

PACEvolve: Enabling Long-Horizon Progress-Aware Consistent Evolution

Minghao Yan^{†12} Bo Peng^{*1} Benjamin Coleman^{*3} Ziqi Chen¹ Zhouhang Xie^{†14} Shuo Chen³ Zhankui He³
 Naveen Sachdeva³ Isabella Ye¹ Weili Wang¹ Chi Wang³ Ed H. Chi³ Fernando Pereira³
 Wang-Cheng Kang³ Derek Zhiyuan Cheng³ Beidou Wang¹

Abstract

Large Language Models (LLMs) have emerged as powerful operators for evolutionary search, yet the design of efficient search scaffolds remains ad hoc. While promising, current LLM-in-the-loop systems lack a systematic approach to managing the evolutionary process. We identify three distinct failure modes: *Context Pollution*, where experiment history biases future candidate generation; *Mode Collapse*, where agents stagnate in local minima due to poor exploration-exploitation balance; and *Weak Collaboration*, where rigid crossover strategies fail to leverage parallel search trajectories effectively. We introduce **Progress-Aware Consistent Evolution** (PACEvolve), a framework designed to robustly govern the agent’s context and search dynamics, to address these challenges. PACEvolve combines hierarchical context management (HCM) with pruning to address context pollution; momentum-based backtracking (MBB) to escape local minima; and a self-adaptive sampling policy that unifies backtracking and crossover for dynamic search coordination (CE), allowing agents to balance internal refinement with cross-trajectory collaboration. We demonstrate that PACEvolve provides a systematic path to consistent, long-horizon self-improvement, achieving state-of-the-art results on LLM-SR and KernelBench, and surpassing the record on Modded NanoGPT.

1. Introduction

Large Language Models (LLMs) are increasingly used by evolutionary processes to optimize challenging scientific

and engineering problems (Novikov et al., 2025; Romera-Paredes et al., 2024; Lange et al., 2024; Cheng et al., 2025b; Lange et al., 2025). They transform evolutionary search by replacing the rigid, random operators of classical algorithms (such as mutation and crossover) with intelligent, context-aware reasoning (Novikov et al., 2025; Romera-Paredes et al., 2024; Lange et al., 2025). Unlike traditional Evolutionary Algorithms (EAs) that evaluate an extensive number of weakly-guided candidates ($> 10^6$ samples) (Fogel, 1988; Holland, 1992), LLM-driven agents leverage in-context evolution history to perform iterative refinement (Novikov et al., 2025). By treating the history as a dynamic knowledge base, these agents can theoretically learn from failures and perform meta-reasoning, shifting the paradigm toward sample-efficient, knowledge-guided optimization (Zhai et al., 2025; Lange et al., 2025).

Our work aims to use these intelligent search priors to unlock state-of-the-art performance in complex research and engineering tasks (Shojaee et al., 2024; Ouyang et al., 2025a; Jordan et al., 2024). However, this new LLM-in-the-loop paradigm introduces significant instability, preventing the search from consistently leveraging the LLM’s reasoning capabilities (Xia et al., 2025; Kim et al., 2025). Rather than steadily improving, these systems suffer from high variance (§2), often failing to produce reliable improvements due to the combined stochasticity of the LLM and the search process (Comanici et al., 2025; Renze, 2024).

Despite many successes in applying LLM-based evolutionary search to diverse tasks (Novikov et al., 2025; Lange et al., 2025), we lack a systematic and principled understanding of how to improve the evolution scaffold, often relying on ad hoc designs. In this work, we aim to answer the central research question:

How should we build an agent scaffold for an LLM-driven evolutionary search process?

We identify three core challenges that hinder the performance of modern LLM-assisted evolutionary agents (§2): First, **Context Pollution** overwhelms the agent history with failed hypotheses due to reward sparsity (Liu et al., 2025a), which degrades the quality of generated ideas (Anthony et al., 2025; Zhu et al., 2025). Second, **Mode Collapse**

^{*}Equal contribution. [†]Work done during an internship at Google. ¹Google ²University of Wisconsin-Madison ³Google DeepMind ⁴University of California San Diego. Correspondence to: Minghao Yan <myan@cs.wisc.edu>, Bo Peng <bopeng0823@gmail.com>, Benjamin Coleman <colemanben@google.com>, Ziqi Chen <chenziqi199692@gmail.com>.

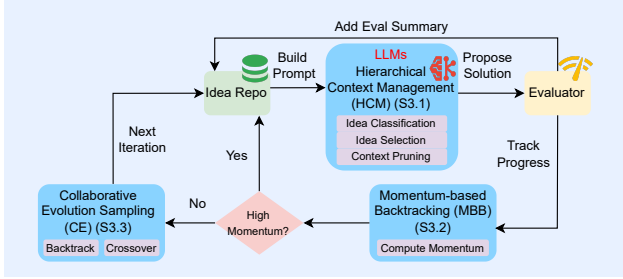


Figure 1. We show the overall workflow of PACEvolve. More details about each module can be found in Figure 2.

occurs when the agent fails to balance exploration and exploitation, leading to stagnation in local minima (Zhang et al., 2025b). Third, **Weak Collaboration** hampers parallel search efficiency because current frameworks lack adaptive mechanisms to transfer knowledge among concurrent processes (Romera-Paredes et al., 2024).

In this paper, we introduce **Progress-Aware Consistent Evolution** (PACEvolve) (Figure 1), a framework that addresses these challenges through a principled, systematic approach. PACEvolve directly tackles the identified challenges via three components (Figure 2): First, we introduce a Hierarchical Context Management (HCM) module to decouple idea generation from selection and mitigate Context Pollution by employing a hierarchical idea memory with context pruning (§3.1). Second, we develop Momentum-based Backtracking (MBB) to combat Mode Collapse, providing a hard escape mechanism that enables the system to break free from local minima (§3.2). Third, we develop Self-adaptive Collaborative Evolution Sampling (CE) to resolve Weak Collaboration; this policy unifies parallel evolution processes by efficiently balancing internal backtracking (deep exploration) with external crossover (knowledge transfer) (§3.3).

We empirically demonstrate that PACEvolve achieves state-of-the-art results, significantly outperforming existing methods on a diverse suite of benchmarks, including Symbolic Regression (LLM-SR), KernelBench, and Modded NanoGPT. Our contributions are summarized as follows:

- We introduce **Hierarchical Context Management (HCM)**, a mechanism that decouples idea generation from selection and applies context pruning. This addresses the challenge of context pollution by ensuring the agent maintains a high signal-to-noise ratio in its evolutionary history, incentivizing diverse idea generation.
- We develop a unified search control policy enabling **Momentum-Based Backtracking (MBB)** and **Self-Adaptive Collaborative Evolution (CE)**. By monitoring search momentum, this policy dynamically balances the trade-off between deep internal exploration (via backtracking) and external knowledge transfer (via crossover), effectively preventing mode collapse across parallel search

processes.

- We demonstrate state-of-the-art empirical performance, significantly outperforming existing methods on diverse and complex benchmarks, including Symbolic Regression (LLM-SR) and KernelBench, and surpassing prior records on Modded NanoGPT.

2. Motivation

Traditional evolutionary algorithms rely on a fixed set of operators, such as mutation and crossover. In contrast, LLM-based search can perform intelligent, context-aware operations, rewriting entire solutions based on a rich prompt that includes past experimental history.

Existing evolutionary agent scaffolds follow an execution-and-reflection paradigm, in which an LLM operates a closed loop comprising idea sampling, execution, feedback collection, and reflection (Novikov et al., 2025; Lange et al., 2025). Though promising, evolutionary agents still yield sub-optimal performance when applied to critical scientific and coding challenges, such as symbolic regression (Shojaee et al., 2024) and kernel design (Ouyang et al., 2025a; Liao et al., 2025). As an example, consider the three independent evolution trajectories on symbolic regression in Figure 3. We observe that if the evolutionary search cannot quickly find a low-NMSE solution, it is unlikely to discover better solutions later. We hypothesize this is due to summarized experiment histories, which serve as context for future iterations, biasing LLMs towards generating similar ideas rather than exploring completely different paths.

In this paper, we conduct a systematic empirical study to identify the key challenges in designing evolutionary agent skeletons, summarized as follows:

1. **Context Pollution disincentivizes diverse candidate generation.** LLM-assisted evolutionary agents rely on their context to guide reflection and idea sampling; context quality is critical to agent performance (Anthony et al., 2025). However, successful discoveries are naturally sparse (Liu et al., 2025a). Consequently, as the agent progresses, the context rapidly saturates with failed attempts (trials yielding no performance gain). As the experimental history grows, a self-reinforcing feedback loop forms, leading LLMs to persist with flawed hypotheses, even in the face of negative results (Anthony et al., 2025; Zhu et al., 2025). This leads to increasingly poor decisions and context explosion, degrading the signal-to-noise ratio. Our empirical observations suggest that these failed trials contribute little to the discovery of innovative and performant solutions (Figure 3). Instead, they cause significant context pollution, diluting the agent’s focus and reducing the probability of finding optimal solutions. LLMs often struggle to balance the refinement

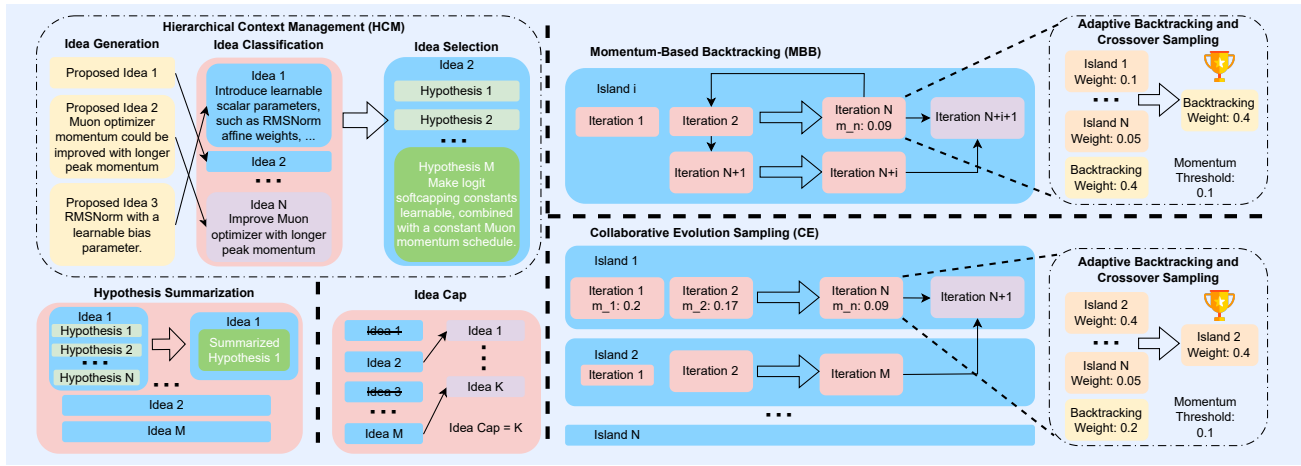


Figure 2. This figure demonstrates the core components of PACEvolve. We decouple idea generation from idea selection to enable easy hierarchical management of idea memory (§3.1). We also design momentum-based self-adaptive backtracking (§3.2) and crossover sampling mechanisms (§3.3) to foster long-horizon reasoning in evolutionary search and escaping local minima.

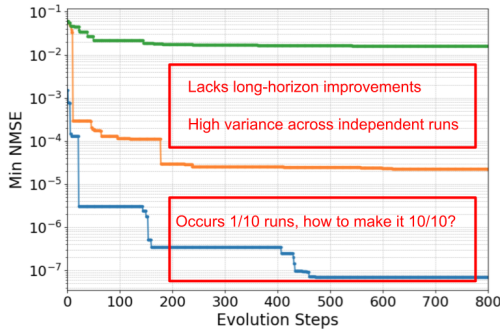


Figure 3. We show three prototypical trajectories from 10 independent trials. If the search process does not converge quickly to a good answer in the first few iterations, it remains in a local minima for the rest of the search. Variance across runs is also large.

of in-context ideas with the need to explore radically different parts of the search space, causing the agent to propose increasingly local candidates (Qin et al., 2025b; Anthony et al., 2025; Agarwal et al.).

2. **Poor exploration-exploitation balance causes Mode Collapse.** To discover innovative solutions, evolutionary agents must explore diverse ideas (Novikov et al., 2025; Real et al., 2017; Hornby et al., 2006). However, our study suggests that evolutionary agents often prefer ideas similar to their context over truly novel ideas. This causes the agent to remain in a self-imposed local minimum (Qin et al., 2025b; Anthony et al., 2025), exploiting known ideas while failing to sufficiently explore new ones (Monea et al., 2024; Chen et al., 2025c; Nie et al., 2024). Similar phenomena are also reported in reinforcement learning agents (Zhu et al., 2025; Zhang et al., 2025b).
3. **Weak Collaboration hurts parallel search efficiency.** Existing paradigms often employ concurrent evolutionary searches (commonly referred to as **multi-island**) to accelerate the evolutionary process (Lange et al., 2025;

Romera-Paredes et al., 2024). The algorithm typically relies on static, periodic crossover patterns where sub-optimal agents are simply replaced with copies of top performers. This rigidity prevents the system from adaptively determining when an agent should incorporate knowledge learned by others.

We show concrete examples of the challenges above in the Appendix A.

3. Progress-Aware Consistent Evolution (PACEvolve)

We introduce PACEvolve, an evolutionary agent framework that explicitly addresses the challenges mentioned above to enable superior solution discovery. Specifically, PACEvolve employs 1) Hierarchical Context Management 3.1 to address Context Pollution; 2) Momentum-Based Backtracking 3.2 to tackle Mode Collapse; and 3) Self-Adaptive Sampling 3.3 to resolve Weak Collaboration. As will be demonstrated in Section 4, by integrating these components, PACEvolve achieves state-of-the-art performance on challenging scientific and real-world engineering tasks. Notations and concepts introduced in this section are summarized in Table 3 in Appendix B.1.

3.1. Hierarchical Context Management (HCM)

To mitigate context pollution while effectively leveraging failed attempts, our HCM features three key designs:

Decomposing High-Level Ideas and Concrete Solutions. Structuring is key to building concise context (Ouyang et al., 2025b; Xu et al., 2025; Wang et al., 2024). PACEvolve disentangles abstract ideas (e.g., “Add Nesterov Momentum to the optimizer”) from specific solutions (e.g., “Try

a specific momentum hyperparameter configuration”) to construct a structured context representation, facilitating context pruning. We implement the structured context using *Macro-Level Conceptual Ideas* to capture global diversity and *Micro-Level Experimental Hypotheses* to refine local details. To support this, we re-architect the search process by decoupling candidate generation into a two-stage process: 1) idea generation and 2) idea selection, supported by a persistent idea pool. The persistent pool acts as an evolving knowledge base for the problem, ensuring that the agent maintains access to a rich, long-term history of conceptual directions. During idea generation, each newly proposed idea runs through an LLM-based classifier to ensure that it is conceptually distinct rather than differing only in minor details. If a conceptual match exists in the persistent pool, the new proposal refines the existing idea; otherwise, it is added as a new entry. We then perform idea selection by granting the agent full access to the knowledge base.

Active Bi-Level Context Pruning To effectively prune unrelated context, we employ a bi-level pruning strategy. At the hypothesis-level, we compress the experimental history associated with each idea. At the idea-level, we identify and actively eliminate ideas with many low-performing hypotheses to encourage the agent to explore innovative directions that are likely to provide high-signal information.

To implement hypothesis-level pruning, we cap the number of hypotheses per idea. Once this limit is reached, a summarization operator is triggered, distilling the accumulated experiment histories into concise key findings. We apply a similar process to cap the number of idea “threads” that are actively considered by the agent, to improve the breadth of the idea pool and force radical exploration. Once the idea cap is reached, the LLM is prompted to discard the least promising directions, thereby encouraging exploration of novel concepts.

Persisting Failures to Permanent Memory One problem with this setup is that a bad idea may be discarded due to low performance, only to be rediscovered and re-explored later. To prevent duplicate solutions and improve sample efficiency, we keep a persistent record of all pruned failures and rejected hypotheses. The agent can effectively filter out known failures by cross-referencing new solutions against this history, ensuring that computational resources are dedicated to exploring novel or high-potential ideas.

In Appendix §B.2, we present the algorithmic details and prompts used for this stage.

3.2. Momentum-Based Backtracking (MBB)

To address mode collapse (Zhang et al., 2025b), where agents over-exploit known solutions at the expense of diversity, we draw inspiration from human problem-solving.

Humans typically pivot when progress stagnates, and we require a similar mechanism for an intelligent evolutionary agent to escape local optima (Balachandran et al., 2025; Yang et al., 2025; Cai et al., 2025a; Liu et al., 2025b; Qin et al., 2025b). The standard approach, fixed-schedule resets (Romera-Paredes et al., 2024; Lange et al., 2025), is inefficient because it ignores the search state. We introduce a **Momentum-Based Backtracking** mechanism (Kingma, 2014) that triggers interventions based on real-time momentum to incentivize exploration and prevent mode collapse.

To effectively detect when a search trajectory stagnates at a local minimum, we require a performance metric that is adaptive to the current scale of the optimization problem. Because the difficulty of achieving further improvements intrinsically increases as the agent approaches the evolutionary target, this metric must adjust for the current difficulty level. To quantify this scale-invariant rate of improvement, we design a new measure that we term **Relative Progress**.

We define the target metric r to be minimized (i.e., r is a lower bound, such as $r = 0$) and let s_t be the best-achieved score at generation t . We define the performance gap G_t as the distance from s_t to the target: $G_t = s_t - r$. The search objective is to drive $G_t \rightarrow 0$. To detect a stagnating trajectory, we track the momentum of its improvement. A simple absolute score formulation ($\Delta t = s_{t-1} - s_t$) is highly dependent on the problem’s scale, motivating the development of a scale-invariant metric: the **Relative Progress** (R_t). When a new best score $s_t < s_{t-1}$ is found at generation t , the Relative Progress R_t is calculated as the fraction of the previous performance gap (G_{t-1}) that has been closed by the new improvement:

$$R_t = \frac{G_{t-1} - G_t}{G_{t-1}} = \frac{(s_{t-1} - r) - (s_t - r)}{s_{t-1} - r} = \frac{s_{t-1} - s_t}{s_{t-1} - r}$$

If no improvement is found ($s_t \geq s_{t-1}$), then $R_t = 0$. This metric is non-negative and represents the fractional reduction in the performance gap, making it inherently adaptive to the current search proximity to r .

We then maintain an Exponentially Weighted Moving Average (EWMA) of this relative progress, which we define as the **Relative Improvement Momentum** (m_t):

$$m_t = \beta \cdot m_{t-1} + (1 - \beta) \cdot R_t$$

where β is the momentum decay factor, which smooths the progress metric over time.

This momentum m_t serves as a direct, adaptive signal of a trajectory’s health. We trigger an intervention when the momentum drops below a predefined stagnation threshold, ϵ_{rel} . When triggered, the agent reverts to an earlier state sampled from a power-law distribution that favors earlier iterations, explicitly unlearning the recent history and reset-

ting the context window. This provides a hard escape from local minima that prompt engineering alone cannot resolve.

3.3. Self-Adaptive Collaborative Evolution (CE)

While the previous sections address individual agent optimization, maximizing search throughput requires parallelization. Existing multi-island frameworks typically employ a static coordination strategy that periodically replaces the worst-performing agents with copies of the best-performing agents (Novikov et al., 2025; Lange et al., 2025; Shojaee et al., 2024). This approach fails to leverage the LLM’s context to make adaptive decisions for knowledge transfer.

To address this, we propose Self-Adaptive Sampling for Collaborative Evolution, which unifies the actions of backtracking and crossover. This framework is designed to trigger automatically when an island stagnates based on momentum and dynamically selects the action that best fosters collaborative evolution between islands.

Once an island i is triggered by MBB (due to $m_{t,i} < \epsilon_{rel}$), we must select an action a from the set $\mathcal{A} = \{\text{Backtrack}\} \cup \{\text{Crossover}_j | j \neq i\}$. The core principle guiding this selection is to prefer the action (backtracking or crossover) that offers the highest potential for global progress.

To design an effective collaboration strategy, we require a global, uniform metric to compare the advancement across all islands. We define the **Absolute Progress** (A_t) for each island as the total fraction of the performance gap it has closed since the beginning of its search:

$$A_t = \frac{s_0 - s_t}{s_0 - r}$$

where s_0 is the initial score of the island and s_t is its current best score. This metric, $A_t \in [0, 1]$, allows us to compare the relative advancement of all islands, regardless of their recent momentum.

We introduce a unified sampling scheme where each action $a \in \mathcal{A}$ is assigned a non-negative weight w_a , based on the $A_{t,i}$ (i.e., absolute progress of island i at time t). The action is then chosen with a probability proportional to its weight, following three key principles:

- **Prioritize High-Reward Knowledge Transfer:** When selecting a crossover partner, the sampling should favor islands j that offer a high potential progress gain ($A_{t,j} > A_{t,i}$).
- **Favor Backtracking for Dominant Agents:** Backtracking should be preferred if the current island i is dominant (i.e., $A_{t,i} \geq A_{t,j}$ for all j), as no other island provides a clear path for improvement.
- **Sensitivity to Global Stagnation:** The decision must be sensitive to progress magnitude when island i and its

best partner j_{best} have similar absolute progress ($A_{t,i} \approx A_{t,\text{best}}$): Similar and low performance (e.g., $A_{t,i} \approx A_{t,\text{best}} \approx 0.1$) indicates shared stagnation, and backtracking should be favored. If both have high progress (e.g., $A_{t,i} \approx A_{t,\text{best}} \approx 0.9$), crossover should be favored, as it suggests potential synergy.

Based on the above principles, we assemble the final sampling probability of choosing any action $a \in \mathcal{A}$:

$$P(a) = \frac{w_a}{w_{BT} + \sum_{j \neq i} w_{C_j}}$$

The details of how the weights are computed, and pseudocode can be found in Appendix B.3. In practice, we set a freeze period at the beginning to build the momentum without backtracking and crossover.

This self-adaptive, momentum-driven framework ensures that the multi-island system efficiently balances internal exploration and external exploitation, thereby maximizing the overall search progress.

4. Experiments

We evaluate PACEvolve using a two-fold approach. First, we validate the effectiveness of our evolutionary algorithm by comparing it directly against state-of-the-art evolutionary frameworks on established benchmarks (Section 4.1). Second, we deploy PACEvolve on complex, open-ended engineering challenges (Section 4.2). Finally, we present an ablation study to quantify the contributions of our core components (Section 4.3).

4.1. Evolutionary Framework Comparison

In this section, we focus on tasks that allow for direct comparison with existing evolutionary search frameworks. We utilize **Symbolic Regression** (Shojaee et al., 2024), which evaluates **scientific reasoning** capability by tasking the agent to recover oscillator acceleration equations from synthetic data, and **KernelBench** (Ouyang et al., 2025a), which evaluates **code optimization** capability by tasking the agent with writing performant custom GPU kernels.

4.1.1. SYMBOLIC REGRESSION

In this experiment, we use the Nonlinear Oscillators task from LLM-SR (Shojaee et al., 2024) to evaluate PACEvolve’s scientific discovery capability. **Nonlinear damped oscillators**, which are widespread in physics and engineering, are described by differential equations that capture the complex interaction among an oscillator’s position, velocity, and the acting forces. We evaluate PACEvolve on discovering a synthetically generated variant by minimizing Normalized Mean Squared Error (NMSE).

Table 1. Performance comparison on LLM-SR task. We report the best, P75, mean, and worst performance across 10 runs in terms of Log10 Normalized Mean Squared Error (Log10 NMSE). Best results are in **bold**, second-best results are underlined.

Method	Best	P75	Mean	Worst
uDSR	-4.06	-4.06	-3.95	-3.73
LLM-SR	-4.80	-4.03	-4.06	-3.62
OpenEvolve	-7.11	-5.79	-5.40	-4.02
CodeEvolve	-7.26	-5.54	-4.97	-3.99
ShinkaEvolve	-6.35	-5.92	-5.35	-3.42
PACEvolve-Single	<u>-8.23</u>	<u>-6.33</u>	<u>-5.87</u>	<u>-4.71</u>
PACEvolve-Multi	-8.24	-7.64	-6.11	-4.73

Experiment Setup We compare PACEvolve against the state-of-the-art evolutionary search frameworks ShinkaEvolve, OpenEvolve, CodeEvolve, LLM-SR, as well as the state-of-the-art non-LLM-based symbolic regression framework uDSR (Landajuela et al., 2022). We run the evolutionary search process for 1000 iterations and repeat each experiment 10 times to obtain a distribution of results. We use the default setup for frameworks that natively support the task (LLM-SR and OpenEvolve). Other baselines and PACEvolve-Multi use a 2-island setup, the default setup ShinkaEvolve used to obtain its state-of-the-art results in Circle Packing. PACEvolve-Single uses a single island setup. We use Gemini 2.5 Pro (Comanici et al., 2025) as the base LLM across all methods and use Gemini 2.5 Flash as the secondary model for frameworks that support model ensembles. For uDSR, we use 10 different random seeds. We report the base-10 logarithm of the Normalized Mean Square Error (NMSE) of the best, worst, P75, and mean across 10 runs for each method. We report the average of the Log-scaled NMSE following (Sharma, 2025), as it better reflects error reduction across independent evolutionary searches. In contrast, the mean NMSE would be dominated by a single bad run (which we also report as the worst Log NMSE).

Results In Table 1, we observe that PACEvolve-Single outperforms every other baseline in the best solution, worst solution, mean log NMSE, as well as P75 log NMSE. When equipped with a multi-island setup, PACEvolve-Multi achieves significant improvements in P75 and mean NMSE, as it discovers 3 solutions with \log_{10} NMSE lower than -8, solidifying our hypothesis that our self-adaptive crossover mechanism can better synergize between islands when at least one of them is performing well.

4.1.2. KERNEL BENCH

KernelBench is a benchmark for developing performant machine learning kernels. In this experiment, we evaluate PACEvolve on KernelBench and demonstrate that

PACEvolve improves on kernels of different difficulties and converges to a higher speedup than other evolutionary search methods.

Evaluation Setup We sample 16 kernels that cover different types of neural network operators and granularity from the KernelBench representative subset, including activation functions (GeLU, Softmax), normalization (RMSNorm, BatchNorm, LayerNorm), operators (MaxPooling, Conv3D, MatMul, Mean Reduction), layers (variants of Conv2D/3D, BMM+Normalization+Residual Add), and models (MLP, RNN, VGG16, AlexNet). A complete list is provided in Appendix C.3. We benchmark the latency on a single A100 40GB GPU using KernelBench’s benchmarking script with modifications to address existing vulnerabilities, such as L2 cache flushing (Ye et al., 2025) and LLM reward hacking (Li et al., 2025a). To mitigate performance variability, we set the GPU frequency to maximum to minimize the impact of DVFS on latency measurements (Le Sueur & Heiser, 2010). We run each method for 1000 iterations per kernel and compare against the PyTorch baseline and the best kernel from the KernelBench Leaderboard v0.1. Due to the higher evaluation cost and the need to cover a wide variety of kernels, we could not repeat the experiments multiple times. To mitigate variance across evolutionary search runs, we compare the individual timing for each of the 16 tested kernels (Table 2). We use Gemini 2.5 Pro across all methods and use Gemini 2.5 Flash as the secondary model for frameworks that support model ensembles.

Results Table 2 shows a per-kernel breakdown of PACEvolve performance. PACEvolve-Single outperforms PyTorch in all but two cases (MLP and Matmul with Large k). The kernels underpinning matrix multiplication in PyTorch are heavily optimized (NVIDIA, 2025), while an MLP can be viewed as a stack of matrix multiplications. PACEvolve-Multi discovered solutions to outperform Torch base in MLP, while achieving near-parity performance with PyTorch on Matrix Multiplication.

Table 2 shows that both the single island (PACEvolve-Single) and the multi-island (PACEvolve-Multi) version of PACEvolve outperform the best existing kernels on KernelBench in all tested cases. In addition, PACEvolve-Single outperforms ShinkaEvolve in all tested kernels, where PACEvolve-Multi further outperforms PACEvolve-Single in 81.25% (13/16) of the tested kernels. When comparing against other evolutionary frameworks, PACEvolve-Multi found equivalent or better kernels compared to ShinkaEvolve and CodeEvolve in 14/16 cases and OpenEvolve in 15/16 cases, clearly demonstrating better framework design despite possible variances in individual runs.

In Appendix C.3, we report the head-to-head win rate across frameworks (Figure 5) and PACEvolve-generated kernels in Appendix D.

Table 2. Kernel performance comparison. Unit: microseconds. PACE-Single and PACE-Multi represent the single and multi-island versions of PACEvolve. We also report the maximum speedup PACEvolve achieves over the PyTorch baseline in the last column. Best results are in **bold**, second-best results are underlined.

Kernel	PyTorch	KBench	Shinka	OpenEvo	CodeEvo	PACE-Single	PACE-Multi	Max Sp.
BatchNorm	1.10×10^0	6.83×10^{-1}	4.76×10^{-1}	1.10×10^0	3.91×10^{-1}	3.99×10^{-1}	3.91×10^{-1}	(2.81x)
ConvDiv	1.11×10^0	9.51×10^{-1}	8.25×10^{-1}	9.37×10^{-1}	8.47×10^{-1}	<u>7.51×10^{-1}</u>	6.09×10^{-1}	(1.82x)
ConvMax	1.22×10^0	1.15×10^0	<u>6.72×10^{-1}</u>	1.02×10^0	1.04×10^0	3.09×10^{-1}	7.13×10^{-1}	(3.95x)
ConvT2D	5.04×10^{-1}	3.35×10^{-1}	<u>2.94×10^{-1}</u>	2.86×10^{-1}	3.64×10^{-1}	<u>2.47×10^{-1}</u>	1.81×10^{-1}	(2.79x)
GELU	8.81×10^{-3}	9.01×10^{-3}	9.32×10^{-3}	<u>8.70×10^{-3}</u>	8.91×10^{-3}	7.88×10^{-3}	<u>8.70×10^{-3}</u>	(1.12x)
MatMul	1.03×10^0	6.49×10^0	1.34×10^0	<u>6.35×10^0</u>	1.15×10^0	1.10×10^0	<u>1.06×10^0</u>	(0.97x)
MaxPool	5.52×10^{-2}	4.83×10^{-2}	4.33×10^{-2}	4.34×10^{-2}	4.71×10^{-2}	<u>4.26×10^{-2}</u>	4.24×10^{-2}	(1.30x)
MLP	3.51×10^{-2}	1.69×10^{-1}	1.06×10^{-1}	<u>2.20×10^{-2}</u>	6.25×10^{-3}	<u>8.63×10^{-2}</u>	2.27×10^{-2}	(1.55x)
RMSNorm	1.03×10^0	4.48×10^0	1.03×10^0	5.64×10^{-1}	4.03×10^{-1}	<u>3.99×10^{-1}</u>	3.97×10^{-1}	(2.59x)
Softmax	1.69×10^{-2}	3.32×10^{-2}	1.62×10^{-2}	1.60×10^{-2}	1.14×10^{-2}	<u>1.36×10^{-2}</u>	1.14×10^{-2}	(1.48x)
VGG16	6.39×10^0	1.16×10^1	<u>2.61×10^0</u>	6.84×10^0	3.87×10^0	2.16×10^0	3.80×10^0	(2.96x)
MeanRedu	2.01×10^{-2}	4.42×10^{-2}	<u>2.01×10^{-2}</u>	1.98×10^{-2}	<u>1.25×10^{-2}</u>	1.35×10^{-2}	1.17×10^{-2}	(1.72x)
RNN	3.49×10^{-2}	1.58×10^1	3.79×10^{-2}	9.39×10^{-2}	<u>3.40×10^{-2}</u>	<u>2.97×10^{-2}</u>	2.58×10^{-2}	(1.35x)
BMMNorm	2.85×10^{-2}	2.60×10^{-2}	1.48×10^{-2}	1.50×10^{-2}	9.52×10^{-3}	1.27×10^{-2}	<u>1.21×10^{-2}</u>	(2.36x)
AlexNet	9.76×10^{-1}	3.31×10^0	6.65×10^{-1}	9.36×10^{-1}	6.02×10^{-1}	<u>5.41×10^{-1}</u>	5.02×10^{-1}	(1.94x)
LayerNorm	1.01×10^1	1.28×10^2	6.14×10^{-1}	9.47×10^{-1}	5.96×10^{-1}	<u>5.90×10^{-1}</u>	5.81×10^{-1}	(17.38x)

4.2. Automated Engineering in Complex Environments

Research innovation often requires interacting with complex environments that are not natively supported by existing evolutionary frameworks. To address this, we develop a versatile integration platform that supports arbitrary tasks within a sandboxed environment, enabling us to extend evolutionary search to any full-stack research and engineering challenges. In this section, we utilize our platform to evaluate PACEvolve on the **Modded NanoGPT** benchmark. We demonstrate how PACEvolve can be used in complex environments to accelerate and automate research and engineering efforts.

Modded NanoGPT In this task, we deploy PACEvolve on the Modded NanoGPT benchmark (Jordan et al., 2024), which represents a higher level of complexity, as the agent can optimize the model architecture, the distributed training pipeline, and improve kernels. Prior tasks, such as Symbolic Regression (Shojaee et al., 2024) and KernelBench (Ouyang et al., 2025a), focus on a single concrete challenge, whereas in Modded NanoGPT (Jordan et al., 2024), PACEvolve is tasked with making any adjustments to improve training efficiency. This includes, but is not limited to, changes to model architecture, more efficient kernels, data processing, and communication.

Evaluation Setup We use the recommended setup in Modded NanoGPT (Jordan et al., 2024) and evaluate how long it takes to reach a validation loss of 3.28 on FineWeb (Penedo et al., 2024) using 8 H100 GPUs. We use Gemini 3 Pro (Team, 2025) as the backbone LLM.

Results PACEvolve discovers improvements upon a heavily optimized state-of-the-art (version 40 of Modded

NanoGPT (Jordan et al., 2024)) from both systems and a modeling perspective. We summarize the innovations found below:

1. PACEvolve discovers an inefficiency in training data loading and proposes to shard data across ranks and pre-load data to minimize GPU idle time. This reduces the training time from **142.8s** to **141.9s** (over 2330 steps).
2. PACEvolve introduces a U-shaped initialization for skip connections where the weights start high (0.8), dip in the middle layers (0.4), and rise again (0.8). This technique reduces the training time from **141.9s** to **141.5s**.
3. PACEvolve optimizes a series of hyperparameters, such as logit softcapping (Team et al., 2024), beta1 in Distributed Adam (Rajbhandari et al., 2020; Loshchilov & Hutter, 2017), alpha in YaRN (Peng et al., 2023), and token smearing lambda. When combined, these hyperparameter updates reduce the training time from **141.5s** to **140.8s**.
4. PACEvolve improves dynamic context window scheduling, training more aggressively on a smaller context window (for 45% of training) while also increasing the maximum context window length for the last 10% of training. This further reduces the training time from **140.8s** to **140.2s**.

While this improvement looks incremental, the Modded NanoGPT benchmark has been heavily optimized from ~ 2700 seconds (version 1) to 142 seconds (version 40). Therefore, any further gain represents a substantial achievement. This result demonstrates PACEvolve’s capability in conducting research tasks in complex environments.

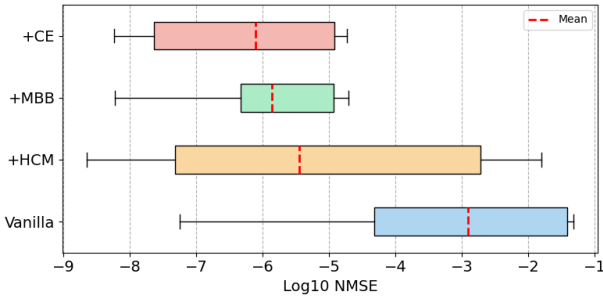


Figure 4. Cumulative Boxplot Comparison of PACEvolve Techniques. The distribution of performance across 10 runs is shown, starting with vanilla append-only context management and progressively adding each optimization technique.

4.3. Ablation studies

We use the Symbolic Regression task to dissect how each component in PACEvolve contributes to overall performance. First, we run the evolutionary search process for 1000 iterations and repeat each experiment 10 times to obtain a distribution of results on the Nonlinear Oscillator task. The vanilla implementation appends a summary of the proposed solution and experiment results of each iteration during the evolution. We then gradually add hierarchical context management, momentum-based backtracking, and adaptive cross-island sampling, isolating the contribution of each technique.

Results Figure 4 demonstrates the progression of improvements after adding each technique. Adding hierarchical context erasure significantly improves the mean and best-performing evolutionary processes. However, the worst-performing processes still made insufficient progress. Hierarchical context erasure only erases low-performing ideas and experiment histories during the process; however, if the best-performing idea in the search process is a local minimum, it is never eliminated and will continue to affect future iterations. This also demonstrates the necessity of backtracking.

After adding momentum-based backtracking to the hierarchical context, we eliminated the low-performing processes. However, this slightly affected the best-performing processes, as backtracking forces the evolutionary search to explore more often and reduces the search budget devoted to exploitative ideas that LLMs prefer without external intervention (Zhu et al., 2025; Anthony et al., 2025).

We then integrate self-adaptive cross-island sampling into backtracking to obtain the best of both worlds. As shown in Figure 4, this preserves backtracking’s advantage of eliminating low-performing processes while significantly improving mean and P75 performance. Our self-adaptive cross-island sampling enables multiple concurrent explorations of the high-reward (in this case, low NMSE) regime, provided that at least one island discovers a promising direction.

5. Related Works

LLMs have redefined how evolutionary search is performed (Holland, 1992; Fogel, 1988; Qian et al., 2024). Moving beyond the traditional requirement of defining a fixed set of mutation and crossover operators (Chen et al., 2023; Real et al., 2017; Lehman et al., 2023; Meyerson et al., 2024), modern evolutionary agents use LLMs prompted with background knowledge and past results to intelligently propose improvements to solutions in each iteration, granting them far more flexibility and reasoning power (Novikov et al., 2025; Lange et al., 2025; Romera-Paredes et al., 2024; Wang et al., 2025). This paradigm has yielded significant successes across various domains, including discrete mathematics (Romera-Paredes et al., 2024; Lange et al., 2025), kernel optimization (Novikov et al., 2025), and general code optimization (Lehman et al., 2023; Meyerson et al., 2024; Cai et al., 2025b; Cheng et al., 2025a).

Current SOTA agents, notably AlphaEvolve and ShinkaEvolve, have focused on improving the quality of the LLM’s context by summarizing past trials (Lange et al., 2025; Novikov et al., 2025) and increasing the number of in-context examples (Assumpção et al., 2025). Specifically, AlphaEvolve introduced a reflection mechanism to derive insights from successful and failed attempts; ShinkaEvolve proposed a method for summarizing past attempts to maintain context within the LLM’s finite window. Our work moves beyond mere context aggregation to focus on context management and dynamic search control. PACEvolve analyzes current failure modes and introduces a principled recipe that significantly outperforms existing approaches.

In addition to evolutionary search processes, other types of iterative search and fine-tuning methods (Chen et al., 2025b; Deng et al., 2025; Chen et al., 2025a; Qin et al., 2025a; Zhang et al., 2025a) that leverage tree (Li et al., 2025b; Li, 2025; Jiang et al., 2025a; Chi et al., 2024) and rubric-based techniques (Gunjal et al., 2025; Huang et al., 2025; Jiang et al., 2025b; Goel et al., 2025) have been developed to improve LLM-based agents in diverse scenarios. Another line of work investigates how to improve LLM reasoning in long-context settings (Li et al., 2024; Bai et al., 2025; Zhou et al., 2025; Motwani et al., 2025). While uncovering similar phenomena (Zhu et al., 2025; Zhang et al., 2025b; Anthony et al., 2025), our work advances the field by identifying the best practices for designing evolutionary search.

6. Conclusion

The rise of LLM-assisted evolutionary search has opened a new frontier in optimization and scientific discovery. In this paper, we introduce PACEvolve, a framework designed to address these core challenges through a principled evolutionary recipe: a Hierarchical Context Management (HCM) mod-

ule to promote idea diversity and prune ineffective histories, Momentum-based Backtracking (MBB) to provide consistent escape from local minima, and a Collaborative Evolution Sampling policy (CE) for self-adaptive multi-island coordination. Through extensive empirical evaluation on diverse, complex benchmarks, we demonstrate that PACEvolve consistently achieves state-of-the-art performance and offers a principled approach to develop robust LLM-in-the-loop evolutionary agents.

References

- Agarwal, D., Majumder, B. P., Adamson, R., Chakravorty, M., Gavireddy, S. R., Parashar, A., Surana, H., Mishra, B. D., McCallum, A., Sabharwal, A., et al. Autodiscovery: Open-ended scientific discovery via bayesian surprise. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*.
- Agrawal, L. A., Tan, S., Soylu, D., Ziemis, N., Khare, R., Opsahl-Ong, K., Singhvi, A., Shandilya, H., Ryan, M. J., Jiang, M., et al. Gepa: Reflective prompt evolution can outperform reinforcement learning. *arXiv preprint arXiv:2507.19457*, 2025.
- Anthony, G.-C., Lin, D., Samiei, M., Precup, D., Richards, B. A., Fergus, R., and Marino, K. Language agents mirror human causal reasoning biases. how can we help them think like scientists? In *Second Conference on Language Modeling*, 2025.
- Assumpção, H., Ferreira, D., Campos, L., and Murai, F. Codeevolve: An open source evolutionary coding agent for algorithm discovery and optimization. *arXiv preprint arXiv:2510.14150*, 2025.
- Bai, Y., Tu, S., Zhang, J., Peng, H., Wang, X., Lv, X., Cao, S., Xu, J., Hou, L., Dong, Y., et al. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 3639–3664, 2025.
- Balachandran, V., Chen, J., Chen, L., Garg, S., Joshi, N., Lara, Y., Langford, J., Nushi, B., Vineet, V., Wu, Y., et al. Inference-time scaling for complex tasks: Where we stand and what lies ahead. *arXiv preprint arXiv:2504.00294*, 2025.
- Cai, H. J., Wang, J., Chen, X., and Dhingra, B. How much backtracking is enough? exploring the interplay of sft and rl in enhancing llm reasoning. *arXiv preprint arXiv:2505.24273*, 2025a.
- Cai, Z., Guo, X., Pei, Y., Feng, J., Chen, J., Zhang, Y.-Q., Ma, W.-Y., Wang, M., and Zhou, H. Flex: Continuous agent evolution via forward learning from experience. *arXiv preprint arXiv:2511.06449*, 2025b.
- Chen, G., Qiao, Z., Chen, X., Yu, D., Xu, H., Zhao, W. X., Song, R., Yin, W., Yin, H., Zhang, L., et al. Iterresearch: Rethinking long-horizon agents via markovian state reconstruction. *arXiv preprint arXiv:2511.07327*, 2025a.
- Chen, J., Coleman, B., and Shrivastava, A. Revisiting consistent hashing with bounded loads. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 3976–3983, 2021.
- Chen, X., Liang, C., Huang, D., Real, E., Wang, K., Pham, H., Dong, X., Luong, T., Hsieh, C.-J., Lu, Y., et al. Symbolic discovery of optimization algorithms. *Advances in neural information processing systems*, 36:49205–49233, 2023.
- Chen, Y., Wang, Y., Zhu, S., Yu, H., Feng, T., Zhang, M., Patwary, M., and You, J. Multi-agent evolve: Llm self-improve through co-evolution. *arXiv preprint arXiv:2510.23595*, 2025b.
- Chen, Z., Qin, X., Wu, Y., Ling, Y., Ye, Q., Zhao, W. X., and Shi, G. Pass@ k training for adaptively balancing exploration and exploitation of large reasoning models. *arXiv preprint arXiv:2508.10751*, 2025c.
- Cheng, A., Liu, S., Pan, M., Li, Z., Wang, B., Krentsel, A., Xia, T., Cemri, M., Park, J., Yang, S., et al. Barbarians at the gate: How ai is upending systems research. *arXiv preprint arXiv:2510.06189*, 2025a.
- Cheng, J., Clark, P., and Richardson, K. Language modeling by language models, 2025b. URL <https://arxiv.org/abs/2506.20249>.
- Chi, Y., Lin, Y., Hong, S., Pan, D., Fei, Y., Mei, G., Liu, B., Pang, T., Kwok, J., Zhang, C., et al. Sela: Tree-search enhanced llm agents for automated machine learning. *arXiv preprint arXiv:2410.17238*, 2024.
- Comanici, G., Bieber, E., Schaekermann, M., Pasupat, I., Sachdeva, N., Dhillon, I., Blistein, M., Ram, O., Zhang, D., Rosen, E., et al. Gemini 2.5: Pushing the frontier with advanced reasoning, multimodality, long context, and next generation agentic capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- Deng, Y., Hsu, I., Yan, J., Wang, Z., Han, R., Zhang, G., Chen, Y., Wang, W., Pfister, T., Lee, C.-Y., et al. Supervised reinforcement learning: From expert trajectories to step-wise reasoning. *arXiv preprint arXiv:2510.25992*, 2025.
- Fogel, D. B. An evolutionary approach to the traveling salesman problem. *Biological Cybernetics*, 60(2):139–144, 1988.

- Goel, S., Hazra, R., Jayalath, D., Willi, T., Jain, P., Shen, W. F., Leontiadis, I., Barbieri, F., Bachrach, Y., Geiping, J., et al. Training ai co-scientists using rubric rewards. *arXiv preprint arXiv:2512.23707*, 2025.
- Gunjal, A., Wang, A., Lau, E., Nath, V., He, Y., Liu, B., and Hendryx, S. Rubrics as rewards: Reinforcement learning beyond verifiable domains. *arXiv preprint arXiv:2507.17746*, 2025.
- Holland, J. H. Genetic algorithms. *Scientific american*, 267 (1):66–73, 1992.
- Hornby, G., Globus, A., Linden, D., and Lohn, J. Automated antenna design with evolutionary algorithms. In *Space 2006*, pp. 7242. 2006.
- Huang, Z., Zhuang, Y., Lu, G., Qin, Z., Xu, H., Zhao, T., Peng, R., Hu, J., Shen, Z., Hu, X., et al. Reinforcement learning with rubric anchors. *arXiv preprint arXiv:2508.12790*, 2025.
- Jiang, M., Lupu, A., and Bachrach, Y. Bootstrapping task spaces for self-improvement. *arXiv preprint arXiv:2509.04575*, 2025a.
- Jiang, Y., Jiang, L., Teney, D., Moor, M., and Brbic, M. Meta-rl induces exploration in language agents. *arXiv preprint arXiv:2512.16848*, 2025b.
- Jordan, K., Bernstein, J., Rappazzo, B., @fern-bear.bsky.social, Vlado, B., Jiacheng, Y., Cesista, F., Koszarsky, B., and @Grad62304977. modded-nanogpt: Speedrunning the nanogpt baseline, 2024. URL <https://github.com/KellerJordan/modded-nanogpt>.
- Kim, Y., Gu, K., Park, C., Park, C., Schmidgall, S., Heydari, A. A., Yan, Y., Zhang, Z., Zhuang, Y., Malhotra, M., Liang, P. P., Park, H. W., Yang, Y., Xu, X., Du, Y., Patel, S., Althoff, T., McDuff, D., and Liu, X. Towards a science of scaling agent systems, 2025. URL <https://arxiv.org/abs/2512.08296>.
- Kingma, D. P. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Landajuela, M., Lee, C. S., Yang, J., Glatt, R., Santiago, C. P., Aravena, I., Mundhenk, T., Mulcahy, G., and Petersen, B. K. A unified framework for deep symbolic regression. *Advances in Neural Information Processing Systems*, 35:33985–33998, 2022.
- Lange, R., Tian, Y., and Tang, Y. Large language models as evolution strategies. In *Proceedings of the Genetic and Evolutionary Computation Conference Companion*, pp. 579–582, 2024.
- Lange, R. T., Imajuku, Y., and Cetin, E. Shinkaevolve: Towards open-ended and sample-efficient program evolution. *arXiv preprint arXiv:2509.19349*, 2025.
- Le Sueur, E. and Heiser, G. Dynamic voltage and frequency scaling: The laws of diminishing returns. In *Proceedings of the 2010 international conference on Power aware computing and systems*, pp. 1–8, 2010.
- Lehman, J., Gordon, J., Jain, S., Ndousse, K., Yeh, C., and Stanley, K. O. Evolution through large models. In *Handbook of evolutionary machine learning*, pp. 331–366. Springer, 2023.
- Li, T., Zhang, G., Do, Q. D., Yue, X., and Chen, W. Long-context llms struggle with long in-context learning. *arXiv preprint arXiv:2404.02060*, 2024.
- Li, X., Sun, X., Wang, A., Li, J., and Shum, C. Cuda-11: Improving cuda optimization via contrastive reinforcement learning. *arXiv preprint arXiv:2507.14111*, 2025a.
- Li, Y. Policy guided tree search for enhanced llm reasoning. *arXiv preprint arXiv:2502.06813*, 2025.
- Li, Y., Gu, Q., Wen, Z., Li, Z., Xing, T., Guo, S., Zheng, T., Zhou, X., Qu, X., Zhou, W., et al. Treepo: Bridging the gap of policy optimization and efficacy and inference efficiency with heuristic tree-based modeling. *arXiv preprint arXiv:2508.17445*, 2025b.
- Liao, G., Qin, H., Wang, Y., Golden, A., Kuchnik, M., Yetim, Y., Ang, J. J., Fu, C., He, Y., Hsia, S., et al. Kernelevolve: Scaling agentic kernel coding for heterogeneous ai accelerators at meta. *arXiv preprint arXiv:2512.23236*, 2025.
- Liu, F., Zhang, Q., Shi, J., Tong, X., Mao, K., and Yuan, M. Fitness landscape of large language model-assisted automated algorithm search. *arXiv preprint arXiv:2504.19636*, 2025a.
- Liu, Z., Wang, J., Song, L., and Bian, J. Sample-efficient llm optimization with reset replay. *arXiv preprint arXiv:2508.06412*, 2025b.
- Loshchilov, I. and Hutter, F. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- Meyerson, E., Nelson, M. J., Bradley, H., Gaier, A., Moradi, A., Hoover, A. K., and Lehman, J. Language model crossover: Variation through few-shot prompting. *ACM Transactions on Evolutionary Learning*, 4(4):1–40, 2024.
- Monea, G., Bosselut, A., Brantley, K., and Artzi, Y. Llms are in-context bandit reinforcement learners. *arXiv preprint arXiv:2410.05362*, 2024.

- Motwani, S. R., Ivanova, A., Cai, Z., Torr, P., Islam, R., Shah, S., de Witt, C. S., and London, C. h1: Bootstrapping llms to reason over longer horizons via reinforcement learning. *arXiv preprint arXiv:2510.07312*, 2025.
- Nie, A., Su, Y., Chang, B., Lee, J. N., Chi, E. H., Le, Q. V., and Chen, M. Evolve: Evaluating and optimizing llms for exploration. *arXiv preprint arXiv:2410.06238*, 2024.
- Novikov, A., Vű, N., Eisenberger, M., Dupont, E., Huang, P.-S., Wagner, A. Z., Shirobokov, S., Kozlovskii, B., Ruiz, F. J., Mehrabian, A., et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint arXiv:2506.13131*, 2025.
- NVIDIA. cublas 13.0, Oct 2025. URL <https://docs.nvidia.com/cuda/cublas/>.
- Ouyang, A., Guo, S., Arora, S., Zhang, A. L., Hu, W., Ré, C., and Mirhoseini, A. Kernelbench: Can llms write efficient gpu kernels?, 2025a. URL <https://arxiv.org/abs/2502.10517>.
- Ouyang, S., Yan, J., Hsu, I., Chen, Y., Jiang, K., Wang, Z., Han, R., Le, L. T., Daruki, S., Tang, X., et al. Reasoningbank: Scaling agent self-evolving with reasoning memory. *arXiv preprint arXiv:2509.25140*, 2025b.
- Penedo, G., Kydlíček, H., Lozhkov, A., Mitchell, M., Raffel, C. A., Von Werra, L., Wolf, T., et al. The fineweb datasets: Decanting the web for the finest text data at scale. *Advances in Neural Information Processing Systems*, 37: 30811–30849, 2024.
- Peng, B., Quesnelle, J., Fan, H., and Shippole, E. Yarn: Efficient context window extension of large language models. *arXiv preprint arXiv:2309.00071*, 2023.
- Qian, C., Xue, K., and Wang, R.-J. Quality-diversity algorithms can provably be helpful for optimization. *arXiv preprint arXiv:2401.10539*, 2024.
- Qin, R., He, W., Huang, W., Zhang, Y., Zhao, Y., Pang, B., Xu, X., Shan, Y., Wu, Y., and Zhang, M. Seer: Online context learning for fast synchronous llm reinforcement learning. *arXiv preprint arXiv:2511.14617*, 2025a.
- Qin, T., Alvarez-Melis, D., Jelassi, S., and Malach, E. To backtrack or not to backtrack: When sequential search limits model reasoning. *arXiv preprint arXiv:2504.07052*, 2025b.
- Rajbhandari, S., Rasley, J., Ruwase, O., and He, Y. Zero: Memory optimizations toward training trillion parameter models. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pp. 1–16. IEEE, 2020.
- Real, E., Moore, S., Selle, A., Saxena, S., Suematsu, Y. L., Tan, J., Le, Q. V., and Kurakin, A. Large-scale evolution of image classifiers. In *International conference on machine learning*, pp. 2902–2911. PMLR, 2017.
- Renze, M. The effect of sampling temperature on problem solving in large language models. In *Findings of the association for computational linguistics: EMNLP 2024*, pp. 7346–7356, 2024.
- Romera-Paredes, B., Barekatin, M., Novikov, A., Balog, M., Kumar, M. P., Dupont, E., Ruiz, F. J., Ellenberg, J. S., Wang, P., Fawzi, O., et al. Mathematical discoveries from program search with large language models. *Nature*, 625 (7995):468–475, 2024.
- Sharma, A. Openevolve: an open-source evolutionary coding agent, 2025. URL <https://github.com/algorithmicsuperintelligence/openevolve>.
- Shojaee, P., Meidani, K., Gupta, S., Farimani, A. B., and Reddy, C. K. Llm-sr: Scientific equation discovery via programming with large language models. *arXiv preprint arXiv:2404.18400*, 2024.
- Team, G. . Gemini 3, Nov 2025. URL <https://blog.google/products/gemini/gemini-3/>.
- Team, G., Riviere, M., Pathak, S., Sessa, P. G., Hardin, C., Bhupatiraju, S., Hussenot, L., Mesnard, T., Shahriari, B., Ramé, A., et al. Gemma 2: Improving open language models at a practical size. *arXiv preprint arXiv:2408.00118*, 2024.
- Wang, R., Shivanna, R., Cheng, D., Jain, S., Lin, D., Hong, L., and Chi, E. Dcn v2: Improved deep & cross network and practical lessons for web-scale learning to rank systems. In *Proceedings of the web conference 2021*, pp. 1785–1797, 2021.
- Wang, Y., Su, S.-R., Zeng, Z., Xu, E., Ren, L., Yang, X., Huang, Z., He, X., Ma, L., Peng, B., et al. Thetaevolve: Test-time learning on open problems. *arXiv preprint arXiv:2511.23473*, 2025.
- Wang, Z. Z., Mao, J., Fried, D., and Neubig, G. Agent workflow memory. *arXiv preprint arXiv:2409.07429*, 2024.
- Xia, P., Zeng, K., Liu, J., Qin, C., Wu, F., Zhou, Y., Xiong, C., and Yao, H. Agent0: Unleashing self-evolving agents from zero data via tool-integrated reasoning. *arXiv preprint arXiv:2511.16043*, 2025.
- Xu, W., Liang, Z., Mei, K., Gao, H., Tan, J., and Zhang, Y. A-mem: Agentic memory for llm agents. *arXiv preprint arXiv:2502.12110*, 2025.

- Yang, X.-W., Zhu, X.-Y., Wei, W.-D., Zhang, D.-C., Shao, J.-J., Zhou, Z., Guo, L.-Z., and Li, Y.-F. Step back to leap forward: Self-backtracking for boosting reasoning of language models. *arXiv preprint arXiv:2502.04404*, 2025.
- Ye, Z., Chen, L., Lai, R., Lin, W., Zhang, Y., Wang, S., Chen, T., Kasikci, B., Grover, V., Krishnamurthy, A., et al. Flash-infer: Efficient and customizable attention engine for llm inference serving. *arXiv preprint arXiv:2501.01005*, 2025.
- Zhai, Y., Tao, S., Chen, C., Zou, A., Chen, Z., Fu, Q., Mai, S., Yu, L., Deng, J., Cao, Z., et al. Agentevolver: Towards efficient self-evolving agent system. *arXiv preprint arXiv:2511.10395*, 2025.
- Zhang, G., Ren, H., Zhan, C., Zhou, Z., Wang, J., Zhu, H., Zhou, W., and Yan, S. Memevolve: Meta-evolution of agent memory systems. *arXiv preprint arXiv:2512.18746*, 2025a.
- Zhang, J., Yu, S., Chong, D., Sicilia, A., Tomz, M. R., Manning, C. D., and Shi, W. Verbalized sampling: How to mitigate mode collapse and unlock llm diversity. *arXiv preprint arXiv:2510.01171*, 2025b.
- Zhou, Z., Qu, A., Wu, Z., Kim, S., Prakash, A., Rus, D., Zhao, J., Low, B. K. H., and Liang, P. P. Mem1: Learning to synergize memory and reasoning for efficient long-horizon agents. *arXiv preprint arXiv:2506.15841*, 2025.
- Zhu, J., Liu, J., Yang, S., Zhang, Q., and He, X. Open benchmarking for click-through rate prediction. In *Proceedings of the 30th ACM international conference on information & knowledge management*, pp. 2759–2769, 2021.
- Zhu, J., Dai, Q., Su, L., Ma, R., Liu, J., Cai, G., Xiao, X., and Zhang, R. Bars: Towards open benchmarking for recommender systems. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 2912–2923, 2022.
- Zhu, K., Liu, Z., Li, B., Tian, M., Yang, Y., Zhang, J., Han, P., Xie, Q., Cui, F., Zhang, W., et al. Where llm agents fail and how they can learn from failures. *arXiv preprint arXiv:2509.25370*, 2025.

A. Failure Analysis

We demonstrate that these challenges actively hinder performance in practice through a case study on a Symbolic Regression task (LLM-SR) (Shojaee et al., 2024). We consider the Nonlinear Oscillators benchmark (see §4.1.1 for details), where the goal is to discover the underlying differential equation governing a system’s motion, which takes the general form: $\ddot{x} + f(t, x, \dot{x}) = 0$. The agent must find the symbolic form of f that minimizes the Normalized Mean Squared Error (NMSE).

Our experiments with the vanilla setup, in which we append the summarized experimental history to the context, reveal that most evolutionary searches become trapped in local minima. This failure stems directly from the challenges identified. Prior works reveal that modern LLMs exhibit reasoning biases; they often persist with a flawed hypothesis even when presented with negative results (Zhu et al., 2025; Anthony et al., 2025).

Insight 1: Hierarchical Context Management is key to enabling innovative thinking

First, to enable easy context management of past trials and results, we introduce a **hierarchical persistent idea memory with context erasure** (§3.1). This technique decouples conceptual idea generation from the selection of one for experimentation. It periodically summarizes or trims ineffective experimental paths, forcing the search to maintain diversity and avoiding context explosion.

Figure 3 shows three independent trajectories in evolution. A common failure mode is that if the evolutionary agent cannot quickly find a solution with low NMSE, it is very unlikely to discover better solutions later in the evolution. We hypothesize that this is due to summarized experiment histories, which serve as context for future iterations, conditioning LLMs to generate similar ideas rather than exploring completely different paths.

Insight 2: Progress-aware regret is the key to escaping local minima

While context erasure helps manage a stagnating line of evolution, it doesn’t fully solve the problem of an island getting stuck in a local minimum (e.g., stagnating trajectories in Figure 3). Recent works also show that LLMs cannot perform exploration-exploitation tradeoff in-context (Nie et al., 2024; Monea et al., 2024) efficiently. To provide a hard escape mechanism, we introduce **momentum-based backtracking**(§3.2). This technique implements progress-aware regret by explicitly pruning evolutionary trajectories: the agent’s context is reverted to a promising ancestor state, effectively removing the conditioning influence of the failed line of inquiry. By refreshing the context and unlearning the detrimental path, we encourage the agent to search in a different direction, facilitating the deep exploration necessary to escape local minima and generate diverse, innovative solutions. In §4.3, we demonstrate how backtracking eliminates trajectories that get stuck in local minima for the entirety of the evolution.

Insight 3: Dynamic cross-island synergy is key to fostering collaborative evolution

The previous techniques are designed to enhance the performance of individual islands, but they do not address the inherent inefficiency of a static multi-island setup. In traditional parallel evolution, knowledge transfer (crossover) is scheduled periodically and uniformly, regardless of the island’s evolutionary progress. We observe that this Static Coordination fails to navigate the fundamental tension between preserving an island’s internal search stability and leveraging external knowledge. When an island is making rapid local progress, forced crossover is disruptive; when an island is stagnating, delayed crossover or backtracking wastes computational steps. This inefficiency limits the potential for coordinated, non-uniform progress.

To address this and scale solutions in a principled manner, we introduce a self-adaptive sampling policy (§3.3) that unifies backtracking and crossover. This policy enables collaborative evolution by allowing islands to dynamically determine whether to backtrack (internal exploration) or perform crossover (external knowledge transfer) based on their current progress and momentum. This mechanism enables islands to learn from each other’s experiences efficiently and ensures that knowledge transfer occurs precisely when it is most beneficial to the collective search effort.

B. Method Details

B.1. Notation Table

Table 3 summarizes all notations used in PACEvolve.

Table 3. Notation table summarizing the mathematical concepts used in PACEvolve.

Symbol	Description	Definition / Context
<i>Evolutionary Search</i>		
Island	A semi-isolated sub-population of candidate solutions that evolves independently.	
Crossover	An operation of combining components from two parent solutions (often from two distinct islands) to produce a new offspring with mixed traits	
t	Current generation index	$t \in \mathbb{N}$
r	Target metric lower bound (e.g., 0 for error)	Constant
s_t	Best-achieved score by an island at generation t	$s_t \in \mathbb{R}$
s_0	Initial score of an island at start of search	
<i>Momentum-Based Backtracking (MBB)</i>		
G_t	Performance Gap: Distance to target	$G_t = s_t - r$
R_t	Relative Progress: Scale-invariant improvement	$R_t = \frac{G_{t-1} - G_t}{G_{t-1}}$
m_t	Relative Improvement Momentum: EWMA of R_t	$m_t = \beta m_{t-1} + (1 - \beta) R_t$
β	Momentum decay factor	$\beta \in [0, 1)$
ϵ_{rel}	Stagnation threshold for triggering intervention	Trigger MBB if $m_t < \epsilon_{rel}$
<i>Collaborative Evolution (CE)</i>		
A_t	Absolute Progress: Total gap closed (global metric)	$A_t = \frac{s_0 - s_t}{s_0 - r}$
i, j	Indices for specific evolutionary islands	
A	Set of available intervention actions	$A = \{\text{Backtrack}\} \cup \{\text{Crossover}_j\} \forall j$
w_a	Sampling weight for action $a \in A$	
$P(a)$	Probability of selecting action a	$\frac{w_a}{\sum_{a' \in A} w_{a'}}$

B.2. Context Management

B.2.1. DECOUPLING GENERATION AND SELECTION

Existing evolutionary frameworks typically prompt the LLM to generate a new candidate solution in a single, context-limited step. This prevents the accumulation of knowledge and hinders the evolutionary agent’s understanding of the problem space. To address this and cultivate Innovative Thinking, we re-architect the search process by decoupling candidate generation into a two-stage process: 1) idea generation and 2) idea selection, supported by a persistent idea pool. The persistent pool acts as an evolving knowledge base for the problem, ensuring that the agent maintains access to a rich, long-term history of conceptual directions, not just execution results. To manage this, we introduce Hierarchical Idea Grouping.

First, the LLM proposes a series of conceptual ideas. The LLM then classifies these proposals, ensuring they are conceptually distinct rather than differing only in minor implementation details. If a conceptual match exists in the persistent pool, the new proposal refines the existing idea; otherwise, it is added as a new entry. Then, we perform idea selection by granting the agent full access to this knowledge base to facilitate high-reward idea selection. In this stage, the LLM selects a conceptual idea to pursue, proposes a concrete experimental hypothesis, and implements it.

After evaluation, the results are summarized and appended as a new hypothesis record to the corresponding conceptual idea. This persistent, hierarchical structure decouples the creative, long-term thinking from the immediate execution, significantly increasing both idea diversity and the sophistication of the agent’s problem-solving knowledge.

B.2.2. CONTEXT PRUNING

Persistent idea memory, while enhancing diversity, introduces new challenges. The append-only nature of the idea memory causes the evolutionary agent to fail to escape local minima due to LLM’s inability to navigate the exploration-exploitation trade-off (Nie et al., 2024). LLMs exhibit a strong bias towards exploiting selected ideas, persisting with them even when

Algorithm 1 Hierarchical Context Management (HCM)

```

1: Input: Idea Pool  $\mathcal{P} \leftarrow \emptyset$ , Global Log  $\mathcal{L} \leftarrow \emptyset$ , Max Ideas  $K_{idea}$ , Max Hypotheses per Idea  $K_{hyp}$ 
2:  $Proposals \leftarrow \text{LLM}_{idea\_gen}(\mathcal{P})$ 
3: for each  $idea \in Proposals$  do
4:    $ID \leftarrow \text{LLM}_{idea\_classification}(idea, \mathcal{P})$ 
5:   if  $ID \in \mathcal{P}.IDs$  then
6:     Refine  $\mathcal{P}_{ID}$  description using  $idea$ 
7:   else
8:      $\mathcal{P} \leftarrow \mathcal{P} \cup \{idea\}$ 
9:   end if
10: end for
11:  $I_{sel}, H_{new} \leftarrow \text{LLM}_{select}(\mathcal{P})$ 
12: if  $H_{new} \notin \mathcal{L}$  then
13:    $Result \leftarrow \text{Execute}(H_{new})$ 
14:    $I_{sel}.history \leftarrow I_{sel}.history \cup \{Result\}$ 
15:    $\mathcal{L} \leftarrow \mathcal{L} \cup \{H_{new}\}$ 
16: end if
17: if  $|I_{sel}.history| > K_{hyp}$  then
18:    $Summary \leftarrow \text{LLM}_{summarize}(I_{sel}.history)$ 
19:    $I_{sel}.history \leftarrow \{Summary\}$ 
20: end if
21: if  $|\mathcal{P}| > K_{idea}$  then
22:    $I_{prune} \leftarrow \text{LLM}_{prune}(\mathcal{P})$ 
23:    $\mathcal{P} \leftarrow \mathcal{P} \setminus \{I_{prune}\}$ 
24: end if

```

only minor or negligible improvement is observed (Anthony et al., 2025).

Since LLMs lack an inherent mechanism, such as curiosity, to drive exploration toward conceptually novel ideas, we introduce two critical context management operators. First, to mitigate growth within an idea, we cap the number of experimental hypotheses per idea. Once this limit is reached, a summarization operator is triggered, distilling the accumulated experiment histories into concise key findings. Second, to manage the breadth of the idea pool and force radical exploration, we limit the total number of ideas. Once this cap is reached, the LLM would discard the least promising conceptual ideas, thereby encouraging exploration of novel concepts. We also maintain a full, separate log of all attempted ideas, asking the LLM to check against this log to avoid retrying the same flawed conceptual hypothesis after pruning. This combined approach ensures the agent’s context remains focused, high-quality, and continuously pushed toward beneficial exploration.

We present the pseudocode for HCM in Algorithm 1:

B.2.3. PROMPT TEMPLATES

The prompt templates for each stage are attached. Texts in **red** represent task-specific information that would be replaced with relevant materials; Texts in **blue** represent contexts that are managed during the evolution.

Idea Generation Prompt Template

We are conducting an evolutionary optimization process for **Task Name**.
BACKGROUND (Replace with task-specific introduction)
 The current state-of-the-art algorithm is as follows:
SoTA Algorithm (Sample SoTA solution tournament style)
 Idea repos contain ideas we have generated so far and the experiments we have run to test these hypotheses.
Idea Repo
 Your Task

When proposing a new design or hyperparameter configuration, you should start by conducting a research brainstorming exercise where you develop 3 different options to explore the design space. Go through each option and provide a comprehensive explanation of the proposed changes, including

- * The underlying rationale and expected impact.
- * The specific reason why you expect this experiment to be worth running.

Final Instructions

You should try to understand why a particular design performed well / poorly, so that you can make a more informed choice for the next set of designs. It is important to strike a good balance between exploring the design space and tuning to identify the best hyperparameter values. One way to tune is to develop hypotheses about how the design might affect performance, and then conduct experiments to test those hypotheses.

You are encouraged to consider configurations we may not have considered before, and you are likely to discover patterns during your search. You are **STRONGLY** encouraged to review your experimental history and refer to the designs we have already tested to ensure that the design you are proposing makes sense in the broader research context (i.e., is not too similar and is informed by past results). If an idea has a rich experimental history but its performance has plateaued, perhaps it's time to switch to a new idea.

Feel free to think outside of the box - you are allowed to propose candidates that are high risk (or even expected to be quality-negative) if it helps you understand the problem better, furthers your long-term research agenda, and (most importantly) increases your chance to produce a final candidate. The goal is to produce a final candidate after roughly X attempts, not to find a locally-optimal solution right away. That said, please be mindful of your compute usage: do not redo configurations you have already tried (the results will almost certainly be the same), and please do your best to deliver a good candidate by the end of your X attempts. We will take the best candidate discovered during your entire tuning process as your final candidate. It does not matter WHEN you find a good candidate, as long as you eventually do so.

Carefully review the idea and experimental history. DO NOT re-propose an idea that has been well tested already.

When there are fewer than or equal to 5 ideas, you should try to think out of the box and come up with at least one idea that does not fall under existing ideas.

You should follow the following format when generating ideas:

Idea 1

Idea: <Your idea here>

Reasoning: <Your reasoning here>

Idea 2

Idea: <Your idea here>

Reasoning: <Your reasoning here>

...

Idea N

Idea: <Your idea here>

Reasoning: <Your reasoning here>

Idea Classification Prompt Template

Below is the database of ideas we have explored:

[Idea Repo](#)

Here is the idea to be classified:

Idea

Your job is to classify the newly generated idea and merge it into the idea repo.

If you believe the following idea is similar or identical to one of the ideas in the database, you should respond in the following format by identifying the ID of the similar idea. For instance, if your idea leverages similar mathematical or computer-science concepts but requires different implementations or hyperparameter tuning, they should be grouped under the same idea.

If your idea is similar but not identical, provide an updated description that combines it with the original, and keep the description concise.

Otherwise, you can reuse the same idea description in the updated description section.

If your idea is a combination of two existing, but very different, ideas, it should be treated as a new idea.

If you find it hard to articulate the new hypothesis together with the original idea in TWO SENTENCES after merging, you should classify the hypothesis as a new idea. Your updated idea description should be NO MORE THAN TWO SENTENCES. But you should also AVOID overly broad terms and BE SPECIFIC about the ideas, while keeping the description concise.

Idea Exists: True

Idea ID: <Idea ID>

Updated description: <Updated description here>

If you believe your idea is new and orthogonal to existing ideas, respond in the following format:

Idea Exists: False

Idea description: <Provide your hypothesis here>

Idea Selection Prompt Template

Background and instructions

You should use the following format for the idea selection part:

Idea ID: <Idea ID>

Experiment description: <Provide a concrete but concise description of the experiment you want to try, DO NOT HALLUCINATE IDEA ID, YOU MUST SELECT ONE FROM the idea repo above; however, you can provide more details about the experiment you want to try based on that idea>

Once your brainstorming and idea-generation process is complete, you are ready to write code. Please follow the guidelines when completing the coding part:

Coding Instructions

History Summarization Prompt Template

Below is the idea we want to summarize: **Idea and Corresponding Experiment Histories**

Your task is to summarize and condense the experimental history into a single bullet point. You should keep the result of the best trial so far, followed by summarizing which ideas work and which do not.

Format like this:

- Results: <Summary of what ideas work and what do not work>

Idea Capping Prompt Template

Below is the database of ideas we have explored: [Idea Repo](#)
 We want to cap the number of ideas under consideration at <Idea Cap>;
 therefore, we need to drop some ideas.
 Your job is to evaluate which ideas we should drop to reduce the number of
 ideas under consideration to <Idea Cap>.
 Your criteria should be as follows:
 1. If thorough experiment results from the experiment history of this and
 other ideas clearly invalidate this idea?
 2. Is there any potential for a breakthrough in performance if we try more
 hypotheses based on this idea?
 3. Has this idea been thoroughly investigated (indicated by experiment
 count), and does the performance clearly underperform the current best
 results in the experiment history?
 4. With the number of experiments performed on this idea, is it possible to
 reduce the gap between the current best-of-class performance of this idea
 and the target metric?
 5. You should prioritize dropping ideas that are either old, lack the
 potential to improve, or have been explored extensively but still do not
 see breakthrough improvements.
 You should output the idea(s) to drop in a Python list format (e.g., if you
 want to drop idea i , you should output $[i]$)
 Format like this:
 Dropping Ideas: <list of idea(s) to drop>

B.3. Action Weighting

In this section, we describe the details of the action-weighting mechanism used in self-adaptive crossover sampling.

The weights are designed to balance the competing principles introduced in §3.3. Let A_i be the absolute progress of our triggered island and $A_{\text{best}} = \max_{j \neq i}(A_j)$ be the progress of the best available partner island, j_{best} .

1. Crossover Weight (w_{C_j}): We define the base utility for crossing over with any island j as the direct performance gain it offers:

$$w_{C_j}^{\text{base}} = \max(0, A_j - A_i)$$

which favors islands with higher absolute progress.

2. Backtrack Weight (w_{BT}): The utility of backtracking is based on two conditions. First, the *dominance* component (w_{BT}^{dom}), which applies if the current island is outperforming all others:

$$w_{BT}^{\text{dom}} = \max(0, A_i - A_{\text{best}})$$

Second, the *low-progress stagnation* component (w_{BT}^{stag}), which applies if two islands are making similarly low progress. We define similarity $S = \max(0, 1 - |A_i - A_{\text{best}}|)$.

$$w_{BT}^{\text{stag}} = S \cdot (1 - A_i) \cdot (1 - A_{\text{best}})$$

This term is high only when S is high (high similarity) and both $(1 - A_i)$ and $(1 - A_{\text{best}})$ are high (low progress). The total backtrack weight is their sum:

$$w_{BT} = w_{BT}^{\text{dom}} + w_{BT}^{\text{stag}}$$

3. Synergy Bonus (w_C^{syn}): Conversely, if two islands are making similarly high progress, they may synergize. We add a bonus to the best partner j_{best} :

$$w_C^{\text{syn}} = S \cdot A_i \cdot A_{\text{best}}$$

This term is high only when S is high (high similarity) and both A_i and A_{best} are high (high progress).

Algorithm 2 Momentum-based Backtracking (MBB) and Collaborative Evolution (CE)

```

1: Input: Island  $i$ , Initial score  $s_0$ , Target lower bound  $r$ , Decay  $\beta$ , Threshold  $\epsilon_{\text{rel}}$ 
2:  $s_{\text{prev}} \leftarrow s_0, m \leftarrow 1.0, G_{\text{prev}} \leftarrow s_0 - r$ 
3:  $s_{\text{curr}} \leftarrow \text{Evaluate Candidate Solution}$ 
4: Update best score  $s_{\text{best}} \leftarrow \min(s_{\text{curr}}, s_{\text{prev}})$ 
   Update Momentum Metrics (§3.2)
5:  $G_{\text{curr}} \leftarrow s_{\text{best}} - r$ 
6: if  $s_{\text{best}} < s_{\text{prev}}$  then
7:    $R_t \leftarrow (s_{\text{prev}} - s_{\text{best}}) / (s_{\text{prev}} - r)$ 
8: else
9:    $R_t \leftarrow 0$ 
10: end if
11:  $m \leftarrow \beta \cdot m + (1 - \beta) \cdot R_t$ 
12:  $s_{\text{prev}} \leftarrow s_{\text{best}}$ 
   Self-Adaptive Sampling (§3.3)
13: if  $m < \epsilon_{\text{rel}}$  then
14:   Calculate Abs. Progress  $A_i \leftarrow (s_0 - s_{\text{best}}) / (s_0 - r)$ 
15:   Fetch neighbor progress  $\{A_j\}_{j \neq i}$ 
16:    $A_{\text{best}} \leftarrow \max_{j \neq i}(A_j)$ 
17:    $S \leftarrow \max(0, 1 - |A_i - A_{\text{best}}|)$ 
18:    $w_{\text{BT}} \leftarrow \max(0, A_i - A_{\text{best}}) + S(1 - A_i)(1 - A_{\text{best}})$ 
19:   for each  $j \neq i$  do
20:      $w_{C_j} \leftarrow \max(0, A_j - A_i)$ 
21:     if  $j = \arg \max(A)$  then
22:        $w_{C_j} \leftarrow w_{C_j} + S \cdot A_i \cdot A_{\text{best}}$ 
23:     end if
24:   end for
25:    $\text{action} \sim P(a) \propto \{w_{\text{BT}}, w_{C_1}, \dots\}$ 
26:   if  $\text{action} = \text{Backtrack}$  then
27:     Revert to previous state  $t' \sim \text{PowerLaw}$ 
28:   else
29:     Crossover with partner  $j$  selected by action
30:   end if
31: end if
    
```

The final weights for all actions are then assembled. For all $j \neq j_{\text{best}}$, the weight is $w_{C_j} = w_{C_j}^{\text{base}}$. For the best partner, $w_{C_{j_{\text{best}}}} = w_{C_{j_{\text{best}}}}^{\text{base}} + w_C^{\text{syn}}$. The backtrack weight is w_{BT} .

The probability of choosing any action $a \in A$ is then:

$$P(a) = \frac{w_a}{w_{\text{BT}} + \sum_{j \neq i} w_{C_j}}$$

This model adaptively balances exploration (favoring high-gain partners) and exploitation (backtracking when a dominant partner is present), while avoiding stagnation (backtracking when all islands show low progress).

We present the pseudocode for momentum-based backtracking and crossover in Algorithm 2:

C. Experiment Details

C.1. Benchmark Task Selection

While evolutionary search can be applied to any task with a defined fitness score, we define four desirable properties for a benchmark task intended to evaluate evolutionary search:

Table 4. This table lists a subset of tasks we support on the PACEvolve platform and whether they satisfy the desirable properties for studying evolutionary search.

	Solution Space	World Knowledge	Reward Smoothness	Eval Cost
Consistent Hashing	✗	✗	✗	✓
Capset	✓	?	✗	✓
DCN Optimizer	?	✗	✓	✗
Feature Cross	✓	✗	✓(Validation Loss)	? (15 mins)
Symbolic Regression	✓	✓(Synthetic)	✓(NMSE)	✓(<1s)
Kernel Bench	? (Kernel Dependent)	?	✓(Latency)	✓(<2 mins)
Modded NanoGPT	✓	?	✓(Latency)	✗(8 GPUs)

- **Complex Solution Space:** Multiple distinct techniques are required to optimize task performance. This property elicits progressive, observable improvements in the evolutionary trajectory, making the comparison of different search strategies more robust.
- **LLM World Knowledge:** To better isolate the effect of different evolutionary recipes, we include tasks where LLM world knowledge plays a less significant role. The task should ideally not admit different optimal solutions for different task instances (e.g., different workloads or datasets in machine learning). This property compels the LLM to engage in iterative exploration rather than merely retrieving a near-optimal solution from its internal knowledge of similar tasks and not produce a near-optimal solution in a few attempts (Wang et al., 2021; Chen et al., 2021).
- **Smooth Reward Landscape:** The reward landscape should not be excessively sparse or non-smooth. In such landscapes, LLMs cannot effectively leverage evolutionary history or build on existing knowledge, thereby reducing the process to a lengthy and exhaustive search (Romera-Paredes et al., 2024).
- **Low-Cost Evaluation:** Task evaluation must be computationally cheap and fast. The inherent stochasticity of LLM decoding and evolutionary search necessitates multiple experimental repetitions to draw robust conclusions when comparing different search methodologies (Zhu et al., 2021; 2022).

These four properties collectively enable the large-scale empirical studies required to understand and improve evolutionary search. Table 4 shows a subset of tasks supported in PACEvolve across various fields, such as algorithm design, combinatorial optimization, machine learning kernel engineering, and deep learning research. They include the following tasks:

- **Consistent Hashing:** We run an evolutionary search to ask the LLM to improve upon the existing SoTA consistent hashing algorithm, especially for heavy workloads (Chen et al., 2021).
- **Capset:** We run an evolutionary search to ask the LLM to improve upon the SoTA construction of Capset (Romera-Paredes et al., 2024).
- **DCN Optimizer Design:** We run an evolutionary search to ask the LLM to improve upon optimizer design for DCN in the FuxiCTR framework (Zhu et al., 2021; Wang et al., 2021).
- **Feature Cross:** We run an evolutionary search to ask the LLM to design a feature set and feature crosses, enabling a Logistic Regression model to achieve performance comparable to state-of-the-art neural networks, such as DCNv2 (Wang et al., 2021).
- **Symbolic Regression:** We integrate LLM-SR into our framework and ask the LLM to fit a synthetically generated oscillator acceleration equation.
- **KernelBench:** We ask the LLM to design and improve the latency of various deep learning kernels.
- **Modded NanoGPT:** We ask the LLM to design and improve the latency of training a GPT2-style model to target validation loss in a distributed training setup with 8 H100 GPUs.

We label KernelBench and Modded NanoGPT as having questionable LLM World Knowledge, since the general strategies for optimizing model architecture and kernels are likely part of the LLM pretraining corpora, despite the benchmarks being created after the knowledge cutoff of the tested LLMs. We label the Solution Space of KernelBench as questionable because we find that the speedups for some discovered kernels can be largely attributed to a single innovation, whereas others benefit from multiple composable innovations.

C.2. LLM-SR

We use the first Nonlinear damped oscillator question in the LLMSR Bench. We note that other problems in the benchmark use a similar format and primarily vary through perturbing term combinations. We conducted ten trials per run on the instance to establish a systematic understanding of inter-run variability. The ground truth instance is: $-1.0267 * x * *3 - 1.0267 * x * \exp(-\text{Abs}(x)) + 0.9480 * \sin(t) - 0.7123 * \sin(v)$

C.3. More KernelBench Results

C.3.1. KERNEL LIST

Table 5 shows the list of kernels we evaluate, and their corresponding difficult levels and problem indices in KernelBench (Ouyang et al., 2025a). We use the prompt from GEPA (Agrawal et al., 2025) as the background and instruction section during idea generation and selection (§B.2.3)

Table 5. Sampled kernel list.

Kernel	Full Kernel Name	Difficulty Level	Problem Index
BatchNorm	BatchNorm	1	33
ConvDiv	Conv3d Divide Max GlobalAvgPool BiasAdd Sum	2	8
ConvMax	Conv3d Max LogSumExp ReLU	2	43
ConvT2D	ConvTranspose2d BiasAdd Clamp Scaling Clamp Divide	2	2
GELU	GELU	1	26
MatMul	Matmul with Large K dimension	1	6
MaxPool	Max Pooling 2D	1	42
MLP	MLP	3	1
RMSNorm	RMSNorm	1	36
Softmax	Softmax	1	23
VGG16	VGG16	3	11
MeanRedu	Mean reduction over a dimension	1	48
RNN	VanillaRNN	3	33
BMMNorm	BMM InstanceNorm Sum ResidualAdd Multiply	2	28
AlexNet	AlexNet	3	5
LayerNorm	LayerNorm	1	40

C.3.2. HEAD TO HEAD COMPARISON

Figure 5 shows that both the single island (PACEvolve-Single) and the multi-island (PACEvolve-Multi) version of PACEvolve outperform the best existing kernels on KernelBench in all tested cases. In addition, PACEvolve-Single outperforms ShinkaEvolve in all tested kernels, which PACEvolve-Multi further outperforms PACEvolve-Single in 81.25% (13/16) of the tested kernels. When comparing against other evolutionary frameworks, PACEvolve-Multi found equivalent or better kernels compared to ShinkaEvolve and CodeEvolve in 14/16 cases and OpenEvolve in 15/16 cases, clearly demonstrating better framework design despite possible variances in individual runs.

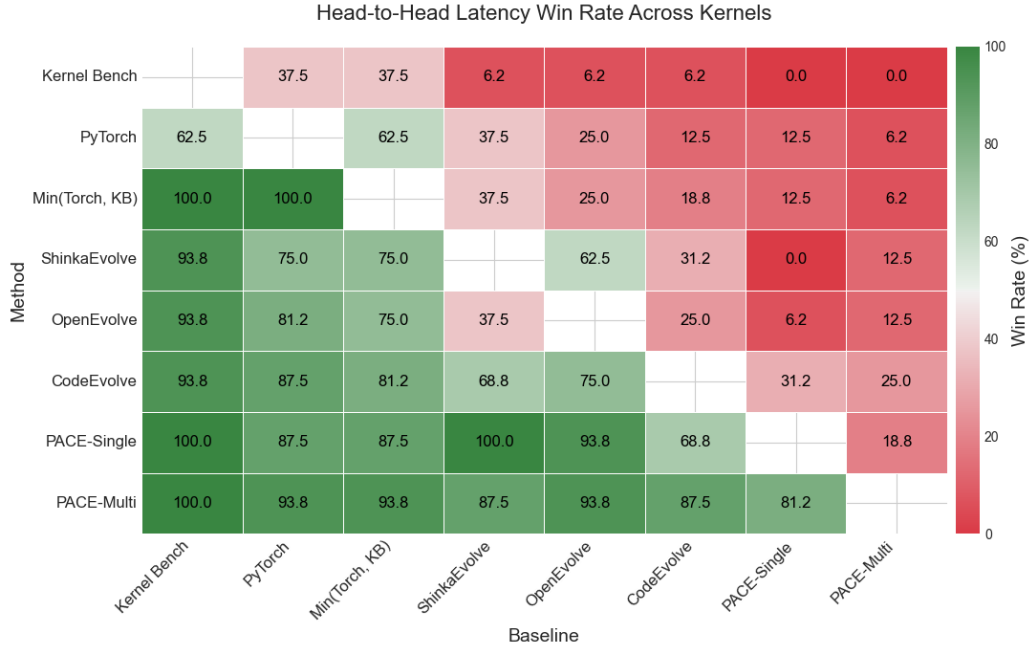


Figure 5. Head-to-head win rate comparison. Win rate percentages are shown for all method pairs, indicating the proportion of kernels where the row method outperformed the column method. Equal counts as a win for both methods; therefore, the heatmap is not strictly symmetric.

D. Discovered Kernels

Here we list the kernels discovered by PACEvolve. We remove PACEvolve-generated comments and inline checks to save space.

D.1. BatchNorm

```
import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

batch_norm_source_ldg_params = """
#include <torch/extension.h>
#include <cuda_fp16.h>
#include <cuda_runtime.h>

__global__ void __launch_bounds__(1024, 1) batch_norm_fp16_compute_kernel(
    const float* __restrict__ input,
    const float* __restrict__ gamma,
    const float* __restrict__ beta,
    const float* __restrict__ running_mean,
    const float* __restrict__ running_var,
    float* __restrict__ output,
    int total_size,
    int channels,
    int height,
    int width,
    float epsilon
) {
```

```

int idx = blockIdx.x * blockDim.x + threadIdx.x;
int size_vec = total_size / 4;

if (idx < size_vec) {
    int original_idx = idx * 4;
    int c = (original_idx / (width * height)) % channels;

    float g = __ldg(&gamma[c]);
    float b = __ldg(&beta[c]);
    float mean = __ldg(&running_mean[c]);
    float var = __ldg(&running_var[c]);

    float inv_std = rsqrtf(var + epsilon);
    float scale = g * inv_std;
    float shift = b - mean * scale;
    const float4 x4 = __ldcg(&(((const float4*)input)[idx]));
    const __half scale_h = __float2half_rn(scale);
    const __half shift_h = __float2half_rn(shift);

    const __half2 scale_h2 = __halves2half2(scale_h, scale_h);
    const __half2 shift_h2 = __halves2half2(shift_h, shift_h);
    const __half2 x_h2_a = __float22half2_rn(make_float2(x4.x, x4.y));
    const __half2 x_h2_b = __float22half2_rn(make_float2(x4.z, x4.w));
    const __half2 y_h2_a = __hfma2(x_h2_a, scale_h2, shift_h2);
    const __half2 y_h2_b = __hfma2(x_h2_b, scale_h2, shift_h2);
    const float2 y_f2_a = __half22float2(y_h2_a);
    const float2 y_f2_b = __half22float2(y_h2_b);
    ((float4*)output)[idx] = make_float4(y_f2_a.x, y_f2_a.y, y_f2_b.x,
        y_f2_b.y);
}
}

std::vector<torch::Tensor> batch_norm_fp16_compute_cuda(
    torch::Tensor input,
    torch::Tensor gamma,
    torch::Tensor beta,
    torch::Tensor running_mean,
    torch::Tensor running_var,
    float epsilon
) {
    const int batch_size = input.size(0);
    const int channels = input.size(1);
    const int height = input.size(2);
    const int width = input.size(3);

    auto output = torch::empty_like(input);
    const int total_size = batch_size * channels * height * width;

    if (total_size == 0) {
        return {output};
    }

    const int threads_per_block = 1024;
    const int num_blocks = (total_size / 4 + threads_per_block - 1) /
        threads_per_block;

    batch_norm_fp16_compute_kernel<<<num_blocks, threads_per_block>>>(
        input.data_ptr<float>(),
        gamma.data_ptr<float>(),
        beta.data_ptr<float>(),
        running_mean.data_ptr<float>(),
        running_var.data_ptr<float>(),
        output.data_ptr<float>(),

```

```

        total_size,
        channels,
        height,
        width,
        epsilon
    );

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) {
        AT_ERROR("CUDA kernel launch failed: ", cudaGetErrorString(err));
    }

    return {output};
}
"""

batch_norm_cpp_source_ldg_params = """
#include <torch/extension.h>
#include <vector>

std::vector<torch::Tensor> batch_norm_fp16_compute_cuda(
    torch::Tensor input,
    torch::Tensor gamma,
    torch::Tensor beta,
    torch::Tensor running_mean,
    torch::Tensor running_var,
    float epsilon
);
"""

batch_norm_ldg_params_module = load_inline(
    name='batch_norm_fp16_compute_ldg_params',
    cpp_sources=batch_norm_cpp_source_ldg_params,
    cuda_sources=batch_norm_source_ldg_params,
    functions=['batch_norm_fp16_compute_cuda'],
    verbose=False,
    extra_cflags=['-O3'],
    extra_cuda_cflags=['-O3', '--use_fast_math']
)

class ModelNew(nn.Module):
    def __init__(self, num_features: int):
        super(ModelNew, self).__init__()
        self.num_features = num_features
        self.weight = nn.Parameter(torch.ones(num_features).contiguous())
        self.bias = nn.Parameter(torch.zeros(num_features).contiguous())
        self.register_buffer('running_mean', torch.zeros(num_features).contiguous())
        self.register_buffer('running_var', torch.ones(num_features).contiguous())
        self.eps = 1e-5
        self.batch_norm_impl = batch_norm_ldg_params_module

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        return self.batch_norm_impl.batch_norm_fp16_compute_cuda(
            x,
            self.weight,
            self.bias,
            self.running_mean,
            self.running_var,
            self.eps
        ) [0]

```

D.2. Conv3d Divide Max GlobalAvgPool BiasAdd Sum

```

import torch
import torch.nn as nn
import torch.nn.functional as F
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

cuda_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <limits>
#include <cuda_fp16.h>

__device__ __forceinline__ __half compute_max_from_coords(
    int d_out, int h_out, int w_out,
    const __half* __restrict__ input_channel,
    int H_in, int W_in,
    int pool_d, int pool_h, int pool_w,
    __half conv_bias_val) {

    int d_start = d_out * pool_d;
    int h_start = h_out * pool_h;
    int w_start = w_out * pool_w;

    const __half* window_base_ptr = input_channel + d_start * H_in * W_in + h_start
        * W_in + w_start;

    const __half2* ptr2 = reinterpret_cast<const __half2*>(window_base_ptr);
    const __half2 conv_bias2 = __half2half2(conv_bias_val);

    __half2 val00 = __hadd2(ptr2[0], conv_bias2);
    __half2 val01 = __hadd2(ptr2[W_in / 2], conv_bias2);
    __half2 val10 = __hadd2(ptr2[H_in * W_in / 2], conv_bias2);
    __half2 val11 = __hadd2(ptr2[H_in * W_in / 2 + W_in / 2], conv_bias2);

    __half max_d0 = __hmax(__hmax(val00.x, val00.y), __hmax(val01.x, val01.y));
    __half max_d1 = __hmax(__hmax(val10.x, val10.y), __hmax(val11.x, val11.y));

    return __hmax(max_d0, max_d1);
}

__global__ void monolithic_fused_kernel_2d_block(
    const __half* __restrict__ input,
    const __half* __restrict__ conv_bias,
    const __half* __restrict__ bias,
    float* __restrict__ output,
    int B, int C,
    int D_in, int H_in, int W_in,
    int D_out, int H_out, int W_out,
    int pool_d, int pool_h, int pool_w,
    float inv_scale_float) {

    const int CHANNELS_PER_BLOCK = 2;
    extern __shared__ __half sdata[];

    int base_bc_idx = blockIdx.x * CHANNELS_PER_BLOCK;
    if (base_bc_idx >= B * C) {
        return;
    }

```



```

int b = base_bc_idx / C;

const __half* input_channel0 = input + (base_bc_idx + 0) * D_in * H_in * W_in;
const __half* input_channel1 = input + (base_bc_idx + 1) * D_in * H_in * W_in;

bool c1_valid = (base_bc_idx + 1 < B * C);

const __half zero_h = __float2half(0.0f);

__half conv_b0 = conv_bias[base_bc_idx % C];
__half conv_b1 = c1_valid ? conv_bias[(base_bc_idx + 1) % C] : zero_h;

__half2 my_sum01 = make_half2(zero_h, zero_h);

int h_out = threadIdx.y;
int w_out = threadIdx.x;

if (h_out < H_out && w_out < W_out) {
    for (int d_out = 0; d_out < D_out; ++d_out) {
        __half v0 = compute_max_from_coords(d_out, h_out, w_out,
            input_channel0, H_in, W_in, pool_d, pool_h, pool_w, conv_b0);
        __half v1 = c1_valid ? compute_max_from_coords(d_out, h_out, w_out,
            input_channel1, H_in, W_in, pool_d, pool_h, pool_w, conv_b1) :
            zero_h;
        my_sum01 = __hadd2(my_sum01, make_half2(v0, v1));
    }
}

#pragma unroll
for (int offset = 16; offset > 0; offset >>= 1) {
    my_sum01 = __hadd2(my_sum01, __shfl_down_sync(0xFFFFFFFF, my_sum01,
        offset));
}

int tid_ld = threadIdx.y * blockDim.x + threadIdx.x;
int lane_id = tid_ld % 32;
int warp_id = tid_ld / 32;
int num_warps = (blockDim.x * blockDim.y) / 32;

__half2* sdata2 = reinterpret_cast<__half2*>(sdata);
if (lane_id == 0) {
    sdata2[warp_id] = my_sum01;
}
__syncthreads();

if (warp_id == 0) {
    __half2 warp_sum01 = (lane_id < num_warps) ? sdata2[lane_id] :
        make_half2(zero_h, zero_h);

    #pragma unroll
    for (int offset = 16; offset > 0; offset >>= 1) {
        warp_sum01 = __hadd2(warp_sum01, __shfl_down_sync(0xFFFFFFFF,
            warp_sum01, offset));
    }

    if (lane_id == 0) {
        float final_val_f = 0.0f;
        const __half inv_scale = __float2half(inv_scale_float);

        __half p_val0 = __hadd(__hmul(warp_sum01.x, inv_scale),
            bias[base_bc_idx % C]);
        final_val_f += __half2float(p_val0);
    }
}

```

```

        if(c1_valid) {
            __half p_val1 = __hadd(__hmul(warp_sum01.y, inv_scale),
                                   bias[(base_bc_idx + 1) % C]);
            final_val_f += __half2float(p_val1);
        }

        atomicAdd(&output[b], final_val_f);
    }
}

void launch_kernel(const torch::Tensor& input, const torch::Tensor& conv_bias,
                  const torch::Tensor& bias, torch::Tensor& output, int pool_d, int pool_h, int
pool_w, float inv_scale) {
    const int B = input.size(0); const int C = input.size(1);
    const int D_in = input.size(2); const int H_in = input.size(3); const int W_in
    = input.size(4);
    const int D_out = D_in / pool_d; const int H_out = H_in / pool_h; const int
    W_out = W_in / pool_w;

    const int CHANNELS_PER_BLOCK = 2;
    const dim3 block_dim(16, 16); // Use a 2D block
    const int block_size_ld = block_dim.x * block_dim.y;
    const int num_warps = block_size_ld / 32;

    const int num_blocks = (B * C + CHANNELS_PER_BLOCK - 1) / CHANNELS_PER_BLOCK;
    const size_t shared_mem_size = num_warps * sizeof(__half2);

    monolithic_fused_kernel_2d_block<<<num_blocks, block_dim, shared_mem_size>>> (
        (const __half*)input.data_ptr(), (const __half*)conv_bias.data_ptr(),
        (const __half*)bias.data_ptr(), output.data_ptr<float>(),
        B, C, D_in, H_in, W_in, D_out, H_out, W_out,
        pool_d, pool_h, pool_w, inv_scale
    );
}

torch::Tensor fused_op(torch::Tensor input, torch::Tensor conv_bias, torch::Tensor
bias, int pool_d, int pool_h, int pool_w, float inv_scale) {
    const int B = input.size(0);
    auto output = torch::zeros({B}, input.options().dtype(torch::kFloat32));

    launch_kernel(input, conv_bias, bias, output, pool_d, pool_h, pool_w,
        inv_scale);

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) { throw std::runtime_error(cudaGetErrorString(err)); }
    return output.view({B, 1, 1, 1});
}

void fused_op_out(torch::Tensor input, torch::Tensor conv_bias, torch::Tensor bias,
torch::Tensor output, int pool_d, int pool_h, int pool_w, float inv_scale) {
    auto output_ld = output.view({-1});
    output_ld.zero_();
    launch_kernel(input, conv_bias, bias, output_ld, pool_d, pool_h, pool_w,
        inv_scale);

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) { throw std::runtime_error(cudaGetErrorString(err)); }
}
"""

cpp_sources = """

```

```

torch::Tensor fused_op(torch::Tensor input, torch::Tensor conv_bias, torch::Tensor
    bias, int pool_d, int pool_h, int pool_w, float inv_scale);
void fused_op_out(torch::Tensor input, torch::Tensor conv_bias, torch::Tensor bias,
    torch::Tensor output, int pool_d, int pool_h, int pool_w, float inv_scale);
"""

fused_op_module = load_inline(
    name='fused_op_module_2d_block',
    cpp_sources=cpp_sources,
    cuda_sources=cuda_source,
    functions=['fused_op', 'fused_op_out'],
    verbose=True
)

class ModelNew(nn.Module):
    def __init__(self, in_channels: int, out_channels: int, kernel_size: int,
        divisor: float, pool_size: tuple, bias_shape: tuple, sum_dim: int):
        super(ModelNew, self).__init__()
        self.conv = nn.Conv3d(in_channels, out_channels, kernel_size,
            bias=True).half()
        self.bias = nn.Parameter(torch.randn(bias_shape).squeeze().half())
        self.divisor = divisor
        self.pool_d, self.pool_h, self.pool_w = pool_size
        self.graph = None
        self.static_input = None
        self.static_output = None
        self.static_inv_scale = 0.0

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x = x.to(self.conv.weight.device).half()
        conv_out_no_bias = F.conv3d(x, self.conv.weight, bias=None,
            stride=self.conv.stride,
            padding=self.conv.padding,
            dilation=self.conv.dilation,
            groups=self.conv.groups)

        if self.training:
            if self.graph is not None:
                self.graph = None
                self.static_input = None
                self.static_output = None

            num_pool_outputs = (conv_out_no_bias.size(2) // self.pool_d) * \
                (conv_out_no_bias.size(3) // self.pool_h) * \
                (conv_out_no_bias.size(4) // self.pool_w)
            inv_scale = 1.0 / (num_pool_outputs * self.divisor) if num_pool_outputs
                > 0 else (1.0 / self.divisor)

            return fused_op_module.fused_op(
                conv_out_no_bias, self.conv.bias, self.bias,
                self.pool_d, self.pool_h, self.pool_w, inv_scale
            )

        if self.graph is None or x.shape != self.static_input.shape:
            num_pool_outputs = (conv_out_no_bias.size(2) // self.pool_d) * \
                (conv_out_no_bias.size(3) // self.pool_h) * \
                (conv_out_no_bias.size(4) // self.pool_w)
            self.static_inv_scale = 1.0 / (num_pool_outputs * self.divisor) if
                num_pool_outputs > 0 else (1.0 / self.divisor)

            output = fused_op_module.fused_op(
                conv_out_no_bias, self.conv.bias, self.bias,
                self.pool_d, self.pool_h, self.pool_w, self.static_inv_scale

```

```
)

self.static_input = x.clone()
self.static_output = torch.empty_like(output)

self.graph = torch.cuda.CUDAGraph()
with torch.cuda.graph(self.graph):
    conv_out_graph = F.conv3d(self.static_input, self.conv.weight,
                              bias=None,
                              stride=self.conv.stride,
                              padding=self.conv.padding,
                              dilation=self.conv.dilation,
                              groups=self.conv.groups)
    fused_op_module.fused_op_out(
        conv_out_graph, self.conv.bias, self.bias, self.static_output,
        self.pool_d, self.pool_h, self.pool_w, self.static_inv_scale
    )
return output

self.static_input.copy_(x)
self.graph.replay()
return self.static_output.clone()
```

D.3. Conv3d Max LogSumExp ReLU

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

cuda_source = """
#include <torch/extension.h>
#include <ATen/cuda/CUDAContext.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>
#include <cmath>

constexpr int CONV_KERNEL_SIZE = 3;
constexpr int CONV_STRIDE = 1;
constexpr int CONV_PADDING = 1;
constexpr int POOL_SIZE = 2;
constexpr int OUT_CHANNELS = 16;
constexpr int CONV_KERNEL_VOL = CONV_KERNEL_SIZE * CONV_KERNEL_SIZE *
    CONV_KERNEL_SIZE;
constexpr int TILE_H_OUT = 4;
constexpr int TILE_W_OUT = 16;
constexpr int TILE_D_CONV = POOL_SIZE; // = 2
constexpr int TILE_H_CONV = TILE_H_OUT * POOL_SIZE; // 4*2 = 8
constexpr int TILE_W_CONV = TILE_W_OUT * POOL_SIZE; // 16*2 = 32
constexpr int TILE_D_IN = TILE_D_CONV + CONV_KERNEL_SIZE - 1; // 2+3-1 = 4
constexpr int TILE_H_IN = TILE_H_CONV + CONV_KERNEL_SIZE - 1; // 8+3-1 = 10
constexpr int TILE_W_IN = TILE_W_CONV + CONV_KERNEL_SIZE - 1; // 32+3-1 = 34

__global__ void fused_conv_maxpool_bias_logsumexp_relu_kernel(
    const half* __restrict__ x, const half* __restrict__ weight, const half*
    __restrict__ bias,
    half* __restrict__ y,
    int N, int C_in, int D_in, int H_in, int W_in,
    int D_out, int H_out, int W_out)
{
    const int w_out_base = blockIdx.x * TILE_W_OUT;
    const int h_out_base = blockIdx.y * TILE_H_OUT;
    const int nd_idx = blockIdx.z;

    if (nd_idx >= N * D_out) return;

    const int n = nd_idx / D_out;
    const int d_out = nd_idx % D_out;

    extern __shared__ half smem[];
    half* sm_weights = smem;
    half* sm_input_x = sm_weights + OUT_CHANNELS * CONV_KERNEL_VOL;
    half* sm_conv_out = sm_input_x + TILE_D_IN * TILE_H_IN * TILE_W_IN;

    half conv_accum[OUT_CHANNELS][4] = {{0.0f}};

    const int d_conv_base = d_out * POOL_SIZE;
    const int h_conv_base = h_out_base * POOL_SIZE;
    const int w_conv_base = w_out_base * POOL_SIZE;

    const int d_in_base = d_conv_base - CONV_PADDING;
    const int h_in_base = h_conv_base - CONV_PADDING;
    const int w_in_base = w_conv_base - CONV_PADDING;

```



```

constexpr int num_weights_per_ic = OUT_CHANNELS * CONV_KERNEL_VOL;

for (int ic = 0; ic < C_in; ++ic) {
    const int num_weights_per_ic_vec = num_weights_per_ic / 2;
    const half2* weight_ic_base_vec = reinterpret_cast<const half2*>(weight +
        (long long)ic * num_weights_per_ic);
    half2* sm_weights_vec = reinterpret_cast<half2*>(sm_weights);
    for (int i = threadIdx.x; i < num_weights_per_ic_vec; i += blockDim.x) {
        sm_weights_vec[i] = weight_ic_base_vec[i];
    }

    half2* sm_input_x_vec = reinterpret_cast<half2*>(sm_input_x);
    const int smem_tile_size_vec = (TILE_D_IN * TILE_H_IN * TILE_W_IN) / 2;

    int initial_flat_idx = threadIdx.x * 2;
    int w_local_curr = initial_flat_idx % TILE_W_IN;
    int h_rem = initial_flat_idx / TILE_W_IN;
    int h_local_curr = h_rem % TILE_H_IN;
    int d_local_curr = h_rem / TILE_H_IN;

    const int stride_flat = blockDim.x * 2;
    const int w_stride = stride_flat % TILE_W_IN;
    const int h_rem_stride = stride_flat / TILE_W_IN;
    const int h_stride = h_rem_stride % TILE_H_IN;
    const int d_stride = h_rem_stride / TILE_H_IN;

    for (int i = threadIdx.x; i < smem_tile_size_vec; i += blockDim.x) {
        int d_in_idx_1 = d_in_base + d_local_curr;
        int h_in_idx_1 = h_in_base + h_local_curr;
        int w_in_idx_1 = w_in_base + w_local_curr;

        half val1 = (half)0.0f;
        if (d_in_idx_1 >= 0 && d_in_idx_1 < D_in &&
            h_in_idx_1 >= 0 && h_in_idx_1 < H_in &&
            w_in_idx_1 >= 0 && w_in_idx_1 < W_in) {
            val1 = x[(((long long)n * C_in + ic) * D_in + d_in_idx_1) * H_in +
                h_in_idx_1 * W_in + w_in_idx_1];
        }

        int d_local_2 = d_local_curr;
        int h_local_2 = h_local_curr;
        int w_local_2 = w_local_curr + 1;
        if (w_local_2 == TILE_W_IN) {
            w_local_2 = 0; h_local_2++;
            if (h_local_2 == TILE_H_IN) { h_local_2 = 0; d_local_2++; }
        }

        int d_in_idx_2 = d_in_base + d_local_2;
        int h_in_idx_2 = h_in_base + h_local_2;
        int w_in_idx_2 = w_in_base + w_local_2;

        half val2 = (half)0.0f;
        if (d_in_idx_2 >= 0 && d_in_idx_2 < D_in &&
            h_in_idx_2 >= 0 && h_in_idx_2 < H_in &&
            w_in_idx_2 >= 0 && w_in_idx_2 < W_in) {
            val2 = x[(((long long)n * C_in + ic) * D_in + d_in_idx_2) * H_in +
                h_in_idx_2 * W_in + w_in_idx_2];
        }

        sm_input_x_vec[i] = make_half2(val1, val2);

        w_local_curr += w_stride;
    }
}

```

```

        if (w_local_curr >= TILE_W_IN) { w_local_curr -= TILE_W_IN;
            h_local_curr++; }

        h_local_curr += h_stride;
        if (h_local_curr >= TILE_H_IN) { h_local_curr -= TILE_H_IN;
            d_local_curr++; }

        d_local_curr += d_stride;
    }
    __syncthreads();

#pragma unroll
for (int oc = 0; oc < OUT_CHANNELS; ++oc) {
    #pragma unroll
    for (int p = 0; p < 4; ++p) { // Each thread computes 4 points
        int point_idx = threadIdx.x * 4 + p;
        int w_conv_local = point_idx % TILE_W_CONV;
        int h_conv_local = (point_idx / TILE_W_CONV) % TILE_H_CONV;
        int d_conv_local = point_idx / (TILE_W_CONV * TILE_H_CONV);

        half accum_val = 0.0f;
        #pragma unroll
        for (int kd = 0; kd < CONV_KERNEL_SIZE; ++kd) {
            #pragma unroll
            for (int kh = 0; kh < CONV_KERNEL_SIZE; ++kh) {
                #pragma unroll
                for (int kw = 0; kw < CONV_KERNEL_SIZE; ++kw) {
                    int d_in_local = d_conv_local + kd;
                    int h_in_local = h_conv_local + kh;
                    int w_in_local = w_conv_local + kw;
                    int input_idx = d_in_local * TILE_H_IN * TILE_W_IN +
                        h_in_local * TILE_W_IN + w_in_local;
                    int weight_idx = oc * CONV_KERNEL_VOL + kd * 9 + kh * 3
                        + kw;
                    accum_val = __hfma(sm_input_x[input_idx],
                        sm_weights[weight_idx], accum_val);
                }
            }
        }
        conv_accum[oc][p] = __hadd(conv_accum[oc][p], accum_val);
    }
    __syncthreads();
}

for (int oc = 0; oc < OUT_CHANNELS; ++oc) {
    for (int p = 0; p < 4; ++p) {
        int point_idx = threadIdx.x * 4 + p;
        int sm_idx = oc * (TILE_D_CONV * TILE_H_CONV * TILE_W_CONV) + point_idx;
        sm_conv_out[sm_idx] = conv_accum[oc][p];
    }
    __syncthreads();
}

if (threadIdx.x < (TILE_H_OUT * TILE_W_OUT)) {
    const int h_out_local = threadIdx.x / TILE_W_OUT;
    const int w_out_local = threadIdx.x % TILE_W_OUT;
    const int h_out = h_out_base + h_out_local;
    const int w_out = w_out_base + w_out_local;

    if (w_out < W_out && h_out < H_out) {
        half max_pooled_vals[OUT_CHANNELS];

```

```

for(int c = 0; c < OUT_CHANNELS; ++c) {
    const int h_conv_start = h_out_local * POOL_SIZE;
    const int w_conv_start = w_out_local * POOL_SIZE;

    const int sm_c_base = c * (TILE_D_CONV * TILE_H_CONV * TILE_W_CONV);
    const int sm_d0_base = sm_c_base;
    const int sm_d1_base = sm_c_base + TILE_H_CONV * TILE_W_CONV;

    half m0 = sm_conv_out[sm_d0_base + (h_conv_start+0)*TILE_W_CONV +
        (w_conv_start+0)];
    half m1 = sm_conv_out[sm_d0_base + (h_conv_start+0)*TILE_W_CONV +
        (w_conv_start+1)];
    half m2 = sm_conv_out[sm_d0_base + (h_conv_start+1)*TILE_W_CONV +
        (w_conv_start+0)];
    half m3 = sm_conv_out[sm_d0_base + (h_conv_start+1)*TILE_W_CONV +
        (w_conv_start+1)];
    half m4 = sm_conv_out[sm_d1_base + (h_conv_start+0)*TILE_W_CONV +
        (w_conv_start+0)];
    half m5 = sm_conv_out[sm_d1_base + (h_conv_start+0)*TILE_W_CONV +
        (w_conv_start+1)];
    half m6 = sm_conv_out[sm_d1_base + (h_conv_start+1)*TILE_W_CONV +
        (w_conv_start+0)];
    half m7 = sm_conv_out[sm_d1_base + (h_conv_start+1)*TILE_W_CONV +
        (w_conv_start+1)];

    half max_plane0 = __hmax(__hmax(m0, m1), __hmax(m2, m3));
    half max_plane1 = __hmax(__hmax(m4, m5), __hmax(m6, m7));
    max_pooled_vals[c] = __hadd(__hmax(max_plane0, max_plane1),
        bias[c]);
}

half m_0_1 = __hmax(max_pooled_vals[0], max_pooled_vals[1]); half m_2_3
    = __hmax(max_pooled_vals[2], max_pooled_vals[3]);
half m_4_5 = __hmax(max_pooled_vals[4], max_pooled_vals[5]); half m_6_7
    = __hmax(max_pooled_vals[6], max_pooled_vals[7]);
half m_8_9 = __hmax(max_pooled_vals[8], max_pooled_vals[9]); half
    m_10_11 = __hmax(max_pooled_vals[10], max_pooled_vals[11]);
half m_12_13 = __hmax(max_pooled_vals[12], max_pooled_vals[13]); half
    m_14_15 = __hmax(max_pooled_vals[14], max_pooled_vals[15]);
half m_0_3 = __hmax(m_0_1, m_2_3); half m_4_7 = __hmax(m_4_5, m_6_7);
half m_8_11 = __hmax(m_8_9, m_10_11); half m_12_15 = __hmax(m_12_13,
    m_14_15);
half m_0_7 = __hmax(m_0_3, m_4_7); half m_8_15 = __hmax(m_8_11,
    m_12_15);
half max_val = __hmax(m_0_7, m_8_15);

if (__hgt(max_val, __float2half(-65000.0f))) {
    float sum_exp_fp32 = 0.0f;
    #pragma unroll
    for(int k = 0; k < OUT_CHANNELS; ++k) {
        sum_exp_fp32 += __half2float(hexp(__hsub(max_pooled_vals[k],
            max_val)));
    }
    half lse_h = __hadd(max_val, hlog(__float2half_rn(sum_exp_fp32)));
    long long y_idx = (long long)n * D_out * H_out * W_out +
        (long long)d_out * H_out * W_out +
        (long long)h_out * W_out + w_out;
    y[y_idx] = __hmax(lse_h, (half)0.0f);
} else {
    long long y_idx = (long long)n * D_out * H_out * W_out +
        (long long)d_out * H_out * W_out +
        (long long)h_out * W_out + w_out;
    y[y_idx] = (half)0.0f;
}

```

```

    }
    }
}

torch::Tensor fused_conv_maxpool_bias_logsumexp_relu_cuda(torch::Tensor x,
    torch::Tensor weight, torch::Tensor bias)
{
    TORCH_CHECK(x.is_cuda() && x.is_contiguous() && x.scalar_type() ==
        torch::kHalf, "Input x invalid");
    TORCH_CHECK(weight.is_cuda() && weight.is_contiguous() && weight.scalar_type()
        == torch::kHalf, "Input weight invalid");
    TORCH_CHECK(bias.is_cuda() && bias.is_contiguous() && bias.scalar_type() ==
        torch::kHalf, "Input bias invalid");

    const int N = x.size(0);
    const int C_in = x.size(1);
    const int D_in = x.size(2);
    const int H_in = x.size(3);
    const int W_in = x.size(4);

    TORCH_CHECK(weight.dim() == 2, "Reshaped weight must be 2D");
    TORCH_CHECK(weight.size(0) == C_in, "Reshaped weight in_channels mismatch");
    TORCH_CHECK(bias.size(0) == OUT_CHANNELS, "Bias size must be 16");

    const int D_conv = (D_in + 2 * CONV_PADDING - CONV_KERNEL_SIZE) / CONV_STRIDE +
        1;
    const int H_conv = (H_in + 2 * CONV_PADDING - CONV_KERNEL_SIZE) / CONV_STRIDE +
        1;
    const int W_conv = (W_in + 2 * CONV_PADDING - CONV_KERNEL_SIZE) / CONV_STRIDE +
        1;

    const int D_out = D_conv / POOL_SIZE;
    const int H_out = H_conv / POOL_SIZE;
    const int W_out = W_conv / POOL_SIZE;

    auto y = torch::zeros({N, 1, D_out, H_out, W_out}, x.options());
    if (y.numel() == 0) return y;

    const dim3 threads_per_block(128);
    const dim3 num_blocks(
        (W_out + TILE_W_OUT - 1) / TILE_W_OUT,
        (H_out + TILE_H_OUT - 1) / TILE_H_OUT,
        (long long)N * D_out
    );

    size_t smem_size = (OUT_CHANNELS * CONV_KERNEL_VOL +
        TILE_D_IN * TILE_H_IN * TILE_W_IN +
        OUT_CHANNELS * TILE_D_CONV * TILE_H_CONV * TILE_W_CONV) *
        sizeof(half);

    fused_conv_maxpool_bias_logsumexp_relu_kernel<<<num_blocks, threads_per_block,
        smem_size>>>(
        (const half*)x.data_ptr<at::Half>(),
        (const half*)weight.data_ptr<at::Half>(),
        (const half*)bias.data_ptr<at::Half>(),
        (half*)y.data_ptr<at::Half>(),
        N, C_in, D_in, H_in, W_in, D_out, H_out, W_out
    );

    C10_CUDA_KERNEL_LAUNCH_CHECK();
    return y;
}

```

```

}
"""

cpp_source = """
#include <torch/extension.h>
torch::Tensor fused_conv_maxpool_bias_logsumexp_relu_cuda(torch::Tensor x,
    torch::Tensor weight, torch::Tensor bias);
"""

vertically_fused_op = load_inline(
    name='vertically_fused_op_v11_inc_addr',
    cpp_sources=cpp_source,
    cuda_sources=cuda_source,
    functions=['fused_conv_maxpool_bias_logsumexp_relu_cuda'],
    verbose=True,
    extra_cuda_cflags=['-U__CUDA_NO_HALF_OPERATORS__',
        '-U__CUDA_NO_HALF_CONVERSIONS__', '-U__CUDA_NO_HALF2_OPERATORS__',
        '--use_fast_math']
)

class ModelNew(nn.Module):
    def __init__(self, in_channels: int, out_channels: int, kernel_size: int,
        stride: int, padding: int):
        super(ModelNew, self).__init__()
        specialized_out_channels = 16
        conv_layer = nn.Conv3d(
            in_channels=in_channels,
            out_channels=specialized_out_channels,
            kernel_size=kernel_size,
            stride=stride,
            padding=padding,
            bias=True
        ).half()

        self.conv_weight = nn.Parameter(conv_layer.weight)
        self.conv_bias = nn.Parameter(conv_layer.bias)
        self.in_channels = in_channels

        self.fused_op =
            vertically_fused_op.fused_conv_maxpool_bias_logsumexp_relu_cuda

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x_half = x.half()
        weight_resaped = self.conv_weight.permute(1, 0, 2, 3,
            4).contiguous().view(self.in_channels, -1)
        output_half = self.fused_op(x_half, weight_resaped, self.conv_bias)
        return output_half.float()

```


D.4. ConvTranspose2d BiasAdd Clamp Scaling Clamp Divide

```

import torch
import torch.nn as nn
import torch.nn.functional as F
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

custom_kernel_source_unified = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>
#include <cmath> // For fminf

union float4_as_half8 {
    float4 f4;
    half2 h2[4];
};

__global__ void custom_kernel_nhwc_unified_po2(
    const float4* __restrict__ input,
    const __half* __restrict__ bias,
    float4* __restrict__ output,
    int size_in_float4,
    int padded_channels,
    float scaling_factor) {

    const half2 h2_zero = __float2half2_rn(0.0f);
    const half2 h2_upper_bound = __float2half2_rn(fminf(1.0f, 1.0f /
        scaling_factor));
    const int out_channels_mask = padded_channels - 1;

    for (int idx = blockIdx.x * blockDim.x + threadIdx.x;
        idx < size_in_float4;
        idx += gridDim.x * blockDim.x) {

        const int elem_idx_start = idx * 8;
        int c = elem_idx_start & out_channels_mask; // Fast modulo

        float4_as_half8 val;
        val.f4 = input[idx];

        __half b0 = bias[c]; c = (c + 1) & out_channels_mask;
        __half b1 = bias[c]; c = (c + 1) & out_channels_mask;
        half2 res0 = __hadd2(val.h2[0], __halves2half2(b0, b1));
        res0 = __hmax2(res0, h2_zero); res0 = __hmin2(res0, h2_upper_bound);
        val.h2[0] = res0;

        __half b2 = bias[c]; c = (c + 1) & out_channels_mask;
        __half b3 = bias[c]; c = (c + 1) & out_channels_mask;
        half2 res1 = __hadd2(val.h2[1], __halves2half2(b2, b3));
        res1 = __hmax2(res1, h2_zero); res1 = __hmin2(res1, h2_upper_bound);
        val.h2[1] = res1;

        __half b4 = bias[c]; c = (c + 1) & out_channels_mask;
        __half b5 = bias[c]; c = (c + 1) & out_channels_mask;
        half2 res2 = __hadd2(val.h2[2], __halves2half2(b4, b5));
        res2 = __hmax2(res2, h2_zero); res2 = __hmin2(res2, h2_upper_bound);
        val.h2[2] = res2;

        __half b6 = bias[c]; c = (c + 1) & out_channels_mask;

```

```

        __half b7 = bias[c];
        half2 res3 = __hadd2(val.h2[3], __halves2half2(b6, b7));
        res3 = __hmax2(res3, h2_zero); res3 = __hmin2(res3, h2_upper_bound);
        val.h2[3] = res3;

        output[idx] = val.f4;
    }
}

torch::Tensor custom_cuda_wrapper_unified(
    torch::Tensor input, torch::Tensor bias, int padded_channels, float
    scaling_factor) {

    const int total_elements = input.numel();
    auto output = torch::empty_like(input,
        input.options().memory_format(at::MemoryFormat::ChannelsLast));
    const int size_in_float4 = total_elements / 8;
    const int block_size = 256;
    const int num_blocks = (size_in_float4 + block_size - 1) / block_size;

    custom_kernel_nhwc_unified_po2<<<(num_blocks, block_size)>>>(
        reinterpret_cast<const float4*>(input.data_ptr<at::Half>()),
        reinterpret_cast<const __half*>(bias.data_ptr<at::Half>()),
        reinterpret_cast<float4*>(output.data_ptr<at::Half>()),
        size_in_float4, padded_channels, scaling_factor);

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) {
        TORCH_CHECK(false, "CUDA kernel launch failed: ", cudaGetErrorString(err));
    }

    return output;
}
"""

custom_cpp_source_unified = "torch::Tensor
    custom_cuda_wrapper_unified(torch::Tensor input, torch::Tensor bias, int
    padded_channels, float scaling_factor);"

custom_op_unified = load_inline(
    name='custom_op_unified',
    cpp_sources=custom_cpp_source_unified,
    cuda_sources=custom_kernel_source_unified,
    functions=['custom_cuda_wrapper_unified'],
    verbose=False,
    extra_cflags=['-O3'],
    extra_cuda_cflags=['-O3', '--use_fast_math', '-arch=sm_80']
)

class ModelNew(nn.Module):
    def __init__(self, in_channels, out_channels, kernel_size, stride, padding,
        output_padding, bias_shape, scaling_factor):
        super(ModelNew, self).__init__()

        self.conv_transpose = nn.ConvTranspose2d(
            in_channels, out_channels, kernel_size, stride=stride,
            padding=padding, output_padding=output_padding, bias=True
        ).half()

        second_bias = nn.Parameter(torch.randn(1, bias_shape[0], 1, 1).half())

        self.scaling_factor = scaling_factor
        self.custom_op = custom_op_unified

```

```

with torch.no_grad():
    combined_bias_data = self.conv_transpose.bias.data +
        second_bias.data.squeeze()

self.out_channels = out_channels
is_po2 = (self.out_channels > 0) and ((self.out_channels &
    (self.out_channels - 1)) == 0)

if not is_po2:
    self.padded_channels = 1 << (self.out_channels - 1).bit_length()
    padded_bias_data = torch.zeros(self.padded_channels,
        dtype=combined_bias_data.dtype, device=combined_bias_data.device)
    padded_bias_data[:self.out_channels] = combined_bias_data
else:
    self.padded_channels = self.out_channels
    padded_bias_data = combined_bias_data

self.bias = nn.Parameter(padded_bias_data)

del self.conv_transpose.bias
self.conv_transpose.bias = None

self.graph = None
self.static_input = None
self.static_output = None

def forward(self, x: torch.Tensor) -> torch.Tensor:
    original_dtype = x.dtype
    if x.is_contiguous(memory_format=torch.channels_last):
        original_format = torch.channels_last
    else:
        original_format = torch.contiguous_format

    x_half = x.half()
    x_nchw = x_half.to(memory_format=torch.channels_last)

    if self.graph is None:
        conv_out_no_bias = self.conv_transpose(x_nchw)
        output_nchw_half = self.custom_op.custom_cuda_wrapper_unified(
            conv_out_no_bias, self.bias, self.padded_channels,
            self.scaling_factor
        )

        self.static_input = torch.empty_like(x_nchw)
        self.graph = torch.cuda.CUDAGraph()
        with torch.cuda.graph(self.graph):
            conv_out_static = self.conv_transpose(self.static_input)
            self.static_output = self.custom_op.custom_cuda_wrapper_unified(
                conv_out_static, self.bias, self.padded_channels,
                self.scaling_factor
            )

        return output_nchw_half.to(dtype=original_dtype,
            memory_format=original_format)

    else:
        self.static_input.copy_(x_nchw)
        self.graph.replay()
        return self.static_output.to(dtype=original_dtype,
            memory_format=original_format)

```

D.5. GELU

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

class ModelNew(nn.Module):
    def __init__(self, num_features: int = 0):
        super(ModelNew, self).__init__()

        gelu_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <cmath>

const float GELU_SCALAR_1 = 0.7978845608028654f; // sqrt(2.0 / PI)
const float GELU_SCALAR_2 = 0.044715f;

__device__ __forceinline__ float gelu_fma(const float x) {
    const float x_sq = x * x;
    const float inner = fmaf(GELU_SCALAR_2 * x_sq, x, x);
    const float v = GELU_SCALAR_1 * inner;
    const float t = tanhf(v);
    const float cdf = fmaf(0.5f, t, 0.5f);
    return x * cdf;
}

__global__ void gelu_grid_stride_float4_fma_kernel(const float* __restrict__ input,
float* __restrict__ output, const int size) {
    const int num_float4_elements = size / 4;
    const int grid_stride = gridDim.x * blockDim.x;
    for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < num_float4_elements; i
+= grid_stride) {
        const float4 in_vec = reinterpret_cast<const float4*>(input)[i];

        float4 out_vec;
        out_vec.x = gelu_fma(in_vec.x);
        out_vec.y = gelu_fma(in_vec.y);
        out_vec.z = gelu_fma(in_vec.z);
        out_vec.w = gelu_fma(in_vec.w);

        reinterpret_cast<float4*>(output)[i] = out_vec;
    }

    const int tail_start_idx = num_float4_elements * 4;
    const int thread_id = blockIdx.x * blockDim.x + threadIdx.x;
    const int tail_idx = tail_start_idx + thread_id;

    if (tail_idx < size) {
        output[tail_idx] = gelu_fma(input[tail_idx]);
    }
}

torch::Tensor gelu_cuda(torch::Tensor input) {
    const auto size = input.numel();
    auto output = torch::empty_like(input);

    if (size == 0) {
        return output;

```

```

    }

    const int block_size = 1024;
    const int num_blocks = 108 * 2;
    gelu_grid_stride_float4_fma_kernel<<<num_blocks, block_size>>>(
        input.data_ptr<float>(),
        output.data_ptr<float>(),
        size
    );

    cudaError_t err = cudaGetLastError();
    TORCH_CHECK(err == cudaSuccess, "CUDA kernel launch failed: ",
        cudaGetErrorString(err));

    return output;
}
"""

    gelu_cpp_source = """
#include <torch/extension.h>

torch::Tensor gelu_cuda(torch::Tensor input);
"""

    self.gelu_cuda_module = load_inline(
        name='gelu_cuda_fixedgrid_A100',
        cpp_sources=gelu_cpp_source,
        cuda_sources=gelu_source,
        functions=['gelu_cuda'],
        verbose=True,
        extra_cuda_cflags=['--use_fast_math', '-O3']
    )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        return self.gelu_cuda_module.gelu_cuda(x)

```


D.6. Matmul with large K dimension

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

matmul_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>

#define THREADS_X 16
#define THREADS_Y 16
#define VEC_C_Y 4
#define VEC_C_X 4
#define BLOCK_ROWS (THREADS_Y * VEC_C_Y) // 16 * 4 = 64
#define BLOCK_COLS (THREADS_X * VEC_C_X) // 16 * 4 = 64
#define K_STEP 64
#define SPLIT_K_FACTOR 128

__global__ void matmul_splitk_atomic_kernel(const float* A, const float* B, float*
    C,
                                     int M, int N, int K) {
    __shared__ float As[BLOCK_ROWS][K_STEP]; // 64 x 64
    __shared__ float Bs[K_STEP][BLOCK_COLS]; // 64 x 64
    const int tx = threadIdx.x; // 0-15
    const int ty = threadIdx.y; // 0-15
    const int bx = blockIdx.x;
    const int by = blockIdx.y;
    const int bz = blockIdx.z;
    const int c_base_row = by * BLOCK_ROWS + ty * VEC_C_Y;
    const int c_base_col = bx * BLOCK_COLS + tx * VEC_C_X;
    float accum[VEC_C_Y][VEC_C_X] = {{0.0f}};
    const int K_per_split = K / SPLIT_K_FACTOR;
    const int k_start_idx = bz * K_per_split;
    const int num_tiles_per_split = K_per_split / K_STEP;

    for (int t = 0; t < num_tiles_per_split; ++t) {
        const int global_k_idx = k_start_idx + t * K_STEP;

        #pragma unroll
        for (int i = 0; i < 4; ++i) {
            int smem_row = ty * 4 + i;
            int gmem_row = by * BLOCK_ROWS + smem_row;

            if (gmem_row < M) {
                reinterpret_cast<float4*>(As[smem_row])[tx] =
                    reinterpret_cast<const float4*>(&A[gmem_row * K +
                        global_k_idx])[tx];
            } else {
                reinterpret_cast<float4*>(As[smem_row])[tx] = make_float4(0.0f,
                    0.0f, 0.0f, 0.0f);
            }
        }

        #pragma unroll
        for (int i = 0; i < 4; ++i) {
            int smem_row = ty * 4 + i;
            int gmem_row = global_k_idx + smem_row;

            if (gmem_row < K) {

```

```

        reinterpret_cast<float4*>(Bs[smem_row])[tx] =
            reinterpret_cast<const float4*>(&B[gmem_row * N + bx *
                BLOCK_COLS])[tx];
    } else {
        reinterpret_cast<float4*>(Bs[smem_row])[tx] = make_float4(0.0f,
            0.0f, 0.0f, 0.0f);
    }
}

__syncthreads();

#pragma unroll
for (int k = 0; k < K_STEP; ++k) {
    float a_reg[VEC_C_Y];
    float b_reg[VEC_C_X];

    #pragma unroll
    for(int i = 0; i < VEC_C_Y; ++i) {
        a_reg[i] = As[ty * VEC_C_Y + i][k];
    }

    #pragma unroll
    for(int j = 0; j < VEC_C_X; ++j) {
        b_reg[j] = Bs[k][tx * VEC_C_X + j];
    }

    #pragma unroll
    for(int i = 0; i < VEC_C_Y; ++i) {
        #pragma unroll
        for(int j = 0; j < VEC_C_X; ++j) {
            accum[i][j] += a_reg[i] * b_reg[j];
        }
    }
}
__syncthreads();
}

#pragma unroll
for(int i = 0; i < VEC_C_Y; ++i) {
    #pragma unroll
    for(int j = 0; j < VEC_C_X; ++j) {
        if (c_base_row + i < M && c_base_col + j < N) {
            atomicAdd(&C[(c_base_row + i) * N + (c_base_col + j)], accum[i][j]);
        }
    }
}
}

torch::Tensor matmul_cuda(torch::Tensor a, torch::Tensor b) {
    a = a.contiguous();
    b = b.contiguous();
    const int M = a.size(0);
    const int K = a.size(1);
    const int N = b.size(1);
    auto c = torch::zeros({M, N}, a.options());
    dim3 threadsPerBlock(THREADS_X, THREADS_Y);
    dim3 numBlocks((N + BLOCK_COLS - 1) / BLOCK_COLS,
        (M + BLOCK_ROWS - 1) / BLOCK_ROWS,
        SPLIT_K_FACTOR);

    matmul_splitk_atomic_kernel<<<numBlocks, threadsPerBlock>>>(
        a.data_ptr<float>(),

```

```

        b.data_ptr<float>(),
        c.data_ptr<float>(),
        M, N, K);

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) {
        throw std::runtime_error(cudaGetErrorString(err));
    }

    return c;
}
"""

matmul_cpp_source = """
torch::Tensor matmul_cuda(torch::Tensor a, torch::Tensor b);
"""

matmul_cuda_module = load_inline(
    name='matmul_cuda_splitk128_atomic',
    cpp_sources=matmul_cpp_source,
    cuda_sources=matmul_source,
    functions=['matmul_cuda'],
    extra_cuda_cflags=['-O3', '--use_fast_math'],
    verbose=False
)

class ModelNew(nn.Module):
    def __init__(self, num_features: int = 0):
        super(ModelNew, self).__init__()
        self.matmul_cuda = matmul_cuda_module

    def forward(self, A: torch.Tensor, B: torch.Tensor) -> torch.Tensor:
        a_cuda = A if A.is_cuda else A.cuda()
        b_cuda = B if B.is_cuda else B.cuda()

        return self.matmul_cuda.matmul_cuda(a_cuda, b_cuda)

```

D.7. Max pooling 2D

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

class ModelNew(nn.Module):
    def __init__(self, kernel_size: int, stride: int, padding: int, dilation: int):
        super(ModelNew, self).__init__()
        self.kernel_size = kernel_size
        self.stride = stride
        self.padding = padding
        self.dilation = dilation
        BLOCK_DIM_X = 32
        BLOCK_DIM_Y = 8
        TILE_DIM_X = 2
        TILE_DIM_Y = 2

        cuda_source = f"""
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <float.h>

#define KERNEL_SIZE {self.kernel_size}
#define STRIDE {self.stride}
#define PADDING {self.padding}
#define DILATION {self.dilation}
#define BLOCK_DIM_X {BLOCK_DIM_X}
#define BLOCK_DIM_Y {BLOCK_DIM_Y}
#define TILE_DIM_X {TILE_DIM_X}
#define TILE_DIM_Y {TILE_DIM_Y}

__global__ void max_pool2d_grid_stride_kernel(
    const float* __restrict__ input,
    float* __restrict__ output,
    int input_height,
    int input_width,
    int output_height,
    int output_width,
    int num_planes) // Total number of planes (batch_size * channels)
{{
    const int w_out_base = (blockIdx.x * blockDim.x + threadIdx.x) * TILE_DIM_X;
    const int h_out_base = (blockIdx.y * blockDim.y + threadIdx.y) * TILE_DIM_Y;
    if (h_out_base >= output_height) {{
        return;
    }}

    for (int plane_idx = blockIdx.z; plane_idx < num_planes; plane_idx += gridDim.z)
    {{
        const float* input_map = input + plane_idx * (long)input_height *
            input_width;
        float* output_map = output + plane_idx * (long)output_height * output_width;

        float max_val[TILE_DIM_Y][TILE_DIM_X];
        #pragma unroll
        for(int i = 0; i < TILE_DIM_Y; ++i) {{
            #pragma unroll
            for(int j = 0; j < TILE_DIM_X; ++j) {{
                max_val[i][j] = -FLT_MAX;
            }}
        }}
    }}
}
"""
        load_inline("max_pool2d_grid_stride_kernel", cuda_source)

```

```

    }}
}}

const int h_starts[TILE_DIM_Y] = {{h_out_base * STRIDE - PADDING,
    (h_out_base + 1) * STRIDE - PADDING}};
const int w_starts[TILE_DIM_X] = {{w_out_base * STRIDE - PADDING,
    (w_out_base + 1) * STRIDE - PADDING}};

#pragma unroll
for (int i = 0; i < KERNEL_SIZE; ++i) {{
    const int h_in[TILE_DIM_Y] = {{h_starts[0] + i * DILATION, h_starts[1]
        + i * DILATION}};

    const bool h_in_valid[TILE_DIM_Y] = {{
        (unsigned)h_in[0] < (unsigned)input_height,
        (unsigned)h_in[1] < (unsigned)input_height
    }};

    #pragma unroll
    for (int j = 0; j < KERNEL_SIZE; ++j) {{
        const int w_in[TILE_DIM_X] = {{w_starts[0] + j * DILATION,
            w_starts[1] + j * DILATION}};

        const bool w_in_valid[TILE_DIM_X] = {{
            (unsigned)w_in[0] < (unsigned)input_width,
            (unsigned)w_in[1] < (unsigned)input_width
        }};

        max_val[0][0] = fmaxf(max_val[0][0], (h_in_valid[0] &&
            w_in_valid[0]) ? input_map[(long)h_in[0] * input_width +
            w_in[0]] : -FLT_MAX);
        max_val[0][1] = fmaxf(max_val[0][1], (h_in_valid[0] &&
            w_in_valid[1]) ? input_map[(long)h_in[0] * input_width +
            w_in[1]] : -FLT_MAX);
        max_val[1][0] = fmaxf(max_val[1][0], (h_in_valid[1] &&
            w_in_valid[0]) ? input_map[(long)h_in[1] * input_width +
            w_in[0]] : -FLT_MAX);
        max_val[1][1] = fmaxf(max_val[1][1], (h_in_valid[1] &&
            w_in_valid[1]) ? input_map[(long)h_in[1] * input_width +
            w_in[1]] : -FLT_MAX);
    }}
}}

#pragma unroll
for(int i = 0; i < TILE_DIM_Y; ++i) {{
    const int h_out = h_out_base + i;
    if (h_out < output_height) {{
        #pragma unroll
        for(int j = 0; j < TILE_DIM_X; ++j) {{
            const int w_out = w_out_base + j;
            if (w_out < output_width) {{
                output_map[(long)h_out * output_width + w_out] =
                    max_val[i][j];
            }}
        }}
    }}
}}
}}

torch::Tensor max_pool2d_grid_stride_cuda(torch::Tensor x) {{
    const int batch_size = x.size(0);

```

```

const int channels = x.size(1);
const int input_height = x.size(2);
const int input_width = x.size(3);

const int output_height = (input_height + 2 * PADDING - DILATION *
    (KERNEL_SIZE - 1) - 1) / STRIDE + 1;
const int output_width = (input_width + 2 * PADDING - DILATION *
    (KERNEL_SIZE - 1) - 1) / STRIDE + 1;

auto options = torch::TensorOptions().dtype(x.dtype()).device(x.device());
auto output = torch::empty({batch_size, channels, output_height,
    output_width}), options);

const int num_planes = batch_size * channels;
if (num_planes == 0) return output;

const int max_z_blocks = 1024;

const dim3 block_dim(BLOCK_DIM_X, BLOCK_DIM_Y, 1);
const dim3 grid_dim(
    (output_width + block_dim.x * TILE_DIM_X - 1) / (block_dim.x * TILE_DIM_X),
    (output_height + block_dim.y * TILE_DIM_Y - 1) / (block_dim.y *
        TILE_DIM_Y),
    (num_planes > max_z_blocks) ? max_z_blocks : num_planes
);

max_pool2d_grid_stride_kernel<<<grid_dim, block_dim>>> (
    x.data_ptr<float>(),
    output.data_ptr<float>(),
    input_height,
    input_width,
    output_height,
    output_width,
    num_planes
);

cudaError_t err = cudaGetLastError();
if (err != cudaSuccess) {{
    // Proper error handling should be here
}}
return output;
}}
"""

cpp_source = """
torch::Tensor max_pool2d_grid_stride_cuda(torch::Tensor x);
"""

self.max_pool2d_module = load_inline(
    name=f'max_pool2d_grid_stride_v3_{BLOCK_DIM_X}x{BLOCK_DIM_Y}_k{self.kernel_size}_s{self.s}',
    cpp_sources=cpp_source,
    cuda_sources=cuda_source,
    functions=['max_pool2d_grid_stride_cuda'],
    verbose=False,
    extra_cuda_cflags=['-O3', '--use_fast_math']
)

def forward(self, x: torch.Tensor) -> torch.Tensor:
    if not x.is_cuda:
        x = x.cuda()
    if not x.is_contiguous(memory_format=torch.contiguous_format):
        x = x.contiguous()

    return self.max_pool2d_module.max_pool2d_grid_stride_cuda(x)

```

D.8. MLP

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

cuda_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>

union float4_longlong2_union {
    float4 f4;
    longlong2 ll2;
};

__global__ void linear_relu_fma_1x4_tile_4x_unroll_4acc(
    const float* __restrict__ input,
    const float* __restrict__ weight,
    const float* __restrict__ bias,
    float* __restrict__ output,
    int batch_size,
    int input_size,
    int output_size)
{
    const int TILE_M = 1;
    const int TILE_N = 4;
    const int batch_idx = blockIdx.y;
    const int block_tile_n = blockIdx.x * TILE_N;
    const int warp_id = threadIdx.y; // 0..3
    const int lane_id = threadIdx.x; // 0..31
    const int output_idx = block_tile_n + warp_id;

    if (batch_idx >= batch_size || output_idx >= output_size) {
        return;
    }

    const float* input_vec = input + batch_idx * input_size;
    const float* weight_row = weight + output_idx * input_size;
    float sum_a = 0.0f;
    float sum_b = 0.0f;
    float sum_c = 0.0f;
    float sum_d = 0.0f;

    const int vectorized_limit = input_size / 4;
    const int stride = blockDim.x * 4;

    if (vectorized_limit > 0) {
        const longlong2* input_ll2 = reinterpret_cast<const longlong2*>(input_vec);
        const longlong2* weight_ll2 = reinterpret_cast<const
            longlong2*>(weight_row);

        float4_longlong2_union in_val, wt_val;

        for (int k = lane_id; k < vectorized_limit; k += stride) {
            in_val.ll2 = input_ll2[k];
            wt_val.ll2 = weight_ll2[k];
            sum_a += in_val.f4.x * wt_val.f4.x + in_val.f4.y * wt_val.f4.y +
                in_val.f4.z * wt_val.f4.z + in_val.f4.w * wt_val.f4.w;

```



```

        if (k + blockDim.x < vectorized_limit) {
            in_val.ll2 = input_ll2[k + blockDim.x];
            wt_val.ll2 = weight_ll2[k + blockDim.x];
            sum_b += in_val.f4.x * wt_val.f4.x + in_val.f4.y * wt_val.f4.y +
                    in_val.f4.z * wt_val.f4.z + in_val.f4.w * wt_val.f4.w;
        }

        if (k + blockDim.x * 2 < vectorized_limit) {
            in_val.ll2 = input_ll2[k + blockDim.x * 2];
            wt_val.ll2 = weight_ll2[k + blockDim.x * 2];
            sum_c += in_val.f4.x * wt_val.f4.x + in_val.f4.y * wt_val.f4.y +
                    in_val.f4.z * wt_val.f4.z + in_val.f4.w * wt_val.f4.w;
        }

        if (k + blockDim.x * 3 < vectorized_limit) {
            in_val.ll2 = input_ll2[k + blockDim.x * 3];
            wt_val.ll2 = weight_ll2[k + blockDim.x * 3];
            sum_d += in_val.f4.x * wt_val.f4.x + in_val.f4.y * wt_val.f4.y +
                    in_val.f4.z * wt_val.f4.z + in_val.f4.w * wt_val.f4.w;
        }
    }

    float sum = sum_a + sum_b + sum_c + sum_d;

    for (int i = vectorized_limit * 4 + lane_id; i < input_size; i += blockDim.x) {
        sum += input_vec[i] * weight_row[i];
    }

    #pragma unroll
    for (int offset = 16; offset > 0; offset /= 2) {
        sum += __shfl_down_sync(0xffffffff, sum, offset);
    }

    if (lane_id == 0) {
        sum += bias[output_idx];
        output[batch_idx * output_size + output_idx] = fmaxf(0.0f, sum);
    }
}

torch::Tensor linear_relu_cuda_1x4_4x_4acc(torch::Tensor input, torch::Tensor
weight, torch::Tensor bias) {
    input = input.contiguous();
    weight = weight.contiguous();
    bias = bias.contiguous();

    const int batch_size = input.size(0);
    const int input_size = input.size(1);
    const int output_size = weight.size(0);

    auto output = torch::empty({batch_size, output_size}, input.options());

    const int TILE_N = 4;
    const int warps_per_block = 4;

    dim3 threads(32, warps_per_block); // 128 threads per block
    dim3 blocks(
        (output_size + TILE_N - 1) / TILE_N,
        batch_size
    );

    linear_relu_fma_1x4_tile_4x_unroll_4acc<<<blocks, threads>>>(
        input.data_ptr<float>(),

```

```

        weight.data_ptr<float>(),
        bias.data_ptr<float>(),
        output.data_ptr<float>(),
        batch_size,
        input_size,
        output_size
    );

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) {
        TORCH_CHECK(false, "CUDA kernel launch failed: ", cudaGetErrorString(err));
    }

    return output;
}
"""

cpp_source = "torch::Tensor linear_relu_cuda_1x4_4x_4acc(torch::Tensor input,
    torch::Tensor weight, torch::Tensor bias);"

fused_op_1x4_4x_4acc = load_inline(
    name='fused_op_1x4_4x_4acc',
    cpp_sources=cpp_source,
    cuda_sources=cuda_source,
    functions=['linear_relu_cuda_1x4_4x_4acc'],
    verbose=True,
    extra_cflags=['-O3'],
    extra_cuda_cflags=['-O3', '--use_fast_math']
)

class ModelNew(nn.Module):
    def __init__(self, input_size, layer_sizes, output_size):
        super(ModelNew, self).__init__()
        self.layers = nn.ModuleList()
        current_input_size = input_size
        for layer_size in layer_sizes:
            self.layers.append(nn.Linear(current_input_size, layer_size))
            current_input_size = layer_size
        self.final_linear = nn.Linear(current_input_size, output_size)

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        for layer in self.layers:
            x = fused_op_1x4_4x_4acc.linear_relu_cuda_1x4_4x_4acc(x, layer.weight,
                layer.bias)
        x = self.final_linear(x)
        return x

```

D.9. RMSNorm

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
import math
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

hybrid_rms_norm_cuda_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>

constexpr int BLOCK_H = 32;
constexpr int BLOCK_W = 32;
constexpr int VEC_SIZE = 4;
constexpr int TILE_W = BLOCK_W * VEC_SIZE;
constexpr int REDUCE_PARTITION_SIZE = 8;
using float4 = float4;

__global__ void rms_norm_nchw_2x8_w_tiling_kernel(const float* __restrict__ x,
    float* __restrict__ out,
    const int N, const int C, const
    int H, const int W,
    const float eps, const float
    inv_c) {

    const int h_pair_idx = blockIdx.y;
    const int tx = threadIdx.x;
    const int ty = threadIdx.y;
    const int C_H_W = C * H * W;
    const int H_W = H * W;
    const int UNROLL_STRIDE = BLOCK_H * 2;
    __shared__ float s_tile[BLOCK_H][TILE_W];

    #pragma unroll
    for (int w_iter = 0; w_iter < 8; ++w_iter) {
        const int w_idx_base = blockIdx.x * (TILE_W * 8) + w_iter * TILE_W;
        if (w_idx_base >= W) continue;

        const int w_offset_base = w_idx_base + tx * VEC_SIZE;

        #pragma unroll
        for (int h_iter = 0; h_iter < 2; ++h_iter) {
            const int global_h_idx = h_pair_idx * 2 + h_iter;
            if (global_h_idx >= N * H) continue;

            const int n_idx = global_h_idx / H;
            const int h_idx = global_h_idx % H;
            float sum_sq[VEC_SIZE] = {0.0f, 0.0f, 0.0f, 0.0f};
            const int w_offset = w_offset_base;

            for (int c_base = ty; c_base < C; c_base += UNROLL_STRIDE) {
                if (w_offset + VEC_SIZE - 1 < W) {
                    float4 val_vec_1 = *reinterpret_cast<const
                        float4*>(&x[n_idx*C_H_W + c_base*H_W + h_idx*W + w_offset]);
                    sum_sq[0] += val_vec_1.x * val_vec_1.x;
                    sum_sq[1] += val_vec_1.y * val_vec_1.y;
                    sum_sq[2] += val_vec_1.z * val_vec_1.z;
                    sum_sq[3] += val_vec_1.w * val_vec_1.w;

                    const int c_base_2 = c_base + BLOCK_H;

```

```

        if (c_base_2 < C) {
            float4 val_vec_2 = *reinterpret_cast<const
                float4*>(&x[n_idx*C_H_W + c_base_2*H_W + h_idx*W +
                    w_offset]);
            sum_sq[0] += val_vec_2.x * val_vec_2.x;
            sum_sq[1] += val_vec_2.y * val_vec_2.y;
            sum_sq[2] += val_vec_2.z * val_vec_2.z;
            sum_sq[3] += val_vec_2.w * val_vec_2.w;
        }
    } else { // Handle ragged edge of W
        #pragma unroll
        for (int i = 0; i < VEC_SIZE; ++i) {
            if (w_offset + i < W) {
                float val1 = x[n_idx*C_H_W + c_base*H_W + h_idx*W +
                    w_offset + i];
                sum_sq[i] += val1 * val1;
                const int c_base_2 = c_base + BLOCK_H;
                if (c_base_2 < C) {
                    float val2 = x[n_idx*C_H_W + c_base_2*H_W + h_idx*W
                        + w_offset + i];
                    sum_sq[i] += val2 * val2;
                }
            }
        }
    }
}

__syncthreads();

*reinterpret_cast<float4*>(&s_tile[ty][tx * VEC_SIZE]) =
    *reinterpret_cast<float4*>(sum_sq);
__syncthreads();

if ((ty % REDUCE_PARTITION_SIZE) == 0 && ty < BLOCK_H) {
    float4* my_sum_vec = reinterpret_cast<float4*>(&s_tile[ty][tx *
        VEC_SIZE]);
    #pragma unroll
    for (int i = 1; i < REDUCE_PARTITION_SIZE; ++i) {
        const float4* other_row = reinterpret_cast<const
            float4*>(&s_tile[ty + i][tx * VEC_SIZE]);
        my_sum_vec->x += other_row->x; my_sum_vec->y += other_row->y;
        my_sum_vec->z += other_row->z; my_sum_vec->w += other_row->w;
    }
}
__syncthreads();

if (ty == 0) {
    float4* final_sum_vec = reinterpret_cast<float4*>(&s_tile[0][tx *
        VEC_SIZE]);
    #pragma unroll
    for (int i = 1; i < BLOCK_H / REDUCE_PARTITION_SIZE; ++i) {
        const float4* partition_sum = reinterpret_cast<const
            float4*>(&s_tile[i * REDUCE_PARTITION_SIZE][tx * VEC_SIZE]);
        final_sum_vec->x += partition_sum->x; final_sum_vec->y +=
            partition_sum->y;
        final_sum_vec->z += partition_sum->z; final_sum_vec->w +=
            partition_sum->w;
    }
    final_sum_vec->x = rsqrtf(final_sum_vec->x * inv_c + eps);
    final_sum_vec->y = rsqrtf(final_sum_vec->y * inv_c + eps);
    final_sum_vec->z = rsqrtf(final_sum_vec->z * inv_c + eps);
    final_sum_vec->w = rsqrtf(final_sum_vec->w * inv_c + eps);
}
__syncthreads();

```

```

const float4 inv_rms_vec = *reinterpret_cast<const
    float4*>(&s_tile[0][tx * VEC_SIZE]);

for (int c_idx = ty; c_idx < C; c_idx += UNROLL_STRIDE) {
    if (w_offset + VEC_SIZE - 1 < W) {
        float4 val_vec1 = *reinterpret_cast<const
            float4*>(&x[n_idx*C_H_W + c_idx*H_W + h_idx*W + w_offset]);
        val_vec1.x *= inv_rms_vec.x; val_vec1.y *= inv_rms_vec.y;
        val_vec1.z *= inv_rms_vec.z; val_vec1.w *= inv_rms_vec.w;
        *reinterpret_cast<float4*>(&out[n_idx*C_H_W + c_idx*H_W +
            h_idx*W + w_offset]) = val_vec1;
        const int c_idx_2 = c_idx + BLOCK_H;
        if (c_idx_2 < C) {
            float4 val_vec2 = *reinterpret_cast<const
                float4*>(&x[n_idx*C_H_W + c_idx_2*H_W + h_idx*W +
                    w_offset]);
            val_vec2.x *= inv_rms_vec.x; val_vec2.y *= inv_rms_vec.y;
            val_vec2.z *= inv_rms_vec.z; val_vec2.w *= inv_rms_vec.w;
            *reinterpret_cast<float4*>(&out[n_idx*C_H_W + c_idx_2*H_W +
                h_idx*W + w_offset]) = val_vec2;
        }
    } else {
        float inv_rms_arr[VEC_SIZE] = {inv_rms_vec.x, inv_rms_vec.y,
            inv_rms_vec.z, inv_rms_vec.w};
        #pragma unroll
        for (int i = 0; i < VEC_SIZE; ++i) {
            if (w_offset + i < W) {
                out[n_idx*C_H_W + c_idx*H_W + h_idx*W + w_offset + i] =
                    x[n_idx*C_H_W + c_idx*H_W + h_idx*W + w_offset + i]
                        * inv_rms_arr[i];
                const int c_idx_2 = c_idx + BLOCK_H;
                if (c_idx_2 < C) {
                    out[n_idx*C_H_W + c_idx_2*H_W + h_idx*W + w_offset +
                        i] = x[n_idx*C_H_W + c_idx_2*H_W + h_idx*W +
                            w_offset + i] * inv_rms_arr[i];
                }
            }
        }
    }
}
}
}
}

```

```
__global__ void rms_norm_nhwc_kernel(const float* __restrict__ x, float*
__restrict__ out,
                                   const int num_vectors, const int C, const
                                   float eps, const float inv_c) {
    const int warp_id = blockIdx.x * (blockDim.x / 32) + (threadIdx.x / 32);
    if (warp_id >= num_vectors) {
        return;
    }

    const int lane_id = threadIdx.x % 32;
    const float* x_vec = x + warp_id * C;
    float* out_vec = out + warp_id * C;

    float sum_sq = 0.0f;
    for (int i = lane_id * VEC_SIZE; i < C; i += 32 * VEC_SIZE) {
        if (i + VEC_SIZE <= C) {
            float4 val = *reinterpret_cast<const float4*>(&x_vec[i]);
            sum_sq += val.x * val.x + val.y * val.y + val.z * val.z + val.w * val.w;
        } else {
```

```

        for (int j = i; j < C; ++j) {
            float val = x_vec[j];
            sum_sq += val * val;
        }
    }

#pragma unroll
for (int offset = 16; offset > 0; offset /= 2) {
    sum_sq += __shfl_xor_sync(0xffffffff, sum_sq, offset);
}
sum_sq = __shfl_sync(0xffffffff, sum_sq, 0);

const float inv_rms = rsqrtf(sum_sq * inv_c + eps);

for (int i = lane_id * VEC_SIZE; i < C; i += 32 * VEC_SIZE) {
    if (i + VEC_SIZE <= C) {
        float4 val = *reinterpret_cast<const float4*>(&x_vec[i]);
        val.x *= inv_rms; val.y *= inv_rms; val.z *= inv_rms; val.w *= inv_rms;
        *reinterpret_cast<float4*>(&out_vec[i]) = val;
    } else {
        for (int j = i; j < C; ++j) {
            out_vec[j] = x_vec[j] * inv_rms;
        }
    }
}
}

torch::Tensor rms_norm_hybrid_cuda(torch::Tensor x, float eps) {
    auto out = torch::empty_like(x);
    const int N = x.size(0);
    const int C = x.size(1);
    const int H = x.size(2);
    const int W = x.size(3);

    if (N*C*H*W == 0) return out;
    const float inv_c = 1.0f / C;

    cudaError_t err;

    if (x.is_contiguous(at::MemoryFormat::Contiguous)) {
        const dim3 block_dim(BLOCK_W, BLOCK_H);
        const dim3 grid_dim( (W + (TILE_W * 8) - 1) / (TILE_W * 8), (N * H + 1) / 2 );

        rms_norm_nchw_2x8_w_tiling_kernel<<<grid_dim, block_dim>>>(&x.data_ptr<float>(), &out.data_ptr<float>(), N, C, H, W, eps, inv_c);
    } else if (x.is_contiguous(at::MemoryFormat::ChannelsLast)) {
        const int num_vectors = N * H * W;
        const int threads_per_block = 256;
        const int warps_per_block = threads_per_block / 32;
        const int blocks = (num_vectors + warps_per_block - 1) / warps_per_block;

        rms_norm_nhwc_kernel<<<blocks, threads_per_block>>>(&x.data_ptr<float>(), &out.data_ptr<float>(), num_vectors, C, eps, inv_c);
    } else {
        AT_ERROR("Unsupported memory format. Input must be contiguous (NCHW) or channels_last (NHWC).");
    }

    err = cudaGetLastError();
}

```

```

    if (err != cudaSuccess) {
        AT_ERROR("CUDA kernel launch failed: ", cudaGetErrorString(err));
    }

    return out;
}
"""

hybrid_rms_norm_cpp_source = """
#include <torch/extension.h>
torch::Tensor rms_norm_hybrid_cuda(torch::Tensor x, float eps);
"""

hybrid_rms_norm_cached = load_inline(
    name='hybrid_rms_norm_v13_2x8_tiling',
    cpp_sources=hybrid_rms_norm_cpp_source,
    cuda_sources=hybrid_rms_norm_cuda_source,
    functions=['rms_norm_hybrid_cuda'],
    verbose=True,
    extra_cuda_cflags=['-std=c++17', '-O3', '-U__CUDA_NO_HALF_OPERATORS__',
                       '-U__CUDA_NO_HALF_CONVERSIONS__', '-U__CUDA_NO_HALF2_OPERATORS__']
)

class ModelNew(nn.Module):
    def __init__(self, num_features: int, eps: float = 1e-5):
        super().__init__()
        self.num_features = num_features
        self.eps = eps
        self.rms_norm_hybrid = hybrid_rms_norm_cached.rms_norm_hybrid_cuda

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        is_supported = (
            x.is_cuda and
            x.dtype == torch.float32 and
            x.dim() == 4 and
            (x.is_contiguous(memory_format=torch.contiguous_format) or
             x.is_contiguous(memory_format=torch.channels_last))
        )

        if not is_supported:
            original_dtype = x.dtype
            original_format = torch.channels_last if
                x.is_contiguous(memory_format=torch.channels_last) else
                torch.contiguous_format

            x_float = x.to(torch.float32)
            variance = x_float.pow(2).mean(dim=1, keepdim=True)
            output = x_float * torch.rsqrt(variance + self.eps)

            return
            output.to(original_dtype).contiguous(memory_format=original_format)

        return self.rms_norm_hybrid(x, self.eps)

```


D.10. Softmax

PACEvolve-generated kernel for Softmax contains a few other cases. Due to space constraints, we show the kernel path that is used in KernelBench evaluation.

```
import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

softmax_cuda_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <cuda_bf16.h>
#include <cfloat>

constexpr int WARP_SIZE = 32;
constexpr int MAX_DIM_SINGLE_PASS = 8192;
constexpr int MAX_DIM_MULTI_ROW = 1024;

__device__ __forceinline__ float warp_reduce_sum(float val) {
    for (int offset = WARP_SIZE / 2; offset > 0; offset >= 1) {
        val += __shfl_down_sync(0xffffffff, val, offset);
    }
    return val;
}

__device__ __forceinline__ float warp_reduce_max(float val) {
    for (int offset = WARP_SIZE / 2; offset > 0; offset >= 1) {
        val = max(val, __shfl_down_sync(0xffffffff, val, offset));
    }
    return val;
}

__global__ __launch_bounds__(1024, 1)
void softmax_online_pass_kernel(const float* __restrict__ input, float*
    __restrict__ output, int batch_size, int dim) {
    const int batch_idx = blockIdx.x;
    if (batch_idx >= batch_size) return;

    const float* row_input = input + batch_idx * dim;
    float* row_output = output + batch_idx * dim;

    const int tid = threadIdx.x;
    const int lane_id = tid % WARP_SIZE;
    const int warp_id = tid / WARP_SIZE;
    const int warps_per_block = blockDim.x / WARP_SIZE;

    extern __shared__ float sdata[];

    float thread_max = -FLT_MAX;
    float thread_sum = 0.0f;

    const int N_vec = dim / 4;
    const float4* input_vec = reinterpret_cast<const float4*>(row_input);

    const int loop_stride = blockDim.x * 2;
    int i = tid;
    for (; i + blockDim.x < N_vec; i += loop_stride) {
        const float4 val0 = input_vec[i];
        const float4 val1 = input_vec[i + blockDim.x];
```

```

float local_max0 = max(max(val0.x, val0.y), max(val0.z, val0.w));
float local_max1 = max(max(val1.x, val1.y), max(val1.z, val1.w));
float combined_max = max(local_max0, local_max1);

if (combined_max > thread_max) {
    thread_sum *= __expf(thread_max - combined_max);
    thread_max = combined_max;
}

thread_sum += __expf(val0.x - thread_max) + __expf(val0.y - thread_max) +
__expf(val0.z - thread_max) + __expf(val0.w - thread_max);
thread_sum += __expf(val1.x - thread_max) + __expf(val1.y - thread_max) +
__expf(val1.z - thread_max) + __expf(val1.w - thread_max);
}
for (; i < N_vec; i += blockDim.x) {
    const float4 val = input_vec[i];
    float local_max = max(max(val.x, val.y), max(val.z, val.w));
    if (local_max > thread_max) { thread_sum *= __expf(thread_max - local_max);
        thread_max = local_max; }
    thread_sum += __expf(val.x - thread_max) + __expf(val.y - thread_max) +
__expf(val.z - thread_max) + __expf(val.w - thread_max);
}
for (int j = N_vec * 4 + tid; j < dim; j += blockDim.x) {
    float val = row_input[j];
    if (val > thread_max) { thread_sum *= __expf(thread_max - val); thread_max
        = val; }
    thread_sum += __expf(val - thread_max);
}

float warp_max = warp_reduce_max(thread_max);
if (lane_id == 0) sdata[warp_id] = warp_max;
__syncthreads();

float block_max = (tid < warps_per_block) ? sdata[tid] : -FLT_MAX;
if (warp_id == 0) block_max = warp_reduce_max(block_max);
if (tid == 0) sdata[0] = block_max;
__syncthreads();
block_max = sdata[0];

thread_sum *= __expf(thread_max - block_max);
float warp_sum = warp_reduce_sum(thread_sum);
if (lane_id == 0) sdata[warp_id] = warp_sum;
__syncthreads();

float block_sum = (tid < warps_per_block) ? sdata[tid] : 0.0f;
if (warp_id == 0) block_sum = warp_reduce_sum(block_sum);
if (tid == 0) sdata[0] = block_sum;
__syncthreads();
block_sum = sdata[0] + 1e-12f;

float4* output_vec = reinterpret_cast<float4*>(row_output);
i = tid;
for (; i + blockDim.x < N_vec; i += loop_stride) {
    const float4 val0 = input_vec[i];
    const float4 val1 = input_vec[i + blockDim.x];
    float4 out_val0, out_val1;

    out_val0.x = __expf(val0.x - block_max) / block_sum; out_val0.y =
__expf(val0.y - block_max) / block_sum;
    out_val0.z = __expf(val0.z - block_max) / block_sum; out_val0.w =
__expf(val0.w - block_max) / block_sum;
    output_vec[i] = out_val0;

```

```

        out_val1.x = __expf(val1.x - block_max) / block_sum; out_val1.y =
            __expf(val1.y - block_max) / block_sum;
        out_val1.z = __expf(val1.z - block_max) / block_sum; out_val1.w =
            __expf(val1.w - block_max) / block_sum;
        output_vec[i + blockDim.x] = out_val1;
    }
    for (; i < N_vec; i += blockDim.x) { // Remainder loop
        const float4 val = input_vec[i];
        float4 out_val;
        out_val.x = __expf(val.x - block_max) / block_sum; out_val.y = __expf(val.y
            - block_max) / block_sum;
        out_val.z = __expf(val.z - block_max) / block_sum; out_val.w = __expf(val.w
            - block_max) / block_sum;
        output_vec[i] = out_val;
    }
    for (int j = N_vec * 4 + tid; j < dim; j += blockDim.x) {
        row_output[j] = __expf(row_input[j] - block_max) / block_sum;
    }
}

torch::Tensor softmax_cuda(torch::Tensor input) {
    const auto batch_size = input.size(0);
    const auto dim = input.size(1);
    auto output = torch::empty_like(input);
    const int threads_per_block = 1024;
    const int warps_per_block = threads_per_block / WARP_SIZE;
    const int blocks_per_grid = batch_size;
    const size_t shared_mem_size = warps_per_block * sizeof(float);
    softmax_online_pass_kernel<<<blocks_per_grid, threads_per_block,
        shared_mem_size>>>(
        input.data_ptr<float>(), output.data_ptr<float>(), batch_size, dim);

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) {
        TORCH_CHECK(false, "CUDA kernel launch failed: ", cudaGetErrorString(err));
    }

    return output;
}
"""

softmax_cpp_source = """
#include <torch/extension.h>
torch::Tensor softmax_cuda(torch::Tensor input);
"""

softmax_module = load_inline(
    name='softmax_cuda_bf16_pass_fixed',
    cpp_sources=softmax_cpp_source,
    cuda_sources=softmax_cuda_source,
    functions=['softmax_cuda'],
    verbose=True,
    extra_cuda_cflags=['-O3', '--use_fast_math']
)

class ModelNew(nn.Module):
    def __init__(self, num_features: int = -1):
        super(ModelNew, self).__init__()
        self.softmax_custom_cuda = softmax_module.softmax_cuda

    def forward(self, x: torch.Tensor) -> torch.Tensor:

```

```
return self.softmax_custom_cuda(x)
```

D.11. VGG16

```

import torch
import torch.nn as nn
import torch.nn.functional as F
from torch.utils.cpp_extension import load_inline
import os
import math
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

cublas_fused_source = """
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>

__global__ void add_bias_relu_kernel_float4(
    const __half* __restrict__ input,
    const __half* __restrict__ bias,
    __half* __restrict__ output,
    const int M,
    const int N)
{
    const int N_vec = N / 8;
    const int total_vecs = M * N_vec;

    for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < total_vecs; i +=
        blockDim.x * gridDim.x) {
        const int col_vec_idx = i % N_vec;
        const float4 input_f4 = reinterpret_cast<const float4*>(input)[i];
        const float4 bias_f4 = reinterpret_cast<const float4*>(bias)[col_vec_idx];

        const __half2* input_h2_vecs = reinterpret_cast<const __half2*>(&input_f4);
        const __half2* bias_h2_vecs = reinterpret_cast<const __half2*>(&bias_f4);

        float4 result_f4;
        __half2* result_h2_vecs = reinterpret_cast<__half2*>(&result_f4);

        const __half zero_h = __float2half(0.0f);
        const __half2 zero_vec = __halves2half2(zero_h, zero_h);

        #pragma unroll
        for (int j = 0; j < 4; ++j) {
            result_h2_vecs[j] = __hmax2(__hadd2(input_h2_vecs[j], bias_h2_vecs[j]),
                zero_vec);
        }
        reinterpret_cast<float4*>(output)[i] = result_f4;
    }
}

__global__ void add_bias_kernel_float4(
    const __half* __restrict__ input,
    const __half* __restrict__ bias,
    __half* __restrict__ output,
    const int M,
    const int N)
{
    const int N_vec = N / 8;
    const int total_vecs = M * N_vec;

    for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < total_vecs; i +=
        blockDim.x * gridDim.x) {

```

```

        const int col_vec_idx = i % N_vec;
        const float4 input_f4 = reinterpret_cast<const float4*>(input)[i];
        const float4 bias_f4 = reinterpret_cast<const float4*>(bias)[col_vec_idx];

        const __half2* input_h2_vecs = reinterpret_cast<const __half2*>(&input_f4);
        const __half2* bias_h2_vecs = reinterpret_cast<const __half2*>(&bias_f4);

        float4 result_f4;
        __half2* result_h2_vecs = reinterpret_cast<__half2*>(&result_f4);

        #pragma unroll
        for (int j = 0; j < 4; ++j) {
            result_h2_vecs[j] = __hadd2(input_h2_vecs[j], bias_h2_vecs[j]);
        }
        reinterpret_cast<float4*>(output)[i] = result_f4;
    }
}

__global__ void add_bias_relu_4d_nhwc_kernel(
    const __half* __restrict__ input,
    const __half* __restrict__ bias,
    __half* __restrict__ output,
    const int C,
    const int total_vectors) // N*H*W*C / 8. Changed from long long to int
{
    const __half zero_h = __float2half(0.0f);
    const __half2 zero_vec = __halves2half2(zero_h, zero_h);
    const int C_vec = C / 8;

    for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < total_vectors; i +=
        blockDim.x * gridDim.x) { // loop variable changed to int
        const int bias_vec_idx = i % C_vec;

        const float4 input_f4 = reinterpret_cast<const float4*>(input)[i];
        const float4 bias_f4 = reinterpret_cast<const float4*>(bias)[bias_vec_idx];

        const __half2* input_h2_vecs = reinterpret_cast<const __half2*>(&input_f4);
        const __half2* bias_h2_vecs = reinterpret_cast<const __half2*>(&bias_f4);

        float4 result_f4;
        __half2* result_h2_vecs = reinterpret_cast<__half2*>(&result_f4);

        #pragma unroll
        for (int j = 0; j < 4; ++j) {
            result_h2_vecs[j] = __hmax2(__hadd2(input_h2_vecs[j], bias_h2_vecs[j]),
                zero_vec);
        }

        reinterpret_cast<float4*>(output)[i] = result_f4;
    }
}

__global__ void add_bias_relu_maxpool_2x2_s2_nhwc_kernel_float4(
    const __half* __restrict__ input,
    const __half* __restrict__ bias,
    __half* __restrict__ output,
    const int C_vec,
    const int W_out,
    const int H_out,
    const int W_in_C_vec,
    const int H_in_W_in_C_vec,
    const int total_output_vectors
) {

```

```

const __half zero_h = __float2half(0.0f);
const __half2 zero_vec = __halves2half2(zero_h, zero_h);

for (int i = blockIdx.x * blockDim.x + threadIdx.x; i < total_output_vectors; i
    += blockDim.x * gridDim.x) {
    const int c_vec_idx = i % C_vec;
    const int spatial_out_idx = i / C_vec;
    const int w_out_idx = spatial_out_idx % W_out;
    const int nh_out_idx = spatial_out_idx / W_out;
    const int h_out_idx = nh_out_idx % H_out;
    const int n_idx = nh_out_idx / H_out;
    const float4 bias_f4 = reinterpret_cast<const float4*>(bias)[c_vec_idx];
    const __half2* bias_h2_vecs = reinterpret_cast<const __half2*>(&bias_f4);
    const int h_in_base = h_out_idx * 2;
    const int w_in_base = w_out_idx * 2;
    const int base_idx_00 = n_idx * H_in_W_in_C_vec + h_in_base * W_in_C_vec +
        w_in_base * C_vec + c_vec_idx;
    const int base_idx_01 = base_idx_00 + C_vec;
    const int base_idx_10 = base_idx_00 + W_in_C_vec;
    const int base_idx_11 = base_idx_10 + C_vec;
    const float4 in_f4_00 = reinterpret_cast<const float4*>(input)[base_idx_00];
    const float4 in_f4_01 = reinterpret_cast<const float4*>(input)[base_idx_01];
    const float4 in_f4_10 = reinterpret_cast<const float4*>(input)[base_idx_10];
    const float4 in_f4_11 = reinterpret_cast<const float4*>(input)[base_idx_11];
    const __half2* in_h2_vecs_00 = reinterpret_cast<const __half2*>(&in_f4_00);
    const __half2* in_h2_vecs_01 = reinterpret_cast<const __half2*>(&in_f4_01);
    const __half2* in_h2_vecs_10 = reinterpret_cast<const __half2*>(&in_f4_10);
    const __half2* in_h2_vecs_11 = reinterpret_cast<const __half2*>(&in_f4_11);
    float4 result_f4;
    __half2* result_h2_vecs = reinterpret_cast<__half2*>(&result_f4);

    #pragma unroll
    for (int j = 0; j < 4; ++j) {
        __half2 v00 = __hmax2(__hadd2(in_h2_vecs_00[j], bias_h2_vecs[j]),
            zero_vec);
        __half2 v01 = __hmax2(__hadd2(in_h2_vecs_01[j], bias_h2_vecs[j]),
            zero_vec);
        __half2 v10 = __hmax2(__hadd2(in_h2_vecs_10[j], bias_h2_vecs[j]),
            zero_vec);
        __half2 v11 = __hmax2(__hadd2(in_h2_vecs_11[j], bias_h2_vecs[j]),
            zero_vec);
        result_h2_vecs[j] = __hmax2(__hmax2(v00, v01), __hmax2(v10, v11));
    }

    reinterpret_cast<float4*>(output)[i] = result_f4;
}

}

torch::Tensor add_bias_fused_dispatch(torch::Tensor input, torch::Tensor bias, bool
    apply_relu) {
    const int M = input.size(0);
    const int N = input.size(1);
    auto output = torch::empty_like(input);
    if (M * N == 0) return output;

    const int block_size = 256;
    const int num_blocks = (M * N / 8 + block_size - 1) / block_size;

    if (apply_relu) {
        add_bias_relu_kernel_float4<<<num_blocks, block_size>>>((const
            __half*)input.data_ptr<at::Half>(), (const
            __half*)bias.data_ptr<at::Half>(), (__half*)output.data_ptr<at::Half>(),
            M, N);
    }
}

```



```

    } else {
        add_bias_kernel_float4<<<num_blocks, block_size>>>((const
            __half*)input.data_ptr<at::Half>(), (const
            __half*)bias.data_ptr<at::Half>(), (__half*)output.data_ptr<at::Half>(),
            M, N);
    }
    return output;
}

torch::Tensor add_bias_relu_forward_float4(torch::Tensor input, torch::Tensor bias)
{ return add_bias_fused_dispatch(input, bias, true); }
torch::Tensor add_bias_forward_float4(torch::Tensor input, torch::Tensor bias) {
    return add_bias_fused_dispatch(input, bias, false); }

torch::Tensor add_bias_relu_4d_nhwc_forward(torch::Tensor input, torch::Tensor
    bias) {
    const int C = input.size(1);
    auto output = torch::empty_like(input, torch::MemoryFormat::ChannelsLast);
    if (input.numel() == 0) return output;

    const int total_vectors = input.numel() / 8; // Changed from long long
    const int block_size = 256;
    const int num_blocks = (total_vectors + block_size - 1) / block_size;

    add_bias_relu_4d_nhwc_kernel<<<num_blocks, block_size>>>(
        (const __half*)input.data_ptr<at::Half>(),
        (const __half*)bias.data_ptr<at::Half>(),
        (__half*)output.data_ptr<at::Half>(),
        C,
        total_vectors);
    return output;
}

void launch_add_bias_relu_maxpool_4d_nhwc(torch::Tensor input, torch::Tensor bias,
    torch::Tensor output) {
    const int N = input.size(0);
    const int C = input.size(1);
    const int H_in = input.size(2);
    const int W_in = input.size(3);
    const int C_vec = C / 8;
    const int H_out = H_in / 2;
    const int W_out = W_in / 2;
    const int W_in_C_vec = W_in * C_vec;
    const int H_in_W_in_C_vec = H_in * W_in_C_vec;
    const int total_output_vectors = N * H_out * W_out * C_vec;
    const int block_size = 256;
    const int num_blocks = (total_output_vectors + block_size - 1) / block_size;
    add_bias_relu_maxpool_2x2_s2_nhwc_kernel_float4<<<num_blocks, block_size>>>(
        (const __half*)input.data_ptr<at::Half>(),
        (const __half*)bias.data_ptr<at::Half>(),
        (__half*)output.data_ptr<at::Half>(),
        C_vec, W_out, H_out, W_in_C_vec, H_in_W_in_C_vec,
        total_output_vectors);
}

torch::Tensor forward_add_bias_relu_maxpool_4d_nhwc(torch::Tensor input,
    torch::Tensor bias) {
    const int N = input.size(0);
    const int C = input.size(1);
    const int H_in = input.size(2);
    const int W_in = input.size(3);
    const int H_out = H_in / 2;
    const int W_out = W_in / 2;

```

```

    auto output = torch::empty({N, C, H_out, W_out}, input.options(),
        torch::MemoryFormat::ChannelsLast);
    if (input.numel() == 0) return output;
    launch_add_bias_relu_maxpool_4d_nhwc(input, bias, output);
    return output;
}

torch::Tensor forward_add_bias_relu_maxpool_flatten_4d_nhwc(torch::Tensor input,
    torch::Tensor bias) {
    const int N = input.size(0);
    const int C = input.size(1);
    const int H_in = input.size(2);
    const int W_in = input.size(3);
    const int H_out = H_in / 2;
    const int W_out = W_in / 2;
    const int flattened_dim = C * H_out * W_out; // Changed from long long
    auto output = torch::empty({N, flattened_dim}, input.options());
    if (input.numel() == 0) return output;
    launch_add_bias_relu_maxpool_4d_nhwc(input, bias, output);
    return output;
}

PYBIND11_MODULE(TORCH_EXTENSION_NAME, m) {
    m.def("forward_relu", &add_bias_relu_forward_float4, "Fused Add Bias and ReLU
        forward (CUDA)");
    m.def("forward_bias_add", &add_bias_forward_float4, "Fused Add Bias forward
        (CUDA)");
    m.def("forward_add_bias_relu_4d_nhwc", &add_bias_relu_4d_nhwc_forward, "Fused Add
        Bias and ReLU for 4D NHWC Tensors (CUDA)");
    m.def("forward_add_bias_relu_maxpool_4d_nhwc",
        &forward_add_bias_relu_maxpool_4d_nhwc, "Fused Add Bias, ReLU, and MaxPool for
        4D NHWC Tensors (CUDA)");
    m.def("forward_add_bias_relu_maxpool_flatten_4d_nhwc",
        &forward_add_bias_relu_maxpool_flatten_4d_nhwc, "Fused Add Bias, ReLU,
        MaxPool, and Flatten for 4D NHWC Tensors (CUDA)");
}
"""

cublas_fused_module = load_inline(
    name='cublas_fused_ops_vgg_int_indexing',
    cpp_sources=[],
    cuda_sources=[cublas_fused_source],
    verbose=True,
    extra_cflags=['-O3'],
    extra_cuda_cflags=['-O3', '--use_fast_math']
)

class FusedLinearLayer(nn.Module):
    def __init__(self, in_features, out_features, has_relu):
        super().__init__()
        self.in_features = in_features
        self.out_features = out_features
        self.has_relu = has_relu
        self.weight = nn.Parameter(torch.Tensor(out_features, in_features))
        self.bias = nn.Parameter(torch.Tensor(out_features))
        self.reset_parameters()

    def reset_parameters(self):
        nn.init.kaiming_uniform_(self.weight, a=math.sqrt(5))
        fan_in, _ = nn.init._calculate_fan_in_and_fan_out(self.weight)
        if fan_in > 0:
            bound = 1 / math.sqrt(fan_in)

```

```

        nn.init.uniform_(self.bias, -bound, bound)

    def forward(self, input):
        matmul_result = F.linear(input, self.weight)
        if self.has_relu:
            return cublas_fused_module.forward_relu(matmul_result, self.bias)
        else:
            return cublas_fused_module.forward_bias_add(matmul_result, self.bias)

class FusedConv2dReLU_NHWC(nn.Module):
    def __init__(self, in_channels, out_channels, kernel_size, padding):
        super().__init__()
        self.weight = nn.Parameter(torch.Tensor(out_channels, in_channels,
            kernel_size, kernel_size))
        self.bias = nn.Parameter(torch.Tensor(out_channels))
        self.padding = padding
        self.reset_parameters()

    def reset_parameters(self):
        nn.init.kaiming_uniform_(self.weight, a=math.sqrt(5))
        fan_in, _ = nn.init._calculate_fan_in_and_fan_out(self.weight)
        if fan_in > 0:
            bound = 1 / math.sqrt(fan_in)
            nn.init.uniform_(self.bias, -bound, bound)

    def forward(self, input):
        conv_output = F.conv2d(input, self.weight, bias=None, stride=1,
            padding=self.padding)
        return cublas_fused_module.forward_add_bias_relu_4d_nhwc(conv_output,
            self.bias)

class FusedConv2dReLUMaxPool_NHWC(nn.Module):
    def __init__(self, in_channels, out_channels, kernel_size, padding):
        super().__init__()
        self.weight = nn.Parameter(torch.Tensor(out_channels, in_channels,
            kernel_size, kernel_size))
        self.bias = nn.Parameter(torch.Tensor(out_channels))
        self.padding = padding
        self.reset_parameters()

    def reset_parameters(self):
        nn.init.kaiming_uniform_(self.weight, a=math.sqrt(5))
        fan_in, _ = nn.init._calculate_fan_in_and_fan_out(self.weight)
        if fan_in > 0:
            bound = 1 / math.sqrt(fan_in)
            nn.init.uniform_(self.bias, -bound, bound)

    def forward(self, input):
        conv_output = F.conv2d(input, self.weight, bias=None, stride=1,
            padding=self.padding)
        return cublas_fused_module.forward_add_bias_relu_maxpool_4d_nhwc(conv_output,
            self.bias)

class FusedConv2dReLUMaxPoolFlatten_NHWC(nn.Module):
    def __init__(self, in_channels, out_channels, kernel_size, padding):
        super().__init__()
        self.weight = nn.Parameter(torch.Tensor(out_channels, in_channels,
            kernel_size, kernel_size))
        self.bias = nn.Parameter(torch.Tensor(out_channels))
        self.padding = padding
        self.reset_parameters()

```

```

def reset_parameters(self):
    nn.init.kaiming_uniform_(self.weight, a=math.sqrt(5))
    fan_in, _ = nn.init._calculate_fan_in_and_fan_out(self.weight)
    if fan_in > 0:
        bound = 1 / math.sqrt(fan_in)
        nn.init.uniform_(self.bias, -bound, bound)

def forward(self, input):
    conv_output = F.conv2d(input, self.weight, bias=None, stride=1,
        padding=self.padding)
    return
        cublas_fused_module.forward_add_bias_relu_maxpool_flatten_4d_nhwc(conv_output,
            self.bias)

class ModelNew(nn.Module):
    def __init__(self, num_features: int):
        super(ModelNew, self).__init__()

        classifier_in_features = 512 * 7 * 7
        classifier_hidden_features = 4096

        self.features = nn.Sequential(
            FusedConv2dReLU_NHWC(3, 64, kernel_size=3, padding=1),
            FusedConv2dReLU_MaxPool_NHWC(64, 64, kernel_size=3, padding=1),

            FusedConv2dReLU_NHWC(64, 128, kernel_size=3, padding=1),
            FusedConv2dReLU_MaxPool_NHWC(128, 128, kernel_size=3, padding=1),

            FusedConv2dReLU_NHWC(128, 256, kernel_size=3, padding=1),
            FusedConv2dReLU_NHWC(256, 256, kernel_size=3, padding=1),
            FusedConv2dReLU_MaxPool_NHWC(256, 256, kernel_size=3, padding=1),

            FusedConv2dReLU_NHWC(256, 512, kernel_size=3, padding=1),
            FusedConv2dReLU_NHWC(512, 512, kernel_size=3, padding=1),
            FusedConv2dReLU_MaxPool_NHWC(512, 512, kernel_size=3, padding=1),

            FusedConv2dReLU_NHWC(512, 512, kernel_size=3, padding=1),
            FusedConv2dReLU_NHWC(512, 512, kernel_size=3, padding=1),
            FusedConv2dReLU_MaxPoolFlatten_NHWC(512, 512, kernel_size=3, padding=1),
        )

        self.classifier = nn.Sequential(
            FusedLinearLayer(classifier_in_features, classifier_hidden_features,
                has_relu=True),
            nn.Dropout(p=0.0),
            FusedLinearLayer(classifier_hidden_features,
                classifier_hidden_features, has_relu=True),
            nn.Dropout(p=0.0),
            FusedLinearLayer(classifier_hidden_features, num_features,
                has_relu=False)
        )

        self.half().cuda().to(memory_format=torch.channels_last)
        self.eval()

        self.batch_size = 128
        self.static_input = torch.randn(self.batch_size, 3, 224, 224,
            device='cuda', dtype=torch.half)
        self.static_input = self.static_input.to(memory_format=torch.channels_last)
        self.graph = torch.cuda.CUDAGraph()
        self.static_output = None
        self._capture_graph()

```

```
def _forward_internal(self, x: torch.Tensor) -> torch.Tensor:
    x = self.features(x)
    x = self.classifier(x)
    return x

def _capture_graph(self):
    s = torch.cuda.Stream()
    s.wait_stream(torch.cuda.current_stream())
    with torch.cuda.stream(s):
        for _ in range(3):
            self._forward_internal(self.static_input)
    torch.cuda.current_stream().wait_stream(s)

    with torch.cuda.graph(self.graph):
        self.static_output = self._forward_internal(self.static_input)

def forward(self, x: torch.Tensor) -> torch.Tensor:
    if x.shape[0] != self.batch_size:
        x_nhwc = x.to(dtype=torch.half, device='cuda',
                      memory_format=torch.channels_last)
        return self._forward_internal(x_nhwc).float().contiguous()

    self.static_input.copy_(x.to(memory_format=torch.channels_last))
    self.graph.replay()
    return self.static_output.clone().float().contiguous()
```

D.12. Mean Reduction over a dimension

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

cuda_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <limits>
#include <c10/cuda/CUDAException.h> // Header for C10_CUDA_CHECK
#include <algorithm>

__device__ __forceinline__ float halfWarpReduceSumXOR(float val) {
    const unsigned int mask = 0xffffffff;
    val += __shfl_xor_sync(mask, val, 8);
    val += __shfl_xor_sync(mask, val, 4);
    val += __shfl_xor_sync(mask, val, 2);
    val += __shfl_xor_sync(mask, val, 1);
    return val;
}

template <bool IsContiguous>
__global__ __launch_bounds__(256, 4) void grid_stride_reduction_kernel(
    const float* __restrict__ input, float* __restrict__ output,
    const int reduce_dim_size, const int outer_size, const int inner_size, const
        int total_outputs,
    const float inv_reduce_dim_size) {

    if constexpr (IsContiguous) {
        const unsigned int threads_per_group = 16;
        const unsigned int group_id = threadIdx.x / threads_per_group;
        const unsigned int local_lane_id = threadIdx.x % threads_per_group;
        const unsigned int groups_per_block = blockDim.x / threads_per_group;

        for (int output_idx = blockIdx.x * groups_per_block + group_id;
            output_idx < total_outputs;
            output_idx += gridDim.x * groups_per_block) {

            const float* input_ptr = input + output_idx * reduce_dim_size;

            const int vec_size = 2;
            const int reduce_dim_vec = reduce_dim_size / vec_size;

            float s0 = 0.0f, s1 = 0.0f, s2 = 0.0f, s3 = 0.0f;
            const int unroll_factor = 4;
            const int loop_stride = threads_per_group * unroll_factor;

            int i = local_lane_id;
            const int unrolled_limit = reduce_dim_vec - (reduce_dim_vec %
                unroll_factor);

            for (; i < unrolled_limit; i += loop_stride) {
                float f1, f2;
                const float2* ptr_a = &reinterpret_cast<const
                    float2*>(input_ptr)[i];
                asm("ld.global.cs.v2.f32 {%0, %1}, [%2];" : "=f"(f1), "=f"(f2) :
                    "l"(ptr_a));
                s0 += f1 + f2;
            }
        }
    }
}

```

```

        const float2* ptr_b = &reinterpret_cast<const float2*>(input_ptr)[i
            + threads_per_group];
        asm("ld.global.cs.v2.f32 {%0, %1}, [%2];" : "=f"(f1), "=f"(f2) :
            "l"(ptr_b));
        s1 += f1 + f2;
        const float2* ptr_c = &reinterpret_cast<const float2*>(input_ptr)[i
            + 2 * threads_per_group];
        asm("ld.global.cs.v2.f32 {%0, %1}, [%2];" : "=f"(f1), "=f"(f2) :
            "l"(ptr_c));
        s2 += f1 + f2;
        const float2* ptr_d = &reinterpret_cast<const float2*>(input_ptr)[i
            + 3 * threads_per_group];
        asm("ld.global.cs.v2.f32 {%0, %1}, [%2];" : "=f"(f1), "=f"(f2) :
            "l"(ptr_d));
        s3 += f1 + f2;
    }

    for (; i < reduce_dim_vec; i += threads_per_group) {
        float f1, f2;
        const float2* ptr = &reinterpret_cast<const float2*>(input_ptr)[i];
        asm("ld.global.cs.v2.f32 {%0, %1}, [%2];" : "=f"(f1), "=f"(f2) :
            "l"(ptr));
        s0 += f1 + f2;
    }

    float thread_sum = s0 + s1 + s2 + s3;
    const int remainder_start = reduce_dim_vec * vec_size;
    for (int j = remainder_start + local_lane_id; j < reduce_dim_size; j +=
        threads_per_group) {
        thread_sum += input_ptr[j]; // Remainder too small to benefit from
        PTX
    }

    float group_sum = halfWarpReduceSumXOR(thread_sum);
    if (local_lane_id == 0) {
        output[output_idx] = group_sum * inv_reduce_dim_size;
    }
}

} else {
    constexpr int ROWS = 8;
    constexpr int COLS = 32;
    extern __shared__ float tile[];

    const int tx = threadIdx.x % COLS;
    const int ty = threadIdx.x / COLS;

    for (int base_output_idx = blockIdx.x * COLS;
        base_output_idx < total_outputs;
        base_output_idx += gridDim.x * COLS) {

        const unsigned int output_idx = base_output_idx + tx;
        if (output_idx >= total_outputs) continue;

        const int outer_idx = output_idx / inner_size;
        const int inner_idx = output_idx % inner_size;
        const float* input_ptr = input + outer_idx * reduce_dim_size *
            inner_size + inner_idx;

        float thread_sum = 0.0f;
        int i = ty;
        const int unrolled_limit = reduce_dim_size - (reduce_dim_size % 4);
        for (; i < unrolled_limit; i += ROWS * 4) {

```



```

        float temp;
        const float* ptr_a = &input_ptr[i * inner_size];
        asm("ld.global.cg.f32 %0, [%1];" : "=f"(temp) : "l"(ptr_a));
        thread_sum += temp;
        const float* ptr_b = &input_ptr[(i + ROWS) * inner_size];
        asm("ld.global.cg.f32 %0, [%1];" : "=f"(temp) : "l"(ptr_b));
        thread_sum += temp;
        const float* ptr_c = &input_ptr[(i + 2*ROWS) * inner_size];
        asm("ld.global.cg.f32 %0, [%1];" : "=f"(temp) : "l"(ptr_c));
        thread_sum += temp;
        const float* ptr_d = &input_ptr[(i + 3*ROWS) * inner_size];
        asm("ld.global.cg.f32 %0, [%1];" : "=f"(temp) : "l"(ptr_d));
        thread_sum += temp;
    }
    for (; i < reduce_dim_size; i += ROWS) {
        float temp;
        const float* ptr = &input_ptr[i * inner_size];
        asm("ld.global.cg.f32 %0, [%1];" : "=f"(temp) : "l"(ptr));
        thread_sum += temp;
    }

    tile[threadIdx.x] = thread_sum;
    __syncthreads();

    if (ty == 0) {
        float final_sum = tile[tx + 0 * COLS] + tile[tx + 1 * COLS] +
                           tile[tx + 2 * COLS] + tile[tx + 3 * COLS] +
                           tile[tx + 4 * COLS] + tile[tx + 5 * COLS] +
                           tile[tx + 6 * COLS] + tile[tx + 7 * COLS];
        output[output_idx] = final_sum * inv_reduce_dim_size;
    }
    __syncthreads();
}
}

torch::Tensor mean_reduction_cuda(torch::Tensor input, int dim) {
    auto input_sizes = input.sizes();
    int ndim = input_sizes.size();
    dim = (dim < 0) ? (dim + ndim) : dim;

    const int reduce_dim_size = input_sizes[dim];
    std::vector<int64_t> output_sizes;
    for (int i = 0; i < ndim; ++i) {
        if (i != dim) {
            output_sizes.push_back(input_sizes[i]);
        }
    }

    auto output = torch::empty(output_sizes, input.options());
    if (reduce_dim_size == 0) {
        output.fill_(std::numeric_limits<float>::quiet_NaN());
        return output;
    }
    if (reduce_dim_size == 1) {
        output.copy_(input.squeeze(dim));
        return output;
    }

    int outer_size = 1;
    for (int i = 0; i < dim; ++i) outer_size *= input_sizes[i];
    int inner_size = 1;
    for (int i = dim + 1; i < ndim; ++i) inner_size *= input_sizes[i];

```

```

const int total_outputs = outer_size * inner_size;
if (total_outputs == 0) return output;

const float inv_reduce_dim_size = 1.0f / static_cast<float>(reduce_dim_size);

if (inner_size == 1) {
    const int block_dim_x = 256;
    const int threads_per_group = 16;
    const int groups_per_block = block_dim_x / threads_per_group;
    const int grid_dim_x = std::min((total_outputs + groups_per_block - 1) /
        groups_per_block, 512);

    grid_stride_reduction_kernel<true><<<grid_dim_x, block_dim_x>>>(
        input.data_ptr<float>(), output.data_ptr<float>(),
        reduce_dim_size, outer_size, inner_size, total_outputs,
        inv_reduce_dim_size
    );
} else {
    const int block_dim_x = 256; // 8 rows x 32 cols
    const int outputs_per_block = 32;
    const int grid_dim_x = std::min((total_outputs + outputs_per_block - 1) /
        outputs_per_block, 512);
    const int shared_mem_size = block_dim_x * sizeof(float);

    grid_stride_reduction_kernel<false><<<grid_dim_x, block_dim_x,
        shared_mem_size>>>(
        input.data_ptr<float>(), output.data_ptr<float>(),
        reduce_dim_size, outer_size, inner_size, total_outputs,
        inv_reduce_dim_size
    );
}

C10_CUDA_CHECK(cudaGetLastError());
return output;
}
"""

cpp_source = """
torch::Tensor mean_reduction_cuda(torch::Tensor input, int dim);
"""

mean_reduction_module = load_inline(
    name='ptx_caching_reduction',
    cpp_sources=cpp_source,
    cuda_sources=cuda_source,
    functions=['mean_reduction_cuda'],
    verbose=True
)

class ModelNew(nn.Module):
    def __init__(self, num_features: int):
        super(ModelNew, self).__init__()
        self.dim = num_features
        self.mean_reduction = mean_reduction_module

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        return self.mean_reduction.mean_reduction_cuda(x.cuda(), self.dim)

```

D.13. RNN

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

try:
    batch_size
except NameError:
    batch_size = 64

rnn_cuda_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>

#define I2H_TILE_H 128
#define H2O_TILE_O 128
#define THREADS_PER_BLOCK 1024

__device__ __forceinline__ float fast_tanhf(float x) {
    float r;
    asm("tanh.approx.f32 %0, %1;" : "=f"(r) : "f"(x));
    return r;
}

__global__ void __launch_bounds__(THREADS_PER_BLOCK, 1) i2h_feature_parallel_kernel(
    const float* __restrict__ input,
    const float* __restrict__ hidden,
    const float* __restrict__ i2h_weight,
    const float* __restrict__ i2h_bias,
    float* __restrict__ new_hidden_out,
    int batch_size,
    int input_size,
    int hidden_size
) {
    const int batch_idx = blockIdx.y;
    const int h_tile_start = blockIdx.x * I2H_TILE_H;
    extern __shared__ float s_mem[];
    const int combined_size = input_size + hidden_size;
    float* s_combined_input = s_mem;
    const float* current_input = input + batch_idx * input_size;
    const float* current_hidden = hidden + batch_idx * hidden_size;

    for (int i = threadIdx.x; i < input_size / 2; i += blockDim.x) {
        reinterpret_cast<float2*>(s_combined_input)[i] = reinterpret_cast<const
            float2*>(current_input)[i];
    }
    if ((input_size % 2 != 0) && (threadIdx.x == 0)) {
        s_combined_input[input_size - 1] = current_input[input_size - 1];
    }

    for (int i = threadIdx.x; i < hidden_size / 2; i += blockDim.x) {
        reinterpret_cast<float2*>(s_combined_input + input_size)[i] =
            reinterpret_cast<const float2*>(current_hidden)[i];
    }
    if ((hidden_size % 2 != 0) && (threadIdx.x == 0)) {
        s_combined_input[input_size + hidden_size - 1] = current_hidden[hidden_size
            - 1];
    }
    __syncthreads();
}

```

```

const int warp_id = threadIdx.x / 32;
const int lane_id = threadIdx.x % 32;
const int h_idx = h_tile_start + warp_id * 4;

if (h_idx >= hidden_size) return;
float sum1 = 0.0f, sum2 = 0.0f, sum3 = 0.0f, sum4 = 0.0f;
const bool process2 = (h_idx + 1 < hidden_size);
const bool process3 = (h_idx + 2 < hidden_size);
const bool process4 = (h_idx + 3 < hidden_size);
const float2* s_combined_input_v2 = reinterpret_cast<const
    float2*>(s_combined_input);
const int combined_size_v2 = combined_size / 2;
const float2* weight_row1_v2 = reinterpret_cast<const float2*>(i2h_weight +
    h_idx * combined_size);
const float2* weight_row2_v2 = process2 ? reinterpret_cast<const
    float2*>(i2h_weight + (h_idx + 1) * combined_size) : nullptr;
const float2* weight_row3_v2 = process3 ? reinterpret_cast<const
    float2*>(i2h_weight + (h_idx + 2) * combined_size) : nullptr;
const float2* weight_row4_v2 = process4 ? reinterpret_cast<const
    float2*>(i2h_weight + (h_idx + 3) * combined_size) : nullptr;

for (int i = lane_id; i < combined_size_v2; i += 32) {
    const float2 s_val = s_combined_input_v2[i];
    sum1 += s_val.x * weight_row1_v2[i].x + s_val.y * weight_row1_v2[i].y;
    if (process2) sum2 += s_val.x * weight_row2_v2[i].x + s_val.y *
        weight_row2_v2[i].y;
    if (process3) sum3 += s_val.x * weight_row3_v2[i].x + s_val.y *
        weight_row3_v2[i].y;
    if (process4) sum4 += s_val.x * weight_row4_v2[i].x + s_val.y *
        weight_row4_v2[i].y;
}

if (combined_size % 2 != 0) {
    if (lane_id == 0) {
        const float s_val_last = s_combined_input[combined_size - 1];
        sum1 += s_val_last * i2h_weight[h_idx * combined_size + combined_size -
            1];
        if (process2) sum2 += s_val_last * i2h_weight[(h_idx + 1) *
            combined_size + combined_size - 1];
        if (process3) sum3 += s_val_last * i2h_weight[(h_idx + 2) *
            combined_size + combined_size - 1];
        if (process4) sum4 += s_val_last * i2h_weight[(h_idx + 3) *
            combined_size + combined_size - 1];
    }
}

#pragma unroll
for (int offset = 16; offset > 0; offset >= 1) {
    sum1 += __shfl_down_sync(0xffffffff, sum1, offset);
    sum2 += __shfl_down_sync(0xffffffff, sum2, offset);
    sum3 += __shfl_down_sync(0xffffffff, sum3, offset);
    sum4 += __shfl_down_sync(0xffffffff, sum4, offset);
}

if (lane_id == 0) {
    new_hidden_out[batch_idx * hidden_size + h_idx] = fast_tanhf(sum1 +
        i2h_bias[h_idx]);
    if (process2) new_hidden_out[batch_idx * hidden_size + h_idx + 1] =
        fast_tanhf(sum2 + i2h_bias[h_idx + 1]);
    if (process3) new_hidden_out[batch_idx * hidden_size + h_idx + 2] =
        fast_tanhf(sum3 + i2h_bias[h_idx + 2]);
}

```

```

        if (process4) new_hidden_out[batch_idx * hidden_size + h_idx + 3] =
            fast_tanhf(sum4 + i2h_bias[h_idx + 3]);
    }
}

__global__ void __launch_bounds__(THREADS_PER_BLOCK, 1) h2o_feature_parallel_kernel(
    const float* __restrict__ new_hidden,
    const float* __restrict__ h2o_weight,
    const float* __restrict__ h2o_bias,
    float* __restrict__ output,
    int batch_size,
    int hidden_size,
    int output_size
) {
    const int batch_idx = blockIdx.y;
    const int o_tile_start = blockIdx.x * H2O_TILE_O;

    extern __shared__ float s_mem[];
    float* s_new_hidden = s_mem;

    const float* current_hidden = new_hidden + batch_idx * hidden_size;
    for (int i = threadIdx.x; i < hidden_size / 2; i += blockDim.x) {
        reinterpret_cast<float2*>(s_new_hidden)[i] = reinterpret_cast<const
            float2*>(current_hidden)[i];
    }
    if ((hidden_size % 2 != 0) && (threadIdx.x == 0)) {
        s_new_hidden[hidden_size - 1] = current_hidden[hidden_size - 1];
    }
    __syncthreads();

    const int warp_id = threadIdx.x / 32;
    const int lane_id = threadIdx.x % 32;
    const int o_idx = o_tile_start + warp_id * 8;

    if (o_idx >= output_size) return;

    float sum1=0, sum2=0, sum3=0, sum4=0, sum5=0, sum6=0, sum7=0, sum8=0;
    const bool p2 = (o_idx + 1 < output_size);
    const bool p3 = (o_idx + 2 < output_size);
    const bool p4 = (o_idx + 3 < output_size);
    const bool p5 = (o_idx + 4 < output_size);
    const bool p6 = (o_idx + 5 < output_size);
    const bool p7 = (o_idx + 6 < output_size);
    const bool p8 = (o_idx + 7 < output_size);
    const float2* s_new_hidden_v2 = reinterpret_cast<const float2*>(s_new_hidden);
    const int hidden_size_v2 = hidden_size / 2;
    const float2* w1 = reinterpret_cast<const float2*>(h2o_weight + o_idx *
        hidden_size);
    const float2* w2 = p2 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        1) * hidden_size) : nullptr;
    const float2* w3 = p3 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        2) * hidden_size) : nullptr;
    const float2* w4 = p4 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        3) * hidden_size) : nullptr;
    const float2* w5 = p5 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        4) * hidden_size) : nullptr;
    const float2* w6 = p6 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        5) * hidden_size) : nullptr;
    const float2* w7 = p7 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        6) * hidden_size) : nullptr;
    const float2* w8 = p8 ? reinterpret_cast<const float2*>(h2o_weight + (o_idx +
        7) * hidden_size) : nullptr;

```

```

for (int i = lane_id; i < hidden_size_v2; i += 32) {
    const float2 s_val = s_new_hidden_v2[i];
    sum1 += s_val.x * w1[i].x + s_val.y * w1[i].y;
    if (p2) sum2 += s_val.x * w2[i].x + s_val.y * w2[i].y;
    if (p3) sum3 += s_val.x * w3[i].x + s_val.y * w3[i].y;
    if (p4) sum4 += s_val.x * w4[i].x + s_val.y * w4[i].y;
    if (p5) sum5 += s_val.x * w5[i].x + s_val.y * w5[i].y;
    if (p6) sum6 += s_val.x * w6[i].x + s_val.y * w6[i].y;
    if (p7) sum7 += s_val.x * w7[i].x + s_val.y * w7[i].y;
    if (p8) sum8 += s_val.x * w8[i].x + s_val.y * w8[i].y;
}

if (hidden_size % 2 != 0) {
    if (lane_id == 0) {
        const float s_val_last = s_new_hidden[hidden_size - 1];
        sum1 += s_val_last * h2o_weight[o_idx * hidden_size + hidden_size - 1];
        if (p2) sum2 += s_val_last * h2o_weight[(o_idx + 1) * hidden_size +
            hidden_size - 1];
        if (p3) sum3 += s_val_last * h2o_weight[(o_idx + 2) * hidden_size +
            hidden_size - 1];
        if (p4) sum4 += s_val_last * h2o_weight[(o_idx + 3) * hidden_size +
            hidden_size - 1];
        if (p5) sum5 += s_val_last * h2o_weight[(o_idx + 4) * hidden_size +
            hidden_size - 1];
        if (p6) sum6 += s_val_last * h2o_weight[(o_idx + 5) * hidden_size +
            hidden_size - 1];
        if (p7) sum7 += s_val_last * h2o_weight[(o_idx + 6) * hidden_size +
            hidden_size - 1];
        if (p8) sum8 += s_val_last * h2o_weight[(o_idx + 7) * hidden_size +
            hidden_size - 1];
    }
}

#pragma unroll
for (int offset = 16; offset > 0; offset >= 1) {
    sum1 += __shfl_down_sync(0xffffffff, sum1, offset);
    sum2 += __shfl_down_sync(0xffffffff, sum2, offset);
    sum3 += __shfl_down_sync(0xffffffff, sum3, offset);
    sum4 += __shfl_down_sync(0xffffffff, sum4, offset);
    sum5 += __shfl_down_sync(0xffffffff, sum5, offset);
    sum6 += __shfl_down_sync(0xffffffff, sum6, offset);
    sum7 += __shfl_down_sync(0xffffffff, sum7, offset);
    sum8 += __shfl_down_sync(0xffffffff, sum8, offset);
}

if (lane_id == 0) {
    output[batch_idx * output_size + o_idx] = sum1 + h2o_bias[o_idx];
    if (p2) output[batch_idx * output_size + o_idx + 1] = sum2 + h2o_bias[o_idx
        + 1];
    if (p3) output[batch_idx * output_size + o_idx + 2] = sum3 + h2o_bias[o_idx
        + 2];
    if (p4) output[batch_idx * output_size + o_idx + 3] = sum4 + h2o_bias[o_idx
        + 3];
    if (p5) output[batch_idx * output_size + o_idx + 4] = sum5 + h2o_bias[o_idx
        + 4];
    if (p6) output[batch_idx * output_size + o_idx + 5] = sum6 + h2o_bias[o_idx
        + 5];
    if (p7) output[batch_idx * output_size + o_idx + 6] = sum7 + h2o_bias[o_idx
        + 6];
    if (p8) output[batch_idx * output_size + o_idx + 7] = sum8 + h2o_bias[o_idx
        + 7];
}

```

```

}

std::vector<torch::Tensor> rnn_forward_feature_parallel_cuda(
    torch::Tensor input,
    torch::Tensor hidden,
    torch::Tensor i2h_weight,
    torch::Tensor i2h_bias,
    torch::Tensor h2o_weight,
    torch::Tensor h2o_bias
) {
    const int batch_size = input.size(0);
    const int input_size = input.size(1);
    const int hidden_size = hidden.size(1);
    const int output_size = h2o_weight.size(0);

    auto new_hidden = torch::empty({batch_size, hidden_size}, input.options());
    auto output = torch::empty({batch_size, output_size}, input.options());

    const dim3 threads(THREADS_PER_BLOCK);

    const dim3 i2h_blocks((hidden_size + I2H_TILE_H - 1) / I2H_TILE_H, batch_size);
    const size_t i2h_shared_mem_size = (input_size + hidden_size) * sizeof(float);
    i2h_feature_parallel_kernel<<<i2h_blocks, threads, i2h_shared_mem_size>>>(
        input.data_ptr<float>(),
        hidden.data_ptr<float>(),
        i2h_weight.data_ptr<float>(),
        i2h_bias.data_ptr<float>(),
        new_hidden.data_ptr<float>(),
        batch_size,
        input_size,
        hidden_size
    );

    const dim3 h2o_blocks((output_size + H2O_TILE_O - 1) / H2O_TILE_O, batch_size);
    const size_t h2o_shared_mem_size = hidden_size * sizeof(float);
    h2o_feature_parallel_kernel<<<h2o_blocks, threads, h2o_shared_mem_size>>>(
        new_hidden.data_ptr<float>(),
        h2o_weight.data_ptr<float>(),
        h2o_bias.data_ptr<float>(),
        output.data_ptr<float>(),
        batch_size,
        hidden_size,
        output_size
    );

    cudaError_t err = cudaGetLastError();
    if (err != cudaSuccess) {
        throw std::runtime_error(cudaGetErrorString(err));
    }

    return {output, new_hidden};
}

"""
rnn_cpp_source = """
#include <torch/extension.h>
#include <vector>

std::vector<torch::Tensor> rnn_forward_feature_parallel_cuda(
    torch::Tensor input,
    torch::Tensor hidden,
    torch::Tensor i2h_weight,
    torch::Tensor i2h_bias,

```

```

        torch::Tensor h2o_weight,
        torch::Tensor h2o_bias
    );
    """

# JIT compilation of the CUDA kernel
rnn_cuda_feature_parallel = load_inline(
    name='rnn_cuda_feature_parallel',
    cpp_sources=rnn_cpp_source,
    cuda_sources=rnn_cuda_source,
    functions=['rnn_forward_feature_parallel_cuda'],
    verbose=True,
    extra_cuda_cflags=['-O3']
)

class ModelNew(nn.Module):
    def __init__(self, input_size: int, hidden_size: int, output_size: int):
        super(ModelNew, self).__init__()
        self.input_size = input_size
        self.hidden_size = hidden_size
        self.output_size = output_size
        self.hidden = torch.randn((batch_size, self.hidden_size))
        self.i2h = nn.Linear(self.input_size + self.hidden_size, self.hidden_size)
        self.h2o = nn.Linear(self.hidden_size, self.output_size)
        self.rnn_cuda = rnn_cuda_feature_parallel

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        if self.hidden.device != x.device:
            self.hidden = self.hidden.to(x.device)

        if self.hidden.shape[0] != x.shape[0]:
            self.hidden = torch.randn((x.shape[0], self.hidden_size),
                                       device=x.device, dtype=x.dtype)

        output, self.hidden = self.rnn_cuda.rnn_forward_feature_parallel_cuda(
            x,
            self.hidden,
            self.i2h.weight,
            self.i2h.bias,
            self.h2o.weight,
            self.h2o.bias
        )
        return output

```


D.14. BMM InstanceNorm Sum ResidualAdd Multiply

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline
import os
import math
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

fused_linear_norm_source = """
#include <torch/extension.h>
#include <cuda.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>

__device__ inline float rsqrt_fast(float x) {
    return rsqrtf(x);
}

__device__ inline float warp_reduce_sum(float val) {
    #pragma unroll
    for (int offset = 16; offset > 0; offset >= 1) {
        val += __shfl_down_sync(0xFFFFFFFF, val, offset);
    }
    return val;
}

__global__ void __launch_bounds__(128, 4) fused_linear_norm_add_mul_kernel(
    const float* __restrict__ x,          // Input tensor (Batch, InFeatures)
    const float* __restrict__ y,          // Second input tensor (Batch,
        OutFeatures)
    const float* __restrict__ linear_w,    // Linear layer weights (OutFeatures,
        InFeatures)
    const float* __restrict__ linear_b,    // Linear layer bias (OutFeatures)
    const float* __restrict__ norm_w,     // InstanceNorm weights (OutFeatures)
    const float* __restrict__ norm_b,     // InstanceNorm bias (OutFeatures)
    float* __restrict__ out,              // Output tensor (Batch, OutFeatures)
    const int batch_size,
    const int in_features,
    const int out_features,
    const float eps
) {
    const int batch_idx = blockIdx.x;
    if (batch_idx >= batch_size) {
        return;
    }

    extern __shared__ float s_mem[];
    float* s_linear_out = s_mem; // Size: out_features
    float* s_reduce_mem = &s_mem[out_features]; // Size for reduction: num_warps * 2

    for (int c_out_idx = threadIdx.x; c_out_idx < out_features; c_out_idx +=
        blockDim.x) {
        float acc[4] = {0.0f, 0.0f, 0.0f, 0.0f};

        const int in_features_vec = in_features / 4;
        const float4* x_vec = reinterpret_cast<const float4*>(&x[batch_idx *
            in_features]);
        const float4* linear_w_vec = reinterpret_cast<const
            float4*>(&linear_w[c_out_idx * in_features]);

        const int unroll_factor = 20;

```

```

int k = 0;
if (in_features_vec >= unroll_factor) {
    for (; k <= in_features_vec - unroll_factor; k += unroll_factor) {
        #pragma unroll
        for(int i=0; i < unroll_factor; ++i) {
            float4 x_v = x_vec[k + i];
            float4 w_v = linear_w_vec[k + i];
            // Cycle through 4 accumulators. The compiler resolves 'i % 4'
            // at compile time.
            acc[i % 4] += x_v.x * w_v.x + x_v.y * w_v.y + x_v.z * w_v.z +
                x_v.w * w_v.w;
        }
    }
}

float remainder_acc = 0.0f;
for (; k < in_features_vec; ++k) {
    float4 x_v = x_vec[k];
    float4 w_v = linear_w_vec[k];
    remainder_acc += x_v.x * w_v.x + x_v.y * w_v.y + x_v.z * w_v.z + x_v.w
        * w_v.w;
}

float total_acc = remainder_acc + acc[0] + acc[1] + acc[2] + acc[3];

s_linear_out[c_out_idx] = total_acc + linear_b[c_out_idx];
}
__syncthreads();

const int warp_id = threadIdx.x / 32;
const int lane_id = threadIdx.x % 32;
float local_sum = 0.0f;
float local_sum_sq = 0.0f;
const int out_features_vec = out_features / 4;
const float4* s_linear_out_vec = reinterpret_cast<const float4*>(s_linear_out);

for (int i = threadIdx.x; i < out_features_vec; i += blockDim.x) {
    float4 val4 = s_linear_out_vec[i];
    local_sum += val4.x + val4.y + val4.z + val4.w;
    local_sum_sq += val4.x * val4.x + val4.y * val4.y + val4.z * val4.z +
        val4.w * val4.w;
}

float warp_sum = warp_reduce_sum(local_sum);
float warp_sum_sq = warp_reduce_sum(local_sum_sq);

if (lane_id == 0) {
    s_reduce_mem[warp_id] = warp_sum;
    s_reduce_mem[warp_id + blockDim.x / 32] = warp_sum_sq;
}
__syncthreads();

if (warp_id == 0) {
    local_sum = (lane_id < blockDim.x / 32) ? s_reduce_mem[lane_id] : 0.0f;
    local_sum_sq = (lane_id < blockDim.x / 32) ? s_reduce_mem[lane_id +
        blockDim.x / 32] : 0.0f;
    warp_sum = warp_reduce_sum(local_sum);
    warp_sum_sq = warp_reduce_sum(local_sum_sq);
}

if (threadIdx.x == 0) {
    float mean = warp_sum / out_features;
    float var = (warp_sum_sq / out_features) - (mean * mean);
}

```

```

        s_reduce_mem[0] = mean;
        s_reduce_mem[1] = rsqrt_fast(var + eps);
    }
    __syncthreads();

    const float mean = s_reduce_mem[0];
    const float inv_std = s_reduce_mem[1];

    for (int i = threadIdx.x; i < out_features_vec; i += blockDim.x) {
        const int offset = i * 4;
        const float4 val4 = *reinterpret_cast<const float4*>(&s_linear_out[offset]);
        const float4 norm_w4 = *reinterpret_cast<const float4*>(&norm_w[offset]);
        const float4 norm_b4 = *reinterpret_cast<const float4*>(&norm_b[offset]);
        const float4 y4 = *reinterpret_cast<const float4*>(&y[batch_idx *
            out_features + offset]);
        const float4 mean4 = {mean, mean, mean, mean};
        const float4 inv_std4 = {inv_std, inv_std, inv_std, inv_std};
        const float4 normalized_val4 = {
            (val4.x - mean4.x) * inv_std4.x, (val4.y - mean4.y) * inv_std4.y,
            (val4.z - mean4.z) * inv_std4.z, (val4.w - mean4.w) * inv_std4.w
        };
        const float4 transformed_val4 = {
            normalized_val4.x * norm_w4.x + norm_b4.x, normalized_val4.y *
            norm_w4.y + norm_b4.y,
            normalized_val4.z * norm_w4.z + norm_b4.z, normalized_val4.w *
            norm_w4.w + norm_b4.w
        };
        const float4 temp_val4 = {
            transformed_val4.x + y4.x, transformed_val4.y + y4.y,
            transformed_val4.z + y4.z, transformed_val4.w + y4.w
        };
        const float4 out4 = {
            temp_val4.x * y4.x, temp_val4.y * y4.y,
            temp_val4.z * y4.z, temp_val4.w * y4.w
        };
        *reinterpret_cast<float4*>(&out[batch_idx * out_features + offset]) = out4;
    }
}

torch::Tensor fused_linear_norm_add_mul(
    torch::Tensor x,
    torch::Tensor y,
    torch::Tensor linear_w,
    torch::Tensor linear_b,
    torch::Tensor norm_w,
    torch::Tensor norm_b,
    double eps
) {
    const int batch_size = x.size(0);
    const int in_features = x.size(1);
    const int out_features = y.size(1);
    auto out = torch::empty_like(y);

    if (batch_size == 0) {
        return out;
    }

    int threads_per_block = 128;
    const int num_blocks = batch_size;
    const int num_warps = threads_per_block / 32;
    const int shared_mem_size = (out_features + num_warps * 2) * sizeof(float);

```

```

        fused_linear_norm_add_mul_kernel<<<num_blocks, threads_per_block,
            shared_mem_size>>>(
            x.data_ptr<float>(),
            y.data_ptr<float>(),
            linear_w.data_ptr<float>(),
            linear_b.data_ptr<float>(),
            norm_w.data_ptr<float>(),
            norm_b.data_ptr<float>(),
            out.data_ptr<float>(),
            batch_size,
            in_features,
            out_features,
            static_cast<float>(eps)
        );

        cudaError_t err = cudaGetLastError();
        if (err != cudaSuccess) {
            AT_ERROR("CUDA kernel launch failed in fused_linear_norm_add_mul: ",
                cudaGetErrorString(err));
        }

        return out;
    }
    """

    fused_linear_norm_cpp_source = """
    torch::Tensor fused_linear_norm_add_mul(
        torch::Tensor x,
        torch::Tensor y,
        torch::Tensor linear_w,
        torch::Tensor linear_b,
        torch::Tensor norm_w,
        torch::Tensor norm_b,
        double eps
    );
    """

    fused_op_reg_pressure_opt = load_inline(
        name='fused_op_reg_pressure_opt',
        cpp_sources=fused_linear_norm_cpp_source,
        cuda_sources=fused_linear_norm_source,
        functions=['fused_linear_norm_add_mul'],
        verbose=True,
        extra_cuda_cflags=['-O3', '--use_fast_math']
    )

    class ModelNew(nn.Module):
        def __init__(self, in_features: int, out_features: int, eps=1e-5, momentum=0.1):
            super(ModelNew, self).__init__()
            self.in_features = in_features
            self.out_features = out_features

            if out_features % 4 != 0:
                raise ValueError(f"out_features must be divisible by 4 for this
                    optimized kernel, but got {out_features}")

            self.linear = nn.Linear(in_features, out_features)
            self.instance_norm = nn.InstanceNorm1d(out_features, eps=eps,
                momentum=momentum, affine=True)
            self.eps = eps
            self.fused_op = fused_op_reg_pressure_opt

        def forward(self, x: torch.Tensor, y: torch.Tensor) -> torch.Tensor:

```

```
if x.dim() != 2 or x.size(1) != self.in_features:
    raise ValueError(f"Input x has wrong dimensions. Expected (batch,
                      {self.in_features}), got {x.shape}")
if y.dim() != 2 or y.size(1) != self.out_features:
    raise ValueError(f"Input y has wrong dimensions. Expected (batch,
                      {self.out_features}), got {y.shape}")

return self.fused_op.fused_linear_norm_add_mul(
    x,
    y,
    self.linear.weight,
    self.linear.bias,
    self.instance_norm.weight,
    self.instance_norm.bias,
    self.eps
)
```

D.15. AlexNet

```

import torch
import torch.nn as nn
import torch.nn.functional as F
import math
from torch.utils.cpp_extension import load_inline
import os
os.environ['TORCH_CUDA_ARCH_LIST'] = '8.0'

fp16_fused_ops_cuda_source = """
#include <torch/extension.h>
#include <cublas_v2.h>
#include <cuda_runtime.h>
#include <cuda_fp16.h>
#include <stdexcept>
#include <iostream>

#define CUDA_CHECK(call)                                \
do {                                                    \
    cudaError_t e = call;                               \
    if (e != cudaSuccess) {                             \
        throw std::runtime_error(cudaGetErrorString(e)); \
    }                                                    \
} while (0)

const char* cublasGetErrorString(cublasStatus_t status) {
    switch(status) {
        case CUBLAS_STATUS_SUCCESS: return "CUBLAS_STATUS_SUCCESS";
        case CUBLAS_STATUS_NOT_INITIALIZED: return "CUBLAS_STATUS_NOT_INITIALIZED";
        case CUBLAS_STATUS_ALLOC_FAILED: return "CUBLAS_STATUS_ALLOC_FAILED";
        case CUBLAS_STATUS_INVALID_VALUE: return "CUBLAS_STATUS_INVALID_VALUE";
        case CUBLAS_STATUS_ARCH_MISMATCH: return "CUBLAS_STATUS_ARCH_MISMATCH";
        case CUBLAS_STATUS_MAPPING_ERROR: return "CUBLAS_STATUS_MAPPING_ERROR";
        case CUBLAS_STATUS_EXECUTION_FAILED: return
            "CUBLAS_STATUS_EXECUTION_FAILED";
        case CUBLAS_STATUS_INTERNAL_ERROR: return "CUBLAS_STATUS_INTERNAL_ERROR";
        case CUBLAS_STATUS_NOT_SUPPORTED: return "CUBLAS_STATUS_NOT_SUPPORTED";
        case CUBLAS_STATUS_LICENSE_ERROR: return "CUBLAS_STATUS_LICENSE_ERROR";
    }
    return "Unknown cuBLAS error";
}

#define CUBLAS_CHECK(call)                                \
do {                                                    \
    cublasStatus_t s = call;                             \
    if (s != CUBLAS_STATUS_SUCCESS) {                     \
        throw std::runtime_error(cublasGetErrorString(s)); \
    }                                                    \
} while (0)

__global__ void fused_relu_maxpool2d_nhwc_fp16_kernel(
    const __half* __restrict__ input,
    __half* __restrict__ output,
    int N, int H_in, int W_in, int C,
    int H_out, int W_out,
    int kernel_size, int stride)
{
    const int total_output_elements = N * H_out * W_out * C;
    const int idx = blockIdx.x * blockDim.x + threadIdx.x;
    const int grid_stride = blockDim.x * gridDim.x;

```

```

for (int i = idx; i < total_output_elements; i += grid_stride) {
    int c = i % C;
    int w_out = (i / C) % W_out;
    int h_out = (i / (C * W_out)) % H_out;
    int n = i / (C * W_out * H_out);
    int h_start = h_out * stride;
    int w_start = w_out * stride;
    half max_val = __float2half(-65504.0f); // Smallest fp16 value

    for (int kh = 0; kh < kernel_size; ++kh) {
        for (int kw = 0; kw < kernel_size; ++kw) {
            int h_in_idx = h_start + kh;
            int w_in_idx = w_start + kw;

            if (h_in_idx < H_in && w_in_idx < W_in) {
                long long input_idx = (long long)n * H_in * W_in * C +
                                      (long long)h_in_idx * W_in * C +
                                      (long long)w_in_idx * C + c;
                half val = input[input_idx];
                val = __hmax(val, __float2half(0.0f));
                max_val = __hmax(max_val, val);
            }
        }
    }
    output[i] = max_val;
}

__device__ unsigned int xorshift32_dev(unsigned int& state) {
    state ^= state << 13;
    state ^= state >> 17;
    state ^= state << 5;
    return state;
}

__device__ float random_uniform_dev(unsigned int& state) {
    return float(xorshift32_dev(state)) / 4294967295.0f;
}

__global__ void __launch_bounds__(512) add_bias_relu_dropout_kernel_fp16(
    const __half* __restrict__ gemm_output,
    const __half* __restrict__ bias,
    __half* __restrict__ final_output,
    int batch_size,
    int out_features,
    float p,
    unsigned long long seed)
{
    const int total_half2_elements = (batch_size * out_features) / 2;
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    const int stride = blockDim.x * gridDim.x;
    unsigned int prng_state = idx + (unsigned int)seed;
    if (prng_state == 0) prng_state = 1;
    const float inv_p = (p > 0.0f) ? (1.0f / (1.0f - p)) : 0.0f;
    const half2 zero_h2 = __float2half2_rn(0.0f);
    const int UNROLL_FACTOR = 4;
    const int unrolled_stride = stride * UNROLL_FACTOR;

    for (; idx + stride * (UNROLL_FACTOR - 1) < total_half2_elements; idx +=
        unrolled_stride) {
        #pragma unroll
        for (int i=0; i < UNROLL_FACTOR; ++i) {

```

```

        int current_idx = idx + i * stride;
        const int col_base = (current_idx * 2) % out_features;
        half2 gemm_val = ((const half2*)gemm_output)[current_idx];
        half2 bias_val = ((const half2*)bias)[col_base / 2];
        half2 biased_val = __hadd2(gemm_val, bias_val);
        half2 relu_val = __hmax2(biased_val, zero_h2);
        if (p > 0.0f) {
            float v1 = __half2float(relu_val.x) *
                (float)(random_uniform_dev(prng_state) > p) * inv_p;
            float v2 = __half2float(relu_val.y) *
                (float)(random_uniform_dev(prng_state) > p) * inv_p;
            relu_val = __floats2half2_rn(v1, v2);
        }
        ((half2*)final_output)[current_idx] = relu_val;
    }
}

for (; idx < total_half2_elements; idx += stride) {
    const int col_base = (idx * 2) % out_features;
    half2 gemm_val = ((const half2*)gemm_output)[idx];
    half2 bias_val = ((const half2*)bias)[col_base / 2];
    half2 biased_val = __hadd2(gemm_val, bias_val);
    half2 relu_val = __hmax2(biased_val, zero_h2);
    if (p > 0.0f) {
        float v1 = __half2float(relu_val.x) *
            (float)(random_uniform_dev(prng_state) > p) * inv_p;
        float v2 = __half2float(relu_val.y) *
            (float)(random_uniform_dev(prng_state) > p) * inv_p;
        relu_val = __floats2half2_rn(v1, v2);
    }
    ((half2*)final_output)[idx] = relu_val;
}
}

__global__ void __launch_bounds__(512) add_bias_and_cast_to_fp32_kernel_fp16(
    const __half* __restrict__ gemm_output,
    const __half* __restrict__ bias,
    float* __restrict__ final_output, // Output is FP32
    int total_elements,
    int out_features)
{
    const int total_half2_elements = total_elements / 2;
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    const int stride = blockDim.x * gridDim.x;
    const int UNROLL_FACTOR = 4;
    const int unrolled_stride = stride * UNROLL_FACTOR;

    for (; idx + stride * (UNROLL_FACTOR - 1) < total_half2_elements; idx +=
        unrolled_stride) {
        #pragma unroll
        for(int i=0; i<UNROLL_FACTOR; ++i) {
            int current_idx = idx + i * stride;
            half2 gemm_val = ((const half2*)gemm_output)[current_idx];
            half2 bias_val = ((const half2*)bias)[((current_idx * 2) %
                out_features) / 2];
            float2 out_f32 = __half22float2(__hadd2(gemm_val, bias_val));
            final_output[current_idx*2] = out_f32.x;
            final_output[current_idx*2 + 1] = out_f32.y;
        }
    }

    for (; idx < total_half2_elements; idx += stride) {

```



```

        const int col_base = (idx * 2) % out_features;
        half2 gemm_val = ((const half2*)gemm_output)[idx];
        half2 bias_val = ((const half2*)bias)[col_base / 2];
        half2 biased_val = __hadd2(gemm_val, bias_val);
        float2 out_f32 = __half22float2(biased_val);
        final_output[idx*2] = out_f32.x;
        final_output[idx*2 + 1] = out_f32.y;
    }
}

static cublasHandle_t cublas_handle;
static bool cublas_handle_initialized = false;

void initialize_cublas() {
    if (!cublas_handle_initialized) {
        CUBLAS_CHECK(cublasCreate(&cublas_handle));
        CUBLAS_CHECK(cublasSetMathMode(cublas_handle, CUBLAS_DEFAULT_MATH));
        cublas_handle_initialized = true;
    }
}

void gemm_fp16(cublasHandle_t handle, torch::Tensor A, torch::Tensor B,
               torch::Tensor C) {
    const int m = B.size(0);
    const int k = B.size(1);
    const int n = A.size(1);

    const float alpha = 1.0f;
    const float beta = 0.0f;

    CUBLAS_CHECK(cublasGemmEx(handle,
                              CUBLAS_OP_N, CUBLAS_OP_N,
                              n, m, k,
                              &alpha,
                              A.data_ptr(), CUDA_R_16F, n,
                              B.data_ptr(), CUDA_R_16F, k,
                              &beta,
                              C.data_ptr(), CUDA_R_16F, n,
                              CUBLAS_COMPUTE_32F,
                              CUBLAS_GEMM_DEFAULT_TENSOR_OP));
}

torch::Tensor cublas_fused_linear_relu_dropout_fp16(
    torch::Tensor input, torch::Tensor weight_t, torch::Tensor bias, float p,
    unsigned long long seed)
{
    const int batch_size = input.size(0);
    const int out_features = bias.size(0);
    initialize_cublas();
    auto gemm_output = torch::empty({batch_size, out_features}, input.options());
    gemm_fp16(cublas_handle, weight_t, input, gemm_output);
    auto final_output = torch::empty({batch_size, out_features}, input.options());
    const int threads_per_block = 512;
    const int total_half2_elements = (batch_size * out_features) / 2;
    const int num_blocks = (total_half2_elements + threads_per_block - 1) /
        threads_per_block;
    add_bias_relu_dropout_kernel_fp16<<<num_blocks, threads_per_block>>>(
        (const __half*)gemm_output.data_ptr(), (const __half*)bias.data_ptr(),
        (__half*)final_output.data_ptr(), batch_size, out_features, p, seed);
    CUDA_CHECK(cudaGetLastError());

    return final_output;
}

```

```

torch::Tensor cublas_fused_flatten_linear_relu_dropout_fp16(
    torch::Tensor input, torch::Tensor weight_t, torch::Tensor bias, float p,
    unsigned long long seed) {
    const int batch_size = input.size(0);
    const int in_features = input.numel() / batch_size;
    auto reshaped_input = input.reshape({batch_size, in_features});
    return cublas_fused_linear_relu_dropout_fp16(reshaped_input, weight_t, bias, p,
        seed);
}

torch::Tensor cublas_fused_linear_fp16_to_fp32(
    torch::Tensor input, torch::Tensor weight_t, torch::Tensor bias)
{
    const int batch_size = input.size(0);
    const int out_features = bias.size(0);
    initialize_cublas();
    auto gemm_output = torch::empty({batch_size, out_features}, input.options());
    gemm_fp16(cublas_handle, weight_t, input, gemm_output);
    auto final_output = torch::empty({batch_size, out_features},
        input.options().dtype(torch::kFloat));
    const int total_elements = batch_size * out_features;
    const int threads_per_block = 512;
    const int total_half2_elements = total_elements / 2;
    const int num_blocks = (total_half2_elements + threads_per_block - 1) /
        threads_per_block;

    add_bias_and_cast_to_fp32_kernel_fp16<<<num_blocks, threads_per_block>>>(
        (const __half*)gemm_output.data_ptr(), (const __half*)bias.data_ptr(),
        (float*)final_output.data_ptr(), total_elements, out_features);
    CUDA_CHECK(cudaGetLastError());

    return final_output;
}

torch::Tensor fused_relu_maxpool2d_fp16(
    torch::Tensor input, int kernel_size, int stride)
{
    CHECK_CUDA(input);
    const int N = input.size(0);
    const int C = input.size(1);
    const int H_in = input.size(2);
    const int W_in = input.size(3);
    const int H_out = (H_in - kernel_size) / stride + 1;
    const int W_out = (W_in - kernel_size) / stride + 1;
    auto output = torch::empty({N, C, H_out, W_out},
        input.options().memory_format(torch::MemoryFormat::ChannelsLast));
    const int total_output_elements = N * H_out * W_out * C;
    if (total_output_elements == 0) return output;
    const int threads_per_block = 512;
    const int num_blocks = (total_output_elements + threads_per_block - 1) /
        threads_per_block;

    fused_relu_maxpool2d_nhwc_fp16_kernel<<<num_blocks, threads_per_block>>>(
        (const __half*)input.data_ptr(),
        (__half*)output.data_ptr(),
        N, H_in, W_in, C, H_out, W_out,
        kernel_size, stride);
    CUDA_CHECK(cudaGetLastError());

    return output;
}
"""

```

```

fp16_fused_ops_cpp_source = """
#include <torch/extension.h>

torch::Tensor cublas_fused_linear_relu_dropout_fp16(torch::Tensor input,
    torch::Tensor weight_t, torch::Tensor bias, float p, unsigned long long seed);
torch::Tensor cublas_fused_flatten_linear_relu_dropout_fp16(torch::Tensor input,
    torch::Tensor weight_t, torch::Tensor bias, float p, unsigned long long seed);
torch::Tensor cublas_fused_linear_fp16_to_fp32(torch::Tensor input, torch::Tensor
    weight_t, torch::Tensor bias);

torch::Tensor fused_relu_maxpool2d_fp16(torch::Tensor input, int kernel_size, int
    stride);
"""

fp16_fused_ops = load_inline(
    name='cublas_fused_ops_fp16_relu_maxpool_fixed',
    cpp_sources=fp16_fused_ops_cpp_source,
    cuda_sources=fp16_fused_ops_cuda_source,
    functions=['cublas_fused_linear_relu_dropout_fp16',
        'cublas_fused_flatten_linear_relu_dropout_fp16',
        'cublas_fused_linear_fp16_to_fp32', 'fused_relu_maxpool2d_fp16'],
    extra_cflags=['-O3'],
    extra_cuda_cflags=['-O3', '--use_fast_math'],
    extra_ldflags=['-lcublas'],
    verbose=False
)

class FusedReLUMaxPool2dFP16(nn.Module):
    def __init__(self, kernel_size, stride):
        super().__init__()
        self.kernel_size = kernel_size
        self.stride = stride

    def forward(self, input):
        return fp16_fused_ops.fused_relu_maxpool2d_fp16(input, self.kernel_size,
            self.stride)

class FusedLinearReLU DropoutFP16(nn.Module):
    def __init__(self, in_features, out_features, p=0.0):
        super().__init__()
        self.in_features = in_features
        self.out_features = out_features
        self.p = p
        self.weight = nn.Parameter(torch.Tensor(in_features, out_features).half())
        self.bias = nn.Parameter(torch.Tensor(out_features).half())
        self.reset_parameters()

    def reset_parameters(self):
        temp_weight = torch.empty(self.out_features, self.in_features)
        nn.init.kaiming_uniform_(temp_weight, a=math.sqrt(5))
        fan_in, _ = nn.init._calculate_fan_in_and_fan_out(temp_weight)
        bound = 1 / math.sqrt(fan_in) if fan_in > 0 else 0
        temp_bias = torch.empty(self.out_features).uniform_(-bound, bound)
        with torch.no_grad():
            self.weight.copy_(temp_weight.T.half())
            self.bias.copy_(temp_bias.half())

    def forward(self, input):
        dropout_p = self.p if self.training else 0.0
        seed = torch.randint(2**63 - 1, (1,)).item()
        return fp16_fused_ops.cublas_fused_linear_relu_dropout_fp16(input,
            self.weight, self.bias, dropout_p, seed)

```

```

class FusedFlattenLinearReLU dropoutFP16(nn.Module):
    def __init__(self, in_channels, h, w, out_features, p=0.0):
        super().__init__()
        self.in_features = in_channels * h * w
        self.out_features = out_features
        self.p = p
        self.weight = nn.Parameter(torch.Tensor(self.in_features,
            out_features).half())
        self.bias = nn.Parameter(torch.Tensor(out_features).half())
        self.reset_parameters()

    def reset_parameters(self):
        temp_weight = torch.empty(self.out_features, self.in_features)
        nn.init.kaiming_uniform_(temp_weight, a=math.sqrt(5))
        fan_in, _ = nn.init._calculate_fan_in_and_fan_out(temp_weight)
        bound = 1 / math.sqrt(fan_in) if fan_in > 0 else 0
        temp_bias = torch.empty(self.out_features).uniform_(-bound, bound)
        with torch.no_grad():
            self.weight.copy_(temp_weight.T.half())
            self.bias.copy_(temp_bias.half())

    def forward(self, input):
        input_contiguous = input.contiguous()
        dropout_p = self.p if self.training else 0.0
        seed = torch.randint(2**63 - 1, (1,)).item()
        return
            fp16_fused_ops.cublas_fused_flatten_linear_relu_dropout_fp16(input_contiguous,
                self.weight, self.bias, dropout_p, seed)

class FusedLinearFP16ToFP32(nn.Module):
    def __init__(self, in_features, out_features):
        super().__init__()
        self.in_features = in_features
        self.out_features = out_features
        self.weight = nn.Parameter(torch.Tensor(in_features, out_features).half())
        self.bias = nn.Parameter(torch.Tensor(out_features).half())
        self.reset_parameters()

    def reset_parameters(self):
        temp_weight = torch.empty(self.out_features, self.in_features)
        nn.init.kaiming_uniform_(temp_weight, a=math.sqrt(5))
        fan_in, _ = nn.init._calculate_fan_in_and_fan_out(temp_weight)
        bound = 1 / math.sqrt(fan_in) if fan_in > 0 else 0
        temp_bias = torch.empty(self.out_features).uniform_(-bound, bound)
        with torch.no_grad():
            self.weight.copy_(temp_weight.T.half())
            self.bias.copy_(temp_bias.half())

    def forward(self, input):
        return fp16_fused_ops.cublas_fused_linear_fp16_to_fp32(input, self.weight,
            self.bias)

class ModelNew(nn.Module):
    def __init__(self, num_features: int = 1000):
        super(ModelNew, self).__init__()
        if num_features % 2 != 0:
            num_features += 1

        num_classes = num_features

        self.features = nn.Sequential(

```

```

        nn.Conv2d(3, 96, kernel_size=11, stride=4, padding=2),
        FusedReLUMaxPool2dFP16(kernel_size=3, stride=2),
        nn.Conv2d(96, 256, kernel_size=5, padding=2),
        FusedReLUMaxPool2dFP16(kernel_size=3, stride=2),
        nn.Conv2d(256, 384, kernel_size=3, padding=1),
        nn.ReLU(inplace=True),
        nn.Conv2d(384, 384, kernel_size=3, padding=1),
        nn.ReLU(inplace=True),
        nn.Conv2d(384, 256, kernel_size=3, padding=1),
        FusedReLUMaxPool2dFP16(kernel_size=3, stride=2),
    )

    self.classifier = nn.Sequential(
        FusedFlattenLinearReLU dropoutFP16(in_channels=256, h=6, w=6,
            out_features=4096, p=0.0),
        FusedLinearReLU dropoutFP16(in_features=4096, out_features=4096, p=0.0),
        FusedLinearFP16ToFP32(in_features=4096, out_features=num_classes),
    )

    self.features.half().to(memory_format=torch.channels_last)

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        x = x.half().to(memory_format=torch.channels_last)
        x = self.features(x)
        x = self.classifier(x)
        return x

```

D.16. LayerNorm

```

import torch
import torch.nn as nn
from torch.utils.cpp_extension import load_inline

layer_norm_source = """
#include <torch/extension.h>
#include <cuda_runtime.h>
#include <cmath>
#include <vector_types.h>

struct Sums {
    float sum;
    float sum_sq;
};

constexpr int BLOCKS_PER_CHANNEL = 64;
constexpr int UNROLL_FACTOR = 16;

__global__ void fused_atomic_stats_kernel(
    const float4* __restrict__ input,
    float* __restrict__ stats,
    int N, int C, int H, int W_div_4) {

    const int total_elements_per_channel_div_4 = N * H * W_div_4;
    const int global_block_idx = blockIdx.x;
    const int feature_idx = global_block_idx / BLOCKS_PER_CHANNEL;
    if (feature_idx >= C) {
        return;
    }

    const int block_in_channel_idx = global_block_idx % BLOCKS_PER_CHANNEL;
    extern __shared__ Sums sdata[];
    float thread_sum = 0.0f;
    float thread_sum_sq = 0.0f;
    const int HW_div_4 = H * W_div_4;
    const int CHW_div_4 = C * HW_div_4;
    const int grid_stride = BLOCKS_PER_CHANNEL * blockDim.x;
    const int unrolled_stride = grid_stride * UNROLL_FACTOR;

    int i = block_in_channel_idx * blockDim.x + threadIdx.x;

    while (i + (UNROLL_FACTOR - 1) * grid_stride <
        total_elements_per_channel_div_4) {
        #pragma unroll
        for (int j = 0; j < UNROLL_FACTOR; ++j) {
            const int current_i = i + j * grid_stride;
            const int n = current_i / HW_div_4;
            const int hw_idx = current_i % HW_div_4;
            const int global_idx = n * CHW_div_4 + feature_idx * HW_div_4 + hw_idx;

            const float4 val4 = input[global_idx];

            thread_sum += val4.x + val4.y + val4.z + val4.w;
            thread_sum_sq += val4.x * val4.x + val4.y * val4.y + val4.z * val4.z +
                val4.w * val4.w;
        }
        i += unrolled_stride;
    }

    for (; i < total_elements_per_channel_div_4; i += grid_stride) {

```

```

    const int n = i / HW_div_4;
    const int hw_idx = i % HW_div_4;
    const int global_idx = n * CHW_div_4 + feature_idx * HW_div_4 + hw_idx;

    const float4 val4 = input[global_idx];

    thread_sum += val4.x + val4.y + val4.z + val4.w;
    thread_sum_sq += val4.x * val4.x + val4.y * val4.y + val4.z * val4.z +
        val4.w * val4.w;
}

sdata[threadIdx.x] = {thread_sum, thread_sum_sq};
__syncthreads();

for (int s = blockDim.x / 2; s > 0; s >= 1) {
    if (threadIdx.x < s) {
        sdata[threadIdx.x].sum += sdata[threadIdx.x + s].sum;
        sdata[threadIdx.x].sum_sq += sdata[threadIdx.x + s].sum_sq;
    }
    __syncthreads();
}

if (threadIdx.x == 0) {
    atomicAdd(&stats[feature_idx * 2], sdata[0].sum);
    atomicAdd(&stats[feature_idx * 2 + 1], sdata[0].sum_sq);
}
}

__global__ void apply_norm_from_stats_kernel(
    const float4* __restrict__ input,
    const float* __restrict__ weight,
    const float* __restrict__ bias,
    const float* __restrict__ stats,
    float4* __restrict__ output,
    int total_elements_div_4,
    int C, int HW,
    float num_elements_per_channel) {

    const int grid_stride = gridDim.x * blockDim.x;
    const int unrolled_stride = grid_stride * UNROLL_FACTOR;
    int idx_div_4 = blockIdx.x * blockDim.x + threadIdx.x;

    while (idx_div_4 + (UNROLL_FACTOR - 1) * grid_stride < total_elements_div_4) {
        #pragma unroll
        for (int j = 0; j < UNROLL_FACTOR; ++j) {
            const int current_idx = idx_div_4 + j * grid_stride;
            const int feature_idx = ((current_idx * 4) / HW) % C;

            const float total_sum = stats[feature_idx * 2];
            const float total_sum_sq = stats[feature_idx * 2 + 1];

            const float m = total_sum / num_elements_per_channel;
            const float var = (total_sum_sq / num_elements_per_channel) - (m * m);
            const float iv = rsqrtf(var + 1e-5f);

            const float w = weight[feature_idx];
            const float b = bias[feature_idx];

            const float4 in4 = input[current_idx];
            float4 out4;
            out4.x = (in4.x - m) * iv * w + b;
            out4.y = (in4.y - m) * iv * w + b;

```

```

        out4.z = (in4.z - m) * iv * w + b;
        out4.w = (in4.w - m) * iv * w + b;
        output[current_idx] = out4;
    }
    idx_div_4 += unrolled_stride;
}

for (; idx_div_4 < total_elements_div_4; idx_div_4 += grid_stride) {
    const int feature_idx = ((idx_div_4 * 4) / HW) % C;

    const float total_sum = stats[feature_idx * 2];
    const float total_sum_sq = stats[feature_idx * 2 + 1];

    const float m = total_sum / num_elements_per_channel;
    const float var = (total_sum_sq / num_elements_per_channel) - (m * m);
    const float iv = rsqrtf(var + 1e-5f);

    const float w = weight[feature_idx];
    const float b = bias[feature_idx];

    const float4 in4 = input[idx_div_4];
    float4 out4;
    out4.x = (in4.x - m) * iv * w + b;
    out4.y = (in4.y - m) * iv * w + b;
    out4.z = (in4.z - m) * iv * w + b;
    out4.w = (in4.w - m) * iv * w + b;
    output[idx_div_4] = out4;
}
}

torch::Tensor layer_norm_cuda(torch::Tensor input, torch::Tensor weight,
    torch::Tensor bias) {
    input = input.contiguous();

    const auto N = input.size(0);
    const auto C = input.size(1);
    const auto H = input.size(2);
    const auto W = input.size(3);
    const auto total_elements = input.numel();
    const auto W_div_4 = W / 4;
    auto output = torch::empty_like(input);
    auto options =
        torch::TensorOptions().device(input.device()).dtype(torch::kFloat32);

    auto stats = torch::zeros({C, 2}, options);
    const int stats_block_size = 256;
    const int stats_num_blocks = C * BLOCKS_PER_CHANNEL;
    const int stats_shared_mem_size = stats_block_size * sizeof(Sums);

    fused_atomic_stats_kernel<<<stats_num_blocks, stats_block_size,
        stats_shared_mem_size>>>(
        reinterpret_cast<const float4*>(input.data_ptr<float>()),
        stats.data_ptr<float>(),
        N, C, H, W_div_4
    );

    const int apply_block_size = 256;
    const int total_elements_div_4 = total_elements / 4;
    const int apply_num_blocks = (total_elements_div_4 + apply_block_size - 1) /
        apply_block_size;
    const float num_elements_per_channel = static_cast<float>(N * H * W);
    apply_norm_from_stats_kernel<<<apply_num_blocks, apply_block_size>>>(

```



```

        reinterpret_cast<const float4*>(input.data_ptr<float>()),
        weight.data_ptr<float>(),
        bias.data_ptr<float>(),
        stats.data_ptr<float>(),
        reinterpret_cast<float4*>(output.data_ptr<float>()),
        total_elements_div_4,
        C, H*W,
        num_elements_per_channel
    );

    auto err = cudaGetLastError();
    TORCH_CHECK(err == cudaSuccess, "CUDA kernel launch error: ",
        cudaGetErrorString(err));

    return output;
}
"""

layer_norm_cpp_source = (
    "torch::Tensor layer_norm_cuda(torch::Tensor input, torch::Tensor weight, "
    "    torch::Tensor bias);"
)

_layer_norm_module = load_inline(
    name="fused_layernorm_2pass_atomic_aos",
    cpp_sources=layer_norm_cpp_source,
    cuda_sources=layer_norm_source,
    functions=["layer_norm_cuda"],
    verbose=False,
    extra_cuda_cflags=["-O3", "--use_fast_math"],
)

class ModelNew(nn.Module):
    def __init__(self, num_features: int):
        super(ModelNew, self).__init__()
        if isinstance(num_features, (list, tuple)):
            channel_count = num_features[0]
        else:
            channel_count = num_features

        self.num_features = channel_count

        self.weight = nn.Parameter(torch.ones(self.num_features))
        self.bias = nn.Parameter(torch.zeros(self.num_features))

        if _layer_norm_module is None:
            raise RuntimeError("CUDA extension for LayerNorm was not compiled "
                               "successfully.")
        self.layer_norm_cuda = _layer_norm_module.layer_norm_cuda

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        if x.dim() != 4:
            raise ValueError(f"Expected 4D input (N, C, H, W), but got {x.dim()}D")

        if x.size(1) != self.num_features:
            raise ValueError(
                f"Expected input to have {self.num_features} features (channels), "
                f"but got {x.size(1)}"
            )

        if x.size(3) % 4 != 0:
            raise ValueError(

```

```
        f"Input width (W) must be divisible by 4 for float4 optimization,\n          but got {x.size(3)}"\n    )\n\n    return self.layer_norm_cuda(x, self.weight, self.bias)
```