

# CUDA getUnrestricted Kernel Execution Time Prediction Report

This report applies the CUDA kernel performance prediction template to the kernel: getUnrestricted, 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

## 1. Effective Parallelism

Launch configuration:

- Grid = (5,5) → 25 blocks
- Block = (32,32) → 1024 threads/block
- THREADSPERDIM = 16
- countx = 5, county = 5
- rows = 100, cols = 10

The kernel uses:

$n = blockIdx.x * THREADSPERDIM + threadIdx.x;$   
 $m = blockIdx.y * THREADSPERDIM + threadIdx.y;$

With  $\text{gridDim.x} = \text{gridDim.y} = 5$  and  $\text{THREADSPERDIM} = 16$ ,  $n$  and  $m$  each range from 0 to 95, but only indices with  $n < \text{countx}$  (5) and  $m < \text{county}$  (5) do work.

This occurs only in block ( $\text{blockIdx.x} = 0$ ,  $\text{blockIdx.y} = 0$ ) and  $\text{threadIdx.x}, \text{threadIdx.y}$  in 0..4.

Thus there are  $5 \times 5 = 25$  active (m,n) pairs per launch; all other threads return immediately. Each active thread performs a full Gram–Schmidt QR solve on a  $100 \times 10$  system (similar to getRestricted).

## 2. FLOP Count

As in getRestricted, for each (m,n) pair we approximate:

- Gram–Schmidt orthogonalization:  $\approx 21,100$  FLOPs
- $R = Q \square X$ ,  $B = Q \square Y$ , and back-substitution:  $\approx 22,100$  FLOPs

Total per-thread (per (m,n)) FLOPs:

$F_{\text{thread}} \approx 43,200$  FLOPs

With 25 active threads:

$F_{\text{total}} = 43,200 \times 25 = 1,080,000 \approx 1.08 \times 10^7$  FLOPs.

## 3. Memory Traffic (Approximate)

Per active thread we reuse the getRestricted estimate:  
 $B_{\text{thread}} \approx 1.99 \times 10^{\text{■}} \text{ bytes}$  (loads/stores of X, Q, R, Y, B).

With 25 active threads:  
 $B_{\text{total}} \approx 1.99 \times 10^{\text{■}} \times 25 = 4.975 \times 10^{\text{■}} \text{ bytes} \approx 4.98 \text{ MB}$  per kernel launch.

This again makes the kernel primarily memory bound.

#### 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:

- $F_{\text{total}} \approx 1.08 \times 10^{\text{■}} \text{ FLOPs}$
- $B_{\text{total}} \approx 4.975 \times 10^{\text{■}} \text{ bytes}$

We compute for each GPU:

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

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

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

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

GPU	$t_{\text{compute}} (\mu\text{s})$	$t_{\text{mem}} (\mu\text{s})$	$t_{\text{body}} (\mu\text{s})$	$t_{\text{total}} (\mu\text{s})$
GTX TITAN Black	0.211	14.81	14.81	$\approx 19.81$
GTX TITAN X	0.176	14.78	14.78	$\approx 19.78$
TITAN V	0.072	7.62	7.62	$\approx 12.62$
RTX 2080 Ti	0.080	8.08	8.08	$\approx 13.08$
RTX 4070	0.037	9.87	9.87	$\approx 14.87$

#### 6. Conclusion

getUnrestricted performs roughly 1.08M FLOPs and moves ~5MB of data per launch, spread over only 25 active (m,n) threads while many threads return immediately.

Given the relatively modest FLOP count and multi-megabyte global memory traffic, the kernel is clearly memory bound. The body time is on the order of 8–15  $\mu\text{s}$ , and adding a ~5  $\mu\text{s}$  CUDA kernel launch overhead yields total times:

- GTX TITAN Black / TITAN X:  $\approx 19.8 \mu\text{s}$
- TITAN V:  $\approx 12.6 \mu\text{s}$
- RTX 2080 Ti:  $\approx 13.1 \mu\text{s}$
- RTX 4070:  $\approx 14.9 \mu\text{s}$

These estimates are approximate but consistent with expectations for a QR-based kernel with limited parallel occupancy and moderate memory bandwidth usage.