

MSRA SH Triton Study Group

1. Naive CUDA Softmax Kernel

2025/07/04

Problem Definition

- Shape: batch size $M \geq 4096$, hidden dim $N = 4096$
- Mathematical

$$\begin{aligned} X &\in \mathbb{R}^{M \times N} \\ Y_{ij} &= \text{softmax}(X_i)_j \\ &= \frac{\exp(X_{ij})}{\sum_{j=1}^N \exp(X_{ij})} \end{aligned}$$

Standard Softmax Operation

- Input is a 32-bit float tensor
- Find maximum value for each row
- Overflow-free exponential calculation
- Sum of exponentials for each row
- Output

$$X \in \mathbb{R}_{\text{float32}}^{M \times N}$$

$$M_i \leftarrow \max_{j=1}^N X_{ij}$$

$$Z_{ij} \leftarrow \exp(X_{ij} - M_i)$$

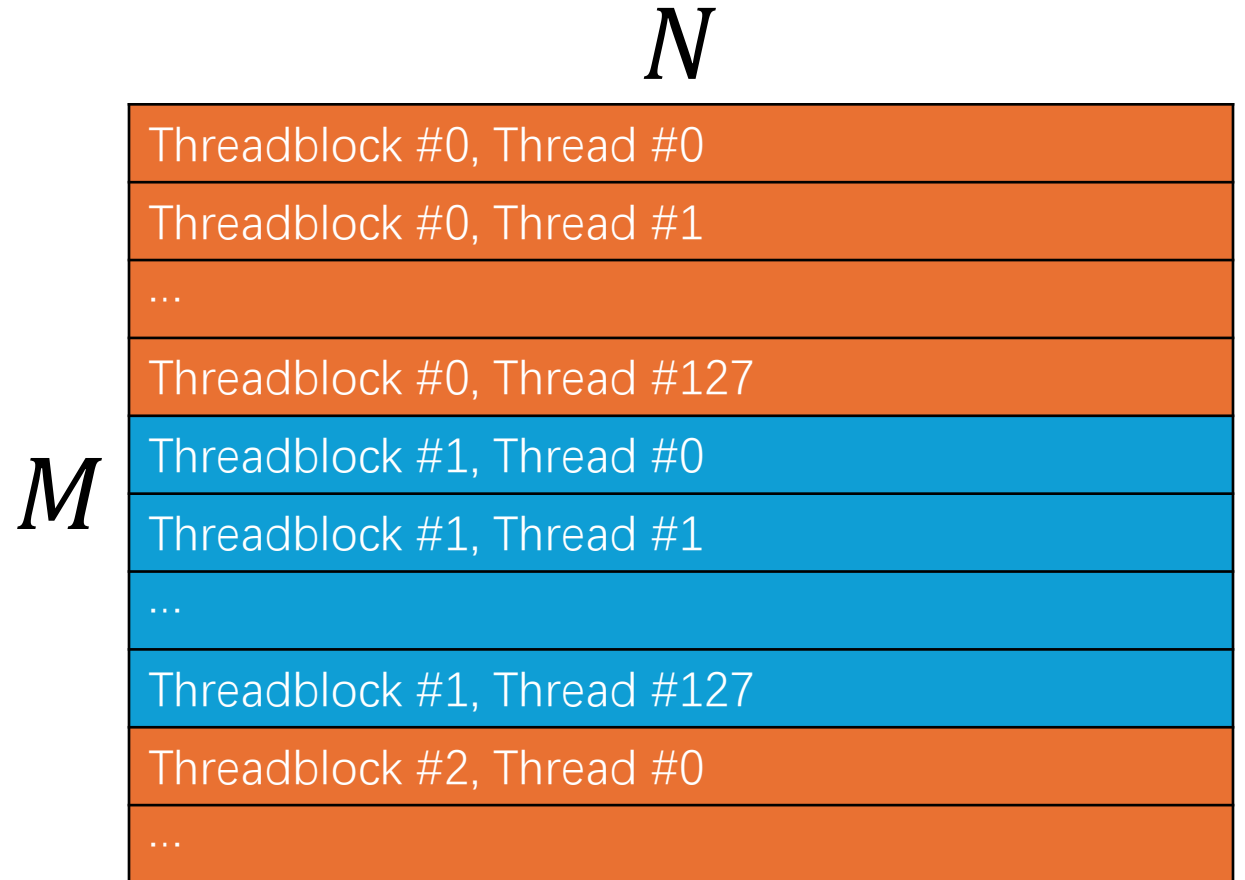
$$S_i \leftarrow \sum_{j=1}^N Z_{ij}$$

$$Y_{ij} \leftarrow \frac{Z_{ij}}{S_i}$$

- Row-wise data dependency

Naïve Softmax Kernel

- 1 thread per row
- For each row i :
 - Loop to find M_i
 - Loop to calc S_i
 - Loop to calc & store Y_{ij}



```

9 // CUDA kernel for naive softmax implementation
10 __global__ void naive_softmax_kernel(float* x, float* y, int batch_size, int hidden_dim) {
11     // Each thread processes one row of the input matrix x
12     int row_idx = blockIdx.x * blockDim.x + threadIdx.x;
13     // Boundary check
14     if (row_idx >= batch_size) return;
15
16     // Calculate the maximum value in the row
17     float max_val = -FLT_MAX;
18     for (int i = 0; i < hidden_dim; i++) {
19         float tmp_val = x[row_idx * hidden_dim + i]; // Read from global memory
20         max_val = max(max_val, tmp_val);
21     }
22
23     // Calculate the sum of exponentials
24     float sum_exp = 0.0f;
25     for (int i = 0; i < hidden_dim; i++) {
26         float tmp_val = x[row_idx * hidden_dim + i]; // Read from global memory
27         sum_exp += expf(tmp_val - max_val);
28     }
29
30     // Write the softmax output
31     for (int i = 0; i < hidden_dim; i++) {
32         float tmp_val = x[row_idx * hidden_dim + i]; // Read from global memory
33         y[row_idx * hidden_dim + i] = expf(tmp_val - max_val) / sum_exp; // Write to global memory
34     }
35 }

```

```
37
38 // C++ function to call the naive softmax kernel
39 torch::Tensor naive_softmax(torch::Tensor X) {
40     cudaSetDevice(X.get_device());
41
42     int batch_size = X.size(0);
43     int hidden_dim = X.size(1);
44     torch::Tensor Y = torch::empty_like(X, X.options());
45
46     // Thread block size
47     const int num_threads = 128;
48     // Grid size (= number of thread blocks)
49     int num_blocks = (batch_size + num_threads - 1) / num_threads;
50
51     // Launch the kernel
52     const dim3 dimBlock(num_threads);
53     const dim3 dimGrid(num_blocks);
54     naive_softmax_kernel<<<dimGrid, dimBlock>>>>(
55         X.data_ptr<float>(), // Pointer to input data
56         Y.data_ptr<float>(), // Pointer to output data
57         batch_size, hidden_dim
58     );
59
60     return Y;
61 }
62
```

Naïve Softmax Kernel

- Obviously slow

```
Batch size: 8765, Hidden dim: 4096  
Torch Softmax Latency: 0.20 ms  
Naive Softmax Latency: 2.16 ms
```

- Low GPU utility
 - Thread block number: $4096 / 128 = 32 < \text{physical SM number}$
- Redundant memory access
 - Register number: $255 < \text{hidden dim}$
 - Repeat load input values from global memory (L2 cache)
 - Non-vectorized memory access: 32 discrete address for 32 threads

Vectorized Memory Access

- GPU memory bus width: 128 bits = 4x float32 values

```
const int MEM_ACCESS_WIDTH = 4;
const int VALS_PER_THREAD = 128;
#pragma unroll
for (int i = 0; i < VALS_PER_THREAD; i += MEM_ACCESS_WIDTH) {
    reinterpret_cast<float4*>(&tmp_val[i])[0] =
        reinterpret_cast<float4*>(&x[row_idx * N + lane_idx * VALS_PER_THREAD + i])[0];
}
```

- Warp-level access consecutive 128 Byte global memory in a cycle

```
#pragma unroll
for (int i = 0; i < VALS_PER_THREAD; i += MEM_ACCESS_WIDTH) {
    reinterpret_cast<float4*>(&tmp_val[i])[0] =
        reinterpret_cast<float4*>(&x[row_idx * N + i * WARP_SIZE + lane_idx * MEM_ACCESS_WIDTH])[0];
}
```

Shuffle Sync

$t = \text{__shfl_xor_sync}(m, \text{laneMask}=1)$

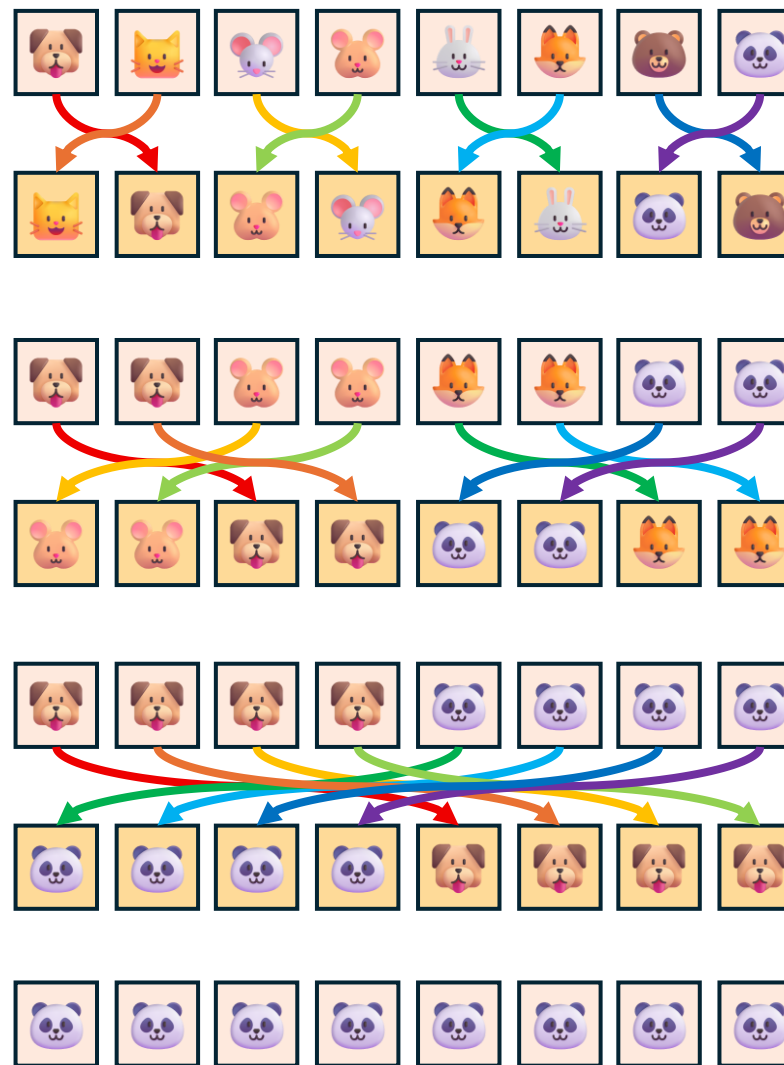
$m = \max(m, t)$

$t = \text{__shfl_xor_sync}(m, \text{laneMask}=2)$

$m = \max(m, t)$

$t = \text{__shfl_xor_sync}(m, \text{laneMask}=4)$

$m = \max(m, t)$



Better Softmax Kernel

- Reasonable fast

```
Batch size: 8765, Hidden dim: 4096  
Torch Softmax Latency: 0.20 ms  
Naive Softmax Latency: 2.16 ms  
Better Softmax Latency: 0.19 ms
```

- High GPU utility
 - Thread block number: $4096 / 4 = 1024 \gg$ physical SM number
- Efficient memory access
 - Register number: $255 \times 32 >$ hidden dim
 - Vectorized memory access
 - Contiguous global memory address in a warp

Homework

- Play with [TritonStudyGroup/1_CUDA_Softmax at main · Starmys/TritonStudyGroup](#)
- Questions:
 - Calculate the theoretical minimum latency for input size [8765, 4096] on single A100 GPU
 - Can we load the entire row into L1 cache or shared memory in naïve softmax?
 - How much acceleration does warp-level contiguous global memory access bring? Why?
 - How does the best softmax implementation changes for inputs of different shapes?

Reading Materials

- [How to Implement an Efficient Softmax CUDA Kernel— OneFlow Performance Optimization | by OneFlow | Medium](#)
- [Using CUDA Warp-Level Primitives | NVIDIA Technical Blog](#)