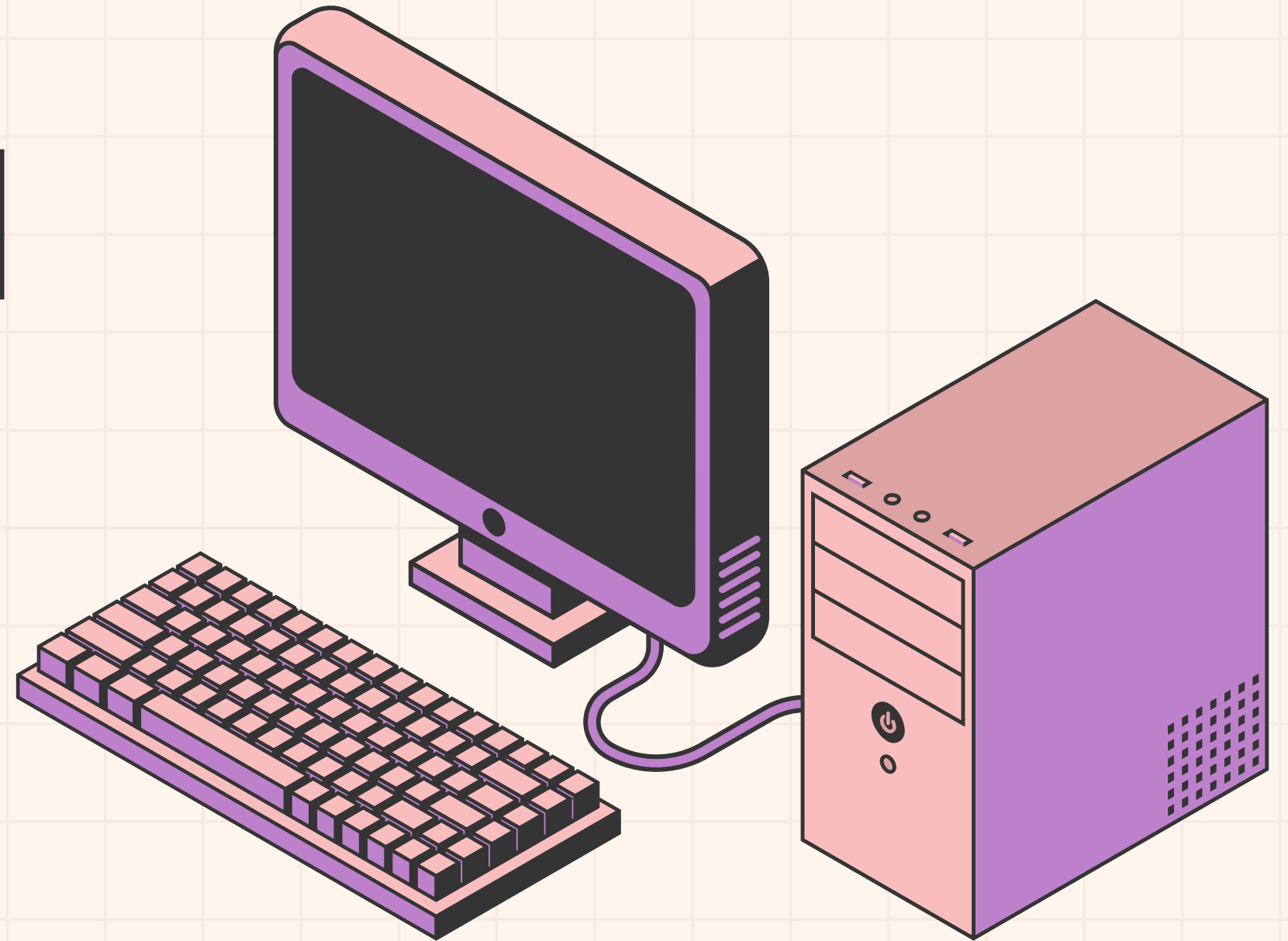# MNIST DIGIT CLASSIFICATION *ACCELERATED WITH CUDA*

Sahrish Mustafa 22iO977
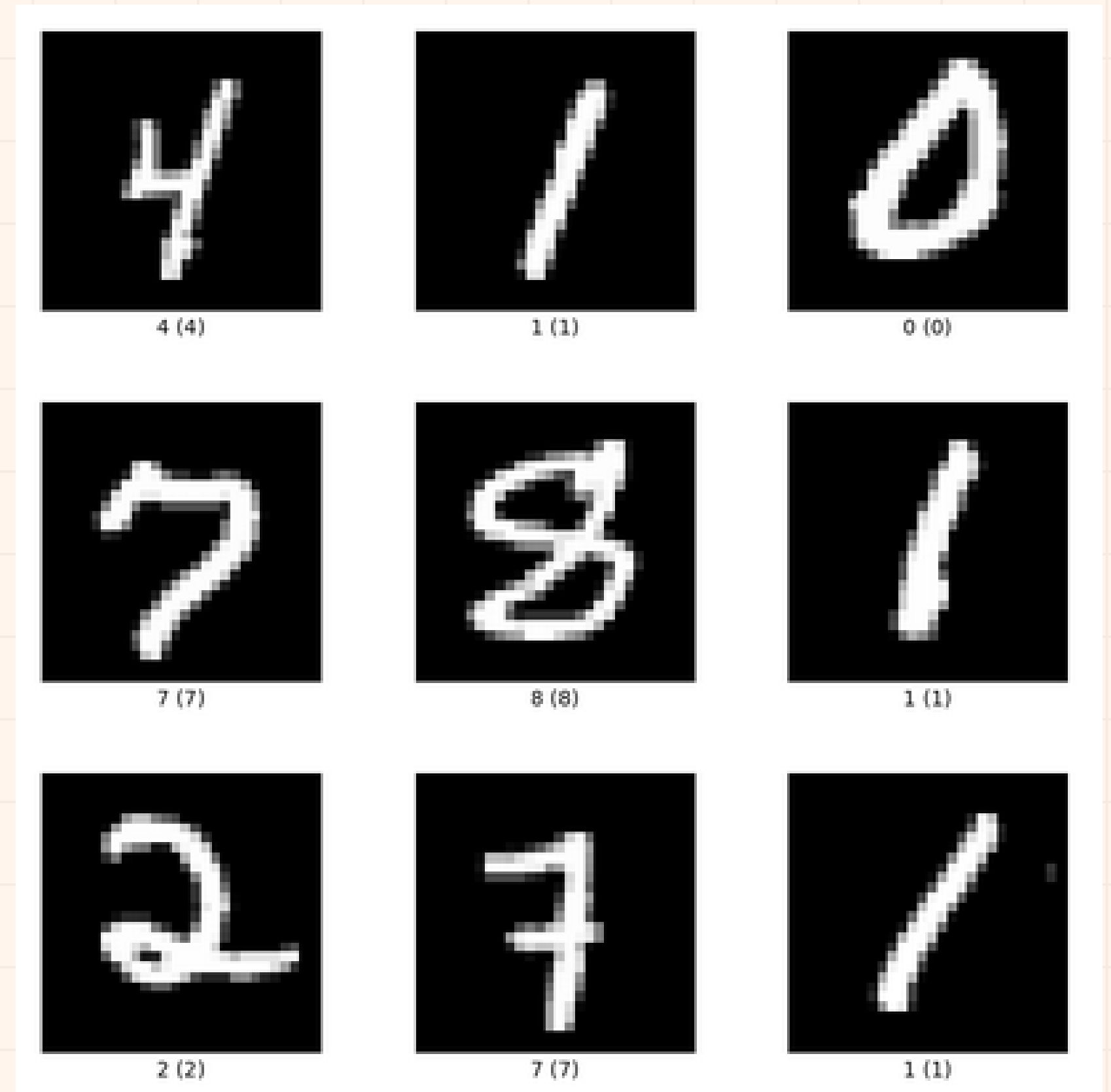
Aalyan Raza Kazmi 22iO833

# MNIST

*(Modified National Institute of Standards and Technology database)*
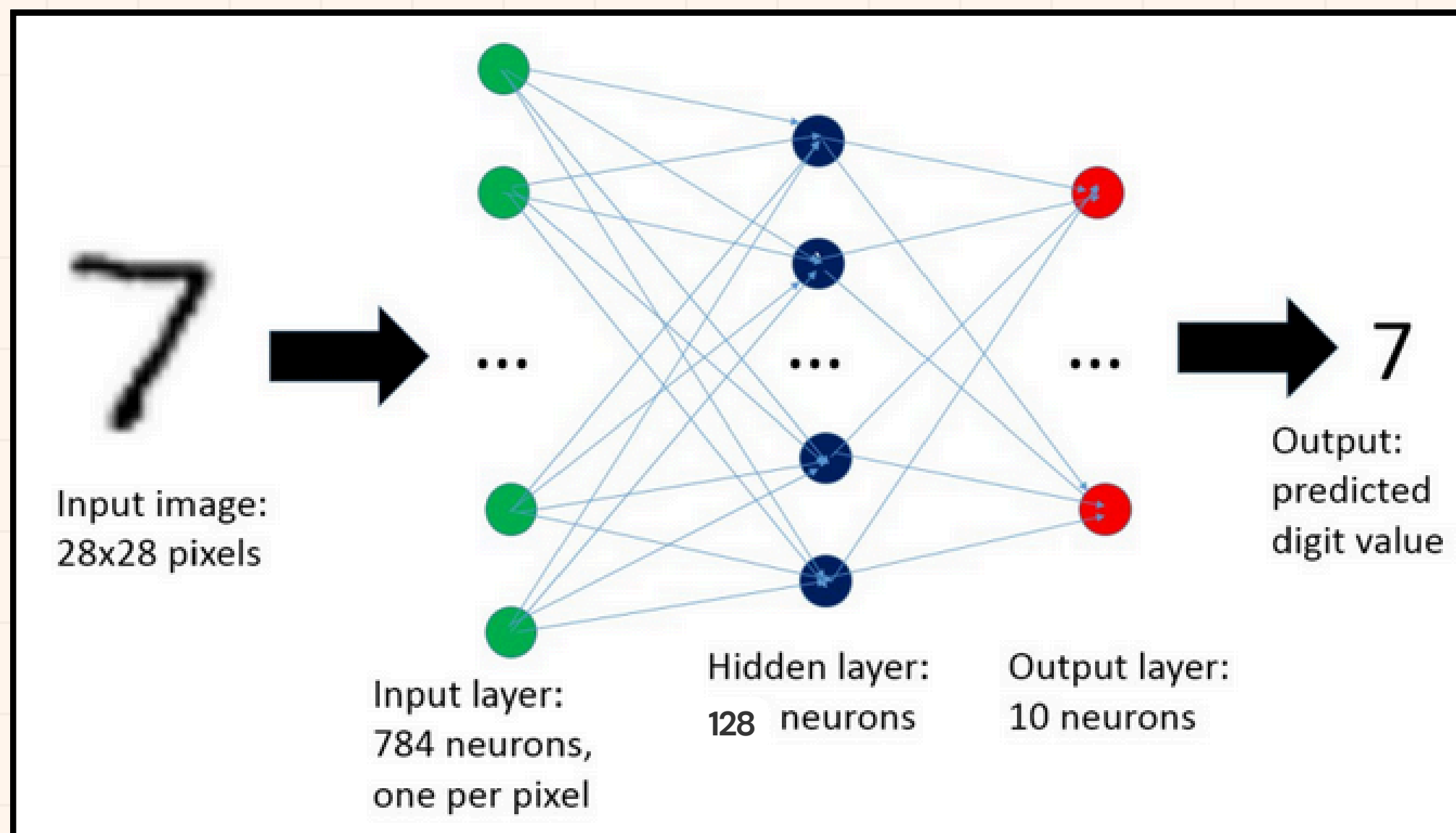
- **Handwritten Digits**
- **70,000 Images**
- **28x28 pixels**
- **10 classes (0-9)**

# NEURAL NETWORK - V1

## INPUT LAYER

- Size: 784 neurons
- 28x28 pixel grayscale images. Flattened, each image becomes a 784-dimensional vector.



Input image: 28x28 pixels

Input layer: 784 neurons, one per pixel

Hidden layer: 128 neurons

Output layer: 10 neurons

Output: predicted digit value

## OUTPUT LAYER

- Size: 10 neurons
- Weights (W2): A matrix of size 10 x 128
- Biases (b2): A vector of size 10
- Activation Function: Softmax (raw outputs to probabilities)

## HIDDEN LAYER

- Size: 128 neurons
- Weights (W1): A matrix of size 128 x 784
- Biases (b1): A vector of size 128
- Activation Function: ReLU (max(0, x))

# ACCELERATION PLAN

- parallelize forward() and backward() by converting loops to kernels running on the GPU

# FORWARD PROPAGATION

## V1

```
1. FUNCTION FORWARD_PROPAGATION(INPUT):
2.     // COMPUTE HIDDEN LAYER
3.     FOR I IN RANGE(HIDDEN_SIZE):
4.         SUM = BIAS1[I]
5.         FOR J IN RANGE(INPUT_SIZE):
6.             SUM += W1[I][J] * INPUT[J]
7.         HIDDEN[I] = RELU(SUM)
8.
9.     // COMPUTE OUTPUT LAYER
10.    FOR I IN RANGE(OUTPUT_SIZE):
11.        SUM = BIAS2[I]
12.        FOR J IN RANGE(HIDDEN_SIZE):
13.            SUM += W2[I][J] * HIDDEN[J]
14.        OUTPUT[I] = SUM  // SOFTMAX APPLIED LATER
```

## V2

```
1. KERNEL FORWARDHIDDEN(W1, B1, INPUT, HIDDEN):
2.     I = THREADIDX.X
3.     IF I < HIDDEN_SIZE:
4.         SUM = B1[I]
5.         FOR J IN RANGE(INPUT_SIZE):
6.             SUM += W1[I][J] * INPUT[J]
7.         HIDDEN[I] = RELU(SUM)
8.
9. KERNEL FORWARDOUTPUT(W2, B2, HIDDEN, OUTPUT):
10.    I = THREADIDX.X
11.    IF I < OUTPUT_SIZE:
12.        SUM = B2[I]
13.        FOR J IN RANGE(HIDDEN_SIZE):
14.            SUM += W2[I][J] * HIDDEN[J]
15.        OUTPUT[I] = SUM  // SOFTMAX ON HOST
```

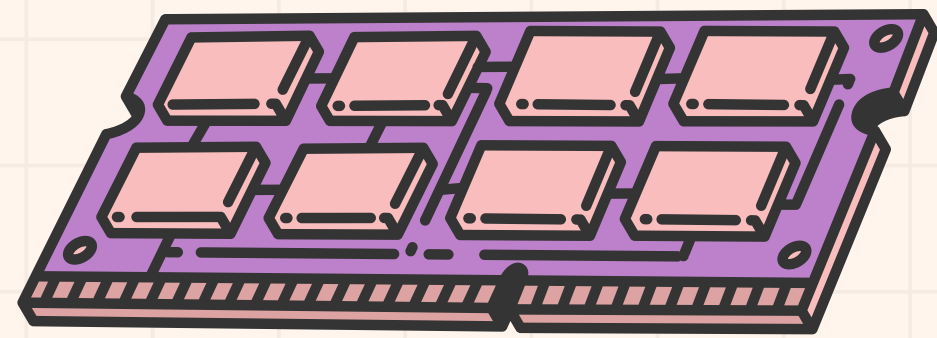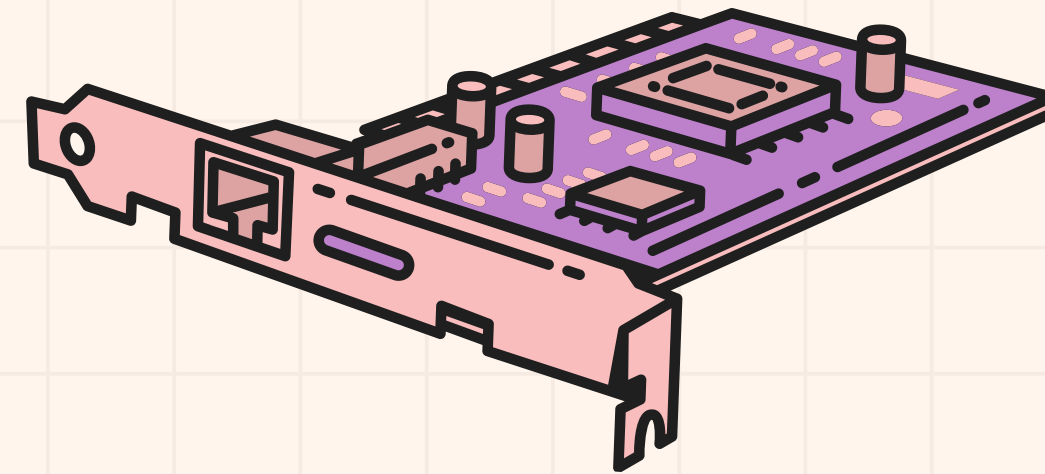| Time | Total Time | Instances | Avg | Med | Min | Max | StdDev | Category | Operation |
|---|---|---|---|---|---|---|---|---|---|
| 46.2% | 16.397 s | 1536592 | 10.671 µs | 416 ns | 351 ns | 126.848 µs | 23.794 µs | MEMORY_OPER | [CUDA memcpy Host-to-Device] |
| 30.7% | 10.896 s | 709194 | 15.363 µs | 1.696 µs | 768 ns | 176.320 µs | 30.173 µs | MEMORY_OPER | [CUDA memcpy Device-to-Host] |
| 18.9% | 6.715 s | 118200 | 56.808 µs | 56.768 µs | 56.415 µs | 65.440 µs | 435 ns | CUDA_KERNEL | forwardHidden(double *, double *, double *, double *) |
| 1.9% | 676.200 ms | 118199 | 5.720 µs | 5.696 µs | 5.663 µs | 6.912 µs | 62 ns | CUDA_KERNEL | forwardOutput(double *, double *, double *, double *) |
| 1.0% | 357.451 ms | 118199 | 3.024 µs | 3.008 µs | 2.943 µs | 4.064 µs | 43 ns | CUDA_KERNEL | updateHiddenLayer(double *, double *, double *, double *) |
| 0.5% | 174.582 ms | 118199 | 1.477 µs | 1.472 µs | 1.440 µs | 1.760 µs | 15 ns | CUDA_KERNEL | computeHiddenGradients(double *, double *, double *, double *) |
| 0.5% | 170.481 ms | 118199 | 1.442 µs | 1.440 µs | 1.407 µs | 1.920 µs | 17 ns | CUDA_KERNEL | updateOutputLayer(double *, double *, double *, double *) |
| 0.4% | 142.075 ms | 118199 | 1.202 µs | 1.216 µs | 1.183 µs | 1.568 µs | 19 ns | CUDA_KERNEL | computeOutputGradients(double *, double *, double *) |

# ACCELERATION PLAN

- TOO MUCH MEMORY TRANSFER TIME!
- …. REDUCE MEMORY TRANSFER (DTH / HTD)
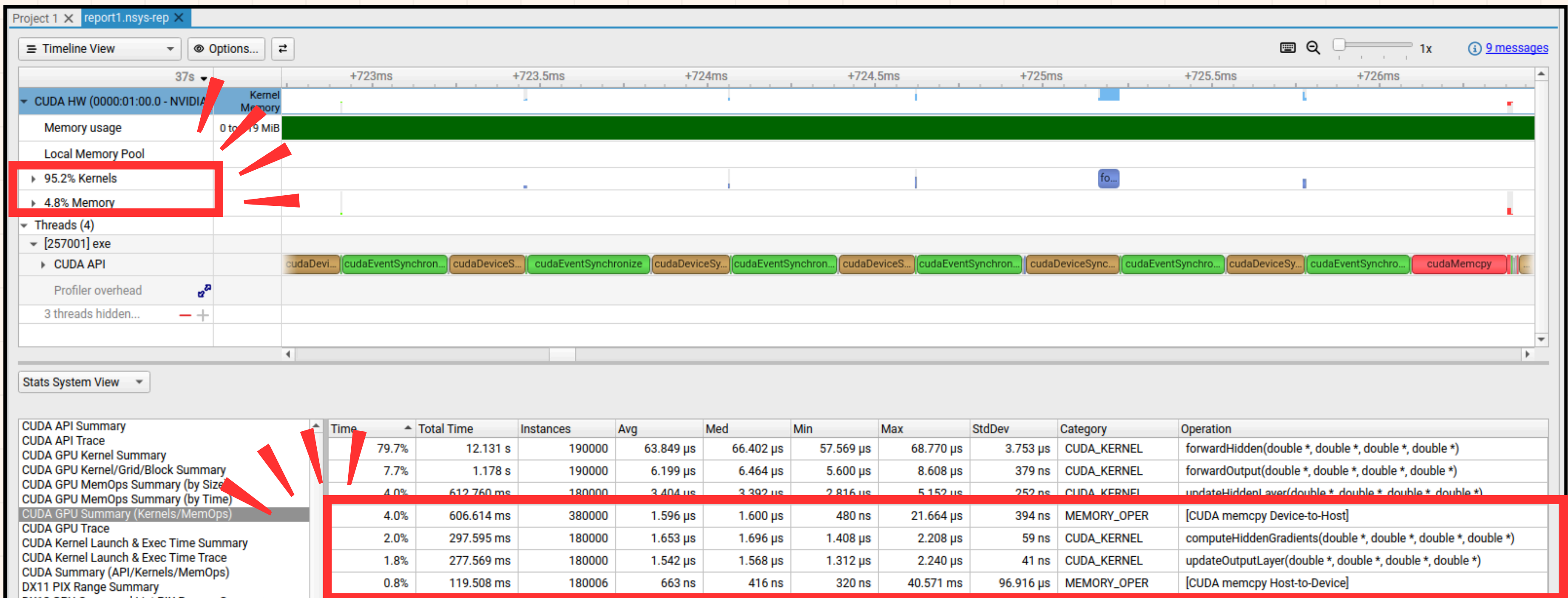
**CPU RAM (HOST STORAGE)**

previously, we were copying each image from the host to the device 60,000 times per epoch during training. **Now, we copy the entire dataset in the beginning only once**

**GPU GLOBAL MEMORY (DEVICE STORAGE)**

| Project 1 × | report1.nsys-rep × |

Timeline View | Options... | ⇄ | 1x | ⓘ 9 messages

| | 37s | +723ms | +723.5ms | +724ms | +724.5ms | +725ms | +725.5ms | +726ms |

CUDA HW (0000:01:00.0 - NVIDIA | Kernel Memory

Memory usage | 0 to 79 MiB

Local Memory Pool

▶ 95.2% Kernels

▶ 4.8% Memory

▼ Threads (4)

▼ [257001] exe

▶ CUDA API | cudaDevi... | cudaEventSynchron... | cudaDeviceS... | cudaEventSynchronize | cudaDeviceSy... | cudaEventSynchron... | cudaDeviceS... | cudaEventSynchron... | cudaDeviceSync... | cudaEventSynchro... | cudaDeviceSy... | cudaEventSynchro... | cudaMemcpy

Profiler overhead

3 threads hidden...

Stats System View

| | Time | Total Time | Instances | Avg | Med | Min | Max | StdDev | Category | Operation |
|---|---|---|---|---|---|---|---|---|---|---|
| CUDA API Summary | 79.7% | 12.131 s | 190000 | 63.849 µs | 66.402 µs | 57.569 µs | 68.770 µs | 3.753 µs | CUDA_KERNEL | forwardHidden(double *, double *, double *, double *) |
| CUDA API Trace | 7.7% | 1.178 s | 190000 | 6.199 µs | 6.464 µs | 5.600 µs | 8.608 µs | 379 ns | CUDA_KERNEL | forwardOutput(double *, double *, double *, double *) |
| CUDA GPU Kernel Summary | 4.0% | 612.760 ms | 180000 | 3.404 µs | 3.392 µs | 2.816 µs | 5.152 µs | 252 ns | CUDA_KERNEL | updateHiddenLayer(double *, double *, double *, double *) |
| CUDA GPU Kernel/Grid/Block Summary | | | | | | | | | | |
| CUDA GPU MemOps Summary (by Size) | 4.0% | 606.614 ms | 380000 | 1.596 µs | 1.600 µs | 480 ns | 21.664 µs | 394 ns | MEMORY_OPER | [CUDA memcpy Device-to-Host] |
| CUDA GPU MemOps Summary (by Time) | 2.0% | 297.595 ms | 180000 | 1.653 µs | 1.696 µs | 1.408 µs | 2.208 µs | 59 ns | CUDA_KERNEL | computeHiddenGradients(double *, double *, double *, double *) |
| CUDA GPU Summary (Kernels/MemOps) | 1.8% | 277.569 ms | 180000 | 1.542 µs | 1.568 µs | 1.312 µs | 2.240 µs | 41 ns | CUDA_KERNEL | updateOutputLayer(double *, double *, double *, double *) |
| CUDA GPU Trace | 0.8% | 119.508 ms | 180006 | 663 ns | 416 ns | 320 ns | 40.571 ms | 96.916 µs | MEMORY_OPER | [CUDA memcpy Host-to-Device] |
| CUDA Kernel Launch & Exec Time Summary | | | | | | | | | | |
| CUDA Kernel Launch & Exec Time Trace | | | | | | | | | | |
| CUDA Summary (API/Kernels/MemOps) | | | | | | | | | | |
| DX11 PIX Range Summary | | | | | | | | | | |

# ACCELERATION PLAN

- **> 13 SECONDS** ON KERNELS
- . . . SPEED UP KERNELS

# OPTIMIZED GPU VERSION - V3

1. Optimizing launch configurations
2. Shared memory
3. Half-precision (FP16) datatype
4. Coalesced memory access
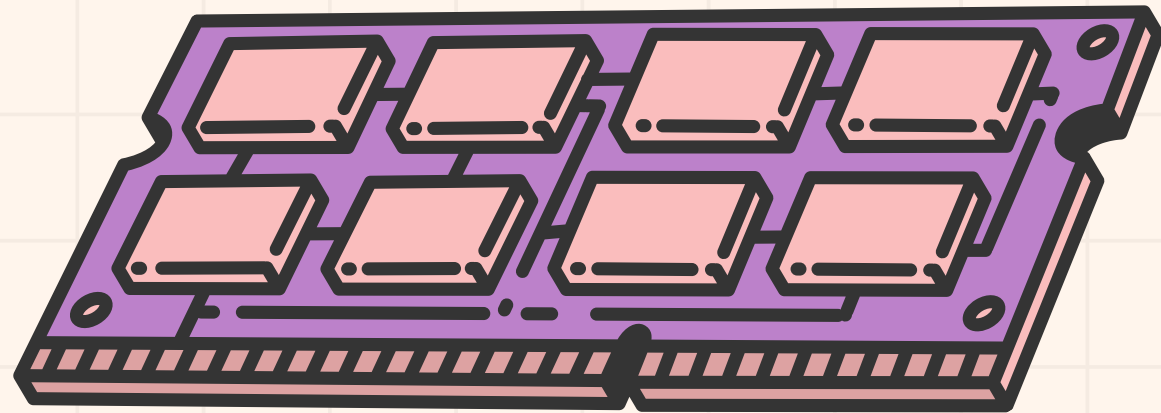
# LAUNCH CONFIGURATIONS

FORWARDHIDDEN<<<4, 32>>>

FORWARDOUTPUT<<<1, 64>>>

COMPUTEHIDDENGRADIENTS<<<16, 4>>>

UPDATEOUTPUTLAYER<<<OUTPUT_SIZE, HIDDEN_SIZE>>>

UPDATEHIDDENLAYER<<<HIDDEN_SIZE, INPUT_SIZE>>>

# OPTIMIZED GPU VERSION – V3

## SHARED MEMORY

## MEMORY COALESCION

Shared memory is used to cache input data, hidden layer activations, and output gradients within each thread block. This reduces global memory access and speeds up forward and backward passes in the neural network.
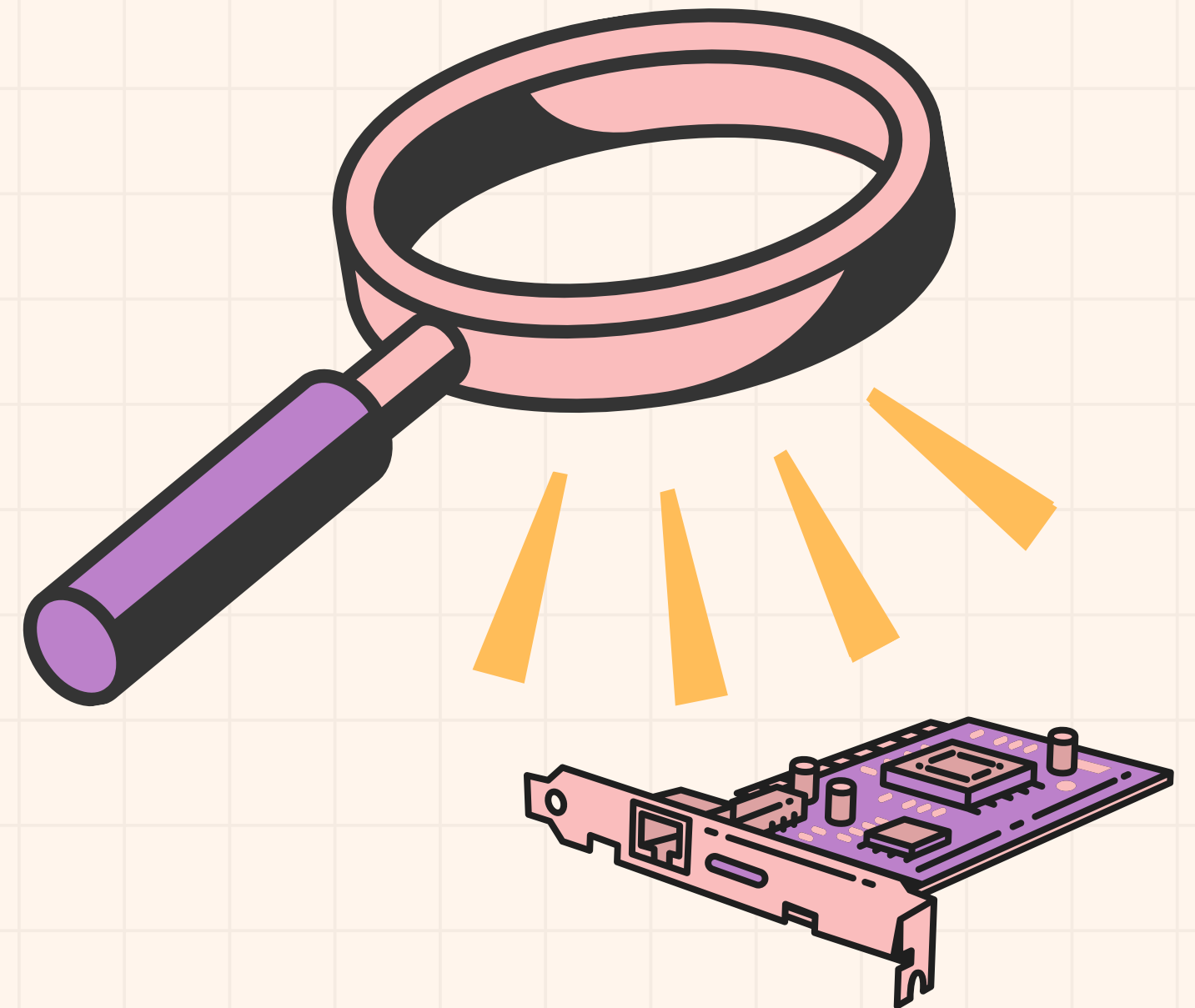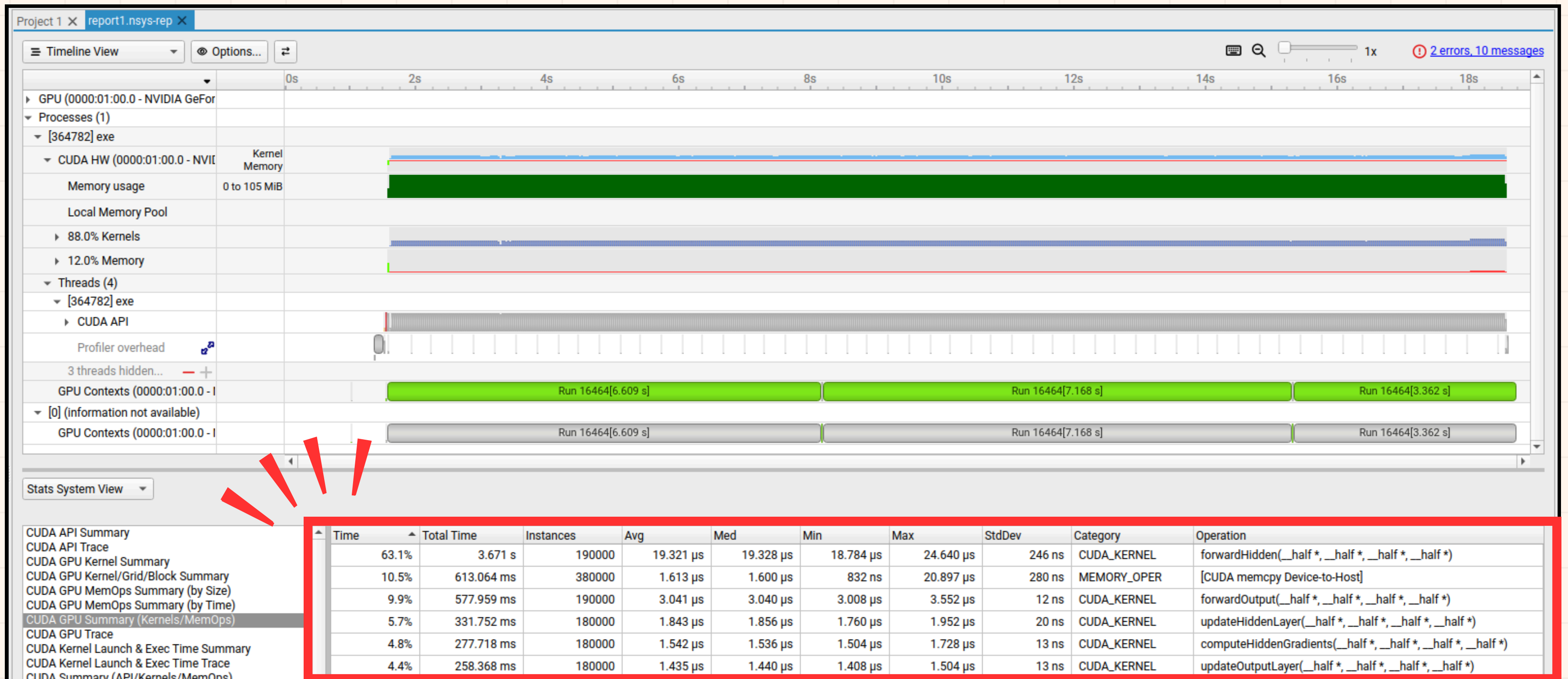
The weight matrix is transposed to ensure memory coalescing during access by parallel threads. This improves global memory access efficiency, reducing latency during matrix multiplications.

# FP 16
## HALF PRECISION

- Uses half type to store inputs, weights, and activations.
- Functions like __float2half() and __half2float() are used to convert between 32-bit and 16-bit values.
- Reduces memory usage and enables faster computation on supported GPUs.
- Ideal for deep learning workloads where precision trade-offs are acceptable.
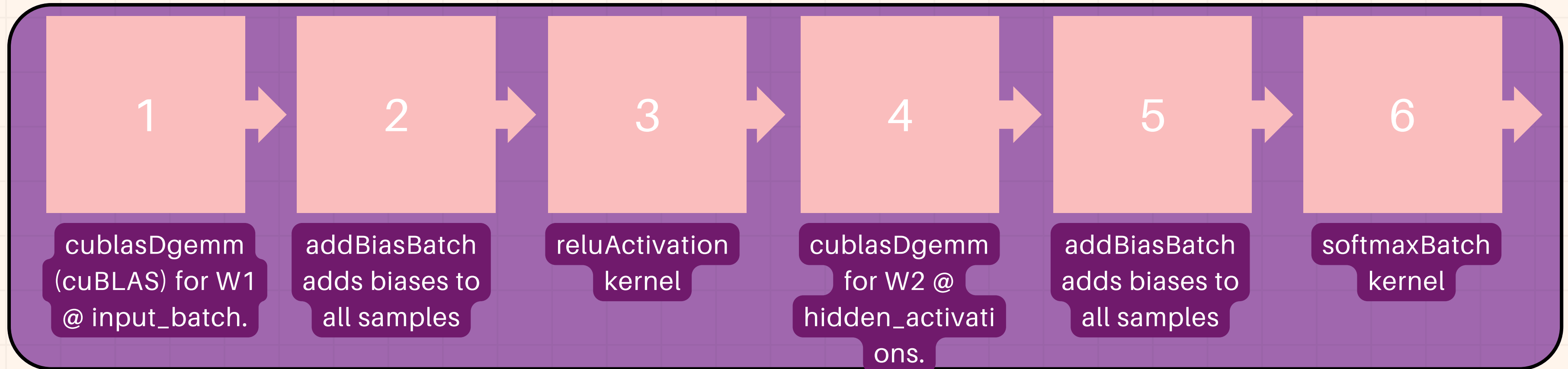
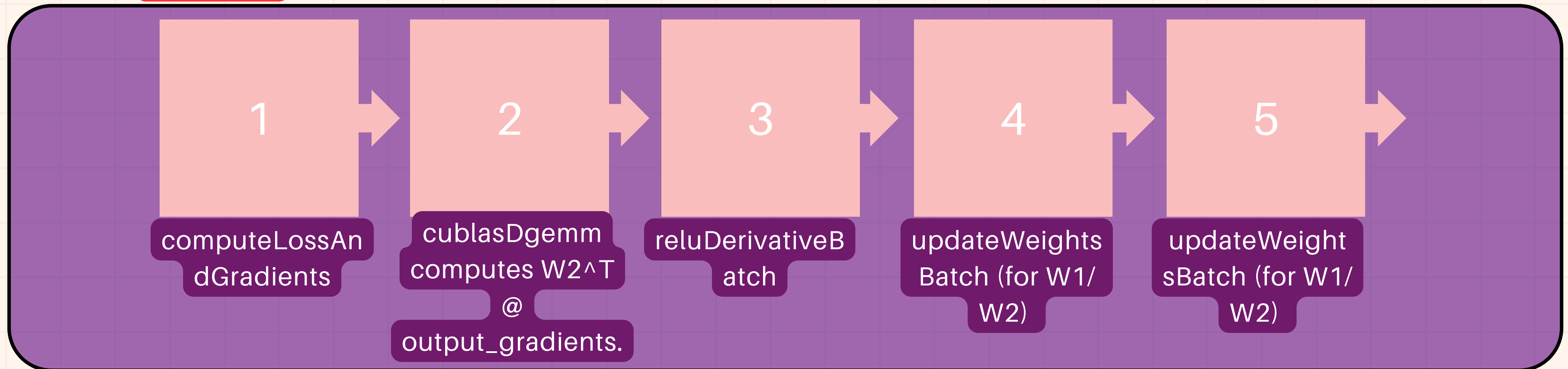| Time | Total Time | Instances | Avg | Med | Min | Max | StdDev | Category | Operation |
|---|---|---|---|---|---|---|---|---|---|
| 63.1% | 3.671 s | 190000 | 19.321 µs | 19.328 µs | 18.784 µs | 24.640 µs | 246 ns | CUDA_KERNEL | forwardHidden(__half *, __half *, __half *, __half *) |
| 10.5% | 613.064 ms | 380000 | 1.613 µs | 1.600 µs | 832 ns | 20.897 µs | 280 ns | MEMORY_OPER | [CUDA memcpy Device-to-Host] |
| 9.9% | 577.959 ms | 190000 | 3.041 µs | 3.040 µs | 3.008 µs | 3.552 µs | 12 ns | CUDA_KERNEL | forwardOutput(__half *, __half *, __half *, __half *) |
| 5.7% | 331.752 ms | 180000 | 1.843 µs | 1.856 µs | 1.760 µs | 1.952 µs | 20 ns | CUDA_KERNEL | updateHiddenLayer(__half *, __half *, __half *, __half *) |
| 4.8% | 277.718 ms | 180000 | 1.542 µs | 1.536 µs | 1.504 µs | 1.728 µs | 13 ns | CUDA_KERNEL | computeHiddenGradients(__half *, __half *, __half *, __half *) |
| 4.4% | 258.368 ms | 180000 | 1.435 µs | 1.440 µs | 1.408 µs | 1.504 µs | 13 ns | CUDA_KERNEL | updateOutputLayer(__half *, __half *, __half *, __half *) |

# ACCELERATION PLAN

- BATCH PROCESSING
- TENSOR CORES

# FORWARD PROPAGATION

**1** — cublasDgemm (cuBLAS) for W1 @ input_batch.

**2** — addBiasBatch adds biases to all samples

**3** — reluActivation kernel

**4** — cublasDgemm for W2 @ hidden_activations.

**5** — addBiasBatch adds biases to all samples

**6** — softmaxBatch kernel

*requires atomic operations

# BACKWARD PROPAGATION

**1** — computeLossAndGradients

**2** — cublasDgemm computes W2^T @ output_gradients.

**3** — reluDerivativeBatch
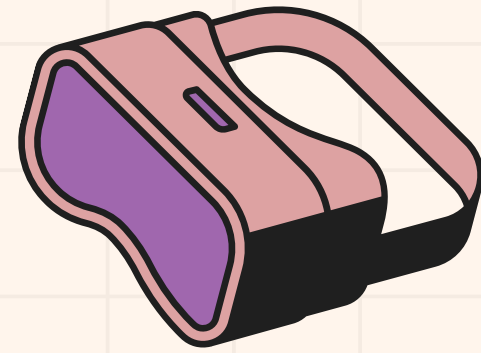
**4** — updateWeightsBatch (for W1/W2)

**5** — updateWeightsBatch (for W1/W2)

CONCLUSION

# THANK YOU