

## Derivation of CUDA Kernel Execution-Time Prediction

This document explains the full reasoning process used to derive the CUDA kernel performance prediction template. The model follows the principles of GPU performance engineering, including FLOP analysis, memory bandwidth limits, and kernel launch latency behavior observed on NVIDIA GPUs.

### 1. Overview of Performance Model

Every CUDA kernel runtime can be approximated using three components:

- Work performed by threads (iterations  $\times$  FLOPs)
- Memory traffic (bytes loaded/stored)
- GPU hardware limits (peak FLOPs, memory bandwidth, launch latency)

This mirrors the NVIDIA Roofline Model and CUDA optimization best practices.

### 2. Workload Derivation from the Kernel

Given a loop of the form:

for (offset = threadIdx.x; offset < K; offset += blockDim.x) { float t = a - b; temp += t \* t; } We derive:

- $\text{iters\_per\_thread} = \text{ceil}(K / \text{blockDim.x})$
- Each iteration performs ~3 FLOPs (subtract, multiply, add)
- Each iteration loads 8 bytes (two float loads)

Threads are organized as:

$\text{threads\_per\_block} = B_x \times B_y \times B_z$

$\text{active\_blocks} = G_x \times G_y \times G_z$  (adjusted by boundary checks)

Total iterations:

$$N_{\text{iters\_total}} = \text{iters\_per\_thread} \times \text{threads\_per\_block} \times \text{active\_blocks}$$

### 3. FLOP Analysis

Total FLOPs =  $N_{\text{iters\_total}} \times \text{FLOPs\_per\_iteration}$

This follows GPU throughput-based modeling. FLOP counts determine the compute-bound time when compared to GPU peak performance.

**Compute-bound time:**

$$t_{\text{compute}} = \text{FLOPs\_total} / \text{Peak\_FP32}$$

### 4. Memory-Traffic Analysis

Bytes per iteration =  $(\# \text{float\_loads} \times 4) + (\# \text{float\_stores} \times 4)$

$$\text{Total memory bytes} = N_{\text{iters\_total}} \times \text{bytes\_per\_iter}$$

**Memory-bound time:**

$$t_{\text{mem}} = \text{bytes\_total} / \text{memory\_bandwidth}$$

Whichever is larger ( $t_{\text{compute}}$  or  $t_{\text{mem}}$ ) determines the kernel body time:

$$t_{\text{body}} = \max(t_{\text{compute}}, t_{\text{mem}})$$

### 5. Kernel Launch Overhead

NVIDIA GPUs exhibit a typical kernel launch latency of ~4–6 microseconds (Turing, Ampere, Ada; verified in NVIDIA Developer Forums).

This overhead dominates for small kernels.

**Final predicted runtime:**

$t_{\text{total}} \approx t_{\text{body}} + t_{\text{launch}}$

## **6. Why This Model Is Correct**

- Matches NVIDIA Roofline methodology
- Matches Nsight Compute performance formulas
- Matches CUDA programming and memory-throughput theory
- Matches empirical kernel launch latency observations

## **7. Template Summary**

1. Compute iterations
2. Compute FLOPs
3. Compute bytes
4. Compute compute-bound and memory-bound times
5. Add launch overhead

**This derivation forms the basis of the prediction template and the GPU-specific runtime predictions generated earlier.**