**GPU WORKLOAD DESIGN**

NVIDIA

# AGENDA

GPU Architecture & Programming Model

Are there opportunities to improve performance?

WHY GPUS?

What is CUDA actually?

**Depending on the context, "CUDA" can refer to different things**
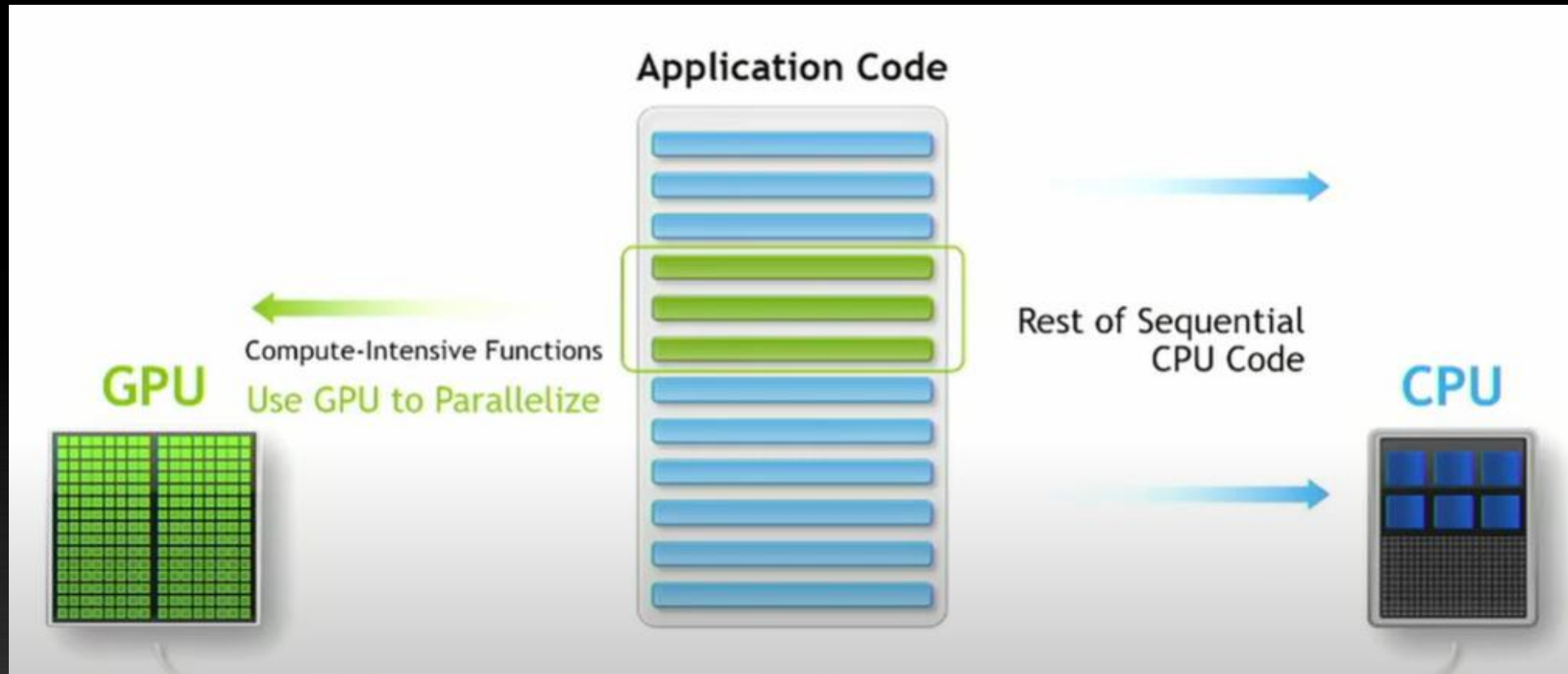
A high-level device architecture

Parallel programming model for architecture with that design

Software platform that extends high-level languages like C to add that programming model

# WE NEED BOTH

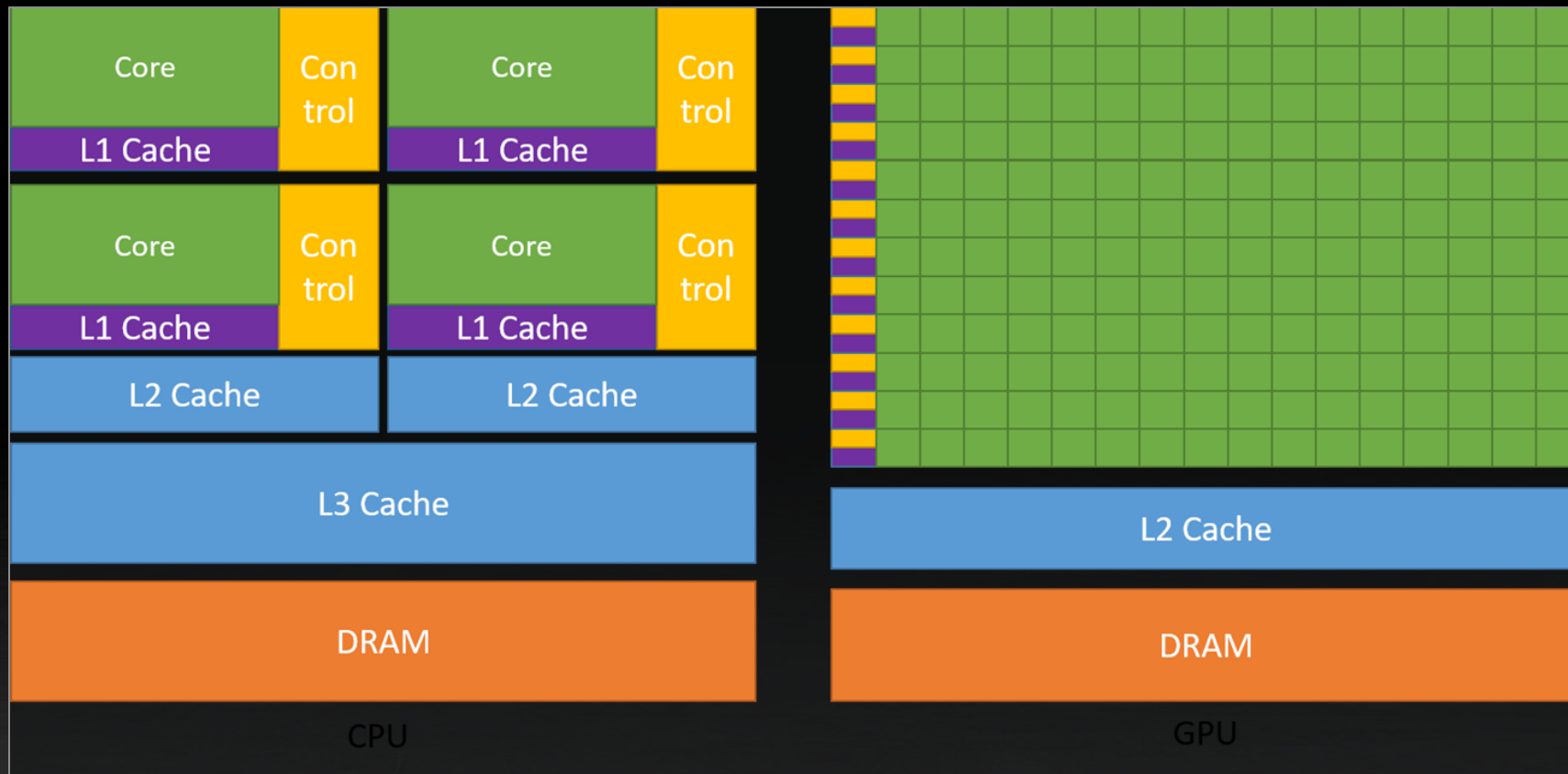Applications contain both sequential and parallel sections

# GPU: throughput-Oriented design

- **Focus**: Process massive numbers of parallel tasks efficiently (high throughput across thousands of threads)

- **Transistor Allocation**: majority dedicated to data processing rather than data caching and flow control

- **Execution:**
  - Many simple cores: No complex out-of-order execution or branch prediction, reducing per-core overhead
  - Optimized for parallel processing
  - Can hide memory access latencies with computation

# CPU: latency-oriented design

- **Focus**: Complete single tasks as quickly as possible

- **Transistor Allocation**: majority dedicated to flow control and large caches

- **Execution:**
  - Small number of high-quality, powerful, and efficient cores
  - Optimized for sequential processing (low latency per thread)
  - Relying on large data caches and complex flow control

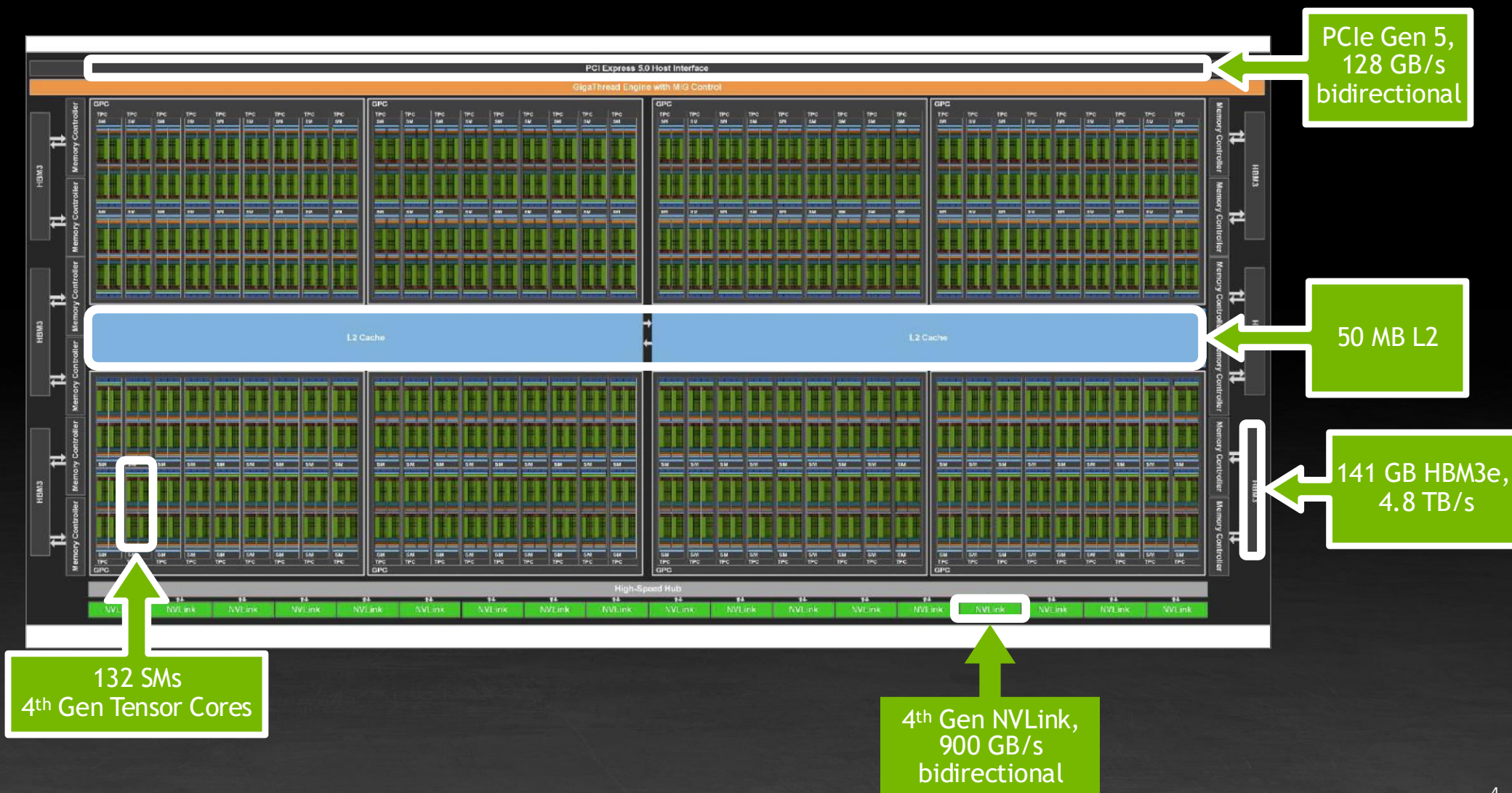**☺ nVIDIA.**

# CPU vs GPU

GPU ARCHITECTURE

# GPU OVERVIEW

## NVIDIA H200 SXM



PCIe Gen 5, 128 GB/s bidirectional

50 MB L2

141 GB HBM3e, 4.8 TB/s

132 SMs 4th Gen Tensor Cores

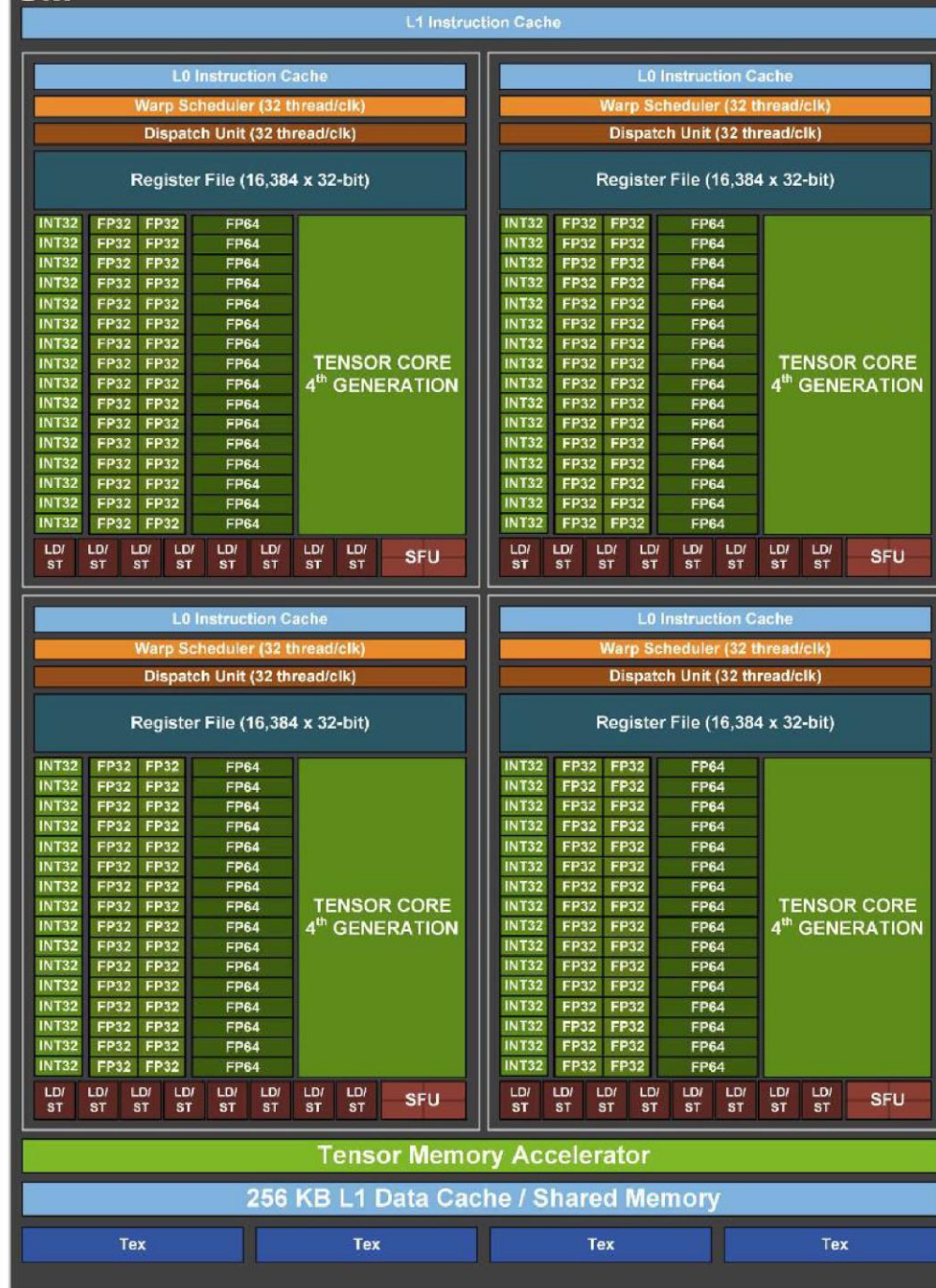4th Gen NVLink, 900 GB/s bidirectional

# STREAMING MULTIPROCESSOR (SM)
## HOPPER ARCHITECTURE

- 128 FP32 cores
- 64 FP64 cores
- 64 INT32 cores
- 4 mixed-precision Tensor Cores
- 16 special function units (transcendentals)
- 4 warp schedulers

- 32 LD/ST units
- 64K 32-bit registers
- 256 KiB unified L1 data cache and shared memory
- Tensor Memory Accelerator (TMA)

# HARDWARE EXECUTION UNITS

- SMs are the fundamental compute units of NVIDIA GPUs, analogous to CPU cores but optimized for massive parallelism

- Key components:
  - CUDA cores: execute scalar arithmetic instructions
  - Tensor cores: operate on entire matrices. These cores are much larger and less numerous than CUDA cores
  - Warp schedulers: Manage thread execution by deciding which group of 32 threads (warps) to run
  - Context switching between warps occurs in one clock cycle

**NVIDIA.**

# CUDA PROGRAMMING MODEL
## SINGLE-PROGRAM MULTIPLE-DATA

- SIMT instructions specify the execution of a **single** thread.

- A SIMT kernel is launched on many threads that execute in parallel.

- Threads use their thread index to work on disjoint data or to enable different execution paths.

- Three key software abstractions enable efficient programming through the CUDA programming model:
  - a hierarchy of **thread groups**,
  - **memory spaces**, and
  - **synchronization**.

**Single-threaded CPU vector addition**

```
for (int i = 0; i <
    N; i++) { c[i] =
    a[i] + b[i];

}
```
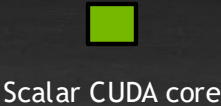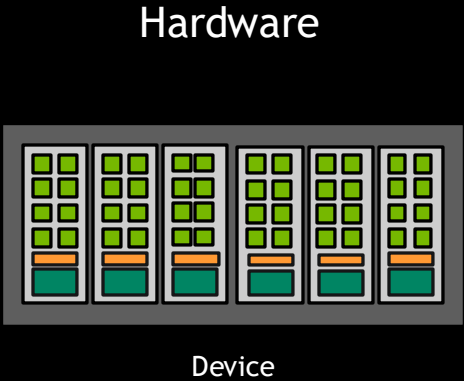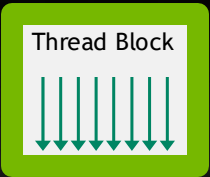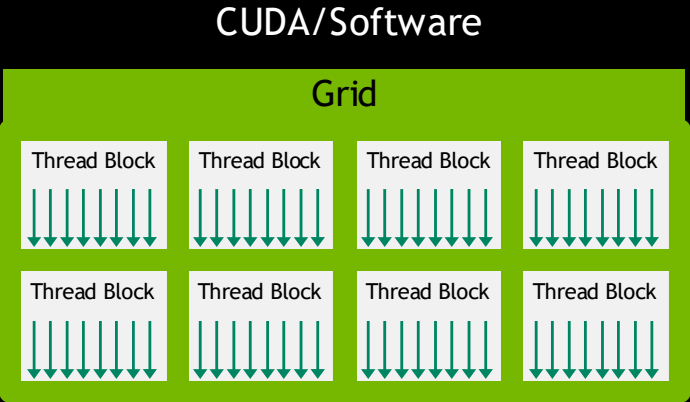
**GPU vector addition**

Thread IDs:

| 0 | 1 | 2 | ... | N-1 | N |

```
int i =
my_global_thread_id();
if (i < N) c[i] = a[i]
+ b[i];
```

NVIDIA.

# THREAD HIERARCHY

## CUDA/Software

### Grid

| Thread Block | Thread Block | Thread Block | Thread Block |
| Thread Block | Thread Block | Thread Block | Thread Block |

### Thread Block

Thread

## Hardware

Device

SM

Scalar CUDA core

- A CUDA kernel is launched on a grid of thread blocks, which are completely independent.

- Thread blocks are executed on SMs.
  - Several concurrent thread blocks can reside on an SM.
  - Thread blocks do not migrate.
  - Each block can be scheduled on any of the available SMs, in any order, concurrently or in series.

- Individual threads execute on scalar CUDA cores.

# HOW IS THIS EXECUTED?

ROOFLINE EXAMPLE STEP THROUGH

# ROOFLINE EXAMPLE WALKTHROUGH

- Step through the CuPy code:

 https://github.com/marcelo-alvarez/datasci211/blob/main/week-2/roofline.py

NVIDIA.

What is just-in-time compilation?

# NUMERICAL COMPUTING IN PYTHON



- Mathematical focus
- Operates on arrays of data
  - *ndarray*, holds data of same type
- Many years of development
- Highly tuned for CPUs

- NumPy like interface
- Trivially port code to GPU
- Copy data to GPU
  - CuPy *ndarray*
- Data interoperability with DL frameworks, RAPIDS, and Numba
- Uses high tuned NVIDIA libraries
- Can write custom CUDA functions

# CUPY

**BEFORE**

**AFTER**

```
import numpy as np

size = 4096
A = np.random.randn(size,size)

Q, R = np.lingalg.qr(A)
```

```
import cupy as cp

size = 4096
A = cp.random.randn(size,size)

Q, R = cp.lingalg.qr(A)
```



*52x Speedup!*

# KERNEL OVERHEAD



## JIT Compilation

- What is the size of A?
- What is the datatype?
- Which GPU-accelerated libraries are available?
- Compiler optimizations for custom kernels

# SETTING UP CUPY SCRIPT

```python
 8    import numpy as np
 9    import cupy as cp
10
11    # KERNEL LOADING: Use CuPy RawModule to JIT-compile CUDA kernel
12    #
13    # RawModule workflow:
14    # 1. Read CUDA source code (roofline_kernel.cuh)
15    # 2. JIT-compile using NVRTC (NVIDIA Runtime Compiler)
16    # 3. Load compiled PTX/CUBIN into current GPU context
17    # 4. Extract kernel function by name
18    #
19    # This ensures the CUDA (roofline.cu) and CuPy implementations use
20    # IDENTICAL device code, enabling fair performance comparison.
21    with open('roofline_kernel.cuh', 'r') as f:
22        _kernel_code = f.read()
23    _module = cp.RawModule(code=_kernel_code)
24    _COMPUTE_K_KERNEL = _module.get_function('compute_k_terms')
25
```

- Loads CUDA kernel code from a file
- Just-in-time compiles it
- Picks the GPU kernel we want in our script

# SOURCE CUDA CODE

```
 1    /*
 2     * Roofline Analysis Polynomial Kernel
 3     *
 4     * Computes b[i] = a[i] + a[i]^2 + ... + a[i]^k for roofline model demonstration.
 5     * Varies K to sweep arithmetic intensity from memory-bound to compute-bound.
 6     * See README.md for detailed explanation.
 7     */
 8
 9    #pragma once
10
11    /*
12     * compute_k_terms: Polynomial evaluation kernel
13     *
14     * Parameters:
15     *   a: Input array (const __restrict__ enables compiler optimizations)
16     *   b: Output array (__restrict__ guarantees no aliasing)
17     *   n: Number of elements
18     *   k: Polynomial degree (controls arithmetic intensity: AI = k/4)
19     *
20     * extern "C" linkage: Required for CuPy RawModule interoperability
21     */
```

- Setup to do our roofline model computation

NVIDIA.

# CUDA CODE IMPLEMENTATION

```
22    extern "C" __global__ void compute_k_terms(const float* __restrict__ a,
23                                                float* __restrict__ b,
24                                                int n,
25                                                int k) {
26        // Calculate global thread ID
27        int i = blockIdx.x * blockDim.x + threadIdx.x;
28
29        // Boundary check: ensure we don't access beyond array bounds
30        if (i < n) {
31            // Load input value from global memory (4-byte read)
32            float x = a[i];
33
34            // Initialize accumulator and running power of x
35            float result = 0.0f;
36            float power = x;  // Start with x^1
37
38            // Compute polynomial sum: x + x^2 + ... + x^k
39            // Loop performs K iterations:
40            //   - Each iteration: 1 addition (result += power) + 1 multiplication (power *= x)
41            //   - Total: 2K FLOPs
42            // All operations use registers only - no memory traffic
43            for (int j = 0; j < k; ++j) {
44                result += power;    // Accumulate current power term
45                power *= x;         // Advance to next power
46            }
47
48            // Write result to global memory (4-byte write)
49            // Total memory traffic: 4 bytes (read) + 4 bytes (write) = 8 bytes
50            b[i] = result;
51        }
52    }
```

- Identify unique thread running in a kernel
- Each thread computes a polynomial sum for one element
- Write results to global memory

NVIDIA.

# CUDA SIDE STEP

```
__global__ void helloWorld()
{
    printf("hello world from device\n");
}
int main(
{
    helloWorld<<<1,1>>>();
    cudaDeviceSynchronize();
    return 0;
}
```

Kernel function

Blocks per grid

Threads per block

Kernel invocation

Host –device synchronization

NVIDIA.

# KERNEL

- Kernels are C++ functions prefixed with `__global__` declaration specifier

- `__global__` prefix defines a function that is called by the host (CPU) and executed by the device (GPU)

- Kernels are executed N times in parallel across N CUDA threads

- Kernels never return a value. Host and device cannot communicate directly. Instead, data needs to be copied back and forth between them

```cpp
__global__ void helloWorld()
{
    printf("hello world from device \n");
}
```

NVIDIA.

# SETUP PYTHON TIMING CODE

```python
27  ∨  def run_sweep(n: int, k: int, a: cp.ndarray, b: cp.ndarray):
28          """Measure kernel performance for a specific K value.
29
30          This function implements GPU benchmarking best practices:
31
32          1. WARMUP RUNS: Execute kernel 3 times before timing to eliminate:
33              – JIT compilation overhead (NVRTC on first launch)
34              – GPU frequency scaling effects
35              – Cache cold–start effects
36
37          2. MULTIPLE TIMED RUNS: Collect 10 timing measurements to:
38              – Compute reliable mean performance
39              – Detect timing variance (system interference)
40
41          3. EVENT–BASED TIMING: Use CUDA events for GPU–side timing:
42              – Eliminates CPU–GPU synchronization overhead
43              – Provides microsecond precision
44              – Correctly measures asynchronous kernel execution
45
46          Parameters:
47              n: Number of array elements
48              k: Polynomial degree (controls arithmetic intensity)
49              a: Input CuPy array (device memory)
50              b: Output CuPy array (device memory)
51
52          Returns:
53              Tuple of (arithmetic_intensity, gflops, bandwidth, percent_peak)
54          """
55
```

- Drop first couple of runs to account for compilation
- Calculate only the GPU compute timing with events

NVIDIA.

# BENCHMARKING PYTHON CODE

```python
56    # KERNEL LAUNCH CONFIGURATION:
57    # Block size: 256 threads (typical for good occupancy on modern GPUs)
58    # Grid size: Ceiling division to cover all elements
59    block = 256
60    grid = (n + block - 1) // block
61
62    # ARITHMETIC INTENSITY CALCULATION:
63    # Per element: k additions + k multiplications = 2k FLOPs
64    # Memory access: 1 read (4 bytes) + 1 write (4 bytes) = 8 bytes
65    # AI = FLOPs / Bytes = 2k / 8 = k/4 FLOPs/byte
66    flops_per_element = 2.0 * k  # k additions + k multiplications
67    bytes_per_element = 8.0       # 1 read (4B) + 1 write (4B)
68    ai = flops_per_element / bytes_per_element
69
70    print(f"\nK={k} (AI={ai:.3f} flops/byte):")
71
72    # WARMUP PHASE: Run kernel 3 times to stabilize GPU state
73    # First few launches may be slower due to:
74    # - NVRTC JIT compilation (CuPy compiles kernels on first use)
75    # - GPU power state transitions (boost clocks)
76    # - Cache warming
77    for w in range(3):
78        # Launch kernel with grid/block dimensions and arguments
79        # Arguments: (a, b, n, k) - must match kernel signature
80        # Note: Python int must be converted to np.int32 for correct type
81        _COMPUTE_K_KERNEL((grid,), (block,), (a, b, np.int32(n), np.int32(k)))
82        cp.cuda.Stream.null.synchronize()  # Wait for completion
83
```

- Set thread block/grid sizes for full GPU use
- Calculate FLOPS/byte
- Launches and syncs the benchmark kernel with CuPy for accurate results

# RUNNING OUR KERNEL

```python
84      # TIMED RUNS PHASE: Measure performance over 10 iterations
85      nruns = 10
86      times = []
87
88      # Create CUDA events for precise GPU timing
89      start = cp.cuda.Event()
90      stop = cp.cuda.Event()
91
92      for i in range(nruns):
93          # Record start event in default stream
94          start.record()
95
96          # Launch kernel
97          _COMPUTE_K_KERNEL((grid,), (block,), (a, b, np.int32(n), np.int32(k)))
98
99          # Record stop event
100         stop.record()
101         stop.synchronize()  # Wait for stop event to complete
102
103         # Compute elapsed time in milliseconds
104         times.append(cp.cuda.get_elapsed_time(start, stop))
105
```

- Tell the kernel how many threads to use, and how to arrange them
- Run our kernel 10 times
- Capture timing for each iteration

```
106     # STATISTICS: Calculate mean and RMS deviation
107     mean_ms = sum(times) / nruns
108
109     # RMS (Root Mean Square) deviation measures timing stability
110     # High RMS (>5%) indicates:
111     # - System interference (other GPU workloads)
112     # - Thermal throttling
113     # - GPU frequency variations
114     sum_sq_diff = sum((t - mean_ms)**2 for t in times)
115     rms = np.sqrt(sum_sq_diff / nruns)
116     rms_percent = (rms / mean_ms) * 100.0
117
118     # PERFORMANCE METRICS:
119     # GFLOPS = (Total FLOPs / 10^9) / (Time in seconds)
120     mean_gflops = (n * flops_per_element / 1e9) / (mean_ms / 1e3)
121
122     # Effective bandwidth = (Total bytes / 10^9) / (Time in seconds)
123     mean_bw = (n * bytes_per_element / 1e9) / (mean_ms / 1e3)
124
125     # Percentage of H100's theoretical peak FP32 performance (67 TFLOPS)
126     percent_peak = 100.0 * mean_gflops / 67000.0
127
128     print(f"  Mean: {mean_ms:.3f} ms ({mean_gflops:.2f} GFLOPS, {mean_bw:.2f} GB/s, {percent_peak:.2f}% peak)", end="")
129
130     # Warn if timing variance exceeds 5% threshold
131     RMS_TOLERANCE = 5.0   # 5% tolerance
132     if rms_percent > RMS_TOLERANCE:
133         print(f" [WARNING: RMS={rms_percent:.2f}% > {RMS_TOLERANCE:.1f}%]")
134     else:
135         print()
136
137     return ai, mean_gflops, mean_bw, percent_peak
138
```

- Calculate and summarize GPU metrics for our kernel

NVIDIA.

# FINALLY, TIME FOR THE MAIN!

```python
140    if __name__ == "__main__":
141        # MAIN EXECUTION: Run roofline analysis sweep
142
143        # Problem size: 2^26 elements = 67M elements
144        # Array size: 67M * 4 bytes/float = 268 MB per array
145        # - Large enough to avoid L2 cache effects (~40 MB on H100)
146        # - Small enough for quick iteration
147        n = 1 << 26  # 67M elements
148
149        # Allocate and initialize GPU arrays
150        # Input:  Fill with 1.01 (avoids overflow for K<1000)
151        # Output: Zero-initialized (will be overwritten by kernel)
152        a = cp.full(n, 1.01, dtype=cp.float32)
153        b = cp.zeros(n, dtype=cp.float32)
154
155        # K values sweep: Chosen to span memory-bound to compute-bound regions
156        # - K=1:    AI=0.25 FLOPs/byte  (deeply memory-bound)
157        # - K=50:   AI=12.5 FLOPs/byte  (near ridge point ~19.7)
158        # - K=1000: AI=250 FLOPs/byte   (deeply compute-bound)
159        k_values = [1, 2, 5, 10, 20, 50, 100, 200, 500, 1000]
160
161        # Write results to CSV for plotting
162        with open("roofline_cupy.csv", "w") as f:
163            # CSV header
164            f.write("k,ai,gflops,bandwidth,percent_peak\n")
165
166            # Run performance sweep across K values
167            for k in k_values:
168                ai, gflops, bw, pct_peak = run_sweep(n, k, a, b)
169
170                # Write results: k, arithmetic_intensity, GFLOPS, bandwidth, percent_peak
171                f.write(f"{k},{ai:.6f},{gflops:.2f},{bw:.2f},{pct_peak:.2f}\n")
```

- Setup problem size
- Increasing k means more floating-point operations per memory access.

- Low k: memory-bound, limited by memory speed.

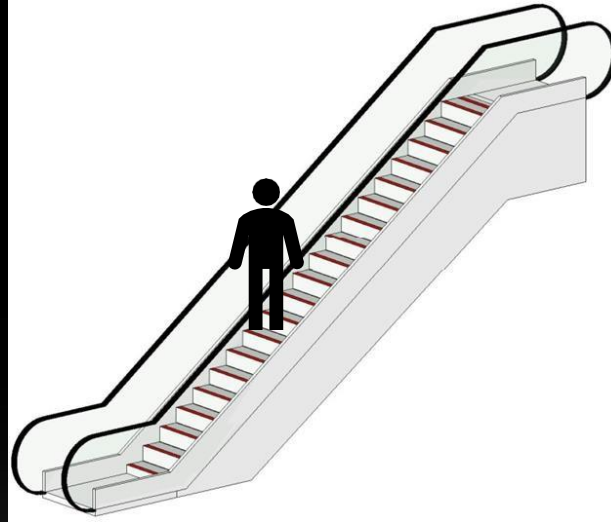- High k: compute-bound, limited by GPU processing speed.

NVIDIA.

LITTLE'S LAW

# LITTLE'S LAW
## FOR ESCALATORS

Our escalator parameters:

- 1 person per step

- A step arrives every 2 seconds
  - **Bandwidth**: 0.5 person/s



- 20 steps tall
  - **Latency** = 40 seconds
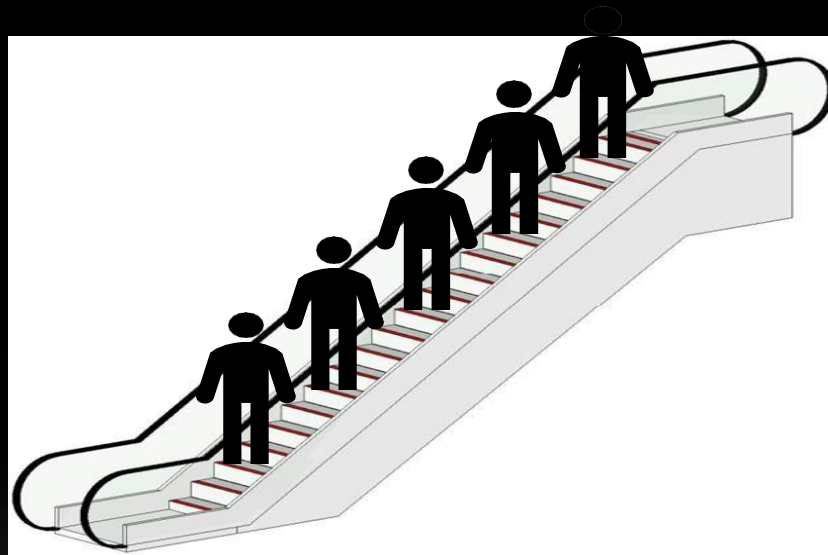
**One person in flight?**

Throughput = 0.025 person/s

# LITTLE'S LAW
## FOR ESCALATORS

Our escalator parameters:

- 1 person per step

- A step arrives every 2 seconds
  - **Bandwidth**: 0.5 person/s



- 20 steps tall
  - **Latency** = 40 seconds

**How many persons do we need in-flight to saturate bandwidth?**

Concurrency = Bandwidth x Latency
= 0.5 persons/s x 40 s
= 20 persons

# LITTLE'S LAW
## FOR GPUS

- How to maximize performance?

  1. Saturate compute units.
  2. Saturate memory bandwidth.

- Need to hide the corresponding latencies to achieve this.

| | |
|---|---|
| FP32 | FP32 |
| FP32 | FP32 |
| FP32 | FP32 |
| FP32 | FP32 |

**FP32 Latency = 24 cycles 8 FP32 ops per cycle**

- Compute latencies.
- Memory access latencies.

- Latencies can be hidden by having more instructions in flight.

**Concurrency = Bandwidth x Latency = 8 x 24 operations in-flight**

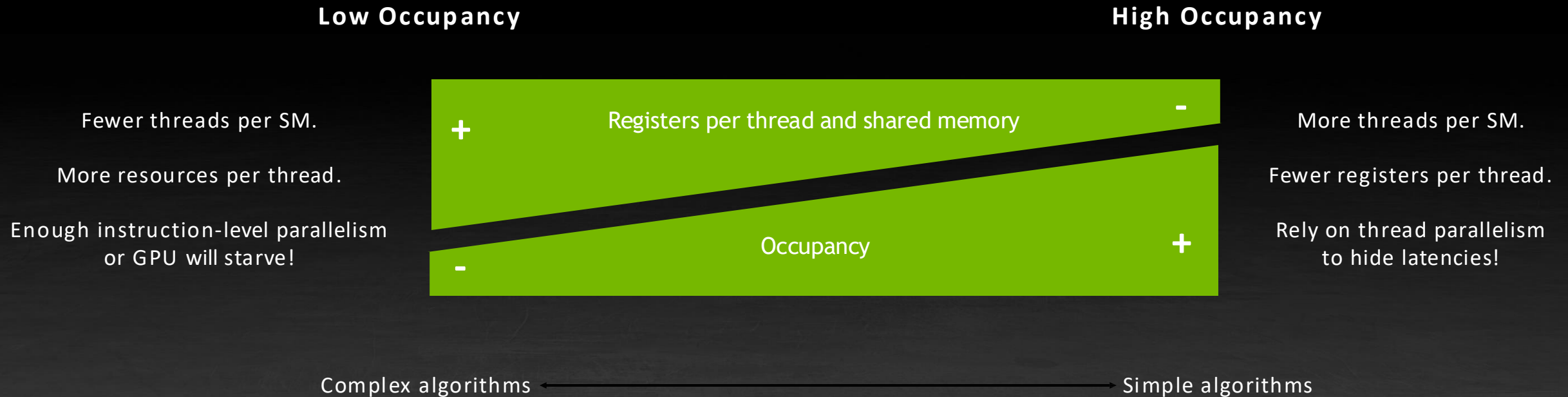NVIDIA.

Is it always this clear cut?

# WHAT OCCUPANCY DO I NEED?
## GENERAL GUIDELINES

**Rule of thumb:** Try to maximize occupancy.

But some algorithms will run better at low occupancy.

More registers and shared memory can allow higher data reuse, higher ILP, higher performance.
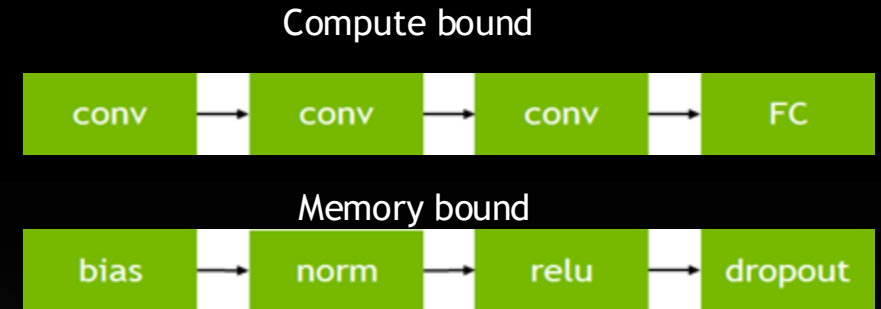
**Low Occupancy**                                    **High Occupancy**

Fewer threads per SM.      + Registers per thread and shared memory -      More threads per SM.

More resources per thread.                               Fewer registers per thread.

Enough instruction-level parallelism
or GPU will starve!      - Occupancy +      Rely on thread parallelism
to hide latencies!

Complex algorithms ←————————————→ Simple algorithms

**OTHER OPPORTUNITIES FOR OPTIMIZATION**

# FINDING PERFORMANCE OPPORTUNITIES

## Models can be data bound by the data pipeline, compute or memory

- GPU utilization as it relates to model code
  - Time being spent on ops in every iteration
  - Time spent on GPU/CPU
  - Data types used for operations
- Bottlenecks could be attributed to
  - Input data pipeline: data loading, preprocessing etc
  - Compute (math) limited operations
  - Memory limited operations
  - Other aspects such as overall system tuning
- Categories of operations in DNNs based on bottleneck
  - Element wise: ReLU, memory bound
  - Reduction: Batch norm, memory bound
  - Dot product: Convolution, math bound

Compute bound

| conv | → | conv | → | conv | → | FC |

Memory bound

| bias | → | norm | → | relu | → | dropout |

Compute heavy ops see speed-ups from GPUs

# DEEP LEARNING OPTIMIZATION
## Performance Analysis at System and DNN Level & Visualization

### System Level Tuning

- System Tuning
  - Thread Synchronization, Multi GPU and node communication
  - Memory management & Kernel profiling
- Leveraging/Optimizing Hardware
- Input Pipeline Optimization
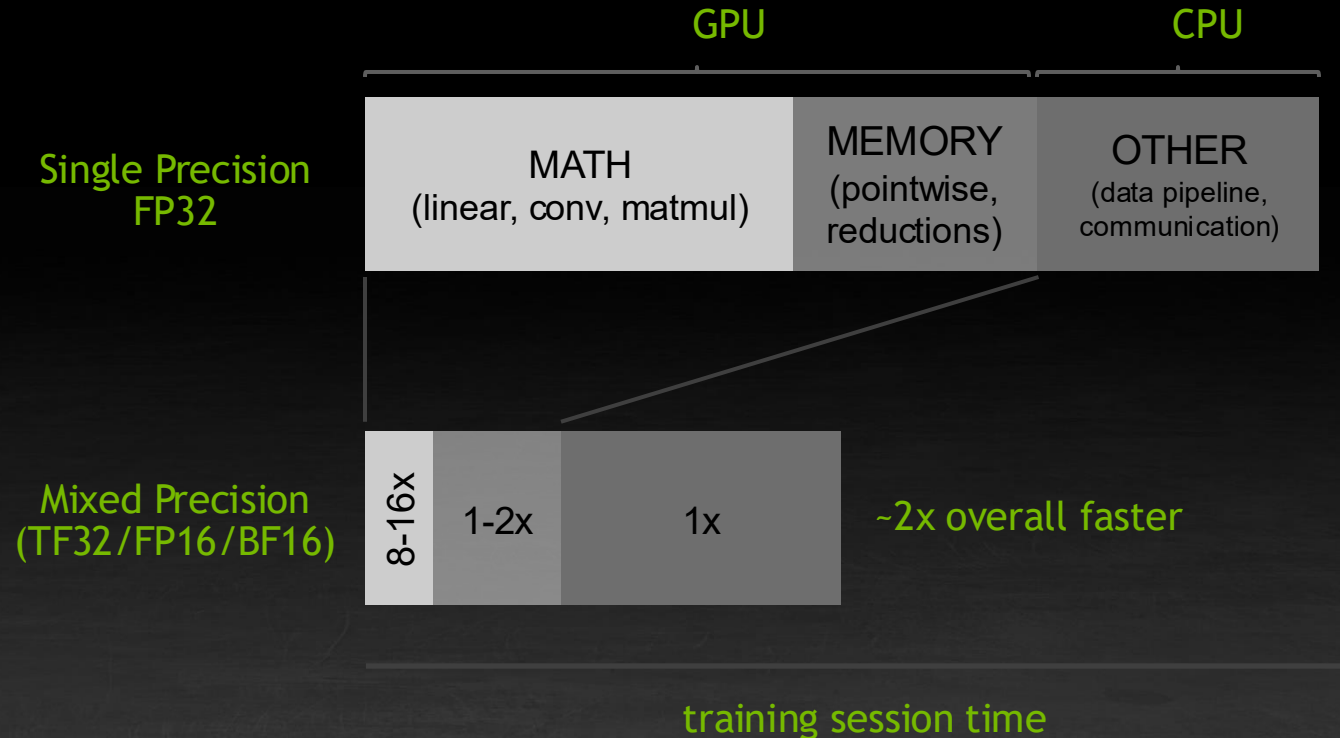- Many others....

### DNN Level Tuning

- Algorithm Techniques & Data Representations
- Pruning
- Calibration
- Quantization
- Many others....

# LIMITS OF PERFORMANCE OPTIMIZATION

## End-to-end perf depends on training composition

▪**Amdahl's Law:**
If you speed up part of your training session (GPU work), then the remaining parts (CPU work) limit your overall performance

# DL PROFILING NEEDS OF DIFFERENT PERSONAS

## Researchers

Fast development of best performant models for research, challenge and domains

## Data Scientists & Applied Researchers

Reduce Training time, focus on data, develop and apply the best models for the applications

## Sysadmins & DevOps

Optimized utilization and uptime, monitor GPU workloads, leverage hardware

# WHICH OPTIMIZATIONS TO FOCUS ON?

## SOLVING THE BOTTLENECKS

- Compute bound
  - Reduce instruction count.
    - E.g., use vector loads/stores.
  - Use tensor cores.
  - Use lower precision arithmetic, fast math intrinsics.

- Bandwidth bound
  - Reduce the amount of data transferred.
    - Optimize memory access patterns.
    - Lower precision datatypes.
    - Kernel fusion.

- Latency bound
  - Increase number of instructions and memory accesses in-flight.
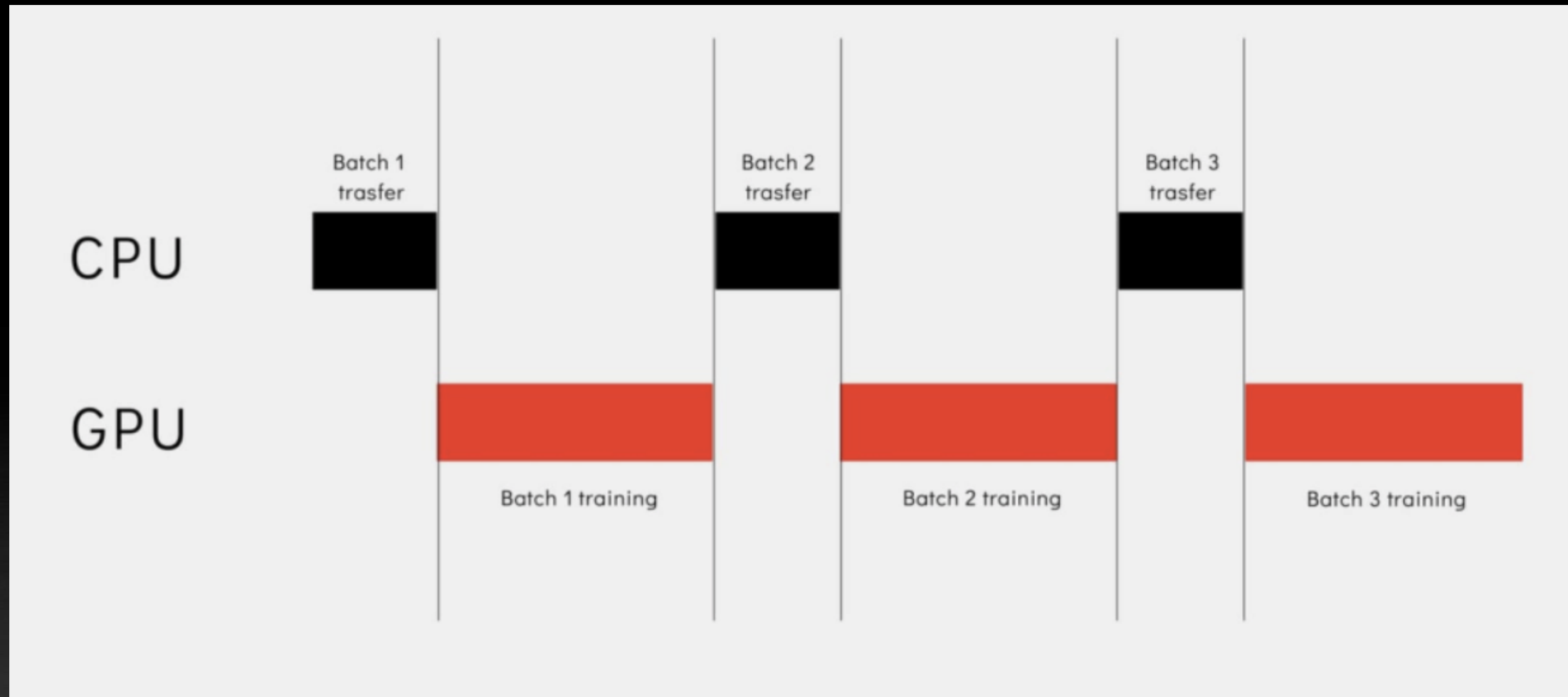  - Increase parallelism, occupancy.

NVIDIA.

PYTORCH PERFORMANCE

# WHAT ARE COMMON PROBLEMS IN PYTORCH

- **I/O and Data:**
  - **GPU starvation**
  - GPU sits idle waiting for the CPU to load and preprocess data.
  - This is common with large data sets that cannot be fully loaded into RAM.

- **GPU utilization:**
- Fully saturated with a heavy model
- Underutilized due to a lack of data

- Memory:
  - `RuntimeError: CUDA out of memory`
  - Forces a reduction in batch size.

# WHAT ARE COMMON PROBLEMS IN PYTORCH?

# WHAT CAN WE DO ABOUT IT?

- **I/O and Data:**
  - Use the fundamental abstractions called <u>Dataset and Dataloader</u>
  - <u>Dataset:</u> This is the base class that represents a set of samples and their labels.
  - Dataloader: Wraps the dataset and makes it efficiently iterable.

- **GPU utilization:**
  - PyTorch Profiler to see Detailed analysis of CUDA operators and kernels.
  - Export results in .jsoninteractive format or visualization with TensorBoard.

- **Memory:**
  - Fuse operations to reduce memory access and kernel launch times
  - Only one kernel is launched for multiple pointwise operations

- 

```
@torch.compile
def gelu(x):
    return x * 0.5 * (1.0 + torch.erf(x / 1.41421))
```

# WHAT CAN WE DO ABOUT IT?

- **GPU utilization:**
  - Automatic Mixed Precision
  - Reduced precision can mean faster computations

- **GPU utilization:**
  - Increase batch size during training
  - Number of training examples utilized in one iteration.
  - Better utilization of GPU parallelism and faster convergence.

- **Memory:**
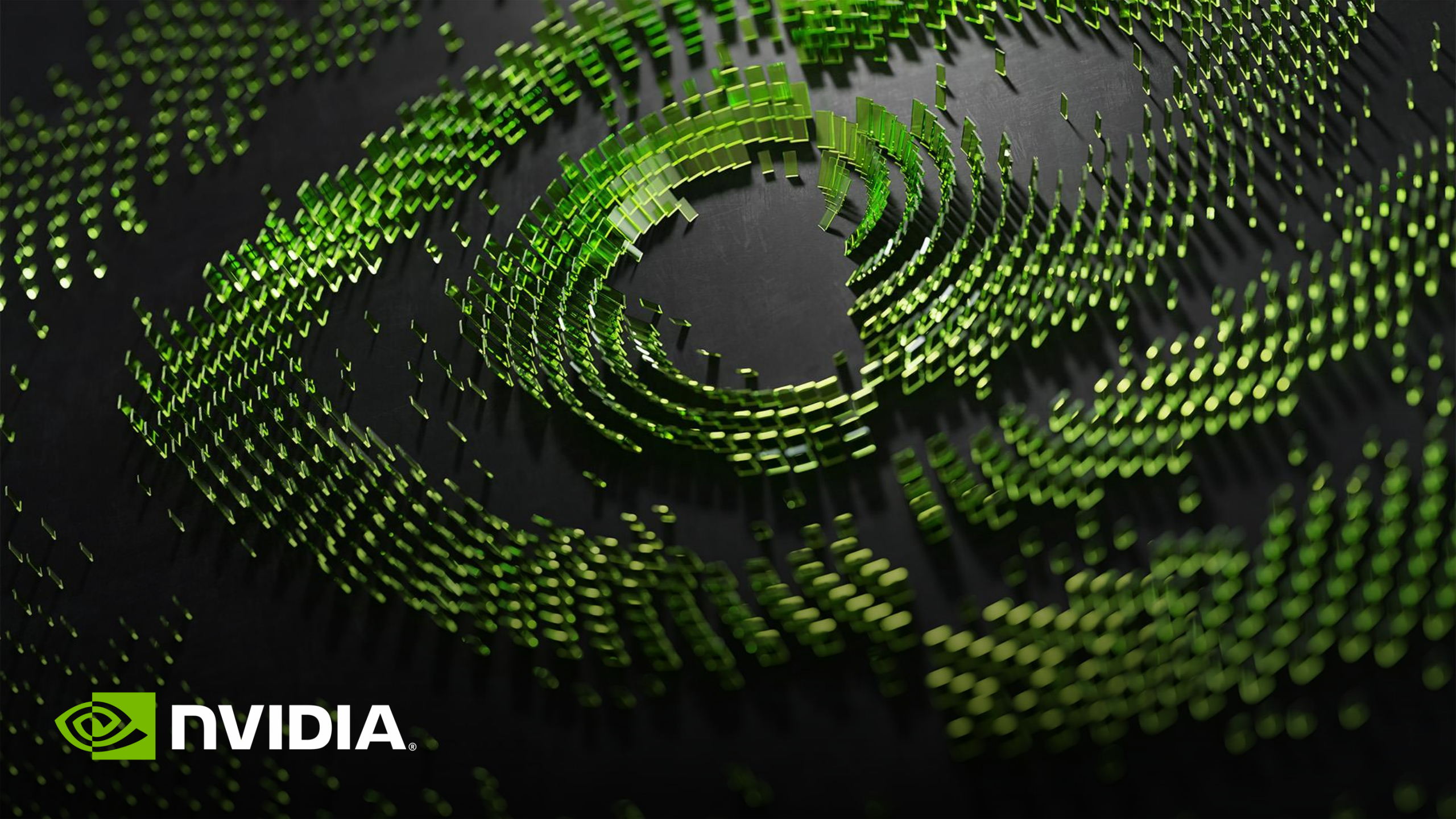  - Memory Pinning for optimizing data transfer between the CPU and GPU

- 

NVIDIA.

# WHAT CAN WE DO ABOUT IT

| Symptom detected by the Profiler | Diagnosis (Bottleneck) | Recommended solution |
|---|---|---|
| High `Self CPU total` % for `DataLoader` | Slow pre-processing and/or data loading on the CPU side | Increase `num_workers` |
| High execution time for `cudaMemcpyAsync` | Slow data transfer between CPU and GPU memory | Enable `pin_memory=True` |

NVIDIA.

# WHERE CAN WE LEARN MORE?

- [PyTorch Performance Tuning Guide](#)
- [PyTorch Optimization Guide](#)
- [NVIDIA Deep Learning Performance Guide](#)
- [Ultimate guide to PyTorch library in Python](#)
- [LLM Training on Grace Hopper](#)