

# CUDA gpuKendall Kernel Execution Time Prediction Report

This report applies the CUDA kernel performance prediction template to the kernel: gpuKendall, launched with Grid = (5,5) and Block = (32,32).

We estimate the expected per-launch execution time on five NVIDIA GPUs:

- GeForce GTX TITAN Black
- GeForce GTX TITAN X
- NVIDIA TITAN V
- GeForce RTX 2080 Ti
- GeForce RTX 4070

```
__global__ void gpuKendall(const float * a, size_t na, const float * b, size_t nb,
                           size_t sampleSize, double * results)
{
    size_t i, j, tests,
          tx = threadIdx.x, ty = threadIdx.y,
          bx = blockIdx.x, by = blockIdx.y,
          rowa = bx * sampleSize, rowb = by * sampleSize;
    float discordant, concordant = 0.f, numer, denom;
    __shared__ float threadSums[NUMTHREADS*NUMTHREADS];

    for(i = tx; i < sampleSize; i += NUMTHREADS) {
        for(j = i+1+ty; j < sampleSize; j += NUMTHREADS) {
            tests = ((a[rowa+j] > a[rowa+i]) && (b[rowb+j] > b[rowb+i]))
                  + ((a[rowa+j] < a[rowa+i]) && (b[rowb+j] < b[rowb+i]))
                  + ((a[rowa+j] == a[rowa+i]) && (b[rowb+j] == b[rowb+i]));
            concordant = concordant + (float)tests;
        }
    }
    threadSums[tx*NUMTHREADS+ty] = concordant;
    __syncthreads();
    // 2D reduction in shared memory...
}
```

## 1. Workload Analysis

Launch configuration:

- Grid = (5,5) → 25 blocks
- Block = (32,32) → 1024 threads/block
- NUMTHREADS = 16 (used for loop strides and shared memory)
- na = nb = 5, sampleSize = 100

Logically, the computation is a pairwise comparison over all distinct index pairs (i,j) with  $0 \leq i < j < \text{sampleSize}$  for each block (bx,by).

Total pairs per block:

$$\begin{aligned} C(\text{sampleSize}, 2) &= \text{sampleSize} \times (\text{sampleSize} - 1) / 2 \\ &= 100 \times 99 / 2 = 4,950 \text{ pairs} \end{aligned}$$

With 25 blocks (all used, since grid = (5,5) and na = nb = 5):

Total pair evaluations across the grid:

$$N_{\text{pairs\_total}} = 4,950 \times 25 = 123,750$$

## 2. FLOP Count (Approximate)

For each pair, the code evaluates three conjunctions of float comparisons and accumulates an integer count into a float:

```
tests = (condition1) + (condition2) + (condition3);
concordant += (float) tests;
```

We approximate this as  $\approx 10$  FLOP-equivalents per pair (comparisons, boolean logic, and the final float add).

Total FLOPs:

$F_{\text{total}} \approx 123,750 \times 10 = 1,237,500 \approx 1.24 \times 10^6$  FLOPs.

## 3. Memory Traffic

For each pair  $(i,j)$  and a given block  $(bx,by)$ , the kernel reads:

- $a[\text{rowa} + i], a[\text{rowa} + j]$
- $b[\text{rowb} + i], b[\text{rowb} + j]$

Assuming these four float values are loaded once each and reused in the three conditions, we get 4 float loads per pair  $\rightarrow 16$  bytes per pair.

Total bytes per block:

$B_{\text{block}} = 4,950 \text{ pairs} \times 16 \text{ bytes} \approx 79,200 \text{ bytes}$

Across 25 blocks:

$B_{\text{total}} \approx 79,200 \times 25 = 1,980,000 \text{ bytes} \approx 1.98 \times 10^6 \text{ bytes} (\sim 1.89 \text{ MiB})$ .

The final write of one double per block ( $\text{results}[by*na + bx]$ ) adds only  $25 \times 8 = 200$  bytes and is negligible in comparison.

Thus, the kernel is predominantly memory- and launch-limited.

## 4. GPU Specifications Used

Approximate FP32 peak performance and memory bandwidth:

GPU	Peak FP32 (FLOPs/s)	Bandwidth (bytes/s)
GTX TITAN Black	5.12e12	3.36e11
GTX TITAN X	6.14e12	3.365e11
TITAN V	1.49e13	6.528e11
RTX 2080 Ti	1.345e13	6.16e11
RTX 4070	2.9e13	5.04e11

## 5. Time Estimates

Using:

- $F_{\text{total}} \approx 1.2375 \times 10^6$  FLOPs

- $B_{\text{total}} \approx 1.98 \times 10^6$  bytes

We compute:

$$t_{\text{compute}} = F_{\text{total}} / \text{Peak\_FP32}$$

$$t_{\text{mem}} = B_{\text{total}} / \text{Bandwidth}$$

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

$t_{\text{total}} \approx t_{\text{body}} + 5 \mu\text{s}$  (for kernel launch overhead).

GPU	<b>t_compute (μs)</b>	<b>t_mem (μs)</b>	<b>t_body (μs)</b>	<b>t_total (μs)</b>
GTX TITAN Black	0.24	5.89	5.89	$\approx 10.89$
GTX TITAN X	0.20	5.88	5.88	$\approx 10.88$
TITAN V	0.08	3.03	3.03	$\approx 8.03$
RTX 2080 Ti	0.09	3.21	3.21	$\approx 8.21$
RTX 4070	0.04	3.93	3.93	$\approx 8.93$

## 6. Conclusion

The gpuKendall kernel performs roughly 1.24M FLOPs and reads about 1.98MB of data per launch, spread across 25 blocks and 16x16 logical threads per block.

Because the memory and compute loads are modest, the kernel body completes in a few microseconds, and total runtime per launch is dominated by the fixed CUDA kernel launch overhead ( $\sim 5 \mu\text{s}$ ).

Predicted per-launch times:

- GTX TITAN Black / TITAN X:  $\approx 10.9 \mu\text{s}$
- TITAN V:  $\approx 8.0 \mu\text{s}$
- RTX 2080 Ti:  $\approx 8.2 \mu\text{s}$
- RTX 4070:  $\approx 8.9 \mu\text{s}$

These values are approximate but consistent with typical behavior of small, pairwise-comparison kernels with limited memory traffic.