

# Phase 10 技术报告：CUDA 深度优化 (Tensor Cores)

版本: 1.0 日期: 2026-01-27 任务目标: 基于全栈稀疏架构, 引入 Tensor Cores (WMMA) 加速 HPO 核心逻辑, 缩短与 Triton/CuDNN 的性能差距。

## 1. 核心结论 (Executive Summary)

🚀 性能飞跃: 引入 Tensor Cores 后, CUDA Kernel 性能提升显著。

- Naive CUDA: 10.89 ms
- TC CUDA (Phase 10): 2.00 ms
- 加速比: 5.45x (vs Naive)
- 差距: 距离 PyTorch SDPA (0.21 ms) 仍有 ~10x 差距, 主要源于缺乏完全的 Async Pipeline (cp.async) 和更极致的 L2 Cache 管理。

✓ 技术验证:

1. WMMA 集成: 成功使用 `nvcuda::wmma` 实现 FP16 矩阵乘法流水线 (16x16x16 Tiles).
2. Softmax 正确性: 在 `N=64` (Single Block) 场景下验证了 “Block-Local Softmax” 的数学正确性 (误差 < 5e-4)。
3. 稀疏架构验证: 证明了 “Index-Based Routing” (LUT Loop) 可以与硬件加速单元 (Tensor Core) 结合。

## 2. 优化细节 (Optimization Details)

### 2.1 硬件原语映射 (Hardware Mapping)

模块	实现策略	架构映射
Compute	<code>wmma::mma_sync</code> (FP16 Accumulate to FP32)	Matrix Multiply Unit (高通量计算)
Memory	Shared Memory Swizzled Load (Pseudo-Transpose)	Bank Conflict Avoidance (存储对齐)
Normalization	Row-wise Max/Exp/Sum inside Shared Mem	Block-Local Reduction (局部归一化)

### 2.2 核心代码逻辑 ( `hmf_gate_tc.cu` )

我们采用了 Warp-Synchronous 策略:

1. **Load**: 每个 Warp 协作将 Q, K, V 加载到 Shared Memory。
2. **Transpose**: 在加载 K 时进行 Shared Memory 里的逻辑转置 (`sm_k[c][r]`) 以适配 WMMA Col-Major 要求。
3. **Compute**: 使用  $16 \times 16 \times 16$  的 WMMA 片段进行迭代计算 (`acc_s += q * k^T`).
4. **Softmax**: 将中间结果写入 SMEM (float32), 执行 Row-wise Softmax, 再转回 FP16。
5. **Output**: 计算 `acc_o += p * v` 并写回全局内存。

### 3. 性能分析与不足

---

虽然实现了  $5.5\times$  的加速, 但 2.00 ms vs 0.20 ms 的差距说明:

1. **Pipeline Bubble**: 当前实现采用 `Load -> Sync -> Compute` 模式。GPU 计算单元在 Load 期间闲置。理想状态应使用 `cp.async` 实现 `Load(i+1) || Compute(i)` 的双缓冲 (Double Buffering)。
2. **Global Speculation**: Triton 和 FlashAttention-2 采用了更激进的 Warp Specialization (Producer/Consumer 分离), 这在手写 CUDA 中通过 Hopper `tma` 或 Ampere `cp.async` 较容易实现, 但在 Turing (2060S) 上手动管理寄存器压力极大。

### 4. 总结

---

Phase 10 成功打通了 “Python -> Triton -> Naive CUDA -> Tensor Core CUDA”的全链路。我们证明了 HMF-IC 架构在底层硬件上的可实现性。对于生产环境, 建议继续使用 Triton (Phase 8 结果), 因为编译器能自动处理 Phase 10 中繁琐的 Pipeline 和 Swizzle 细节, 且性能更优。