# NVIDIA SYSTEM PROFILING

Raghav, Tianhao, Vasu, Ram

# Intro

# What is a System?



Stage 1
Transfer data to GPU and allocate device memory

Stage 2
Launch Kernel

Stage 3
CUDA API calls

Stage 4
Repeat until all kernels have executed

Stage 5
Copy all data back to CPU

## Profiling

A session in which a programmer can collect data about the performance of the application.

## Sampling

Collecting the stacks of threads that are active in order to build an understand of the time spent in different kernel functions.

## Testing

Provides sophisticated information about API functions that have been called and the amount of time spent.

# What is Nsight Systems?

- Overall performance, system level profiling, and analysis
- Shows runtime of Cuda API calls like memcpy and synchronize and kernels
- Helps identify issues across entire systems not just within kernels
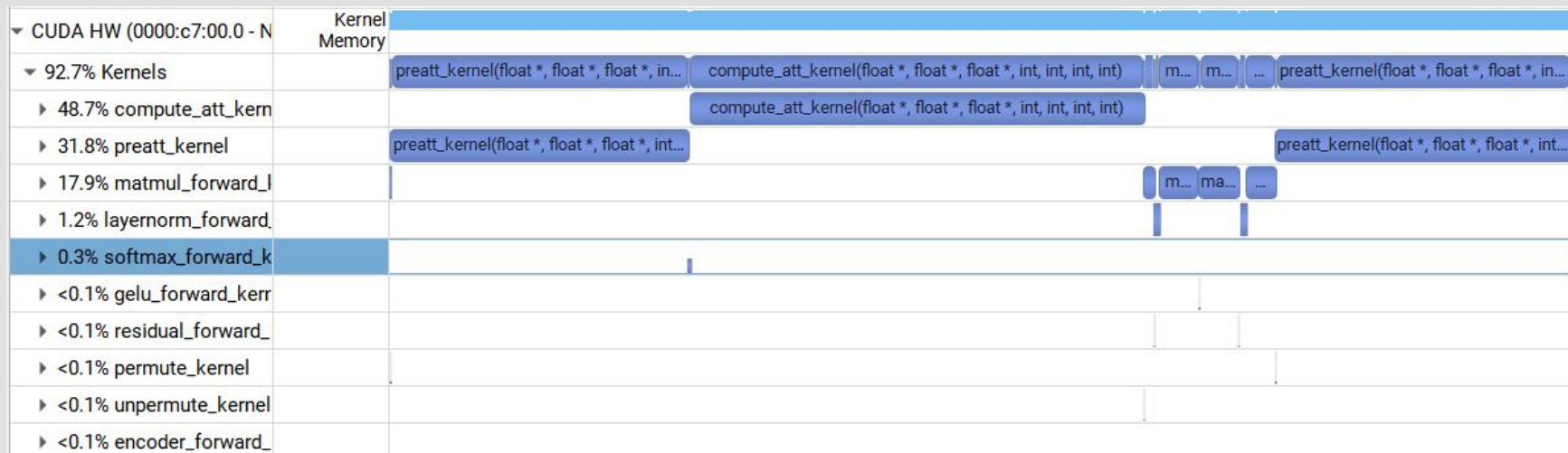- Shows timeline of events

# Nsight Systems

- Does GPU kernel level performance analysis
- Allows for more focused profiling of GPU kernels
- Gives detailed metrics of each kernel
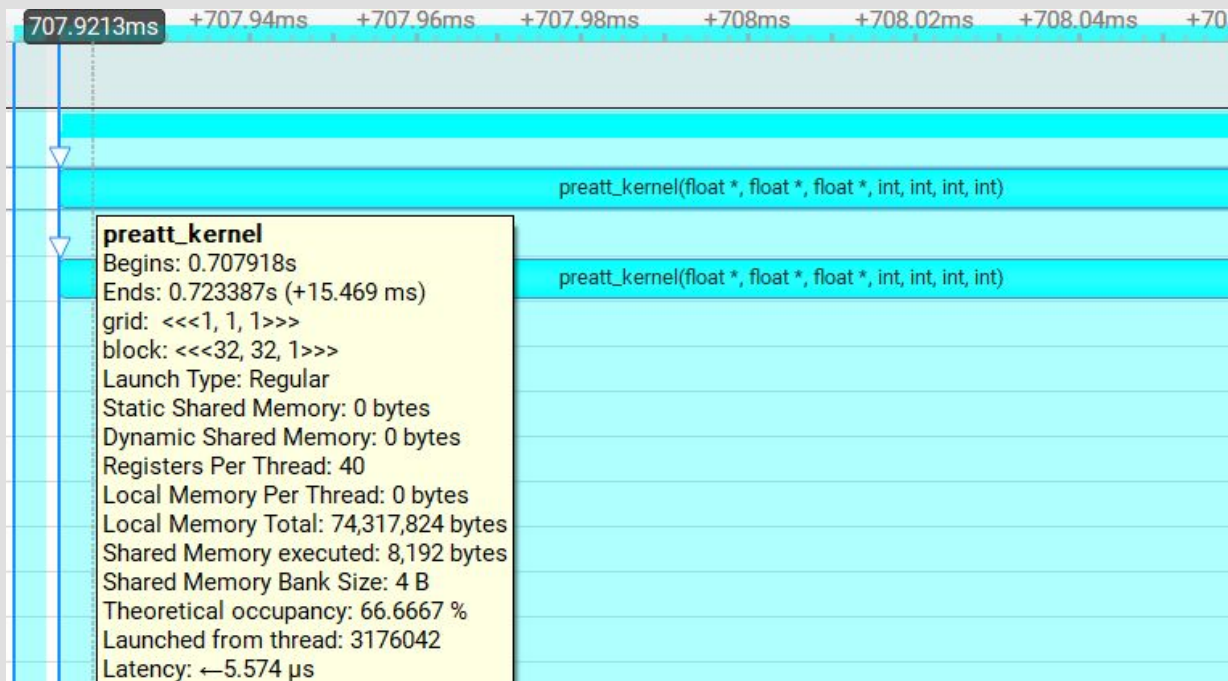- Gives hints of solving bottlenecks

# Nsight Compute

# Application Overview

# Examples from GPT-2

# Before M3 optimizations

| Time | Total Time | Instances | Avg | Med | Min | Max | StdDev | Category | Operation |
|---|---|---|---|---|---|---|---|---|---|
| 39.4% | 492.071 ms | 371 | 1.326 ms | 26.600 µs | 972 ns | 23.795 ms | 4.933 ms | CUDA_API | cudaDeviceSynchronize |
| 22.8% | 284.942 ms | 12 | 23.745 ms | 23.741 ms | 23.711 ms | 23.789 ms | 21.810 µs | CUDA_KERNEL | compute_att_kernel(float *, float *, float *, int, int, int, int) |
| 14.9% | 185.817 ms | 12 | 15.485 ms | 15.483 ms | 15.468 ms | 15.510 ms | 13.325 µs | CUDA_KERNEL | preatt_kernel(float *, float *, float *, int, int, int, int) |
| 12.6% | 157.374 ms | 13 | 12.106 ms | 10.950 µs | 391 ns | 146.605 ms | 40.477 ms | CUDA_API | cudaFree |
| 3.5% | 44.320 ms | 3 | 14.773 ms | 8.040 ms | 21.199 µs | 36.258 ms | 19.034 ms | CUDA_API | cudaMemcpy |
| 2.9% | 36.165 ms | 2 | 18.082 ms | 18.082 ms | 1.792 µs | 36.163 ms | 25.570 ms | MEMORY_OPER | [CUDA memcpy Host-to-Device] |
| 1.1% | 14.112 ms | 153 | 92.238 µs | 110.807 µs | 3.757 µs | 281.175 µs | 55.779 µs | CUDA_API | cudaMalloc |
| 0.6% | 7.960 ms | 1 | 7.960 ms | 7.960 ms | 7.960 ms | 7.960 ms | 0 ns | MEMORY_OPER | [CUDA memcpy Device-to-Host] |
| 0.6% | 7.282 ms | 25 | 291.287 µs | 290.754 µs | 287.938 µs | 296.257 µs | 2.444 µs | CUDA_KERNEL | layernorm_forward_kernel(float *, float *, float *, const fl |
| 0.5% | 5.721 ms | 49 | 116.758 µs | 83.040 µs | 22.752 µs | 2.478 ms | 345.510 µs | CUDA_KERNEL | matmul_forward_kernel(float *, const __half *, const __h |
| 0.2% | 2.513 ms | 1 | 2.513 ms | 2.513 ms | 2.513 ms | 2.513 ms | 0 ns | CUDA_API | cudaFreeHost |
| 0.2% | 2.165 ms | 367 | 5.900 µs | 3.858 µs | 3.236 µs | 167.222 µs | 9.133 µs | CUDA_API | cudaLaunchKernel |
| 0.2% | 1.972 ms | 98 | 20.126 µs | 9.632 µs | 5.312 µs | 482.241 µs | 48.424 µs | CUDA_KERNEL | convertFP32to16(const float *, __half *, int, int, int, int) |
| 0.2% | 1.946 ms | 49 | 39.707 µs | 37.184 µs | 9.792 µs | 572.386 µs | 78.506 µs | CUDA_KERNEL | transpose(const float *, float *, int, int) |
| 0.1% | 1.487 ms | 12 | 123.880 µs | 123.872 µs | 123.008 µs | 124.704 µs | 537 ns | CUDA_KERNEL | softmax_forward_kernel(float *, float, const float *, int, ir |
| 0.1% | 1.366 ms | 1 | 1.366 ms | 1.366 ms | 1.366 ms | 1.366 ms | 0 ns | CUDA_API | cudaMallocHost |
| 0.1% | 1.130 ms | 99 | 11.413 µs | 4.448 µs | 2.496 µs | 406.881 µs | 42.428 µs | MEMORY_OPER | [CUDA memset] |

# Overall Kernel Runtimes

707.9213ms    +707.94ms    +707.96ms    +707.98ms    +708ms    +708.02ms    +708.04ms    +70

preatt_kernel(float *, float *, float *, int, int, int, int)

preatt_kernel(float *, float *, float *, int, int, int, int)

**preatt_kernel**
Begins: 0.707918s
Ends: 0.723387s (+15.469 ms)
grid: <<<1, 1, 1>>>
block: <<<32, 32, 1>>>
Launch Type: Regular
Static Shared Memory: 0 bytes
Dynamic Shared Memory: 0 bytes
Registers Per Thread: 40
Local Memory Per Thread: 0 bytes
Local Memory Total: 74,317,824 bytes
Shared Memory executed: 8,192 bytes
Shared Memory Bank Size: 4 B
Theoretical occupancy: 66.6667 %
Launched from thread: 3176042
Latency: ←5.574 µs

# Individual Kernel Breakdown

| Start | Duration | Name | Result | CorrID | Pid | Tid | T-Pri | Thread Name |
|---|---|---|---|---|---|---|---|---|
| 0.226729s | 147.162 ms | cudaFree | 0 | 597 | 500821 | 500821 | 20 | test_gpt2 |
| 0.49506s | 36.545 ms | cudaMemcpy | 0 | 1129 | 500821 | 500821 | 20 | test_gpt2 |
| 1.10111s | 31.373 ms | cudaDeviceSynchronize | 0 | 1476 | 500821 | 500821 | 20 | test_gpt2 |
| 0.933391s | 23.781 ms | cudaDeviceSynchronize | 0 | 1372 | 500821 | 500821 | 20 | test_gpt2 |
| 0.840914s | 23.770 ms | cudaDeviceSynchronize | 0 | 1316 | 500821 | 500821 | 20 | test_gpt2 |
| 0.65599s | 23.768 ms | cudaDeviceSynchronize | 0 | 1204 | 500821 | 500821 | 20 | test_gpt2 |
| 0.97964s | 23.765 ms | cudaDeviceSynchronize | 0 | 1400 | 500821 | 500821 | 20 | test_gpt2 |
| 0.74849s | 23.761 ms | cudaDeviceSynchronize | 0 | 1260 | 500821 | 500821 | 20 | test_gpt2 |
| 0.887186s | 23.757 ms | cudaDeviceSynchronize | 0 | 1344 | 500821 | 500821 | 20 | test_gpt2 |
| 0.702251s | 23.755 ms | cudaDeviceSynchronize | 0 | 1232 | 500821 | 500821 | 20 | test_gpt2 |
| 1.07207s | 23.744 ms | cudaDeviceSynchronize | 0 | 1456 | 500821 | 500821 | 20 | test_gpt2 |
| 0.794731s | 23.739 ms | cudaDeviceSynchronize | 0 | 1288 | 500821 | 500821 | 20 | test_gpt2 |

# CUDA API Call Durations

# memory metrics

Begins: 0.612214s
Ends: 0.648377s (+36.163 ms)
HtoD memcpy 497,903,616 bytes
Source memory kind: Pageable
Destination memory kind: Device
Throughput: 12.8228 GiB/s
Launched from thread: 3176042
Latency: ←123.096 µs
Correlation ID: 1129
Stream: Default stream 7

**NVIDIA A40**

| Chip Name | GA102 |
|---|---|
| SM Count | 84 |
| L2 Cache Size | 6.00 MiB |
| Memory Bandwidth | 648.29 GiB/s |

▾ 8.5% Memory
  2.5% Memset
  79.9% HtoD memcpy
  17.6% DtoH memcpy

# Suspicious Tensor Core Conversions

cudaMalloc 1: 4 microseconds

cudaMalloc 2: 111 microseconds

FP32to16 call 1: 5 microseconds

FP32to16 call 2: 31 microseconds

# Solution: Overlapping!

cudaMallocAsync 2 //use a stream
cudaMalloc 1
FP32to16 call 1
Synchronize
FP32to16 call 2

# Questions?

**Nsight Compute:**
https://developer.nvidia.com/tools-overview/nsight-compute/get-started