TILELANG: BRIDGE PROGRAMMABILITY AND PERFORMANCE IN MODERN NEURAL KERNELS

Anonymous authors

000

001

002 003 004

006

008

010 011

012

013

014

015

016

017

018

019

021

023

025 026

027

029

031

033

034

037

040

041

042

043

044

046

047

048

051

052

Paper under double-blind review

ABSTRACT

Achieving high performance in modern AI increasingly requires kernels codesigned with underlying hardware, but writing efficient kernels remains challenging due to hardware-level complexity and limited fine-grained control in compilers like Triton. In this paper, we introduce TILELANG, a programmable tile-level system that provides explicit primitives for memory placement, data movement, and parallel scheduling. Using a unified fused tile-level dataflow graph (FTG), TILE-LANG streamlines kernel development by unifying tile recommendation, which guides developers with hardware-aware defaults, and tile inference, which automates completion through constraint propagation. TILELANG makes it easy to express a wide range of AI algorithms in under 70 lines of Python code, reducing code size by up to 85.5% compared to manual implementations. Evaluations show that TILELANG achieves up to a 10.59× speedup over Triton on NVIDIA H100 and up to 11.56× on AMD GPUs, bridging programmability and performance.

1 Introduction

The rapid progress of modern neural networks has driven a growing demand for highly optimized compute kernels, particularly for memory-bound operations such as attention. In recent years, modern attention algorithms such as Multi-Head Attention(MHA) (Vaswani et al., 2017), Multi-Head Latent Attention (MLA) (Liu et al., 2024), Gated Query Attention (GQA) (Ainslie et al., 2023), and Linear Attention (Gu & Dao, 2023; Dao & Gu, 2024; Sun et al., 2023; Yang et al., 2024), increasingly demand fine-grained control over memory hierarchy, scheduling, and data movement to fully utilize hardware capabilities. However, existing systems like Triton (Tillet et al., 2019) lack programmable abstractions to support this level of control. For instance, FlashMLA relies on carefully pipelined computations and shared memory reuse, but Triton gives programmers no direct control over tile reuse or pipeline scheduling, restricting performance optimization.

As a result, developers often face a steep tradeoff between achieving peak performance and maintaining programmability: they must either manually write complex CUDA kernels or sacrifice significant performance due to abstraction mismatches. As illustrated in Figure 1, the Triton implementation of MLA requires only 130 lines of code, whose convenience comes at a steep cost—its performance reaches only 14.2% of the hand-written CUDA version (DeepSeek, 2025) (\sim 500 lines) on NVIDIA H100 GPUs. Bridging the gap between programmability and performance requires addressing two key challenges. First, a programming model must give developers precise control over data movement and computation, enabling direct interaction with hardware resources. Second, a compiler must efficiently lower these high-level programs to GPU code, mapping abstractions onto hard-

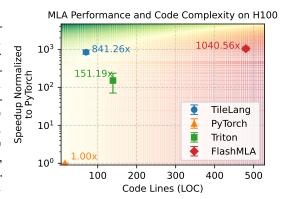


Figure 1: Performance vs. code size trade-off for MLA kernels on NVIDIA H100. Points closer to the top-left indicate better balance between performance and implementation simplicity. The annotated speedup values indicate performance gains over the PyTorch implementation.

ware resources without adding programming complexity. Solving both challenges is essential to balance developer productivity with near-peak hardware performance on modern accelerators.

We introduce TILELANG, a controllable programming system for modern neural workloads. TILELANG provides programmable tile abstractions that let developers express and optimize low-level kernel behaviors in a high-level, composable way. Unlike existing compilers such as Triton, which rely on opaque optimization passes, TILELANG gives developers explicit control over memory, data movement, layout, and parallel execution. Specifically, developers can allocate buffers in different hardware memory levels (alloc_shared, alloc_fragment), orchestrate data transfers (copy), define custom memory layouts (annotate_layout, use_swizzle), and fine-tune parallelism and pipelining strategies (Parallel, Pipelined).

Under programmable tile abstractions, TILELANG programs can be represented as a unified fused tile-level dataflow graph (FTG). By operating on this FTG, TILELANG enables fine-grained reasoning and optimization of AI kernels, guiding developers from high-level design choices to fully specified, hardware-efficient kernel configurations. It introduces two complementary techniques. First, *tile recommendation* analyzes the FTG along with partially specified configurations to provide hardware-aware defaults for tile shapes, memory placement, and warp partitions, offering developers high-quality starting points that can be accepted, adjusted, or further tuned. Second, tile inference propagates shape and layout constraints across the FTG to complete the remaining configurations based on the partially annotated operators. It also automatically aligns buffer shapes, layouts, and memory allocations both downstream and upstream. This cooperative workflow narrows the design space while producing consistent, efficient kernel configurations, reducing manual effort without sacrificing fine-grained control.

As shown in Figure 1, TILELANG achieves on average 5.56× the performance of Triton and approaches the hand-written CUDA version in performance, while requiring less than 16% of the code size of the manual kernel and even fewer LOCs than Triton. This highlights TILELANG's ability to attain a more favorable balance between programmability and performance, offering both high efficiency and low development effort. We also implement other modern AI kernels—including Dequantize Matmul (Wang et al., 2024), Multi-Head Attention (MHA) (Vaswani et al., 2017), and Block-Sparse Attention (BSA) (Guo et al., 2024). Despite its deliberately streamlined interface, TILELANG achieves state-of-the-art throughput across heterogeneous GPUs, delivering speed-ups of up to 10.59× over Triton on an NVIDIA H100 and 11.56× on an AMD MI300X (AMD, 2024).

Our contributions are twofold: (1) programmable tile abstractions that let developers directly control and interact with hardware; and (2) tile recommendation and inference that guide developers with hardware-aware defaults and automatically complete configurations over a unified FTG graph. We believe TILELANG improve both the productivity and performance of modern AI kernel development.

2 Related work

AI kernel programming and optimizations. To simplify the development of AI kernels, libraries like FlashAttention-3 (Shah et al., 2024), CUTLASS (NVIDIA, 2019), and ThunderKittens (Spector et al., 2025) rely on manual or template-driven designs. Triton (Tillet et al., 2019) provides a high-level Python DSL but restricts control over critical performance paths. Cypress (Yadav et al., 2025) introduces a task-based programming model with sequential semantics. Frameworks such as PyTorch (Paszke et al., 2019), Graphene (Hagedorn et al., 2023), MLIR (Lattner et al., 2021), and Welder (Shi et al., 2023) take a compiler-centric approach. Unlike these works, TileLang is a tile-level programmable language that automates layout and low-level configuration while giving users fine-grained control. Its flexible tile programming abstraction can help researchers obtain kernels for a broad range of AI operations, and enable advanced optimizations like software pipelining (Cheng et al., 2025) and warp specialization (Huang et al., 2023).

Cost modeling. TANGRAM (Gao et al., 2019) optimizes dataflow across scheduling layers, along with a performance modeling tool extended by SET (Cai et al., 2023) with Resource Allocation Trees. KPerfIR (Guan et al., 2025) adds instrumentation for profiling and pipeline reordering in Triton. ML-based predictors like Path Forward (Li et al., 2023) and NEUSIGHT (Lee et al., 2025) also exist. In contrast, TILELANG's tile-level analytical cost model uniquely captures both computation and data movement at tile granularity, supporting fusion-aware scheduling with high accuracy and usability.

An extended discussion of related work is in Appendix A, where we discussed several classic works such as Tensor-level IRs (e.g., XLA (Google, 2019)), the polyhedral model (Griebl et al., 1998), loop synthesizers (e.g., Halide (Ragan-Kelley et al., 2013)), TVM (Chen et al., 2018), and CuTe library (NVIDIA, 2019).

112 113

3 Programming Model

114 115

108

109

110

111

3.1 TILE LANGUAGE

117118119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

116

Tile declarations. TILELANG elevates a *tile*—a hyper-rectangular slice of a tensor—to a first-class citizen. A tile may be owned by a warp, a thread block, or any programmer-defined parallel unit, and can be reshaped or re-partitioned at compile time. In the FlashMLA kernel, the global matrices are consumed in tiles whose extents are parameterized by **block_H**, **block_N**, and related symbolic sizes. The **T.Kernel** structure establishes the kernel's launch configuration (e.g. **bx**, **by**, and the thread count), enabling both index derivation for each thread block and subsequent compiler analyses such as memory-access coalescing and loop tiling.

Tile placement. A distinguishing feature of TILELANG is the ability to map every tile buffer to a concrete level of the target accelerator's memory hierarchy via user-visible intrinsics, rather than relying on opaque compiler heuristics. **T.alloc_shared** reserves storage in low-latency, software-managed shared memory on NVIDIA GPUs (or an architecturally analogous space on other devices). **T.alloc_fragment** places accumulator tiles in the register file. Although registers are scarcer than shared memory, their single-cycle latency is indispensable for performance-critical reductions. During compilation, a layout-inference pass distributes these register tiles across threads while respecting register-pressure constraints and bank conflicts.

Tile operators and schedulable primitives. Table 1 in Appendix C showcases the representative subset of core building blocks that orchestrate computation and movement among tiles. Fundamental operators (**T.copy**, **T.gemm**, **T.reduce**) act on tile operands directly, allowing the programmer to express dense linear algebra, pointwise transforms, and reductions without resorting to scalarized loops. Orthogonal *scheduling primitives* expose fine-grained control over parallelism (**T.Parallel**), pipelining (**T.Pipelined**), and memory layout (**T.annotate_layout**, **T.use_swizzle**).

137 138 139 140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

156

157

158

159

160

161

3.2 A FLASH MULTI-HEAD LATENT ATTENTION EXAMPLE

By fusing high-level expressiveness with architecture-aware orchestration, TILELANG succinctly captures sophisticated AI algorithms such as FlashMLA (Liu et al., 2024) while fully harnessing the performance envelope of modern GPU architectures. Figure 2 illustrates TILELANG's developer-compiler co-optimization model: the *developer* specifies key decisions—such as tile configuration, launch grid (block_H, block_N), buffer placement (T.alloc shared, T.alloc_fragment), swizzled layouts (T.annotate layout), and warp-level collaboration (T.Parallel). The compiler then infers the remaining low-level details, including latency-hiding pipelines, conflict-free memory layouts, and instruction selection for peak hardware performance. To balance flexibility with automation, TILELANG offers two developer-facing facilities. First, tile recommendation (Sec. 4.2) supplies hardware-aware defaults that serve as high-quality starting points. Second, tile inference (Sec. 4.3) analytically

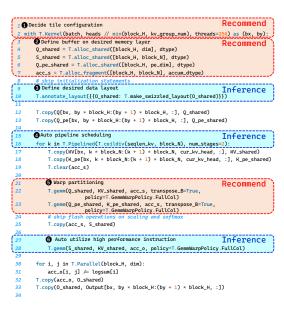


Figure 2: FlashMLA TILELANG kernel example and language features

propagates user-provided or recommended hints to complete the schedule and guarantee consistency. Working in concert, these facilities deliver near-optimal performance with limited manual tuning.

4 SCHEDULING GUIDANCE AND AUTOMATION

4.1 Two-stage framework

Optimization space. High-performance kernel design in TILELANG begins with a tile-level program, represented as a fused tile-level graph (FTG) capturing dataflow and tiling structure—each node represents a tile operator and each edge encodes a data dependency. By operating on this unified graph, TILELANG exposes and reasons about hardware-aware optimizations across six key dimensions: *tile size* (affecting shared memory and register usage), *memory placement* (selecting appropriate memory scope), *warp partitioning* (how threads collaborate and bind within a block), *memory layout* (how tile data is organized across memory levels), *software pipelining* (overlapping compute and data transfer, e.g., via TMA), and *tensorization* (mapping operations to CUDA or Tensor Cores).

Tile recommendation and inference. To efficiently explore the optimization space, TILELANG adopts a unified two-stage workflow over the FTG. In the first stage, *tile recommendation* analyzes the FTG to provide hardware-aware defaults for partially annotated operators, covering dimensions such as initial tile shapes, memory placement, and warp partitioning (Section 4.2). These recommendations shape the memory footprint, compute partitioning, and thread collaboration, providing high-quality starting points. In the second stage, leveraging the context from recommendation, *tile inference* propagates constraints through the FTG, automatically inferring the remaining configuration, including tile size, memory layout, software pipelining, and tensorization. It ensures consistency, compatibility, and hardware efficiency (Section 4.3). Together, these stages unify developer guidance and automated completion: recommendation narrows the design space with informed hints, while inference finalizes fully specified, hardware-efficient kernels with minimal manual effort.

Running example. Taking MLA as an example (Figure 2), TILELANG first performs tile recommendation as illustrated in Figure 3. Tile operators in the FTG expose tunable parameters—such as tile size, memory placement, and warp-partitioning strategies—serving as the user interface for these optimization knobs. For instance, in the first **T.gemm** operator (Figure 3), memory placement annotations specify Q and KV tiles in shared memory, while S resides in registers. The S tile is further partitioned across columns using the "policy=FullCol" warp-partitioning strategy. These decisions directly shape the memory footprint and influence data access patterns across the FTG. The cost model analyzes the FTG to estimate memory traffic, guiding the search toward configurations that minimize data movement.

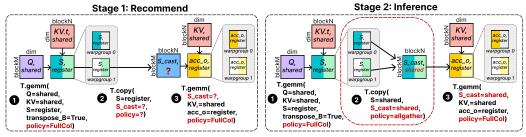


Figure 3: Two-stage workflow of optimizing MLA example.

Tile inference completes the configuration by operating over the FTG. For example, once S (output of the first **T**.gemm in Figure 3(1) and S_cast (input of the second **T**.gemm in Figure 3(3) are fixed in location, shape, and partitioning in the first step, inference automatically determines the tile placement and partitioning (e.g., all-gather or scatter) of **copy** (Figure 3(2)) in the second step, ensuring consistency without manual effort. Beyond copy decisions, inference also derives memory layouts by mapping multi-dimensional indices to physical addresses, explicitly considering vectorization, coalescing, and bank conflicts. Finally, it automates software pipelining and tensorization, ensuring that the resulting kernel configuration is efficient on the underlying hardware. TILELANG also provides platform-specific recommendations and inference (see Appendix D).

4.2 TILE RECOMMENDATION

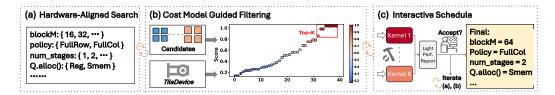


Figure 4: Tile recommendation with cost model. In (b), we show a scatter plot of candidate schedules. The x-axis orders candidates by their $\frac{1}{\text{predicted latency}}$, and the y-axis shows their normalized scores defined as $\frac{\text{latency of best candidate}}{\text{latency of current candidate}}$ which lies in the range (0,1].

Roofline-based cost model. As outlined in Figure 4, TILELANG uses a static roofline-based cost model to evaluate candidate configurations, which include tile shapes, memory placement strategies, and warp partitioning. The cost model operates directly on the fused tile graph (FTG): each FTG under a given configuration is lowered into an intermediate representation (IR), a structured, tile-oriented compute plan that explicitly encodes compute and memory access patterns per tile. From this IR, the model statically extracts two key quantities: total memory traffic at each memory level, and total floating-point operations for each compute type. These quantities are used in a roofline formulation that assumes perfect overlap between computation and memory transfers, ignoring pipeline prologues/epilogues. The execution time is estimated as:

$$Time = \max_{i,j} \left(\frac{MemoryTraffic_i}{Bandwidth_i}, \frac{Computation_j}{Performance_j} \right) + t_{intrinsic}$$
 (1)

where i indexes levels of the memory hierarchy (e.g., HBM, L2, L1), and j indexes compute unit types (e.g., tensor cores, vector cuda cores, special function units (SFUs)). The term $t_{\rm intrinsic}$ accounts for inherent overheads such as kernel launch latency and loop prologue and epilogue costs. This model provides a tight performance upper bound and allows rapid evaluation across large configuration spaces without actual execution or runtime profiling.

Based on the cost model, TILELANG generates actionable recommendations for kernel tuning, including tile shapes, memory placement, and warp partitioning. These recommendations form an interactive baseline: developers can accept, adjust, or iteratively refine them across multiple rounds. This human-in-the-loop workflow balances automation with expert insight, slashing tuning effort while preserving full design control.

Tile size. TILELANG presents a ranked shortlist of tile shapes that are multiples of the device's native tensor-core fragments and respect register and shared-memory limits. Each candidate shows predicted arithmetic intensity, memory traffic, and roofline utilisation. Developers can accept the top choice, pin alternatives for later benchmarking, or adjust dimensions manually.

Memory placement. Given a chosen tile shape, TILELANG enumerates legal bindings of operands and temporaries to registers or shared memory, flagging options that exceed capacity. Each binding includes estimated pipeline stalls and effective bandwidth, letting developers quickly explore tradeoffs and commit or refine placements.

Warp partition. To ensure sufficient thread-level parallelism, TILELANG proposes warp partitions that evenly cover the output tile and match the SM topology. With predicted occupancy and compute—memory overlap, developers can select, benchmark, or override, retaining full control while benefiting from data-driven guidance.

4.3 TILE INFERENCE

Layout inference. While memory placement and computation partitioning in Section 4.2 decide where tensors reside and how computation is split, layout inference determines how multi-dimensional indices are converted into physical memory addresses—taking into account vectorization, memory coalescing, and bank conflict avoidance. In other words, layout is not about which memory scope

is used, but how data is accessed within that scope. Once placement and partitioning are fixed, the system can then infer an appropriate layout to ensure efficient low-level memory access.

TILELANG supports high-level indexing into multi-dimensional arrays (e.g., A[i, k]), which is eventually lowered to physical memory addresses through a hierarchy of abstractions. At the physical level, layouts are modeled as linear address expressions of the form $\sum_i y_i s_i$, where y_i is the index along dimension i, and s_i is its stride. To capture such mappings, TILELANG introduces a composable Layout algebra based on IterVar—a loop iterator that carries range and stride information. This allows layout transformations (e.g., transposes) to be expressed as algebraic mappings, such as Lambda i, j: (j, i). Formally, a layout becomes a function $f: \mathbb{K}^n \to \mathbb{K}^m$, converting high-level indices into memory addresses. Additionally, TILELANG defines Fragment layouts—a specialized extension where $f: \mathbb{K}^n \to \mathbb{K}^2$, mapping each index to a thread's register ID and its local offset. This enables precise modeling of intra-thread register allocation. Although a buffer of size N theoretically allows O(N!) memory layouts, the set of feasible layouts is significantly constrained by hardware. Global memory prefers coalesced access, shared memory requires bank conflict avoidance, and Tensor Core instructions impose strict layout requirements. To explore these constraints, TILELANG employs a greedy strategy that derives valid layouts by enforcing layout rules on selected tile operators.

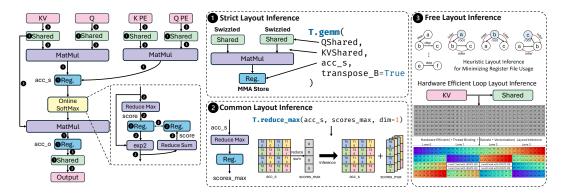


Figure 5: Layout Inference mechanism in TILELANG

We propose a hierarchical layout inference algorithm that operates over an FTG. As illustrated in Figure 5, FlashMLA can be represented as a FTG, where nodes are tile-level operators (e.g., matmul, softmax) and edges encode data dependencies. The graph captures how Q, K, and V tiles are loaded into shared memory, attention scores computed and normalized in registers, and final outputs written back. This structure makes memory movement and parallelism explicit, enabling layout inference and efficient scheduling.

Our goal is to synthesize memory layouts that optimize low-level execution efficiency while preserving high-level tensor semantics. The inference process is modeled as a constraint propagation algorithm (Algorithm 1 in Appendix E) that iteratively traverses the FTG and incrementally refines the layout mapping \mathcal{L} until convergence. As illustrated in Figure 5, the algorithm integrates three complementary inference strategies: (1) Strict Layout Inference (Fig.51) enforces operator-specific constraints for hardware-sensitive primitives such as tensor core GEMM, including swizzled shared memory layouts and MMA-aligned register allocations; (2) Common Layout Inference (Fig.52) propagates layout decisions through structurally aligned operators (e.g., reductions), ensuring consistent thread bindings and register reuse; and (3) Free Layout Inference (Fig. 53) handles the remaining unconstrained layouts by partitioning them into subgraphs via connected component analysis. For each subgraph, the partitioning scheme with the lowest register usage is selected. This step also determines the loop layout using the hardware cost model, which specifies thread binding and vectorization length to maximize memory coalescing and minimize bank conflicts. This unified inference pipeline supports composable, performance-portable layout generation and seamlessly bridges high-level loop indexing with low-level memory organization.

Pipeline inference. TILELANG automatically infers a pipelined schedule from a sequential program. As shown in Figure 6 (a), operations like **copy** and **gemm** are overlapped to increase parallelism. The system analyzes dependencies in the FTG and generates a structured pipeline that preserves

execution correctness, exposing only a single num_stages parameter to users. Additionally, TILELANG applies Warp Specialization to fully exploit asynchronous copy instructions on Hopper GPUs, inserting synchronization barriers where necessary to maintain correct data dependencies.

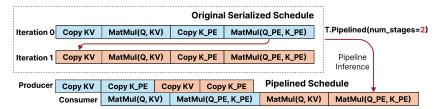


Figure 6: Pipeline Inference mechanism in TILELANG

Instruction inference. In TILELANG, while low-level hardware instructions such as dp4a or mma can be manually invoked via source injection or inline PTX (NVIDIA, 2021), choosing the most appropriate instruction based on input shapes and data types can be challenging. To address this, TILELANG integrates with high-level Tile Libraries like NVIDIA's cute (NVIDIA, 2019) and AMD's ck (AMD, 2025), which abstract hardware-specific details and automatically choose efficient instructions based on input configurations. These libraries expose standardized tile-based APIs (e.g., t1::gemm_ss), and TILELANG supports their invocation via a unified T.call_extern interface, simplifying development while ensuring performance portability.

5 EVALUATION

TILELANG is realised as a Pythonic DSL whose compiler lowers high-level tile programs to hardware-specialised kernels through a modular IR and code-generation pipeline.

5.1 EXPERIMENTAL SETUP

Hardware platforms. We assess the performance of TILELANG on two leading GPU architectures: NVIDIA and AMD, which dominate contemporary accelerator ecosystems. Our evaluation employs state-of-the-art hardware, including the NVIDIA H100 (80GB) (NVIDIA, 2023) and the AMD Instinct MI300X (192GB) (AMD, 2023). The NVIDIA H100 leverages CUDA 12.8, while the AMD MI300X utilizes ROCm 6.2.0. Both GPUs are benchmarked under the Ubuntu 20.04 operating system to ensure consistency in environmental configurations.

AI kernels. To evaluate system performance, we analyze nine representative operators: (1) GEMM, (2) fused dequantized GEMM ($W_{\rm INT4}A_{\rm FP16}$), (3) Attention, (4) Multi-Head Latent Attention, (5) Block Sparse Attention, (6) 2D Convolution, (7) Chunk Gated Delta Net, (8) Vertical Slash Sparse Attention, and (9) Attention Sink. Shape configurations are provided in Appendix H.

Baselines. Our comparative analysis considers the following baselines: (1) PyTorch Inductor—torch.matmul for GEMM, SDPA (PyTorch, 2023) for attention, and other operators compiled via Inductor; (2) Triton implementations, including GemLite (Mobius ML, 2024) and MLA from SGLang (Zheng et al., 2024); (3) ThunderKittens (TK) (Spector et al., 2025)—a template-based framework for high-performance AI kernels on NVIDIA GPUs; and (4) Highly optimized libraries, including CUTLASS (NVIDIA, 2019) and Composable Kernel (AMD, 2025) for GEMM, Marlin (Frantar et al., 2025) for dequantized GEMM, FlashAttention-V3 (Dao, 2023) for MHA, AITER (AMD, 2025) for MLA, and Block Sparse Attention (Guo et al., 2024) for sparse attention.

We evaluate kernel performance versus code complexity (Section 5.2) and present ablation results (Section 5.3), with cost model and tuning time analyses in Appendices F and G.

5.2 Kernel Performance

Matrix Multiplication. TILELANG achieves high performance with low code complexity across diverse GEMM configurations, demonstrating $1.18-1.40\times$ speedup over PyTorch on NVIDIA H100, while maintaining competitive performance $(0.94-1.05\times)$ on AMD MI300X. It also delivers $1.08-1.43\times$ speedup over Triton with minimal kernel code, enabled by automated inference that abstracts low-level hardware details such as TMA and pipeline scheduling. Compared with TK,

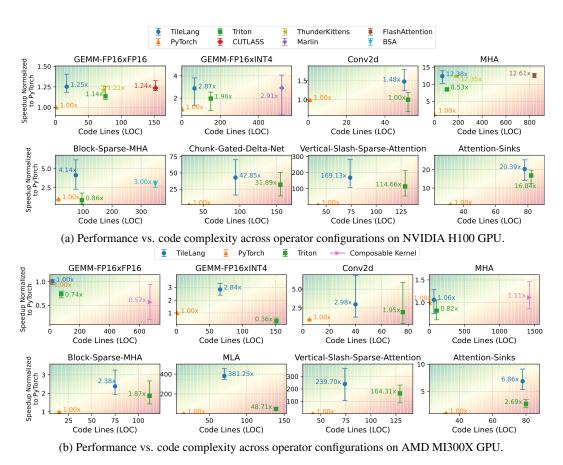


Figure 7: Performance and code complexity on an NVIDIA H100 GPU and AMD MI300X GPU. The y-axis denotes the speedup relative to PyTorch, while the x-axis indicates lines of code (LOC). Ideal solutions appear toward the top-left corner.

TILELANG achieves $0.99-1.11\times$ speedups while reducing code complexity by 77%. Its cost-model guidance and automated tile inference eliminate manual tuning. TK depends on curated CUDA templates, limiting it to NVIDIA GPUs, whereas TILELANG supports multiple hardware backends.

Low-Bit Matmul. For $W_{\rm INT4}A_{\rm FP16}$ GEMM, TILELANG achieves 1.35–3.81× speedups over PyTorch and up to 1.55× over Triton on H100, while outperforming the specialized Marlin kernel with far simpler code. On MI300X, it delivers up to 11.56× over Triton. These gains arise because TILELANG exposes low-level memory, dequantization, and layout controls that Triton hides.

Convolution. On H100, TILELANG achieves 1.24– $1.79\times$ and 1.10– $1.97\times$ speedups over PyTorch and Triton, respectively, with reduced code complexity. These gains come from its instruction inference mechanism, which maps data movement efficiently to TMA im2col. On MI300X, the improvements are even larger, reaching 1.29– $6.80\times$ over PyTorch and 1.02– $3.10\times$ over Triton.

Flash Attention. TILELANG achieves efficient attention computation with concise code across sequence lengths. On H100 and MI300X, it delivers 1.08–1.58× and 1.22–1.37× speedups over Triton, while matching the performance of FlashAttention-V3 (0.98× and 0.96× on average). These results stem from TILELANG's ability to infer and apply platform-specific partitioning and pipelining strategies that exploit specialized compute units. TILELANG achieves up to 1.10× speedup over TK while significantly reducing code complexity (from 185 lines to 66), highlighting its programmability. By combining tile-level guidance with automated inference, TILELANG streamlines kernel development. This is particularly valuable for complex attention operators, where TK often requires extensive manual tuning of tiling, warp partitioning, layout, and pipelines.

Flash MLA. As showin in Figure 1, TILELANG achieves 4.06–10.59× speedups over Triton on H100, with substantially reduced code complexity. It matches the latency of the specialized FlashMLA

kernel while reducing code complexity by $6.86\times$. On MI300X, TILELANG delivers $5.64\text{--}12.97\times$ gains over Triton and slightly outperforms the hand-tuned ROCm library AITER $(1.05\times)$. These improvements arise from warp specialization and automated TMA mapping.

Block Sparse Attention. TILELANG achieves acceleration of $3.42-7.87 \times$ and $1.22-1.37 \times$ over Triton with less code on H100 and MI300X, respectively. On H100, it matches BlockSparse (BSA) latency $(0.91-1.82 \times)$ while greatly reducing complexity. Implementing block-sparse MHA requires only adding two lines to the standard MHA code (Appendix I.5).

Chunk Gated Delta Net. On H100, TILELANG achieves $15.88-70.35\times$ speedups over PyTorch by fusing complex operations into a single kernel. Compared to Triton, it attains $1.10-1.45\times$ speedups with 39% fewer lines of code. These gains come from automated tile recommendation and inference, which optimize memory placement and partitioning for efficient hardware utilization.

Vertical Slash Sparse Attention. TILELANG delivers $108.55-280.41 \times$ and $105.01-363.53 \times$ speedups over PyTorch on H100 and MI300X, largely by fusing the sparse attention operation into a single efficient kernel. Compared to Triton, it achieves $1.16-1.97 \times$ and $1.19-1.60 \times$ speedups on H100 and MI300X, respectively, while cutting the code size by roughly half.

Attention Sinks. For attention with the sinks mechanism, TILELANG achieves $14.21-25.57\times$ and $5.35-9.11\times$ speedups over PyTorch on H100 and MI300X, respectively, enabled by TILELANG's FTG-based fusion into a single optimized kernel. Against Triton, it reaches $1.13-1.30\times$ on H100 and $2.32-2.69\times$ on MI300X. The attention-sink variant can be implemented with only minor changes to standard MHA, demonstrating TILELANG's ability to support diverse attention types with minimal development effort.

5.3 ABLATION STUDIES

To evaluate the impact of TILELANG's automatic optimization, we conduct an ablation study on FlashMLA. Starting from a baseline that uses manually crafted scheduling heuristics (TL-Heuristic), we progressively enable (i) cost model—guided tiling (+Tile), which adjusts the compute-to-memory ratio and optimizes cache utilization; (ii) cost model—guided memory placement (+Alloc), which selects the optimal memory region for each buffer to maximize cache efficiency and avoid register spilling; and (iii) warp partitioning (+Partition), which improves workload balance within warps. Performance is measured at each stage relative to the PyTorch baseline.

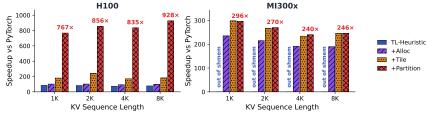


Figure 8: Ablation study for FlashMLA on both H100 and MI300X GPUs.

As shown in Figure 8, each of the evaluated optimizations provides measurable performance gains, validating their effectiveness. Our analysis also highlights architecture-specific behaviors (illustrated in Appendix D). On NVIDIA H100, warp partition contributes the most, achieving a $4.4\times$ speedup over +Alloc by aligning task partitioning with Tensor Core shapes and reducing register spilling, which balances workloads across warps. In contrast, on AMD MI300X, which features abundant registers but limited shared memory, +Alloc delivers the largest improvement. Here, caching query vectors (Q) in registers instead of shared memory leverages the large register file to enhance performance.

6 Conclusion

TILELANG offers a controllable tile-level programming model with graph-based optimizations via tile recommendation and inference. By combining automated configuration with fine-grained developer control, it streamlines kernel development and delivers significant speedups. It enables rapid experimentation with emerging AI algorithms, such as custom attention, sparsity, and quantization. TILELANG also lowers barriers for systems-aware research across diverse hardware platforms.

7 REPRODUCIBILITY STATEMENT

We provide a detailed description of our experimental setup in Section 5.1. Operator shapes used in our benchmarks are drawn from widely adopted, real-world AI models (e.g., GPT-OSS, DeepSeek V3, Qwen3-Next). A list of these operator configurations is included in Appendix H, and the corresponding TILELANG code of kernels used in evaluation is provided in Appendix I. The system implementation and scripts for reproducing our experiments will be made publicly available after the review process, ensuring full reproducibility while maintaining anonymity.

REFERENCES

- Joshua Ainslie, James Lee-Thorp, Michiel De Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*, 2023.
- AMD. Amd cdnaTM 3 architecture. Technical report, AMD, 2023. URL https://www.amd.com/content/dam/amd/en/documents/instinct-tech-docs/white-papers/amd-cdna-3-white-paper.pdf.
- AMD. White paper: Introducing amd cdnaTM 3 architecture, 2024. URL https://www.amd.com/content/dam/amd/en/documents/instinct-tech-docs/white-papers/amd-cdna-3-white-paper.pdf.
 - AMD. ROCm aiter, 2025. https://github.com/ROCm/aiter.
- AMD. AMD Composable Kernel, 2025. https://github.com/ROCm/composable_kernel.
 - Jingwei Cai, Yuchen Wei, Zuotong Wu, Sen Peng, and Kaisheng Ma. Inter-layer scheduling space definition and exploration for tiled accelerators. In *Proceedings of the 50th Annual International Symposium on Computer Architecture*, pp. 1–17, 2023.
 - Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. {TVM}: An automated {End-to-End} optimizing compiler for deep learning. In 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18), pp. 578–594, 2018.
 - Yu Cheng, Lei Wang, Yining Shi, Yuqing Xia, Lingxiao Ma, Jilong Xue, Yang Wang, Zhiwen Mo, Feiyang Chen, Fan Yang, et al. Pipethreader: Software-defined pipelining for efficient dnn execution. In 19th USENIX Symposium on Operating Systems Design and Implementation (OSDI 25), 2025.
 - Tri Dao. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv* preprint arXiv:2307.08691, 2023.
 - Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv* preprint arXiv:2405.21060, 2024.
 - DeepSeek. FlashMLA, 2025. https://github.com/deepseek-ai/FlashMLA.
 - Elias Frantar, Roberto L Castro, Jiale Chen, Torsten Hoefler, and Dan Alistarh. Marlin: Mixed-precision auto-regressive parallel inference on large language models. In *Proceedings of the 30th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming*, pp. 239–251, 2025.
- Mingyu Gao, Xuan Yang, Jing Pu, Mark Horowitz, and Christos Kozyrakis. Tangram: Optimized coarse-grained dataflow for scalable nn accelerators. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 807–820, 2019.
 - Google. XLA, 2019. https://www.tensorflow.org/xla.

- Martin Griebl, Christian Lengauer, and Sabine Wetzel. Code generation in the polytope model. In
 Proceedings. 1998 International Conference on Parallel Architectures and Compilation Techniques (Cat. No. 98EX192), pp. 106–111. IEEE, 1998.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv* preprint arXiv:2312.00752, 2023.
 - Yue Guan, Yuanwei Fang, Keren Zhou, Corbin Robeck, Manman Ren, Zhongkai Yu, Yufei Ding, and Adnan Aziz. Kperfir: Towards a open and compiler-centric ecosystem for gpu kernel performance tooling on modern ai workloads. 2025.
 - Junxian Guo, Haotian Tang, Shang Yang, Zhekai Zhang, Zhijian Liu, and Song Han. Block Sparse Attention. https://github.com/mit-han-lab/Block-Sparse-Attention, 2024.
 - Bastian Hagedorn, Bin Fan, Hanfeng Chen, Cris Cecka, Michael Garland, and Vinod Grover. Graphene: An ir for optimized tensor computations on gpus. In *Proceedings of the 28th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3*, pp. 302–313, 2023.
 - Guyue Huang, Yang Bai, Liu Liu, Yuke Wang, Bei Yu, Yufei Ding, and Yuan Xie. Alcop: Automatic load-compute pipelining in deep learning compiler for ai-gpus. *Proceedings of Machine Learning and Systems*, 5:680–694, 2023.
 - Hyoukjun Kwon, Prasanth Chatarasi, Vivek Sarkar, Tushar Krishna, Michael Pellauer, and Angshuman Parashar. Maestro: A data-centric approach to understand reuse, performance, and hardware cost of dnn mappings. *IEEE micro*, 40(3):20–29, 2020.
 - Chris Lattner, Mehdi Amini, Uday Bondhugula, Albert Cohen, Andy Davis, Jacques Pienaar, River Riddle, Tatiana Shpeisman, Nicolas Vasilache, and Oleksandr Zinenko. Mlir: Scaling compiler infrastructure for domain specific computation. In 2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO), pp. 2–14. IEEE, 2021. https://mlir.llvm.org/.
 - Seonho Lee, Amar Phanishayee, and Divya Mahajan. Forecasting gpu performance for deep learning training and inference. In *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 1*, pp. 493–508, 2025.
 - Ying Li, Yifan Sun, and Adwait Jog. Path forward beyond simulators: Fast and accurate gpu execution time prediction for dnn workloads. In *Proceedings of the 56th Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 380–394, 2023.
 - Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, et al. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*, 2024.
 - Mobius ML. Gemlite: Lightweight generative machine learning engine. https://github.com/mobiusml/gemlite, 2024. Accessed: [Insert Access Date].
 - NVIDIA. NVIDIA cutlass, 2019. https://github.com/NVIDIA/cutlass.
- NVIDIA. PTX ISA, 2021. https://docs.nvidia.com/cuda/parallel-thread-execution/.
 - NVIDIA. Nvidia h100 tensor core gpu architecture. Technical report, NVIDIA Corporation, 2023. URL https://resources.nvidia.com/en-us-tensor-core.
- Angshuman Parashar, Priyanka Raina, Yakun Sophia Shao, Yu-Hsin Chen, Victor A Ying, Anurag Mukkara, Rangharajan Venkatesan, Brucek Khailany, Stephen W Keckler, and Joel Emer. Timeloop: A systematic approach to dnn accelerator evaluation. In 2019 IEEE international symposium on performance analysis of systems and software (ISPASS), pp. 304–315. IEEE, 2019.
- Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. https://pytorch.org/.

- PyTorch. torch.nn.functional.scaled_dot_product_attention, 2023. URL https://docs.pytorch.org/docs/stable/generated/torch.nn.functional.scaled_dot_product_attention.html.
- Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. Halide: A language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. In *Proceedings of the 34th ACM SIGPLAN Conference on Programming Language Design and Implementation*, PLDI '13, pp. 519–530, New York, NY, USA, 2013. ACM. ISBN 978-1-4503-2014-6. doi: 10.1145/2491956.2462176. URL http://doi.acm.org/10.1145/2491956.2462176.
- Jay Shah, Ganesh Bikshandi, Ying Zhang, Vijay Thakkar, Pradeep Ramani, and Tri Dao. Flashattention-3: Fast and accurate attention with asynchrony and low-precision. *arXiv* preprint *arXiv*:2407.08608, 2024.
- Yining Shi, Zhi Yang, Jilong Xue, Lingxiao Ma, Yuqing Xia, Ziming Miao, Yuxiao Guo, Fan Yang, and Lidong Zhou. Welder: Scheduling deep learning memory access via tile-graph. In *17th USENIX Symposium on Operating Systems Design and Implementation (OSDI 23)*, pp. 701–718, 2023.
- Benjamin F Spector, Simran Arora, Aaryan Singhal, Daniel Y Fu, and Christopher Ré. Thunderkittens: Simple, fast, and adorable ai kernels. In *Proceedings of the 13th International Conference on Learning Representations*, 2025. URL https://openreview.net/forum?id=0fJfVOSUra.https://github.com/HazyResearch/ThunderKittens.
- Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *arXiv preprint arXiv:2307.08621*, 2023.
- Philippe Tillet, H. T. Kung, and David Cox. *Triton: An Intermediate Language and Compiler for Tiled Neural Network Computations*, pp. 10–19. Association for Computing Machinery, New York, NY, USA, 2019. ISBN 9781450367196. URL https://doi.org/10.1145/3315508.3329973.
- Nicolas Vasilache, Oleksandr Zinenko, Theodoros Theodoridis, Priya Goyal, Zachary DeVito, William S. Moses, Sven Verdoolaege, Andrew Adams, and Albert Cohen. Tensor comprehensions: Framework-agnostic high-performance machine learning abstractions. *CoRR*, abs/1802.04730, 2018. URL http://arxiv.org/abs/1802.04730.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- Lei Wang, Lingxiao Ma, Shijie Cao, Quanlu Zhang, Jilong Xue, Yining Shi, Ningxin Zheng, Ziming Miao, Fan Yang, Ting Cao, et al. Ladder: Enabling efficient {Low-Precision} deep learning computing through hardware-aware tensor transformation. In 18th USENIX Symposium on Operating Systems Design and Implementation (OSDI 24), pp. 307–323, 2024.
- Rohan Yadav, Michael Garland, Alex Aiken, and Michael Bauer. Task-based tensor computations on modern gpus. *Proceedings of the ACM on Programming Languages*, 9(PLDI):396–420, 2025.
- Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. Parallelizing linear transformers with the delta rule over sequence length. *Advances in neural information processing systems*, 37: 115491–115522, 2024.
- Lianmin Zheng, Liangsheng Yin, Zhiqiang Xie, Chuyue Livia Sun, Jeff Huang, Cody Hao Yu, Shiyi Cao, Christos Kozyrakis, Ion Stoica, Joseph E Gonzalez, et al. Sglang: Efficient execution of structured language model programs. *Advances in Neural Information Processing Systems*, 37: 62557–62583, 2024.

- Our appendix is organized as follows:
- Appendix A: Extended discussion of related work.
- Appendix B: Details of the MLA algorithm.
- Appendix C: Semantics of a partial list of TILELANG primitives.
- 654 Appendix D: Platform-specific scheduling

- Appendix E: Layout inference algorithm of TILELANG.
- Appendix F: Evaluation of the cost model.
- Appendix G: Tuning time measurements.
- Appendix H: Operator shapes used in our benchmark.
 - Appendix I: TileLang code of kernels used in the evaluation.

A EXTENDED DISCUSSION OF RELATED WORK

Tensor-level IRs, The polyhedral model, and Loop synthesizers. Traditional approaches address program optimization at different abstraction levels: Tensor-level IRs (e.g., XLA (Google, 2019)) lower tensor programs via pattern-matched templates (e.g., LLVM, CUDA). Polyhedral models (Griebl et al., 1998) (e.g., TC (Vasilache et al., 2018)) automate affine loop transforms, mainly for DNN layers. Loop synthesizers (e.g., Halide (Ragan-Kelley et al., 2013)) generate loop nests guided by user-defined schedules. TILELANG targets a distinct programming model and control granularity. It differs fundamentally by introducing tiles as first-class programming units. It offers programmable control over fusion strategies, memory hierarchy, and parallelism. This enables developers to design fused kernels with both high performance and portability across hardware.

TVM (Chen et al., 2018). TILELANG builds upon TVM's IR and arithmetic passes. However, unlike TVM's schedule-driven loop generation from high-level compute definitions, TILELANG offers explicit, tile-level programmability and control over memory, fusion, and parallelism. This enables much finer-grained kernel customization beyond what TVM can achieve. For instance, TVM cannot fully express advanced algorithms like FlashAttention (FA) or Multi-Level Attention (MLA), which demand precise management of memory hierarchy and execution order—capabilities that TileLang supports.

Warp Partition. Warp Partition (WP) is a key component of TILELANG's execution model, building directly on the tile abstraction. Given a specified tile size, WP allows further partitioning of the tile along each dimension across multiple warps. For example, consider a GEMM operation C = A@B, where $A \in \mathbb{R}^{M \times K}$, $B \in \mathbb{R}^{K \times N}$, and $C \in \mathbb{R}^{M \times N}$. The output tile C can be partitioned along either the M or N axis, corresponding to full-row or full-column warp-partitioning strategies, respectively. By giving users explicit control over warp partitioning, TILELANG enables fine-grained management of resources such as register usage within each warp. This, in turn, allows users to better control the performance of operations. Such flexibility is crucial for mapping computations efficiently to hardware, especially when optimizing diverse and performance-sensitive kernel workloads.

CuTe library. While both TILELANG and CuTe (NVIDIA, 2019) share this high-level goal, their underlying mechanisms differ: CuTe relies on shape/stride pairs, whereas TILELANG encodes the mappings using explicit arithmetic expressions. This arithmetic formulation offers advantages by more directly capturing index transformations and enabling more flexible, composable manipulations, allowing for clear definition and description in the DSL frontend.

The roofline-guided cost model. Several analytical modeling approaches have been proposed, such as the nested-loop-based modeling in Timeloop (Parashar et al., 2019), the data-centric representation in Maestro (Kwon et al., 2020). In contrast to these methods, our work leverages a tile-level programming abstraction, which naturally lends itself to a tile-centric cost model. This enables us to accurately capture both computation and data movement at the tile granularity, while maintaining simplicity and enhanced support for modeling operator fusion. This design strikes a balance between accuracy and usability, making it effective for guiding schedule selection without introducing excessive complexity.

B MLA ALGORITHM

Instead of storing full-sized key and value matrices, MLA projects input token embeddings into a lower-dimensional latent space using a down-projection matrix:

$$\mathbf{z}_t = \mathbf{x}_t \mathbf{W}_{\text{down}}$$
.

The latent vector \mathbf{z}_t is then used to reconstruct the key and value representations:

$$\mathbf{k}_t = \mathbf{z}_t \mathbf{W}_{\mathsf{up}}^K, \quad \mathbf{v}_t = \mathbf{z}_t \mathbf{W}_{\mathsf{up}}^V.$$

To incorporate positional information, Rotary Positional Embedding (RoPE) is applied to the reconstructed keys and queries:

$$\mathbf{k}_t^{\mathrm{rot}} = \mathrm{RoPE}(\mathbf{k}_t).$$

Queries are also compressed using a similar process to reduce activation memory:

$$\mathbf{q}_t = \mathbf{z}_t^Q \mathbf{W}_{\mathrm{up}}^Q.$$

MLA further enhances computational efficiency through a technique known as *matrix absorption*, which reorders matrix multiplications to optimize performance. This approach enables the key and value inputs to share the same latent representation \mathbf{z}_t , thereby reducing redundancy and memory usage. In the adopted configuration, MLA employs a single shared key-value (KV) head, with a head dimension of 512.

C PARTIAL LIST OF TILELANG PRIMITIVES

Table 1: A partial list of primitives supported by TILELANG.

	Dataflow Centric Tile Operators	Scheduling Primitives			
сору	data movement among hierarchy memory.	Parallel	Parallelization of loop iterations over threads.		
gemm	matrix multiplication on different GPUs.	Pipelined	Enables pipelining to overlap data transfers with computation.		
reduce	reduction operator (e.g., sum, min, max) exploiting warp/block-level parallelism.	annotate_layout	Definition of custom memory layouts to minimize bank conflicts and optimize thread binding.		
atomic	atomic operations to ensure thread-safe updates in shared or global memory.	use_swizzle	Improves L2 cache locality via swizzled access patterns.		

D PLATFORM-SPECIFIC SCHEDULING

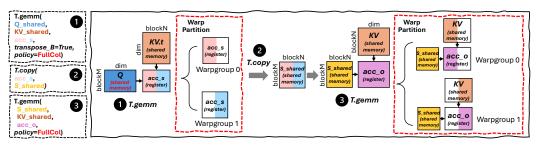
TILELANG also provides platform-specific recommendations and inference. Taking MLA as an example, we illustrate how TILELANG performs tile recommendation and inference based on the code shown in Figure 2.

On the H100, each SM features 228 KiB of shared memory and a 256 KiB register file, whereas the MI300X provides 64 KiB of Local Data Share (LDS) and a total of 512 KiB in registers. Given these architectural differences, TILELANG first recommends different tile configurations.

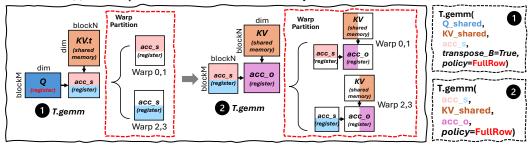
As shown in Figure 9, for memory placement, users may initially allocate the Q tile to shared memory on the MI300X. However, this approach fails due to the limited capacity of shared memory. TILELANG detects this constraint and instead recommends placing both Q and acc_s in registers. In contrast, on the H100, both tiles fit comfortably in shared memory and are placed there accordingly. For Software Pipelining, TILELANG disables pipelining on the MI300X to support larger tile sizes and reduce register pressure, whereas on the H100, pipelining is enabled to maximize pipeline overlap. Tile sizes are also adjusted accordingly to fit each platform's resource constraints. For Warp Partitioning, users may initially adopt a default policy for the two gemm operators, which often leads to sub-optimal performance. TILELANG addresses this by analyzing the underlying hardware and recommending platform-specific partitioning strategies, as illustrated in Figure 3. On the H100, both gemm operators use the FullCol scheme, partitioning acc_s and acc_o vertically to match the Tensor Core shape. In contrast, TILELANG applies a FullRow policy on the MI300X, partitioning tiles horizontally.

Table 2: Comparison of specifications between NVIDIA H100 SXM and AMD MI300X.

NVIDIA H100 SXM	AMD MI300X
1.83 GHz	2.10 GHz
3.35 TB/s	5.30 TB/s
9.45 TB/s	16.63 TB/s
30.92 TB/s	81.72 TB/s
132 SMs	304 CUs
228 KiB	64 KiB
256 KiB	512 KiB
989 TFLOPs	1307 TFLOPs
	1.83 GHz 3.35 TB/s 9.45 TB/s 30.92 TB/s 132 SMs 228 KiB 256 KiB



(a) The Warp Partition schedule recommended by TILELANG for MLA on H100.



(b) The Warp Partition schedule recommended by TILELANG for MLA on MI300X.

Figure 9: Cooperative workflow between tile-recommendation and inference stages on NVIDIA H100 and AMD MI300X GPUs.

E LAYOUT INFERENCE ALGORITHM

```
Algorithm 1 Hardware-Aware Layout Inference over Fused Tile Graph
Require: Fused Tile Graph \mathcal{G} = (V, E), where V is the set of tile operators
Ensure: Final layout mapping \mathcal{L} for all buffers
 1: Initialize layout map \mathcal{L} \leftarrow \emptyset
 2: Initialize constraint worklist Q \leftarrow SEEDSTRICTCONSTRAINTS(\mathcal{G})
 3:
     while not Q.EMPTY() do
           v \leftarrow \mathcal{Q}.POP()
 4:
           C_v \leftarrow \text{COLLECTCONSTRAINTS}(v, \mathcal{L})
 5:
           \Delta \mathcal{L} \leftarrow \text{InferLayout}(v, \mathcal{C}_v)
 6:
 7:
           if \Delta \mathcal{L} \neq \emptyset then
 8:
                \mathcal{L} \leftarrow \mathcal{L} \cup \Delta \mathcal{L}
                 Q \leftarrow Q \cup \text{Neighbors}(v)
 9:
10: return \mathcal{L}
```

F EFFECTIVENESS OF THE COST MODEL

 TILELANG employs an analytical cost model to prune suboptimal candidates and prioritize high-potential ones. This approach yields schedules that match or closely approach the performance of the best results of brute-force search or exhaustive autotuning, while requiring orders of magnitude less tuning effort. For example, on GEMM-FP16×FP16 workloads derived from models such as LLaMA-70B, TILELANG prunes 95% of candidate schedules, retaining only the top 5%. Despite this aggressive pruning, it achieves on average 98.47% of the performance (in TFLOPS) of the best configurations found by exhaustive search, substantially reducing compilation time with only negligible performance loss. This also serves as an evaluation of the effectiveness of our roofline-based analytical cost model.

\overline{M}	N	K	Predicted-TopX / Best (TFLOPS)
512	1024	8192	100.0%
512	12288	12288	99.9%
512	28672	8192	98.7%
2048	12288	49152	100.0%
4096	1024	7168	100.0%
4096	14336	14336	100.0%
4096	28672	8192	99.6%
8192	8192	28672	100.0%
8192	28672	8192	100.0%
16384	1024	7168	98.4%

Table 3: Accuracy of our analytical cost model: predicted top-5% schedules retain over 98% of the best performance while pruning 95% of candidate schedules.

G TUNING TIME

As demonstrated in Table 4, TILELANG leverages a hardware-aware recommendation mechanism to efficiently automate the design of high-performance computational kernels. The system achieves average tuning durations of approximately 10 seconds across both NVIDIA H100 and AMD MI300X accelerators. For the most complex operations, extended tuning times average 13.69 seconds on H100 and 15.08 seconds on MI300X, reflecting the scalability of our approach under computationally intensive workloads.

Table 4: Average Tuning Times for Different Operators

Operation	GEMM	DequantGEMM	FlashMHA	FlashMLA	FlashBSA
H100 Time (s)	9.05	9.15	13.48	13.69	13.46
MI300 Time (s)	10.99	11.10	14.67	15.03	15.08

H OPERATOR SHAPES IN OUR BENCHMARK

Table 5: Matrix shapes in our FP16 Matmul evaluation

	D0	D1	D2	D3
M	8192	8192	8192	8192
N	1024	8192	28672	8192
K	8192	8192	8192	28672

Table 6: Matrix shapes in our Fused Dequantize-Matmul evaluation

	M0	M1	M2	M3
M	1	1	1	1
N	1024	8192	28672	8192
K	8192	8192	8192	28672

Table 7: FlashAttention and Block Sparse Attention(with 50%, 90% sparsity) shapes in our evaluation

	FA0	FA1	FA2	FA3	FA4	FA5	FA6	FA7
batch	64	64	64	64	64	64	64	64
nheads	64	64	64	64	64	64	64	64
seq_len	1024	2048	4096	8192		2048	4096	8192
head_dim	128	128	128	128	128	128	128	128
causal	false	false	false	false	true	true	true	true

Table 8: FlashMLA shapes in our evaluation

	FMLA0	FMLA1	FMLA2	FMLA3
batch	64	64	64	64
nheads	128	128	128	128
seq_len	1024	2048	4096	8192
head_dim	512	512	512	512
pe_dim	64	64	64	64
causal	false	false	false	false

Table 9: Convolution-2D shapes in our evaluation

	Conv0	Conv1	Conv2	Conv3	Conv4	Conv5	Conv6	Conv7
N	128	128	128	128	128	128	128	128
C	2048	512	512	512	256	1024	512	64
Н	7	7	14	7	14	14	28	56
W	7	7	14	7	14	14	28	56
F	512	2048	512	512	256	256	128	64
K	1	1	3	3	3	1	1	1
S	1	1	2	1	1	1	1	1
D	1	1	1	1	1	1	1	1
P	0	0	1	1	1	0	0	0
G	1	1	1	1	1	1	1	1
НО	7	7	7	7	14	14	28	56
WO	7	7	7	7	14	14	28	56
Count	2	3	1	2	5	5	3	1



Table 10: Chunk-Gated-Delta-Net kernel shapes in our evaluation

	CGDN0	CGDN1	CGDN2	CGDN3	CGDN4	CGDN5
batch	1	1	1	64	64	64
nheads	32	32	32	32	32	32
seq_len	16384	32768	65536	1024	2048	4096
head_dim	128	128	128	128	128	128

Table 11: Vertical Slash Sparse Attention shapes in our evaluation

	VSSA0	VSSA1	VSSA2	VSSA3
batch	1	1	1	1
nheads	1	1	1	1
seq_len	8192	16384	32768	65536
head_dim	64	64	64	64
vertical size	1000	1000	800	1000
slash size	600	200	600	600

Table 12: Attention Sink shapes in our evaluation

	Sink0	Sink1	Sink2	Sink3
batch	1	1	1	1
nheads	64	64	64	64
kv_heads	8	8	8	8
seq_len	1024	2048	4096	8192
kv_seq_len	1024	2048	4096	8192
head_dim	64	64	64	64
casual	true	true	true	true

I KERNEL IMPLEMENTATIONS

972

973 974

975

987 988 989

990

1023

1024 1025

I.1 MATRIX MULTIPLICATION (MATMUL)

```
976 1
          @tilelang.jit
          def Matmul(A: T.Tensor, B: T.Tensor, C: T.Tensor):
977 2
            with T.Kernel(N // block_N, M // block_M,
978 4
              threads=threads) as (bx, by):
               A_shared = T.alloc_shared(block_M, block_K)
B_shared = T.alloc_shared(block_K, block_N)
979 5
980 7
               C_local = T.alloc_fragment(block_M, block_N)
981 8
               T.clear(C_local)
               for k in T.Pipelined(K // block_K, num_stages=2):
982 10
                    T.copy(A[by * block_M, k * block_K], A_shared)
T.copy(B[k * block_K, bx * block_N], B_shared)
983 ...
984 13
                    T.gemm(A_shared, B_shared, C_local)
985 <sup>14</sup> <sub>15</sub>
               T.copy(C_local, C[by * block_M, bx * block_N])
986
```

Figure 10: Kernel Implementation of Matrix Multiplication.

I.2 DEQUANTIZED MATRIX MULTIPLICATION

```
991
        @tilelang.jit
        def dequantize_gemv(A: T.Tensor, B: T.Tensor, C: T.Tensor):
992 2
             with T.Kernel(T.ceildiv(N, n_partition), M, threads=(reduce_thread, n_partition)) as (bx, by):
993
             A_local = T.alloc_local((micro_size_k,), in_dtype)
             B_quant_local = T.alloc_local([micro_size_k_compressed], storage_dtype)
994 5
             B_dequantize_local = T.alloc_local([micro_size_k], in_dtype)
995 7
             accum_res = T.alloc_local((1,), accum_dtype)
996 8
             reduced_accum_res = T.alloc_local((1,), accum_dtype)
997 10
             T.clear(accum_res)
998 11
             for ko in T.serial(T.ceildiv(K, block_K)):
    12
                 for v in T.vectorized(micro_size_k):
999 13
                     A_local[v] = A[by, ko * block_K + kr * micro_size_k + v]
1000_{15}^{14}
                 for v in T.vectorized(micro_size_k_compressed):
100116
                     B_quant_local[v] = B[
    bx * n_partition + ni,
1002_{18}^{17}
                          ko * (reduce_thread * micro_size_k_compressed) +
100319
                          kr * micro_size_k_compressed + v,
                     ]
1004_{21}^{20}
100522
                 T.call_extern(
100624
                      "fast_decode_int4",
                     T.address_of(B_quant_local[0]),
100725
                     T.address_of(B_dequantize_local[0]),
                     dtype=in_dtype,
100827
1009_{29}^{28}
                 for ki in T.serial(micro_size_k):
101030
                     accum_res[0] += A_local[ki] * B_dequantize_local[ki]
1011_{32}^{31}
             with T.attr(
101233
                     T.comm_reducer(lambda x, y: x + y, [T.Cast(accum_dtype, 0)]),
101334
                      "reduce_scope"
                     T.reinterpret(T.uint64(0), dtype="handle"),
101436
101537
                 T.evaluate(
                     T.tvm_thread_allreduce(
101639
                          T.uint32(1),
    40
                          accum_res[0],
101741
                          True.
101842
                          reduced_accum_res[0],
    43
                          kr.
101944
                          dtype="handle",
                     ))
1020_{46}^{45}
             if kr == 0:
102147
                 C[by, bx * n_partition + ni] = reduced_accum_res[0]
1022
```

Figure 11: Implementation of Weight-Only Quantization ($W_{\text{FP4_E2M1}}A_{\text{FP16}}$) Matmul using TILE-LANG, showcasing support for mixed-precision computations via a simple form.

I.3 FLASH ATTENTION IMPLEMENTATION

1026

1027

```
1028 1
         @tilelang.jit
1029 2
                     _attention(Q: T.Tensor, K: T.Tensor, V: T.Tensor, Output: T.Tensor):
         def flash
              with T.Kernel (
1030 4
                    \begin{tabular}{ll} T.ceildiv(seq\_len, block\_M), heads, batch, threads=threads) & as (bx, by, bz): \\ Q\_shared = T.alloc\_shared([block\_M, dim], dtype) \\ \end{tabular} 
1031\frac{5}{6}
                   K_shared = T.alloc_shared([block_N, dim], dtype)
                   V_shared = T.alloc_shared([block_N, dim], dtype)
10327
                   O_shared = T.alloc_shared([block_M, dim], dtype)
1033 0
                   acc_s = T.alloc_fragment([block_M, block_N], accum_dtype)
103410
                   acc_s_cast = T.alloc_fragment([block_M, block_N], dtype)
103511
                   acc_o = T.alloc_fragment([block_M, dim], accum_dtype)
                   scores_max = T.alloc_fragment([block_M], accum_dtype)
                   scores_max_prev = T.alloc_fragment([block_M], accum_dtype)
scores_scale = T.alloc_fragment([block_M], accum_dtype)
scores_sum = T.alloc_fragment([block_M], accum_dtype)
103613
103715
103816
                   logsum = T.alloc_fragment([block_M], accum_dtype)
103918
                   T.copy(Q[bz, bx * block_M: (bx + \frac{1}{1}) * block_M, by, :], Q_shared)
                   T.fill(acc_o, 0)
1040_{20}^{19}
                   T.fill(logsum, 0)
104121
                   T.fill(scores_max, -T.infinity(accum_dtype))
1042_{23}^{22}
                   loop_range = (
104324
                        T.min(T.ceildiv(seq_len, block_N), T.ceildiv(
1044_{26}^{25}
                             (bx + 1) * block_M, block_N)) if is_causal else T.ceildiv(seq_len, block_N))
104527
                   for k in T.Pipelined(loop_range, num_stages=num_stages):
                        T.copy(K[bz, k * block_N:(k + 1) * block_N, by, :], K_shared)
104629
                        if is causal:
                             for i, j in T.Parallel(block_M, block_N):
104730
                                 acc_s[i, j] = T.if_then_else(bx \star block_M + i >= k \star block_N + j, 0,
104832
                                                                   -T.infinity(acc_s.dtype))
1049_{34}^{33}
                        else:
                             T.clear(acc_s)
105035
                        T.gemm(Q_shared, K_shared, acc_s, transpose_B=True, policy=T.GemmWarpPolicy.FullRow)
1051<sup>36</sup><sub>37</sub>
                        T.copy(scores_max, scores_max_prev)
                        T.fill(scores_max, -T.infinity(accum_dtype))
105238
                        T.reduce_max(acc_s, scores_max, dim=1, clear=False)
1053 39 40
                       for i in T.Parallel(block_M):
                            scores_scale[i] = T.exp2(scores_max_prev[i] * scale - scores_max[i] * scale)
105441
                       for i, j in T.Parallel(block_M, block_N):
1055_{43}^{42}
                            acc_s[i, j] = T.exp2(acc_s[i, j] * scale - scores_max[i] * scale)
                       T.reduce_sum(acc_s, scores_sum, dim=1)
for i in T.Parallel(block_M):
105644
                             logsum[i] = logsum[i] * scores_scale[i] + scores_sum[i]
105746
                        T.copy(acc_s, acc_s_cast)
                       for i, j in T.Parallel(block_M, dim):
    acc_o[i, j] *= scores_scale[i]
1058<sup>47</sup>
48
105949
                        T.copy(V[bz, k * block_N:(k + 1) * block_N, by, :], V_shared)
1060_{51}^{50}
                        T.gemm(acc_s_cast, V_shared, acc_o, policy=T.GemmWarpPolicy.FullRow)
                   for i, j in T.Parallel(block_M, dim):
106152
                        acc_o[i, j] /= logsum[i]
1062_{54}^{53}
                   T.copy(acc_o, O_shared)
                   T.copy(O_shared, Output[bz, bx * block_M: (bx + \frac{1}{1}) * block_M, by, :])
1063
```

Figure 12: Implementation of Flash Attention with TILELANG.

```
1080
          I.4 FLASHMLA IMPLEMENTATION
1081
1082 i
          @tilelang.jit
         def flash_mla(
1083 2
                   Q: T.Tensor([batch, heads, dim], dtype),
1084 4
                   Q_pe: T.Tensor([batch, heads, pe_dim], dtype),
KV: T.Tensor([batch, seqlen_kv, kv_head_num, dim], dtype),
1085 \frac{5}{6}
                   K_pe: T.Tensor([batch, seqlen_kv, kv_head_num, pe_dim], dtype),
                   Output: T. Tensor([batch, heads, dim], dtype),
10867
         ):
1087 <sup>8</sup><sub>9</sub>
              with T.Kernel(batch, heads // min(block_H, kv_group_num), threads=256) as (bx, by):
                   Q_shared = T.alloc_shared([block_H, dim], dtype)
S_shared = T.alloc_shared([block_H, block_N], dtype)
108810
108911
                   Q_pe_shared = T.alloc_shared([block_H, pe_dim], dtype)
                   KV_shared = T.alloc_shared([block_N, dim], dtype)
109013
                   K_pe_shared = T.alloc_shared([block_N, pe_dim], dtype)
109115
                   O_shared = T.alloc_shared([block_H, dim], dtype)
                   acc_s = T.alloc_fragment([block_H, block_N], accum_dtype)
1092^{16}
                   acc_o = T.alloc_fragment([block_H, dim], accum_dtype)
scores_max = T.alloc_fragment([block_H], accum_dtype)
109318
                   scores_max_prev = T.alloc_fragment([block_H], accum_dtype)
1094_{20}^{19}
                   scores_scale = T.alloc_fragment([block_H], accum_dtype)
109521
                   scores_sum = T.alloc_fragment([block_H], accum_dtype)
1096_{23}^{22}
                   logsum = T.alloc_fragment([block_H], accum_dtype)
109724
                   cur_kv_head = by // (kv_group_num // block_H)
1098<sup>25</sup><sub>26</sub>
                   T.use_swizzle(10)
                   T.copy(Q[bx, by * VALID_BLOCK_H: (by + 1) * VALID_BLOCK_H, :], Q_shared)
109927
                   T.copy(Q_pe[bx, by * VALID_BLOCK_H:(by + 1) * VALID_BLOCK_H, :], Q_pe_shared)
1100_{29}^{28}
                   T.fill(acc_o, 0)
1101^{30}
                   T.fill(logsum, 0)
                   T.fill(scores_max, -T.infinity(accum_dtype))
110232
                   loop_range = T.ceildiv(seqlen_kv, block_N)
1103<sup>33</sup><sub>34</sub>
                   for k in T.Pipelined(loop_range, num_stages=2):
110435
                        T.copy(KV[bx, k * block_N:(k + 1) * block_N, cur_kv_head, :], KV_shared)
1105_{37}^{36}
                        T.copy(K_pe[bx, k * block_N:(k + 1) * block_N, cur_kv_head, :], K_pe_shared)
                        T.clear(acc_s)
110638
                        T.gemm (
1107<sup>39</sup><sub>40</sub>
                            Q_shared, KV_shared, acc_s, transpose_B=True, policy=T.GemmWarpPolicy.FullCol)
                        T.gemm(
110841
                            Q_pe_shared,
1109_{43}^{42}
                            K_pe_shared,
                             acc s,
                             transpose_B=True,
111044
                            policy=T.GemmWarpPolicy.FullCol)
1111_{46}^{-1}
                       T.copy(scores_max, scores_max_prev)
                        T.fill(scores_max, -T.infinity(accum_dtype))
1112<sup>47</sup>
48
                        T.reduce_max(acc_s, scores_max, dim=1, clear=False)
111349
                       for i in T.Parallel(block_H):
1114_{51}^{50}
                            scores_scale[i] = T.exp2(scores_max_prev[i] * scale - scores_max[i] * scale)
                       for i, j in T.Parallel(block_H, block_N):
111552
                            acc_s[i, j] = T.exp2(acc_s[i, j] * scale - scores_max[i] * scale)
1116_{54}^{53}
                       T.reduce_sum(acc_s, scores_sum, dim=1)
                       T.copy(acc_s, S_shared)
111755
                       for i in T.Parallel(block_H):
                       logsum[i] = logsum[i] * scores_scale[i] + scores_sum[i]
for i, j in T.Parallel(block_H, dim):
1118_{57}^{56}
                       acc_o(i, j] *= scores_scale[i]
T.gemm(S_shared, KV_shared, acc_o, policy=T.GemmWarpPolicy.FullCol)
111958
1120<sub>60</sub>
                   for i, j in T.Parallel(block_H, dim):
    acc_o[i, j] /= logsum[i]
112161
                   T.copy(acc_o, O_shared)
112263
                   T.copy(O_shared, Output[bx, by * VALID_BLOCK_H: (by + 1) * VALID_BLOCK_H, :])
1123
```

Figure 13: Implementation of FlashMLA with TILELANG.

I.5 BLOCK SPARSE ATTENTION IMPLEMENTATION

1134

1135

```
1136<sub>1</sub>
          @tilelang.jit
1137 2
          def blocksparse_attn(Q: T.Tensor, K: T.Tensor, V: T.Tensor, BlockMask: T.Tensor, Output: T.Tensor):
               with T.Kernel (
1138 4
                     \begin{tabular}{ll} T.ceildiv(seq\_len, block\_M), heads, batch, threads=threads) & as (bx, by, bz): \\ Q\_shared = T.alloc\_shared([block\_M, dim], dtype) \\ \end{tabular} 
1139_{6}^{5}
                    K_shared = T.alloc_shared([block_N, dim], dtype)
                    V_shared = T.alloc_shared([block_N, dim], dtype)
11407
                    O_shared = T.alloc_shared([block_M, dim], dtype)
1141 8
                    acc_s = T.alloc_fragment([block_M, block_N], accum_dtype)
114210
                    acc_s_cast = T.alloc_fragment([block_M, block_N], dtype)
114311
                    acc_o = T.alloc_fragment([block_M, dim], accum_dtype)
                    scores_max = T.alloc_fragment([block_M], accum_dtype)
                    scores_max_prev = T.alloc_fragment([block_M], accum_dtype)
scores_scale = T.alloc_fragment([block_M], accum_dtype)
scores_sum = T.alloc_fragment([block_M], accum_dtype)
114413
114515
1146^{16}
                    logsum = T.alloc_fragment([block_M], accum_dtype)
114718
                    T.copy(Q[bz, bx * block_M: (bx + \frac{1}{1}) * block_M, by, :], Q_shared)
                    T.fill(acc_o, 0)
1148_{20}^{19}
                    T.fill(logsum, 0)
114921
                   T.fill(scores_max, -T.infinity(accum_dtype))
1150_{23}^{22}
                   loop_range = (
115124
                         T.min(T.ceildiv(seq_len, block_N), T.ceildiv(
1152_{26}^{25}
                              (bx + 1) * block_M, block_N)) if is_causal else T.ceildiv(seq_len, block_N))
115327
                    for k in T.Pipelined(loop_range, num_stages=num_stages):
                         if BlockMask[bz, bx, by, k]:
    T.copy(K[bz, k * block_N: (k + 1) * block_N, by, :], K_shared)
115429
1155<sup>30</sup>
                              if is causal:
                                   for i, j in T.Parallel(block_M, block_N):
115632
                                        acc_s[i, j] = T.if_then_else(bx * block_M + i >= k * block_N + j, 0,
1157<sup>33</sup><sub>34</sub>
                                                                          -T.infinity(acc_s.dtype))
                              else:
115835
                                   T.clear(acc_s)
1159<sup>36</sup><sub>37</sub>
                              T.gemm(Q_shared, K_shared, acc_s, transpose_B=True, policy=T.GemmWarpPolicy.FullRow)
                              T.copy(scores_max, scores_max_prev)
116038
                              T.fill(scores_max, -T.infinity(accum_dtype))
1161<sup>39</sup><sub>40</sub>
                              T.reduce_max(acc_s, scores_max, dim=1, clear=False)
                             for i in T.Parallel(block_M):
116241
                             scores_scale[i] = T.exp2(scores_max_prev[i] * scale - scores_max[i] * scale)
for i, j in T.Parallel(block_M, block_N):
1163_{43}^{42}
                              acc_s[i, j] = T.exp2(acc_s[i, j] * scale - scores_max[i] * scale)
T.reduce_sum(acc_s, scores_sum, dim=1)
116444
                              for i in T.Parallel(block_M):
116546
                                   logsum[i] = logsum[i] * scores_scale[i] + scores_sum[i]
                              T.copy(acc_s, acc_s_cast)
for i, j in T.Parallel(block_M, dim):
1166<sup>47</sup>
48
116749
                                   acc_o[i, j] *= scores_scale[i]
1168<sup>50</sup>
                              T.copy(V[bz, k * block_N: (k + 1) * block_N, by, :], V_shared)
                              T.gemm(acc_s_cast, V_shared, acc_o, policy=T.GemmWarpPolicy.FullRow)
116952
                    for i, j in T.Parallel(block_M, dim):
1170_{54}^{53}
                         acc_o[i, j] /= logsum[i]
                    T.copy(acc_o, O_shared)
117155
                    T.copy(O_shared, Output[bz, bx * block_M:(bx + 1) * block_M, by, :])
1172
```

Figure 14: Implementation of Block Sparse Flash Attention with TILELANG.