range(10):
194         with torch.no_grad():
195             _ = model(x)
196
197     # Measure performance
198     torch.cuda.synchronize()
199     start_event = torch.cuda.Event(enable_timing=True)
200     end_event = torch.cuda.Event(enable_timing=True)
201
202     n_runs = 100
203     start_event.record()
204     with torch.no_grad():
205         for _ in range(n_runs):
206             output = model(x)
207     end_event.record()
208
209     torch.cuda.synchronize()
210     avg_time = start_event.elapsed_time(end_event) / n_runs
211
212     print(f"\nPerformance: {avg_time:.3f} ms per forward
213         pass")

```

```
209     print(f"    Expected speedup: ~3.42x with all
210           optimizations")
211     else:
212     print("\n⚠ Running on CPU - only 🟢 memory
213           contiguity active")
214     print("    (CUDA graphs and static tensor reuse require
215           GPU)")
216     with torch.no_grad():
217         output = model(x)
218     print(f"\n    Output shape: {output.shape}")
```

## F.2 3DCONV

Table 16: Reference code and CUDA-L1 generation for Conv3D class

**Conv3D | Reference Code - Simple baseline implementation**

```

1  import torch
2  import torch.nn as nn
3
4  class Model(nn.Module):
5      """
6      Model that performs a 3D convolution, applies Group
7      Normalization, minimum, clamp, and dropout.
8      """
9      def __init__(self, in_channels, out_channels, kernel_size,
10                 groups, min_value, max_value, dropout_p):
11         super(Model, self).__init__()
12         self.conv = nn.Conv3d(in_channels, out_channels,
13                               kernel_size)
14         self.norm = nn.GroupNorm(groups, out_channels)
15         self.dropout = nn.Dropout(dropout_p)
16         self.min_value = min_value
17         self.max_value = max_value
18
19     def forward(self, x):
20         x = self.conv(x)
21         x = self.norm(x)
22         x = torch.min(x, torch.tensor(self.min_value))
23         x = torch.clamp(x, min=self.min_value,
24                        max=self.max_value)
25         x = self.dropout(x)
26         return x
27
28     # Hyperparameters
29     batch_size = 128
30     in_channels = 3
31     out_channels = 16
32     depth, height, width = 16, 32, 32
33     kernel_size = 3
34     groups = 8
35     min_value = 0.0
36     max_value = 1.0
37     dropout_p = 0.2
38
39     def get_inputs():
40         return [torch.randn(batch_size, in_channels, depth,
41                              height, width)]
42
43     def get_init_inputs():
44         return [in_channels, out_channels, kernel_size, groups,
45               min_value, max_value, dropout_p]

```

**Conv3D | Fully Optimized Code - All optimizations enabled (120x faster)**

```

1  import torch
2  import torch.nn as nn
3
4  # Hyperparameters
5  batch_size = 128
6  in_channels = 3
7  out_channels = 16
8  depth, height, width = 16, 32, 32

```

```

9 kernel_size = 3
10 groups = 8
11 min_value = 0.0
12 max_value = 1.0
13 dropout_p = 0.2
14
15 class ModelNew(nn.Module):
16     def __init__(self, in_channels, out_channels, kernel_size,
17                 groups, min_value, max_value, dropout_p):
18         super(ModelNew, self).__init__()
19         # Store the original layers for parameter compatibility
20         self.conv = nn.Conv3d(in_channels, out_channels,
21                               kernel_size)
22         self.norm = nn.GroupNorm(groups, out_channels)
23         self.dropout = nn.Dropout(dropout_p)
24         self.min_value = min_value
25         self.max_value = max_value
26         self.dropout_p = dropout_p
27
28         # ● TECH 1: Mathematical Short-Circuit Optimization
29         # Detects when min_value=0.0 to skip entire computation
30         self.use_optimized_path = (min_value == 0.0)
31
32         # ● TECH 4: Pre-computed Convolution Parameters
33         # Extract and store conv parameters once during
34         # initialization
35         if isinstance(kernel_size, int):
36             self.kernel_size = (kernel_size, kernel_size,
37                                 kernel_size)
38         else:
39             self.kernel_size = kernel_size
40             self.stride = self.conv.stride
41             self.padding = self.conv.padding
42             self.dilation = self.conv.dilation
43
44         # ● TECH 4: Pre-compute output dimensions for
45         # standard input
46         self.out_depth = ((depth + 2 * self.padding[0] -
47                             self.dilation[0] * (self.kernel_size[0] - 1) - 1)
48                             // self.stride[0]) + 1
49         self.out_height = ((height + 2 * self.padding[1] -
50                             self.dilation[1] * (self.kernel_size[1] - 1) - 1)
51                             // self.stride[1]) + 1
52         self.out_width = ((width + 2 * self.padding[2] -
53                             self.dilation[2] * (self.kernel_size[2] - 1) - 1)
54                             // self.stride[2]) + 1
55
56         # Standard output shape for the default batch size
57         self.standard_shape = (batch_size, out_channels,
58                                 self.out_depth, self.out_height, self.out_width)
59
60         # ● TECH 2: Pre-allocated Zero Tensors
61         # Create zero tensors once to avoid allocation overhead
62         if self.use_optimized_path:
63             self.register_buffer('zero_output_float32',
64                                 torch.zeros(self.standard_shape,
65                                              dtype=torch.float32),
66                                 persistent=False)
67             self.register_buffer('zero_output_float16',
68                                 torch.zeros(self.standard_shape,
69                                              dtype=torch.float16),
70                                 persistent=False)

```

```

57         self.register_buffer('zero_output_bfloat16',
58                               torch.zeros(self.standard_shape,
59                                             dtype=torch.bfloat16),
60                               persistent=False)
61     def calculate_output_shape(self, input_shape):
62         """Calculate the output shape of the convolution
63         operation."""
64         batch_size, _, d, h, w = input_shape
65
66         # 🟡 TECH 4: Use precomputed parameters
67         # Avoid repeated attribute lookups
68         out_d = ((d + 2 * self.padding[0] - self.dilation[0] *
69                  (self.kernel_size[0] - 1) - 1) // self.stride[0])
70                 + 1
71         out_h = ((h + 2 * self.padding[1] - self.dilation[1] *
72                  (self.kernel_size[1] - 1) - 1) // self.stride[1])
73                 + 1
74         out_w = ((w + 2 * self.padding[2] - self.dilation[2] *
75                  (self.kernel_size[2] - 1) - 1) // self.stride[2])
76                 + 1
77
78         return (batch_size, self.conv.out_channels, out_d,
79                 out_h, out_w)
80
81     def forward(self, x):
82
83         # 🔵 TECH 1: Mathematical Short-Circuit - Main
84         # optimization
85         # Skip all computation when we know result will be
86         # zeros
87         if not self.use_optimized_path:
88             # Standard path for non-optimized cases
89             x = self.conv(x)
90             x = self.norm(x)
91             x = torch.minimum(x, torch.tensor(self.min_value,
92                                               device=x.device))
93             x = torch.clamp(x, min=self.min_value,
94                             max=self.max_value)
95             x = self.dropout(x)
96             return x
97
98         # Optimized path when min_value == 0.0
99         # Since min(x, 0) followed by clamp(0, 1) always
100         # produces zeros
101
102         # 🟢 TECH 3: Direct Shape Matching
103         # Fast path for standard input dimensions
104         if x.shape == (batch_size, in_channels, depth, height,
105                       width):
106
107             # 🟣 TECH 2: Use pre-allocated tensors
108             # Return pre-allocated zeros matching input dtype
109             if x.dtype == torch.float32:
110                 return self.zero_output_float32
111             elif x.dtype == torch.float16:
112                 return self.zero_output_float16
113             elif x.dtype == torch.bfloat16:
114                 return self.zero_output_bfloat16
115             else:
116                 # Fallback for other dtypes
117                 return torch.zeros(self.standard_shape,
118                                   device=x.device, dtype=x.dtype)
119         else:

```

```
103         # For non-standard input shapes, calculate output
           shape
104         output_shape = self.calculate_output_shape(x.shape)
105         return torch.zeros(output_shape, device=x.device,
           dtype=x.dtype)
106
107     def get_inputs():
108         return [torch.randn(batch_size, in_channels, depth,
           height, width)]
109
110     def get_init_inputs():
111         return [in_channels, out_channels, kernel_size, groups,
           min_value, max_value, dropout_p]
112
113     # Color Legend:
114     # ● TECH 1: Mathematical Short-Circuit (Blue) - Skips
           computation when min_value=0
115     # ● TECH 2: Pre-allocated Tensors (Purple) - Pre-allocates
           zero tensors
116     # ● TECH 3: Direct Shape Matching (Green) - Fast path for
           standard shapes
117     # ● TECH 4: Pre-computed Parameters (Orange) - Pre-computes
           conv parameters
```