TRITONRL: TRAINING LLMS TO THINK AND CODE TRITON WITHOUT CHEATING

Anonymous authors

Paper under double-blind review

ABSTRACT

With the rapid evolution of large language models (LLMs), the demand for automated, high-performance system kernels has emerged as a key enabler for accelerating development and deployment. We introduce TRITONRL, a domainspecialized LLM for Triton kernel generation, trained with a novel reinforcement learning (RL) framework that enables robust and automated kernel synthesis. Unlike CUDA, which benefits from abundant programming data, high-performance Triton kernels are scarce and typically require costly crawling or manual authoring. Furthermore, reliable evaluation methods for validating Triton kernels remain underdeveloped and even hinder proper diagnosis of base model performance. Our approach addresses these challenges end-to-end with a fully opensource recipe: we curate datasets from KernelBook, enhance solution quality via DeepSeek-assisted distillation, and fine-tune Qwen3-8B to retain both reasoning ability and Triton-specific correctness. We further introduce hierarchical reward decomposition and data mixing to enhance RL training. With correct reevaluations of existing models, our experiments on KernelBench demonstrate that TRITONRL achieves state-of-the-art correctness and speedup, surpassing all other Triton-specific models and underscoring the effectiveness of our RL-based training paradigm.

1 Introduction

The exponential growth in demand for GPU computing resources has driven the need for highly optimized GPU kernels that improve computational efficiency, yet with the emergence of numerous GPU variants featuring diverse hardware specifications and the corresponding variety of optimization kernels required for each, developing optimized kernels has become an extremely time-consuming and challenging task. In response to this need, there is growing interest in leveraging large language models (LLMs) for automated kernel generation. While there have been attempts introducing inference frameworks that utilize general-purpose models, such as OpenAI models and DeepSeek, for generating kernels (Ouyang et al., 2025; Lange et al., 2025; Li et al., 2025a; NVIDIA Developer Blog, 2025), they often struggle with even basic kernel implementations, thereby highlighting the critical need for domain-specific models specifically tailored for kernel synthesis.

As the need for specialized models for kernel generation has emerged, several works have focused on fine-tuning LLMs for CUDA or Triton. In the CUDA domain, recent RL-based approaches include Kevin-32B Baronio et al. (2025), which progressively improves kernels using execution feedback as reward signals, and CUDA-L1, which applies contrastive RL to DeepSeek-V3 Li et al. (2025c). While these large models (32B-671B parameters) achieve strong CUDA performance, their training costs remain prohibitively expensive. To address these limitations, researchers have proposed specialized 8B Triton models, including KernelLLM Fisches et al. (2025) (supervised training on torch compiler-generated code) and AutoTriton Li et al. (2025b), fine-tuned via LLM distillation and RL using execution feedback. Though these smaller models outperform their base models, significant room remains for improving efficiency and accuracy compared to larger counterparts.

Furthermore, there is a common issue reported across kernel generation works, reward hacking (Baronio et al., 2025; Li et al., 2025b). Due to the scarcity of high-quality kernel examples compared to other programming languages, most approaches rely on RL training using runtime measurements and correctness rewards from unit tests after kernel execution. However, models frequently learn to

056

060

061

062 063 064

065

066

067

068

069

071

072

073

074 075

076

077

078

079

081

082

084

085

087

090

091

092

094

096

098

099

100

101

102

103

104 105

106

107

Figure 1: TRITONRL components and workflow.

exploit unit test loopholes, such as direct use of high-level PyTorch modules, rather than generating proper code, and this phenomenon is particularly prevalent in smaller models (8B and below) (Baronio et al., 2025). This issue fundamentally undermines the core objective of developing more efficient custom kernels to replace existing pre-optimized libraries, while current approaches predominantly rely on simple rule-based syntax verification whose effectiveness remains uncertain.

In this paper, we present TRITONRL, an 8B-scale LLM specialized for Triton programming that achieves state-of-the-art performance in both correctness and runtime speedup, while effectively mitigating reward hacking. To enable high-quality Triton kernel generation with small models (up to 8B), we design a training pipeline with the following key contributions:

- Simplified dataset curation with distillation: Instead of large-scale web crawling, we build on the curated *KernelBook* problems. Their solutions are refined and augmented through DeepSeek-R1 distillation (Guo et al., 2025), providing high-quality supervision for SFT of our base model Qwen3-8B (Team, 2025).
- Fine-grained and robust verification: We incorporate enhanced rule-based checks (e.g., nn.Module) together with LLM-based judges (e.g., Qwen3-235B-Instruct) to construct verifiable rewards. This enables reliable diagnosis across commercial and open-source models, while preventing reward hacking that arises from naive syntax-only verification.
- Hierarchical reward decomposition with data mixing: Our RL stage decomposes rewards into multiple dimensions (e.g., correctness, efficiency, style) and applies token-level credit assignment. Combined with strategic data mixing across SFT and RL, this yields better kernel quality, generalization, and robustness.
- Comprehensive evaluation and open-sourcing: Through rigorous validity analysis that filters
 out syntactically or functionally invalid code, we reveal true performance differences among
 models. Ablation studies further confirm the effectiveness of our hierarchical reward design and
 data mixing. At the 8B scale, TRITONRL surpasses existing Triton-specific LLMs, including
 KernelLLM (Fisches et al., 2025) and AutoTriton (Li et al., 2025b). We fully release our datasets,
 recipes, pretrained checkpoints, and evaluation framework to ensure reproducibility and foster
 future research.

2 TritonRL

In this section, we present TRITONRL, a specialized model designed for Triton programming. Our objective is to develop a model that can generate Triton code that is both correct and highly optimized for speed, outperforming the reference implementation. To achieve this, we adopt a two-stage training strategy: we first apply supervised fine-tuning (SFT) to instill fundamental Triton syntax and kernel optimization skills, followed by reinforcement learning (RL) with fine-grained verifiable rewards to further refine the model for correctness and efficiency. We will first detail the SFT procedure, and subsequently present the design of the reinforcement learning framework.

2.1 TRITON KNOWLEDGE DISTILLATION VIA SUPERVISED FINE-TUNING

Recent work (Fisches et al., 2025; Li et al., 2025b) has demonstrated that large language models (LLMs) exhibit weak Triton programming ability at the 8B scale, struggling both with syntax and

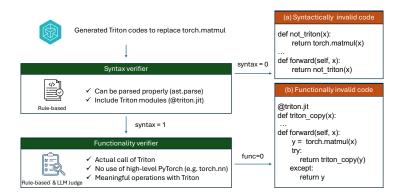


Figure 2: Illustration of the flow of our robust verifier incorporating syntax and functionality checkers and the examples of invalid Triton codes. (a) invalid syntax: the code lacks any Triton blocks and consists solely of PyTorch code. (b) invalid functionality: the code include dummy Triton code that just copies data without meaningful operation delegating core operation (matrix multiplication) to PyTorch modules (torch.matmul).

with performance-oriented design patterns. Effective dataset curation is therefore essential, not only to expose Triton-specific primitives and coding structures but also to preserve the reasoning traces that guide kernel optimization. To address this, our pipeline follows three key steps: (i) data augmentation, (ii) synthesis of reasoning traces paired with corresponding code, and (iii) construction of high-quality training pairs for supervised fine-tuning.

- (i) Data augmentation: We start from the problem sets in KernelBook (Paliskara & Saroufim, 2025), which provides curated pairs of PyTorch programs and equivalent Triton kernels. To enrich this dataset, we augment the tasks with additional variations (e.g., diverse input shapes), thereby exposing the model to broader performance scenarios.
- (ii) **Data synthesis:** To obtain diverse reasoning traces that guide correct and efficient Triton generation, we employ DeepSeek-R1 (Guo et al., 2025) to jointly synthesize reasoning steps and Triton implementations. For each task wrapped with instruction and Pytorch reference, multiple candidate kernels are collected, each paired with an explicit reasoning trace. This yields a dataset of $\mathcal{D} = \{(q, o_i)\} = \{(\text{task query, Triton code with CoT})\}$. See concrete template in Appendix E.3).
- (iii) Supervised fine-tuning: In the SFT stage, the model is trained to produce valid Triton code as well as reproduce the associated reasoning traces from the instruction. This distillation process transfers essential Triton programming patterns while reinforcing reasoning ability.

2.2 REINFORCEMENT LEARNING WITH HIERARCHICAL REWARD DECOMPOSITION

While supervised fine-tuning (SFT) equips the base model with basic Triton syntax and kernel optimization abilities, the resulting code may still contain errors or lack efficiency. To further improve the quality of Triton code generation of TRITONRL, we train the model via reinforcement learning (RL) with verifiable rewards, which incentivize model to generate more correct and efficient Triton code yielding higher rewards. In RL, designing effective reward feedback is essential since crude reward designs not perfectly aligned with original objectives of tasks often lead to reward hacking, guiding models to exploit loopholes. To address this, we first introduce robust and fine-grained verifiers that rigorously assess the quality of Triton code in diverse aspects, forming the basis for constructing reward functions. Building on these verifiers, we present a GRPO-based RL framework with hierarchical reward decomposition that provides targeted feedback for reasoning traces and Triton code, thereby improving the correctness and efficiency of generated Triton kernels.

Fine-Grained Verification for High-Quality Triton Code. We denote i-th generated output sample o_i for a given prompt or task q. Recall that q and o_i include Pytorch reference $q^{\rm ref}$ and Triton code $o_i^{\rm code}$ that is executable on some proper input x, i.e., $q^{\rm ref}(x)$, $o^{\rm code}(x)$. We introduce fine-grained verifiers v that comprehensively evaluate different aspects of code quality as follows.

- Robust verifier: Sequential verifier to check output is a valid Triton code without non-nonsensical hacking. See Figure 2 for an illustration
 - syntax: A binary verifier that assesses whether code o^{code} is valid Triton syntax. We use a rule-based linter to verify the presence of Triton kernels annotated with @triton.jit.
 - func: A binary functionality verifier to detect whether o^{code} constitutes a valid Triton kernel. Syntax checks alone are insufficient, since models may output code that superficially passes verification but defers computations to high-level PyTorch modules (e.g., torch.nn, @) or hardcodes constants, as in Figure 2 (b). To address this, we combine a rule-based linter—which ensures Triton kernels are invoked and flags reliance on PyTorch modules—with an LLM-based judge that evaluates semantic correctness against task specifications.
- compiled: A binary verifier that checks whether the generated Triton code can be successfully
 compiled without errors.
- correct: A binary verifier that evaluates whether the generated Triton code produces correct outputs by compiling and comparing its results against those of the reference PyTorch code using provided test input x.

$$correct(q, o) = compile(o^{code}) \cdot \mathbb{1}[o^{code}(x) == q^{ref}(x)]$$

• speedup: A scalar score that quantifies the execution time improvement of the generated Triton code relative to the reference PyTorch implementation. For a prompt q with reference Pytorch code q^{ref} and corresponding triton code o^{code} generated, speed-up is defined as, given test input x,

$$\operatorname{speedup}(q,o) = \frac{\tau(o^{\operatorname{code}},x)}{\tau(q^{\operatorname{ref}},x)} \cdot \operatorname{correct}(q,o),$$

where $\tau(\cdot, x)$ measure the runtimes of given code and input x.

For notational simplicity, for a given prompt q and output sample o_i , we define $v(q, o_i)$ as v_i any verification function v throughput the paper.

While prior works (Li et al., 2025b; Baronio et al., 2025) addressed reward hacking with rule-based linters, such methods remain vulnerable to loopholes. Our verifier combines rule-based and LLM-based checks to capture both syntactic and semantic errors, offering stronger guidance during training (see Appendix E.2 for examples and Section 3 for evaluation). Building on these fine-grained verifiers, we develop an RL framework that delivers targeted feedback to both reasoning traces and Triton code, improving kernel correctness and efficiency.

Hierarchical Reward Decomposition. Training LLMs with long reasoning traces remains challenging because providing appropriate feedback across lengthy responses is difficult. When a single final reward is uniformly applied to all tokens, it fails to distinguish between those that meaningfully contribute to correctness or efficiency and those that do not. This issue is particularly pronounced in kernel code generation, where reasoning traces often outline complex optimization strategies for GPU operations. Even if the plan itself is well-formed and could yield significant speedups, errors in the subsequent Triton implementation may cause the entire response to be penalized, which was also identified by (Qu et al., 2025). This prevents the model from effectively learning good optimization strategies, conflating high-quality reasoning with poor execution.

To address this, we propose a GRPO with hierarchical reward decomposition for Triton code generation. Specifically, Triton code generation o_i can be viewed as two-level hierarchy action pairs,

- o_i^{plan} : CoT reasoning traces correspond to high-level planning actions, providing abstract kernel optimization strategies, such as tiling or shared memory.
- o_i^{code} : final Triton code answers correspond to low-level coding actions that execute the plan given by the previous reasoning traces.

The key idea is to assign different reward credit for different class of output tokens between plan and code. By jointly optimizing rewards for both levels, we can train the model to better align its

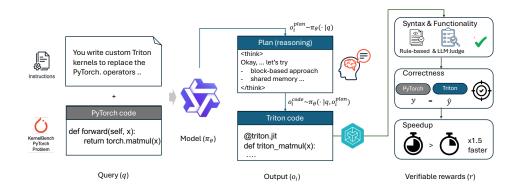


Figure 3: An example LLM output for Triton code generation, showing a reasoning trace (plan) and the generated Triton kernel code conditioned on the plan.

reasoning with the desired code output as follows:

$$\mathcal{J}_{\text{GRPO}}(\theta) = \mathbb{E}_{q \sim P(Q), \{o_i = (o_i^{\text{plan}}, o_i^{\text{code}})\}_{i=1}^G \sim \pi_{\theta_{old}}(\cdot | q)} \left[\frac{1}{G} \sum_{i=1}^G \alpha \mathcal{F}_{\text{GRPO}}^{\text{plan}}(\theta, i) + \mathcal{F}_{\text{GRPO}}^{\text{code}}(\theta, i) \right]$$
(1)

where π_{θ} and $\pi_{\theta_{old}}$ are the policy model and reference model, q denotes a prompt given to the model, defining a task to implement in Triton, and $o_i = (o_i^{\mathrm{plan}}, o_i^{\mathrm{code}})$ represents i-th response generated by the model for q in the group G. Here we denote $o_{i,t}^{\mathrm{plan}}$ and $o_{i,t}^{\mathrm{code}}$ as t-th token of output plan and code. The GRPO losses $\mathcal{F}^{\mathrm{plan}}(\theta)$ and $\mathcal{F}^{\mathrm{code}}(\theta)$ are computed over the tokens in generated plans and Triton codes, respectively, and $\alpha \in [0,1]$ is a weighting factor that balances the training speed of planning and coding policy. Here, α is set to a small value (e.g., 0.1) so that the planning distribution is updated slowly, allowing the coding policy sufficient time to learn correct implementations conditioned on those high-level plans. The detailed formulation of each loss component is as follows.

$$\mathcal{F}_{\text{GRPO}}^{\text{plan}} = \left[\frac{1}{|\sigma_{i}^{\text{plan}}|} \sum_{t=0}^{|\sigma_{i}^{\text{plan}}|} \left\{ \min \left[\frac{\pi_{\theta}(\sigma_{i,t}^{\text{plan}}|q, \sigma_{i,c+t}^{\text{plan}})}{\pi_{\theta_{old}}(\sigma_{i,t}^{\text{plan}}|q, \sigma_{i,c+t}^{\text{plan}})} \hat{A}_{i,t}^{\text{plan}}, \text{clip} \left(\frac{\pi_{\theta}(\sigma_{i,t}^{\text{plan}}|q, \sigma_{i,c+t}^{\text{plan}})}{\pi_{\theta_{old}}(\sigma_{i,t}^{\text{plan}}|q, \sigma_{i,c+t}^{\text{plan}})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t}^{\text{plan}} \right] \right\} \right],$$

$$\mathcal{F}_{\text{GRPO}}^{\text{code}} = \left[\frac{1}{|\sigma_{i}^{\text{code}}|} \sum_{t=0}^{|\sigma_{i}^{\text{code}}|} \left\{ \min \left[\frac{\pi_{\theta}(\sigma_{i,t}^{\text{code}}|q, \sigma_{i}^{\text{plan}}, \sigma_{i,c+t}^{\text{code}})}{\pi_{\theta_{old}}(\sigma_{i,t}^{\text{code}}|q, \sigma_{i}^{\text{plan}}, \sigma_{i,c+t}^{\text{code}})}, 1 - \epsilon, 1 + \epsilon \right) \hat{A}_{i,t}^{\text{code}} \right] \right\} \right],$$

$$(2)$$

where $\hat{A}_{i,t}^{\text{plan}}$ and $\hat{A}_{i,t}^{\text{code}}$ are the group-wise advantages for plan and code tokens, computed as $A_{i,t} = r_i - \frac{1}{G} \sum_{j=1}^G r_j$, with rewards for plan and code tokens defined as:

$$r_i^{\text{plan}} = \text{syntax}_i \cdot \text{func}_i \cdot \text{speedup}_i, \quad r_i^{\text{code}} = \text{syntax}_i \cdot \text{func}_i \cdot \text{correct}_i.$$
 (3)

Here, we note that the syntax and functionality checks serve as necessary conditions for correctness and speedup evaluations. If either the syntax or functionality check fails, the generated code is deemed invalid, and both correctness and speedup rewards are set to zero.

By assigning speedup-based rewards to plan tokens and correctness-based rewards to code tokens, our approach provides targeted feedback, allowing for more fine-grained credit assignment: reasoning traces are encouraged to propose optimization strategies that yield efficient kernels, while code generation is guided to produce valid implementations that faithfully realize these plans. Moreover, by choosing small values for α , we slow down the training of plan tokens relative to code tokens to ensure the model to learn effective coding skills first for given plan distributions. This helps avoid overly penalizing the reasoning trace due to code generation instability in early training, thus preserving promising optimization plans that may yield higher speedup once the code generation stabilizes. In our experiments, we explore different configurations of α and reward functions to analyze the benefits of our reward decomposition.

Data Mixing Optimization for Reinforcement Learning. In our task, we have a training dataset \mathcal{D}_{train} that is made of examples from KernelBook. We have an evaluation dataset \mathcal{D}_{test} from

KernelBench. We regard the problem of selecting the right examples to post-train our LLM model as a data mixing optimization problem, and we formulate the problem as follows.

Given a training dataset \mathcal{D}_{train} and a test dataset \mathcal{D}_{test} , either dataset can be divided into m subsets. we suppose the probability of selecting m training data subsets forms a probability simplex $\mathbf{p} \in \Delta$, which represents the probability of drawing different subsets in the training set. We define the reward function $R^i_{test}(\mathbf{p}) = \frac{1}{|\mathcal{D}^i_{test}|} \sum_{q \in \mathcal{D}^i_{test}} \max_{o \in \pi^p(q)} r(q, o)$, where π^p is the policy that is trained under the data mixture \mathbf{p} . In this problem, we use a static mixture \mathbf{p} for all post-training steps. Following the proposal by Chen et al. (2024), we can express a data mixture optimization problem as the following optimization problem:

$$\underset{\boldsymbol{p} \in \Delta}{\text{maximize}} \sum_{j=1}^{m} R_{test}^{j}(\boldsymbol{p}) \text{ s.t. } R_{test}^{j}(\boldsymbol{p}) = \zeta(S, \boldsymbol{p}), \text{ where } S \in \mathbb{R}^{m \times m}.$$

$$\tag{4}$$

We define S as the reward interaction matrix, where S_{ij} captures the effect of post-training on the i-th subset in \mathcal{D}_{train} and then using the learned policy π to generate corresponding kernels for the j-th \mathcal{D}_{test} . We use a general function ζ to leverage this information to predict the reward obtained on the j-th subset in \mathcal{D}_{test} .

We make two crucial changes to prior data mixing work. First, unlike the typical data mixing work that assumes ζ can be parametrized as a linear function, we do not make specific assumptions (Chen et al., 2023; Xie et al., 2023; Fan et al., 2023). Second, we do not create random subsets of our data. Instead, we note that KernelBench (Ouyang et al., 2025) has three complexity levels. We use an LLM labeler to create difficulty labels for each data point in \mathcal{D}_{train} . Since we focus on difficulty level 1 and 2, we create two subsets for training and test: $\mathcal{D}_{train} = \{\mathcal{D}_{train}^{L1}, \mathcal{D}_{train}^{L2}\}$. Since the problem is under-constrained without a known parametric form of ζ , we simply evaluate three candidate initializations, $\boldsymbol{p} \in [1,0], [0,1], [0.5,0.5]$, and choose the best-performing mixture rather than fully modeling S and ζ or solving for \boldsymbol{p} exhaustively.

3 EXPERIMENTS

This section provides the detailed recipe of training and evaluation of TRITONRL, followed by the main results and ablation studies.

3.1 TRAINING AND EVALUATION SETUPS

Data Preparation. For both SFT and RL, we use 11k tasks from KernelBook (Paliskara & Saroufim, 2025). We expand each task with five reasoning traces and corresponding Triton implementations generated by DeepSeek-R1 (Guo et al., 2025), yielding 58k <task query, Triton code with CoT> pairs. Prompts adopt a one-shot format, where the reference PyTorch code is given and the model is asked to produce an optimized Triton alternative (examples in Appendix E.3). To support curriculum in RL training, we further label tasks into three difficulty levels using Qwen3-235B-Instruct (Team, 2025), following the task definitions in Ouyang et al. (2025), yielding 11k <task query, level> pairs. See detailed generation and classification in Appendix G and E.1.

Training Configuration. We begin by fine-tuning the base model Qwen3-8B on Level-1 tasks. After SFT, we move to RL training on the same KernelBook tasks, but without output labels—rewards are computed directly from code execution—so the RL dataset consists only of task instructions. An example of such instructions is shown in Appendix E.3. We implement training under the VeRL framework (Sheng et al., 2025), starting from Level-1 tasks and gradually progressing to higher levels as performance improves (though current results use only Level-1). Hyperparameters for both SFT and RL are provided in Appendix G. By default, we use the reward function r from equation 3, setting $\alpha=0.1$ unless specified otherwise.

Evaluation Benchmarks. We evaluate TRITONRL on KernelBench (Li et al., 2025a)¹. KernelBench offers an evaluation framework covering 250 tasks, divided into Level 1 (100 single-kernel

¹We use the Triton backend version of KernelBench from https://github.com/ ScalingIntelligence/KernelBench/pull/35.

Model	#Params	LEVEL1 (ROBUST VERIFIER)					LEVEL1 (W/O ROBUST VERIFIER)		
	"I di dilis	valid	compiled / correct	fast ₁ / fast ₂	mean speedup	valid	compiled / correct		
Qwen3 (base)	8B	73.0	40.0 / 14.0	0.0 / 0.0	0.03	54.0	52.0 / 15.0		
Qwen3	14B	82.0	65.0 / 17.0	0.0 / 0.0	0.04	66.0	71.0 / 16.0		
Qwen3	32B	75.0	61.0 / 16.0	2.0 / 0.0	0.06	52.0	50.0 / 15.0		
KernelLLM	8B	42.0	40.0 / 20.0	0.0 / 0.0	0.05	100.0	98.0 / 29.0		
AutoTriton	8B	97.0	78.0 / 50.0	2.0 / 1.0	0.25	100.0	95.0 / 70.0		
TRITONRL (ours)	8B	99.0	82.0 / 56.0	5.0 / 1.0	0.33	99.0	83.0 / 58.0		
w/o RL (SFT only)	8B	97.0	88.0 / 44.0	4.0 / 2.0	0.33	98.0	93.0 / 47.0		
Claude-3.7	-	99.0	99.0 / 53.0	3.0 / 1.0	0.32	100.0	100.0 / 64.0		

Table 1: Main results on KernelBench Level 1. All metrics are reported as pass@10 (%). Our model achieves the best results among models with fewer than 32B parameters. The left block reports evaluation with the robust verifier (syntax + functionality). The right block (w/o robust verifier) lacks functionality checks, leading to misleading correctness estimates.

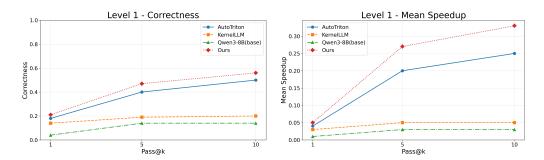


Figure 4: Pass@k correctness and mean speedup for k=1,5,10 on KernelBench Level 1 tasks. We adopted our robust verifier to check validity.

tasks, such as convolution), Level 2 (100 simple fusion tasks, such as conv+bias+ReLU), and Level 3 (50 full architecture tasks, such as MobileNet), to assess LLM proficiency in generating efficient CUDA kernels. We conduct experiments mainly on the Level 1 and Level 2 tasks from KernelBench. The prompts used for these benchmarks are provided in Appendix F.1.

Metrics. We evaluate the performance of LLMs for generating Triton code in terms of (1) Validity (syntax and functionality); (2) Correctness (compilation and correct output); (3) Speedup (relative execution time improvement). We report $fast_1$ and $fast_2$ to indicate the model's ability to generate Triton code that is at least as fast as or twice as fast as the reference PyTorch implementation, respectively. The formal definition of metrics is provided in Appendix D. We measure the pass@k metrics for each aspect, which indicates the ratio of generating at least one successful solution among k sampled attempts. We use k=10 as a default unless specified. We test both Triton codes and reference PyTorch codes on an NVIDIA L40S.

Baselines. We compare TRITONRL with several baselines, including KernelLLM (Fisches et al., 2025) and AutoTriton (Li et al., 2025b), which are fine-tuned LLMs specifically for Triton programming. We also include our base model Qwen3-8B (Team, 2025) without any fine-tuning, fine-tuned Qwen3-8B only after SFT, and larger Qwen3 models (e.g., Qwen3-14B and Qwen3-32B). Additionally, we evaluate Claude-3.7 (Anthropic, 2025) with unknown model size. For large model classes beyond 100B (e.g., GPT-OSS 120B (OpenAI, 2025), DeepSeek-R1-0528 (NVIDIA Developer Blog, 2025)), we report the numbers to the Appendix H.1 as a reference.

3.2 Main Experiment Results

Overall Performance on Level 1 Tasks. The left side of Table 1 presents the performance comparison results for pass@10 evaluated with robust verifiers (syntax and functionality) on Kernel-Bench Level 1 tasks. TRITONRL consistently outperforms most baseline models with < 32B parameter sizes in terms of validity, correctness, and speedup. In particular, TRITONRL surpasses AutoTriton, which also leverages SFT and RL, by achieving higher correctness and speedup, underscoring the advantages of our hierarchical reward assignment. Notably, the correctness metric

improves from 44% (SFT only) to 55% after RL, indicating that RL provides substantial gains in addition to supervised fine-tuning. Furthermore, TRITONRL achieves on-par performance to much larger models, highlighting the efficiency of our approach in enabling smaller models to excel in specialized code generation tasks. We further evaluated inference scaling by varying the number of sampled attempts (k=1,5,10), as illustrated in Figure 4. The correctness of TRITONRL increases with more samples, whereas KernelLLM and Qwen3-8B show limited improvement, suggesting that TRITONRL generates a more diverse set of codes and benefits from additional sampling during inference. Additional pass@1 and pass@5 results are provided in Appendix H.2.

Effectiveness of Validity Reward. analyze the validity of Triton codes generated by fine-tuned models to understand the types of errors each model is prone to. To examine the effects of fine-tuning on validity, we also include the base models of TRITONRL and the baseline models. In Figure 5, although both AutoTriton and TRITONRL achieve relatively high rates of validity compared to KernelLLM. a more detailed breakdown reveals that AutoTriton exhibits a much higher proportion of functionally invalid codes. Interestingly, the base model of AutoTriton shows a low rate of functionally and syntactically invalid codes, indicating that the finetuning process of AutoTriton may have led to learn functionally invalid codes. In contrast, TRITONRL generates significantly fewer invalid codes in terms of both syntax and functionality after fine-tuning, demonstrating the effectiveness of our robust verification in enhancing code quality.

Figure 5: Side-by-side comparison of base models and their post-trained variants. Successful training should reduce invalid syntax errors (yellow) and functional invalidity (red), while increasing correctness (green) from individual base model. Only TRITONRL shows this outcome.

Moreover, Table 1 highlights how heavily prior models relied on cheating shortcuts. Without functionality verification (w/o robust), AutoTriton's correctness jumps from 50% to 70%, revealing its tendency to exploit reward-driven shortcuts rather than produce true Triton code. In contrast, TRITONRL shows only a slight increase (56% to 58%), suggesting it learned to generate genuine code.

Effectiveness of Hierarchical Reward Decomposition. We analyze the impact of different reward design choices during RL training on the performance. Denoting our default configurations of reward decomposition, where reward function is set as equation 3 for $\alpha=0.1$, by λ^* , we compare our choice of reward decomposition with the following reward configurations:

- $\pmb{\lambda^1}$: $r_i^{\text{plan}} = r_i^{\text{code}} = \text{syntax}_i \cdot \text{func}_i \cdot \text{correct}_i \text{ and } \alpha = 1.0 \text{ (uniform correctness)}.$
- λ^2 : $r_i^{\mathrm{plan}} = r_i^{\mathrm{code}} = \mathrm{syntax}_i \cdot \mathrm{func}_i \cdot \mathrm{speedup}_i$ and $\alpha = 1.0$ (uniform speedup).
- λ^3 : Set reward function as equation 3 with $\alpha = 1.0$.

Note that λ^1 applies the correctness reward uniformly to all tokens, similar to the approach in Li et al. (2025b), while λ^2 does the same with the speedup reward. λ^3 and λ^* use the same default reward decomposition described in equation 3, but λ^3 assigns equal weight to both plan and code token, resulting in both being trained at the same rate, while λ^* assigns a smaller weight to the plan tokens. We train TRITONRL for these configurations under the same hyperparameter settings and the performance of each configuration evaluated on KernelBench Level 1 tasks.

The result in Table 2 shows that reward design has a substantial impact on model performance. The default configuration, λ^* , yields the best overall results, suggesting that while correctness is crucial for generating valid Triton code, incorporating speedup feedback into high-level planning helps the model uncover more efficient implementations without compromising correctness. Additionally,

performance gap between λ^* and λ^3 shows that slower updates of plan tokens enhances speedup and correctness, which highlights the importance of balanced training dynamics between plan and code tokens to prevent premature convergence.

Effectiveness of Data Mixing. In order to solve the data mixture optimization problem, given the limited budget, we explore three common parametrization of p: training solely on level 1 tasks, solely on level 2 tasks, and a mixture of both level 1 and level 2 tasks with equal ratio. We present the result in Table 3. We can see that there is a non-trivial interaction effect

Reward	valid	compiled / correct	fast ₁ / fast ₂
λ^{\star}	99.0	82.0 / 55.0	5.0 / 1.0
λ^1	95.0	75.0 / 38.0	2.0 / 1.0
λ^2	94.0	87.0 / 47.0	3.0 / 1.0
λ^3	93.0	64.0 / 33.0	2.0 / 1.0

Table 2: Ablation study on reward design of TRITONRL on KernelBench Level 1 tasks.

between post-training on different subsets of \mathcal{D}_{train} . We obtain notably outperforming model in terms of correctness and fast₁ in Level 1 tasks if we choose p = [1,0], while the performance on Level 2 is not improved even if p = [0,1]. A plausible explanation is that training on L2 tasks is not as effective as L1 tasks for boosting evaluation performance on both L1 and L2 tasks, for the model checkpoint we obtained through SFT (as confirmed in Table 4). This observation provides a direction for future work where we use adaptive p for different the post-training steps.

\mathcal{D}_{train}			\mathcal{D}_{test} Level1			\mathcal{D}_{test} Level2	
Subset	\boldsymbol{p}	valid	compiled / correct	$fast_1 / fast_2$	valid	compiled / correct	fast ₁ / fast ₂
Level 1	[1, 0]	99.0	82.0 / 56.0	5.0 / 1.0	66.0	29.0 / 7.0	0.0 / 0.0
Level 1+2	[0.5, 0.5]	99.0	92.0 / 43.0	2.0 / 1.0	74.0	35.0 / 8.0	0.0 / 0.0
Level 2	[0, 1]	100.0	97.0 / 49.0	3.0 / 1.0	57.0	37.0 / 6.0	0.0 / 0.0

Table 3: Ablation study on data mixture for RL training of TRITONRL, where the performance is evaluated on KernelBench level 1 and level 2 tasks.

Limited Performance on Fusion Tasks. We evaluated TRITONRL and baseline models on KernelBench Level 2 tasks, which involve fused implementations such as Conv+ReLU. As shown in Table 4, TRITONRL outperforms other 8B-scale Triton-specific models in correctness and speedup, achieving performance comparable to Claude 3.7. Nevertheless, all models, including TRITONRL, show a marked drop from Level 1 to Level 2, highlighting the greater difficulty of generating fully valid Triton code for fusion tasks. This gap reflects the complexity and advanced optimizations required, underscoring substantial room for improvement.

Model	#Params		LEVEL2	LEVEL2 (W/O ROBUST	
	# 2 41 4111 5	valid	compiled / correct	fast ₁ / fast ₂	compiled / correct
Qwen3 (base)	8B	56.0	1.0 / 0.0	0.0 / 0.0	52.0 / 11.0
Qwen3	14B	35.0	24.0 / 1.0	0.0 / 0.0	94.0 / 65.0
Qwen3	32B	31.0	16.0 / 0.0	0.0 / 0.0	73.0 / 22.0
KernelLLM	8B	0.0	0.0 / 0.0	0.0 / 0.0	96.0 / 3.0
AutoTriton	8B	70.0	3.0 / 0.0	0.0 / 0.0	97.0 / 76.0
TRITONRL (ours)	8B	69.0	29.0 / 7.0	0.0 / 0.0	88.0 / 42.0
w/o RL (SFT only)	8B	67.0	32.0 / 6.0	0.0 / 0.0	98.0 / 41.0
Claude-3.7	-	34.0	34.0 / 12.0	1.0 / 0.0	98.0 / 60.0

Table 4: Main results on KernelBench Level 2. The left block reports evaluation with the robust verifier (syntax + functionality). The right block (w/o robust verifier) lacks functionality checks, leading to misleading correctness estimates, more severe than Level 1 evaluation.

4 Conclusion

In this work, we introduce TRITONRL, a specialized LLM for Triton code generation, trained with a novel RL framework featuring robust verifiable rewards and hierarchical reward assignment. Our experiments on KernelBench show that TRITONRL surpasses existing fine-tuned Triton models in validity, correctness, and efficiency. Ablation studies demonstrate that both robust reward design and hierarchical reward assignment are essential for achieving correctness and efficiency. We believe TRITONRL marks a significant advancement toward fully automated and efficient GPU kernel generation with LLMs.

REPRODUCIBILITY STATEMENT

Our code-base is built upon publicly available frameworks (Verl (Sheng et al., 2025). Section 3.1 and the Appendix G F describe the experimental settings in detail.

REFERENCES

- Anthropic. Claude 3.7 sonnet system card, 2025. URL https://www.anthropic.com/claude-3-7-sonnet-system-card.
- Carlo Baronio, Pietro Marsella, Ben Pan, and Silas Alberti. Multi-turn training for cuda kernel generation. Cognition AI Blog. URL: https://cognition.ai/blog/kevin-32b, 2025. Accessed on May 06, 2025.
 - Mayee Chen, Nicholas Roberts, Kush Bhatia, Jue Wang, Ce Zhang, Frederic Sala, and Christopher Ré. Skill-it! a data-driven skills framework for understanding and training language models. *Advances in Neural Information Processing Systems*, 36:36000–36040, 2023.
 - Mayee F Chen, Michael Y Hu, Nicholas Lourie, Kyunghyun Cho, and Christopher Ré. Aioli: A unified optimization framework for language model data mixing. *arXiv preprint arXiv:2411.05735*, 2024.
 - Simin Fan, Matteo Pagliardini, and Martin Jaggi. Doge: Domain reweighting with generalization estimation. *arXiv preprint arXiv:2310.15393*, 2023.
 - Zacharias Fisches, Sahan Paliskara, Simon Guo, Alex Zhang, Joe Spisak, Chris Cummins, Hugh Leather, Joe Isaacson, Aram Markosyan, and Mark Saroufim. Kernelllm, 5 2025. URL https://huggingface.co/facebook/Kernelllm. Corresponding authors: Aram Markosyan, Mark Saroufim.
 - Jiaxuan Gao, Shusheng Xu, Wenjie Ye, Weilin Liu, Chuyi He, Wei Fu, Zhiyu Mei, Guangju Wang, and Yi Wu. On designing effective rl reward at training time for llm reasoning. *arXiv preprint arXiv:2410.15115*, 2024.
 - Etash Guha, Ryan Marten, Sedrick Keh, Negin Raoof, Georgios Smyrnis, Hritik Bansal, Marianna Nezhurina, Jean Mercat, Trung Vu, Zayne Sprague, Ashima Suvarna, Benjamin Feuer, Liangyu Chen, Zaid Khan, Eric Frankel, Sachin Grover, Caroline Choi, Niklas Muennighoff, Shiye Su, Wanjia Zhao, John Yang, Shreyas Pimpalgaonkar, Kartik Sharma, Charlie Cheng-Jie Ji, Yichuan Deng, Sarah Pratt, Vivek Ramanujan, Jon Saad-Falcon, Jeffrey Li, Achal Dave, Alon Albalak, Kushal Arora, Blake Wulfe, Chinmay Hegde, Greg Durrett, Sewoong Oh, Mohit Bansal, Saadia Gabriel, Aditya Grover, Kai-Wei Chang, Vaishaal Shankar, Aaron Gokaslan, Mike A. Merrill, Tatsunori Hashimoto, Yejin Choi, Jenia Jitsev, Reinhard Heckel, Maheswaran Sathiamoorthy, Alexandros G. Dimakis, and Ludwig Schmidt. Openthoughts: Data recipes for reasoning models, 2025.
 - Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in Ilms via reinforcement learning. *ArXiv preprint*, abs/2501.12948, 2025. URL https://arxiv.org/abs/2501.12948.
 - Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V. Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, Yuling Gu, Saumya Malik, Victoria Graf, Jena D. Hwang, Jiangjiang Yang, Ronan Le Bras, Oyvind Tafjord, Chris Wilhelm, Luca Soldaini, Noah A. Smith, Yizhong Wang, Pradeep Dasigi, and Hannaneh Hajishirzi. Tulu 3: Pushing frontiers in open language model post-training, 2025.
 - 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. 2025.
- Cong Duy Vu Le, Jinxin Chen, Zihan Li, Hongyu Sun, Yuan Liu, Ming Chen, Yicheng Zhang, Salam Zhihong Zhang, Hong Wang, Sheng Yang, et al. Coderl: Mastering code generation through pretrained models and deep reinforcement learning. arXiv preprint arXiv:2207.01780, 2022. URL https://ar5iv.labs.arxiv.org/html/2207.01780.

544

545

546

547 548

549

550551

552

553

554

555 556

558

559

560

561 562

563

564

565

566

567 568

569

570 571

572

573

574 575

576

577

578 579

580

581

582

583

584 585

586

587

588 589

590

592

- Jianling Li, Shangzhan Li, Zhenye Gao, Qi Shi, Yuxuan Li, Zefan Wang, Jiacheng Huang, Haojie Wang, Jianrong Wang, Xu Han, et al. Tritonbench: Benchmarking large language model capabilities for generating triton operators. *arXiv preprint arXiv:2502.14752*, 2025a.
 - Shangzhan Li, Zefan Wang, Ye He, Yuxuan Li, Qi Shi, Jianling Li, Yonggang Hu, Wanxiang Che, Xu Han, Zhiyuan Liu, and Maosong Sun. Autotriton: Automatic triton programming with reinforcement learning in llms. *arXiv preprint arXiv:2507.05687*, 2025b. URL https://arxiv.org/pdf/2507.05687.
 - Xiaoya Li, Xiaofei Sun, Albert Wang, Jiwei Li, and Chris Shum. Cuda-11: Improving cuda optimization via contrastive reinforcement learning. *arXiv preprint arXiv:2507.14111*, 2025c.
 - Yujia Li, David Choi, Junyoung Chung, Nate Glaese, Rew Beattie, Markus Pex, Huanling Wu, Edward Zielinski, Quandong Ma, Timo Wicke, et al. Competition-level code generation with alphacode. *Science*, 378(6624):1092-1097, 2022. URL https://www.researchgate.net/publication/366137000_Competition-level_code_generation_with_AlphaCode.
 - Blog. Automating **NVIDIA** Developer generation with gpu kernel deepseek-r1 and inference **NVIDIA** Developer time scaling. Blog, https://developer.nvidia.com/blog/ February URL 2025. automating-gpu-kernel-generation-with-deepseek-rl-and-inference-time-scaling/. Accessed on May 20, 2025.
 - OpenAI. gpt-oss-120b & gpt-oss-20b model card, 2025. URL https://arxiv.org/abs/2508.10925.
 - 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.
 - Sahan Paliskara and Mark Saroufim. Kernelbook, 5 2025. URL https://huggingface.co/datasets/GPUMODE/KernelBook.
 - Yuxiao Qu, Anikait Singh, Yoonho Lee, Amrith Setlur, Ruslan Salakhutdinov, Chelsea Finn, and Aviral Kumar. Learning to discover abstractions for LLM reasoning. In *ICML 2025 Workshop on Programmatic Representations for Agent Learning*, 2025. URL https://openreview.net/forum?id=zwEUO0KT8G.
 - 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.
 - Archit Sharma, Sedrick Keh, Eric Mitchell, Chelsea Finn, Kushal Arora, and Thomas Kollar. A critical evaluation of ai feedback for aligning large language models. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), *Advances in Neural Information Processing Systems*, volume 37, pp. 29166–29190. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/33870b3e099880cd8e705cd07173ac27-Paper-Conference.pdf.
 - Guangming Sheng, Chi Zhang, Zilingfeng Ye, Xibin Wu, Wang Zhang, Ru Zhang, Yanghua Peng, Haibin Lin, and Chuan Wu. Hybridflow: A flexible and efficient rlhf framework. In *Proceedings of the Twentieth European Conference on Computer Systems*, pp. 1279–1297, 2025.
 - Joar Skalse, Nikolaus Howe, Dmitrii Krasheninnikov, and David Krueger. Defining and characterizing reward gaming. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), Advances in Neural Information Processing Systems, volume 35, pp. 9460-9471. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/3d719fee332caa23d5038b8a90e81796-Paper-Conference.pdf.

 Kimi Team, Angang Du, Bofei Gao, Bowei Xing, Changjiu Jiang, Cheng Chen, Cheng Li, Chenjun Xiao, Chenzhuang Du, Chonghua Liao, Chuning Tang, Congcong Wang, Dehao Zhang, Enming Yuan, Enzhe Lu, Fengxiang Tang, Flood Sung, Guangda Wei, Guokun Lai, Haiqing Guo, Han Zhu, Hao Ding, Hao Hu, Hao Yang, Hao Zhang, Haotian Yao, Haotian Zhao, Haoyu Lu, Haoze Li, Haozhen Yu, Hongcheng Gao, Huabin Zheng, Huan Yuan, Jia Chen, Jianhang Guo, Jianlin Su, Jianzhou Wang, Jie Zhao, Jin Zhang, Jingyuan Liu, Junjie Yan, Junyan Wu, Lidong Shi, Ling Ye, Longhui Yu, Mengnan Dong, Neo Zhang, Ningchen Ma, Qiwei Pan, Qucheng Gong, Shaowei Liu, Shengling Ma, Shupeng Wei, Sihan Cao, Siying Huang, Tao Jiang, Weihao Gao, Weimin Xiong, Weiran He, Weixiao Huang, Weixin Xu, Wenhao Wu, Wenyang He, Xianghui Wei, Xianqing Jia, Xingzhe Wu, Xinran Xu, Xinxing Zu, Xinyu Zhou, Xuehai Pan, Y. Charles, Yang Li, Yangyang Hu, Yangyang Liu, Yanru Chen, Yejie Wang, Yibo Liu, Yidao Qin, Yifeng Liu, Ying Yang, Yiping Bao, Yulun Du, Yuxin Wu, Yuzhi Wang, Zaida Zhou, Zhaoji Wang, Zhaowei Li, Zhen Zhu, Zheng Zhang, Zhexu Wang, Zhilin Yang, Zhiqi Huang, Zihao Huang, Ziyao Xu, Zonghan Yang, and Zongyu Lin. Kimi k1.5: Scaling reinforcement learning with llms, 2025.

Qwen Team. Qwen3 technical report, 2025. URL https://arxiv.org/abs/2505.09388.

Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information Processing Systems*, 36:69798–69818, 2023.

Yuyu Zhang, Jing Su, Yifan Sun, Chenguang Xi, Xia Xiao, Shen Zheng, Anxiang Zhang, Kaibo Liu, Daoguang Zan, Tao Sun, et al. Seed-coder: Let the code model curate data for itself. *arXiv* preprint arXiv:2506.03524, 2025.

A LLM USAGE

We used an LLM to improve the writing by correcting grammar in our draft. It was not used to generate research ideas.

B RELATED WORK

B.1 LLM FOR KERNEL GENERATION

The exponential growth in demand for GPU computing resources has driven the need for highly optimized GPU kernels that improves computational efficiency. However, writing efficient GPU kernels is a complex and time-consuming task that requires specialized knowledge of GPU architectures and programming models. This has spurred significant interest in leveraging Large Language Models (LLMs), for automated kernel generation, especially for CUDA and Triton (Shao et al., 2024; Ouyang et al., 2025; Li et al., 2025a; NVIDIA Developer Blog, 2025). While these general-purpose models excel at a variety of programming tasks, they often struggle with custom kernel generation, achieving low success rates on specialized gpu programming tasks (Ouyang et al., 2025), highlighting the need for domain-specific models tailored to kernel synthesis.

For CUDA kernel generation, Ouyang et al. (2025) introduced KERNELBENCH, an open-source framework for evaluating LMs' ability to write fast and correct kernels on a suite of 250 carefully selected PyTorch ML workloads. Furthermore, Lange et al. (2025) presented an agentic framework, which leverages LLMs to translate PyTorch code into CUDA kernels and iteratively optimize them using performance feedback. Additionally, several works have focused on fine-tuning LLMs tailored for CUDA kernel generation. For example, Kevin-32B (Baronio et al., 2025) is a 32B parameter model fine-tuned via multi-turn RL to enhance kernel generation through self-refinement, and CUDA-L1 (Li et al., 2025c) applies contrastive reinforcement learning to DeepSeek-V3-671B, achieving notable speedup improvements in CUDA optimization tasks.

Another line of research focuses on Triton kernel generation. Li et al. (2025a) introduced TRITON-BENCH, providing evaluations of LLMs on Triton programming tasks and highlighting the challenges of Triton's domain-specific language and GPU programming complexity. To further enhance LLMs' capabilities in Triton programming, Fisches et al. (2025) has introduced KernelLLM, a fine-tuned model of Llama3.1-8B-Instruct via supervised fine-tuning with Pytorch and Triton code pairs in KernelBook Paliskara & Saroufim (2025), but its performance is limited by the quality of training data. Similarly, Li et al. (2025b) introduced AutoTriton, a model fine-tuned specifically for Triton programming from Seed-Coder-8B-Reasoning Zhang et al. (2025), which achieves improved performance via SFT and RL with verifiable rewards based on correctness and rule-based Triton syntax verification, which may have limited improvement in runtime efficiency due to correctness-focused rewards. Both KernelLLM and AutoTriton are concurrent works developed alongside our work, and we provide a detailed comparison in Section 3.2.

B.2 REINFORCEMENT LEARNING WITH VERIFIABLE REWARDS

Reinforcement Learning (RL) has become a key technique for training Large Language Models (LLMs), especially in domains where verifiable reward signals are available. Unlike supervised fine-tuning (SFT), which relies on curated examples, RL enables models to learn through trial and error, guided solely by reward feedback. This makes the design of accurate reward functions critical, as the model's behavior is shaped entirely by the reward signal. As a result, RL with verifiable rewards (RLVR) (Lambert et al., 2025; Team et al., 2025; Guo et al., 2025) has gained significant traction in applications like mathematics and code generation (Shao et al., 2024; Li et al., 2022), where external verification is feasible through solution correctness or unit test outcomes.

In math and coding applications, the reward can be directly computed solely based on the final outcomes when ground-truth answers or unit tests are available. For tasks where validation is not available or noisy, rule-based verification or LLM-based judges can be employed to verify the quality of generated content (Guha et al., 2025; Guo et al., 2025). For coding tasks, unit tests are commonly used to measure whether generated code meets the specified requirements (Le et al., 2022; Ouyang et al., 2025). However, unit tests often fail to cover edge cases or fully capture the problem require-

ments, leading to potential "reward hacking" (Skalse et al., 2022) where the model generates code that passes the tests but does not genuinely solve the task (Sharma et al., 2024; Gao et al., 2024). Such reward hacking has been observed in kernel generation tasks, where models produce superficially correct codes passing unit tests by using high-level Pytorch modules instead of implementing custom kernels. To address this, some works Li et al. (2025b); Baronio et al. (2025) have introduced rule-based verification, which checks kernel syntax or use of specific high-level modules.

C NOTATIONS

The following notations will be used throughout this paper. For notational simplicity, we denote any function $f(q, o_i)$ as f_i when the context is clear.

- q: prompt given to the model, defining a task to implement in Triton
- o: output sequence generated by the model, which includes both reasoning trace and Triton code
- π_{θ} : policy model with parameters θ
- G: group size for GRPO
- o_i : *i*-th sample in the group G, which includes a reasoning trace that provides the "plan" for Triton code optimization and implementation and the final "Triton code", i.e. $o_i = \{o_{i,\text{plan}}, o_{i,\text{triton}}\}$
- T_i^c : set of token indices corresponding to token class $c \in \{\text{plan, triton}\}$ in the i-th sample
- $r^c(q, o_i) = r_i^c$: reward function for token class $c \in \{\text{plan}, \text{triton}\}.$
- $\hat{A}_t^c(q, o_i) = r^c(q, o_i) \frac{1}{G} \sum_{j=1}^G r^c(q, o_j)$: token-level advantage of the t-th token of the i-th sample belonging to token class $c \in \{\text{plan, triton}\}$, shortened as $\hat{A}_{i,t}^c$.

D METRICS

We provide the formal definitions of the evaluation metrics used in this paper. Given a set of N tasks $\{q_n\}_{n=1}^N$ and k samples $\{o_i\}_{i=1}^k$ generated by the model for each task, we define the following metrics:

$$\begin{aligned} & \text{valid} = \frac{1}{N} \sum_{n=1}^{N} \max_{i \in [k]} \mathbb{I}(\text{syntax}(q_n, o_i) \cdot \text{func}(q_n, o_i) = 1) \\ & \text{compiled} = \frac{1}{N} \sum_{n=1}^{N} \max_{i \in [k]} \mathbb{I}(\text{syntax}(q_n, o_i) \cdot \text{func}(q_n, o_i) \cdot \text{compiled}(q_n, o_i) = 1) \\ & \text{correct} = \frac{1}{N} \sum_{n=1}^{N} \max_{i \in [k]} \mathbb{I}(\text{syntax}(q_n, o_i) \cdot \text{func}(q_n, o_i) \cdot \text{correct}(q_n, o_i) = 1) \\ & \text{fast}_p = \frac{1}{N} \sum_{n=1}^{N} \max_{i \in [k]} \mathbb{I}(\text{syntax}(q_n, o_i) \cdot \text{func}(q_n, o_i) \cdot \text{correct}(q_n, o_i) \cdot \text{speedup}(q_n, o_i) > p) \\ & \text{mean_speedup} = \frac{1}{N} \sum_{n=1}^{N} \max_{i \in [k]} (\text{syntax}(q_n, o_i) \cdot \text{func}(q_n, o_i) \cdot \text{correct}(q_n, o_i) \cdot \text{speedup}(q_n, o_i)) \end{aligned}$$

E DATA CURATION AND EXAMPLES

E.1 DATA MIXING SUBSET CREATION

We labeled difficulty level of 11k PyTorch reference codes in KernelBook based on the complexity of kernel implementation using Qwen3-235B-Instruct (Team, 2025). For each given PyTorch reference code, we prompt Qwen3-235B-Instruct (temperature=0.7, top_p=0.8) to label the difficulty level of replacing the PyTorch reference with Triton code as follows:

Instruction (input) example

```
''' <PyTorch reference code> '''
```

Assign a kernel implementation complexity level (1, 2, or 3) of the provided reference Py-Torch architecture according to the criteria below:

- Level 1: Single primitive operation. This level includes the foundational building blocks of AI (e.g. convolutions, matrix-vector and matrix-matrix multiplications, losses, activations, and layer normalizations). Since PyTorch makes calls to several well-optimized and often closed-source kernels under-the-hood, it can be challenging for LMs to outperform the baseline for these primitive operations. However, if an LM succeeds, the open-source kernels could be an impactful alternative to the closed-source (e.g., CuBLAS [27]) kernels.
 Level 2: Operator sequences. This level includes AI workloads containing multiple primitive operations, which can be fused into a single kernel for improved performance (e.g.,
- itive operations, which can be fused into a single kernel for improved performance (e.g., a combination of a convolution, ReLU, and bias). Since compiler-based tools such as the PyTorch compiler are effective at fusion, it can be challenging for LMs to outperform them. However, LMs may propose more complex algorithms compared to compiler rules. Level 3: This level includes architectures that power popular AI models, such as AlexNet and MiniGPT, collected from popular PyTorch repositories on GitHub.

E.2 INVALID TRITON CODE EXAMPLES

Here, we provide some examples of invalid Triton code generated by the baseline models.

Example 1: Triton syntax failure (syntax = 0) – no Triton kernel. This is not a Triton kernel because it does not include any custom @triton.jit kernel.

```
import torch
import torch.nn as nn
import triton
import triton.language as tl
import torch.nn.functional as F
# To replace with Triton, a custom kernel would be needed (not provided
   here)
class ModelNew(nn.Module):
  def __init__(self) -> None:
      super().__init__()
   def forward(self, x):
     return F.avg_pool2d(x, kernel_size=3, stride=None, padding=0)
def get_inputs():
   a = torch.randn(1, 128).cuda()
  b = torch.randn(1, 128).cuda()
  return [a, b]
def get_init_inputs():
   return []
```

Example 2: Functional failure (syntax = 1, func = 0) – use of high-level PyTorch module. This implementation contains a Triton kernel, but the main computation (convolution) is still handled by the PyTorch nn.Conv3d module. The Triton kernel is only used for a simple addition, rather than implementing the full intended operation.

```
import torch
import torch.nn as nn
import triton
import triton.language as tl
```

```
810
       @triton.jit
811
       def add_bias_kernel(
812
          out_ptr,
813
          bias_ptr,
          n_elements,
814
          out_C,
815
          out_D,
816
          out_W,
817
          out_H,
818
          BLOCK_SIZE: tl.constexpr,
819
          pid = tl.program_id(0)
820
          block_start = pid * BLOCK_SIZE
821
          offsets = block_start + tl.arange(0, BLOCK_SIZE)
822
          mask = offsets < n_elements</pre>
823
          out = tl.load(out_ptr + offsets, mask=mask, other=0.0)
824
825
          stride_channel = out_D * out_W * out_H
826
          channel_idx = (offsets // stride_channel) % out_C
827
          bias = tl.load(bias_ptr + channel_idx, mask=mask, other=0.0)
828
          out = out + bias
829
          tl.store(out_ptr + offsets, out, mask=mask)
830
831
832
       def triton_add_bias(out: torch.Tensor, bias: torch.Tensor):
          assert out.is_cuda and bias.is_cuda, "Tensors_must_be_on_CUDA."
833
834
          out = out.contiguous()
835
          bias = bias.contiquous()
836
          n_elements = out.numel()
837
          BLOCK_SIZE = 128
838
          grid = lambda meta: (
839
             (n_elements + meta["BLOCK_SIZE"] - 1) // meta["BLOCK_SIZE"],
840
          )
841
842
          batch_size, out_channels, D, W, H = out.shape
843
          add_bias_kernel[grid](
844
             out,
845
             bias,
846
             n_elements,
847
             out_channels,
848
             D,
             W,
849
             Η.
850
             BLOCK_SIZE=BLOCK_SIZE
851
852
          return out
853
854
855
       class ModelNew(nn.Module):
856
          def __init__(
857
             self,
             in_channels: int,
858
             out_channels: int,
859
             kernel_size: int,
860
             stride: int = 1,
861
             padding: int = 0,
862
             dilation: int = 1,
             groups: int = 1,
863
             bias: bool = False
```

882

883

884

885

886

```
864
865
             super(ModelNew, self).__init__()
866
             self.conv3d = nn.Conv3d(
867
                in channels,
                out_channels,
868
                 (kernel_size, kernel_size, kernel_size),
869
                stride=stride,
870
                padding=padding,
871
                dilation=dilation,
872
                groups=groups,
                bias=bias
873
874
875
          def forward(self, x: torch.Tensor) -> torch.Tensor:
876
             out = self.conv3d(x)
             if self.conv3d.bias is not None:
877
                out = triton_add_bias(out, self.conv3d.bias)
878
             return out
879
880
```

Example 3: Functional failure (syntax = 1, func = 0) – hardcoded output and no meaningful computation. While the Triton kernel is syntactically correct, but it doesn't actually implement the intended operation (Group Normalization). The kernel doesn't compute mean or variance, which are essential for GroupNorm. To implement real GroupNorm, you'd need to compute pergroup statistics and normalize accordingly. Also, it only loads the input tensor and writes it back unchanged.

```
887
       import torch
888
       import torch.nn as nn
889
       import triton
890
       import triton.language as tl
891
892
       @triton.jit
893
      def groupnorm_kernel(
894
          x_ptr, # Pointer to x tensor
895
          y_ptr, # Pointer to y tensor (not used here)
896
          out_ptr, # Pointer to output tensor
897
          n_elements, # Total number of elements
          BLOCK_SIZE: tl.constexpr,
898
      ):
899
          # Each program handles a contiguous block of data of size BLOCK_SIZE
900
         block_start = tl.program_id(0) * BLOCK_SIZE
901
          # Create a range of offsets [0..BLOCK_SIZE-1]
902
          offsets = block_start + tl.arange(0, BLOCK_SIZE)
          # Mask to ensure we don't go out of bounds
903
         mask = offsets < n_elements</pre>
904
          # Load input value
905
          x = tl.load(x_ptr + offsets, mask=mask, other=0.0)
906
          # Compute x squared
907
          x\_squared = x * x
          # Store the result
908
          tl.store(out_ptr + offsets, x, mask=mask)
909
910
911
      def triton_groupnorm(x: torch.Tensor, y: torch.Tensor):
          assert x.is_cuda and y.is_cuda, "Tensors must be on CUDA."
912
          x = x.contiguous()
913
          y = y.contiguous()
914
915
          # Prepare output tensor
916
          out = torch.empty_like(x)
917
          # Number of elements in the tensor
```

```
918
          n_{elements} = x.numel()
919
         BLOCK_SIZE = 128 # Tunable parameter for block size
920
          # Determine the number of blocks needed
921
          grid = lambda meta: ((n_elements + meta["BLOCK_SIZE"] - 1) // meta["
922
             BLOCK_SIZE"],)
923
924
          # Launch the Triton kernel
         groupnorm_kernel[grid](x, y, out, n_elements, BLOCK_SIZE=BLOCK_SIZE)
925
         return out
926
927
928
      class ModelNew(nn.Module):
929
         def __init__(self, num_features: int, num_groups: int) -> None:
930
             super().__init__()
             self.num_features = num_features
931
             self.num_groups = num_groups
932
933
          def forward(self, x: torch.Tensor) -> torch.Tensor:
934
             # Use Triton kernel for elementwise operations
             x_{triton} = triton_{groupnorm}(x, x)
             # Manually compute mean and variance (as Triton kernel only handles
936
937
             # Actual GroupNorm logic would go here
938
             # For this example, we return the Triton processed tensor
939
             return x_triton
940
```

E.3 SFT AND RL DATASET CONSTRUCTION WITH KERNELBOOK

To synthesize SFT dataset, we extract 11,621 PyTorch reference codes from KernelBook, executable without errors, such as

```
import torch
import torch.nn as nn

class Model(nn.Module):

def __init__(self):
    super(Model, self).__init__()

def forward(self, neighbor):
    return torch.sum(neighbor, dim=1)

def get_inputs():
    return [torch.rand([4, 4, 4, 4])]

def get_init_inputs():
    return [[], {}]
```

For each given PyTorch reference code, we construct an instruction for DeepSeek-R1 to generate CoTs and Triton kernels as:

Instruction (input) example

Your task is to write custom Triton kernels to replace as many PyTorch operators as possible in the given architecture, aiming for maximum speedup. You may implement multiple custom kernels, explore operator fusion (such as combining matmul and relu), or introduce algorithmic improvements (like online softmax). You are only limited by your imagination. You are given the following architecture:

```
''' <PyTorch reference code> '''
```

You have to optimize the architecture named Model with custom Triton kernels. Optimize the architecture named Model with custom Triton kernels! Name your optimized output architecture ModelNew. Output the new code in codeblocks. Please generate real code, NOT pseudocode, make sure the code compiles and is fully functional. Just output the new model code, no other text, and NO testing code! Before writing a code, reflect on your idea to make sure that the implementation is correct and optimal.

Given the instruction for each PyTorch reference code, we collect (CoT, Triton kernel code) pairs fom DeepSeek-R1 and construct outputs for SFT by concatenating the pairs as follows:

```
Triton kernel with CoT (output) example

<think>
CoT
</think>
...
<Triton kernel code>
```

In this manner, for each Pytorch reference code in KernelBook, we construct 5 (input, output) SFT samples.

For RL training, we use the same instruction input as a prompt for the same set of Pytorch reference codes in KernelBook, without the output synthesized by DeepSeek-R1 because RL training only requires reward feedback, which can be directly obtained from executing the generated Triton code.

F EVALUATION AND EXAMPLES

F.1 EVALUATION WITH KERNELBENCH

To evaluate the trained models, we construct prompts for 250 tasks in KernelBench. Similar to KernelBook, KernelBench provides a reference PyTorch code for each task. For each given reference PyTorch code, we construct a prompt with one simple example pair of (PyTorch code, Triton kernel code), similarly to the one-shot prompting format in KernelBench. Here, we use the following PyTorch and Triton codes for a simple add operation as an example:

```
1013
       ### PyTorch reference code ###
1014
       import torch
       import torch.nn as nn
1015
       import torch.nn.functional as F
1016
1017
       class Model(nn.Module):
1018
          def __init__(self) -> None:
1019
             super().__init__()
1020
          def forward(self, a, b):
1021
             return a + b
1022
1023
      def get_inputs():
1024
          # randomly generate input tensors based on the model architecture
1025
          a = torch.randn(1, 128).cuda()
          b = torch.randn(1, 128).cuda()
```

```
1026
          return [a, b]
1027
1028
      def get_init_inputs():
1029
          # randomly generate tensors required for initialization based on the
              model architecture
1030
          return []
1031
1032
       ### Triton kernel code ###
1033
1034
      import torch
       import torch.nn as nn
1035
       import torch.nn.functional as F
1036
       import triton
1037
      import triton.language as tl
1038
      @triton.jit
1039
      def add_kernel(
1040
          x_ptr, # Pointer to first input
1041
          y_ptr, # Pointer to second input
1042
          out_ptr, # Pointer to output
1043
          n_elements, # Total number of elements in input/output
1044
          BLOCK_SIZE: tl.constexpr,
1045
          # Each program handles a contiguous block of data of size BLOCK_SIZE
1046
         block_start = tl.program_id(0) * BLOCK_SIZE
1047
          # Create a range of offsets [0..BLOCK_SIZE-1]
1048
          offsets = block_start + tl.arange(0, BLOCK_SIZE)
1049
          # Mask to ensure we don't go out of bounds
         mask = offsets < n_elements</pre>
1050
          # Load input values
1051
          x = tl.load(x_ptr + offsets, mask=mask, other=0.0)
1052
          y = tl.load(y_ptr + offsets, mask=mask, other=0.0)
1053
          # Perform the elementwise addition
1054
          out = x + y
          # Store the result
1055
          tl.store(out_ptr + offsets, out, mask=mask)
1056
1057
      def triton_add(x: torch.Tensor, y: torch.Tensor):
1058
          This function wraps the Triton kernel call. It:
1059
           1. Ensures the inputs are contiguous on GPU.
1060
           2. Calculates the grid (blocks) needed.
1061
           3. Launches the Triton kernel.
1062
1063
          assert x.is_cuda and y.is_cuda, "Tensors_must_be_on_CUDA."
1064
          x = x.contiguous()
          y = y.contiguous()
1065
1066
          # Prepare output tensor
1067
          out = torch.empty_like(x)
1068
          # Number of elements in the tensor
1069
          n_{elements} = x.numel()
1070
          BLOCK_SIZE = 128 # Tunable parameter for block size
1071
1072
          # Determine the number of blocks needed
1073
          grid = lambda meta: ((n_elements + meta["BLOCK_SIZE"] - 1) // meta["
             BLOCK_SIZE"],)
1074
1075
          # Launch the Triton kernel
1076
          add_kernel[grid](x, y, out, n_elements, BLOCK_SIZE=BLOCK_SIZE)
1077
          return out
1078
1079
      class ModelNew(nn.Module):
          def __init__(self) -> None:
```

```
super().__init__()

def forward(self, a, b):
    # Instead of "return a + b", call our Triton-based addition
    return triton_add(a, b)
```

Prompt (input) example

1080

1081 1082

1083

1084 1085

1087 1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

110511061107

11081109

You write custom Triton kernels to replace the pytorch operators in the given architecture to get speedups.

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

Here's an example to show you the syntax of inline embedding custom Triton kernels in torch. The example given architecture is:

```
''' <Example PyTorch code> '''
```

The example new architecture with custom Triton kernels looks like this:

```
''' <Example Triton kernel code> '''
```

You are given the following architecture:

```
''' <PyTorch reference code> '''
```

Optimize the architecture named Model with custom Triton kernels! Name your optimized output architecture ModelNew. Output the new code in codeblocks. Please generate real code, NOT pseudocode, make sure the code compiles and is fully functional. Just output the new model code, no other text, and NO testing code!

F.2 Example of Triton codes with speedup > 1

Problem 12 in KernelBench Level 1: diagonal matrix multiplication with x11 speedup.

```
1110
       import torch
1111
       import triton
1112
      import triton.language as tl
1113
       @triton.jit
1114
      def fused_diag_matmul_kernel(
1115
        vec_a_ptr, # Pointer to the diagonal vector A (N,)
1116
        mat_b_ptr, # Pointer to the input dense matrix B (N, M)
1117
        output_ptr, # Pointer to the output matrix (N, M)
1118
        # --- Matrix dimensions ---
        N, # The size of the diagonal vector A
1119
        {\rm M}, # The number of columns in matrix {\rm B}
1120
        TOTAL_ELEMENTS, # Total number of elements in the output matrix (N * M)
1121
        # --- Tuning parameters ---
1122
        BLOCK_SIZE: tl.constexpr, # Number of elements each Triton program will
            handle
1123
       ):
1124
        ,, ,, ,,
1125
        Computes C = diag(A) * B, where A is a vector representing the diagonal.
1126
        This kernel treats the operation as a simple element-wise multiplication
1127
        where each row of B is scaled by the corresponding element of A.
1128
        11 11 11
        # 1. CALCULATE OFFSETS: Determine which elements this program instance
1129
            will process.
1130
        pid = tl.program_id(axis=0)
1131
        block_start_offset = pid * BLOCK_SIZE
1132
        offsets = block_start_offset + tl.arange(0, BLOCK_SIZE)
1133
        # 2. CREATE BOUNDARY-CHECK MASK: Prevent out-of-bounds memory access.
```

```
1134
        mask = offsets < TOTAL_ELEMENTS</pre>
1135
1136
        # 3. COMPUTE 2D INDICES FROM 1D OFFSETS: Find the row for vector A.
1137
        row indices = offsets // M
1138
        # 4. LOAD DATA EFFICIENTLY: Read from global memory into registers.
1139
        vec_a_elements = tl.load(vec_a_ptr + row_indices, mask=mask, other=0.0)
1140
        mat_b_elements = tl.load(mat_b_ptr + offsets, mask=mask, other=0.0)
1141
1142
        # 5. PERFORM THE CORE COMPUTATION: The actual fused operation.
        output = vec_a_elements * mat_b_elements
1143
1144
        # 6. STORE THE RESULT: Write the computed values back to global memory.
1145
        tl.store(output_ptr + offsets, output, mask=mask)
1146
```

Problem 61 in KernelBench Level 1: 3D Transposed Convolution with x1.5 speedup.

```
1148
       import torch
1149
       import torch.nn as nn
1150
       import triton
1151
      import triton.language as tl
1152
1153
       @triton.jit
1154
      def conv_transpose3d_kernel(
1155
         input_ptr, weight_ptr, output_ptr,
1156
          B, IC, OC, D_in, H_in, W_in, D_out, H_out, W_out,
1157
          input_batch_stride, input_channel_stride, input_d_stride,
              input_h_stride, input_w_stride,
1158
          weight_in_channels_stride, weight_out_channels_stride, weight_d_stride
1159
              , weight_h_stride, weight_w_stride,
1160
          output_batch_stride, output_channel_stride, output_d_stride,
1161
             output_h_stride, output_w_stride,
1162
         BLOCK_SIZE: tl.constexpr,
      ):
1163
          # Each block handles BLOCK_SIZE output elements
1164
         pid = tl.program_id(0)
1165
         n_elements = B * OC * D_out * H_out * W_out
1166
         block_start = pid * BLOCK_SIZE
1167
         offsets = block_start + tl.arange(0, BLOCK_SIZE)
1168
1169
          # Mask to avoid out-of-bounds access in last block
1170
         mask = offsets < n_elements</pre>
1171
1172
          # Precompute for flattening and indexing
          OC_D_outH_outW_out = OC * D_out * H_out * W_out
1173
         D_outH_outW_out = D_out * H_out * W_out
1174
1175
          # Decompose flat index into (batch, channel, depth, height, width)
1176
         b_idx = offsets // OC_D_outH_outW_out
         residual = offsets % OC_D_outH_outW_out
1177
         oc_idx = residual // D_outH_outW_out
1178
          spatial_idx = residual % D_outH_outW_out
1179
1180
         d_idx = spatial_idx // (H_out * W_out)
1181
         hw_idx = spatial_idx % (H_out * W_out)
         h_idx = hw_idx // W_out
1182
          w_idx = hw_idx % W_out
1183
1184
          # Accumulator for output value (per element)
1185
         acc = tl.zeros((BLOCK_SIZE,), dtype=tl.float32)
1186
1187
          # Loop over all input channels and 3x3x3 kernel positions
         for ic in range(IC):
```

```
1188
             for kd in range(3):
1189
                for kh in range(3):
1190
                   for kw in range(3):
1191
                       # Compute corresponding input position for each output
                       d_{in} = d_{idx} - kd
1192
                      h_{in} = h_{idx} - kh
1193
                      w_in = w_idx - kw
1194
1195
                       # Check input bounds to avoid invalid memory access
1196
                       in\_bounds = (
                          (d_{in} >= 0) & (d_{in} < D_{in}) &
1197
                          (h_{in} >= 0) & (h_{in} < H_{in}) &
1198
                          (w_in \ge 0) & (w_in < W_in)
1199
                       )
1200
                       # Compute input tensor offset
1201
                       input\_offsets = (
1202
                          b_idx * input_batch_stride +
1203
                          ic * input_channel_stride +
1204
                          d_in * input_d_stride +
1205
                          h_in * input_h_stride +
1206
                          w_in * input_w_stride
1207
1208
                       # Load input values with masking (zeros for out-of-bounds)
1209
                       input_val = tl.load(input_ptr + input_offsets, mask=
1210
                           in_bounds, other=0.0)
1211
                       # Compute weight tensor offset (flipped in transpose)
1212
                       weight_offsets = (
1213
                          oc_idx * weight_out_channels_stride +
1214
                          ic * weight_in_channels_stride +
1215
                          kd * weight_d_stride +
                          kh * weight_h_stride +
1216
                          kw * weight_w_stride
1217
1218
1219
                       weight_val = tl.load(weight_ptr + weight_offsets)
1220
1221
                       # Accumulate product into output accumulator
                       acc += input_val * weight_val
1222
1223
          # Store the final output value
1224
          output_offsets = (
1225
             b_idx * output_batch_stride +
             oc_idx * output_channel_stride +
1226
             d_idx * output_d_stride +
1227
             h_idx * output_h_stride +
1228
             w_idx * output_w_stride
1229
1230
          # Store with mask for thread-safety
1231
          tl.store(output_ptr + output_offsets, acc, mask=mask)
1232
1233
1234
       def triton_conv_transpose3d(x: torch.Tensor, weight: torch.Tensor) ->
1235
          torch.Tensor:
          B, IC, D_{in}, H_{in}, W_{in} = x.shape
1236
          OC = weight.shape[1]
1237
1238
          # Output dimensions for stride=1, padding=0, kernel=3
1239
          D_out = D_in + 2
1240
          H_out = H_in + 2
          W_out = W_in + 2
1241
```

```
1242
          # Allocate output tensor
1243
          output = torch.empty((B, OC, D_out, H_out, W_out), device=x.device,
1244
              dtype=x.dtype)
1245
          # Launch grid size: enough blocks to cover all output elements
1246
          BLOCK_SIZE = 256
1247
          grid = (triton.cdiv(B * OC * D_out * H_out * W_out, BLOCK_SIZE),)
1248
1249
          # Launch the kernel
1250
          conv_transpose3d_kernel[grid](
             x, weight, output,
1251
             B, IC, OC, D_in, H_in, W_in, D_out, H_out, W_out,
1252
             x.stride(0), x.stride(1), x.stride(2), x.stride(3), x.stride(4),
1253
             weight.stride(0), weight.stride(1), weight.stride(2), weight.stride
1254
                 (3), weight.stride(4),
             output.stride(0), output.stride(1), output.stride(2), output.stride
1255
                 (3), output.stride(4),
1256
             BLOCK_SIZE
1257
1258
1259
          return output
1260
1261
       class ModelNew(nn.Module):
1262
          def ___init___(
1263
             self,
1264
             in_channels: int,
1265
             out_channels: int,
             kernel_size: int,
1266
             stride: int = 1,
1267
             padding: int = 0,
1268
             output_padding: int = 0,
1269
             groups: int = 1,
             bias: bool = False
1270
          ):
1271
             super().__init__()
1272
1273
             # This implementation is hardcoded for kernel_size=3, stride=1, no
1274
                 padding, no bias
             assert (
1275
                kernel_size == 3 and stride == 1 and
1276
                padding == 0 and output_padding == 0 and
1277
                groups == 1 and not bias
1278
1279
             self.in_channels = in_channels
1280
             self.out_channels = out_channels
1281
1282
             # Weight shape: [in_channels, out_channels, 3, 3, 3]
1283
             self.weight = nn.Parameter(torch.empty(in_channels, out_channels,
1284
                 3, 3, 3))
1285
             # Initialize with Kaiming uniform
1286
             nn.init.kaiming_uniform_(self.weight, a=5**0.5)
1287
1288
          def forward(self, x: torch.Tensor) -> torch.Tensor:
1289
             return triton_conv_transpose3d(x, self.weight)
1290
```

Problem 35 in KernelBench Level 1: Group Normalization with x1.1 speedup.

1291

```
import torch
import torch.nn as nn
import triton
import triton
import triton.language as tl
```

```
1296
1297
      @triton.jit
1298
      def group_norm_reduce_kernel(
1299
          x_ptr,
          mean_group_ptr,
1300
          var_group_ptr,
1301
          N, C, H, W, num_groups, M,
1302
          BLOCK_SIZE_REDUCE: tl.constexpr
1303
      ):
1304
          # Program IDs for parallel execution across batch and groups
          pid0 = tl.program_id(0) # batch index
1305
          pid1 = tl.program_id(1) # group index
1306
1307
          group\_size = M * H * W
1308
          group_start = pid0 * (C * H * W) + pid1 * group_size
1309
          # Use float accumulators for numerical stability
1310
          sum1 = 0.0
1311
          sum2 = 0.0
1312
1313
          # Loop over the group in chunks for better cache usage
1314
          for i in range(0, group_size, BLOCK_SIZE_REDUCE):
             offset = i + tl.arange(0, BLOCK_SIZE_REDUCE)
1315
             mask = offset < group_size # avoid out-of-bounds loads</pre>
1316
1317
             offsets = group_start + offset
1318
1319
             # Efficient memory access with masking and zero-padding
             chunk = tl.load(x_ptr + offsets, mask=mask, other=0.0)
1320
1321
             # Accumulate sum and squared sum
1322
             sum1 += tl.sum(chunk, axis=0)
1323
             sum2 += tl.sum(chunk * chunk, axis=0)
1324
          # Numerically stable variance computation
1325
          mean = sum1 / group_size
1326
          var = (sum2 - mean * sum1) / group_size
1327
1328
          # Store per-group mean and variance
          mean_idx = pid0 * num_groups + pid1
1329
          tl.store(mean_group_ptr + mean_idx, mean)
1330
          tl.store(var_group_ptr + mean_idx, var)
1331
1332
1333
       @triton.jit
1334
      def group_norm_forward_kernel(
         x_ptr,
1335
          out_ptr,
1336
          mean_group_ptr,
1337
          var_group_ptr,
1338
          weight_ptr,
1339
         bias_ptr,
         N, C, H, W, num_groups, M, eps,
1340
          BLOCK_SIZE: tl.constexpr
1341
1342
          # 1D parallelism across total number of elements
1343
         pid = tl.program_id(0)
          total_elements = N * C * H * W
1344
1345
          block_start = pid * BLOCK_SIZE
1346
          offsets = block_start + tl.arange(0, BLOCK_SIZE)
1347
          mask = offsets < total_elements # bounds checking
1348
1349
          # Flattened indexing to recover n, c from offset
          total2 = H * W
```

```
1350
          total3 = C * total2
1351
1352
         n = offsets // total3
1353
         rest = offsets - n * total3
          c = rest // total2
1354
          group_index = c // M # which group the channel belongs to
1355
1356
          # Memory load with masking
1357
          x\_offsets = offsets
1358
          x_val = tl.load(x_ptr + x_offsets, mask=mask)
1359
          # Load per-sample, per-group mean and variance
1360
          mean_idx = n * num_groups + group_index
1361
          mean_val = tl.load(mean_group_ptr + mean_idx)
1362
          var_val = tl.load(var_group_ptr + mean_idx)
1363
          # Load affine transformation parameters per channel
1364
          gamma = tl.load(weight_ptr + c)
1365
          beta = tl.load(bias_ptr + c)
1366
1367
          # Avoid negative variance (numerical safety)
1368
          var_val = tl.maximum(var_val, 0)
          std = tl.sqrt(var_val + eps)
1369
1370
          # Normalize and apply affine transformation
1371
          normalized = (x_val - mean_val) / std
1372
          out_val = normalized * gamma + beta
1373
          # Store result with mask to handle edge threads
1374
          tl.store(out_ptr + x_offsets, out_val, mask=mask)
1375
1376
1377
       class ModelNew(nn.Module):
1378
          def __init__(self, num_features, num_groups):
             super().__init__()
1379
             self.num_groups = num_groups
1380
             self.num_features = num_features
1381
1382
             # Learnable affine parameters
             self.weight = nn.Parameter(torch.ones(num_features))
1383
             self.bias = nn.Parameter(torch.zeros(num_features))
1384
1385
          def forward(self, x):
1386
             x = x.contiguous() # ensure contiguous memory layout for Triton
1387
             N, C, H, W = x.shape
1388
             M = C // self.num_groups # channels per group
1389
1390
             # Allocate buffers for per-group statistics
1391
             mean_group = torch.empty((N, self.num_groups), device=x.device)
1392
             var_group = torch.empty((N, self.num_groups), device=x.device)
1393
             # Launch reduction kernel: one thread per (N, group)
1394
             grid_reduce = (N, self.num_groups)
1395
             group_norm_reduce_kernel[grid_reduce](
1396
                x, mean_group, var_group,
1397
                N, C, H, W, self.num_groups, M,
                BLOCK_SIZE_REDUCE=1024 # large block size for better throughput
1398
1399
1400
             # Allocate output tensor
1401
             out = torch.empty_like(x)
1402
             total_elements = N * C * H * W
1403
             # Launch forward kernel: 1D block across all elements
```

```
1404
             grid_forward = (triton.cdiv(total_elements, 1024),)
1405
             group_norm_forward_kernel[grid_forward](
1406
                x, out, mean_group, var_group,
                self.weight, self.bias,
1407
                N, C, H, W, self.num_groups, M, 1e-5,
1408
                BLOCK_SIZE=1024 # tuneable block size
1409
1410
1411
             return out
1412
```

1414

1415 1416

1417

1418

1419

1420

1421

1422

1423

1424 1425

1426 1427

1428 1429

1430

1431

1432

1433

1434

1435

1436

1437

1438

1439 1440

1441

1442

1443

1444

1445 1446

1447 1448

1449

1450

1451

1452

HYPERPARAMETERS OF SFT AND RL TRAINING

For data preparation, we use temperature=0.6, top_p=0.95 for DeepSeek-R1 (Guo et al., 2025). To label difficulty, we further label tasks into three difficulty levels using Qwen3-235B-Instruct (Team, 2025) (temperature=0.7, top_p=0.8).

For SFT, we training for 2 epochs with a batch size of 16, a learning rate of 1×10^{-5} , and a maximum sequence length of 12,288 tokens. We train the model for 2 epochs using the VeRL framework (Sheng et al., 2025) with a batch size of 32, a learning rate of 1×10^{-6} , the maximum prompt length 2,048, and the maximum response length 16,384. We use 8 NVIDIA A100 80GB GPUs for both SFT and RL training.

Η ADDITIONAL EXPERIMENT RESULTS

H.1 OVERALL PERFORMANCE COMPARISON

LEVEL1 LEVEL1 (W/O ROBUST) Model #Params compiled / correct compiled / correct valid fast₁ / fast₂ valid fast₁ / fast₂ Owen3 (base) 8B73.0 40.0 / 14.0 0.0 / 0.054.0 52.0 / 15.0 0.0 / 0.014B 65.0 / 17.0 0.070.066.0 71.0 / 16.0 0.0 / 0.0Owen3 82.0 Owen3 32B 75.061.0 / 16.02.0 / 0.052.0 50.0 / 15.0 2.0 / 0.0KernelLLM 8B 42.0 40.0 / 20.0 0.0 / 0.0100.0 98.0 / 29.0 2.0 / 0.078.0 / 50.0 95.0 / 70.0 AutoTriton 8B 97.0 2.0 / 1.0100.0 8.0 / 1.0TRITONRL (ours) 82.0 / 56.0 5.0/1.083.0 / 58.0 5.0 / 1.0 8B 99.0 99.0 w/o RL (SFT only) 8B 97.0 88.0 / 44.0 4.0 / 2.098.0 93.0 / 47.0 5.0 / 2.0GPT-oss 120B 100.0 100.0 / 74.0 7.0 / 2.0100.0 100.0 / 78.0 7.0 / 2.0 Claude-3.7 99.0 / 53.03.0 / 1.0100.0 100.0 / 64.0 8.0 / 1.099.0 6.0 / 2.0 DeepSeek-R1 685B 100.0 100.0 / 66.0 100.0 100.0 / 72.06.0 / 2.0

Table 5: Main results on KernelBench Level 1. All metrics are reported in terms of pass@10 (%). We obtained the best result in model parameter sizes < 32B. The left side shows the results where the validity of generated codes is verified using the robust verifier, checking both syntax and functionality. The right side (w/o robust) shows the results without the robust verifier, where functionality is not checked.

H.2 PASS@K RESULTS

In addition to the pass@10 results shown in the main text, we also report pass@1 and pass@5 results for KernelBench Level 1 and Level 2 tasks in Table 7 and Table 8. These results further demonstrate the strong performance of TRITONRL in generating valid, correct, and efficient Triton code across various pass@k metrics.

Model	#Params		LEVEL2		LEVEL2 (W/O ROBUST)				
Widdel	"I di dilis	valid	compiled / correct	fast ₁ / fast ₂	valid	compiled / correct	fast ₁ / fast ₂		
Qwen3 (base)	8B	56.0	1.0 /0.0	0.0 / 0.0	94.0	52.0 / 11.0	1.0 / 0.0		
Qwen3	14B	35.0	24.0 / 1.0	0.0 / 0.0	100.0	94.0 / 65.0	14.0 / 1.0		
Qwen3	32B	31.0	16.0 / 0.0	0.0 / 0.0	90.0	73.0 / 22.0	7.0 / 1.0		
KernelLLM	8B	0.0	0.0 / 0.0	0.0 / 0.0	98.0	96.0 / 3.0	3.0 / 1.0		
AutoTriton	8B	70.0	3.0 / 0.0	0.0 / 0.0	100.0	97.0 / 76.0	15.0 / 0.0		
TRITONRL (ours)	8B	69.0	41.0 / 10.0	0.0 / 0.0	100.0	88.0 / 42.0	13.0 / 1.0		
w/o RL (SFT only)	8B	67.0	32.0 / 6.0	0.0 / 0.0	100.0	98.0 / 41.0	11.0 / 1.0		
GPT-oss	120B	39.0	38.0 / 12.0	0.0 / 0.0	100.0	99.0 / 74.0	23.0 / 1.0		
Claude-3.7	-	34.0	34.0 / 12.0	1.0 / 0.0	100.0	98.0 / 60.0	18.0 / 1.0		
DeepSeek-R1	685B	30.0	29.0 / 10.0	0.0 / 0.0	100.0	98.0 / 72.0	25.0 / 3.0		

Table 6: Main results on KernelBench level 2 tasks. All metrics are reported in terms of pass@10 (%). We obtained the best result in model parameter sizes < 32B. The left side shows the results where the validity of generated codes is verified using the robust verifier, checking both syntax and functionality. The right side (w/o robust) shows the results without the robust verifier, where functionality is not checked.

Model	PASS@1				PASS@5			PASS@10		
	valid	compiled / correct	fast ₁ / fast ₂	valid	compiled / correct	fast ₁ / fast ₂	valid	compiled / correct	fast ₁ / fast ₂	
Qwen3 (8B)	24.0	14.0 / 4.0	0.0 / 0.0	54.0	33.0 / 14.0	0.0 / 0.0	73.0	40.0 / 14.0	0.0 / 0.0	
Qwen3 (14B)	25.0	22.0 / 10.0	0.0 / 0.0	66.0	53.0 / 14.0	0.0 / 0.0	82.0	65.0 / 17.0	0.0 / 0.0	
Qwen3 (32B)	19.0	16.0 / 4.0	0.0 / 0.0	52.0	41.0 / 14.0	2.0 / 0.0	75.0	61.0 / 16.0	2.0 / 0.0	
KernelLLM	32.0	29.0 / 14.0	0.0 / 0.0	41.0	39.0 / 19.0	0.0 / 0.0	42.0	40.0 / 20.0	0.0 / 0.0	
AutoTriton	61.0	42.0 / 18.0	0.0 / 0.0	94.0	73.0 / 40.0	1.0 / 1.0	97.0	78.0 / 50.0	2.0 / 1.0	
TRITONRL (ours)	70.0	51.0 / 21.0	0.0 / 0.0	97.0	80.0 / 47.0	3.0 / 1.0	99.0	82.0 / 56.0	5.0 / 1.0	
w/o RL (SFT only)	48.0	37.0 /12.0	0.0 / 0.0	94.0	77.0 / 31.0	2.0 / 1.0	97.0	88.0 / 44.0	4.0 / 2.0	
GPT-oss	82.0	82.0 / 38.0	3.0 / 1.0	100.0	100.0 / 64.0	7.0 / 2.0	100.0	100.0 / 74.0	7.0 / 2.0	
Claude-3.7	78.0	73.0 / 25.0	1.0 / 1.0	99.0	98.0 / 40.0	1.0 / 1.0	99.0	99.0 / 53.0	3.0 / 1.0	
DeepSeek-R1	87.0	86.0 / 29.0	1.0 / 0.0	99.0	98.0 / 51.0	4.0 / 2.0	100.0	100.0 / 66.0	6.0 / 2.0	

Table 7: Pass@k performance comparison for k=1,5,10 on KernelBench Level 1 tasks.

Model	PASS@1			PASS@5			PASS@10		
	valid	compiled / correct	fast ₁ / fast ₂	valid	compiled / correct	fast ₁ / fast ₂	valid	compiled / correct	fast ₁ / fast ₂
Qwen3 (8B)	9.0	0.0 / 0.0	0.0 / 0.0	30.0	1.0 / 0.0	0.0 / 0.0	56.0	1.0 / 0.0	0.0 / 0.0
Qwen3 (14B)	4.0	2.0 / 1.0	0.0 / 0.0	23.0	17.0 / 1.0	0.0 / 0.0	35.0	24.0 / 1.0	0.0 / 0.0
Qwen3 (32B)	4.0	3.0 / 0.0	0.0 / 0.0	19.0	9.0 / 0.0	0.0 / 0.0	31.0	16.0 / 0.0	0.0 / 0.0
KernelLLM	0.0	0.0 / 0.0	0.0 / 0.0	0.0	0.0 / 0.0	0.0 / 0.0	0.0	0.0 / 0.0	0.0 / 0.0
AutoTriton	21.0	0.0 / 0.0	0.0 / 0.0	57.0	1.0 / 0.0	0.0 / 0.0	70.0	3.0 / 0.0	0.0 / 0.0
TRITONRL (ours)	16.0	10.0 / 0.0	0.0 / 0.0	56.0	20.0 / 4.0	0.0 / 0.0	71.0	29.0 / 7.0	0.0 / 0.0
w/o RL (SFT only)	15.0	8.0 / 0.0	0.0 / 0.0	51.0	25.0 / 5.0	0.0 / 0.0	67.0	32.0 / 6.0	0.0 / 0.0
GPT-oss	9.0	8.0 / 2.0	0.0 / 0.0	30.0	28.0 / 7.0	0.0 / 0.0	39.0	38.0 / 12.0	0.0 / 0.0
Claude-3.7	15.0	13.0 / 1.0	0.0 / 0.0	31.0	30.0 / 10.0	1.0 / 0.0	34.0	34.0 / 12.0	1.0 / 0.0
DeepSeek-R1	6.0	4.0 / 0.0	0.0 / 0.0	23.0	22.0 / 4.0	0.0 / 0.0	30.0	29.0 / 10.0	0.0 / 0.0

Table 8: Pass@k performance comparison for k = 1, 5, 10 on KernelBench Level 2 tasks.