

---

# Hawkeye: Hardware-Aware GPU Kernel Optimization with Minimal Supervision

---

Anonymous Authors<sup>1</sup>

## Abstract

Achieving peak GPU kernel performance increasingly relies on architecture-specific optimizations targeting new hardware features. While AI coding agents show promise in generating performant kernels, they lack the necessary context to effectively implement and stack hardware-specific optimizations, especially on newer GPU architectures. We propose HAWKEYE (**H**ardware-**A**ware **K**ernel Optimization), an open-source framework that grounds autonomous kernel generation in a minimal and comprehensive taxonomy with only one unit test per optimization strategy per target architecture. Supporting a new accelerator therefore requires only 10 expert-written unit tests per architecture (one per recurring optimization strategy) that generalize across downstream workloads, rather than hand writing a new kernel for each workload and precision. HAWKEYE effectively scales the test-time compute of coding agents with this minimal expert supervision to enable kernel generation that consistently leverages hardware-specific features, approaching and even surpassing expert-written PyTorch or Triton in BF16 and emerging low precision (FP8, NVFP4, MXFP4) across Ampere, Hopper, Blackwell, and MI350 GPUs. HAWKEYE demonstrates that minimally supervised coding agents can exploit architecture-specific hardware features and reduce the overhead of supporting emerging hardware accelerators.

## 1 Introduction

Modern machine learning systems depend on GPU kernels carefully co-designed with the underlying architecture (Dao et al., 2022). As GPUs evolve, the number of architecture-specific instruction families grows substantially (Figure 1a) (Yadav et al., 2025). The underlying programming models have also become increasingly hardware-

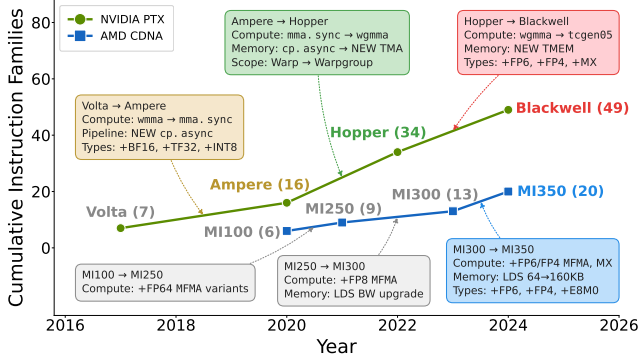
specific, with new tensor core instructions, functional units, and warp specialization strategies introduced with each generation. Tensor core paths provide roughly 94% of Hopper FLOPs (NVIDIA, 2022) and over 99% of Blackwell FLOPs (NVIDIA, 2024a), making hardware co-design essential for peak kernel performance.

The cost and complexity of supporting new hardware can be seen in real open-source projects. Architecture-specific code in NVIDIA’s CUTLASS library has grown 29× across recent versions and is now the majority of the codebase (per-version breakdown in Appendix A.6). Manually porting projects like FlashAttention (Dao et al., 2022), FlashInfer (Ye et al., 2025) and ThunderKittens (Spector et al., 2024) across GPU generations or vendors takes months to years of expert effort (porting timelines from Git history in Appendix A.6). The alternative of using high-level DSLs like Triton (Tillet et al., 2019) takes months to support new hardware (Appendix A.6) and still struggles with new hardware features (Issues, 2025). Hardware-specific kernels that achieve optimal performance on one hardware generation do not automatically achieve optimal performance on the next. Closing this gap, or providing *day-zero porting support*, means retargeting high-performance kernels to a newly released chip in hours rather than months while maintaining peak utilization on the new hardware.

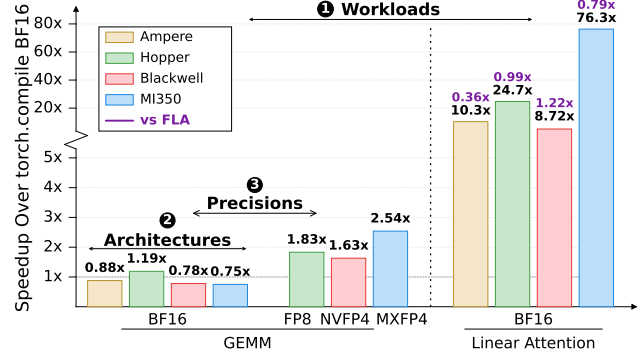
Coding agents accelerate these manual porting timelines with test-time feedback loops (Ouyang et al., 2025; Lange et al., 2025), reinforcement learning from runtime rewards (Baronio et al., 2025; Dai et al., 2026), and retrieval over CUDA documentation or hand-authored skills (Dong et al., 2026; Hugging Face, 2026). However, these approaches lack the context structure to enable architecture-specific kernel optimizations. In practice, kernel engineers add raw hardware context like programming model documentation (NVIDIA, 2024b; AMD, 2024) or kernels from older generations, which can pollute context without preventing common hardware-awareness failures (Section 4.4). For example, frontier coding agents can rarely implement Blackwell-specific optimizations such as `tcgen05` tensor core instructions or TMEM accumulators. Instead, they often fall back to Hopper-era `wgmma` code paths, which cannot fully utilize Blackwell’s new hardware and leave much

---

<sup>1</sup>Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.



(a) The number of unique assembly instruction families grows with each new GPU generation across NVIDIA and AMD product lines, as architects expose more hardware-specific optimizations to the programmer.



(b) HAWKEYE shows hardware awareness from unit tests enables agents to port kernels across workloads, architectures and vendors, and 16, 8, and 4-bit precision (vs. torch.compile and Triton FLA BF16).

Figure 1. Each GPU generation introduces new hardware-specific instruction families that widen the gap between peak hardware capability and what software actually delivers. HAWKEYE brings hardware awareness to coding agents through a minimal taxonomy of expert-authored unit tests, enabling autonomous kernel porting across workloads, GPU architectures and vendors, and precisions.

of the chip idle (Appendix B).

We define this missing capability as *hardware awareness*, the ability to identify hardware-specific optimizations and correctly implement kernels that utilize them. Without it, agents fall into predictable failure modes by generating hardware-unaware kernels that reuse optimization patterns from the wrong GPU generation, misconfigure architecture-specific data layouts, and abandon specialized low-precision paths for naive code that leaves tensor cores idle (Appendix D.2). These failures compound on newer and non-NVIDIA architectures (Appendix B), exactly where day-zero porting support is most valuable. This motivates our central question: *How can we make coding agents hardware-aware with minimal expert intervention?*

HAWKEYE introduces a generalizable (Section 4.1), minimal (Section 4.4), and comprehensive (Appendix E.3) taxonomy of unit tests that enables coding agents to scale test-time compute more effectively (Section 4.3) and generate hardware-aware kernels. Each unit test pairs an optimization strategy with a profiling metric, and the agent composes them into high-performance kernels for downstream workloads. In contrast with baselines of hardware-aware agent context, HAWKEYE’s taxonomy has better context efficiency (Section 4.4) by organizing rows of recurring optimization strategies (MMA, async pipelines, shared-memory swizzling, and so on) and columns of target architectures, with one expert-authored kernel per cell. Adding a new accelerator therefore costs only 10 unit tests, regardless of the number of downstream workloads the agent later optimizes.

We evaluate HAWKEYE on porting PyTorch workloads to high-performance kernels across NVIDIA Ampere, Hopper, Blackwell, and AMD MI350, and across BF16, FP8, NVFP4, and MXFP4 precisions (Section 4). On established workloads, where torch.compile dispatches to expert-

tuned vendor libraries like cuBLAS and cuDNN, HAWKEYE matches or exceeds it in both BF16 and low precision, including in formats PyTorch cannot natively run. On emerging attention variants where torch.compile cannot fuse non-standard scans and gates, HAWKEYE reaches an 18.9× geomean speedup. Against expert-authored Triton kernels from the Flash Linear Attention library (Yang and Zhang, 2024), HAWKEYE approaches or exceeds FLA on Linear Attention across every architecture, including 1.22× on Blackwell and 1.00× on MI350.

In summary, our contributions are as follows.

- 1. An extensible taxonomy of unit tests for hardware-aware kernel optimization.** The taxonomy is a minimal, comprehensive, and generalizable expert input needed to create hardware-aware kernels for a target accelerator, requiring only 10 unit tests per architecture (one per recurring optimization strategy) that generalize to downstream workloads.
- 2. Scaling test-time compute to generate hardware-aware kernels from minimal supervision.** HAWKEYE is an open-source framework that scales coding agents to compose the per-architecture taxonomy unit tests into performant GPU kernels with better context efficiency than baselines. We release the taxonomy, kernels, and experimental harness for all four GPU architectures.
- 3. Enabling kernel porting across architectures, vendors, and precisions.** Across NVIDIA Ampere, Hopper, Blackwell, and AMD MI350 in BF16, FP8, NVFP4, and MXFP4, HAWKEYE reaches 1.13× BF16 and 1.28× low-precision geomean speedup against torch.compile on established workloads, beats it by 18.9× geomean on emerging attention variants, and approaches or exceeds expert-written Triton kernels.

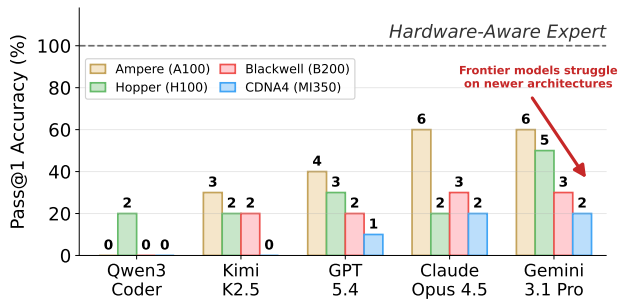


Figure 2. Pass@1 of frontier LLMs on the taxonomy unit tests. Accuracy drops on newer GPU architectures, indicating a lack of architecture-specific domain knowledge.

## 2 HAWKEYE

HAWKEYE is a framework that bootstraps coding agents with a taxonomy of architecture-specific unit tests and enables agents to generate hardware-aware kernels for a diverse set of workloads. The methodology walks through the steps of a kernel engineer automating day-zero kernel optimization (Lin, 2025; Liao et al., 2025), namely the alternative frameworks that lead to coding agents (Section 2.1) and the context for these agents (Section 2.2). The progression of context structures motivates the HAWKEYE unit test taxonomy (Section 2.3) and kernel generation (Section 2.4).

### 2.1 Context Frameworks for Kernel Optimization

Frontier coding agents lack hardware awareness on day-zero GPU architectures, and the gap widens with each new generation (Figure 2, Appendix B.1). Recent work has explored several context-engineering frameworks for kernel optimization (evaluated in Table 3). Iterative refinement (Ouyang et al., 2025) uses hardware specifications and runtime feedback over candidate kernels. GEPA-style prompt evolution (Agrawal et al., 2026) self-maintains a markdown reflection across attempts. Agent skills (Hugging Face, 2026) curate static programming model guides for specific optimizations. None inject architecture-specific knowledge that the agent can reliably retrieve and compose, since iterative feedback adds no external knowledge between trajectories, self-rewriting reflections cannot contribute new knowledge, and skills lack extensibility and detailed implementation guidance (Appendix D.4). The open design question is what hardware-aware content the agent context should hold.

### 2.2 A Progression of Hardware-Aware Context

On day-zero of a new architecture, kernel engineers choose context to put in the agent’s context window. We walk through a progression of prior approaches, ordered from portable algorithmic representations to curated architecture-specific instruction wrappers. Each exposes a structural deficit that HAWKEYE addresses (Section 4.4).

Table 1. Hardware-aware context structures against HAWKEYE. Turn-0 tokens are in the agent prompt. Workspace tokens are on disk and accessed through read tools.

Context structure	Description	Turn-0	Workspace
Triton FLA	DSL kernel	3.4k	13.6k
Raw docs	ISA reference	2.1k	1.22M
Augmented docs	Curated ISA index	5k	1.22M
Hw-aware kernel	Native example	26k	104k
Vendor FMHA	Production kernel	2.3k	287k
PTX abstractions	PTX functions	0.9k	17k
<b>HAWKEYE</b>	<b>Unit tests + guides</b>	<b>11k</b>	<b>43k</b>

**Algorithm without hardware syntax.** High-level GPU DSLs like Triton are designed for portability by expressing kernels through abstractions that the compiler lowers for each target. Using DSL kernels as context, for example the FLA library (Yang and Zhang, 2024) at 3.4k tokens (Table 1, token counts via (tiktoken, 2022)), gives an agent the algorithm but abstracts away the native CUDA or HIP instructions (Appendix D.3). *The agent has the algorithm but no architecture-specific syntax to implement it, and falls back to older generation tensor-core paths.*

**Hardware documentation.** An alternative provides the chip’s full programming model documentation (NVIDIA, 2024b; AMD, 2024), with a 2k-token index and complete reference (1.22M tokens on NVIDIA, 256k on AMD) on disk. Every primitive is present, but the corpus pollutes the context window and retrieval anchors in single chapters (Appendix D.4). Curating the documentation into a 5–7k-token index of structured excerpts removes the navigation cost, but the depth required still floods the context with unusable instruction semantics (Appendix D.5). *The agent finds the right primitives but cannot compose them, because there is no feedback on optimization effectiveness and the cross-primitive recipe is split across separate sections.*

**Few-shot hardware-aware kernels.** A more involved approach replaces verbose documentation with hand-written kernel examples. A 26–35k-token GEMM kernel in context demonstrates MMA and asynchronous data movement primitives, but entangles algorithmic shape with descriptor strides and scheduler choices. Porting to new algorithms requires changing these optimizations in coordination (Appendix D.6). Production kernels like CUTLASS (NVIDIA, 2025) and Composable Kernel (Liu et al.) fused multi-head attention stack every optimization at once across 287–854k tokens, but bury the primitive optimization recipes under many layers of abstractions (Appendix D.7). *The agent has a working kernel composition but cannot decompose it because it has no granular abstraction boundary to get usable insights into individual optimizations.*

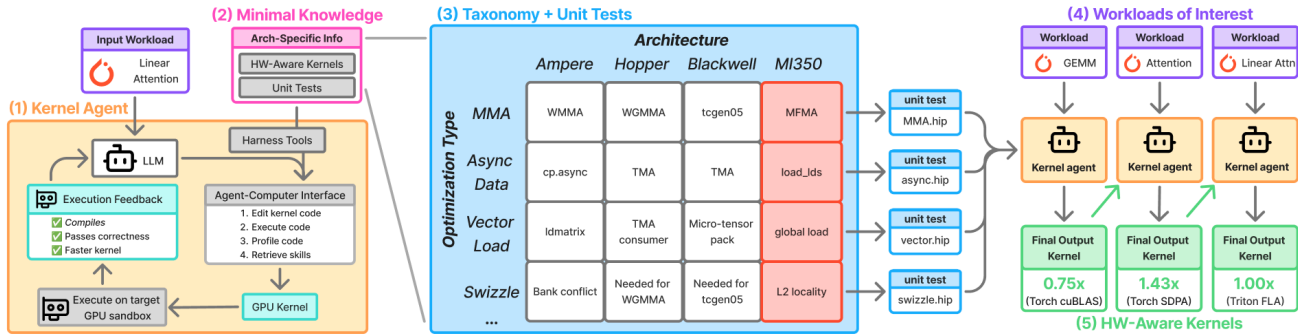


Figure 3. HAWKEYE framework. A kernel agent writes an optimized GPU kernel for a PyTorch workload, iteratively editing, executing, and profiling in a GPU sandbox. The agent is bootstrapped with the minimal knowledge needed for hardware awareness, consisting of 10 unit tests grounded in the taxonomy and previously generated kernels on the same architecture. Generalizing to different workloads of interest produces hardware-aware kernels evaluated against expert baselines.

**Curated abstractions for important optimizations.** An engineer can extract PTX-level abstractions from these production kernels and put these minimal API wrappers in context. At under 1k tokens, the agent receives correct syntax for every primitive, but composing them into a coherent kernel requires reasoning about their interactions, which the wrappers alone do not encode (Appendix D.8). The agent implements each optimization correctly but cannot stack them, since the wrappers encode syntax in isolation.

### 2.3 The Hawkeye Taxonomy of Unit Tests

We organize hardware-aware optimizations as a taxonomy whose rows are recurring optimization strategies and columns are target GPU architectures (Table 2). The optimization strategy rows mirror patterns from production projects like FlashAttention-2 (Dao et al., 2022), ThunderKittens (Spector et al., 2024), and HipKittens (Hu et al., 2025) (Appendices A.4 and A.5). The patterns that recur across workloads and architectures form ten rows that span compute, data movement, and scheduling. The columns in this paper are NVIDIA Ampere, Hopper, Blackwell, and AMD MI350. Each cell captures the concrete syntax, memory semantics, and thread scope required to realize the row’s strategy on the column’s architecture (Appendix A).

Adding a new accelerator to the taxonomy costs only one column of unit tests, bounding the per-architecture human effort to  $O(\text{optimization strategies})$  rather than  $O(\text{workloads} \times \text{precisions})$ . Each unit test is the minimal abstraction that pairs a human-authored solution kernel with the profiling metric that verifies the optimization. Each kernel is wrapped as a callable function with a usage guide, so the agent can use it as a syntax example or compose fragments into a larger kernel. This matters for instructions postdating frontier-model pretraining, where the unit test is the only signal that the feature exists (Appendix C.2).

### 2.4 Hardware-Aware Kernel Generation

Figure 3 shows the kernel generation loop where the agent receives a PyTorch reference for the target workload, the

column of unit tests for the target architecture (Section 2.3), and previously composed hardware-aware kernels for other workloads on the same architecture. The agent iterates through a closed loop of editing the kernel, compiling for the target, executing in the sandbox, validating against the reference, and measuring performance and per-row profiling metrics. Over the course of many turns, the agent identifies the bottleneck the next optimization should address, consults the unit test that targets it, and stacks hardware-aware optimizations into one complete kernel. The per-row counter lets the agent verify that the targeted hardware feature fired before moving on, so the trajectory accumulates composed optimizations (Appendix E.1). Correct kernels that improve over the previous best are added to the workspace and condition the agent on subsequent workloads.

## 3 Experimental Setup

**Hardware and precisions.** We evaluate on NVIDIA A100-40GB (Ampere), H100 (Hopper), B200 (Blackwell), and AMD MI350 (CDNA4) using CUDA with inline PTX and HIP. Generated kernels are disallowed from using external libraries. For each architecture we run BF16 and the lowest natively-supported low-precision floating point format, namely FP8 (E4M3) on Hopper, NVFP4 on Blackwell, and MXFP4 (Rouhani et al., 2023) on MI350. Inputs are generated in the target precision to avoid dequantization overhead.

**Established workloads.** GEMM, Conv2D, and standard attention cover the compute-intensive subset of KernelBench (Ouyang et al., 2025). We compare against torch.compile with max-autotune, which dispatches to cuBLAS (NVIDIA, 2026), cuDNN (Chetlur et al., 2014), and FlashAttention. GEMM uses  $M=N=K=4096$  (KernelBench standard). Conv2D uses the first AlexNet (Krizhevsky et al., 2012) layer with  $N=256$ ,  $3 \rightarrow 96$  channels,  $224 \times 224$  input,  $11 \times 11$  kernel,  $S=4$ , and  $P=2$ . Attention uses the FlashAttention (Dao et al., 2022) configuration  $(B, H, S, D)=(4, 8, 4096, 64)$ .

Table 2. The HAWKEYE taxonomy is a minimal, comprehensive, and generalizable context structure for autonomous kernel optimization. Each row is a general optimization strategy and each unit test cell shows its architecture-specific details. Three representative rows are shown here, and the full taxonomy of 10 rows and profiling breakdown can be found in Appendix C.1.

	Ampere (A100)	Hopper (H100)	Blackwell (B200)	MI350
<b>MMA Unit</b>	wmma + loading ldmatrix/.trans	wgmma + TMA 128B-swizzle loads	tcgen05 + TMA bulk loads + TMEM accum.	MFMA 32x32x16, AGPR accumulators
<b>Shared Memory Layout</b>	XOR swizzle (s32/64/256) for conflict-free	128B column-group XOR swizzle matching wgmma descriptors	Parameterized 64B offset_swizzle for tcgen05 descriptors	64-bank LDS XOR swizzle for conflict-free
<b>Async Pipeline</b>	cp.async pipeline, commit_group	TMA bulk.tensor + mbarrier pipeline	TMA 2D + mbarrier pipeline, phase-parity	Triple-buffered LDS pipeline

**Emerging workloads.** ReBased Linear Attention (Aksenov et al., 2024), DeltaNet (Yang et al., 2024), and Forgetting Attention (Lin et al., 2025) offer architectural advantages over the standard transformer but lack vendor kernel support. We compare against autotuned expert Triton kernels in the Flash Linear Attention (FLA) library (Yang and Zhang, 2024) as a human upper bound. All attention variants are causal at  $(B, H, S, D)=(4, 8, 4096, 64)$ .

**Agent baselines.** All agents use Gemini 3.1 Pro in OpenHands (Wang et al., 2025) and have a budget of 100 *productive turns* of kernel edits, evaluations, and taxonomy or deduplicated kernel-pool reads. We compare HAWKEYE against (1) 10x10 iterative refinement loop (Ouyang et al., 2025), (2) a coding agent with the HuggingFace CUDA-kernels skill (Hugging Face, 2026) and programming guides, (3) GEPA (Agrawal et al., 2026) with a self-rewriting markdown reflection, and (4) a coding agent given programming-model docs or other kernels as context. Full prompt templates and workspace layouts are in Appendix G.

**Metrics.** Correctness uses an oracle in-precision naive kernel (Appendix G.3). Performance is in TFLOPS (workload’s theoretical FLOP count divided by latency from Triton `do_bench`) using  $2MNK$  for GEMM,  $im2col$  count for Conv2D,  $2BHS^2D$  for the causal attentions, and  $4BHS^2D$  for DeltaNet. Speedups are calculated as TFLOPS ratios against `torch.compile` in Table 3 and expert Triton FLA in Figure 4.

## 4 Results

We evaluate HAWKEYE on the established and emerging workloads in Table 3. Section 4.1 covers the porting axes across architectures, vendors, and precisions. Section 4.2 compares against expert-authored Triton kernels on emerging attention variants. Section 4.3 measures how kernel quality scales with the agent’s turn budget and why HAWKEYE trajectories continue to improve while alternatives plateau. Section 4.4 compares against alternative hardware-aware context structures and walks through why HAWKEYE overcomes their shortcomings.

### 4.1 Porting Across Architectures, Vendors, and Precisions

**Across GPU architectures (Ampere, Hopper, Blackwell, and MI350).** HAWKEYE uses the per-architecture unit test column as context to port the PyTorch reference across architectures. On Blackwell, the agent directly composes `tcgen05.mma` and TMEM accumulators, while translating a Hopper `wgmma`-based GEMM with static conversion tools bypasses these primitives entirely (Appendix A.1). On established workloads this gives  $0.88\times$  BF16 GEMM,  $1.69\times$  Conv2D, and  $0.98\times$  standard attention geomean against `torch.compile` across four chips. On emerging variants where `torch.compile` cannot fuse non-standard scans and gates, HAWKEYE reaches  $19.70\times$ ,  $26.47\times$ , and  $12.89\times$  geomean speedup on causal Linear Attention, DeltaNet, and Forgetting Attention.

**Across vendors (NVIDIA and AMD).** HAWKEYE enables cross-vendor porting that requires translation across programming models and core hardware primitives (Davis et al., 2025), and we compare to the autonomous NVIDIA to AMD source-level translation tool Hipify (AMD ROCm Team, 2024). On GEMM, Conv2D, and attention, Hipify MI350 kernels have a geomean  $0.05\times$  performance against `torch.compile` (Appendix A.2), while HAWKEYE uses the MI350 column of unit tests to reach  $1.26\times$  speedup.

**Across precisions (BF16, FP8, NVFP4, and MXFP4).** Roofline analysis (Appendix F.5) shows the low-precision HAWKEYE kernels increasing both the arithmetic intensity and achievable ceiling, lifting GEMM throughput by  $1.5\times$  on H100,  $2.1\times$  on B200, and  $3.4\times$  on MI350. Because the unit tests are specific to precision, translations are clean trajectories rather than retrofits of BF16 kernels. Across the established workloads in low precision, HAWKEYE reaches a  $1.28\times$  geomean against `torch.compile` in formats PyTorch cannot run natively. Some low precision kernels show a slowdown versus BF16 because they introduce intermediate quantization overhead to pass correctness.

Table 3. HAWKEYE reaches speedups against torch.compile on established and emerging workloads results across four GPU architectures and four precisions, with the largest gains on newer generations and lower precisions. X denotes incorrect or runtime errors.

Method	Workload	Ampere (A100)		Hopper (H100)				Blackwell (B200)				MI350				Geomean	
		BF16		BF16		FP8		BF16		NVFP4		BF16		MXFP4		BF16	Low Prec.
		TFLOPS	Spdup	TFLOPS	Spdup	TFLOPS	Spdup	TFLOPS	Spdup	TFLOPS	Spdup	TFLOPS	Spdup	TFLOPS	Spdup	Spdup	Spdup
10×10 Iterative Refinement	GEMM	96.55	0.49x	X	X	X	X	X	X	10.53	0.01x	359.71	0.50x	26.87	0.01x	0.07x	0.01x
	Conv2d	9.57	0.57x	13.85	0.58x	13.95	0.58x	3.64	0.11x	3.68	0.12x	28.24	0.61x	7.20	0.15x	0.39x	0.22x
	Attention	2.89	0.04x	0.33	0.00x	X	X	3.33	0.01x	X	X	0.74	0.00x	2.13	0.01x	0.01x	0.01x
	Lin. Attn.	X	X	X	X	0.12	0.02x	X	X	X	X	X	X	X	X	0.01x	0.01x
	Delta Net	X	X	X	X	0.66	11.63x	X	X	0.50	13.06x	X	X	X	X	0.01x	1.15x
	FoX	6.18	2.02x	15.33	2.37x	X	X	2.45	0.13x	X	X	2.02	0.37x	11.01	2.00x	0.69x	0.06x
Kernel Optimization Agent Skill	GEMM	147.40	0.75x	701.00	<b>1.26x</b>	X	X	X	X	16.95	0.02x	201.11	0.28x	236.82	0.11x	0.23x	0.03x
	Conv2d	17.20	1.02x	3.23	0.13x	X	X	60.10	1.82x	5.03	0.16x	81.92	1.76x	73.89	1.49x	0.81x	0.13x
	Attention	X	X	X	X	0.06	0.00x	0.13	0.00x	X	X	X	X	X	X	0.01x	0.01x
	Lin. Attn.	22.51	6.15x	12.40	1.79x	X	X	3.21	0.10x	X	X	X	X	X	X	0.32x	0.01x
	Delta Net	X	X	X	X	X	X	X	X	0.01	0.26x	X	X	X	X	0.01x	0.03x
	FoX	27.20	8.91x	0.38	0.06x	X	X	X	X	X	X	0.01	0.00x	0.01	0.00x	0.09x	0.01x
GEPA Reflective Prompt Evolution	GEMM	161.40	0.82x	698.10	1.26x	4.70	0.00x	X	X	4.93	0.01x	291.52	0.40x	187.65	0.09x	0.25x	0.02x
	Conv2d	16.50	0.98x	3.50	0.15x	X	X	X	X	3.63	0.12x	19.85	0.43x	67.53	1.37x	0.16x	0.12x
	Attention	22.20	0.34x	X	X	222.40	<b>1.50x</b>	X	X	X	X	X	X	X	X	0.02x	0.05x
	Lin. Attn.	15.90	4.34x	X	X	X	X	X	X	X	X	X	X	X	X	0.05x	0.01x
	Delta Net	X	X	X	X	X	X	X	X	0.30	7.84x	X	X	X	X	0.01x	0.09x
	FoX	39.90	13.07x	57.75	8.95x	90.20	13.88x	0.01	0.00x	2.39	0.13x	37.75	6.85x	X	X	1.68x	0.26x
HAWKEYE	GEMM	173.79	<b>0.88x</b>	660.40	1.19x	1018.96	<b>0.78x</b>	760.66	<b>0.78x</b>	1577.64	<b>1.61x</b>	544.96	<b>0.75x</b>	1844.80	<b>0.88x</b>	<b>0.88x</b>	<b>1.03x</b>
	Conv2d	20.82	<b>1.24x</b>	38.57	<b>1.61x</b>	47.44	<b>1.97x</b>	71.78	<b>2.17x</b>	62.66	<b>2.04x</b>	87.53	<b>1.88x</b>	125.60	<b>2.54x</b>	<b>1.69x</b>	<b>2.17x</b>
	Attention	116.38	<b>0.88x</b>	337.34	<b>1.13x</b>	387.46	1.31x	455.44	<b>0.66x</b>	253.15	<b>0.54x</b>	372.47	<b>1.43x</b>	307.90	<b>1.17x</b>	<b>0.98x</b>	<b>0.94x</b>
	Lin. Attn.	75.79	<b>10.35x</b>	342.30	<b>24.70x</b>	346.07	<b>25.20x</b>	548.54	<b>8.72x</b>	378.26	<b>3.03x</b>	326.46	<b>67.57x</b>	368.59	<b>76.29x</b>	<b>19.70x</b>	<b>17.99x</b>
	Delta Net	1.91	<b>33.50x</b>	2.51	<b>25.50x</b>	1.77	<b>31.20x</b>	2.33	<b>36.20x</b>	1.24	<b>32.40x</b>	3.90	<b>15.88x</b>	3.60	<b>15.43x</b>	<b>26.47x</b>	<b>24.99x</b>
	FoX	104.13	<b>17.05x</b>	245.49	<b>19.02x</b>	246.17	<b>18.94x</b>	440.20	<b>11.89x</b>	287.47	<b>7.75x</b>	39.47	<b>7.16x</b>	49.80	<b>9.05x</b>	<b>12.89x</b>	<b>10.99x</b>

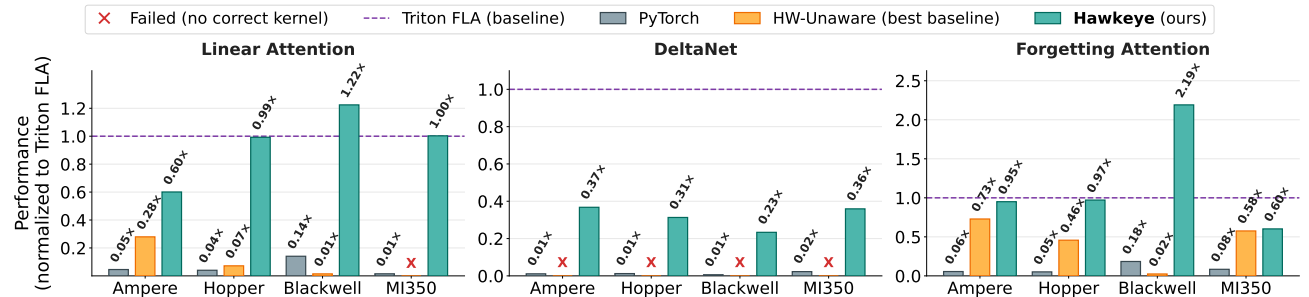


Figure 4. On attention variants, HAWKEYE approaches or exceeds performance against PyTorch, the best baseline from Table 3, and expert-written Triton FLA kernels (Yang and Zhang, 2024) (dotted line). HAWKEYE achieves 0.92× geomean against FLA on causal Linear Attention and 2.19× speedup on Blackwell Forgetting Attention. There is still headroom on the expert DeltaNet kernel.

### 4.2 Comparing Against Expert Triton Kernels

HAWKEYE approaches or exceeds Triton kernels in the Flash Linear Attention library (Yang and Zhang, 2024), which is the closest proxy to an expert upper bound. HAWKEYE reaches geomeans of 0.92×, 0.31×, and 1.05× FLA on causal Linear Attention, DeltaNet, and Forgetting Attention (Figure 4), including 1.22× and 1.00× on Linear Attention on Blackwell and MI350 and 2.19× on Forgetting Attention on Blackwell where the Triton backend is still maturing (Ren et al., 2026). Approaching this performance autonomously, in minutes per kernel rather than days of expert effort, is a substantial productivity win over hand-

written kernels. On DeltaNet, trajectory analysis shows that the agent attempts the full hardware-aware optimization stack but struggles to compose unit tests into a complete kernel, and falls back to simpler implementations to pass correctness (Appendix E.2). While the taxonomy is complete, closing this gap likely requires exposing more insight to the agent on composing unit tests into optimized kernels.

### 4.3 Scaling Performance with Test-Time Compute

Figure 5 shows that scaling test-time compute with HAWKEYE generates the most performant kernels across architectures. In observed trajectories, early turns retrieve unit

330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383  
384

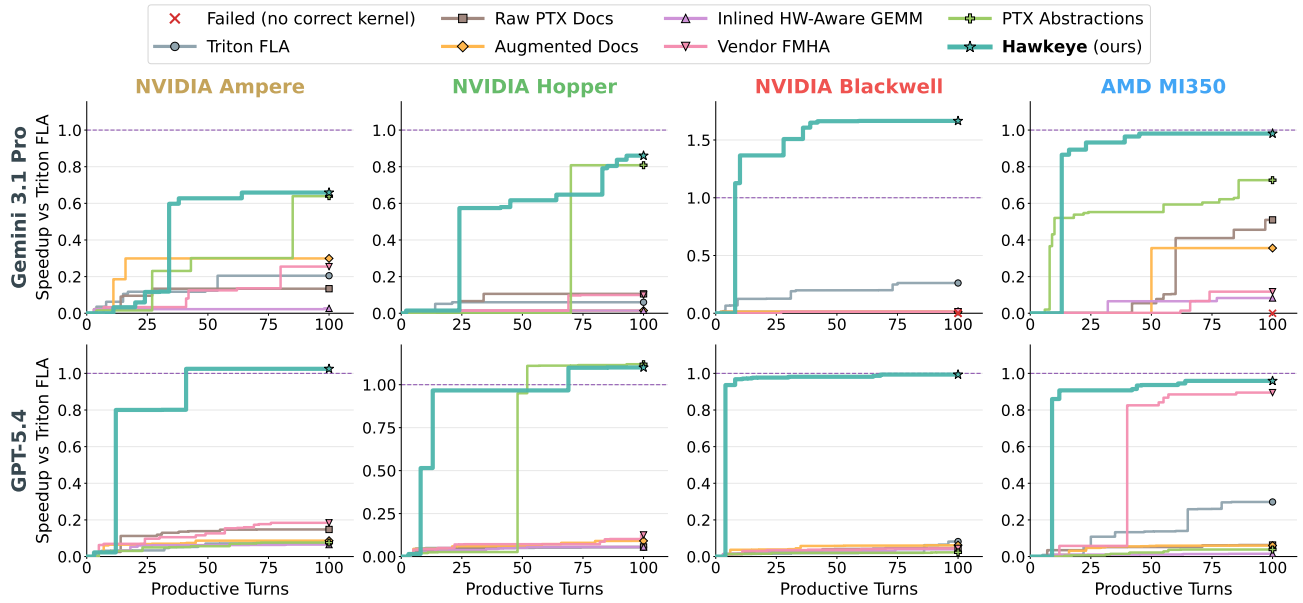


Figure 5. Scaling agent test-time compute with HAWKEYE yields strong speedups on causal Linear Attention with Gemini 3.1 Pro and GPT-5.4 against baseline context structures, all measured against expert-written Triton FLA (dashed line). HAWKEYE is the only method that keeps producing faster kernels as productive turns grow (kernel edits, evaluations, and abstraction reads, see Section 3), reaching  $0.78\times-1.65\times$  FLA across architectures while baselines plateau at lower speedups.

tests and plan stacking strategies (Appendix E.2). The jump arrives when optimizations can be composed into a coherent kernel. The depth of the required stack explains the curve shapes across architectures. On MI350 the agent reaches  $0.90\times$  expert in 23 productive turns because MFMA, LDS, and wave-occupancy compose into a thin stack. On Blackwell the agent needs roughly 50 productive turns to reach  $5\times$  expert because TMEM accumulators, ring-buffered TMA, cluster-multicast loads, and a phase-bit protocol must all converge. Past the transition, additional turns invent cross-optimization patterns no unit test describes, like ping-pong TMEM and a reversed causal chunk schedule that close the gap on Blackwell (Appendix E.3).

The unit test profiling metrics enable the MI350 trajectory to debug tightly coupled individual optimizations and avoid regression from partial composition (Figure 6). On Hopper, a producer-consumer split without the ring buffer regresses throughput by half, and the per-optimization signal lets the agent revert (Appendix E.1). Baseline context structures plateau because each fails to implement hardware-aware optimizations in distinct ways (Appendix D.2).

#### 4.4 Evaluating Context Efficiency Against Other Hardware-Aware Context Structures

A context structure is efficient for GPU kernel optimization when it is minimal enough to avoid context pollution and detailed enough to enable optimization stacking. We evaluate six baselines along this spectrum (Table 1), and HAWKEYE addresses each baseline’s unique failure modes.

**Comparing to algorithm without hardware syntax.** DSL-level conditioning shows an effective algorithm but no path to the architecture’s tensor-core primitives. On Blackwell, the trajectory plateaus at 220.8 TFLOPS ( $0.26\times$  expert) by falling back to Ampere-era MMA syntax without reaching `tcgen05` or its TMEM accumulator. Ampere correctly uses `wmma` syntax, but gets to  $0.27\times$  because it cannot orchestrate asynchronous data movement (Appendix D.3). HAWKEYE closes this gap because each unit test shows the architecture-specific syntax for optimizations, so the agent never has to fall back on syntax from older generations.

**Comparing to hardware documentation.** Raw and augmented documentation show that the right primitives are present in the corpus but cannot be composed. The Blackwell raw-docs trajectory stops at 5.6 TFLOPS ( $0.013\times$  expert) because retrieval anchors in the legacy `wgmma` chapter and never reaches `tcgen05`. Curated documentation removes the navigation cost, but Hopper never beats the naive baseline because information across primitives like `wgmma`, TMA, and `mbarrier` is split across excerpts (Appendix D.5). HAWKEYE closes this gap because each unit test shows a minimal but complete set of optimizations, so the agent easily retrieves syntax and composition rather than reassembling them across documentation excerpts.

**Comparing to few-shot hardware-aware kernels.** Inlined kernel examples and vendor production references show that finished compositions cannot be decomposed into the primitives they are built up from. An inlined hardware-aware

<pre> 385 // Q: reload from gmem each iter 386 load_to_LDS(Q[q_start..+64]); 387 // GEMM1: S = Q @ K^T 388 // 3 of 4 waves idle in MFMA 389 if (warp_id == 0) 390     for (k, r, c) 391         s[r*4+c] = mfma_16x16x16( 392             scalar_gather(LDS_Q), 393             scalar_gather(LDS_K), 394             s[r*4+c]); 395 // causal: per-elem branch 396 // fires on every tile 397 for_each_elem (r,c,e): 398     if (col &gt; row) v = 0; </pre>	<pre> 385 // Q: load once, re-read from LDS 386 load_to_LDS(Q); 387 a = LDS_Q[row*16 + lane]; 388 // GEMM1: S = Q @ K^T 389 // all 4 waves cooperate, but 390 // small tile -&gt; 4 outer steps 391 for (k = 0; k &lt; 4; k++) 392     for (r, c in 2x2) 393         C[r][c] = mfma_16x16x16( 394             a, b, C[r][c]); 395 // causal: per-elem branch 396 // fires on every tile 397 for_each_elem: 398     if (col &gt; row) v = 0; </pre>	<pre> 385 // Q: load once, kept in VGPRs 386 q_regs = load_Q_to_REGS(); 387 // GEMM1: S = Q @ K^T 388 // 4 waves, 4x larger MFMA tile 389 for (ki = 0; ki &lt; D; ki += 16) 390     s = mfma_32x32x16( 391         q_regs[ki/16], 392         LDS_K[cur][ki], s); 393 // causal: tile-level fast path, 394 // only diagonal tile pays 395 if (kv_start + BN &lt;= q_start) 396     write_W(); 397 else 398     write_W_causal(); </pre>
---	--	---

(a) Vendor FMHA (39 TFLOPS)

(b) Raw PTX Docs (170 TFLOPS)

(c) HAWKEYE (326 TFLOPS)

Figure 6. HAWKEYE wins by composing hardware-native primitives end-to-end, while baselines stop part way through the stack. Inner KV-tile loop of three causal Linear Attention kernels (Gemini 3.1 Pro, MI350, BF16). The Vendor FMHA baseline only utilizes a quarter of the matrix cores, and the Raw PTX Docs baseline stays on the smaller MFMA shape and reloads operands through shared memory. HAWKEYE stacks all three optimizations and reaches 98% of the hand-tuned expert.

GEMM exposes the right tensor-core and asynchronous-copy primitives but entangles them with the source workload’s tile geometry and scheduling choices (Appendix D.6). Vendor production kernels like CUTLASS and Composable Kernel are written for human reuse and their abstractions hide the primitive recipes the agent needs to extract (Appendix D.7). HAWKEYE closes this gap because the taxonomy shows composable optimization boundaries, so the agent reuses an abstraction the engineer authored rather than inferring it from a finished kernel.

**Comparing to curated abstractions for important optimizations.** The PTX-abstractions show that primitives alone do not enable multi-optimization composition. Without the HAWKEYE unit tests and corresponding profiling metrics, the Hopper kernel only reaches 240.9 TFLOPS with annotations where the agent recognizes the producer-consumer split as the next move but cannot implement it. The MI350 control case shows a 35% gap to HAWKEYE with Gemini and a 28x gap with GPT-5.4 (Appendix D.8). HAWKEYE closes this gap because the unit tests guide optimization stacking that primitives alone cannot encode to turn correct syntax into composed kernels.

## 5 Related Work

**Training-based methods.** KevinRL (Baronio et al., 2025) and CudaAgent (Dai et al., 2026) fine-tune coding agents on runtime rewards, but rely on the target architectures and precisions being in the training distribution. HAWKEYE does not require this and instead exposes architecture-specific optimizations as a structured taxonomy of expert-authored unit tests, inspectable and extensible to new accelerators.

**Search-based methods.** K-Search (Cao et al., 2026) and AI CUDA Engineer (Lange et al., 2025) explore the kernel space with a world model and evolutionary discovery, while CudaForge (Zhang et al., 2025) and KernelEvolve (Liao et al., 2025) augment search with runtime feedback or re-

trieval. Without per-architecture grounding, these methods struggle to utilize the required primitives. HAWKEYE grounds test-time search in architecture-specific unit tests so the agent reaches the right primitives directly.

**Context engineering methods.** KernelBlaster (Dong et al., 2026) retrieves cross-task examples into a continual memory, while SwizzlePerf (Tschand et al., 2025) and Hugging-Face CUDA skills (Hugging Face, 2026) ship hand-written context for specific optimizations. These inject hardware-aware context but are limited to a static set of optimizations or architectures. HAWKEYE’s taxonomy is structured per architecture, scales to downstream workloads, and extends with one new column of unit tests (comparison in Table 5).

## 6 Discussion

**Limitations.** HAWKEYE still requires an expert to author each new accelerator’s column of unit tests. Low-precision kernels depend on quantization tolerance (Xiao et al., 2023), and attention variants often force mixed precision due to compounding loss. On heavily vendor-tuned workloads, the agent does not always reach hand-written performance because composing isolated unit tests into complete kernels remains hard. We view this as an opportunity for future work in autonomous kernel composition now that agents can implement hardware-specific optimizations.

**Conclusion.** HAWKEYE is a taxonomy of expert-authored domain knowledge that uses only 10 unit tests to ground coding agents in the architecture’s programming model. Scaling test-time compute against this minimal supervision, HAWKEYE generates hardware-aware kernels that match or exceed expert-written PyTorch and Triton in BF16 and low precision across NVIDIA Ampere, Hopper, Blackwell, and AMD MI350. We release the taxonomy, kernels, and evaluation harness, and hope this shifts the human’s role from writing kernels to architecting structured knowledge that lets coding agents scale autonomous kernel optimization.

## References

- Lakshya A Agrawal, Shangyin Tan, Dilara Soyulu, Noah Ziem, Rishi Khare, Krista Opsahl-Ong, Arnav Singhvi, Herumb Shandilya, Michael J Ryan, Meng Jiang, Christopher Potts, Koushik Sen, Alexandros G. Dimakis, Ion Stoica, Dan Klein, Matei Zaharia, and Omar Khattab. Gepa: Reflective prompt evolution can outperform reinforcement learning, 2026. URL <https://arxiv.org/abs/2507.19457>.
- Yaroslav Aksenov, Nikita Balagansky, Sofia Maria Lo Cicero Vaina, Boris Shaposhnikov, Alexey Gorbatovski, and Daniil Gavrilov. Linear transformers with learnable kernel functions are better in-context models, 2024. URL <https://arxiv.org/abs/2402.10644>.
- AMD. *AMD Instinct MI-Series ISA Reference Guides (CDNA1-CDNA4)*, 2024. URL <https://developer.amd.com/>. Instruction families counted from MFMA and related entries in each CDNA architecture ISA document.
- AMD ROCm Team. Hipify: Convert cuda to portable hip c++ code. <https://github.com/ROCm/HIPIFY>, 2024. Accessed: 2026-03-30.
- Carlo Baronio, Pietro Marsella, Ben Pan, Simon Guo, and Silas Alberti. Kevin: Multi-turn rl for generating cuda kernels, 2025. URL <https://arxiv.org/abs/2507.11948>.
- Shiyi Cao, Ziming Mao, Joseph E. Gonzalez, and Ion Stoica. K-search: Llm kernel generation via co-evolving intrinsic world model, 2026. URL <https://arxiv.org/abs/2602.19128>.
- Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cudnn: Efficient primitives for deep learning, 2014. URL <https://arxiv.org/abs/1410.0759>.
- Weinan Dai, Hanlin Wu, Qiyang Yu, Huan-ang Gao, Jiahao Li, Chengquan Jiang, Weiqiang Lou, Yufan Song, Hongli Yu, Jiaze Chen, et al. Cuda agent: Large-scale agentic rl for high-performance cuda kernel generation. *arXiv preprint arXiv:2602.24286*, 2026.
- Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. *Advances in neural information processing systems*, 35:16344–16359, 2022.
- Joshua Hoke Davis, Pranav Sivaraman, Joy Kitson, Konstantinos Parasyris, Harshitha Menon, Isaac Minn, Giorgis Georgakoudis, and Abhinav Bhatele. Taking gpu programming models to task for performance portability. In *Proceedings of the 39th ACM International Conference on Supercomputing*, pages 776–791, 2025.
- Kris Shengjun Dong, Sahil Modi, Dima Nikiforov, Sana Damani, Edward Lin, Siva Kumar Sastry Hari, and Christos Kozyrakis. Kernelblaster: Continual cross-task cuda optimization via memory-augmented in-context reinforcement learning, 2026. URL <https://arxiv.org/abs/2602.14293>.
- William Hu, Drew Wadsworth, Sean Siddens, Stanley Winata, Daniel Y Fu, Ryann Swann, Muhammad Osama, Christopher Ré, and Simran Arora. Hipkittens: Fast and furious amd kernels. *arXiv preprint arXiv:2511.08083*, 2025.
- Hugging Face. Custom kernels for all from codex and claude. <https://huggingface.co/blog/custom-cuda-kernels-agent-skills>, 2026. Blog post.
- Triton Issues. <https://github.com/triton-lang/triton/issues/7392>, 2025. GitHub issue, triton-lang/triton, opened July 4, 2025.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- Robert Tjarko Lange, Qi Sun, Aaditya Prasad, Maxence Faldor, Yujin Tang, and David Ha. Towards robust agentic cuda kernel benchmarking, verification, and optimization, 2025. URL <https://arxiv.org/abs/2509.14279>.
- Gang Liao, Hongsen Qin, Ying Wang, Alicia Golden, Michael Kuchnik, Yavuz Yetim, Jia Jiunn Ang, Chunli Fu, Yihan He, Samuel Hsia, et al. Kernelevolve: Scaling agentic kernel coding for heterogeneous ai accelerators at meta. *arXiv preprint arXiv:2512.23236*, 2025.
- Tao Lin. Measuring automated kernel engineering. <https://metr.org/blog/2025-02-14-measuring-automated-kernel-engineering/>, 02 2025.
- Zhixuan Lin, Evgenii Nikishin, Xu Owen He, and Aaron Courville. Forgetting transformer: Softmax attention with a forget gate. *arXiv preprint arXiv:2503.02130*, 2025.
- Chao Liu, Jing Zhang, Letao Qin, Qianfeng Zhang, Liang Huang, Shaojie Wang, Anthony Chang, Chunyu Lai, Illia Silin, Adam Osewski, Poyen Chen, Rosty Geyyer, Hanwen Chen, Tejash Shah, Xiaoyan Zhou, and Jianfeng Yan. Composable kernel. URL [https://github.com/ROCm/composable\\_kernel](https://github.com/ROCm/composable_kernel).

- 495 NVIDIA. NVIDIA H100 Tensor Core GPU Architecture.  
 496 Whitepaper V1.01, 2022. URL <https://resources.nvidia.com/en-us-tensor-core/gtc22-whitepaper-hopper>.  
 497  
 498  
 499 NVIDIA. NVIDIA Blackwell Architecture Technical Brief.  
 500 Technical brief, 2024a. URL <https://www.nvidia.com/en-us/data-center/technologies/blackwell-architecture/>.  
 501  
 502  
 503 NVIDIA. *Parallel Thread Execution ISA Reference Manual*,  
 504 2024b. URL <https://docs.nvidia.com/cuda/parallel-thread-execution/>. Versions 6.0–  
 505 8.6. Instruction family counts tallied from per-version  
 506 changelogs.  
 507  
 508  
 509 NVIDIA. CUTLASS: CUDA Templates for Linear Algebra  
 510 Subroutines. <https://github.com/NVIDIA/cutlass>, 2025. Accessed: 2025-03-29.  
 511  
 512  
 513 NVIDIA. cuBLAS library documentation. NVIDIA Developer Documentation, 2026. URL <https://docs.nvidia.com/cuda/cublas/>.  
 514  
 515  
 516 Anne Ouyang, Simon Guo, Simran Arora, Alex L Zhang,  
 517 William Hu, Christopher Ré, and Azalia Mirhoseini. Kernelbench: Can llms write efficient gpu kernels? *arXiv preprint arXiv:2502.10517*, 2025.  
 518  
 519  
 520 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer,  
 521 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.  
 522  
 523  
 524 PyTorch Contributors. Pytorch source code: test\_fp8.py  
 525 (line 442). [https://github.com/pytorch/pytorch/blob/696b6fa86f8f6443724c7d03756238a032264892/test/inductor/test\\_fp8.py#L442](https://github.com/pytorch/pytorch/blob/696b6fa86f8f6443724c7d03756238a032264892/test/inductor/test_fp8.py#L442), 2024. Commit 696b6fa.  
 526  
 527  
 528 Manman Ren, Nick Riasanovsky, Neil Dhar, Hongtao Yu,  
 529 Jie Liu, Partha Kanuparth, and Shane Nay. Warp specialization in triton: Design and roadmap. <https://pytorch.org/blog/warp-specialization-in-triton-design-and-roadmap/>, January  
 530 2026. PyTorch Blog.  
 531  
 532  
 533 Bitar Darvish Rouhani, Ritchie Zhao, Ankit More, Mathew  
 534 Hall, Alireza Khodamoradi, Summer Deng, Dhruv Choudhary, Marius Cornea, Eric Dellinger, Kristof Denolf, et al. Microscaling data formats for deep learning. *arXiv preprint arXiv:2310.10537*, 2023.  
 535  
 536  
 537 Benjamin F Spector, Simran Arora, Aaryan Singhal,  
 538 Daniel Y Fu, and Christopher Ré. Thunderkittens: Simple, fast, and adorable ai kernels. *arXiv preprint arXiv:2410.20399*, 2024.  
 539  
 540  
 541 tiktoken. tiktoken: A fast bpe tokeniser for use with openai’s  
 542 models. <https://github.com/openai/tiktoken>, 2022.  
 543  
 544  
 545 Philippe Tillet, Hsiang-Tsung Kung, and David Cox. Triton: an intermediate language and compiler for tiled neural network computations. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages*, pages 10–19, 2019.  
 546  
 547  
 548 Arya Tschand, Muhammad Awad, Ryan Swann, Kesavan Ramakrishnan, Jeffrey Ma, Keith Lowery, Ganesh Dasika, and Vijay Janapa Reddi. Swizzleperf: Hardware-aware llms for gpu kernel performance optimization, 2025. URL <https://arxiv.org/abs/2508.20258>.  
 549  
 550  
 551 Xingyao Wang, Boxuan Li, Yufan Song, Frank F. Xu, Xiangru Tang, Mingchen Zhuge, Jiayi Pan, Yueqi Song, Bowen Li, Jaskirat Singh, Hoang H. Tran, Fuqiang Li, Ren Ma, Mingzhang Zheng, Bill Qian, Yanjun Shao, Niklas Muennighoff, Yizhe Zhang, Binyuan Hui, Junyang Lin, Robert Brennan, Hao Peng, Heng Ji, and Graham Neubig. Openhands: An open platform for ai software developers as generalist agents, 2025. URL <https://arxiv.org/abs/2407.16741>.  
 552  
 553  
 554 Nina Wiedemann, Quentin Leboutet, Michael Paulitsch, Diana Wofk, and Benjamin Ummenhofer. Kernelfoundry: Hardware-aware evolutionary gpu kernel optimization, 2026. URL <https://arxiv.org/abs/2603.12440>.  
 555  
 556  
 557 Samuel Williams, Andrew Waterman, and David Patterson. Roofline: an insightful visual performance model for multicore architectures. *Communications of the ACM*, 52(4):65–76, 2009.  
 558  
 559  
 560 Guangxuan Xiao, Ji Lin, Mickael Seznec, Hao Wu, Julien Demouth, and Song Han. SmoothQuant: Accurate and efficient post-training quantization for large language models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett, editors, *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 38087–38099. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/xiao23c.html>.  
 561  
 562  
 563 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.  
 564  
 565  
 566 Songlin Yang and Yu Zhang. Fla: A triton-based library for hardware-efficient implementations of linear attention mechanism, January 2024. URL <https://github.com/fla-org/flash-linear-attention>.

550 Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and  
551 Yoon Kim. Parallelizing linear transformers with the delta  
552 rule over sequence length. *Advances in neural informa-*  
553 *tion processing systems*, 37:115491–115522, 2024.

554 Zihao Ye, Lequn Chen, Ruihang Lai, Wuwei Lin, Yineng  
555 Zhang, Stephanie Wang, Tianqi Chen, Baris Kasikci,  
556 Vinod Grover, Arvind Krishnamurthy, et al. Flashin-  
557 fer: Efficient and customizable attention engine for llm  
558 inference serving. *Proceedings of Machine Learning and*  
559 *Systems*, 7, 2025.

560 Zijian Zhang, Rong Wang, Shiyang Li, Yuebo Luo, Mingyi  
561 Hong, and Caiwen Ding. Cudaforge: An agent framework  
562 with hardware feedback for cuda kernel optimization,  
563 2025. URL [https://arxiv.org/abs/2511.0](https://arxiv.org/abs/2511.01884)  
564 [1884](https://arxiv.org/abs/2511.01884).

565  
566  
567  
568  
569  
570  
571  
572  
573  
574  
575  
576  
577  
578  
579  
580  
581  
582  
583  
584  
585  
586  
587  
588  
589  
590  
591  
592  
593  
594  
595  
596  
597  
598  
599  
600  
601  
602  
603  
604

## Appendix

## Table of Contents

605	<b>A</b>	<b>Hardware-specific kernel code anatomy</b>	<b>13</b>
606		A.1 MMA across architectures	13
607		A.2 Async data movement across architectures	13
608		A.3 Kernel across precisions	15
609		A.4 Taxonomy derivation from FlashAttention-2	15
610		A.5 Taxonomy derivation from ThunderKittens	17
611		A.6 Growth of hardware-specific code	17
612	<b>B</b>	<b>Frontier LLMs on the taxonomy</b>	<b>19</b>
613		B.1 Pass@1 across architectures	19
614		B.2 Failure-mode buckets	19
615	<b>C</b>	<b>Structure of the Hawkeye taxonomy and agent workspace</b>	<b>20</b>
616		C.1 Full taxonomy	20
617		C.2 Anatomy of a taxonomy cell	21
618		C.3 Workspace layout and retrieval	23
619	<b>D</b>	<b>Comparative analysis of kernel optimization systems and context structures</b>	<b>23</b>
620		D.1 Related autonomous kernel optimization systems	23
621		D.2 Three-cluster failure-mode framework	24
622		D.3 Triton FLA	25
623		D.4 Raw PTX and ROCm docs	26
624		D.5 Augmented documentation	27
625		D.6 Inlined GEMM exemplar	28
626		D.7 Vendor FMHA	30
627		D.8 PTX abstractions ablation	31
628	<b>E</b>	<b>Optimization stacking and composition</b>	<b>33</b>
629		E.1 Closed-loop pattern	33
630		E.2 Trajectories across workloads	35
631		E.3 Composition across kernels	36
632		E.3.1 Cross-workload reuse	36
633		E.3.2 Cross-architecture composition	38
634		E.3.3 Cross-precision cascade	38
635		E.3.4 Invented patterns	39
636	<b>F</b>	<b>Generalization across porting axes</b>	<b>40</b>
637		F.1 Cross-workload	40
638		F.2 Cross-architecture	40
639		F.3 Cross-precision cascade	41
640		F.4 Problem-size sweep	42
641		F.5 Roofline analysis	42
642	<b>G</b>	<b>Experimental infrastructure</b>	<b>43</b>
643		G.1 Prompt templates	43
644		G.2 Evaluation pipeline	43
645		G.3 Precision and correctness tolerance	45
646		G.4 Compilation and profiling	45

```

660 template<int ScaleD, int ScaleA, int ScaleB, int TransA, int TransB>
661 __device__ __forceinline__ void mma_async_m64n128k16(
662     float d[8][8], uint64_t desc_a, uint64_t desc_b) {
663     asm volatile(
664         "{\n"
665         "wmma.mma_async.sync.aligned.m64n128k16.f32.bf16.bf16 "
666         "{%0, %1, ..., %63}," // 64 f32 accumulator registers
667         "%64," // desc_a (SMEM descriptor)
668         "%65," // desc_b (SMEM descriptor)
669         "%66, %67, %68, %69, %70;\n" // ScaleD, ScaleA, ScaleB, TransA, TransB
670         "}\n"
671         : /* d[0][0] ... d[7][7] as "+f" outputs */
672         : "l"(desc_a), "l"(desc_b),
673         "n"(int32_t(ScaleD)), "n"(int32_t(ScaleA)),
674         "n"(int32_t(ScaleB)), "n"(int32_t(TransA)), "n"(int32_t(TransB));
675     }
676 
```

Figure 7. Hopper WGMMMA `wmma.mma_async.sync.aligned.m64n128k16` BF16 wrapper.

```

677 __device__ __forceinline__
678 fp32x16 bf16_32x32x16(bf16x8 a, bf16x8 b, fp32x16 c) {
679     return __builtin_amdgcn_mfma_f32_32x32x16_bf16(
680         a, b, c,
681         0, 0, 0 // cbsz, abid, blgp
682     );
683 }
684 
```

Figure 8. CDNA4 / MI350 `v_mfma_f32_32x32x16_bf16` wrapper.

## A Hardware-specific kernel code anatomy

The same algorithm requires materially different code on each GPU architecture and each precision, and hardware-specific code dominates production kernel libraries. Sections A.1–A.3 show side-by-side excerpts from HAWKEYE API wrappers, Section A.4 derives the taxonomy rows from FlashAttention-2, and Section A.6 quantifies the cost of supporting new hardware.

### A.1 MMA across architectures

Four HAWKEYE API wrappers realize the same logical  $16 \times 16 \times 16$  BF16 matmul, with the issue scope widening from a warp on Ampere (`mma.sync`) to a warpgroup on Hopper (Figure 7, 128 threads issuing asynchronously into 64 accumulator registers via 64-bit SMEM descriptors and `wmma.commit_group/wmma.wait_group`) to a full CTA on Blackwell (`tcgen05.mma` with TMEM accumulators), while CDNA4 holds at the wavefront (Figure 8, 64 threads producing a  $32 \times 32$  tile in 16 AGPR accumulators per lane with `cbsz`, `abid`, and `blgp` immediates and `s_waitcnt` completion). Each step moves operand staging deeper into the memory hierarchy and pushes tile-shape selection out of the mnemonic and into a runtime descriptor (full Ampere and Blackwell wrappers in the code release).

### A.2 Async data movement across architectures

Each generation introduces a new load mechanism, shifting both the unit of transfer and the addressing entity. Ampere uses per-thread `cp.async` with software-pipelined `commit/wait`. Hopper offloads the tile to a hardware DMA engine via a host-side `CUTensorMap`, with completion tracked as an `mbarrier` byte count (Figure 9), and adds a `multicast::cluster` variant that broadcasts one global read into every CTA in a thread-block cluster. Blackwell keeps the same TMA unit and adds 2-CTA cooperative loads where a single tile feeds a `tcgen05.mma.cta_group::2` spanning the cluster pair. CDNA4 uses a descriptor-driven DMA that flows from HBM into LDS without touching a VGPR, with completion tracked by scalar wait counts (Figure 10). The agent learns each from one API wrapper paired with one unit test exercising the matching `config.json` counter (full Ampere and Blackwell wrappers in the code release).

```

715
716
717
718
719 template<int col_group_elems = 64>
720 __device__ __forceinline__ void load(
721     void* dst,
722     void const* src_tma_map,
723     uint64_t* bar,
724     int row_offset,
725     int col_offset) {
726     uint64_t tma_ptr = reinterpret_cast<uint64_t>(src_tma_map);
727     uint32_t mbar_ptr = static_cast<uint32_t>(__cvta_generic_to_shared(bar));
728     uint32_t dst_ptr = static_cast<uint32_t>(__cvta_generic_to_shared(dst));
729     asm volatile(
730         "cp.async.bulk.tensor.5d.shared::cluster.global.tile."
731         "mbarrier::complete_tx::bytes"
732         " [%0], [%1, {%3, %4, %5, 0, 0}], [%2];"
733         : : "r"(dst_ptr), "l"(tma_ptr), "r"(mbar_ptr),
734         "n"(0), "r"(row_offset), "r"(col_offset / col_group_elems)
735         : "memory");
736 }

```

Figure 9. Hopper TMA tile load. One thread instructs the DMA engine to copy a full 2D tile, with mbarrier byte-count completion.

```

745 __device__ __forceinline__
746 i32x4 make_srsrc(const void* ptr, uint32_t range_bytes,
747     uint32_t row_stride_bytes = 0) {
748     auto as_u64 = (uint64_t)reinterpret_cast<uintptr_t>(ptr);
749     // config = 0x110000: NUM_FORMAT=float, DATA_FORMAT=32bit, raw buffer
750     buffer_resource rsrc = {as_u64, range_bytes, 0x110000};
751     /* optional swizzle-enable bits in row_stride_bytes ... */
752     return *reinterpret_cast<const i32x4*>(&rsrc);
753 }
754 __device__ __forceinline__
755 void buffer_load_lds_dwordx4(i32x4 srd, uint32_t lds_base,
756     int voffset, int soffset) {
757     uint32_t lds_addr = __builtin_amdgcn_readfirstlane(lds_base);
758     asm volatile(
759         "s_mov_b32 m0, %0\n"
760         "buffer_load_dwordx4 %1, %2, %3 offen lds"
761         : : "s"(lds_addr), "v"(voffset), "s"(srd), "s"(soffset)
762         : "memory", "m0");
763 }

```

Figure 10. CDNA4 SRD-based global-to-LDS DMA. Each lane provides only an offset, the base lives in SGPRs, and data flows directly from HBM to LDS without touching VGPRs.

```

770 // Step 1: S = Q_consumer @ K^T (64x128) using FP8 WGMMA
771 // FP8 WGMMA K=32, HEAD_DIM/32 = 4 iterations
772 float s_acc[BN / 16][8];
773 memset(s_acc, 0, sizeof(s_acc));
774
775 warpgroup::fence();
776 #pragma unroll
777 for (int k_it = 0; k_it < HEAD_DIM / WGMMA_K_FP8; ++k_it) { // WGMMA_K_FP8 = 32
778     uint64_t dq = wgmma::make_descriptor(
779         &sQ_base[k_it * WGMMA_K_FP8], 16, 512, wgmma::Swizzle::B64);
780     uint64_t dk = wgmma::make_descriptor(
781         &sK[k_it * WGMMA_K_FP8], 16, 512, wgmma::Swizzle::B64);
782     wgmma::mma_async_e4m3<BN, 1, 1, 1>(s_acc, dq, dk); // 3 immediates, not 5
783 }
784 warpgroup::commit();
785 warpgroup::wait<0>();

```

Figure 11. Hopper FP8 DeltaNet inner loop. Three immediates (ScaleD, ScaleA, ScaleB) since FP8 wgmma does not expose transpose control,  $K=32$  per instruction collapses the inner  $k_{it}$  loop by  $2\times$ , and the descriptor stride for the swizzle group is B64 instead of B128.

### A.3 Kernel across precisions

Moving from BF16 to FP8 to NVFP4 (Blackwell) or MXFP4 (CDNA4) is a cascade through the kernel rather than a datatype substitution. Element size halves, so MMA  $K$  per instruction doubles, the SMEM column group (128B / element\_size) shrinks, the swizzle pattern and TMA tensor map geometry change to match, the minimum  $BK$  shifts to fill one swizzle group, and FP4 grows an entire new operand class for block scales (TMEM on Blackwell, LDS on CDNA4) with its own staging code. Precision-invariant steps such as applying  $\beta$  and the causal mask, writing  $W$ , and the second matmul in BF16 on the dequantized intermediate are unchanged, so low-precision kernels are actually mixed-precision in practice. The taxonomy carries a dedicated quantized\_mma cell per architecture so the agent learns the cascade end to end before composing it into a workload.

Figures 11 and 12 show the inner loop of the same DeltaNet kernel under this cascade, contrasting NVIDIA’s descriptor side channel with AMD’s instruction-level scaling. On Hopper the WGMMA template collapses from five immediates to three because FP8 exposes no transpose control, and the descriptor stride moves to a B64 swizzle group. CDNA4 reuses the  $32\times 32\times 64$  MFMA family and selects FP4 via `cbsz : 4` and `blgp : 4`, with MXFP4 block scaling exposed as a separate `v_mfma_scale_f32_32x32x64_f8f6f4` instruction that reaches roughly 87% of unscaled BF16 throughput on MI350. The Hopper BF16 baseline and the Blackwell NVFP4 path appear in the code release.

### A.4 Taxonomy derivation from FlashAttention-2

We derive the taxonomy rows by decomposing a public FlashAttention-2 implementation into optimization regions, with each recurring region mapping to exactly one row. Figure 13 shows the MMA Unit region (Row 1). Rows 2 (Shared Memory Layout via `ldmatrix`), 3 (Async Pipeline via `cp.async.cg commit/wait`), 4 (Warp Reduction for online softmax), and 5 (Vectorized Memory via `float4`-equivalent loads) follow the same decomposition and appear in the code release.

The MMA region pins the kernel’s compute-side primitive. Both attention matmuls compile to Ampere `mma.sync.aligned.m16n8k16`, and that one instruction choice propagates outward to dictate thread count, output tile size, and operand format across the rest of the kernel. The Async Pipeline region pins the data-side primitive the same way (`cp.async.cg` bypassing L1, `commit/wait` forming a software pipeline that overlaps the next K/V load with the current matmul, illustrated for each architecture in Section A.2). Each row in the taxonomy isolates one such region so the agent can vary it without disturbing the rest.

Every row pairs to a profiling counter that confirms the optimization fired. The MMA row reads `sm_inst_executed_pipe_tensor`, Shared Memory Layout reads shared-memory efficiency and bank-conflict counters, Async Pipeline reads DRAM bandwidth and the barrier-stall percentage, Vectorized Memory reads `lltex_t_sectors_pipe_lsu_mem_global_op_ld`, and Warp Reduction reads warp execution efficiency. The mapping appears alongside the row in Table 4, and Section C.2 shows how each counter is wired into the unit-test artifact.

```

825
826
827
828
829 // ---- FP4 MFMA: QK = Q @ K^T with MXFP4 block scaling ----
830 // One scaled call. Q and K are raw FP4 with E8M0 per-block scales.
831 fp32x16 qk;
832 {
833     uint8_t* q_lds = (uint8_t*)(smem + MK3_Qfp4);
834     uint8_t* k_lds = (uint8_t*)(smem + MK3_Kfp4);
835     constexpr int HF = D_C / 2;
836     constexpr int scales_per_row = D_C / 32; // 1 E8M0 per 32 elements
837
838     intx4_t a_op = *(const intx4_t*)(q_lds + (lm + my_row)*HF + my_kh*16);
839     intx4_t b_op = *(const intx4_t*)(k_lds + (ln + my_row)*HF + my_kh*16);
840
841     // Load E8M0 block scales for Q (A operand) and K (B operand)
842     int sq = (int)Q_sc[(cs + lm + my_row) * scales_per_row + my_kh];
843     int sk = (int)K_sc[(cs + ln + my_row) * scales_per_row + my_kh];
844
845     qk = {};
846     qk = mfma_fp4::mfma_scale_32x32x64(a_op, b_op, qk, sq, sk);
847
848     #pragma unroll
849     for (int i = 0; i < 16; i++) qk[i] *= scale; // post-multiply by 1/sqrt(D)
850 }
851
852
853
854
855
856
857
858

```

Figure 12. CDNA4 MXFP4 DeltaNet inner loop. Same MFMA mnemonic family as BF16 (`v_mfma_..._f8f6f4`) with FP4 mode selected via `cbisz:4 blgp:4` immediates and a separate scaled variant `v_mfma_scale_f32_32x32x64_f8f6f4` that takes E8M0 scale registers.  $K=64$ .

```

859 #pragma unroll
860 for (int j = 0; j < kWarpTileSeqLenK; ++j) {
861     HMMA16816(R_S[0][j][0], R_S[0][j][1],
862             R_Q[0][0][0], R_Q[0][0][1],
863             R_Q[0][0][2], R_Q[0][0][3],
864             R_K[j][0], R_K[j][1],
865             R_S[0][j][0], R_S[0][j][1]);
866 }
867 #pragma unroll
868 for (int j = 0; j < kWarpTileHeadDimV; ++j) {
869     HMMA16816(R_O[0][j][0], R_O[0][j][1],
870             R_S[0][w][0], R_S[0][w][1],
871             R_S[0][w + 1][0], R_S[0][w + 1][1],
872             R_V[j][0], R_V[j][1],
873             R_O[0][j][0], R_O[0][j][1]);
874 }
875
876
877
878
879

```

Figure 13. Row 1, MMA Unit. FlashAttention-2 HMMA16816 calls for  $QK^T$  and PV.

```

880 struct producer {
881     __device__ static void setup(producer_setup_args<layout> args) {
882         warpgroup::decrease_registers<40>();
883     }
884     __device__ static void load(producer_load_args<layout> args) {
885         if (warpgroup::laneid() == 0) {
886             tma::expect(args.inputs_arrived, args.input);
887             for (int i = 0; i < M_BLOCK; i++)
888                 tma::load_async(args.input.a[i], args.globals.A,
889                               {args.common.coord.x+i, args.iter}, args.inputs_arrived);
889             for (int i = 0; i < N_BLOCK; i++)
890                 tma::load_async(args.input.b[i], args.globals.B,
891                               {args.iter, args.common.coord.y+i}, args.inputs_arrived);
891         }
892     }
893 };
894 struct consumer {
895     __device__ static void setup(consumer_setup_args<layout> args) {
896         warpgroup::increase_registers<232>();
897         kittens::warp::zero(args.state.accum);
898     }
899     __device__ static void compute(consumer_compute_args<layout> args) {
900         warpgroup::mma_AB(args.state.accum,
901                          args.input.a[warpgroup::groupid()],
902                          reinterpret_cast<wide_tile&>(args.input.b));
903         warpgroup::mma_async_wait();
904         if (warp::laneid() == 0) arrive(args.inputs_finished);
905     }
906 };
    
```

Figure 14. ThunderKittens Hopper BF16 GEMM. Producer warpgroup streams A/B tiles in via TMA, consumer warpgroup issues WGMMMA, with explicit register reallocation between the two roles.

The remaining rows (Quantized MMA, Producer-Consumer, Epilogue Pipeline, Multi-unit Coordination, Persistent Scheduling) cover kernel regions that FA-2 does not reach. We derive these in Appendix A.5 from the Hopper and Blackwell GEMMs in ThunderKittens, with HipKittens informing the CDNA4 cells.

## A.5 Taxonomy derivation from ThunderKittens

The remaining rows draw deep inspiration from ThunderKittens (Spector et al., 2024) on NVIDIA and its AMD counterpart HipKittens (Hu et al., 2025) on CDNA4. We analyze the BF16, FP8, and NVFP4 production GEMMs in ThunderKittens for Hopper and Blackwell to identify the optimizations they rely on for peak throughput, and the operators they are built from. Figures 14 and 15 show the central regions of each.

The recurring regions we extract this way inform both the row inventory and the reference solution kernel in each Hopper, Blackwell, or CDNA4 cell. FA-2 grounds rows 1–5, ThunderKittens and HipKittens ground rows 6–10, and the taxonomy keeps that lineage explicit so adding a new architecture column reduces to porting the same set of regions rather than rediscovering them.

## A.6 Growth of hardware-specific code

Supporting a new GPU is dominated by architecture-specific code rather than algorithm. Figure 16 shows the per-architecture cost growing within one project, where NVIDIA’s CUTLASS library has expanded 29× across recent versions and architecture-specific instruction wrappers, descriptor builders, and pipeline templates now make up the majority of the codebase. The growth is not in cross-architecture algorithm code, since each new generation introduces a new tier of hardware concepts (warpgroup, cluster, TMEM, AGPR, SRD) that requires its own wrapper layer and does not amortize against earlier generations.

Figure 17 shows the same cost across the open-source ecosystem. Manually porting projects like FlashAttention, FlashInfer,

935  
936  
937  
938  
939  
940  
941  
942  
943  
944  
945  
946  
947  
948  
949  
950  
951  
952  
953  
954  
955  
956  
957  
958  
959  
960  
961  
962  
963  
964  
965  
966  
967  
968  
969  
970  
971  
972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989

```

d_tt_t d_tt[C::MMA_PIPE_DEPTH];
#pragma unroll
for (int i = 0; i < C::MMA_PIPE_DEPTH; i++)
    d_tt[i] = tm_alloc.template allocate<d_tt_t>((i+warpgroup::warpid())*C::Nb);

for (int task_iter = 0; ; task_iter++) {
    for (int idx = 0; idx < iters_per_task; idx++) {
        tma::cluster::load_async(a_smem[input_ring][i], g.a,
            {(tile_coord.x*2+cta_rank)*C::NUM_CONSUMERS+i, idx},
            inputs_arrived[input_ring], (uint16_t)(1<<cta_rank), 0);
        wait(inputs_arrived[input_ring], get_phasebit<0>(bitfield, input_ring));
        if (idx == 0) mm2_ABt (d_tt[task_iter%C::MMA_PIPE_DEPTH],
            a_smem[input_ring][warpgroup::warpid()],
            b_smem[input_ring], inputs_finished[input_ring]);
        else mma2_ABt (d_tt[task_iter%C::MMA_PIPE_DEPTH],
            a_smem[input_ring][warpgroup::warpid()],
            b_smem[input_ring], inputs_finished[input_ring]);
        input_ring = ring_advance<C::LOAD_PIPE_DEPTH>(input_ring);
    }
    detail::tcgen05::commit<C::CLUSTER_SIZE>(outputs_arrived[warpgroup::warpid()]);
    if (!schedule.success) break;
}
    
```

Figure 15. ThunderKittens Blackwell BF16 GEMM. A 2-CTA cooperative MMA (mm2\_ABt/mma2\_ABt) accumulates into tensor memory (d\_tt); TMA loads carry a cluster multicast mask.

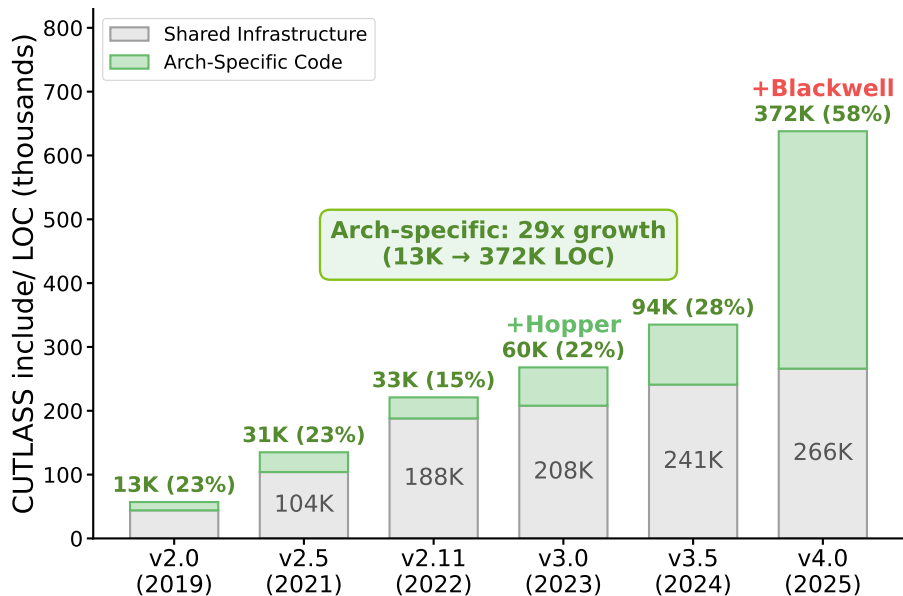


Figure 16. Architecture-specific code in NVIDIA’s CUTLASS library (NVIDIA, 2025) has grown 29× across recent versions and now constitutes the majority of the codebase. The growth is in per-architecture instruction wrappers, descriptor builders, and pipeline templates rather than in cross-architecture algorithm code.

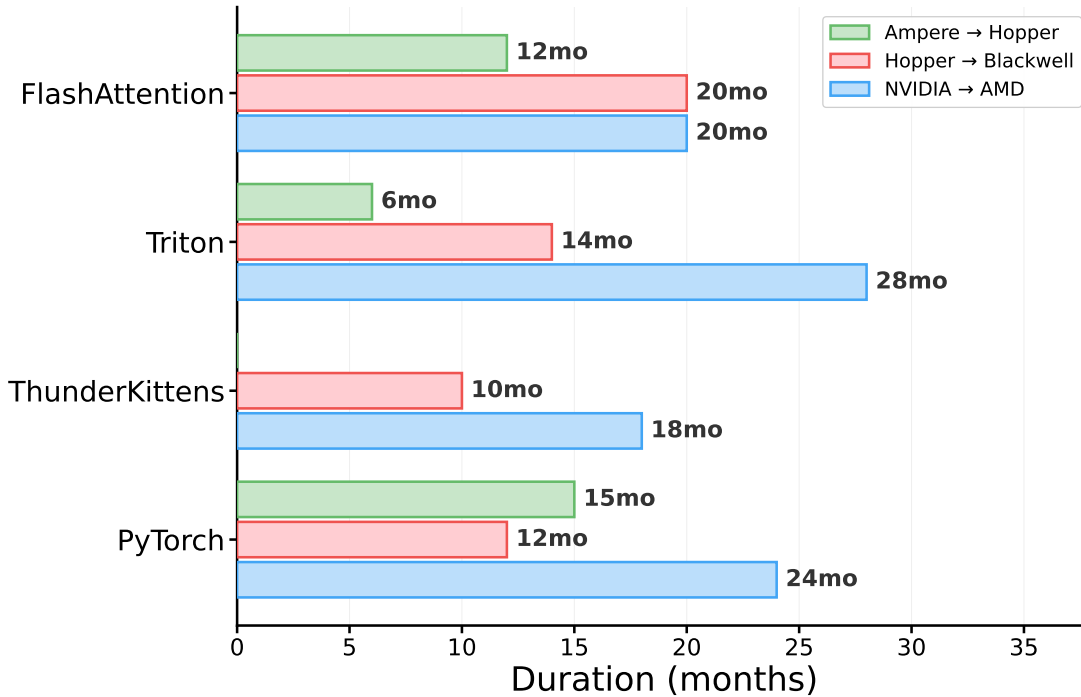


Figure 17. Open-source porting timelines computed from Git history. Manually porting projects like FlashAttention (Dao et al., 2022), FlashInfer (Ye et al., 2025), Triton (Tillet et al., 2019), ThunderKittens (Spector et al., 2024; Hu et al., 2025), and PyTorch (Paszke et al., 2019) between GPU generations or between vendors has historically taken months to years of expert effort. The trend does not flatten with experience because each new generation introduces a new tier of hardware concepts for which prior porting work transfers only partially.

Triton, ThunderKittens, and PyTorch between GPU generations or between vendors has historically taken motivated expert teams months to years, and the trend does not flatten with experience because each new tier of hardware concepts transfers only partially from prior porting work. Together the two figures bound the manual cost any day-zero porting approach must compete against.

## B Frontier LLMs on the taxonomy

The unit tests double as a measurement instrument for hardware awareness in frontier LLMs, motivating the gap (Section 1), pinning model selection (Section 3), and setting up the per-baseline failure clusters (Section D.2).

### B.1 Pass@1 across architectures

Per-architecture Pass@1 across the five frontier models is reported in Figure 2. A cell counts as passed when the first sampled completion compiles, produces an output that matches the expert solution kernel for that cell, and reaches at least 99% of the expert’s value on the cell’s targeted profiling counter (thresholds from `config.json`, Section C.2, with the 1% slack absorbing measurement noise). The threshold means a passing kernel matches or exceeds expert hardware utilization rather than merely returning a correct numerical result. All five models trace the same Ampere > Hopper > Blackwell ≥ MI350 ordering, losing roughly a third of the Ampere score moving to Hopper and roughly half of the Hopper score moving to Blackwell or MI350, and restricting Pass@1 to cells satisfiable by an Ampere-style implementation collapses the gradient to within 5 percentage points across architectures.

### B.2 Failure-mode buckets

Across failed unit tests, three buckets dominate, each tagged with the corresponding agent-loop cluster from Section D.2.

**Bucket 1, invalid syntax for the right strategy (missing API ground truth, cluster 2).** The model identifies the correct strategy but produces invalid syntax, with malformed `wgmma` descriptor immediates, wrong operand types where a `tmem`

```

1045 // Strategy is right (warp-group MMA over a shared-memory descriptor),
1046 // but the descriptor encoding is from an obsolete preview spec.
1047 uint64_t desc = make_smem_desc(
1048     smem_ptr,
1049     /*lbo*/ 64,    // <-- should be encoded as 64 >> 4 = 4
1050     /*sbo*/ 128); // model emits raw byte counts
1051 // nvcc accepts the literal but the GPU traps at issue.
1052 wgmma::mma_async<128, ...>(&desc, &acc[0]);

```

Figure 18. Bucket 1 example. Hopper wgmma descriptor with wrong stride field width.

```

1055 // Blackwell cell. tcgen05.mma + tmem is required.
1056 // Model outputs Hopper warp-group MMA instead.
1057 asm volatile("wgmma.mma_async.sync.aligned.m64n128k16.f32.bf16.bf16 "
1058             "{%0,%1,...}, %16, %17, p, 1, 1, 1, 0;\n" ...);
1059 // On sm_100a this either fails to encode or runs on the legacy
1060 // tensor-core path. tmem stays empty and the cell metric stays at 0.

```

Figure 19. Bucket 3 example. Hopper-style wgmma emitted on a Blackwell tensorcore\_mma cell that targets tcgen05.

register handle is required, and missing `.kind` qualifiers on Blackwell `tcgen05` variants. This is the most common failure on Hopper and Blackwell, illustrated in Figure 18. The code looks plausible until a real GPU run.

**Bucket 2, naive kernel with no targeted-metric improvement (missing algorithm structure, cluster 1).** The model falls back to a functionally correct scalar kernel with hand-rolled FMA loops on cells that exist precisely to exercise tensor cores, so `sm__inst_executed_pipe_tensor_op_hmma` stays at zero and the cell fails. Bucket 2 dominates on Ampere and CDNA4, where the API surface is small enough that syntax is right but the algorithmic mapping is wrong.

**Bucket 3, optimization intended for the wrong architecture (missing optimization pattern, cluster 3).** The model reaches for an optimization pattern from its corpus that does not match the cell’s architecture, emitting Hopper `wgmma` on a Blackwell test, Ampere `cp.async` on a Hopper test, or Ampere-style swizzle on an MI350 LDS layout. Bucket 3 is what the architecture-recency gradient predicts most strongly (Figure 19). The more recent the SM tier, the more often the model substitutes a same-name pattern from the previous tier.

The three buckets map cleanly to agent-loop clusters. Bucket 1 maps to Cluster 2 (missing API ground truth), patched by an API library that pins the encoding. Bucket 2 maps to Cluster 1 (missing algorithm structure), patched by a per-problem algorithmic recipe. Bucket 3 maps to Cluster 3 (missing optimization pattern), patched by per-arch optimization-pattern modules.

**Model selection and GPT replication.** A 30-turn agent loop on Hopper BF16 GEMM picks Gemini 3.1 Pro as the backbone for primary trajectories given its highest unguided ceiling and largest hardware-aware delta. We replicate the Linear-Attention-scaling sweep with GPT-5.4 across six baseline contexts and four architectures at 100 productive turns per run, and both models show the same plateau ordering. On GPT, HAWKEYE’s full configuration reaches  $\geq 0.66$  of expert TFLOPS on Ampere,  $\geq 0.86$  on Hopper, and  $\geq 0.91$  on MI350. On Gemini, the same architectures reach  $\geq 0.94$ ,  $\geq 0.99$ , and  $\geq 0.91$ .

## C Structure of the Hawkeye taxonomy and agent workspace

This appendix specifies what the agent reads when it consumes the HAWKEYE taxonomy, backing Section 2 via Sections C.2 and C.3.

### C.1 Full taxonomy

The complete HAWKEYE taxonomy is reproduced in Table 4. Rows are recurring optimization strategies and columns are target GPU architectures (Ampere, Hopper, Blackwell, MI350). Each cell records the architecture-specific syntax, the bottleneck utilization counter the unit test verifies, and the end-to-end speedup the optimization buys over a naive baseline.

Section C.2 expands what ships inside a single cell, and Section C.3 describes how the agent retrieves it at inference time.

Table 4. Hardware-aware optimization taxonomy. Each row is a general optimization strategy and each cell shows its architecture-specific details, paired with a unit test and expert-authored solution kernel. Every cell reports a bottleneck utilization metric (M=MMA, D=DRAM, G=global memory, S=shared memory, W=warp/wave efficiency;  $0 \rightarrow \infty$  marks features absent in the naive kernel) and the resulting end-to-end speedup over a naive baseline.

	Ampere (A100)	Hopper (H100)	Blackwell (B200)	MI350
<b>MMA Unit</b>	mma.sync.m16n8k8 BF16→FP32, ldmatrix.x2/.trans loads M↑8× 3.48×	WGMMMA m64n64k16 BF16 + TMA 128B-swizzle loads + mbarrier M↑2× 157.48×	tcgen05.mma BF16→FP32, TMA bulk loads, TMEM accumulators M∞ 7.17×	MFMA 32×32×16 BF16, AGPR accumulators M∞ 108×
<b>Quantized Precision</b>	No floating-point quantization support	FP8 E4M3 WGMMMA m64n64k32 + TMA, per-row/col scaling M↑2× 19.53×	NVFP4 tcgen05.mma mxF4nvf4 block-scale + 2-CTA cluster M∞ 6.73×	On-the-fly MXFP4 quantization, MFMA matrix core M∞ 6.7×
<b>Shared Mem Layout</b>	XOR swizzle family (s32/s64/s256) for conflict-free ldmatrix S∞ 1.05×	128B column-group XOR swizzle matching WGMMMA descriptors S∞ 2.59×	Parameterized 64B offset_swizzle for tcgen05 descriptors S↑1238.6× 1.62×	64-bank LDS XOR swizzle for bank-conflict-free transpose S↑1.3× 1.5×
<b>Vectorized Memory</b>	uint4 128-bit vectorized BF16 scale-and-add, grid-stride loop G↑8× 1.26×	uint4 128-bit vectorized BF16 scale-and-add, grid-stride loop G↑8× 1.50×	uint4 128-bit vectorized BF16 scale-and-add, grid-stride loop G↑4× 3.26×	128-bit global loads/stores via buffer_load_dwordx4 G↑16× 3.7×
<b>Async Pipeline</b>	3-stage cp.async pipeline, commit_group/wait_group D↑5.2× 5.22×	TMA cp.async.bulk.tensor + mbarrier 3-stage pipeline D↑2.27× 2.21×	TMA 2D + mbarrier 3-stage pipeline, phase-parity tracking D↑4.09× 3.89×	Triple-buffered LDS pipeline overlapping loads with MFMA D↑1.2× 5.1×
<b>Producer/Consumer</b>	Warp-specialized WMMA, warp 0 loads A/B, warps 1-7 compute, double-buffered SMEM D↑1.2× 5.22×	WGMMMA + TMA with warpgroup register reallocation D∞ 823.28×	tcgen05.mma + TMA, producer/consumer warp specialization D∞ 834.81×	Double-buffered 4-wave MFMA, load/compute overlap D↑11.4× 3.4×
<b>Epilogue Pipeline</b>	mma.sync BF16 GEMM, ReLU fused on FP32 accumulator registers D↑7.7× 7.70×	WGMMMA tiled GEMM, fused GeLU epilogue in FP32 accumulators D↑46.2× 14.62×	tcgen05 GEMM, fused GeLU epilogue from TMEM D↑7.5× 9.71×	Fused bias + GeLU + E8M0 scale + FP8 single pass M↑7.9× 4.0×
<b>Warp/Wave Reduction</b>	Hierarchical block_reduce_sum via __shfl_down_sync shuffles W↑64× 1.11×	block_reduce_sum via __shfl_down_sync + atomic finalize W↑64× 1.00×	block_reduce_sum via __shfl_down_sync + atomic finalize W↑64× 1.17×	64-lane __shfl_xor butterfly replacing global atomics W↑4.0× 7.3×
<b>Multi-Unit Coordination</b>	Cooperative grid_sync() two-phase reduction, L2 partial-sum reuse D↑8× 14.54×	2-CTA cluster, DSMEM (mapa.shared::cluster + ld.shared) exchange D↑2× 1.05×	2-CTA cluster, DSMEM (mapa.shared::cluster) cross-CTA exchange D↑2× 1.12×	XCD-aware blockIdx remapping with L2 swizzle D↑2.8× 1.9×
<b>Persistent Scheduling</b>	mma.sync PTX, SWIZZLE_GROUP=4 L2-friendly blockIdx remap M↑2× 2.14×	Persistent WGMMMA + TMA, atomicAdd tile dispatch + L2 swizzle M↑2× 22.15×	Persistent tcgen05 + TMA, atomicAdd tile dispatch + L2 swizzle M∞ 3.83×	Persistent MFMA wavefronts with atomicAdd tile dispatch M∞ 6.9×

## C.2 Anatomy of a taxonomy cell

A taxonomy cell sits at the intersection of one row (an optimization strategy) and one column (an architecture), and it ships four artifacts. A `naive_kernel.cu` or `naive_kernel.hip` establishes the baseline the optimization has to beat by intentionally avoiding the targeted hardware feature so its profiling counter rides at the floor. An `optimized_kernel.cu` is the expert-authored reference solution that the agent can read as a syntax example, invoke directly as a callable, or lift fragments from when composing into a larger kernel. A `config.json` pairs the cell to a profiling counter and expected gain, and a guide encodes the protocol the agent must follow. The Blackwell async-pipeline cell illustrates the shape (Figures 20 and 21). The optimized kernel is a 3-stage TMA ring with mbarrier multi-phase coordination paired to `dram__throughput.avg.pct_of_peak_sustained_elapsed`, and the guide encodes the phase-bit protocol with initial value `0xFFFF0000` (bits 31-16 producer phase, bits 15-0 consumer phase). Without per-module guides preloaded, the Blackwell agent caps performance at ~240 TFLOPS and at one point hallucinates the non-existent intrinsic `__pack_half2` at turn 32 (Section D.8).

The CDNA4 column shares the four-artifact shape with extensions `.hip` and flags `-offload-arch=gfx950`, and it exposes the AGPR/VGPR split where accumulators live in a separate  $256 \times 32$ -bit register file requiring explicit `v_accvgpr_read_b32` and `v_accvgpr_write_b32` transfers. Its `config.json` pairs the cell to three counters at once, with `SQ_INSTS_VALU_MFMA_BF16` expected to rise and `SQ_INSTS_VALU` and `SQ_WAVES` expected to fall, so the agent can detect MFMA adoption versus a scalar-FMA fallback, and the guide commits to typed wrappers like

```

1155
1156
1157
1158
1159
1160 constexpr int NUM_STAGES = 3;
1161
1162 __device__ void mbarrier_arrive_expect_tx(int mbar, int tx_bytes) {
1163     asm volatile("mbarrier.arrive.expect_tx.release.cta.shared::cta.b64 "
1164                 "_ , [%0], %1;" :: "r"(mbar), "r"(tx_bytes) : "memory");
1165 }
1166 __device__ void tma_load_2d(int dst, const void* tmap,
1167                             int x, int y, int mbar) {
1168     asm volatile("cp.async.bulk.tensor.2d.shared::cta.global"
1169                 ".mbarrier::complete_tx::bytes [%0], [%1, {%2, %3}], [%4];"
1170                 :: "r"(dst), "l"(tmap), "r"(x), "r"(y), "r"(mbar) : "memory");
1171 }
1172 // Prologue: prefetch first NUM_STAGES-1 tiles
1173 for (int s = 0; s < min(NUM_STAGES - 1, num_blocks_k); s++) {
1174     if (threadIdx.x == 0) {
1175         tma_load_2d(smem_addr_s, tensorMapA, s*BK, m_block*BM, mbar_s);
1176         mbarrier_arrive_expect_tx(mbar_s, tile_bytes);
1177     }
1178 }

```

Figure 20. Blackwell async-pipeline cell, optimized\_kernel.cu. 3-stage TMA prologue with mbarrier-based completion.

```

1190 {
1191     "test_name": "async_data",
1192     "task_description": "Asynchronous data transfer with TMA and 3-stage
1193                         pipelining on Blackwell (SM100)",
1194     "compilation": { "flags": "-arch=sm_100a -O3 -lcuda" },
1195     "profiling": {
1196         "metrics": [{
1197             "name": "dram_throughput.avg.pct_of_peak_sustained_elapsed",
1198             "description": "DRAM bandwidth utilization percentage",
1199             "direction": "higher_is_better"
1200         }]
1201     }
1202 }

```

Figure 21. Blackwell async-pipeline config.json. The dram\_throughput counter is what the harness reports back to the agent so it can verify the TMA ring actually saturated DRAM rather than some unrelated change accounting for the runtime improvement.

```

1210 workspace/
1211   task_prompt.md, workload_implementation_guide.md, spec.json
1212   eval.py, bf16_torch_baseline.py, bf16_expert_baseline.py
1213   naive_kernel.cu, optimized_kernel.cu      # agent edits this
1214   blackwell-apis/                          # local API library
1215     blackwell.cuh, general.md              # umbrella + arch overview
1216     mma/, tmem/, tma/, barrier/, cluster/  # each has <name>.cuh + guide.md
1217     epilogue/, pipeline/, clc/, tile_scheduler/, attention_ops/, nvfp4/
1218   taxonomy/                                # column of unit tests
1219     01_tensorcore_mma/
1220     naive_kernel.cu, optimized_kernel.cu, config.json
1221     02_quantized_mma/ ... 10_persistent_scheduling/
1222   kernel_pool/                             # prior outer-loop kernels
1223     gemm_bf16.cu, attention_bf16.cu, ...

```

Figure 22. Workspace tree for the HAWKEYE method on Blackwell. The taxonomy/ column ships only the architecture-specific cells the agent should consult, and the kernel\_pool/ accumulates correct kernels from prior outer-loop runs.

mfma::bf16\_32x32x16 over raw builtins to avoid silent miscasts. All four artifacts are load-bearing. The naive baseline is what makes the targeted counter readable, the optimized kernel is the only executable syntax example for instructions that postdate the model’s training corpus, the config is what tells the agent which counter signals success, and the guide is what encodes the protocol the syntax alone does not. The harness always reports the counter alongside TFLOPS so a runtime improvement is verified to have come from exercising the targeted feature (full per-row mapping in Section A.4 and Table 4).

### C.3 Workspace layout and retrieval

Each method runs in a clean workspace visible to the agent through the shell, with HAWKEYE exposing the architecture column of the taxonomy and the accumulated kernel pool from prior outer-loop runs (Figure 22). The generic and KernelBench workspaces drop taxonomy/ and kernel\_pool/, GEPA replaces them with a single skill.md, and weaker hardware-aware workspaces swap taxonomy/ for docs/ or a single prior\_arch\_kernel.cu.

Across instrumented HAWKEYE trajectories, retrievals concentrate on taxonomy/ and apis/, and no substitute carries comparable retrieval-frequency mass through to peak performance (Section D.2). Three Blackwell baselines show what fills the void instead. ptx\_abstractions strips the per-module guide.md files, so the agent cats the headers and re-derives missing protocols inside the kernel (Section D.8). raw\_ptx\_docs replaces <arch>-apis/ with the full 1.22 M-token PTX ISA reference, and the agent issues 17 distinct wgmma-pattern searches across 221 turns. vendor\_fmha replaces both directories with the CUTLASS or Composable Kernel source tree; retrievals collapse into filesystem navigation, and the kernel body remains an empty stub for 207 turns.

## D Comparative analysis of kernel optimization systems and context structures

This appendix positions HAWKEYE comparatively at two levels: a landscape view of recent autonomous GPU kernel optimization systems (Section D.1), and a controlled comparison against six alternative context structures evaluated under a common agent harness and turn budget (Sections D.2–D.8). For the controlled comparison, we build a master baseline matrix across two backbone LLMs and four architectures, frame a three-cluster failure-mode taxonomy, walk each baseline through its plateau pathology, and ablate the optimization-pattern modules to isolate their contribution.

### D.1 Related autonomous kernel optimization systems

Table 5 compares HAWKEYE to recent autonomous GPU kernel optimization systems across implementation language, optimization method, hardware-specific support, required domain expertise, target GPU architectures, and supported precisions.

## Hawkeye: Hardware-Aware GPU Kernel Optimization with Minimal Supervision

System	Lang.	Method	HW-Specific Opt.	Domain Expertise	GPU Archs.	Precisions
KernelBench	CUDA	Beam search	✗	Basic HW specs	A	FP32
KevinRL	CUDA	RL self-feedback	✗	✗	H	FP32
AI CUDA Eng.	CUDA	Feedback loop	✗	✗	A	BF16, FP32
KernelEvolve	CUDA	Tree search + RAG	✗	RAG on docs	A H M	FP32
KernelFoundry (Wiedemann et al., 2026)	SYCL, CUDA	Quality-diversity search	Some	Templated kernels	A	FP32
KernelBlaster	CUDA	In-context RL	✗	✗	H	FP16, FP32
CudaAgent	CUDA	Planning + RL	✗	✗	H	FP32
CudaForge	CUDA	Feedback loop	✗	✗	A	FP32
SwizzlePerf	Triton	Feedback loop	1 opt type	Hand-written docs	M	BF16, FP32
K-Search	Triton, CUDA	World model search	Some	Past examples	H B	BF16, FP8
HF CUDA Skills	CUDA	Coding agent	Hand-written	Hand-written skills	A H	BF16, FP16, FP32
<b>HAWKEYE</b>	<b>CUDA, HIP</b>	<b>Coding agent with HW-awareness</b>	✓	<b>Taxonomy + unit tests</b>	<b>A H B M</b>	<b>BF16, FP8, NVFP4, MXFP4</b>

Table 5. Comparison of autonomous GPU kernel optimization systems. Architecture support: **A** = Ampere, **H** = Hopper, **B** = Blackwell, **M** = MI350. Prior systems either lack hardware-specific optimizations entirely (✗), achieve them partially through hand-written or outdated approaches, or require high expert effort that does not scale. Most prior systems also stay within standard precisions (FP32, BF16, FP16). HAWKEYE is the first framework to evaluate hardware-aware kernels across four frontier GPU architectures and four precisions including FP8, NVFP4, and MXFP4.

### D.2 Three-cluster failure-mode framework

The baselines span the practitioner spectrum from converged skill-evolution memory to the full vendor ISA reference. We evaluate all six on causal Parallel ReBased Linear Attention (henceforth *Linear Attention* for brevity) across four architectures (Ampere A100, Hopper H100, Blackwell B200, CDNA4 MI350) with two backbones (Gemini 3.1 Pro and GPT-5.4) for  $6 \times 4 \times 2 = 48$  runs at a fixed 100-productive-turn budget, and Table 6 reports the best TFLOPS for every (baseline, backbone, architecture) cell together with the total trajectory token cost. Per-baseline failure-cluster assignments are summarized in the prose below.

Table 6. Master baseline matrix. “Best TFLOPS” is the highest correct kernel observed in the trajectory at the 100-productive-turn budget. “-” marks runs where no correct kernel was produced (status NO\_CORRECT\_KERNEL or best\_tflops = null). Tokens are cumulative trajectory tokens (millions). Expert (peak) reference for Parallel ReBased Linear Attention BF16: Ampere  $\approx 132$ , Hopper  $\approx 298$ –300, Blackwell  $\approx 424$ , CDNA4  $\approx 333$ .

Baseline	Best TFLOPS (Gemini / GPT)				Tokens (M)	
	Ampere	Hopper	Blackwell	CDNA4	Gem	GPT
Triton FLA	27.2 / 4.4	17.9 / 16.2	110.4 / 34.6	- / 0.3	8.0	7.0
ptx_abstractions	84.7 / 10.1	240.9 / 17.5	- / 9.7	241.8 / 11.2	13.6	6.0
Raw PTX docs	17.8 / 19.6	31.6 / 9.2	5.6 / 19.8	169.7 / 16.0	10.4	8.2
Augmented PTX docs	39.7 / 10.0	4.4 / 27.2	5.6 / 26.2	118.4 / 20.0	9.2	7.0
inlined_gemm	3.6 / 8.7	4.4 / 16.6	- / 15.0	27.7 / 2.1	37.4	13.9
Vendor FMHA	33.8 / 18.0	29.7 / 37.6	5.7 / 18.9	39.2 / 0.8	16.6	7.3
<b>HAWKEYE full</b>	<b>87.3 / 135.8</b>	<b>256.4 / 327.9</b>	<b>702.0 / 418.9</b>	<b>326.5 / 319.2</b>	<b>25.0</b>	<b>12.1</b>

Every non-HAWKEYE baseline collapses on at least one architecture under at least one backbone (collapse defined as best TFLOPS at most  $0.1 \times$  expert), and each collapse is tied to a specific deficit in the context. The patterns partition into three structural clusters matching the missing-context types in Section 2.

**Cluster 1. Missing algorithm structure.** The agent lacks the high-level loop, pipeline, or composition pattern the architecture demands, so the kernel compiles and is correct but has no architecture-aware structure (a scalar  $O(S^2)$  loop on a tensor-core GPU, a single-warpgroup where the chip needs a producer-consumer split, or an empty stub). Triton FLA and inlined\_gemm land here, with Triton FLA Blackwell collapsing to 110.4 / 34.6 TFLOPS against the  $\sim 424$  expert and

```

1320 uint32_t bfl6_vals = __pack_half2(__float2bfloat16(val_x),
1321                                 __float2bfloat16(val_y));
1322 asm volatile(
1323     "stmatrix.sync.aligned.m8n8.x4.shared.b16 [%0], %1;\n"
1324     ":: "r"(smem_addr), "r"(bfl6_vals));
1325 // error: identifier "__pack_half2" is undefined

```

Figure 23. `ptx_abstractions` on Blackwell at turn 32. The agent hallucinates a non-existent CUDA intrinsic name and never recovers across 199 turns.

```

1329 T201 decision: "The text mentions Figure 149 for the D accumulator
1330 layout, but I can't see the image. However, I know from standard
1331 WGMMMA layouts that for m64n64k16 with f32, the D accumulator is
1332 32 registers per thread, laid out as [row_tile, col_tile][i]."

```

```

1333 T202 result: 77.9 TFLOPS reported, correctness = 1.92%
1334             (98% of output elements violate tolerance)

```

Figure 24. Raw PTX docs on Hopper at turn 201 (the Figure 149 hallucination). The agent infers the WGMMMA accumulator layout from text alone, producing 1.92% correctness.

`inlined_gemm` Hopper to 4.4 / 16.6 against  $\sim 298$ .

**Cluster 2. Missing API ground truth.** The agent knows the optimization exists but cannot construct a correct invocation, so descriptor immediates, barrier byte counts, encodings, and swizzle modes are subtly wrong and the kernel compiles to silent garbage or fails on a hallucinated intrinsic. Raw PTX docs is the canonical case. The Hopper turn-201 “Figure 149” hallucination (Figure 24) is the cleanest example, where the agent infers a WGMMMA accumulator layout from text alone and produces 1.92% correctness at 77.9 TFLOPS. `ptx_abstractions` on Blackwell exhibits the brittle compile-floor variant where the agent invents `__pack_half2` at turn 32 (Figure 23), retries variants for 167 turns because the `epilogue/guide.md` with the correct `stmatrix::store_b16_x4` signature is not preloaded, and terminates at `best_tflops = null`.

**Cluster 3. Missing optimization patterns.** The agent has correct primitive APIs and the right algorithm but does not know which combinations compose into the throughput-critical pattern, so the kernel is correct but capped and either contains an explicit comment acknowledging the missing optimization or a hand-rolled re-derivation. `ptx_abstractions` on Hopper and CDNA4, augmented PTX docs, and vendor FMHA all land here. The Hopper smoking gun is the literal “// Single warpgroup for now” in `ptx_abstractions` at turn 150 (Section D.8), and the CDNA4 smoking gun is the LDS XOR swizzle re-derived at six call sites.

**Cross-model robustness.** The two backbones disagree quantitatively but agree structurally, with cluster assignments stable across both models. Every baseline that collapses under Gemini also collapses under GPT in the same mode, whether compile floor, hallucinated intrinsic, or plateaued correct kernel, and HAWKEYE is the only context where both LLMs scale to expert-relative parity on every architecture.

**Token-cost framing.** Across the Linear-Attention sweep, baselines spend 120k–2216k trajectory tokens per TFLOPS of best result while HAWKEYE spends 24–50k, with the token band structurally bounded on both sides. Below  $\sim 5$ k turn-0 tokens (skill evolution, `ptx_abstractions`, raw PTX docs) the agent has primitives but no organizing principle and clusters 2 and 3 dominate. Above  $\sim 25$ k turn-0 tokens (`inlined_gemm` on Hopper at 32k and Blackwell at 35k) the exemplar’s algorithmic shape is hard-coded into the syntax and cluster 1 dominates. Above  $\sim 100$ k workspace tokens (raw PTX docs at 1.22M), doc-grep dominates and cluster 2 returns. HAWKEYE’s 11.8k–20.1k turn-0 band sits in the trough.

### D.3 Triton FLA

The `triton` baseline gives the agent a Triton expert reference for causal Parallel ReBased Linear Attention plus the standard Triton DSL, and the agent must emit CUDA or HIP. This is the cleanest missing-algorithm-structure probe, since the algorithm is fully described but every architecture-specific primitive is hidden behind `tl.dot` and the lowering tables.

```

1375 wmma::fragment<wmma::matrix_a, 16, 16, 16, bf16, wmma::row_major> q_frag;
1376 wmma::fragment<wmma::matrix_b, 16, 16, 16, bf16, wmma::col_major> k_frag;
1377 wmma::fragment<wmma::accumulator, 16, 16, 16, float> s_frag;
1378
1379 wmma::load_matrix_sync(q_frag, &sq[...], HEAD_DIM);
1380 wmma::load_matrix_sync(k_frag, &sk[...], HEAD_DIM);
1381 wmma::mma_sync(s_frag, q_frag, k_frag, s_frag);

```

Figure 25. Triton FLA on Blackwell at turn 149. The agent uses the `nvcuda::wmma` fragment API on every NVIDIA arch because that is the only tensor-core path the Triton reference exposes, leaving `tcgen05.mma`, `TMEM`, and `CLC` unused.

```

1385 T49: "The mma.sync.aligned.m16n8k16.row.col.f32.bf16.bf16.f32
1386 instruction compiles successfully. I will use this to replace wmma."
1387
1388 T93: 21.6 TFLOPS reported, correctness = 3.57%
1389 (the kernel compiles but ldmatrix lane-mapping is wrong)
1390
1391 T107--T131: sustained 21--22 TFLOPS at constant correctness ~1.84%
1392
1393 T232: revert to T30 (WMMA kernel at 17.9 TFLOPS)

```

Figure 26. Triton FLA on Hopper, turns 49–232. The agent identifies that WMMA is the ceiling and tries to hand-roll `mma.sync` with `ldmatrix`, but cannot derive correct lane-to-element addressing without a worked example and reverts to WMMA after 44 turns of incorrect kernels.

The Blackwell run lands on `nvcuda::wmma` at turn 16 and iterates that same kernel for 220 turns, peaking at 110.4 TFLOPS at turn 149 (0.26× the 424.7 TFLOPS expert) on the WMMA fragment API that the Triton lowering compiles to on every NVIDIA arch (Figure 25). `tl.dot` would expand on Blackwell to `tcgen05.mma.cta_group:1.kind:f16` with `TMEM`, `CLC`, and multicast TMA, but Triton hides those and the workspace has no syntactic anchor, so across 237 turns and 9.04 M tokens the kernel never advances past WMMA.

Hopper reproduces the same WMMA fallback (peaks at 17.9, 29.2, 106.3 TFLOPS vs 298 expert), and the agent recognizes the ceiling at turn 31 (“WMMA is slow”) but its 44-turn attempt to port to `mma.sync` with `ldmatrix` regresses to 1.84% correctness (Figure 26) before reverting to the 17.9 TFLOPS WMMA kernel for the remaining 150 turns. Ampere is the muted case because its tensor-core path is `nvcuda::wmma` (27.2 and 36.3 TFLOPS vs 132 expert), which shows the DSL’s implicit fallback target rather than DSL competence on Ampere. CDNA4 never passes the correctness gate (~72% within tolerance vs the required 1.0), the running log at turns 190–193 reads “Still failing with `opaque_abort`” with no further hypothesis, and the MFMA intrinsic `__builtin_amdgcn_mfma_f32_16x16x16bf16_1k` appears nowhere in the Triton reference for the agent to copy.

GPT-5.4 lands in the same WMMA fallback at lower absolute peaks (Blackwell 34.6, Hopper 16.2, Ampere 4.4, CDNA4 0.2 TFLOPS), confirming the failure is structural rather than a Gemini idiosyncrasy. The agent receives a complete-looking kernel but cannot see which primitives the DSL silently drops on each architecture, and HAWKEYE surfaces those primitives as named cells with composition order baked in, recovering 5.0× on Blackwell against the same agent and model (Table 6).

#### D.4 Raw PTX and ROCm docs

The `raw_ptx_docs` baseline gives the agent the full PTX ISA (1.22 M tokens) or ROCm ISA (~0.31 M tokens) on disk, accessed by a productive-turn-counted `read_context` tool. The bottleneck is navigation within budget, not corpus content.

Hopper peaks at 31.6 TFLOPS at turn 82 (0.11× expert) by re-issuing the same WGMMA-pattern `grep` across 7 distinct turns (T7, 9, 11, 53, 90, 92, 95), each costing a productive turn while growing trajectory tokens without growing kernel state. The agent retrieves the right symbol `cp.async.bulk.tensor`, the right operand (`tensormap`), and the right reason (swizzle), then declines to construct one and falls back to Ampere-style `cp.async.cg`. Another run reaches 88.6 TFLOPS at turn 118 (0.30×) by adding a WGMMA-shaped block tile but never wires up TMA.

Blackwell is the same pathology made worse, with peaks at 5.6 TFLOPS (turn 4) and 12.0 TFLOPS (turn 84) versus the 424

```

1430 T17: "I will read the PTX ISA for wgmma first, because it's simpler and
1431     I'm more familiar with it. WGMMMA is still usable on Blackwell."
1432
1433 T23: "The tcgen05.mma instruction is quite complex. It requires TMEM
1434     allocation, descriptors, and mbarrier synchronization. Since I have
1435     100 turns, I can try WGMMMA first and switch to tcgen05 if needed."
1436
1437 T64: "The kernel is now correct! However the performance is very poor
1438     (0.9 TFLOPS, 0.00x expert). This is because the kernel is doing a
1439     lot of scalar operations inside the warp-group loop."
1440
1441 T124: "The current kernel is correct but very slow (0.9 TFLOPS vs 424 for
1442     the expert)." -- agent does not pivot to tcgen05

```

Figure 27. Raw PTX docs on Blackwell. The agent reasons that `tcgen05` is too complex, commits to WGMMMA from prior knowledge, and across 221 turns attributes the resulting 0.9 TFLOPS to “scalar operations” rather than the wrong instruction family.

```

1445 WGMMMA-pattern read_context() calls on the PTX ISA, Blackwell:
1446 T17, T19, T21, T29, T33, T35, T37, T39, T64, T66, T68, T74, T83, T85,
1447 T87, T127, T129 -- 17 distinct productive turns
1448
1449 A representative grep:
1450 read_context("blackwell/full_docs/ptx_isa.md", pattern=r"^### .*wgmma")
1451     -> 0 hits (wgmma is buried in section 9.7.15.5.2, not a heading)
1452
1453 T19: "The pattern '^### .*wgmma' didn't match anything. Let's try a
1454     simpler pattern. Let's search for wgmma.mma_async."

```

Figure 28. Raw PTX docs on Blackwell, the WGMMMA grep loop. The PTX ISA’s structure is not indexed for the symbols the agent searches, so 17 productive turns are spent re-issuing pattern variants and never reach `tcgen05`.

expert (0.013× and 0.028×). At turn 23 the agent reasons that `tcgen05` is too complex and commits to a WGMMMA-first plan it never abandons (Figure 27), then re-issues the same WGMMMA-pattern grep across 17 distinct turns (Figure 28), and by turn 124 reverts to the naive kernel (`best_turn = 4`). The string-match bake-off favors the larger WGMMMA, `cp.async`, and WMMA chapter cluster over the smaller `tcgen05` addendum. The corpus is correct and the navigation policy is wrong.

CDNA4 is the inverted case and the strongest `raw_ptx_docs` result (the workspace name is kept for consistency; on CDNA4 the corpus is the ROCm ISA rather than PTX), peaking at 169.7 TFLOPS at turn 240 (0.51× expert) because navigation cost is monotone in corpus size. The CDNA4 ISA is ~0.31 M tokens vs PTX’s 1.22 M, the MFMA section is one contiguous block, and the agent finds `__builtin_amdgcn_mfma_f32_16x16x16bf16_1k` at turn ~170. Even on the small corpus, navigation still bites in the long tail, with the same run missing `global_load_lds` (in the memory-instruction chapter while the agent anchors in the MFMA chapter) and falsely reporting at turn 197 that the symbol does not exist. GPT-5.4 reproduces the same plateau at lower peaks (Hopper 9.2, Blackwell 19.8, Ampere 19.6, CDNA4 16.0 TFLOPS), with the CDNA4-easiest-to-navigate ranking surviving the model swap, so the bottleneck is navigation rather than content. HAWKEYE closes this by replacing free-form retrieval with named cells the agent addresses by recipe, recovering the largest single-cell gap in the matrix on Blackwell (46× behind `raw_ptx_docs`, with 12.36 M post-peak tokens spent for zero improvement, Table 11).

### D.5 Augmented documentation

The `augmented_ptx_docs` baseline gives the agent curated, hyperlinked PTX or HIP ISA with TL;DR blocks and call-sequence templates per primitive. Curation amortizes navigation but does not provide the cross-primitive recipe, and the result is the most expensive baseline (13.6 M Gemini Hopper tokens, 13.0–13.4 M on Blackwell, 17.7 M on CDNA4) and one of the worst-performing (4.4 / 5.6 / 118.4 Gemini TFLOPS). The failure mode is composition, not navigation or syntax.

The Hopper Gemini run (100 productive turns, 13.97 M tokens, peak 4.4 TFLOPS at turn 4 vs 298 expert) is the cleanest case. The agent identifies WGMMMA, TMA, and `mbarrier` primitives correctly through turn 16 and confirms the

```

1485 T78: attempt wgmma -> correctness fail (harness_abort) -> revert
1486 T85: attempt wgmma -> correctness fail (harness_abort) -> revert
1487 T95: attempt wgmma -> correctness fail (harness_abort) -> revert
1488 T111: attempt wgmma -> correctness fail (harness_abort) -> revert
1489 ...
1490 T199: best_turn still 4 (the naive kernel), 13.97M tokens spent

```

Figure 29. Augmented PTX docs on Hopper, the canonical-sequence over-commitment trap. The augmented wgmma doc presents one “canonical” call sequence with no fallback path, so each failed attempt reverts to the naive kernel rather than degrading to a simpler tensor-core path, and the 4.4 TFLOPS naive scalar survives 195 turns.

```

1495 T163 compile error:
1496   ptxas /tmp/...: error: Arguments mismatch for instruction 'tcgen05.mma'
1497
1498 The augmented blackwell/tensor_cores_and_tmem.md TL;DR lists the API
1499 names (tcgen05.alloc, tcgen05.mma, tcgen05.commit, tcgen05.ld,
1500 tcgen05.dealloc) but does not inline the descriptor encoding:
1501   - bits 0--2 : op type
1502   - bits 3--5 : shape selector
1503   - bits 6--8 : data type
1504   - bits 9--15: trans / smulo / block-scaled flags
1505 The agent must read_context() into the doc for these tables, and each
1506 attempt rebuilds the bit layout from scratch.

```

Figure 30. Augmented PTX docs on Blackwell at turn 163. The agent has the tcgen05 API names from the TL;DR but reconstructs the instruction descriptor bit layout by trial and error, producing repeated ptxas argument-mismatch errors and reverting to the naive kernel.

right WGMMA shape at turn 39, but cannot compose them. The augmented tensor\_cores\_and\_warpgroup.md frames the WGMMA call sequence as the canonical path, so the agent commits to that sequence and reverts to the naive kernel after every failed attempt rather than degrading to a simpler shared-memory + WMMA path (Figure 29). The missing piece is the cross-primitive recipe pairing mbarrier.arrive.expect\_tx with cp.async.bulk.tensor....mbarrier::complete\_tx::bytes on the same barrier with phase-bit flips, which the augmented docs treat as separate topics, and the run terminates at best\_turn = 4.

Blackwell is the same shape with a tcgen05 ceiling. Both runs terminate at best\_turn = 3-4 (peak 5.6 TFLOPS, the naive kernel). The agent writes tcgen05.mma, TMEM alloc/dealloc, and cluster TMA boilerplate but cannot decide cta\_group::1 vs ::2, TMEM column allocation 128 vs 256 vs 512, or phase-bit tracking style. The augmented doc gives the API names but not the descriptor encoding tables at turn-0 inlined level, so each compile fails on argument-mismatch errors the agent cannot resolve from the available context (Figure 30), and every tcgen05.fence placement creates a race that reverts the kernel.

CDNA4 (118.4 TFLOPS at turn 175, 17.7 M tokens, 354 raw turns, the only matrix run to hit MaxIterationsReached) is the most diagnostic case because the augmented HIP doc describes v\_mfma\_f32\_16x16x16bf16\_1k and v\_mfma\_f32\_32x32x8bf16\_1k side by side without committing to one. The agent peaks at turn 175 with 16x16, then wastes 179 turns (~9.8 M tokens) re-attempting the 32x32 shape. GPT-5.4 lands on shared-memory tiled kernels without tensor cores at all (Hopper 27.2, Blackwell 26.2, CDNA4 20.0 TFLOPS), the same cluster expressed at a different floor. Across all three architectures the bottleneck is composition rather than navigation, since each excerpt teaches one primitive correctly but none describes the producer-consumer recipe pairing TMA into an mbarrier programmed by arrive.expect\_tx, decremented by the TMA via complete\_tx::bytes, and gated by a phase bit that flips per ring slot. That recipe is what HAWKEYE’s pipeline/ module encodes (Section D.8), and augmented\_ptx\_docs spends +14% of raw\_ptx\_docs’s Hopper tokens for one-seventh the throughput (Table 6).

## D.6 Inlined GEMM exemplar

The inlined\_gemm baseline preloads a hand-written GEMM as one monolithic file (3380 Blackwell lines, ~2500 on Hopper and CDNA4, 189 on Ampere). The exemplar’s density couples algorithmic and hardware concerns tightly enough that the agent treats it as the algorithm and cannot decouple them. This is a missing-algorithm-structure probe.

```

1540 static constexpr int BM = 256;           // per cluster (128 per CTA)
1541 static constexpr int BN = 256;
1542 static constexpr int BK = 64;
1543 static constexpr int SUPERGROUP_SIZE = 4;
1544
1545 static constexpr int CLUSTER_SIZE      = 2;
1546 static constexpr int NUM_CONSUMERS     = 2; // non-overlap mode
1547 static constexpr int NUM_PRODUCERS     = 1;
1548 static constexpr int NUM_WARPGROUPS   = NUM_CONSUMERS + NUM_PRODUCERS; // 3
1549 static constexpr int LOAD_PIPE_DEPTH  = 4;
1550 static constexpr int EPI_PIPE_DEPTH   = 8; // 8 chunks of 32 cols
1551 static constexpr int CTA_M            = BM / CLUSTER_SIZE; // 128
1552 static constexpr int A_ROWS           = CTA_M; // 128
1553 static constexpr int B_ROWS           = BN / CLUSTER_SIZE; // 128
1554 static constexpr int K_CHUNKS        = BK / 16; // 4

```

Figure 31. Blackwell BF16 GEMM exemplar configuration.

```

1555
1556
1557
1558 // inlined_gemm helper (immutable, 2500-line monolithic file):
1559 uint64_t desc_a = make_smem_desc(&sA[0], 16, 1024, true); // <-- 1024
1560 wgmma64<1, 1, 1, 0, 1>(s_acc, sQ + k_it * 16, sK + k_it * 16);
1561
1562 // T28 agent reasoning:
1563 //   The stride of sQ is HEAD_DIM elements = 64 elements = 128 bytes.
1564 //   But make_smem_desc uses 1024 as the stride multiplier.
1565 //   1024 bytes = 512 elements. wgmma64 hardcodes 1024.
1566 //   I need to change the stride in wgmma64 or write my own wrapper.

```

Figure 32. Inlined GEMM on Hopper at turn 28. The `wgmma64` helper hardcodes a 1024-byte descriptor stride for the GEMM tile, but causal Parallel ReBased Linear Attention needs 128 bytes. The helper is inlined and immutable, and 180 turns of from-scratch wrapper attempts produce only fresh PTX errors.

The Blackwell exemplar is a 470-line 2-CTA cluster BF16 GEMM with the configuration in Figure 31.

This configuration is correct for square BF16 GEMM but incorrect for causal Linear Attention, where the tile is  $S$ -by- $D$ ,  $D=64$  is fixed, and the inter-chunk linear-state recurrence has no GEMM analog. Gemini Blackwell (100 productive turns, 17.27 M tokens, peak 5.8 TFLOPS at turn 173) is the exemplar trap. The agent copies the GEMM TMEM/fence/dealloc and `mbarrier.expect_tx` scaffolding verbatim, crops BM/BN/BK while keeping the GEMM’s stride immediates, and falls back to a scalar inner loop for the recurrence. The resulting kernel compiles and is correct (every running line comes from the naive reference), but TMEM is allocated and immediately deallocated, the compute path bypasses it, and the producer/consumer scaffolding adds cost without throughput. Other runs yield `best_tflops = null`.

Hopper exposes the same coupling at the helper-function level. The exemplar ships a `wgmma64` helper with a hardcoded 1024-byte stride for descriptor B, correct for the 64-column GEMM tile but wrong for causal Linear Attention with `HEAD_DIM=64` and 128-row K tiles. The agent identifies the mismatch at turn 28 (Figure 32) but cannot edit the inlined helper without rewriting all 2500 lines, and from turn 29 to 208 either the wrong stride silently produces garbage or a from-scratch WGMMMA wrapper introduces fresh PTX errors. The run plateaus at `best_turn = 4`, 4.4 TFLOPS, 24.18 M tokens.

The exemplar entangles algorithmic shape, tile geometry, descriptor strides, warpgroup split, and persistent CLC scheduling, and porting to Linear Attention requires moving four of those five in coordination without an abstraction boundary that lets the agent move one lever at a time. GPT-5.4 refuses the WGMMMA and `tcgen05` paths entirely, applies shared-memory tiling to the naive kernel and stops (Blackwell 15.0, Hopper 16.6, CDNA4 2.1 TFLOPS, Table 6). This is the dual of the Triton FLA failure, with hardware paths present but no algorithmic decoupling.

```

1595 __global__ void naive_rebased_attention_kernel(
1596     const bf16* Q, const bf16* K, const bf16* V,
1597     bf16* O, int S) {
1598 }
1599 void optimized_rebased_attention(
1600     const bf16* Q, const bf16* K, const bf16* V,
1601     bf16* O, int B, int H, int S, int D) {
1602 }

```

Figure 33. Vendor FMHA on Blackwell at turn 64. After failing to extract `tcgen05.mma` patterns from CUTLASS source, the agent reverts both kernel bodies to literal empty stubs and the run freezes at the 5.7 TFLOPS naive baseline for 207 more turns.

```

1607 T2:    "**Examining the Kernel Code**\n\nOkay, so I need to get a clear
1608       picture of what's going on with these kernels..."
1609 T41:   "**Analyzing the Implementation for Causal Linear Attention on
1610       Blackwell**\n\nOkay, so we need to get a hardware-aware kernel..."
1611 T78:   "**My Approach to Implementing a Hardware-Aware Kernel for Causal
1612       Linear Attention**\n\nOkay, so the challenge is clear..."
1613 T117:  "Here's my analysis of the provided files, keeping in mind my
1614       expert perspective:\n\n**Kernel Optimization Comparison**..."
1615 T229:  "**Navigating the PTX Maze for Blackwell Kernel Optimization**
1616       \n\nOkay, here's the situation. I'm deep in the weeds..."

```

Figure 34. Reasoning openers from `vendor_fmha/blackwell/conversation.jsonl`, turns 2, 41, 78, 117, and 229. Every multi-turn reasoning block opens with a navigation verb. The action the reasoning produces is another `read_context` call, not a kernel edit.

## D.7 Vendor FMHA

The `vendor_fmha` baseline gives the agent CUTLASS Blackwell and Hopper FMHA (~854 KB) and Composable Kernel `ck_tile/01_fmha` (~287 KB) on the workspace. Both are aggressively abstracted with CuTe layouts, collective builders, and warp-specialized templates that present optimization patterns as objects of composition rather than as exemplars of primitives, so the agent loses turns to navigation. This is a missing-optimization-pattern probe.

The Blackwell run is canonical (Gemini, 328 raw turns, 100 productive, 13.10 M tokens, peak 5.7 TFLOPS at turn 121). The agent identifies the right customization point at turn 11 (replace `softmax_step` with the Taylor feature map), spends turns 25–66 discovering ground-truth violations (CUTLASS not on include path, `blackwell-apis` missing), then reverts both kernel bodies to literal empty stubs at turn 64 (Figure 33) and freezes at 5.7 TFLOPS for the remaining 207 turns. Another run fights 408 raw turns and peaks at 24.1 TFLOPS (0.043×).

CDNA4 reduces `vendor_fmha` to `raw_ptx_docs` after 27 turns of CK navigation. The agent `grep`-cycles for the CK kernel entry point through turns 7–22, then at turn 27 searches CK for “MFMA” and gets zero hits because the calls are hidden inside `ck_tile`, and at turn 28 abandons CK and hand-rolls raw HIP MFMA. Peak is 39.2 TFLOPS (0.118×) at turn 237 from `__builtin_amdgcn_mfma_f32_16x16x16bf16_1k` without CK pipelining. The Hopper CUTLASS run (309 raw / 100 productive turns, 13.77 M tokens, peak 29.7 TFLOPS at turn 260, 0.10×) is the only `vendor_fmha` run with a non-trivial kernel. The agent abandons the include-path attempt at turn 75 and hand-rolls a chunked  $O(S^2)$  kernel, but lacks the inter-chunk linear-state accumulator that distinguishes causal Linear Attention from FMHA.

The structural symptom across architectures is that every multi-turn reasoning block opens with a navigation verb rather than a kernel edit (Figure 34).

The action each turn produces is another `read_context` call against vendor source rather than a kernel edit. Production GPU libraries are written to be reused, not read, and CUTLASS’s `collective::Builder` and CK’s `ck_tile::FmhaFwdKernel` hide primitive recipes under reusable layers that are inverted for an LLM that needs to extract those recipes. GPT-5.4 follows the same shape (Hopper 37.6, Blackwell 18.9, CDNA4 0.8 TFLOPS, Table 6), and HAWKEYE’s `pipeline/` and `epilogue/` modules expose the same primitives as named single-purpose recipes that the agent does not have to navigate to.

D.8 PTX abstractions ablation

We start from the full HAWKEYE taxonomy and surgically remove the optimization-pattern modules while keeping the primitive-instruction wrappers, so the agent retains thin C++ wrappers around tensor-core, memory, and synchronization intrinsics but loses every multi-instruction composition pattern. Removed modules encode ring-buffered async pipelines, producer-consumer warpgroup splits, decoupled epilogues, persistent tile schedulers, and composed attention helpers. Ampere and Hopper lose only math helpers (swizzle), Blackwell loses epilogue, pipeline, clc, tile\_scheduler, attention\_ops, swizzle, nvfp4, and CDNA4 loses memory, gemm, attention, dpp, fp4, conv. The prediction is a small gap on Ampere and Hopper and a large gap on Blackwell and CDNA4, which Tables 7 and 8 confirm.

Table 7. ptx\_abstractions versus full HAWKEYE on causal Parallel ReBased Linear Attention, Gemini 3.1 Pro backbone, 100-productive-turn budget.

Arch	ptx_abstractions	HAWKEYE full	Gap	Best turn (ptx)
Ampere	84.7 TFLOPS	87.3 TFLOPS	+3%	205
Hopper	240.9 TFLOPS	256.4 TFLOPS	+6%	150
Blackwell	null (compile-fail)	702.0 TFLOPS	catastrophic	—
CDNA4	241.8 TFLOPS	326.5 TFLOPS	+35%	256

Table 8. ptx\_abstractions versus full HAWKEYE under the GPT-5.4 backbone. The catastrophic Blackwell collapse is reproduced under GPT (9.7 TFLOPS, 44× behind HAWKEYE), and a Hopper collapse appears (17.5 vs 327.9) that is more severe than under Gemini.

Arch	ptx_abstractions	HAWKEYE full	Gap	Best turn (ptx)
Ampere	10.1 TFLOPS	135.8 TFLOPS	13.4×	104
Hopper	17.5 TFLOPS	327.9 TFLOPS	18.8×	124
Blackwell	9.7 TFLOPS	418.9 TFLOPS	43.4×	171
CDNA4	11.2 TFLOPS	319.2 TFLOPS	28.6×	114

Under Gemini, ptx\_abstractions on Ampere and Hopper recovers most of HAWKEYE’s performance with a 6–35% residual reconstruction gap, but under GPT the gap explodes by an order of magnitude on every architecture, so the optimization-pattern modules act as an LLM-capability multiplier rather than a fixed bonus.

Gemini ptx\_abstractions on Hopper reaches 240.9 TFLOPS at turn 150 with a single-warpgroup body that contains the literal admission “// Single warpgroup for now” (Figure 35). The agent recognizes that the producer-consumer split is the next move but cannot synthesize the ring-buffer phase-bit protocol from barrier/ within budget. The kernel is correct but single-warpgroup, single-stage, and serializes both WGMMAs against the same TMA chain, plateauing at ~241 TFLOPS while HAWKEYE reaches 256.4 by composing through pipeline/, a single markdown file the agent burns 100 turns rediscovering 60% of.

CDNA4 is the cleanest A/B because both kernels compile and produce real performance with no compile-floor confounder, so the 84.7 TFLOPS (35%) gap is purely the six excluded modules. The agent re-derives the LDS XOR swizzle (col / 8) ^ (row % 8) at six call sites (Figure 36), each of which would invoke a single helper from HAWKEYE’s memory/ module. The kernel is also single-role (all 512 threads do compute and memory together) where HAWKEYE overlaps K and V loads with S and O via memory/’s async\_pipeline, so the 35% gap is a clean lower bound on the six removed modules.

The kernel is also single-role (all 512 threads do compute and memory together) while HAWKEYE overlaps K/V loads with S and O via memory/’s async\_pipeline. The 35% gap is a clean lower bound on the six removed modules.

Blackwell layers a brittle floor on top of this ceiling, with Gemini hallucinating \_\_pack\_half2 at turn 32 and terminating at best\_tflops = null, and under GPT the floor softens to 9.7 TFLOPS at turn 171 (vs 418.9 HAWKEYE) but composition still fails. Abstraction-budget scaling curves over  $K \in \{1, 2, 4, 6\}$  optimization-pattern modules rise monotonically from ptx\_abstractions ( $K=0$ ) to full HAWKEYE ( $K=6$ ) on Blackwell under both backbones, with the steepest slope in the first three modules, so the contribution is a property of the context rather than the LLM. The wins are not from unit tests, intrinsic wrappers, or the workload guide, all of which are retained in ptx\_abstractions and still cap at ~241 TFLOPS on Hopper. The only structural difference is the optimization-pattern modules, and under GPT those produce a 13–43× gap on every architecture, with the cleanest single-arch comparisons being CDNA4 under Gemini (35%) and CDNA4 under GPT (28.6×). Tables 9 and 10 trace turn-by-turn comparisons, and Table 11 reports the two most catastrophic post-peak runs.

```

1705
1706
1707
1708
1709
1710 if (wg_id == 0) {
1711     // Single warpgroup for now
1712     warpgroup::reg_alloc<232>();
1713     // Load Q via TMA, then loop over K, V tiles
1714     for (int s_block = 0; s_block <= t_block; s_block++) {
1715         // ... TMA loads K, V into smem ...
1716         barrier::wait(&smem.full_k, phase_k); phase_k ^= 1;
1717         barrier::wait(&smem.full_v, phase_v); phase_v ^= 1;
1718
1719         warpgroup::fence(); // S = Q @ K^T
1720         for (int k = 0; k < HEAD_DIM / 16; k++)
1721             wgmma_bf16::mma_async<BK, ...>(s_acc, dq, dk);
1722         warpgroup::commit();
1723         warpgroup::wait<0>();
1724
1725         // ... causal mask + feature map ...
1726
1727         warpgroup::fence(); // O += S @ V
1728         for (int k = 0; k < BK / 16; k++)
1729             wgmma_bf16::mma_async<HEAD_DIM, ...>(o_acc, ds, dv);
1730         warpgroup::commit();
1731         warpgroup::wait<0>();
1732     }
1733 }

```

Figure 35. ptx\_abstractions on Hopper at turn 150 (the smoking gun). The agent has manually unpacked the producer-consumer roles into a single warpgroup with explicit warpgroup::fence/commit/wait calls between WGMMAs. The Hawkeye pipeline/module ships these three calls as a single ring\_advance helper.

```

1746 Line 54: int swizzle_col0 = (col0 / 8) ^ (row0 % 8);
1747 Line 83: int swizzle_col = (col_start / 8) ^ (row % 8);
1748 Line 121: int swizzle_col = (global_col / 8) ^ (global_row % 8);
1749 Line 134: int swizzle_col_s = (col_s_start / 8) ^ (row_s % 8);
1750 Line 166: int swizzle_col = (global_col / 8) ^ (global_row % 8);
1751 Line 175: int o_swizzle_col0 = (o_col0 / 8) ^ (o_row0 % 8);

```

Figure 36. ptx\_abstractions on CDNA4 at turn 256. Same XOR swizzle re-derived at six call sites. With HAWKEYE’s memory/module each site invokes a single helper.

Table 9. Trajectory comparison on Blackwell, same agent and model. Tokens are cumulative trajectory tokens. HAWKEYE climbs monotonically through the rung ladder. raw\_ptx\_docs burns 11.6 M tokens after turn 25 with zero improvement.

Turn	Raw PTX docs		HAWKEYE	
	Tokens	Best TFLOPS	Tokens	Best TFLOPS
1	35k	—	102k	—
25	850k	5.6	1.98M	<b>576.1</b>
50	2.06M	5.6	4.99M	576.1
100	4.72M	5.6	12.97M	701.0
123	5.78M	5.6	17.30M	<b>702.0</b>
150	7.38M	5.6	22.77M	702.0
200	10.36M	5.6	—	—
221	12.43M	5.6	—	—

Raw PTX docs after turn 4 is textbook context pollution. The agent issues 17 distinct wgmma-pattern greps and accumulates each result in trajectory text. HAWKEYE climbs monotonically from T25 to T123 (576 to 702 TFLOPS), then plateaus with bounded post-peak exploration.

Table 10. Trajectory comparison on Hopper. By turn 100 (half the budget), HAWKEYE is at 184 TFLOPS, already 76% of ptx\_abstractions’s eventual best.

Turn	ptx_abstractions		HAWKEYE	
	Tokens	Best TFLOPS	Tokens	Best TFLOPS
1	46k	—	84k	—
100	5.36M	<25	9.52M	<b>183.8</b>
150	8.10M	240.9	18.99M	193.0
168	—	—	22.94M	<b>256.4</b>
203	11.16M	240.9	—	—

Table 11. Two baselines where post-peak tokens are pure context pollution. The trajectory cost grows but the kernel does not.

Run	Total tokens	Best TFLOPS	Post-peak waste
augmented_ptx_docs on CDNA4	17.7M	118.4 at T175	9.8M tokens after T175, no gain
raw_ptx_docs on Blackwell	12.4M	5.6 at T4	12.36M tokens after T4, no gain

Across the sweep, baselines spend 120k to 2216k trajectory tokens per TFLOPS of best result while HAWKEYE spends 24 to 50k.

## E Optimization stacking and composition

This appendix shows what HAWKEYE agents do with the taxonomy, anchored on the trajectory waterfall (Figure 37).

### E.1 Closed-loop pattern

Across Figure 37 the scalar-to-hardware-aware transition is discontinuous. Hopper GEMM sits at ~23 TFLOPS until one step turns on WGMMMA, TMA, mbarrier, and swizzle simultaneously and reaches 493 TFLOPS, Hopper FP8 GEMM jumps from 5.8 to 988 TFLOPS in the same step (170×), Hopper attention from 0.15 to 69.0 TFLOPS in one step, and Blackwell causal Linear Attention from 446.2 to 557.0 TFLOPS between turns 49 and 52 on top of an already-converged async pipeline. The discontinuity is structural: the Hopper features are tightly coupled, so a partial stack earns none of the speedup. WGMMMA descriptors need swizzle to read coherent data, TMA needs mbarrier for completion signaling, WGMMMA needs the warpgroup fence/commit/wait protocol to issue, and WGMMMA needs TMA to stay fed. The unit tests and kernel pool expose the full stack in one place, so the agent crosses the gap in a single phase transition rather than incrementally.

Within each trajectory the agent runs the same five-step inner loop on every turn.

- Profile.** Run `evaluate_kernel` to obtain TFLOPS, SM-busy, barrier-stall, tensor-core-issue, bank-conflict, and L2-hit-rate counters from the previous kernel snapshot.
- Identify the bottleneck.** Compare the counters to the targeted-metric column of the relevant taxonomy cell to name the next bottleneck (TC underutilized, barriers oversynchronized, wrong swizzle, etc.).

## Hawkeye: Hardware-Aware GPU Kernel Optimization with Minimal Supervision

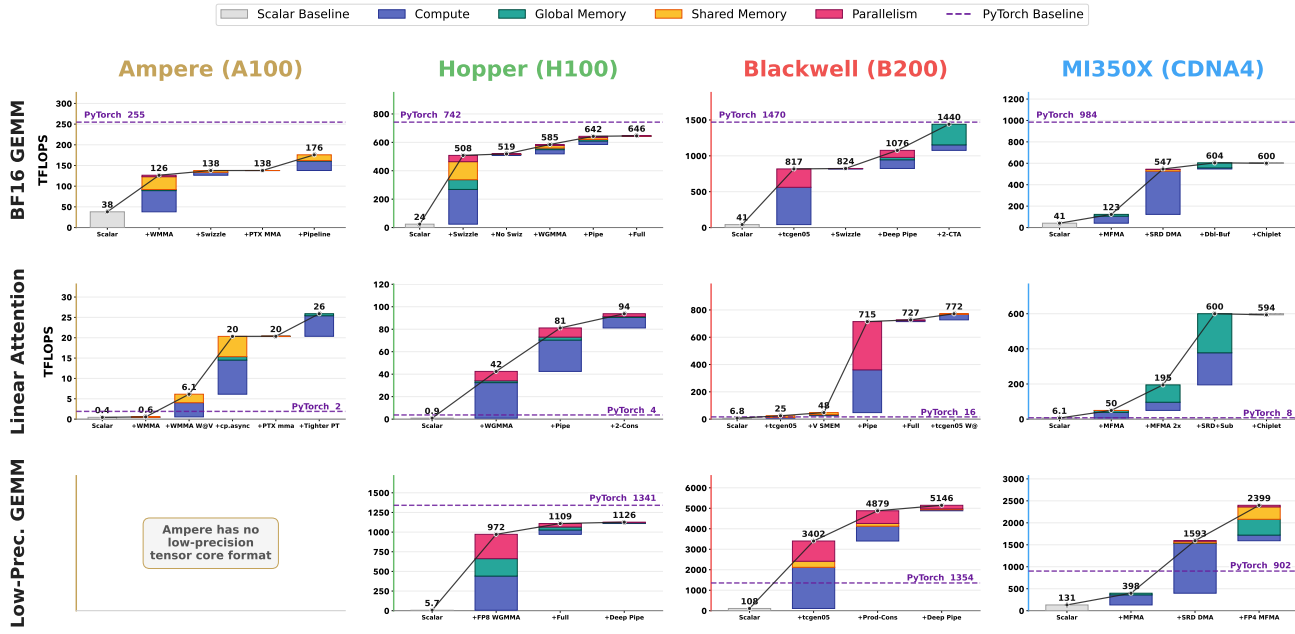


Figure 37. Optimization stacking trajectories across three NVIDIA architectures (Ampere, Hopper, Blackwell) and four workloads (BF16 GEMM, BF16 attention, BF16 causal Parallel ReBased Linear Attention, BF16/low-precision GEMM). Each bar is a cumulative optimization step in the agent’s trajectory toward the final kernel; colored segments attribute the throughput delta to a specific hardware-utilization metric drawn from the Hawkeye taxonomy. The waterfalls let us read off not just *whether* the kernel got faster but *which* part of the hardware was newly exercised.

3. **Consult the cell.** Read the matching unit-test cell (`naive_kernel.cu`, `optimized_kernel.cu`, `guide.md`) plus any API headers it pulls in.
4. **Modify the kernel.** Lift one or two primitives out of the cell and stack them onto the current kernel via `file_editor` or `terminal_writes`. Almost every modification is local.
5. **Re-evaluate.** Run `evaluate_kernel` again, observe the metric delta, and either accept the change (advance to the next rung) or revert (debug). Compile failures and correctness failures route directly back to step 3 with a more focused query.

The loop mirrors what a human kernel author runs, and the taxonomy makes step 3 cheap because the agent reads one cell whose `guide.md` already names the metric being chased.

Table 12 extracts ten representative inflection turns from the canonical Blackwell BF16 causal-Linear-Attention trajectory (Gemini 3.1 Pro, 111 turns, 50 productive, peak 557.0 TFLOPS at turn 52,  $\sim 5.0\times$  over the expert reference and  $\sim 74\times$  over the Torch baseline). The trajectory’s `decision` field records the agent’s own naming of the bottleneck and the cell it consulted, and `eval_result.tflops` records the metric delta on the next `evaluate_kernel` call.

Every productive jump (T49, T52) is preceded by a bottleneck that is identified and one cell that the agent consults to fix it, while the long T59–T103 ring-buffer debug arc oscillates across `pipeline/`, `barrier/`, and `general.md` because no single cell describes `barrier::init` parity. Profitable phase transitions live where one cell explains one bottleneck.

Optimization is non-monotonic, with roughly a third of trajectories regressing at least once. On Hopper GEMM, CTA swizzle on a shallow `QSIZE=1` pipeline drops 504 TFLOPS to 493 (2%) and the same swizzle adds 9% at `QSIZE=3`, on Ampere attention an early shuffle softmax with the wrong reduction width sits at 2 TFLOPS until the agent finds the 4-lane reduction and jumps to 22 TFLOPS in two steps, and T58 of the Blackwell trajectory above regresses from 557.0 to 335.5 because the producer-consumer split was introduced without the ring buffer. The barrier-stall, tensor-core-issue, and bank-conflict counters let the agent tell a productive transition from a misleading one, since scalar Hopper GEMM reads 73% SM-busy from scalar FMAs and only the tensor-core counter exposes the truth. The optimization moves live on a partial-order lattice where some edges are blocked and others commute, with marginal TFLOPS dropping sharply after three

## Hawkeye: Hardware-Aware GPU Kernel Optimization with Minimal Supervision

Turn	Bottleneck	Cell consulted	TFLOPS
T01	Rung 0 hangs / all-zero output	prior-run replay	—
T05	Minimal correct <code>tcgen05</code> kernel	<code>tmem.cuh</code>	—
T10	TMA tensor-map signature	<code>tma.cuh</code>	—
T22–T34	TMA coords + 128B swizzle illegal mem	<code>tma, swizzle, descriptor</code>	—
T44	Barrier phase deadlock	<code>pipeline phase-bit guide</code>	—
T49	Rung 0 correct	—	446.2
T50	Triangular work imbalance	<code>reverse query-chunk iteration</code>	—
T52	Schedule reversed	—	<b>557.0</b>
T53	Single-WG: no TMA/MMA overlap	<code>06_producer_consumer</code>	335.5
T59–T75	Ring-buffer phase-bit semantics	<code>pipeline/get_phasebit</code>	—
T87–T103	<code>barrier::init_array</code> arity	<code>general.md</code>	—
T111	Run ends, T52 best	—	557.0

Table 12. Representative inflection turns from the Blackwell BF16 causal Parallel ReBased Linear Attention trajectory (Gemini 3.1 Pro, 111 turns, peak 557.0 TFLOPS at T52). *Cell consulted* names the unit-test cell or API header the agent cited; TFLOPS shows the next correctness-passing evaluation. Dashes mark debug turns where the kernel still failed.

or four cells are stacked, and the rest of this section instantiates this closed loop on four saved trajectories.

### E.2 Trajectories across workloads

We walk through four saved Gemini 3.1 Pro trajectories, one per architecture on Parallel ReBased Linear Attention. The same closed loop recurs across all four, and only the activated cells differ.

**Ampere BF16 Parallel ReBased Linear Attention (Gemini 3.1 Pro, 63 turns, 41 productive, peak 99.9 TFLOPS at T59).** The trajectory walks Rungs 0–6 cleanly, activating `mma/`, `async_copy/`, `swizzle/`, and the Taylor fragment. The biggest jump is at T53, where a register-chained  $W \cdot V$  more than doubles throughput from 45.0 to 99.3 TFLOPS in one edit, and the final 99.9 TFLOPS is  $57\times$  over Torch.

Turn	Bottleneck	Cell consulted	TFLOPS
T16	WMMA fragment store NaN: PAD too small for HD=64	<code>mma, swizzle</code>	7.8
T23	Warps racing on shared K,V buffers	<code>async_copy</code> ( <code>__syncthreads</code> after continue)	26.5
T28	Synchronous Q/K/V loads kill MMA throughput	<code>ampere::async_copy::load_tile_padded</code>	43.2
T38	Single-buffered K,V; no overlap	<code>double-buffered cp.async</code>	45.0
T53	$W$ materialized in SMEM, $W \cdot V$ scalar	<code>mma_ptx::wv_chain_k8</code> (register-chained MMA)	<b>99.3</b>
T59	Triangular work imbalance for causal mask	<code>causal-mask exit</code> (Rung 5)	<b>99.9</b>

Table 13. Ampere BF16 Parallel ReBased Linear Attention trajectory inflection turns.

**Hopper BF16 Parallel ReBased Linear Attention (Gemini 3.1 Pro, 113 turns, 50 productive, peak 146.9 TFLOPS at T88).** The kernel sits at 4.4 TFLOPS for 81 productive turns while the agent works through Rung 1 TMA bugs, and the WGMMA replacement at T87 lifts throughput  $33\times$  to 146.9 TFLOPS in one edit. Cells activated are `tma/`, `swizzle/`, `wgmma/wgmma_bf16`, and `wgmma_descriptor/`.

Turn	Bottleneck	Cell consulted	TFLOPS
T6	Per-thread strided GMEM reads	<code>naive baseline</code> → <code>tma</code> for Rung 1	4.4
T74	128B swizzle offset reads; TMA coords as element offsets	<code>swizzle::offset</code> , <code>tma::create_tensor_map</code>	4.4
T88	Two scalar GEMMs; no tensor cores	<code>wgmma::mma_async</code> + register-A WGMMA for $P \cdot V$	<b>146.9</b>

Table 14. Hopper BF16 Parallel ReBased Linear Attention trajectory inflection turns.

**Blackwell BF16 Parallel ReBased Linear Attention (Gemini 3.1 Pro, 111 turns, 50 productive, peak 557.0 TFLOPS at T52).** The trajectory activates `tcgen05_mma/`, `tmem/`, `tma/`, `descriptor/`, `swizzle/`, `epilogue/`,

pipeline/, and barrier/, and the peak of 557.0 TFLOPS at T52 is 5.0× over expert and 74× over Torch. Turns T53–T111 attempt Rungs 2 and 3 without converging.

**CDNA4 BF16 Parallel ReBased Linear Attention (Gemini 3.1 Pro, 35 turns, 23 productive, peak 299.8 TFLOPS at T27).** This is the shortest run because the activated cells are smaller. A correct MFMA baseline reaches 292.6 TFLOPS at T6 (0.88× expert), incremental LDS double-buffering at T11 lifts throughput to 296.1, and amdgpu\_waves\_per\_eu(4, 4) at T27 reaches 299.8 TFLOPS.

Turn	Bottleneck	Cell consulted	TFLOPS
T6	Naive global loads; no MFMA	mfma::compute_2x2_bf16, Q-in-register	292.6
T11	Single LDS buffer; no overlap	cdna4-apis/l ds double-buffer	296.1
T27	Underutilized occupancy on chiplets	amdgpu_waves_per_eu(4, 4)	<b>299.8</b>

Table 15. CDNA4 BF16 Parallel ReBased Linear Attention trajectory inflection turns.

Table 16. Cross-workload trajectories. The same closed-loop pattern from the four Linear-Attention vignettes recurs in each of these runs; full per-turn data is in the run directories.

Workload	Architecture	Precision	Turns	Peak TFLOPS
Linear Attention (v4)	Hopper	BF16	231	174.5
Linear Attention (v4)	CDNA4	BF16	201	316.9
Linear Attention (v4)	Ampere	BF16	157	128.7
Linear Attention	CDNA4	MXFP4	96	384.6
Conv2D	Blackwell	BF16	65	37.5
GEMM	CDNA4	BF16	63	643.4
GEMM	CDNA4	MXFP4	140	2341.4
DeltaNet	CDNA4	MXFP4	86	4.6
Forgetting Attention	CDNA4	BF16	113	25.4

The same closed loop recurs across every workload, architecture, and precision in the saved corpus, with activated cells differing predictably. Ampere stacks mma/ + async\_copy/ + swizzle/, Hopper adds wgmma/ + tma/ + wgmma\_descriptor/, Blackwell adds tcgen05\_mma/ + tmem/ + pipeline/ + epilogue/, and CDNA4 stacks mfma/ + lds/ + global\_load\_lds/ + wave/. When the workload extends a familiar pattern, the agent adds a thin specialization on top of the shared stack and converges. When it requires composing across cells in a way no single cell describes, the closed loop oscillates without converging.

### E.3 Composition across kernels

HAWKEYE agents compose unit-test cells rather than authoring one-off arch-specific blobs. HAWKEYE produces 51 hardware-aware kernels across four architectures, six workloads, and two precision tiers (Figure 38).

Reading the matrix surfaces three structural patterns. An arch-specific bundle of cells co-occurs in nearly every kernel for that architecture (Hopper WGMMA + TMA + mbarrier + 128B swizzle, Blackwell tcgen05 + TMEM + TMA + producer-consumer, CDNA4 MFMA + LDS swizzle + global\_load\_lds + wave occupancy), forming the shared infrastructure layer that Section E.3.1 traces in detail. The remaining cells (warp or wave reduction, persistent scheduling, epilogue pipeline) attach to that infrastructure in workload-dependent combinations rather than uniformly, with reduction concentrating in attention-family kernels, persistent scheduling in GEMM and Conv2D, and epilogue pipelining where a fused activation or precision conversion lives at the kernel boundary. The matrix is dense off the diagonal in every architecture column, so no HAWKEYE kernel reaches expert performance with a single primitive in isolation, and the subsubsections below trace the recurring multi-cell patterns along the workload, architecture, precision, and invention axes.

#### E.3.1 CROSS-WORKLOAD REUSE

A small set of arch-specific primitives forms an 80–90% shared substrate per architecture (Ampere mma.sync/WGMMMA + three-stage cp.async, Hopper WGMMA + TMA + mbarrier + 128B swizzle, Blackwell adds tcgen05 + TMEM, CDNA4 mfma\_32x32x16\_bf16 + global\_load\_lds + LDS swizzle), with only a thin workload-specific layer on top. The Hopper WGMMA QK<sup>T</sup> inner loop is byte-identical between attention and Linear Attention (Figure 39), with only the tile dims and K iteration count changing.

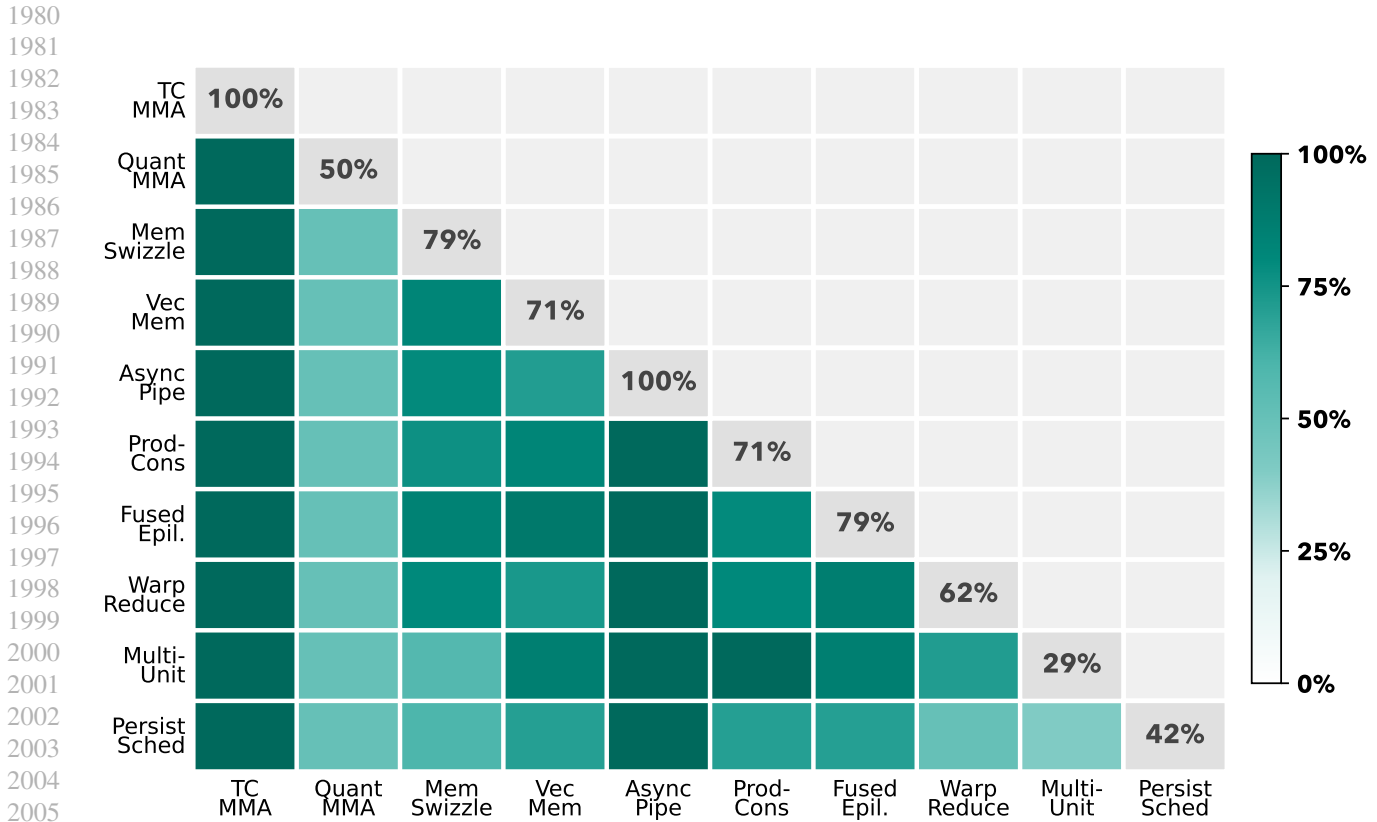


Figure 38. Co-occurrence of unit tests in HAWKEYE kernels across 6 workloads, 2 precisions, and 4 architectures. The dense off-diagonal pattern indicates that performant kernels compose several primitives drawn from across the taxonomy rather than relying on any single optimization in isolation.

```

// attention/optimized.cu (lines 219-229)
hopper::wgmma::fence_sync();
hopper::wgmma::wgmma128<0,1,1,0,0>(s_acc, sQ_base + 0*WGMMMA_K,
                                   sK_slot + 0*WGMMMA_K, 1024);

#pragma unroll
for (int k = 1; k < K_ITERS_S; ++k)
    hopper::wgmma::wgmma128<1,1,1,0,0>(s_acc,
                                        sQ_base + k*WGMMMA_K,
                                        sK_slot + k*WGMMMA_K, 1024);
hopper::wgmma::commit_group_sync();

// linear_attention/optimized.cu (lines 248-254)
wgmma::fence_sync();
bf16* sK_slot = &s.K[qi * BN * HD];
wgmma::wgmma<BN,1,1,1,0,0>(sa, &sQ_my[0*WK], &sK_slot[0*WK]);
wgmma::wgmma<BN,1,1,1,0,0>(sa, &sQ_my[1*WK], &sK_slot[1*WK]);
wgmma::wgmma<BN,1,1,1,0,0>(sa, &sQ_my[2*WK], &sK_slot[2*WK]);
wgmma::wgmma<BN,1,1,1,0,0>(sa, &sQ_my[3*WK], &sK_slot[3*WK]);
wgmma::commit_group_sync();
    
```

Figure 39. Hopper QK<sup>T</sup> inner loop. Same fence-commit-wait protocol, same descriptor-based operand binding; only the tile dims and the K iteration count change.

```

2035 // --- Ampere: WMMA fragment API ---
2036 #pragma unroll
2037 for (int ki = 0; ki < QKT; ki++) {
2038     Fbc fK; wmma::load_matrix_sync(fK, kc + ni*16*PAD + ki*16, PAD);
2039     wmma::mma_sync(sf, fQ[ki], fK, sf);
2040 }
2041 // --- Hopper: WGMMMA over SMEM descriptors ---
2042 wgmma::fence_sync();
2043 for (int ki = 0; ki < 4; ki++)
2044     wgmma::wgmma<BN,1,1,1,0,0>(sa, &sQ_my[ki*WK], &sK_slot[ki*WK]);
2045 wgmma::commit_group_sync();
2046 // --- Blackwell: tcgen05 MMA into TMEM ---
2047 tcgen05_mma::smem_fence();
2048 tcgen05_mma::k_loop_ss<K_CHUNKS_QK, CHUNK_SIZE, CHUNK_SIZE>(
2049     tmem_qk_ab[next_side], q_desc, k_desc, idesc_qk, /*acc=*/true);
2050 tcgen05_mma::commit<1>(&smem.bar_qk_done[next_side]);
2051 // --- CDNA4: MFMA with LDS-resident operands ---
2052 #pragma unroll
2053 for (int ki = 0; ki < BK; ki += 16) {
2054     bf16x8 k_op;
2055     k_rd.read_lo<BK>(sm + cur_off, ki, k_op);
2056     sync::waitcnt_lgkmcnt<0>();
2057     s_acc = mfma::bf16_32x32x16(q_cached[ki/16], k_op, s_acc);
2058 }

```

Figure 40. Parallel ReBased Linear Attention  $QK^T$  inner loop across Ampere, Hopper, Blackwell, and CDNA4. The four invocations share nothing at the instruction level, yet the agent reuses the algorithmic skeleton with the cell composition shifting to match the available primitives.

The same pattern recurs on Blackwell ping-pong TMEM (linear and Forgetting Attention), Ampere WMMA (GEMM and attention), and CDNA4 `a_subtile / b_subtile` readers (GEMM and Linear Attention). On Blackwell, three cells (`01_tensorcore_mma`, `05_async_pipeline`, `06_producer_consumer`) appear in every attention-family kernel, with the workload-specific layer (Taylor hi/lo split for Linear Attention,  $\beta$ -prescaling for DeltaNet, cumsum-bias factoring for Forgetting Attention, `im2col` for Conv2D) sitting on top.

### E.3.2 CROSS-ARCHITECTURE COMPOSITION

Linear Attention is the cleanest cross-arch test since the algorithm is fixed (Taylor map, causal mask,  $W \cdot V$  matmul) but the cell composition shifts at every step. Ampere  $\rightarrow$  Hopper activates `producer_consumer` and switches to WGMMMA, Hopper  $\rightarrow$  Blackwell adds 2-CTA cluster multicast and TMEM with `tcgen05` replacing WGMMMA, and Blackwell  $\rightarrow$  CDNA4 drops producer-consumer and substitutes LDS double-buffering with `t & 1` stage indexing. The four  $QK^T$  inner loops below (Figure 40) share nothing at the instruction level yet collapse to the same algorithmic skeleton.

The activation set grows monotonically: Ampere (5 cells)  $\rightarrow$  Hopper (6, adding producer-consumer)  $\rightarrow$  Blackwell (adding warp-reduction cell 08 for TMEM row reduction) before re-shuffling for CDNA4 (drops producer-consumer, keeps cross-lane reduction). Forgetting Attention shows the same monotone-then-shuffle pattern.

### E.3.3 CROSS-PRECISION CASCADE

A precision change cascades through cells 02 (`quantized_mma`), 03 (`swizzle_stride`), 05 (`descriptor_stride`), and 07 (epilogue dequant). On Hopper FP8 the WGMMMA K grows  $16 \rightarrow 32$ , the descriptor stride drops  $1024 \rightarrow 512$  B, the TMA col-group doubles, QSIZE halves to 2, and the epilogue gains `cvt.rn.satfinite.e4m3x2` (Figure 41). On Blackwell NVFP4 the agent pairs BF16 MMA with offline FP4 dequant, running a preprocessing kernel that converts FP4  $Q/K$  to BF16 with FP8 E4M3 per-16-block scales and pre-scales K by  $\beta$ , then reuses the BF16 DeltaNet pipeline verbatim. On CDNA4 MXFP4 the agent uses native `mfma_scale_32x32x64` with the E8M0 scale tile transposed to  $[K/32, M]$  so 512 threads populate it coalesced. Each precision-specific cell exposes the bundle as one change.

```

2090 static constexpr int WGMMMA_K_FP8 = 32; // K=32 instead of 16
2091 static constexpr int COL_GROUP_FP8 = 64; // 128 bytes / 1 byte per element
2092
2093 #pragma unroll
2094 for (int k_it = 0; k_it < HEAD_DIM / WGMMMA_K_FP8; ++k_it) {
2095     uint64_t dq = wgmma::make_descriptor(
2096         &sQ_base[k_it * WGMMMA_K_FP8], 16, 512, wgmma::Swizzle::B64);
2097     uint64_t dk = wgmma::make_descriptor(
2098         &sK[k_it * WGMMMA_K_FP8], 16, 512, wgmma::Swizzle::B64);
2099     wgmma::mma_async_e4m3<BN, 1, 1, 1>(s_acc, dq, dk);
2100 }
    
```

Figure 41. Hopper FP8 DeltaNet cascade. The K-doubling and the col-group-doubling are simultaneous, and the descriptor stride drops from 1024 to 512.

```

2104 // Without FoX algebra: every QK entry pays cq[i] - cq[j]
2105 // for (i,j) in tile: s[i,j] -= cq[q_row + i] - cq[k_row + j];
2106 // row_max = max_j s[i,j];
2107 // p[i,j] = exp(s[i,j] - row_max);
2108 //
2109 // With FoX algebra: factor cq out of the bias, track effective_max instead.
2110 // Producer fills cum_K (length BK) per chunk via a small TMA stream.
2111 float effective_max = row_max - cq; // one scalar update per row
2112 // Inside the softmax: subtract cum_K[j] from s[i,j] (already partial),
2113 // and subtract the per-row effective_max instead of (row_max + cq).
2114 p[i][j] = __expf(s[i][j] - cum_K[j] - effective_max);
    
```

Figure 42. Blackwell Forgetting-Attention FoX-algebra trick. The cumulative-decay bias  $cq$  is factored out of the per-element subtraction and folded into a row-wise `effective_max` that the softmax consumes.

The CDNA4 MXFP4 row also illustrates how the precision cell catches a correctness-only bug. The first kernel hit peak FLOPs but failed correctness because the E8M0 scale was read from a row-major  $[M, K/32]$  layout with lane-cache-line collisions, and the transposed  $[K/32, M]$  fix reaches 87% of unscaled peak and reuses across all five other CDNA4 MXFP4 workloads. The cascade is asymmetric across vendors (Hopper FP8 inverts column-group, Blackwell NVFP4 dequantizes offline, CDNA4 MXFP4 commits to native FP4 MFMA), and each precision-specific cell hides that asymmetry behind a single named change.

#### E.3.4 INVENTED PATTERNS

Several agent-invented patterns compose on top of cells without appearing in any single cell. The Ampere INT8 GEMM pre-transposes  $B$  to  $BT[N,K]$  for col-major WMMA fragments with padded SMEM stride 144, and the Hopper FP8 DeltaNet keeps  $Q$  and  $K$  as FP8 e4m3 in 8 KB SMEM with the descriptor stride dropped from 1024 to 512. The CDNA4 GEMM MXFP4 kernel uses a transposed  $[K/32, M]$  LDS layout for E8M0 block scales so 512 threads populate the tile coalesced and the lane-0/32 bank conflict disappears.

The 702 TFLOPS Blackwell Linear-Attention kernel composes five inventions on top of four cells. (a) A three-warpgroup split (compute, MMA-scheduler, TMA-producer) with asymmetric register allocations 256/40/40 instead of the canonical two-warpgroup pattern. (b) A QK lookahead that pre-issues  $QK[k+1]$  before waiting for  $W\_ready[k]$ . (c) A hi/lo BF16 split of the Taylor result for  $\sim 24$ -bit precision recovery. (d) A reverse chunk-order schedule ( $qc = (\text{num\_chunks} - 1) - \text{blockIdx.y}$ ) that puts heavy causal chunks first for load balance. (e) A Q-aliased-as-O-staging trick that frees  $\sim 16$  KB SMEM by having the epilogue stage O through the Q tile region after Q is no longer live.

The Blackwell Forgetting-Attention kernel adds an algebraic invention on top of `08_warp_reduction` (Figure 42). It factors `cumsum_q` out of the additive cumulative-log-decay bias and tracks `effective_max = row_max - cq` in TMEM instead of subtracting the per-position decay term inside the softmax loop, collapsing one full TMEM read per softmax row into a single scalar update.

```

2145 hopper::wgmma::fence_sync();
2146 hopper::wgmma::wgmma<WGMMMA_N_S, 0, 1, 1, 0, 0>(
2147     s_acc, &sQ_my[0 * WGMMMA_K], &s.K[qidx * BN * HD + 0 * WGMMMA_K]);
2148 #pragma unroll
2149 for (int k = 1; k < K_ITERS_S; ++k) {
2150     hopper::wgmma::wgmma<WGMMMA_N_S, 1, 1, 1, 0, 0>(
2151         s_acc, &sQ_my[k * WGMMMA_K], &s.K[qidx * BN * HD + k * WGMMMA_K]);
2152 }
2153 hopper::wgmma::commit_group_sync();
2154 // Prefetch cumK that overlaps both the WGMMMA group and the prior PV-half2.
2155 hopper::mbarrier::try_wait_parity(&s.ck_full[qidx], phase_ck[qidx]);
2156 phase_ck[qidx] ^= 1;
2157 float pre_ck[NSTRIPES][4];
2158 #pragma unroll
2159 for (int g = 0; g < NSTRIPES; ++g) {
2160     int col_base = g * WGMMMA_K + col_lane_pre;
2161     pre_ck[g][0] = s.cumK[qidx * BN + col_base + 0];
2162     /* ... three more cum-K loads ... */
2163 }
2164 hopper::wgmma::wait_group_sync<0>();
    
```

Figure 43. Forgetting Attention  $QK^T$  core, Hopper. One  $m64n256k16$  WGMMMA chain across  $K\_ITERS\_S=4$  K-itors, with SMEM-prefetched cumK overlapping the WGMMMA pipeline.

## F Generalization across porting axes

HAWKEYE kernels generalize along four orthogonal porting axes (Sections F.1–F.4), with Section F.5 placing the absolute numbers against the hardware ceiling.

### F.1 Cross-workload

Fix a chip, vary the workload, and the hardware-touching cells stay byte-identical across workloads in the same arithmetic family. Five Hopper BF16 kernels (attention, causal Linear Attention, DeltaNet, Forgetting Attention,  $D=128$  Linear Attention) share a one-producer / two-consumer warpgroup partition, TMA loads, an mbarrier-driven K/V ring, the fence\_sync/wgmma/commit\_group\_sync/wait\_group\_sync pipeline, and 128B XOR swizzle, with workload-specific code at 10–20% of the kernel and the inner  $Q \cdot K^T$  loop byte-identical up to BN and the accumulator name (the same comparison demonstrated in Section E.3.1). The same shape recurs on Blackwell across attention, Linear Attention, Forgetting Attention, Conv2D, and DeltaNet with a Blackwell-specific ping-pong TMEM ( $NUM\_QK\_SIDES = 2$ ) holding pre-issued  $QK[0]$  and lookahead  $QK[k+1]$  while  $PV[k]$  drains, and on CDNA4 across six workloads sharing `mfma::bf16_32x32x16`, `global_load_lds` SRD-DMA, `a_subtile/b_subtile` LDS staging, `waitcnt_lgkmcnt` barriers, and `dpp` cross-lane reduction. Shared infrastructure is 80–90% of file size on Hopper and Blackwell and 70–85% on CDNA4, ending at the QK accumulator where online softmax, the Taylor map, and the cumulative-decay broadcast become workload-specific. DeltaNet adds an intra-tile MMA for the recurrent  $W$ -carry, and Conv2D replaces the QK descriptor build with an `im2col` index walk.

### F.2 Cross-architecture

Fix the algorithm, vary the chip. The four Linear-Attention  $QK^T$  inner loops in Figure 40 differ along five orthogonal axes (tensor-core call, operand source, sync primitive, pipeline depth, cooperation unit). We triangulate with Forgetting Attention (Hopper, Blackwell, CDNA4), which adds a per-position cumulative-log-decay term  $\text{cumlog}_f[i] - \text{cumlog}_f[j]$  before softmax and exposes how each chip handles a broadcast on top of the matmul. Ampere is omitted because the broadcast does not fit SMEM at  $D=64$ ,  $S=4096$ . On Hopper the broadcast lives in an SMEM gather from `s.cumK` into registers that hides under the WGMMMA pipe (Figure 43). On Blackwell the consumer warpgroup defers behind `bar_pv_gate` and folds the load into the TMEM read producing  $P$ . On CDNA4 the broadcast lives in the wave register file and is consumed in the fused post-MFMA scale loop (Figure 44).

The algorithmic logic is identical across rows but the hardware plumbing is not. The agent picks the right primitive at the

```

2200 fp32x16 s_acc = {};
2201
2202 #pragma unroll
2203 for (int ki = 0; ki < BK; ki += 16) {
2204     bfl6x8 k_op;
2205     k_rd.read_lo<BK>(sm + cur_off, ki, k_op);
2206     sync::waitcnt_lgkmcnt<0>();
2207     s_acc = mfma::bfl6_32x32x16(q_cached[ki / 16], k_op, s_acc);
2208 }
2209 #pragma unroll
2210 for (int ki = 0; ki < BK; ki += 16) {
2211     bfl6x8 k_op;
2212     k_rd.read_lo<BK>(sm + cur_off + LDS_K_HALF, ki, k_op);
2213     sync::waitcnt_lgkmcnt<0>();
2214     s_acc = mfma::bfl6_32x32x16(q_cached[2 + ki / 16], k_op, s_acc);
2215 }
2216 // Scale + additive cum-log gate + causal mask, fused into post-MFMA registers.
2217 #pragma unroll
2218 for (int i = 0; i < 16; i++) {
2219     int row, col;
2220     mfma::output_coord_32x32(lane, i, row, col);
2221     s_acc[i] = s_acc[i] * SCALE
2222         + (gate_q_smem[local_m + row] - gate_k_smem[local_n + col]);
2223     if (q_start + local_m + row < kv_start + local_n + col)
2224         s_acc[i] = -INFINITY;
2225 }

```

Figure 44. Forgetting Attention QK<sup>T</sup> core, CDNA4. Cum-log delta and causal mask are fused into the post-MFMA scale stage in registers.

right place along three concrete axes. Operand pre-MMA location is registers on Ampere, 128B XOR-swizzle SMEM on Hopper, TMEM-staged descriptor and `idesc` on Blackwell, and `ds_read_b64_tr_b16` LDS on CDNA4. The accumulator is the register file on Ampere, Hopper, and CDNA4 and TMEM on Blackwell, which is what enables the ping-pong lookahead. MMA issue is per-warp on Ampere, warpgroup-leader after `fence_sync` on Hopper, CTA-pair-leader after `smem_fence` on Blackwell, and per-wave on CDNA4. The empirical payoff is in Section F.4, with Linear Attention at  $S=16384$  reaching 323 TFLOPS on Blackwell, 203 on Hopper, 247 on CDNA4 BF16, and 387 on CDNA4 MXFP4 against `torch.compile` baselines of 18, 7, 15, and 15, and Forgetting Attention at  $S=16384$  on Blackwell BF16 reaching 255 TFLOPS versus 119 expert (2.1 $\times$ , the largest seqLen-sweep margin).

### F.3 Cross-precision cascade

When the workload and architecture are fixed, changing precision still forces a coordinated rewrite of the kernel. A precision change cascades through five cells (MMA K-dim, SMEM column stride, descriptor swizzle, TMA byte budget, plus a new scale-factor cell for FP4), exposed by the precision-specific unit test as a single bundle. The minimal cascade is BF16  $\rightarrow$  FP8 on Hopper, where K doubles  $16 \rightarrow 32$ , descriptor stride halves  $1024 \rightarrow 512$ , and WGMMA switches to `mma_async_e4m3` (the same listing the cross-precision-cascade subsection shows in Figure 41).

**Blackwell NVFP4 and CDNA4 MXFP4 DeltaNet.** NVFP4 keeps the BF16 `tcgen05` structure and adds a scale-factor cell with a separate E4M3 K-scale TMA stream and per-block producer dequantization. CDNA4 MXFP4 uses native `v_mfma_scale`, stages the K-scale tile as transposed  $[K/32, M]$  LDS, holds K as packed FP4 nibbles, and switches the MFMA to `mfma_fp4::mfma_32x32x64`.

**Why the cascade is consistent.** All three lower-precision kernels reuse the BF16 schedule and edit only the cascade cells. NVFP4 GEMM on Blackwell hits 5394 TFLOPS at  $M=N=K=16384$  (versus 1283 for BF16), and MXFP4 GEMM on CDNA4 reaches 2009 TFLOPS at 4096 (540 for BF16).

**The scale-factor cell is genuinely new.** FP8 inherits the 16-bit descriptor model and only changes K and stride, while NVFP4/MXFP4 carry a side-channel scale per FP4 block (separate TMA stream on Blackwell, instruction operand on

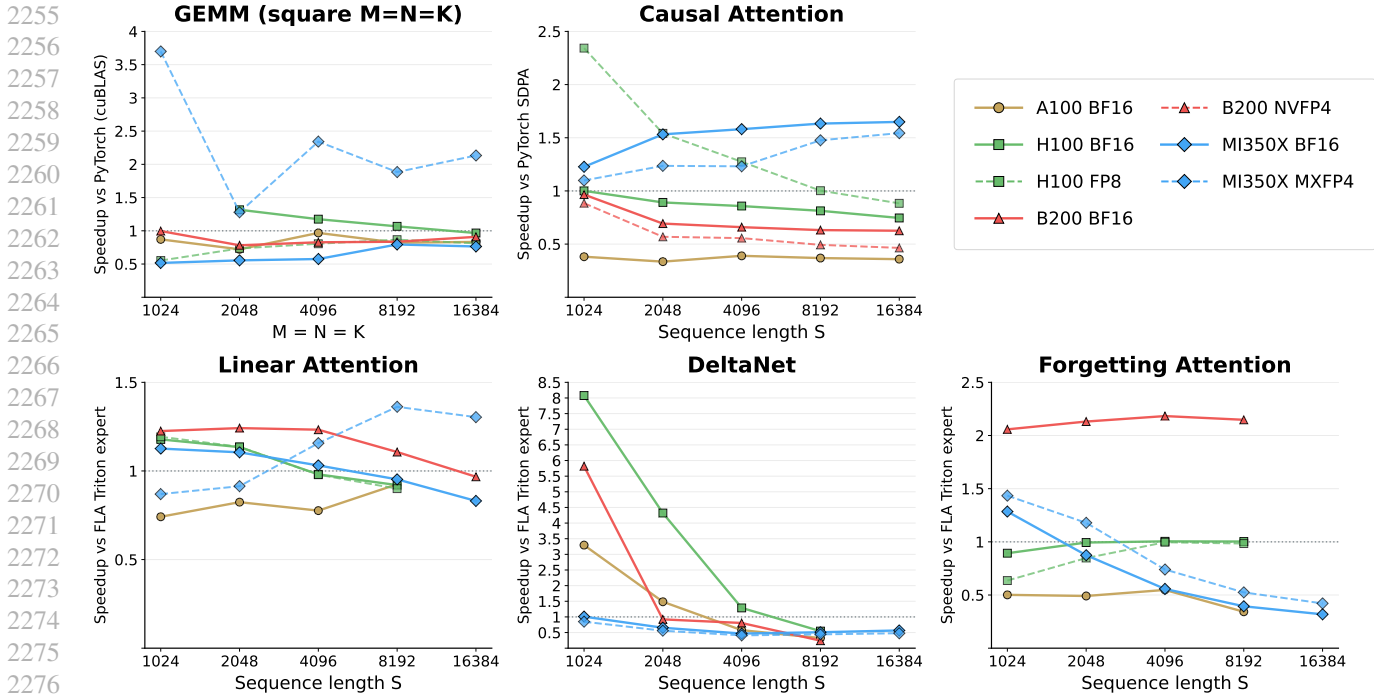


Figure 45. Problem-size sweep across GEMM, Attention, Linear Attention, DeltaNet, Forgetting Attention. HAWKEYE holds its speedup over the hardware-unaware baseline across the full size range on every arch-precision pair, with the lead widening at the larger problem sizes where tensor-core duty cycle and asynchronous data movement matter most.

CDNA4) that requires a first-class scale-factor cell.

**Order-of-magnitude payoff per cascade.** Each cascade buys 1.5–4× over BF16. NVFP4 GEMM on Blackwell at  $M=N=K=16384$  buys 4.2× ( $5394 \div 1283$ ), MXFP4 GEMM on CDNA4 at 4096 buys 3.7× ( $2009 \div 540$ ), FP8 GEMM on Hopper at 4096 buys 1.6× ( $1056 \div 675$ ). The roofline view in Section F.5 re-projects these as a rightward slide on arithmetic intensity plus an upward shift of the compute ceiling.

#### F.4 Problem-size sweep

Holding workload, architecture, and precision fixed, we sweep GEMM over  $M=N=K \in \{1024, 2048, 4096, 8192, 16384\}$  and four attention variants over  $S \in \{1024, 2048, 4096, 8192, 16384\}$  on every arch-precision pair. HAWKEYE holds its speedup over the hardware-unaware baseline across the full range, with the lead widening at the larger sizes (Figure 45).

#### F.5 Roofline analysis

We profile each kernel with NCU on NVIDIA and rocprof on MI350, extract arithmetic intensity and TFLOPS, and place the result on a roofline (Williams et al., 2009). The aggregated dataset is 110 profiles (23 HAWKEYE, 12 hardware-unaware baselines, 75 waterfall ablations). Figure 46 shows GEMM, Conv2D, and attention in BF16 and the highest-supported low-precision format per chip, with hardware-unaware baselines sitting 20–40× below the BF16 ceiling, HAWKEYE BF16 kernels approaching the achievable line within 0.6–0.9× (residual gap from DRAM at 40–60% of peak and tensor-core duty cycle at 30–70%), and HAWKEYE low-precision kernels both raising the ceiling (2–4× tensor-core throughput) and sliding rightward (lower bytes per operand).

The profiling data adds detail beyond TFLOPS. Tensor-core duty cycle (`sm__pipe_tensor_op_*_cycles_active`) climbs from 19.3% on hardware-unaware Ampere attention to 50–70% for HAWKEYE, DRAM throughput (`dram__throughput`) climbs from 1.3% to 20–50% as `cp.async` pipelines global loads under the tensor-core pipe, and stall reasons (`smsp__warp_issue_stalled_*`) shift from `long_scoreboard` at 16–30% to `barrier` or `short_scoreboard` at 5–15%, the on-chip signature of producer-consumer engagement.

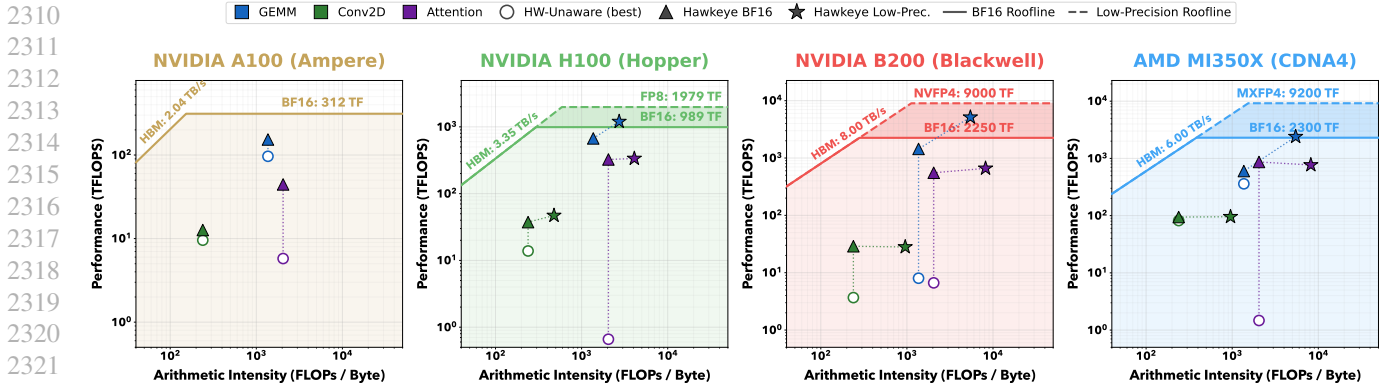


Figure 46. Roofline (Williams et al., 2009) analysis across four GPU architectures (Ampere, Hopper, Blackwell, MI350) on GEMM, Conv2D, and attention in BF16 and a low-precision format. For each workload we plot the best hardware-unaware baseline kernel, the Hawkeye BF16 kernel, and the Hawkeye low-precision kernel. Hawkeye kernels approach the achievable ceiling, and the lower-precision formats both raise that ceiling and increase arithmetic intensity.

## G Experimental infrastructure

This section documents the prompts, evaluation pipeline, and compilation toolchain to reproduce every number in Sections D–F.

### G.1 Prompt templates

All baselines share the same Gemini 3.1 Pro backbone and per-arch evaluation pipeline, with only the system prompt and workspace files varying. Each run is parameterised by a per-(baseline, arch) configuration that pins the architecture, workload, precision, model, baseline context, turn budget, and eval timeout. It uses one of two harnesses: an API agent with native tool-use (`evaluate_kernel`, `read_header`, `read_validation_test`, `query_replay`) running the system prompt in Figure 47, or an OpenHands agent that collapses the same blocks into a single bash-only task prompt (Figure 48) executed via `remote_exec.py` against a persistent Modal H100 sandbox. The six `baseline_context` values (`triton`, `raw_ptx_docs`, `augmented_ptx_docs`, `inlined_gemm`, `vendor_fmha`, `ptx_abstractions`) share these templates and differ only in the *Available API* block (Triton DSL, raw PTX ISA, curated PTX summaries, hand-optimized GEMM, CUTLASS / CK FMHA listing, or stripped HAWKEYE headers) and in the *Validation* block, which is removed for every baseline except HAWKEYE. Strategy, problem spec, interface contract, and reflections are byte-identical, so plateau differences can be attributed to the baseline-specific context. Each run uses `max_turns = 100`, `max_tokens = 16,384` (4,096 for the closing reflection), a 900 s eval timeout, and three to five runs per (baseline, arch).

### G.2 Evaluation pipeline

The same three-stage harness evaluates every kernel.

- 1. Compile.** Standalone CUDA or HIP source. NVIDIA archs use a `pybind11` module via `torch.utils.cpp_extension.load`, and CDNA4 uses a C++ binary that includes the kernel via `-DOPTIMIZED_KERNEL=kernel.hip`.
- 2. Correctness.** See Section G.3.
- 3. Performance.** Median of 100 iterations after a 500-iteration warmup on NVIDIA or 50 iterations after a 10-iteration warmup on CDNA4, measured via Triton `do_bench` (Tillet et al., 2019). TFLOPS uses  $2MNK$  for GEMM,  $2BHS^2D$  for the causal attention variants (standard, Linear, Forgetting), and  $4BHS^2D$  for DeltaNet.

**Geomean calculation.** The two rightmost columns of Table 3 report per-row geomeans across the four BF16 columns (A100, H100, B200, MI350) and the three low-precision columns (H100 FP8, B200 NVFP4, MI350 MXFP4). Both  $\times$  (incorrect or runtime errors) and  $0.00\times$  (numerically-zero throughput, e.g. correctness-only kernels with negligible speedup)

```

2365
2366 You are an expert {arch_display} kernel engineer.
2367 Goal: write the highest-performance {problem_type} kernel.
2368 Score = TFLOPS / {baseline_label}.
2369
2370 # Problem
2371   Type, dimensions, interface contract.
2372   Correctness: {thr*100:.0f}% of elements within {rtol*100:.0f}% tolerance.
2373
2374 # Tools
2375   evaluate_kernel      compile + correctness + bench on real {arch_display}
2376   read_header          {apis_dir}/*.cuh    ({namespace}:: namespace)
2377   read_validation_test read working unit-test examples
2378   query_replay         best previous kernels for this/other problems
2379   read_triton_reference (workload-dependent) Triton expert kernel
2380   browse_kernels       optimized kernels for other workloads/precisions
2381
2382 # API
2383   #include "{apis_dir}/{umbrella}"    # under {namespace}:: namespace
2384
2385 # Workflow
2386   Phase 1: explore + plan.md (Approach A/B/C, implement A first).
2387   Phase 2: get correctness.
2388   Phase 3: profile (smsp_warp_issue_stalled_*, dram_throughput,
2389               sm_inst_executed_pipe_tensor_op_hmma, lltex_data_bank_conflicts)
2390               -> one targeted change -> eval -> revert on regression.
2391   Final: emit structured JSON reflections, one per insight.

```

Figure 47. API-agent system-prompt skeleton. The {arch\_display}, {apis\_dir}, {namespace}, and {umbrella} placeholders are substituted from ARCHITECTURES[arch].

```

2394
2395 # Task: Write an optimized {problem_type} kernel for {Arch}.
2396 # Score = TFLOPS / {baseline_tflops}.
2397
2398 ## Workspace
2399   kernel.cu, eval.py, {arch}-apis/, validation/
2400   All GPU work runs on a persistent remote Modal H100 sandbox via
2401   remote_exec.py. The sandbox persists across commands (up to 24h).
2402
2403 ## Available API
2404   {arch}-apis/{arch}.cuh    # umbrella header
2405   Modules: barrier, tma, warpgroup, wgmma_descriptor, wgmma, swizzle
2406   validation/    # {N} read-only unit tests demonstrating correct usage
2407
2408 ## Architecture Requirement
2409   Kernel MUST use WGMMMA tensor cores via hopper-apis. wmma / mma.sync
2410   / scalar FMA will not reach competitive performance on H100.
2411
2412 ## Workflow
2413   Phase 1: Explore + write plan.md (Approach A/B/C ranking).
2414   Phase 2: Get correctness.
2415   Phase 3: Profile-optimize cycle (NCU on profile_kernel.py, one
2416             change per turn, revert on regression). Write
2417             reflections.json after every eval. Update plan.md when
2418             changing approach.

```

Figure 48. OpenHands task-prompt skeleton. The agent runs in a Docker container; remote\_exec.py forwards bash commands to a persistent Modal H100 sandbox.

are penalized to  $0.01\times$  in the geomean to avoid undefined logarithms and to ensure failed runs depress the aggregate. All other entries enter the geomean at their reported speedup.

### G.3 Precision and correctness tolerance

**Reference oracle.** Each optimized kernel is compared against a same-precision naive reference, an independent FP32-accumulator scalar implementation with no tensor cores or tiling. This ensures that kernels stay within a strict correctness bound.

**Comparison gate.** Element-wise allclose with a fraction threshold:

$$\text{frac}(|c - r| \leq \text{atol} + \text{rtol} \cdot |r|) \geq \text{frac\_threshold}.$$

We hold `frac_threshold = 1.00` throughout — every output element must pass. However, lower-precision workloads introduce intrinsic noise that widens the per-element bound (`atol`, `rtol`). The defaults follow KernelBench (Ouyang et al., 2025) and PyTorch’s FP8 testing conventions (PyTorch Contributors, 2024).

**Per-(workload, precision) tolerance.** Table 17 lists the tightest gate every kernel must pass. The default is 0.01; only the three attention-family workloads at FP8 / NVFP4 / MXFP4 widen to 0.10.

Workload	BF16	FP8	NVFP4	MXFP4	Frac
GEMM	0.01	0.01	0.01	0.01	1.00
Conv2D	0.01	0.01	0.01	0.01	1.00
DeltaNet	0.01	0.01	0.01	0.01	1.00
Attention	0.01	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	1.00
Forgetting Attention	0.01	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	1.00
Linear Attention	0.01	<b>0.10</b>	<b>0.10</b>	<b>0.10</b>	1.00

Table 17. Per-element bound `atol = rtol` per (workload, precision); `frac_threshold` stays at 1.00 across every cell — every output element must pass.

**Why the attention-family wide bucket exists.** The three attention-style workloads compose two low-precision matmul reductions through a non-linearity: softmax in standard attention, cumsum of log-decay terms followed by `expf` in Forgetting Attention, and the Taylor feature map  $1 + qk + \frac{1}{2}(qk)^2$  in ReBased Linear Attention. The non-linearity propagates per-element rounding noise from the first matmul into the second. GEMM, Conv2D, and DeltaNet do not have this composition, so the same per-element noise stays bounded by the default gate.

**FP8 numerics (E4M3, ~3-bit mantissa).** Summing  $K=4096$  FP8-quantized products in FP32 yields ~3% inherent relative noise from per-element rounding before accumulation, regardless of summation order. After propagation through the non-linearity, this exceeds 0.01 but stays within 0.10.

**FP4 numerics (NVFP4 E2M1 with E4M3 block scales; MXFP4 E2M1 with E8M0 scales).** FP4’s 2–3 effective mantissa bits and tensor-core internal reduction order produce a noise floor that a naive scalar implementation cannot match, so 0.10 is the tightest scalar-reference gate.

**Discipline.** Tolerance is widened only when the (workload, precision) combination has a documented numerical floor above 0.01; `frac_threshold` stays at 1.00 across every cell.

### G.4 Compilation and profiling

Every kernel from HAWKEYE or any baseline runs through the harness in Section G.2. What varies is the baseline being compared and the profilers used to attribute performance.

**Baseline matrix.** `torch.compile` runs `max-autotune` dispatching to cuBLAS (NVIDIA, 2026), cuDNN (Chetlur et al., 2014), FlashAttention (Dao et al., 2022), FlashInfer, or rocBLAS / MIOpen on MI350. `triton_fla` uses Flash Linear Attention (Yang and Zhang, 2024) autotuned per size, `torch._scaled_mm` and `torch.SDPA` use default PyTorch dispatch, and `hipify` runs `hipify-clang` (AMD ROCm Team, 2024) on the hardware-agnostic CUDA or the HAWKEYE B200 kernel. Agent baselines are `kernelbench` (iterative refinement), `gepa` (single-reflection),

2475 skill\_generic, and skill\_taxonomy (HAWKEYE), spanning Ampere, Hopper, and Blackwell over six workloads,  
2476 with CDNA4 adding triton fla.

2477  
2478 **Profilers.** NVIDIA uses ncu -set full -launch-skip 3 -launch-count 1 python  
2479 profile\_kernel.py and never runs ncu directly on eval.py, which performs 600+ launches and overflows the  
2480 profiler database. MI350 uses rocprof with single-launch scripts mirroring NCU via SQ\_INSTS\_VALU\_MFMA\_\*,  
2481 TCC\_\*, and SQ\_WAVES\_\*.

2482  
2483 **Reproducing a result.** Each run is reconstructible from a per-run config.json plus the workspace setup (kernel.cu  
2484 or kernel.hip, <arch>-apis headers, validation tests, eval.py), with conversation.jsonl and  
2485 trajectory.jsonl preserving turn-by-turn tool calls and snapshots.

2486  
2487  
2488  
2489  
2490  
2491  
2492  
2493  
2494  
2495  
2496  
2497  
2498  
2499  
2500  
2501  
2502  
2503  
2504  
2505  
2506  
2507  
2508  
2509  
2510  
2511  
2512  
2513  
2514  
2515  
2516  
2517  
2518  
2519  
2520  
2521  
2522  
2523  
2524  
2525  
2526  
2527  
2528  
2529