CUDA-L1: IMPROVING CUDA OPTIMIZATION VIA CONTRASTIVE REINFORCEMENT LEARNING

Anonymous authorsPaper under double-blind review

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.

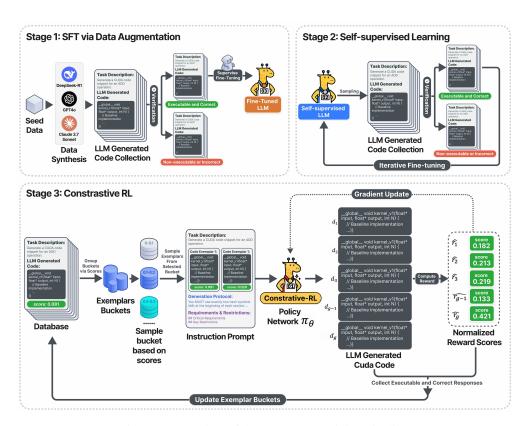


Figure 1: Overview of the CUDA-L1 training pipeline.

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 improving through distinguishing between effective and ineffective optimization strategies. Contrastive-RL simultaneously optimizes the foundation 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 coevolutionary 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 enhances the model's CUDA programming capabilities:

111 112 113

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

114 115

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.¹

116 117 118

2. Correctness: A CUDA code is correct if it produces equivalent outputs to the reference implementation across 1000 random test inputs.²

119

3. **Success**: A CUDA code is successful if it is both executable and correct.

120

2.2 SFT VIA DATA AUGMENTATION

121 122

123

124

125

126

130

134

135

136

137

138

139

140

144

145

146

147

148

149

150

151

152

153

154

155

156

157

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.

127 128 129

2.3 Self-supervised Learning

131 132 133

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.

141 142 143

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 RE-INFORCE (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 REIN-FORCE, 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 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

158 159 160

 $^{^{1}}$ This threshold is reasonable since code with $1000 \times$ slower performance contradicts our speedup optimization goals.

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

165

166

167

168

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:

170 171 172

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.

178 179

177

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).

181

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:

187 188 I) Task Descrpition: A detailed description of the computational problem, including input/output specifications, performance requirements, and optimization objectives.

189 190 191

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.

192

III) Generation Protocol: Explicit instructions defining the required output format and com-

193 194

IV) Requirements and Restrictions: Requirements and restrictions to prevent reward hacking in RL.

195 196

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

197 199

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.

200 201

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.

202 203

III) Code Implementation: The complete CUDA kernel implementation incorporating optimizations.

204

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

205 206

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 addresses 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.

212 213 214

215

211

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:

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

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 τ 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 μ_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, scoreweighted 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:

$$\mathbf{r}_{\text{single-run}}(d) = \frac{t_{q_i}}{t_d} \tag{2}$$

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.

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

2	7	0
2	7	1
2	7	2
2	7	3
2	7	4
2	7	5
2	7	6
2	7	7
2	7	8
2	7	9
2	8	0
2	8	1
2	8	2
2	8	3
2	8	4
2	8	5

Configuration	Method	Mean	Max	75%	50%	25%	Success↑ # out of total	Speedup ↑ >1.01x out of total
	All	3.12	120	2.25	1.42	1.17	249/250	226/250
D.C. L	Level 1	2.78	65.8	1.75	1.28	1.12	99/100	80/100
Default	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
	All	2.77	69.0	2.55	1.72	1.14	249/250	203/250
T 1 C '1	Level 1	3.04	59.7	2.71	1.99	1.41	99/100	89/100
Torch Compile	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
	All	2.88	80.1	2.48	1.67	1.13	249/250	200/250
Tl. C:1- PO	Level 1	3.38	55.3	3.02	2.29	1.61	99/100	90/100
Torch Compile RO	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
	All	2.81	97.9	1.83	1.20	0.954	249/250	147/229
CUDA C. I	Level 1	3.18	59.6	2.09	1.38	1.04	99/100	68/88
CUDA Graph	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.

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 justin-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, providing 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 im-

			l
Met	thods	Model	Mean
		Llama 3.1-405B	0.23
Van	:II.a	DeepSeek-V3	0.34
vun	ши	DeepSeek-R1	0.88
		OpenAI-O1	0.73
		Llama 3.1-405B	1.18
г	,	DeepSeek-V3	1.32
Evo	ive	DeepSeek-R1	1.41
		OpenAI-O1	1.35
		Stage 1	1.14
		Stage 1+2	1.36
CU	DA-L1	Stage 1+2+GRPO	2.41
COL	DA-LI	3 stages - random	2.14
		- island	3.21
		 bucket 	3.12

339 340 341

342

343

344

345

346

347

348

349

350

351

352

353 354

355

356

357

358

359

360

361

362

363

364

366

367

368 369

370

371

372

373

374

375

376

377

Table 2: Model performances on KernelBench All Level.

Max

3.14

2.96

14.4

12.4

18.4

52.4

44.2

63.9

32.7

48.3

84.6

64 5

126

75%

0.63

0.76

1.00

1.00

1.03

1.32

1.45

1.38

1.00

1.41

1.83

1.62

2.21

2.25

50%

0

0.75

0.55

1.00

1.03

1.17

1.16

1.00

1.09

1.33

1.21

1.40

1.42

25%

0

0

1.00

1.00

1.00

1.00

0.96

1.00

1.11

1.09

1.16

Success1

out of 250

99

179

141

247

247

248

247

240

247

247

241

249

249

Speedup1

>1.01 # out of 250

18

14

88

113

162

158

50

175

207

186

223

226

provements 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 2.77–2.88× mean speedup ratios, while demonstrating 2.81× mean improvement over CUDA Graph baseline with notable 97.9× 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 2.78–3.38× over different baselines. Level 2 (operator sequences) shows CUDA-L1's strongest performance with 3.55× improvement over Default baseline. Level 3 (complex ML tasks) reveals interesting baseline-dependent effectiveness: CUDA-L1 achieves 2.96× improvement over Default baseline, but shows reduced effectiveness against Torch compilation baselines (only 1.62-1.98× 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 OpenAIo1, 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 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
- 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

Configuration	GPU Device	Mean	Max	75%	50%	25%	Success ↑	Speedup 1
Comiguration	GI O DEVICE	Wican	IVIAA	13 /6	30 /0	23 /0	# out of 250	>1.01x
	A100	3.12	120	2.25	1.42	1.17	249	226 /250
	3090	2.51	114	1.57	1.18	1.03	242	201/250
Default	H100	3.85	368	1.76	1.32	1.09	250	218/250
	H20	2.38	63.7	1.81	1.34	1.11	247	226 /250
	L40	3.13	182	1.88	1.31	1.08	248	215/250
	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
Torch Compile	H100	2.74	49.7	2.83	1.92	1.11	250	195/250
	H20	2.89	49.4	3.21	2.04	1.19	247	209 /250
	L40	2.85	96.9	2.43	1.82	1.13	248	199/250
	A100	2.88	80.1	2.48	1.67	1.13	249	200/250
	3090	2.61	72.9	2.29	1.48	1.00	242	172/250
Torch Compile RO	H100	2.77	61.2	2.78	1.61	1.00	247	187/250
	H20	2.82	52.1	3.18	1.64	1.06	247	192/250
	L40	2.89	90.9	2.54	1.72	1.08	248	193/250
	A100	2.81	97.9	1.83	1.20	0.954	229	147/229
	3090	3.34	156	1.94	1.28	0.997	206	148/206
CUDA Graph	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	229	133/229
	L40	3.98	275	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.

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—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 Vũ, 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 Ilms write efficient gpu kernels? *arXiv preprint arXiv:2502.10517*, 2025.

- Bernardino Romera-Paredes, Mohammadamin Barekatain, 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.

DETAILS FOR CUDA-L1

594

595

596

597

598

600

601

607

608 609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624 625

626 627

630

631

632

633

634

635

636

637

638

639 640

641 642

643 644

645

646

647

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 PSUDO CODE FOR STAGE2: SELF-SUPERVISED LEARNING

Self-supervised Learning Algorithm

1: Initialize finetuned model M_0 after SFT stage with parameters $\theta_{\rm sft}$ 2: for i=1 to $N_{\text{iterations}}$ do 3: Generate batch of CUDA codes $C_i = \{c_1, ..., c_k\}$ using model M_{i-1}

```
4:
         Evaluate each c \in C_i for:
 5:
             1. Executability (compiles and runs)
             2. Correctness (produces expected output)
 6:
         Filter successful codes: C_i^{\text{success}} = \{c \in C_i | \text{executable} \land \text{correct}\}
 7:
         if C_i^{\text{success}} \neq \emptyset then
 8:
             Compute gradient update \nabla \theta using C_i^{\mathrm{success}}
 9:
10:
             Update model: \theta_i \leftarrow \theta_{i-1} + \eta \nabla \theta
11:
         else
            \theta_i \leftarrow \theta_{i-1} (no update)
12:
```

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 π_{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:

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

The complete GRPO objective optimizes the policy model by maximizing:

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

where:

- π_{θ} is the policy model being optimized
- $\pi_{\theta_{old}}$ is the old policy model from the previous iteration
- ε is the parameter for clipping
- β 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 $18\times$. 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)
```

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:

- 1. **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.
- 2. 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.
- 3. **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.
- 4. **Bucketized Variance Control**: We partition all Score_{single-run}(d) measurements into 7 buckets and compute bucket-wise averages. Evaluations with inter-bucket variance exceeding 0.005 are discarded.
- 5. **Robust Central Tendency**: The final reward uses the median of bucket averages, which proves more stable than the mean against outlier effects:

$$r(d) = median(\{Bucket_k\}_{k=1}^7)$$
(5)

804

808

- Conservative Rounding: We apply conservative rounding to speedup ratios (i.e., Score(d)), truncating to two decimal places while biasing toward unity (e.g., 1.118 → 1.11, 0.992 → 1.00).
- 7. **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:

$$\begin{split} r_{\text{normalized}} &= \frac{r - \mu}{\sigma} \\ r_{\text{smooth}} &= \text{clip}(r_{\text{normalized}}, -k, k) \end{split} \tag{6}$$

where μ and σ 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 $1.5\times$ 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.

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)

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.

C.1 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 $\mathbf{64}\times$ 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

31	6	ì	4
31	6	ì	5
3	6	ò	6
31	6	ì	7
3	6		8
	31	36 36	36: 36: 36: 36:

873 874 875

877 878 879

882 883

880

885

887 889

890

895 897

901 902 903

899

900

905 906 907

908

904

909 910 911

912 914

915 916

Configuration	CUDA Graphs	Memory Contiguity	Static Tensor Reuse	Speedup	Bottleneck
CUDA + Memory + Static	✓	1	1	3.42×	LSTM computation
CUDA + Memory	1	1	х	2.96×	Memory allocation
CUDA + Static	1	×	1	2.84×	Memory layout
CUDA Only	1	×	х	2.77×	Memory overhead
Memory + Static	×	1	1	1.00×	Kernel launch overhead
Memory Only	х	1	х	1.00×	Kernel launch overhead
Static Only	х	×	1	1.00×	Kernel launch overhead
Baseline	х	×	×	1.00×	Kernel launch overhead

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

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× 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 recording 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.

3D TRANSPOSED CONVOLUTION: 120× FASTER C.3

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.

Configuration	Math Short-Circuit	Pre-allocated Tensors	Direct Shape Match	Pre-computed Params	Speedup
Math + PreAlloc + Shape + Params	✓	✓	✓	1	120.9×
Math + Shape + Params	✓	×	✓	1	32.8×
Math + PreAlloc + Params	1	1	×	1	30.6×
Math + Params	✓	×	×	1	30.2×
Math + PreAlloc	✓	1	×	×	29.2×
Math Only	✓	×	×	×	29.2×
Math + Shape	1	х	1	х	28.6×
Math + PreAlloc + Shape	1	1	1	х	28.6×
PreAlloc + Shape + Params	×	1	1	1	1.0×
PreAlloc + Shape	x	1	1	х	1.0×
Shape + Params	×	х	1	1	1.0×
PreAlloc + Params	×	1	х	1	1.0×
Params Only	×	×	×	1	1.0×
Shape Only	х	×	1	×	1.0×
PreAlloc Only	×	1	×	х	1.0×
Baseline	×	×	×	×	1.0×

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

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 improvements (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 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

Task for CUDA Optimization

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:

- Produce **identical** results to the reference PyTorch implementation.
- Demonstrate **speed improvements** on GPU.
- Maintain stability for large input values.

Reference Implementation (exact copy)

```
1 import torch
2 import torch.nn as nn
3
   class Model(nn.Module):
5
       Simple model that performs a Softmax activation.
6
7
       def __init__(self):
9
           super(Model, self).__init__()
10
       def forward(self, x: torch.Tensor) -> torch.Tensor:
11
12
           Applies Softmax activation to the input tensor.
13
14
               x (torch.Tensor): Input tensor of shape
15
                    (batch_size, num_features).
16
               torch. Tensor: Output tensor with Softmax
17
                   applied, same shape as input.
18
19
           return torch.softmax(x, dim=1)
20
   batch\_size = 16
21
   dim = 16384
22
23
   def get_inputs():
25
       x = torch.randn(batch_size, dim)
       return [x]
26
27
   def get_init_inputs():
28
29
       return [] # No special initialization inputs needed
```

Table 8: Prompt illustration for data augmentation in Section ??. For each KernelBench task (soft-max shown here for illustration), the prompt is fed to each of six LLM models—GPT-40, OpenAI-01, 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

Task for CUDA Optimization

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.

Reference Code

Previous Cuda Implementations with Scores

Generation Protocol

You MUST use exactly two hash symbols (##) at the beginning of each section.

- **## Performance Analysis**: Compare code snippets above and articulate on :
 - 1. Which implementations demonstrate superior performance and why?
 - 2. What particular optimization strategies exhibit the greatest potential for improvement?
 - 3. What are the primary performance limitations in the implementation?
 - 4. What CUDA-specific optimization techniques remain unexploited?
 - 5. Where do the most significant acceleration opportunities exist?
- ## Algorithm Design: Describe your optimization approach
- ## Code Implementation: Provide your improved CUDA kernel

Requirements and Restrictions

Critical Requirements:

- 1. Functionality must match the reference implementation exactly. Failure to do so will result in a score of 0.
- 2. Code must compile and run properly on modern NVIDIA GPUs

Key Restrictions:

- Do not cache or reuse previous results the code must execute fully on each run.
- 2. Keep hyperparameters unchanged (e.g., batch size, dimensions, etc.) as specified in the reference.

E CASE STUDY: CODE SNIPPETS BEFORE AND AFTER OPTIMIZATIONS

Tech + Desc	Before optimization	After optimization
Memory Layout Optimization	- Non-contiguous memory access	- Ensuring contiguous memory layout
-	<pre>! ``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>	<pre>1 ``Python 2 def matrix_multiply_optimized(A, B): 3 # Ensure contiguous memory layout for efficient access patterns 4 A = A.contiguous() if not A.is_contiguous() else A 5 B = B.contiguous() if not B.is_contiguous() else B B.contiguous() return C 6 C = torch.mm(A, B) 7 return C 8 ```</pre>
computations. Memory Coa-		
escing	- Uncoalesced memory access	- Coalesced memory access with loop unrolling
Memory coa- lescing optimizes GPU memory access by ensur- ing threads in a warp access contiguous mem- ory locations, reducing memory transactions and increasing band- width utilization.	<pre>2global void uncoalesced_kernel(float* input, float* output) { int tid = threadIdx.x; 4 int stride = blockDim.x; 5 6 // Each thread accesses non-contiguous memory locations 7 for (int i = 0; i < 1024; i++) { output[tid + i * stride] = input[tid + i</pre>	<pre>2global void coalesced_kernel(float* input, float* output)</pre>
		<pre>batch_output[i+1] = batch_input[i+1] * 2.0f;</pre>
		<pre>batch_output[i+2] = batch_input[i+2] * 2.0f;</pre>
		<pre>16</pre>
		2.0f; 18 }
		19 }
Warp-Level	- Each thread independently calculates min value	- Using warp-level operations for parallel reduction
Optimizations Warp-Level Optimizations leverage the CUDA execution model where threads execute in groups of 32 (warps) to improve par- allel efficiency through collabo- trative operations and memory access patterns.	11 }	<pre>1 '''cuda 2 3global void min_kernel_after(const float* input, float* output, int size) { 4 int idx = blockIdx.x * blockDim.x + threadIdx.x; 5 int lane_id = threadIdx.x * 32; // Thread's position within warp 6 int warp_id = threadIdx.x / 32; // Warp number within the block 7 8 float min_val = lel0f; 9 if (idx < size) { 10</pre>

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

```
Tech + Desc
                                                      Before optimization
                                                                                                                                            After optimization
               Memory Hierarchy Opti-
1134
                                                      - Using global memory directly
                                                                                                                                           \hbox{-} \ Using \ memory \ hierarchy \ (shared, constant, registers)\\
               mization
1135
                                                     ı '''cuda
                                                                                                                                           2 __constant__ float c_depthwise_weight[3*3*3];
1136
               Memory Hierarchy Opti-
                                                                                                                                           // Constant memory for weights;
3 __constant__ float o_pointwise_weight[3*64];
                                                             depthwise_separable_conv_kernel_unoptimized(
               mization involves strate-
gically utilizing different
levels of GPU memory
(registers, shared memory,
                                                           const float* input, const float*
  depthwise_weight, const float*
1137
                                                                                                                                          s _global__void
  depthwise_separable_conv_kernel_optimized(
  const float* input, float* output, /* other
    parameters */) {
                                                            pointwise_weight,
float* output, /* other parameters */) {
1138
1139
                constant memory) to mini-
                                                            int out_y = blockIdx.y * blockDim.y +
                mize global memory access
                                                           threadIdx.y;
int out_x = blockIdx.x * blockDim.x +
threadIdx.x;
1140
               latency and maximize data
                                                                                                                                                 // Shared memory for input tile with padding
                                                                                                                                                __shared__ float
    shared_input[3][SHARED_MEM_HEIGHT]
[SHARED_MEM_STRIDE];
1141
                                                           // Each thread directly accesses global
  memory for each computation
for (int oc = 0; oc < out_channels; oc++) {
  float result = 0.0f;
  for (int ic = 0; ic < in_channels; ic++)</pre>
1142
                                                                                                                                                 // Collaborative loading of input data to
1143
                                                                                                                                                 shared memory
// [shared memory loading code...]
__syncthreads();
1144
                                                                     float depthwise_result = 0.0f;
1145
                                                                                                                                                 // Register caching for intermediate results
float depthwise_results[3]; // Store in
    registers
1146
                                                                                                                                                memory
for (int c = 0; c < in_channels; ++c) {
  float sum = 0.0f;
  // Fully unrolled convolution using
  shared memory
  sum +=
                                                                                                                                                 // Compute using shared memory and constant
1148
                                                    18
1149
1150
                                                                                                                                                     shared_input[c][sm_y_base][sm_x_base]
    * c_depthwise_weight[c*9 + 0];
    // [more unrolled operations...]
depthwise_results[c] = sum; // Store in
1151
1152
1153
                                                                                                                  * 9 + ky *
3 + kx];
1154
                                                                                                                                                 // Cache output values in registers
                                                                                                                                                 float output_cache[32];
1155
                                                                                                                                                 // Compute pointwise convolution using
1156
                                                                                                                                                 registers and constant memory
for (int i = 0; i < oc_limit; ++i) {
  output_cache[i] = depthwise_results[0] *
    c_pointwise_weight[i * 3 + 0] +
    depthwise_results[1] *
                                                                      result += depthwise_result >
                                                                          pointwise_weight[oc * in_channels
+ ic];
1157
1158
                                                                 output[((batch_idx * out_channels + oc)
    * out_height + out_y) * out_width +
    out_x] = result;
                                                                                                                                                                           c_pointwise_weight[i *
3 + 1] +
depthwise_results[2] *
1159
1160
                                                                                                                                                                                  c_pointwise_weight[i *
3 + 2];
1161
1162
                                                                                                                                                 // Coalesced write to global memory
for (int i = 0; i < oc_limit; ++i) {
   output[output_idx] = output_cache[i];</pre>
1163
1164
1165
               Asynchronous Execution
1166
                                                      - Sequential execution
                                                                                                                                           - Asynchronous execution with custom stream
                                                     ' ' 'Python 2 def forward(self, x):
                Asynchronous
                                        Execu-
1167
                tion in CUDA allows
                                                                                                                                           2 def forward(self, x):
                                                       # Operations execute in the default stream,
blocking sequentially
result = self.conv_transpose3d(x)
return result
                                                                                                                                                 # Create dedicated compute stream
               operations to be queued and executed concurrently on
1168
                                                                                                                                                 # Greater destrocked compute Stream
self.compute_stream =
    torch.cuda.Stream(priority=-1) # High
    priority stream
               separate streams, enabling
overlapping computation
with memory transfers for
1169
1170
                                                                                                                                                 # Execute operations asynchronously in the
    custom stream
with torch.cuda.stream(self.compute_stream):
                improved GPU utilization.
1171
                                                                                                                                                     result = self._optimized_cuda_forward(x,
1172
                                                                                                                                                           x.dtype)
1173
                                                                                                                                            computation continues in background return result
1174
1175
1176
```

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

1238

```
1190
1191
                  Tech + Desc
                                                                 Before optimization
                                                                                                                                                                       After optimization
                 Memory Access Optimization
1192
                                                                 - Naive matrix multiplication with poor memory access
                                                                                                                                                                       - Using shared memory tiling and register blocking
                 Memory Access Optimiza-
tion in CUDA improves
performance by organizing data
1193
                                                                                                                                                                     2 // Before optimization - Naive matrix
1194
                                                               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 +
                  access patterns to maximize
1195
                  cache utilization and minimize
                 memory latency through tech-
niques like tiling, coalescing,
                                                                       threadIdx.y;
int col = blockIdx.x * blockDim.x +
threadIdx.x;
1196
                  and shared memory usage.
                                                                                                                                                                              const int ty = threadIdx.y;
1197
                                                                      if (row < M && col < N) {
   float sum = 0.0f;
   for (int k = 0; k < K; ++k) {
      sum += A[row * K + k] * B[col * K +</pre>
1198
                                                                                                                                                                              const int row = by \star 8 + ty;
const int col = bx \star 32 + tx;
1199
                                                                                                                                                                             // Register accumulation
float sum00 = 0.0f, sum01 = 0.0f;
float sum10 = 0.0f, sum11 = 0.0f;
1200
                                                                                         k];
1201
                                                                             C[row * N + col] = sum;
                                                                                                                                                                             // Shared memory tiles with padding to
  avoid bank conflicts
  _shared__ float As[8][33];
  _shared__ float Bs[32][33];
1202
1203
                                                                                                                                                                             // Loop over tiles
for (int tile = 0; tile < (K + 31) / 32;
    ++tile) {
    // Collaborative loading of tiles into
        shared memory
    if (row < M && tile * 32 + tx < K)
        As[ty][tx] = A[row * K + tile * 32 +
        tx];
else</pre>
1204
                                                                                                                                                                    23
1206
                                                                                                                                                                    24
25
1207
                                                                                                                                                                                   else
As[ty][tx] = 0.0f;
1208
1209
                                                                                                                                                                                   if (col < N && tile * 32 + ty < K)
Bs[ty][tx] = B[col * K + tile * 32 +
1210
                                                                                                                                                                                                ty];
1211
                                                                                                                                                                    31
32
33
34
                                                                                                                                                                                        Bs[ty][tx] = 0.0f;
1212
                                                                                                                                                                                    __syncthreads();
1213
                                                                                                                                                                                   // Compute partial dot products using shared memory #pragma unroll 8
1214
                                                                                                                                                                                    #pragma unroll 8
for (int k = 0; k < 32; ++k) {
  float a0 = As[ty][k];
  float a1 = As[ty + 4][k];
  float b0 = Bs[k][tx];
  float b1 = Bs[k][tx + 16];</pre>
1215
1216
1217
                                                                                                                                                                                         sum00 += a0 * b0;
sum01 += a0 * b1;
sum10 += a1 * b0;
sum11 += a1 * b1;
1218
1219
1221
                                                                                                                                                                                    __syncthreads();
1222
                                                                                                                                                                             // Write results to global memory
if (row < M && col < N) C[row * N + col] =
    sum00;
if (row < M && col + 16 < N) C[row * N +
    col + 16] = sum01;
if (row + 4 < M && col < N) C[(row + 4) * N
    + col] = sum10;
if (row + 4 < M && col < N) C[(row + 4) * N
    + col] = sum10;
1223
1224
1225
1226
                                                                                                                                                                    57
1227
1228
                  Operation Fusion
                                                               - Separate operations
                                                                                                                                                                      - Fused operations with JIT
1229
                                                              Operation Fusion combines
                 Operation Fusion combines multiple consecutive operations into a single optimized kernel to reduce memory transfers and improve computational efficiency on CUDA devices.
1230
                                                                                                                                                                      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)
1231
1232
                                                                                                                                                                            return torch.max(x, dim=1)[0]
1233
                                                                                                                                                                      9 def forward(self, x):
0 x = F.max_pool3d(x,
1234
                                                                                                                                                                                   kernel_size=self.pool_kernel_size,
1235
                                                                                                                                                                                   stride=self.pool_stride)
turn self.fused_post_process(x,
self.subtract.view(1, -1, 1, 1, 1))
1236
1237
```

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

12451246

1289

```
1247
1248
               Tech + Desc
                                                     Before optimization
                                                                                                                                         After optimization
1249
               Optimized Thread Block
Configuration
                                                     - Basic thread block configuration
                                                                                                                                        - Carefully tuned thread block configuration
1250
               Optimized Thread Block
Configuration involves care-
fully selecting grid and block
dimensions for CUDA ker-
1251
                                                   2 block_dim = (16, 16) # Simple square thread
                                                                                                                                       2 block_dim = (32, 8) # Rectangular block
                                                                                                                                       optimized for matrix multiplication grid_dim = (math.ceil(N / 32), math.ceil(M /
                                                           block
1252
                                                   3 grid_dim = (math.ceil(N / 16), math.ceil(M /
16))
                                                   s kernel(grid=grid_dim, block=block_dim,
    args=[A.data_ptr(), B.data_ptr(),
    C.data_ptr(), M, N, K])
1253
               nels to maximize parallelism,
                                                                                                                                       memory access efficiency,
and computational through-
1254
               put based on the hardware architecture and algorithm
1255
                characteristics.
1256
              Branchless Implementation
                                                     - With branches
1257
               Branchless
                               implementation
                                                   '''cuda
2 if (val > 1.0f) {
                                                                                                                                        ı '''cuda
               replaces conditional statements
with mathematical operations
                                                                                                                                       2 output = fmaxf(-1.0f, fminf(1.0f, val));
1258
                                                      output = 1.0f;
} else if (val < -1.0f) {
output = -1.0f;
               to avoid branch divergence and
               improve GPU performance.
1260
                                                          output = val;
1261
1262
              Shared Memory Usage
                                                    - Each thread reads diagonal element from global memory
                                                                                                                                        - Using shared memory to cache diagonal elements
1263
               Shared memory in CUDA
                                                   1 '''cuda
2 __global
               allows threads within
                                           the
                                                                                                                                       2 __global
1264
               same block to efficiently share
data, reducing global memory
                                                         global__vold
diag_matmul_kernel_unoptimized(const
float* A, const float* B, float* C, int N,
int M) {
  int row = blockIdx.y * blockDim.y +
        threadIdx.y;
  int col = blockIdx.x * blockDim.x +
        threadIdx.x;
                                                                                                                                               diag_matmul_kernel_optimized(const float*
A, const float* B, float* C, int N, int M)
1265
               accesses and improving perfor-
                                                                                                                                             const int BLOCK_SIZE_Y = 8;
__shared__ float A_shared[BLOCK_SIZE_Y]; //
Shared memory for diagonal elements
               mance for algorithms with data
1266
1267
                                                                                                                                              int row = blockIdx.y * blockDim.y +
1268
                                                                                                                                             threadIdx.y;
int col = blockIdx.x * blockDim.x +
                                                          if (row < N && col < M) {
                                                               // Each thread loads the same diagonal
element multiple times from global
memory
C[row * M + col] = A[row] * B[row * M +
1269
                                                                                                                                                    threadIdx.x;
1270
                                                                                                                                             // Load diagonal elements into shared
                                                                                                                                             memory (once per row in block)
if (threadIdx.x == 0 && row < N) {
    A_shared[threadIdx.y] = A[row];
                                                                     col];
1271
1272
1273
                                                                                                                                              __syncthreads(); // Ensure all threads see
                                                                                                                                      15
                                                                                                                                             if (row < N && col < M) {
    // Use cached diagonal element from
        shared memory
    C[row * M + col] = A_shared[threadIdx.y]
        * B[row * M + col];</pre>
1275
1276
1277
1278
1279
               Minimal Synchronization
                                                     - Default synchronization behavior
                                                                                                                                        - Minimal Synchronization
1280
               Minimal Synchronization reduces overhead by minimizing
                                                    ı '''Python
                                                   def forward(self, x):

3  # Create dedicated stream for computation

4  with torch.cuda.stream(self.compute_stream):
1281
              the number of synchronization
points between CPU and GPU
operations, allowing asyn-
chronous execution through
dedicated CUDA streams.
1282
                                                                                                                                                  # Operations run asynchronously in this
1283
                                                         result = self.conv_transpose3d(x)
                                                      return result
                                                                                                                                                   1284
                                                                                                                                                   result =
                                                                                                                                             self.conv_transpose3d(x_optimized)
# Implicit synchronization only happens
when result is used
1285
1286
                                                                                                                                      when result
1287
1288
```

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

```
Tech + Desc
                                                                                                                   Before optimization
                                                                                                                                                                                                                                                                                                         After optimization
                                Thread Coarsening
                                                                                                                   - Each thread processes one feature element
                                                                                                                                                                                                                                                                                                        - Each thread processes two feature elements at once
                                Thread Coarsening is an optimization technique where each thread processes multiple
                                                                                                               data elements instead of just
                               one, increasing arithmetic intensity and reducing thread overhead.
                                                                                                                                                                                                                                                                                                                                 threads x];
1316
                                                                                                                                                                                                                                                                                                                   1318
1319
1320
                               Asynchronous Execution
                                                                                                                   - Sequential execution
                                                                                                                                                                                                                                                                                                        - Asynchronous execution with custom stream
1326
                                                                                                                  ı '''Python
                                                                                                                                                                                                                                                                                                       ı '''Python
                                                                                                              Asynchronous Execution in
                                                                                                                                                                                                                                                                                                      gython
gyth
                                CUDA allows operations to
be queued and executed con-
currently on separate streams,
enabling overlapping computa-
                                                                                                                                                                                                                                                                                                                                 torch.cuda.Stream(priority=-1) # High
                                                                                                                                                                                                                                                                                                                                 priority stream
                               tion with memory transfers for improved GPU utilization.
                                                                                                                                                                                                                                                                                                                     # Execute operations asynchronously in the
                                                                                                                                                                                                                                                                                                                    custom stream
with torch.cuda.stream(self.compute_stream):
    result = self._optimized_cuda_forward(x,
                                                                                                                                                                                                                                                                                                                                          x.dtype)
                                                                                                                                                                                                                                                                                                                    # Control returns immediately while computation continues in background
                                                                                                                                                                                                                                                                                                                      return result
```

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
1356
1357
        1 import torch
1358
         import torch.nn as nn
1359
1360
         4
            class Model(nn.Module):
1361
                def __init__(self, input_size, hidden_size, num_layers,
1362
                    output_size, dropout=0.0):
         6
1363
                     Initialize the LSTM model.
1364
1365
                     super(Model, self).__init__()
         9
1366
                     # Initialize hidden state with random values
         10
1367
                     self.h0 = torch.randn((num_layers, batch_size,
                        hidden_size))
1368
                     self.c0 = torch.randn((num_layers, batch_size,
         12
1369
                        hidden_size))
1370
                     self.lstm = nn.LSTM(input_size, hidden_size,
         13
1371
                        num_layers, batch_first=True, dropout=dropout,
1372
                        bidirectional=False)
                     self.fc = nn.Linear(hidden_size, output_size)
         14
1373
         15
1374
                def forward(self, x):
1375
         17
1376
                     Forward pass through the LSTM model.
         18
         19
1377
                     self.h0 = self.h0.to(x.device)
         20
1378
         21
                     self.c0 = self.h0.to(x.device) # BUG: This should be
1379
                        self.c0.to(x.device)
1380
         22
1381
                     # Forward propagate LSTM
                     out, state = self.lstm(x, (self.h0, self.c0)) # shape
1382
                        of out: (batch_size, seq_length, hidden_size)
1383
         25
1384
                     # Decode the hidden state of the last time step
         26
1385
                     out = self.fc(out[:, -1, :]) # shape of out:
1386
                         (batch_size, output_size)
1387
         28
         29
                     return state[0]
1388
1389
         31 # Test code
1390
         32 batch_size = 10
         sequence_length = 512
1391
            input_size = 128
1392
         35 hidden_size = 256
1393
         36 num_layers = 6
1394
         37 output_size = 10
1395
        38 	ext{ dropout} = 0.0
1396
        39
         40 def get_inputs():
1397
                return [torch.randn(batch_size, sequence_length,
         41
1398
                    input_size)]
1399
         42
1400
         43 def get_init_inputs():
1401
               return [input_size, hidden_size, num_layers, output_size,
         44
1402
                    dropout]
1403
```

```
1404
        LSTM | Fully Optimized Code - All optimizations enabled (3.4x faster)
1405
1406
            import torch
1407
            import torch.nn as nn
1408
            import torch.cuda as cuda
         4
1409
          5
            class ModelNew(nn.Module):
1410
                 def __init__(self, input_size, hidden_size, num_layers,
1411
                    output_size, dropout=0.0):
1412
1413
                     Initialize the LSTM model with three core optimization
          8
1414
                         techniques.
         9
1415
                     Color coding:
        10
1416
                     - BLUE: CUDA Graphs optimization
1417
                          GREEN: Memory Contiguity optimization
1418
1419
                          ORANGE: Static Tensor Reuse optimization
         14
1420
                     super(ModelNew, self).__init__()
         15
1421
         16
1422
                     # Initialize hidden states as buffers
         17
1423
                     self.register_buffer('h0', torch.randn((num_layers,
         18
1424
                         batch_size, hidden_size)))
                     self.register_buffer('c0', torch.randn((num_layers,
         19
1425
                         batch_size, hidden_size)))
1426
         20
1427
                     # Use PyTorch's optimized LSTM implementation
         21
1428
         22
                     self.lstm = nn.LSTM(
1429
                         input_size=input_size,
         23
                         hidden_size=hidden_size,
         24
1430
         25
                         num_layers=num_layers,
1431
                         batch_first=True,
1432
                         dropout=dropout,
         27
1433
                         bidirectional=False
         28
1434
         29
1435
         31
                     self.fc = nn.Linear(hidden_size, output_size)
1436
         32
1437
                     # CUDA GRAPHS: Variables for graph capture and
1438
                         replay
1439
                     self.graph = None
         34
1440
         35
                     self.graph_ready = False
1441
         36
                     self.input_shape = None
         37
1442
                     # STATIC TENSOR REUSE: Pre-allocated tensors for
1443
                         graph execution
1444
                     self.static_input = None
         39
1445
         40
                     self.static_output = None
1446
         41
1447
                     # CUDA GRAPHS: Streams for graph operations
         42
1448
                     self.graph_stream = None
         43
1449
         44
1450
                     # Track if we're running on CUDA
         45
                     self.is_cuda_available = torch.cuda.is_available()
1451
         46
         47
1452
                 def _initialize_cuda_resources(self):
         48
1453
                     """ CUDA GRAPHS: Initialize CUDA stream for graph
         49
1454
                         operations"""
1455
         50
                     if self.graph_stream is None:
1456
         51
                         self.graph_stream = cuda.Stream()
1457
         52
```

```
1458
         53
                 def _capture_graph(self, x, result):
1459
         54
1460
                        CUDA GRAPHS: Capture the computation graph for
1461
                         replay
1462
                        STATIC TENSOR REUSE: Create static tensors for
1463
                         graph capture
1464
1465
                       STATIC TENSOR REUSE: Clone tensors for static
1466
                         allocation
1467
                     self.static_input = x.clone()
                     self.static_output = result.clone()
1468
         60
1469
                     # CUDA GRAPHS: Capture the computation graph
1470
         62
                     with torch.cuda.stream(self.graph_stream):
         63
1471
                          self.graph = cuda.CUDAGraph()
         64
1472
                          with cuda.graph(self.graph):
         65
1473
                              # Operations to capture in the graph
         66
1474
                              static_out, _ = self.lstm(self.static_input,
                                  (self.h0, self.c0))
1475
         68
1476
                              # MEMORY CONTIGUITY: Ensure contiquous
1477
                                  memory layout
1478
                              static_last = static_out[:, -1, :].contiguous()
         70
1479
         71
1480
                              self.static_output.copy_(self.fc(static_last))
1481
                     # Wait for graph capture to complete
1482
         74
         75
                     torch.cuda.synchronize()
1483
         76
1484
         77
                     # Mark graph as ready for use
1485
                     self.graph_ready = True
         78
1486
         79
         80
                 def _standard_forward(self, x):
1487
                     """Standard forward pass with memory contiguity
         81
1488
                         optimization"""
1489
         82
1490
                       MEMORY CONTIGUITY: Ensure input is contiguous
         83
1491
                     if not x.is_contiguous():
         84
1492
         85
                          x = x.contiguous()
1493
         87
                     # Forward pass through LSTM
1494
                     out, \_ = self.lstm(x, (self.h0, self.c0))
         88
1495
         89
1496
                       MEMORY CONTIGUITY: Make last output contiguous
         90
1497
                         for optimal memory access
1498
                     last_out = out[:, -1, :].contiguous()
1499
         92
1500
         93
                     return self.fc(last_out)
1501
         94
                 def forward(self, x):
         95
1502
1503
                     Forward pass through the LSTM model with three
         97
1504
                         optimization techniques.
1505
         98
1506
         99
                     Optimization flow:
1507
                         CUDA GRAPHS: Check if we can use the captured
         100
                         graph (fast path)
1508
1509
                        STATIC TENSOR REUSE: Use pre-allocated tensors
        101
                         for graph replay
1510
```

```
1512
                         MEMORY CONTIGUITY: Ensure optimal memory layout
1513
                         throughout
1514
         103
1515
         104
1516
                          CUDA GRAPHS: Fast path - use captured graph if
1517
                         available
1518
                     if (x.is_cuda and self.graph_ready and x.shape ==
         106
1519
                         self.input_shape):
1520
         107
                          # STATIC TENSOR REUSE: Copy to pre-allocated
1521
                              tensor with non-blocking transfer
1522
                          self.static_input.copy_(x, non_blocking=True)
         109
1523
         110
1524
                          # O CUDA GRAPHS: Replay the captured graph
1525
                          self.graph.replay()
         112
1526
         113
1527
                          # Return the output from static buffer
         114
                          return self.static_output.clone()
1528
         116
1529
                     # Standard execution path
1530
         118
                     with torch.no_grad():
1531
                          result = self._standard_forward(x)
         119
1532
         120
1533
                             CUDA GRAPHS: Initialize graph on first CUDA
1534
                              input.
1535
                          if x.is cuda and self.is cuda available and not
                              self.graph_ready:
1536
         123
                              try:
1537
                                   # Store the current input shape
         124
1538
                                  self.input_shape = x.shape
         125
1539
        126
1540
                                  # CUDA GRAPHS: Initialize CUDA resources
         127
1541
         128
                                  self._initialize_cuda_resources()
1542
         129
1543
                                    CUDA GRAPHS + STATIC TENSOR REUSE:
         130
1544
                                      Capture the graph
         131
                                  self._capture_graph(x, result)
1545
1546
                              except Exception as e:
1547
                                  # If graph capture fails, continue without
         134
1548
                                  self.graph_ready = False
1549
        136
1550
        137
                          return result
1551
        138
1552
            # Hyperparameters from the reference implementation
         139
1553
        140 batch_size = 10
1554
         141
            sequence_length = 512
            input_size = 128
         142
1555
        143 hidden_size = 256
1556
        num_layers = 6
1557
        output_size = 10
1558
        146 dropout = 0.0
1559
        147
        148
            def get_inputs():
1560
                 return [torch.randn(batch_size, sequence_length,
        149
1561
                     input_size)]
1562
         150
1563
         151
            def get_init_inputs():
1564
                 return [input_size, hidden_size, num_layers, output_size,
         152
                    dropout]
1565
```

```
1566
         153
1567
         154
            # Example usage demonstrating the three techniques
1568
         155
             if __name__ == "__main__":
1569
                 import time
         156
1570
         157
                 print(" BLUE: CUDA Graphs optimization")
1571
         158
1572
                 print(" GREEN: Memory Contiguity optimization")
         159
1573
                 print(" ORANGE: Static Tensor Reuse optimization")
         160
1574
                 print("=" * 60)
         161
1575
         162
                 # Create model
         163
1576
                 model = ModelNew(*get_init_inputs())
         164
1577
         165
                 model.eval()
1578
         166
1579
                 # Test input
         167
1580
                 x = get_inputs()[0]
         168
         169
1581
                 # Move to GPU if available
         170
1582
                 if torch.cuda.is available():
1583
                     model = model.cuda()
1584
         173
                      x = x.cuda()
1585
         174
         175
                      print("Running on CUDA - all three optimizations
1586
                          active")
1587
         176
1588
                      # First run - captures graph
         177
1589
                      print("\n First forward pass: Capturing CUDA
1590
                          graph...")
1591
                      with torch.no_grad():
         179
1592
         180
                          output = model(x)
                      print(f" Output shape: {output.shape}")
         181
1593
                      print(f"
                                Graph ready: {model.graph_ready}")
         182
1594
         183
1595
                      # Subsequent runs - uses captured graph
         184
1596
                      print("\n Subsequent passes: Using captured graph
         185
1597
                          with")
1598
                      print(" static tensor reuse and memory
         186
1599
                          contiguity")
1600
         187
1601
         188
                      # Warmup
         189
                      for _ in range(10):
1602
                          with torch.no_grad():
         190
1603
                              _{-} = model(x)
         191
1604
         192
1605
         193
                      # Measure performance
1606
                      torch.cuda.synchronize()
         194
                      start_event = torch.cuda.Event(enable_timing=True)
         195
1607
         196
                      end_event = torch.cuda.Event(enable_timing=True)
1608
         197
1609
         198
                      n_runs = 100
1610
         199
                      start_event.record()
1611
                      with torch.no_grad():
         200
         201
                          for _ in range(n_runs):
1612
         202
                               output = model(x)
1613
                      end_event.record()
         203
1614
         204
1615
         205
                      torch.cuda.synchronize()
1616
                      avg_time = start_event.elapsed_time(end_event) / n_runs
         206
1617
         207
                      print(f"\nPerformance: {avg_time:.3f} ms per forward
         208
1618
                        pass")
1619
```

```
1620
                      print(f" Expected speedup: ~3.42x with all
1621
                          optimizations")
1622
         210
1623
         211
                 else:
1624
                      print("\n! Running on CPU - only  memory
1625
                          contiguity active")
                      print(" (CUDA graphs and static tensor reuse require
1626
         213
                          GPU)")
1627
         214
1628
         215
                      with torch.no_grad():
1629
                      output = model(x)
print(f"\n Output shape: {output.shape}")
         216
1630
         217
1631
1632
```

F.2 3DConv

1674

1675 1676

1677

1719

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

```
Conv3D | Reference Code - Simple baseline implementation
1678
1679
            import torch
1680
            import torch.nn as nn
1681
1682
            class Model(nn.Module):
1683
          5
                Model that performs a 3D convolution, applies Group
1684
                    Normalization, minimum, clamp, and dropout.
1685
1686
                def __init__(self, in_channels, out_channels, kernel_size,
          8
1687
                    groups, min_value, max_value, dropout_p):
1688
          9
                     super(Model, self).__init__()
                     self.conv = nn.Conv3d(in_channels, out_channels,
         10
1689
                        kernel_size)
1690
                     self.norm = nn.GroupNorm(groups, out_channels)
         11
1691
                     self.dropout = nn.Dropout(dropout_p)
1692
         13
                     self.min_value = min_value
1693
                     self.max_value = max_value
         14
         15
1694
         16
                def forward(self, x):
1695
                     x = self.conv(x)
         17
1696
                     x = self.norm(x)
         18
1697
                     x = torch.min(x, torch.tensor(self.min_value))
         19
1698
                     x = torch.clamp(x, min=self.min_value,
         20
                        max=self.max_value)
1699
         2.1
                     x = self.dropout(x)
1700
                     return x
         22
1701
         23
1702
            # Hyperparameters
1703
         25 batch_size = 128
         in_{channels} = 3
1704
         27 out_channels = 16
1705
         28 depth, height, width = 16, 32, 32
1706
         29 kernel_size = 3
1707
         30 groups = 8
         min_value = 0.0
1708
         max_value = 1.0
1709
         dropout_p = 0.2
1710
         34
1711
        35 def get_inputs():
1712
                return [torch.randn(batch_size, in_channels, depth,
         36
1713
                    height, width)]
         37
1714
         38 def get_init_inputs():
1715
                return [in_channels, out_channels, kernel_size, groups,
1716
                    min_value, max_value, dropout_p]
1717
1718
```

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

```
1720
1721
            import torch
1722
         2
            import torch.nn as nn
1723
         4 # Hyperparameters
1724
         5 batch_size = 128
1725
         6 	 in_channels = 3
1726
         7 out_channels = 16
1727
        8 depth, height, width = 16, 32, 32
```

```
1728
        9 kernel_size = 3
1729
         10 groups = 8
1730
         min_value = 0.0
1731
         12 max_value = 1.0
            dropout_p = 0.2
1732
         13
         14
1733
            class ModelNew(nn.Module):
         15
1734
                 def __init__(self, in_channels, out_channels, kernel_size,
         16
1735
                    groups, min_value, max_value, dropout_p):
1736
         17
                     super(ModelNew, self).__init__()
                     # Store the original layers for parameter compatibility
1737
         18
                     self.conv = nn.Conv3d(in_channels, out_channels,
         19
1738
                         kernel_size)
1739
                     self.norm = nn.GroupNorm(groups, out_channels)
1740
                     self.dropout = nn.Dropout(dropout_p)
1741
                     self.min_value = min_value
         22
                     self.max_value = max_value
1742
         23
         24
                     self.dropout_p = dropout_p
1743
1744
                         TECH 1: Mathematical Short-Circuit Optimization
         26
1745
                     # Detects when min_value=0.0 to skip entire computation
1746
                     self.use_optimized_path = (min_value == 0.0)
         28
1747
         29
1748
                     # TECH 4: Pre-computed Convolution Parameters
         30
1749
         31
                     # Extract and store conv parameters once during
1750
                         initialization
1751
                     if isinstance(kernel_size, int):
                         self.kernel_size = (kernel_size, kernel_size,
         33
1752
                             kernel_size)
1753
         34
                     else:
1754
                         self.kernel_size = kernel_size
1755
                     self.stride = self.conv.stride
1756
         37
                     self.padding = self.conv.padding
                     self.dilation = self.conv.dilation
         38
1757
         39
1758
                       TECH 4: Pre-compute output dimensions for
1759
                         standard input
1760
                     self.out_depth = ((depth + 2 * self.padding[0] -
1761
                         self.dilation[0] * (self.kernel_size[0] - 1) - 1)
1762
                         // self.stride[0]) + 1
                     self.out_height = ((height + 2 * self.padding[1] -
1763
         42
                         self.dilation[1] * (self.kernel_size[1] - 1) - 1)
1764
                         // self.stride[1]) + 1
1765
                     self.out\_width = ((width + 2 * self.padding[2] -
         43
1766
                         self.dilation[2] * (self.kernel_size[2] - 1) - 1)
1767
                         // self.stride[2]) + 1
1768
         45
                     # Standard output shape for the default batch size
1769
                     self.standard_shape = (batch_size, out_channels,
1770
                         self.out_depth, self.out_height, self.out_width)
1771
         47
1772
                     # TECH 2: Pre-allocated Zero Tensors
         48
1773
         49
                     # Create zero tensors once to avoid allocation overhead
1774
         50
                     if self.use_optimized_path:
                         self.register_buffer('zero_output_float32',
1775
         51
                                             torch.zeros(self.standard_shape,
         52
1776
                                                 dtype=torch.float32),
1777
                                             persistent=False)
         53
1778
         54
                         self.register_buffer('zero_output_float16',
1779
                                             torch.zeros(self.standard_shape,
         55
                                                 dtype=torch.float16),
1780
                                             persistent=False)
1781
```

```
1782
         57
                         self.register_buffer('zero_output_bfloat16',
1783
                                             torch.zeros(self.standard_shape,
         58
1784
                                                  dtype=torch.bfloat16),
1785
         59
                                             persistent=False)
1786
         61
                 def calculate_output_shape(self, input_shape):
1787
                     """Calculate the output shape of the convolution
         62
1788
                        operation."""
1789
                     batch_size, _, d, h, w = input_shape
1790
1791
                     # TECH 4: Use precomputed parameters
         65
1792
                     # Avoid repeated attribute lookups
                     out_d = ((d + 2 * self.padding[0] - self.dilation[0] *
1793
                         (self.kernel_size[0] - 1) - 1) // self.stride[0])
1794
                         + 1
1795
                     out_h = ((h + 2 * self.padding[1] - self.dilation[1] *
1796
                         (self.kernel_size[1] - 1) - 1) // self.stride[1])
1797
                         + 1
                     out_w = ((w + 2 * self.padding[2] - self.dilation[2] *
1798
         69
                         (self.kernel_size[2] - 1) - 1) // self.stride[2])
1799
                         + 1
1800
         70
1801
                     return (batch_size, self.conv.out_channels, out_d,
1802
                         out_h, out_w)
1803
         72
                 def forward(self, x):
1804
                       TECH 1: Mathematical Short-Circuit - Main
1805
                         optimization
1806
                     # Skip all computation when we know result will be
1807
                         zeros
1808
                     if not self.use_optimized_path:
1809
                         # Standard path for non-optimized cases
         77
1810
         78
                         x = self.conv(x)
                         x = self.norm(x)
         79
1811
                         x = torch.minimum(x, torch.tensor(self.min_value,
         80
1812
                             device=x.device))
1813
                         x = torch.clamp(x, min=self.min_value,
1814
                             max=self.max_value)
                         x = self.dropout(x)
1815
         82
                         return x
1816
         84
1817
                     # Optimized path when min_value == 0.0
         85
1818
                     \# Since min(x, 0) followed by clamp(0, 1) always
         86
1819
                         produces zeros
1820
                       TECH 3: Direct Shape Matching
1821
                     # Fast path for standard input dimensions
1822
                     if x.shape == (batch_size, in_channels, depth, height,
         90
1823
                         width):
1824
                           TECH 2: Use pre-allocated tensors
         91
1825
                         # Return pre-allocated zeros matching input dtype
         92
1826
                         if x.dtype == torch.float32:
         93
1827
                             return self.zero_output_float32
1828
         95
                         elif x.dtype == torch.float16:
                             return self.zero_output_float16
1829
         96
                         elif x.dtype == torch.bfloat16:
         97
1830
         98
                             return self.zero_output_bfloat16
1831
                         else:
         99
1832
         100
                              # Fallback for other dtypes
1833
                             return torch.zeros(self.standard_shape,
        101
                                 device=x.device, dtype=x.dtype)
1834
                     else:
1835
```

```
1836
                       # For non-standard input shapes, calculate output
1837
1838
                       output_shape = self.calculate_output_shape(x.shape)
       104
1839
                       return torch.zeros(output_shape, device=x.device,
       105
1840
                          dtype=x.dtype)
       106
1841
           def get_inputs():
       107
1842
               return [torch.randn(batch_size, in_channels, depth,
       108
1843
                  height, width)]
1844
       109
           def get_init_inputs():
1845
       110
               return [in_channels, out_channels, kernel_size, groups,
1846
                  min_value, max_value, dropout_p]
1847
1848
           # Color Legend:
       113
1849
              TECH 1: Mathematical Short-Circuit (Blue) - Skips
1850
               computation when min_value=0
1851
             TECH 2: Pre-allocated Tensors (Purple) - Pre-allocates
1852
               zero tensors
1853
             TECH 3: Direct Shape Matching (Green) - Fast path for
1854
              standard shapes
1855
       117 # TECH 4: Pre-computed Parameters (Orange) - Pre-computes
1856
              conv parameters
1857
1858
```