

# CUDA Programming

Rajat Ghosh

[rghosh08.github.io](https://github.com/rghosh08)

<https://github.com/rghosh08/cuda>

NUTANIX



# Roadmap

- Preliminary
- Getting Started
- Matrix Multiplication
- Alexnet
- Map Reduce
- Similarity Search
- Checksum
- LLM Inference Paged Attention

# **Compute Unified Device Architecture (CUDA)**

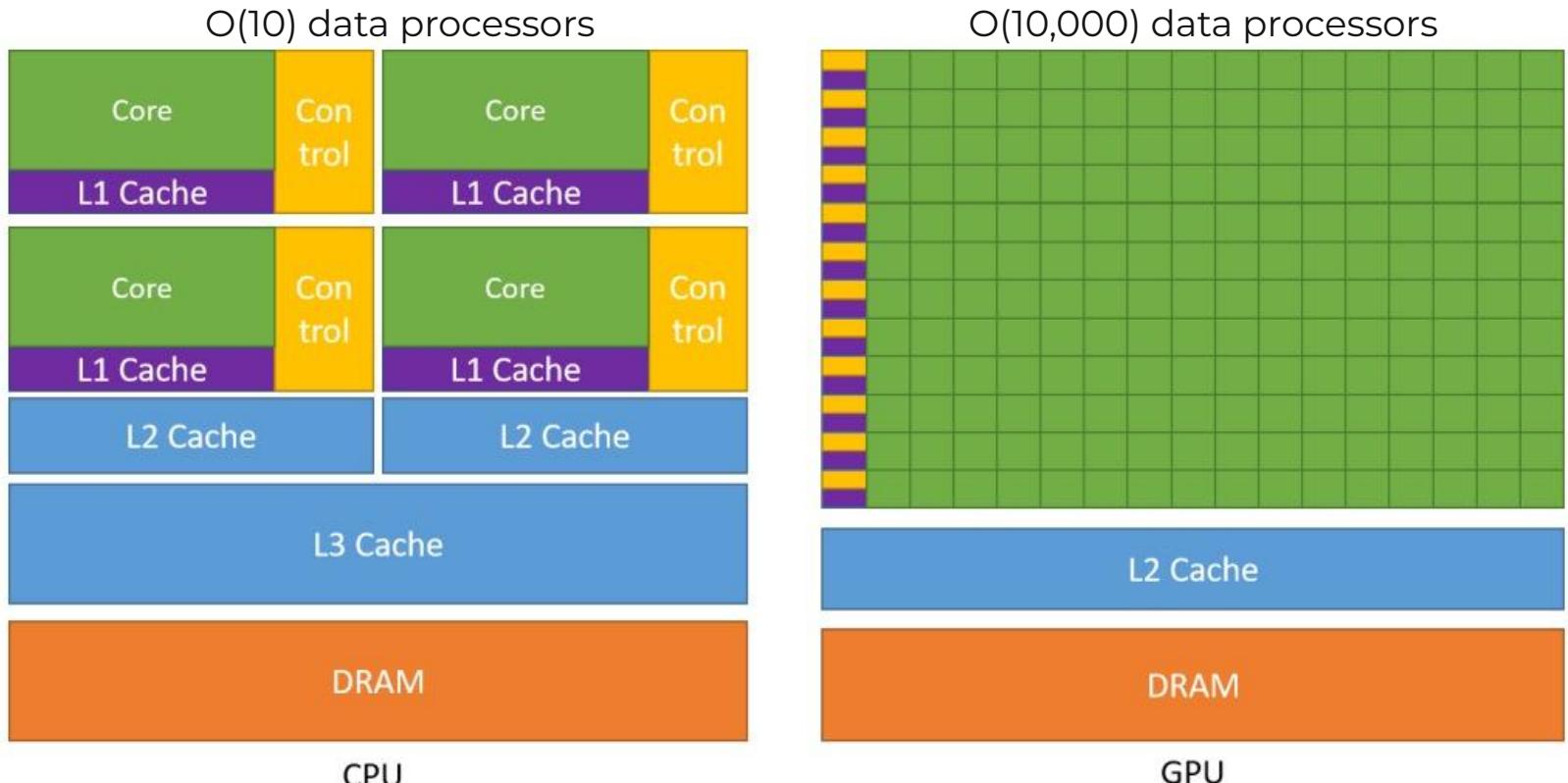
## **Before CUDA (Pre-2006)**

GPU for graphics programming with OpenGL/DirectX

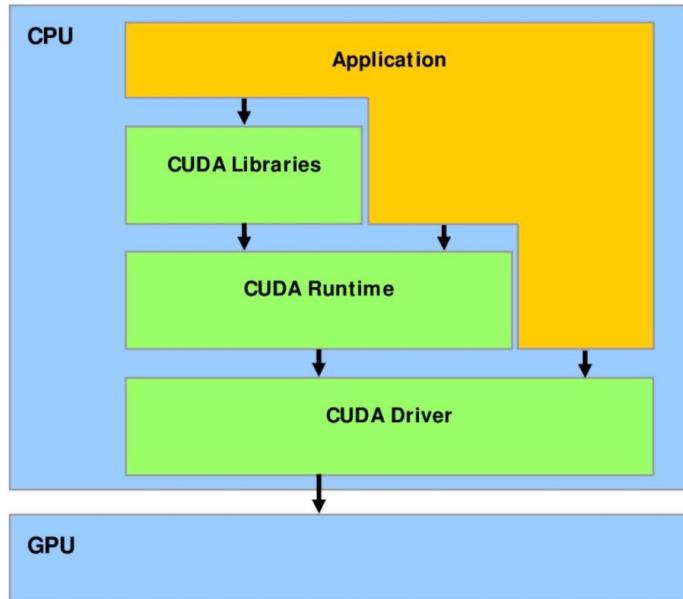
## **After CUDA**

GPU for general-purpose parallel computing with C/C++ API

# GPU vs GPU



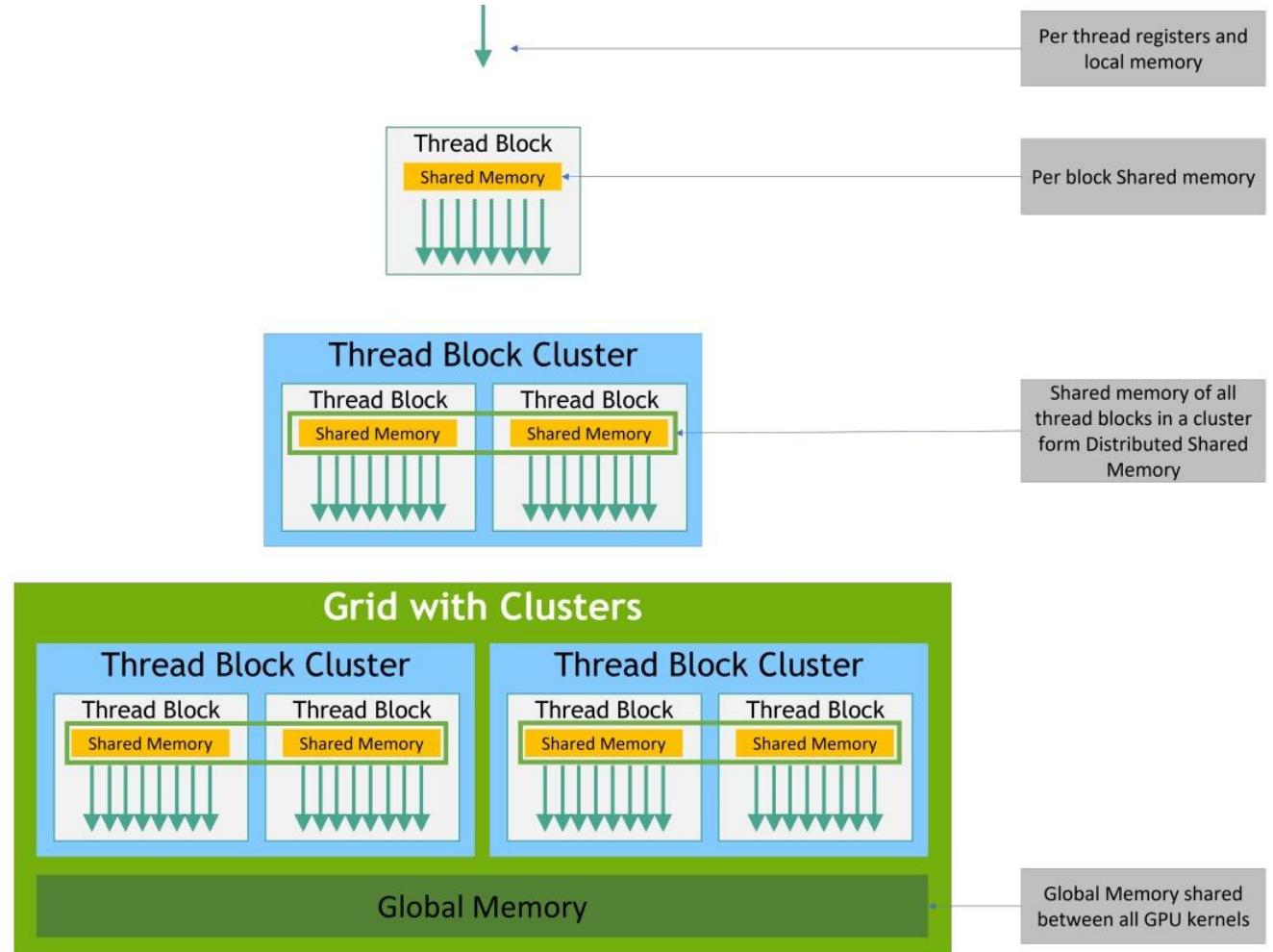
# Programming Model



```
// Kernel definition
__global__ void VecAdd(float* A, float* B, float* C)
{
    int i = threadIdx.x;
    C[i] = A[i] + B[i];
}

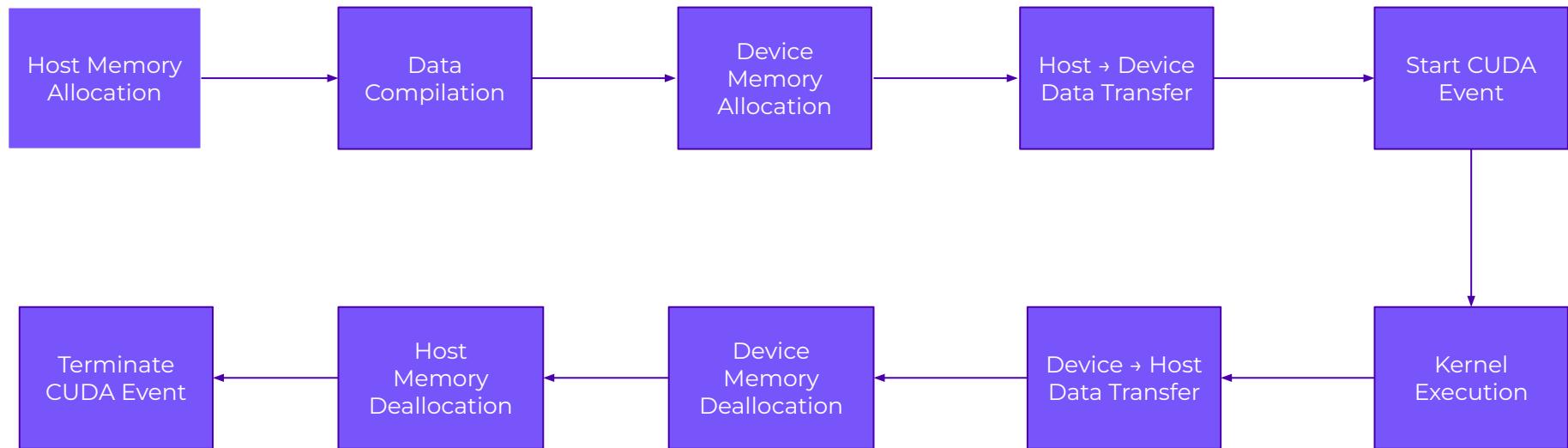
int main()
{
    ...
    // Kernel invocation with N threads
    VecAdd<<<1, N>>>(A, B, C);
    ...
}
```

# Hierarchy



# CUDA Operations Model

- Break a task into tiny operations
- Run each tiny operations on a single GPU thread
- Coalesce the results from single threads into a final answer
- Host CPU as an orchestrator



# Experimental Setup

# CUDA Device Info

NVIDIA-SMI 535.183.01			Driver Version: 535.183.01		CUDA Version: 12.2		
GPU	Name		Persistence-M	Bus-Id	Disp.A	Volatile	Uncorr. ECC
Fan	Temp	Perf	Pwr:Usage/Cap		Memory-Usage	GPU-Util	Compute M.
							MIG M.
0	NVIDIA A10G		On	00000000:00:1E.0	Off	0%	0
0%	34C	P8	24W / 300W		0MiB / 23028MiB	0%	Default
							N/A
Processes:							
GPU	GI	CI	PID	Type	Process name	GPU Memory Usage	
ID		ID					
No running processes found							

# CUDA Host Info

Specification	Value
Physical CPU Cores	8
Total Threads	16
Memory Size	64 GB
Disk Space	1.1 TB
OS	Ubuntu 20.04.6 LTS

# Getting Started

# Synchronization Issue

```
1 #include <stdio.h>
2
3 // Kernel function to generate a markdown-formatted table showing thread indices
4 __global__ void generateThreadTable() {
5     // Calculate global thread ID
6     int threadId = threadIdx.x + blockIdx.x * blockDim.x;
7
8     // Print table header only from the first thread to avoid duplicate headers
9     if (threadId == 0) {
10         printf("| blockIdx.x | threadIdx.x | blockDim.x | threadId |\n");
11         printf("|-----|-----|-----|-----|\n");
12     }
13
14     // Synchronize all threads to ensure the header prints before any thread data
15     __syncthreads();
16
17     // Each thread prints its own indices and global ID in table format
18     printf(" | %-10d | %-11d | %-10d | %-8d |\n",
19           blockIdx.x, threadIdx.x, blockDim.x, threadId);
20 }
21
22 int main() {
23     // Launch kernel with 2 blocks, each containing 4 threads
24     generateThreadTable<<<2, 4>>>();
25
26     // Wait for all GPU tasks to complete before ending the program
27     cudaDeviceSynchronize();
28
29     return 0;
30 }
```

1	0	4	4
blockIdx.x	threadIdx.x	blockDim.x	threadId
1	1	4	5
1	2	4	6
1	3	4	7
0	0	4	0
0	1	4	1
0	2	4	2
0	3	4	3

# Fixed by Changing Execution Order

```
1 #include <stdio.h>
2
3 // Kernel function to generate thread data (no header printing here)
4 __global__ void generateThreadTable() {
5     int threadIdx = threadIdx.x + blockIdx.x * blockDim.x;
6
7     printf("| %-10d | %-11d | %-10d | %-8d |\n",
8           blockIdx.x, threadIdx.x, blockDim.x, threadIdx);
9 }
10
11 int main() {
12     // Print table header from CPU (host)
13     printf("| blockIdx.x | threadIdx.x | blockDim.x | threadIdx |\n");
14     printf("|-----|-----|-----|-----|\n");
15
16     // Launch kernel
17     generateThreadTable<<<2, 4>>>();
18
19     // Wait for kernel to complete
20     cudaDeviceSynchronize();
21
22     return 0;
23 }
```

blockIdx.x	threadIdx.x	blockDim.x	threadId
1	0	4	4
1	1	4	5
1	2	4	6
1	3	4	7
0	0	4	0
0	1	4	1
0	2	4	2
0	3	4	3

# Matrix Multiplication

Task: Benchmark CUDA for Massive Scale Matrix Multiplication

```

1 #include <cuda_runtime.h>
2 #include <stdio.h>
3
4 #define N 65536
5
6 __global__ void matmul(float *A, float *B, float *C, int width) {
7     int row = blockIdx.y * blockDim.y + threadIdx.y;
8     int col = blockIdx.x * blockDim.x + threadIdx.x;
9
10    if (row < width && col < width) {
11        float val = 0;
12        for (int k = 0; k < width; k++) {
13            val += A[row * width + k] * B[k * width + col];
14        }
15        C[row * width + col] = val;
16    }
17 }
18
19 int main() {
20     size_t bytes = N * N * sizeof(float);
21
22     float *h_A = (float*)malloc(bytes);
23     float *h_B = (float*)malloc(bytes);
24     float *h_C = (float*)malloc(bytes);
25
26     for (int i = 0; i < N * N; i++) {
27         h_A[i] = 1.0f;
28         h_B[i] = 2.0f;
29     }
30
31     float *d_A, *d_B, *d_C;
32     cudaMalloc(&d_A, bytes);
33     cudaMalloc(&d_B, bytes);
34     cudaMalloc(&d_C, bytes);
35
36     cudaMemcpy(d_A, h_A, bytes, cudaMemcpyHostToDevice);
37     cudaMemcpy(d_B, h_B, bytes, cudaMemcpyHostToDevice);
38
39     dim3 threads(16, 16);
40     dim3 blocks((N + threads.x - 1) / threads.x, (N + threads.y - 1) / threads.y);
41
42     cudaEvent_t start, stop;
43     cudaEventCreate(&start);
44     cudaEventCreate(&stop);
45
46     cudaEventRecord(start);
47     matmul<<blocks, threads>>>(d_A, d_B, d_C, N);
48     cudaEventRecord(stop);
49     cudaEventSynchronize(stop);
50
51     float milliseconds = 0;
52     cudaEventElapsedTime(&milliseconds, start, stop);
53
54     cudaMemcpy(h_C, d_C, bytes, cudaMemcpyDeviceToHost);
55
56     // Compute FLOPS
57     double total_ops = 2.0 * N * N * N;
58     double seconds = milliseconds / 1000.0;
59     double flops = total_ops / seconds;
60
61     printf("Execution time: %.4f ms\n", milliseconds);
62     printf("Performance: %.4f GFLOPs\n", flops / 1e9);
63     printf("Sample output C[0]: %f\n", h_C[0]);
64
65     cudaFree(d_A);
66     cudaFree(d_B);
67     cudaFree(d_C);
68     free(h_A);
69     free(h_B);
70     free(h_C);
71
72     cudaEventDestroy(start);
73     cudaEventDestroy(stop);
74
75 }

```

ID	Description
1	Host Memory Allocation
2	Data Compilation
3	Device Memory Allocation
4	Host→Device Memory Copy
5	Execution Invocation
6	Kernel Execution
7	Device→Host Memory Copy
8	Device Memory Deallocation
9	Host Memory Deallocation

# Result

Matrix Size	CUDA	CPP	CPP-OpenMP	CUDA Speedup
1024x1024	1.1 ms	6.9 s	79 ms	~70x
2048x2048	8.3 ms	181.9 s	7.2 s	~800x
4096x4096	66 ms	2391.27 s	117 s	~1772x
8192x8192	527 ms	RIP	2,130 s	~4042x
16384x16384	4.3 s	RIP	RIP	CUDA survives 268M
32768x32768	36 s	RIP	RIP	CUDA survives 1B
65536x65536	OOM on 24GB VRAM	RIP	RIP	Memory Bound

# Deep Neural Network Training with Alexnet

Task: Replicate Alexnet Training on CUDA

# AlexNet with 60M Parameters (SOTA in 2012)

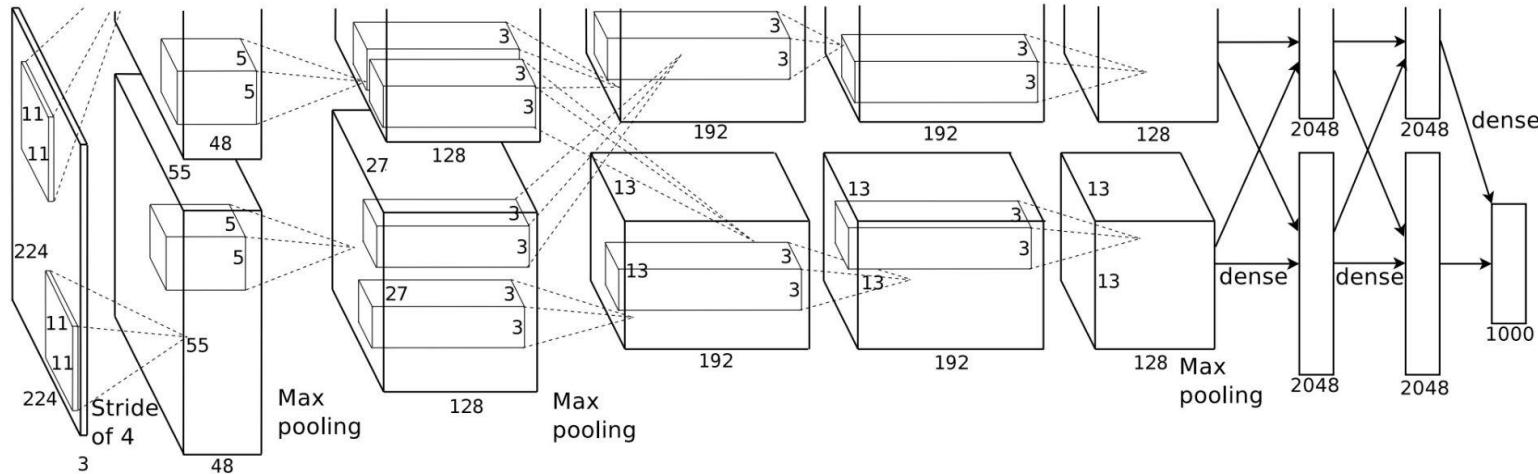


Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

# Architecture Implemented

Parameters	Original	Implemented
Input Size	227 x 227 x 3	227 x 227 x 3
No of Classes	1000	1000
Batch Size	128	512
Epoch	90	10
Training Samples	1,281,167	5120
Interactions per Epoch	10,009	10
Set Up	2x NVIDIA GTX 580	1x NVIDIA A10
Architecture	5 Conv + 3 FC	3 Conv + 3 FC

# Data Generator on Host

```
51 // Simple data loader
52 class SimpleDataLoader {
53 private:
54     std::vector<float> batch_input;
55     std::vector<int> batch_labels;
56     std::mt19937 rng;
57     int current_batch;
58     int total_batches;
59
60 public:
61     SimpleDataLoader() : rng(std::random_device{}()), current_batch(0) {
62         total_batches = Config::SAMPLES_PER_EPOCH / Config::BATCH_SIZE;
63
64         size_t input_size = Config::BATCH_SIZE * Config::INPUT_CHANNELS *
65                         Config::INPUT_HEIGHT * Config::INPUT_WIDTH;
66         batch_input.resize(input_size);
67         batch_labels.resize(Config::BATCH_SIZE);
68
69         std::cout << "DataLoader initialized: " << total_batches << " batches per epoch" << std::endl;
70     }
71
72     bool hasNextBatch() { return current_batch < total_batches; }
73     void resetEpoch() { current_batch = 0; }
74     int getCurrentBatch() const { return current_batch; }
75     int getTotalBatches() const { return total_batches; }
76
77     void generateBatch(float** d_input, int** d_labels) {
78         if (!hasNextBatch()) return;
79
80         // Generate random data
81         std::uniform_real_distribution<float> input_dist(-1.0f, 1.0f);
82         std::uniform_int_distribution<int> label_dist(0, Config::NUM_CLASSES - 1);
83
84         for (size_t i = 0; i < batch_input.size(); ++i) {
85             batch_input[i] = input_dist(rng);
86         }
87
88         for (size_t i = 0; i < batch_labels.size(); ++i) {
89             batch_labels[i] = label_dist(rng);
90         }
91
92         // Copy to GPU
93         CHECK_CUDA(cudaMemcpy(*d_input, batch_input.data(),
94                               batch_input.size() * sizeof(float), cudaMemcpyHostToDevice));
95         CHECK_CUDA(cudaMemcpy(*d_labels, batch_labels.data(),
96                               batch_labels.size() * sizeof(int), cudaMemcpyHostToDevice));
97
98         current_batch++;
99     }
100 };
```

# Convolution on CUDA

```
// CUDA kernels
__global__ void simple_conv2d_kernel(float* input, float* weight, float* bias, float* output,
    int batch_size, int in_channels, int out_channels,
    int input_h, int input_w, int kernel_size, int stride, int padding) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    int output_h = (input_h + 2 * padding - kernel_size) / stride + 1;
    int output_w = (input_w + 2 * padding - kernel_size) / stride + 1;
    int total_outputs = batch_size * out_channels * output_h * output_w;

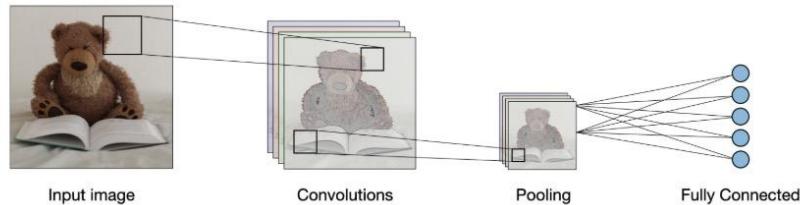
    if (idx < total_outputs) {
        int b = idx / (out_channels * output_h * output_w);
        int oc = (idx / (output_h * output_w)) % out_channels;
        int oh = (idx / output_w) % output_h;
        int ow = idx % output_w;

        float sum = bias[oc];

        // Simple convolution (unoptimized but functional)
        for (int ic = 0; ic < in_channels; ic++) {
            for (int kh = 0; kh < kernel_size; kh++) {
                for (int kw = 0; kw < kernel_size; kw++) {
                    int ih = oh * stride - padding + kh;
                    int iw = ow * stride - padding + kw;

                    if (ih >= 0 && ih < input_h && iw >= 0 && iw < input_w) {
                        int input_idx = b * (in_channels * input_h * input_w) +
                            ic * (input_h * input_w) + ih * input_w + iw;
                        int weight_idx = oc * (in_channels * kernel_size * kernel_size) +
                            ic * (kernel_size * kernel_size) + kh * kernel_size + kw;
                        sum += input[input_idx] * weight[weight_idx];
                    }
                }
            }
        }

        output[idx] = fmaxf(0.0f, sum); // ReLU activation
    }
}
```



$$Y[i,j] = \sum_{m} \sum_{n} X[i \times s + m - p, j \times s + n - p] \cdot W[m,n] + b$$

# Simple Max Pooling Kernel

```
141 __global__ void simple_maxpool_kernel(float* input, float* output,
142                                         int batch_size, int channels, int input_h, int input_w,
143                                         int kernel_size, int stride) {
144     int idx = blockIdx.x * blockDim.x + threadIdx.x;
145     int output_h = (input_h - kernel_size) / stride + 1;
146     int output_w = (input_w - kernel_size) / stride + 1;
147     int total_outputs = batch_size * channels * output_h * output_w;
148
149     if (idx < total_outputs) {
150         int b = idx / (channels * output_h * output_w);
151         int c = (idx / (output_h * output_w)) % channels;
152         int oh = (idx / output_w) % output_h;
153         int ow = idx % output_w;
154
155         float max_val = -INFINITY;
156
157         for (int kh = 0; kh < kernel_size; kh++) {
158             for (int kw = 0; kw < kernel_size; kw++) {
159                 int ih = oh * stride + kh;
160                 int iw = ow * stride + kw;
161                 int input_idx = b * (channels * input_h * input_w) +
162                               c * (input_h * input_w) + ih * input_w + iw;
163                 max_val = fmaxf(max_val, input[input_idx]);
164             }
165         }
166
167         output[idx] = max_val;
168     }
169 }
```

$$Y[i, j] = \max(X[i \times s + m, j \times s + n])$$

$m, n$

$$\text{ReLU}(x) = \max(0, x)$$

$$y = X \odot m$$

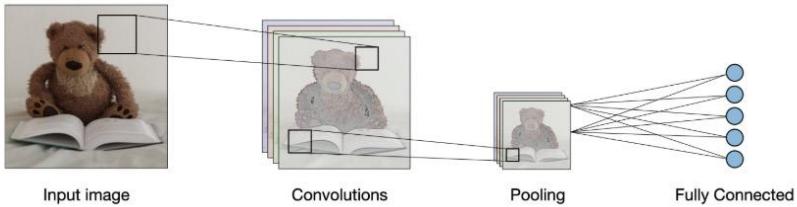
$$L = -\sum_{i=1}^n y_i \log(\hat{y}_i)$$

$$\theta_{t+1} = \theta_t - \eta/m \sum_{i=1}^m \nabla_{\theta} L(\theta_t; x^{(i)}, y^{(i)})$$

```

171 __global__ void relu_kernel(float* input, float* output, int size) {
172     int idx = blockIdx.x * blockDim.x + threadIdx.x;
173     if (idx < size) {
174         output[idx] = fmaxf(0.0f, input[idx]);
175     }
176 }
177
178 __global__ void dropout_kernel(float* input, float* output, int size, unsigned long long seed) {
179     int idx = blockIdx.x * blockDim.x + threadIdx.x;
180     if (idx < size) {
181         // Simple deterministic dropout for demo
182         if ((idx + (int)seed) % 2 == 0) {
183             output[idx] = input[idx] * 2.0f;
184         } else {
185             output[idx] = 0.0f;
186         }
187     }
188 }
189
190 __global__ void cross_entropy_loss_kernel(float* logits, int* labels, float* loss,
191                                            float* grad_output, int batch_size, int num_classes) {
192     int batch_idx = blockIdx.x * blockDim.x + threadIdx.x;
193
194     if (batch_idx < batch_size) {
195         int label = labels[batch_idx];
196         int offset = batch_idx * num_classes;
197
198         // Find max for numerical stability
199         float max_val = logits[offset];
200         for (int i = 1; i < num_classes; i++) {
201             max_val = fmaxf(max_val, logits[offset + i]);
202         }
203
204         // Compute softmax denominator
205         float sum_exp = 0.0f;
206         for (int i = 0; i < num_classes; i++) {
207             sum_exp += expf(logits[offset + i] - max_val);
208         }
209
210         // Loss
211         float log_prob = logits[offset + label] - max_val - logf(sum_exp);
212         atomicAdd(loss, -log_prob / batch_size);
213
214         // Gradient
215         for (int i = 0; i < num_classes; i++) {
216             float softmax_val = expf(logits[offset + i] - max_val) / sum_exp;
217             grad_output[offset + i] = (softmax_val - (i == label ? 1.0f : 0.0f)) / batch_size;
218         }
219     }
220 }
221
222 __global__ void sgd_update_kernel(float* weights, float* gradients, float learning_rate, int size) {
223     int idx = blockIdx.x * blockDim.x + threadIdx.x;
224     if (idx < size) {
225         weights[idx] -= learning_rate * gradients[idx];
226     }
227 }
```

# Forward Pass and Loss Function



```
float computeLoss() {
    float d_loss;
    CHECK_CUDA(cudaMalloc(&d_loss, sizeof(float)));
    CHECK_CUDA(cudaMemset(d_loss, 0, sizeof(float)));

    dim3 block(256);
    dim3 grid((Config::BATCH_SIZE + block.x - 1) / block.x);

    cross_entropy_loss_kernel<<<grid, block>>>(
        d_fc3_out, d_labels, d_loss, d_grad_fc3_out, Config::BATCH_SIZE, 1000);

    float host_loss;
    CHECK_CUDA(cudaMemcpy(&host_loss, d_loss, sizeof(float), cudaMemcpyDeviceToHost));
    CHECK_CUDA(cudaFree(d_loss));

    return host_loss;
}
```

```
public:
    void forward() {
        const float alpha = 1.0f, beta = 0.0f;

        // Conv1 + ReLU + Pool1 (simplified using custom kernels)
        dim3 block(256);
        int conv1_outputs = Config::BATCH_SIZE * 96 * 55 * 55;
        dim3 grid_conv1((conv1_outputs + block.x - 1) / block.x);

        simple_conv2d_kernel<<<grid_conv1, block>>>(
            d_input, d_conv1_w, d_conv1_b, d_conv1_out,
            Config::BATCH_SIZE, 3, 96, 227, 227, 11, 4, 2);

        // Pool1
        int pool1_outputs = Config::BATCH_SIZE * 96 * 27 * 27;
        dim3 grid_pool1((pool1_outputs + block.x - 1) / block.x);
        simple_maxpool_kernel<<<grid_pool1, block>>>(
            d_conv1_out, d_pool1_out, Config::BATCH_SIZE, 96, 55, 55, 3, 2);

        // FC1 (flatten pool1 and use cuBLAS)
        int fc1_input_size = 96 * 27 * 27;
        CHECK_CUBLAS(cublasSgemm(cublas, CUBLAS_OP_T, CUBLAS_OP_N,
            4096, Config::BATCH_SIZE, fc1_input_size,
            &alpha, d_fc1_w, fc1_input_size, d_pool1_out, fc1_input_size,
            &beta, d_fc1_out, 4096));

        // Add bias and apply ReLU + Dropout
        for (int b = 0; b < Config::BATCH_SIZE; ++b) {
            CHECK_CUBLAS(cublasSaxpy(cublas, 4096, &alpha, d_fc1_b, 1, d_fc1_out + b * 4096, 1));
        }

        dim3 grid_fc1((Config::BATCH_SIZE * 4096 + block.x - 1) / block.x);
        relu_kernel<<<grid_fc1, block>>>(d_fc1_out, d_fc1_out, Config::BATCH_SIZE * 4096);
        dropout_kernel<<<grid_fc1, block>>>(d_fc1_out, d_fc1_out, Config::BATCH_SIZE * 4096, dropout_seed++);

        // FC2
        CHECK_CUBLAS(cublasSgemm(cublas, CUBLAS_OP_T, CUBLAS_OP_N,
            4096, Config::BATCH_SIZE, 4096,
            &alpha, d_fc2_w, 4096, d_fc1_out, 4096,
            &beta, d_fc2_out, 4096));

        for (int b = 0; b < Config::BATCH_SIZE; ++b) {
            CHECK_CUBLAS(cublasSaxpy(cublas, 4096, &alpha, d_fc2_b, 1, d_fc2_out + b * 4096, 1));
        }

        relu_kernel<<<grid_fc1, block>>>(d_fc2_out, d_fc2_out, Config::BATCH_SIZE * 4096);
        dropout_kernel<<<grid_fc1, block>>>(d_fc2_out, d_fc2_out, Config::BATCH_SIZE * 4096, dropout_seed++);

        // FC3
        CHECK_CUBLAS(cublasSgemm(cublas, CUBLAS_OP_T, CUBLAS_OP_N,
            1000, Config::BATCH_SIZE, 4096,
            &alpha, d_fc3_w, 4096, d_fc2_out, 4096,
            &beta, d_fc3_out, 1000));

        for (int b = 0; b < Config::BATCH_SIZE; ++b) {
            CHECK_CUBLAS(cublasSaxpy(cublas, 1000, &alpha, d_fc3_b, 1, d_fc3_out + b * 1000, 1));
        }
    }
```

# Backward Pass

```
void backward() {
    // Simplified backward pass - only update FC3 weights
    const float alpha = 1.0f, beta = 0.0f;

    // FC3 weight gradients
    CHECK_CUBLAS(cublasSgemm(cublas, CUBLAS_OP_N, CUBLAS_OP_T,
        4096, 1000, Config::BATCH_SIZE,
        &alpha, d_fc2_out, 4096, d_grad_fc3_out, 1000,
        &beta, d_grad_fc3_w, 4096));

    // FC3 bias gradients
    CHECK_CUBLAS(cublasSgemv(cublas, CUBLAS_OP_T,
        Config::BATCH_SIZE, 1000,
        &alpha, d_grad_fc3_out, Config::BATCH_SIZE,
        d_ones, 1, &beta, d_grad_fc3_b, 1));
}

void updateWeights() {
    dim3 block(256);

    // Update FC3 weights
    dim3 grid_w((1000 * 4096 + block.x - 1) / block.x);
    sgd_update_kernel<<<grid_w, block>>>(d_fc3_w, d_grad_fc3_w, learning_rate, 1000 * 4096);
    dim3 grid_b((1000 + block.x - 1) / block.x);
    sgd_update_kernel<<<grid_b, block>>>(d_fc3_b, d_grad_fc3_b, learning_rate, 1000);
}
```

# Cuda vs CPP

X

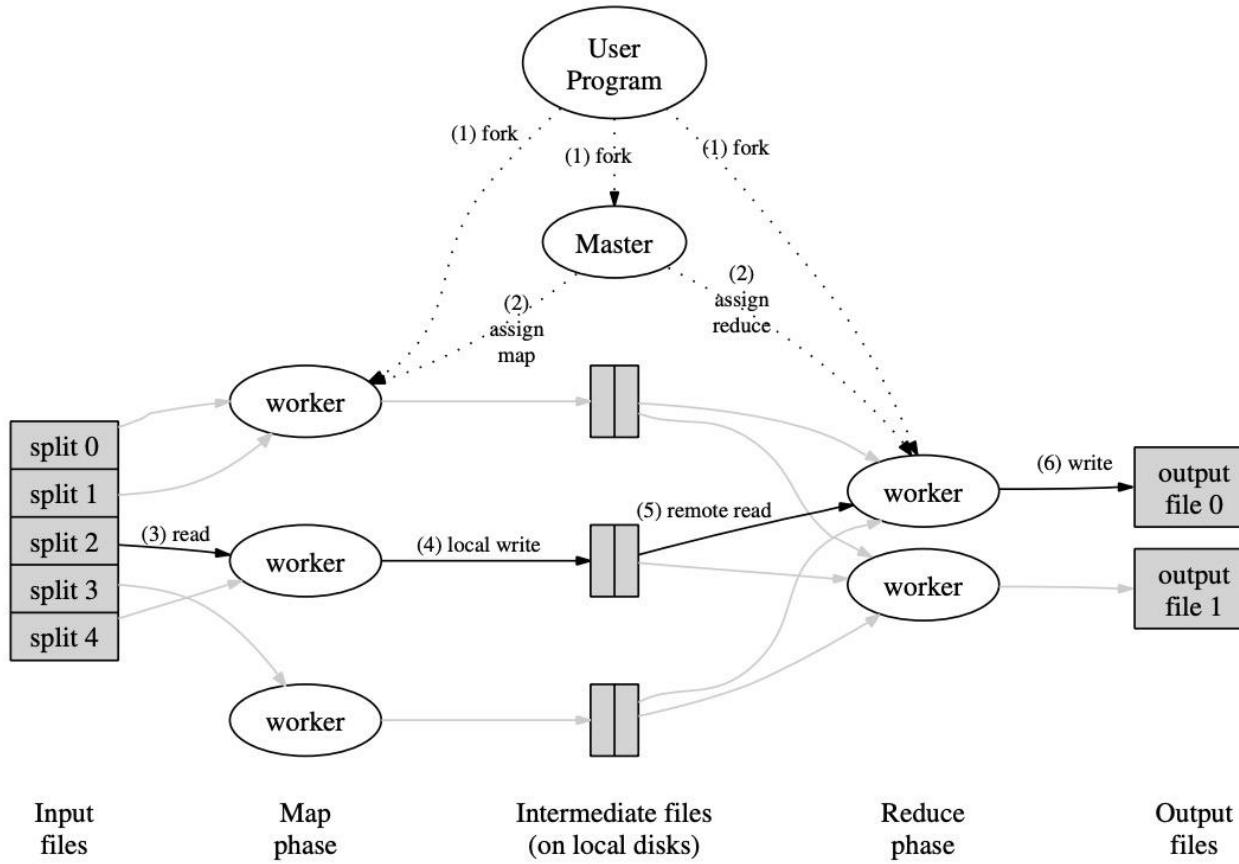
CUDA	CPP-OpenMP
<p>==== Training Summary ==== Total Training Time: 5 minutes Average Time per Epoch: 0.5 minutes ==== Loss Progression ==== Epoch 1: 9.5615 (33.0000s) Epoch 2: 8.1356 (33.0000s) Epoch 3: 7.8586 (33.0000s) Epoch 4: 7.8395 (33.0000s) Epoch 5: 7.7629 (33.0000s) Epoch 6: 7.7826 (33.0000s) Epoch 7: 7.7677 (33.0000s) Epoch 8: 7.7502 (33.0000s) Epoch 9: 7.7596 (33.0000s) Epoch 10: 7.7541 (33.0000s) ==== Performance Metrics ==== Initial Loss: 9.5615 Final Loss: 7.7541 Peak GPU Memory: 2486 MB</p> <p>==== Training Summary ==== Total Training Time: 20 minutes Average Time per Epoch: 2.0 minutes ==== Loss Progression ==== Epoch 1: 23.0079 (125.0000s) Epoch 2: 9.5457 (118.0000s) Epoch 3: 6.9092 (125.0000s) Epoch 4: 6.9085 (125.0000s) Epoch 5: 6.9091 (125.0000s) Epoch 6: 6.9095 (118.0000s) Epoch 7: 6.9080 (125.0000s) Epoch 8: 6.9094 (127.0000s) Epoch 9: 6.9234 (125.0000s) Epoch 10: 6.9079 (126.0000s) ==== Performance Metrics ==== Initial Loss: 23.0079 Final Loss: 6.9079 OpenMP Threads Used: 16 CPU Architecture: Multi-core parallel training ==== Performance Comparison ==== CPU Training (this run): 20 minutes Expected GPU Training: 4-6 minutes (A10 GPU) CPU vs GPU Speedup: ~4x faster on GPU ==== OpenMP Optimization Analysis ==== • Convolution parallelized across output elements</p>	

# Energy-Aware Training

<b>Epoch</b>	<b>Temperature (°C)</b>
10	46
30	48
100	52
300	66
1000	86

# Map Reduce

Task: Benchmark massive scale term frequency computation on CUDA



<https://static.googleusercontent.com/media/research.google.com/en//archive/mapreduce-osdi04.pdf>

# Kernel Execution for Term Frequency Computation

```
1 // Custom MapReduce kernels
2 __global__ void map_kernel(int* input, KeyValue* mapped, int n) {
3     int idx = blockIdx.x * blockDim.x + threadIdx.x;
4     if (idx < n) {
5         mapped[idx].key = input[idx];
6         mapped[idx].value = 1;
7     }
8 }
9
10 __global__ void reduce_kernel(KeyValue* input, KeyValue* output, int* counts, int n) {
11     int idx = blockIdx.x * blockDim.x + threadIdx.x;
12
13     if (idx < n) {
14         int key = input[idx].key;
15         atomicAdd(&counts[key], 1);
16     }
17 }
18
19 __global__ void compact_results(int* counts, KeyValue* output, int max_keys, int* result_count) {
20     int idx = blockIdx.x * blockDim.x + threadIdx.x;
21
22     if (idx < max_keys && counts[idx] > 0) {
23         int pos = atomicAdd(result_count, 1);
24         output[pos].key = idx;
25         output[pos].value = counts[idx];
26     }
27 }
28 }
```

Generate key value pairs {'integer:1'}

Generate value sum for each integer

Skipping zero frequency integers

# Host Orchestration in Term Frequency Computation

```
void custom_mapreduce(int* data, int n, KeyValue* results, int* num_results) {
    KeyValue *d_mapped, *d_results;
    int *d_data, *d_counts, *d_result_count;

    // Allocate memory
    CUDA_CHECK(cudaMalloc(&d_mapped, n * sizeof(KeyValue)));
    CUDA_CHECK(cudaMalloc(&d_results, MAX_WORDS * sizeof(KeyValue)));
    CUDA_CHECK(cudaMalloc(&d_data, n * sizeof(int)));
    CUDA_CHECK(cudaMalloc(&d_counts, MAX_WORDS * sizeof(int)));
    CUDA_CHECK(cudaMalloc(&d_result_count, sizeof(int)));

    // Copy data to device
    CUDA_CHECK(cudaMemcpy(d_data, data, n * sizeof(int), cudaMemcpyHostToDevice));
    CUDA_CHECK(cudaMemset(d_counts, 0, MAX_WORDS * sizeof(int)));
    CUDA_CHECK(cudaMemset(d_result_count, 0, sizeof(int)));

    // Map phase
    int grid_size = (n + BLOCK_SIZE - 1) / BLOCK_SIZE;
    map_kernel<<<grid_size, BLOCK_SIZE>>>(d_data, d_mapped, n);
    CUDA_CHECK(cudaDeviceSynchronize());

    // Reduce phase
    reduce_kernel<<<grid_size, BLOCK_SIZE>>>(d_mapped, d_results, d_counts, n);
    CUDA_CHECK(cudaDeviceSynchronize());

    // Compact results
    int count_grid = (MAX_WORDS + BLOCK_SIZE - 1) / BLOCK_SIZE;
    compact_results<<<count_grid, BLOCK_SIZE>>>(d_counts, d_results, MAX_WORDS, d_result_count);
    CUDA_CHECK(cudaDeviceSynchronize());

    // Copy results back
    CUDA_CHECK(cudaMemcpy(num_results, d_result_count, sizeof(int), cudaMemcpyDeviceToHost));
    CUDA_CHECK(cudaMemcpy(results, d_results, (*num_results) * sizeof(KeyValue), cudaMemcpyDeviceToHost));

    // Cleanup
    cudaFree(d_mapped);
    cudaFree(d_results);
    cudaFree(d_data);
    cudaFree(d_counts);
    cudaFree(d_result_count);
}
```



# Host Orchestration with Thrust

```
void thrust_mapreduce(int* data, int n, KeyValue* results, int* num_results) {
    thrust::device_vector<int> d_words(data, data + n);
    thrust::device_vector<int> d_counts(d_words.size(), 1);

    // Sort by key for grouping
    thrust::sort_by_key(d_words.begin(), d_words.end(), d_counts.begin());

    // Reduce by key
    thrust::device_vector<int> unique_words(n);
    thrust::device_vector<int> word_counts(n);

    auto end = thrust::reduce_by_key(d_words.begin(), d_words.end(), d_counts.begin(),
                                    unique_words.begin(), word_counts.begin());

    *num_results = end.first - unique_words.begin();

    // Copy results
    for (int i = 0; i < *num_results; i++) {
        results[i].key = unique_words[i];
        results[i].value = word_counts[i];
    }
}
```

# Result

X

Size	Time	Memory Usage
1M	149 ms	0.0037 GB   21.98 GB
10M	159 ms	0.037 GB   21.98 GB
100M	250 ms	0.372 GB   21.98 GB
1B	1.2 s	3.72 GB   21.98 GB
1.2B	1.5 s	4.47 GB   21.98 GB
1.5B	OOM	5.58 GB   21.98 GB
2B	OOM	7.45 GB   21.98 GB
3B	OOM	11.2 GB   21.98 GB
5B	OOM	18.2 GB   21.98 GB
10 B	Failure	37.2 GB   21.98 GB

# Massive Scale Similarity Search

Task: Benchmark massive scale distance computation on CUDA

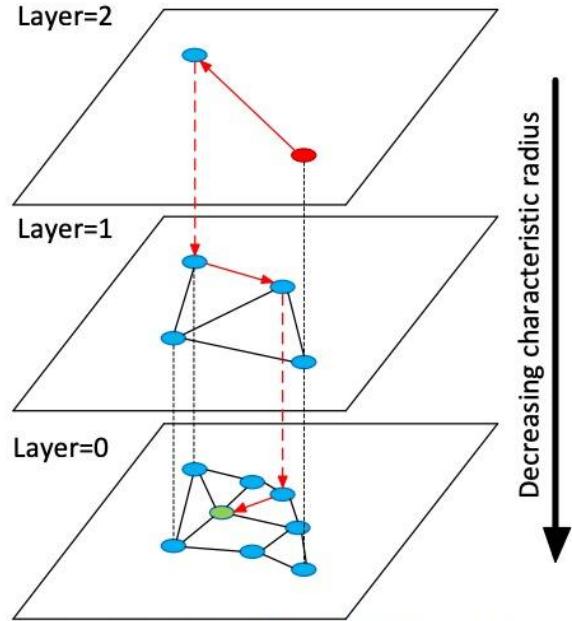


Fig. 1. Illustration of the Hierarchical NSW idea. The search starts from an element from the top layer (shown red). Red arrows show direction of the greedy algorithm from the entry point to the query (shown green).

$$L = k\text{-argmin}_{i=0:\ell} \|x - y_i\|_2,$$

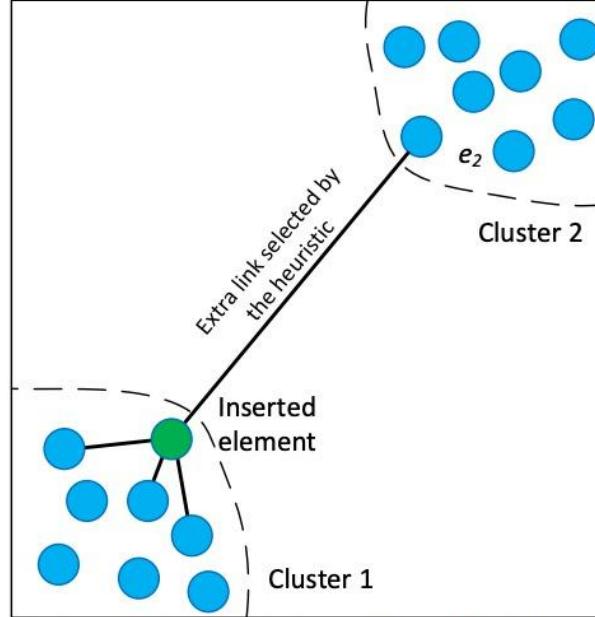


Fig. 2. Illustration of the heuristic used to select the graph neighbors for two isolated clusters. A new element is inserted on the boundary of Cluster 1. All of the closest neighbors of the element belong to the Cluster 1, thus missing the edges of Delaunay graph between the clusters. The heuristic, however, selects element  $e_2$  from Cluster 2, thus, maintaining the global connectivity in case the inserted element is the closest to  $e_2$  compared to any other element from Cluster 1.

# L2 Norm Computation on GPU

```
__global__ void compute_distances(float *nodes, float *query, float *distances, int dim, size_t num_nodes) {
    size_t idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < num_nodes) {
        float distance = 0.0f;
        for (int d = 0; d < dim; d++) {
            float diff = nodes[idx * dim + d] - query[d];
            distance += diff * diff;
        }
        distances[idx] = sqrtf(distance);
    }
}
```

# Host Orchestration

```
int main(int argc, char* argv[]) {
    if (argc != 3) {
        fprintf(stderr, "Usage: %s <num_nodes> <dim>\n", argv[0]);
        return EXIT_FAILURE;
    }

    size_t num_nodes = atol(argv[1]);
    int dim = atoi(argv[2]);

    printf("Allocating host memory...\n");
    float *nodes = (float*)malloc(num_nodes * dim * sizeof(float));
    float *query = (float*)malloc(dim * sizeof(float));
    float *distances = (float*)malloc(num_nodes * sizeof(float));

    printf("Initializing data...\n");
    for (size_t i = 0; i < num_nodes * dim; i++) nodes[i] = rand() / (float)RAND_MAX;
    for (int i = 0; i < dim; i++) query[i] = rand() / (float)RAND_MAX;

    float *d_nodes, *d_query, *d_distances;

    printf("Allocating GPU memory...\n");
    cudaMalloc(&d_nodes, num_nodes * dim * sizeof(float));
    cudaMalloc(&d_query, dim * sizeof(float));
    cudaMalloc(&d_distances, num_nodes * sizeof(float));

    printf("Transferring data to GPU...\n");
    cudaMemcpy(d_nodes, nodes, num_nodes * dim * sizeof(float), cudaMemcpyHostToDevice);
    cudaMemcpy(d_query, query, dim * sizeof(float), cudaMemcpyHostToDevice);

    printGPUMemoryUsage("After allocations & transfers");

    int threadsPerBlock = 256;
    int blocksPerGrid = (num_nodes + threadsPerBlock - 1) / threadsPerBlock;

    printf("Starting CUDA kernel execution...\n");
    auto start = std::chrono::high_resolution_clock::now();

    compute_distances<<<blocksPerGrid, threadsPerBlock>>>(d_nodes, d_query, d_distances, dim, num_nodes);
    cudaDeviceSynchronize();

    auto end = std::chrono::high_resolution_clock::now();
    std::chrono::duration<double>, std::milli> elapsed = end - start;

    printf("CUDA kernel execution completed.\n");
    printf("Kernel Execution Time: %.3f ms\n", elapsed.count());
```

<b>Size</b>	<b>CUDA</b>	<b>CUDA Memory</b>	<b>CPP OpenMP</b>
1M	3 ms	0.73 GB	15 ms
10M	24 ms	5 GB	140 ms
20 M	49 ms	9.8 GB	279 ms
30 M	73 ms	14.5 GB	418 ms
40M	97 ms	19 GB	555 ms
45M	109 ms	21.8 GB	624 ms

# Checksum Computation

Task: Benchmark massive scale checksum computation for 16MB Data on CUDA

## Mathematical Definition

**Checksum** =  $\sum_{i=0 \text{ to } n-1} \text{data}[i] = \text{data}[0] + \text{data}[1] + \dots + \text{data}[n-1]$

## Parallel Reduction Algorithm

### Block-Level Reduction

For a block with B threads and shared memory array S:

#### Algorithm: Parallel Reduction within Block

1. **Require:** Block size B =  $2^k$  for some integer k
2. **Require:** Shared memory array S[0..B-1]
3. **Require:** Thread ID tid  $\in [0, B-1]$
4. Initialize:  $S[\text{tid}] \leftarrow \text{data}[\text{globalIdx}]$
5. **syncthreads()**
6. **for** stride = B/2; stride > 0; stride = stride/2 **do**
7.   **if** tid < stride **then**
8.      $S[\text{tid}] \leftarrow S[\text{tid}] + S[\text{tid} + \text{stride}]$
9.   **end if**
10. **syncthreads()**
11. **end for**
12. **if** tid = 0 **then**
13.   **atomicAdd**(globalChecksum, S[0])
14. **end if**

### Tree-Based Reduction Pattern

The reduction follows a binary tree pattern:

**Level 0:**  $S[i] = \text{data}[i] \quad \forall i \in [0, B-1]$

**Level 1:**  $S[i] = S[i] + S[i + B/2] \quad \forall i \in [0, B/2-1]$

**Level 2:**  $S[i] = S[i] + S[i + B/4] \quad \forall i \in [0, B/4-1]$

:

**Level k:**  $S[0] = S[0] + S[1]$

### Multi-Block Global Reduction

For G blocks, each contributing partial sum  $P_j$ :

$P_j = \sum_{i=j \cdot B \text{ to } (j+1) \cdot B - 1} \text{data}[i] \quad \text{for block } j$

**Checksum** =  $\sum_{j=0 \text{ to } G-1} P_j = \sum_{j=0 \text{ to } G-1} \sum_{i=j \cdot B \text{ to } (j+1) \cdot B - 1} \text{data}[i]$

### Complexity Analysis

#### Time Complexity

$T_{\text{block}} = O(\log_2 B)$  per block

$T_{\text{total}} = O(\log_2 B)$  (parallel across blocks)

```
// CUDA Kernel: Compute simple checksum
__global__ void computeChecksum(unsigned char* data, size_t size, unsigned int* checksum) {
    __shared__ unsigned int shared_sum[BLOCK_SIZE];

    int idx = threadIdx.x + blockIdx.x * blockDim.x;
    int tid = threadIdx.x;

    // Initialize shared memory
    shared_sum[tid] = (idx < size) ? data[idx] : 0;
    __syncthreads();

    // Parallel reduction within block
    for (int stride = BLOCK_SIZE / 2; stride > 0; stride /= 2) {
        if (tid < stride)
            shared_sum[tid] += shared_sum[tid + stride];
        __syncthreads();
    }

    // Atomic add block result to global checksum
    if (tid == 0)
        atomicAdd(checksum, shared_sum[0]);
}
```

```

int main() {
    const size_t dataSize = 1 << 24; // 16 MB
    unsigned char* h_data = (unsigned char*)malloc(dataSize);

    // Initialize data (simulate read from storage)
    for (size_t i = 0; i < dataSize; ++i)
        h_data[i] = rand() % 256;

    unsigned char* d_data;
    unsigned int* d_checksum;
    unsigned int h_checksum = 0;

    cudaMalloc(&d_data, dataSize);
    cudaMalloc(&d_checksum, sizeof(unsigned int));
    cudaMemcpy(d_data, h_data, dataSize, cudaMemcpyHostToDevice);
    cudaMemset(d_checksum, 0, sizeof(unsigned int));

    dim3 block(BLOCK_SIZE);
    dim3 grid((dataSize + BLOCK_SIZE - 1) / BLOCK_SIZE);

    cudaEvent_t start, stop;
    cudaEventCreate(&start);
    cudaEventCreate(&stop);

    cudaEventRecord(start);
    computeChecksum<<<grid, block>>>(d_data, dataSize, d_checksum);
    cudaEventRecord(stop);
    cudaEventSynchronize(stop);

    float milliseconds;
    cudaEventElapsedTime(&milliseconds, start, stop);

    cudaMemcpy(&h_checksum, d_checksum, sizeof(unsigned int), cudaMemcpyDeviceToHost);

    printf("Checksum: %u\n", h_checksum);
    printf("Processing Time: %.4f ms\n", milliseconds);

    // Cleanup
    cudaFree(d_data);
    cudaFree(d_checksum);
    free(h_data);
    cudaEventDestroy(start);
    cudaEventDestroy(stop);

    return 0;
}

```

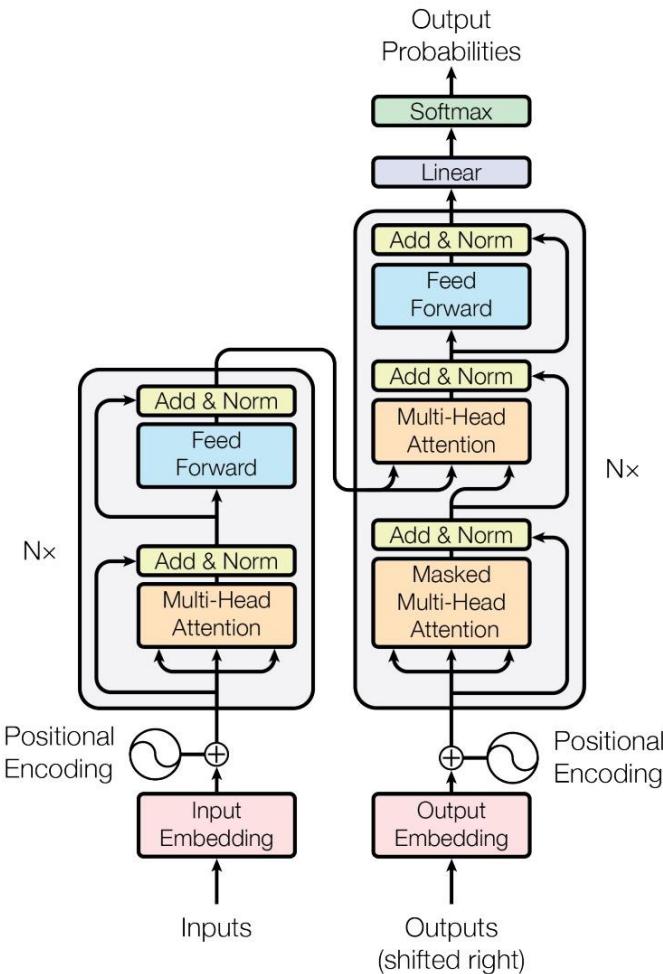
// Assuming BLOCK\_SIZE = 256:  
 dataSize = 16,777,216 bytes  
 BLOCK\_SIZE = 256 threads per block  
 grid\_size = (16,777,216 + 256 - 1) / 256 =  
 16,777,471 / 256 = 65,536 blocks



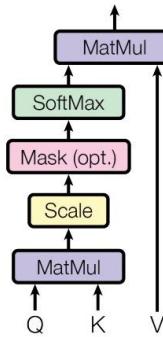
<b>Block Size</b>	<b>CUDA</b>	<b>CPP</b>	<b>CUDA Speed up</b>
256	0.3512 ms	13.4741 ms	38x
512	0.3891 ms	14.2124 ms	36x
1024	0.5335 ms	16.5069 ms	30x

# LLM Inference with Paged Attention

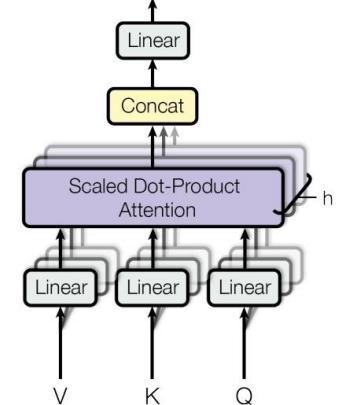
Task: Implement Paged Attention from Scratch on CUDA



### Scaled Dot-Product Attention



### Multi-Head Attention



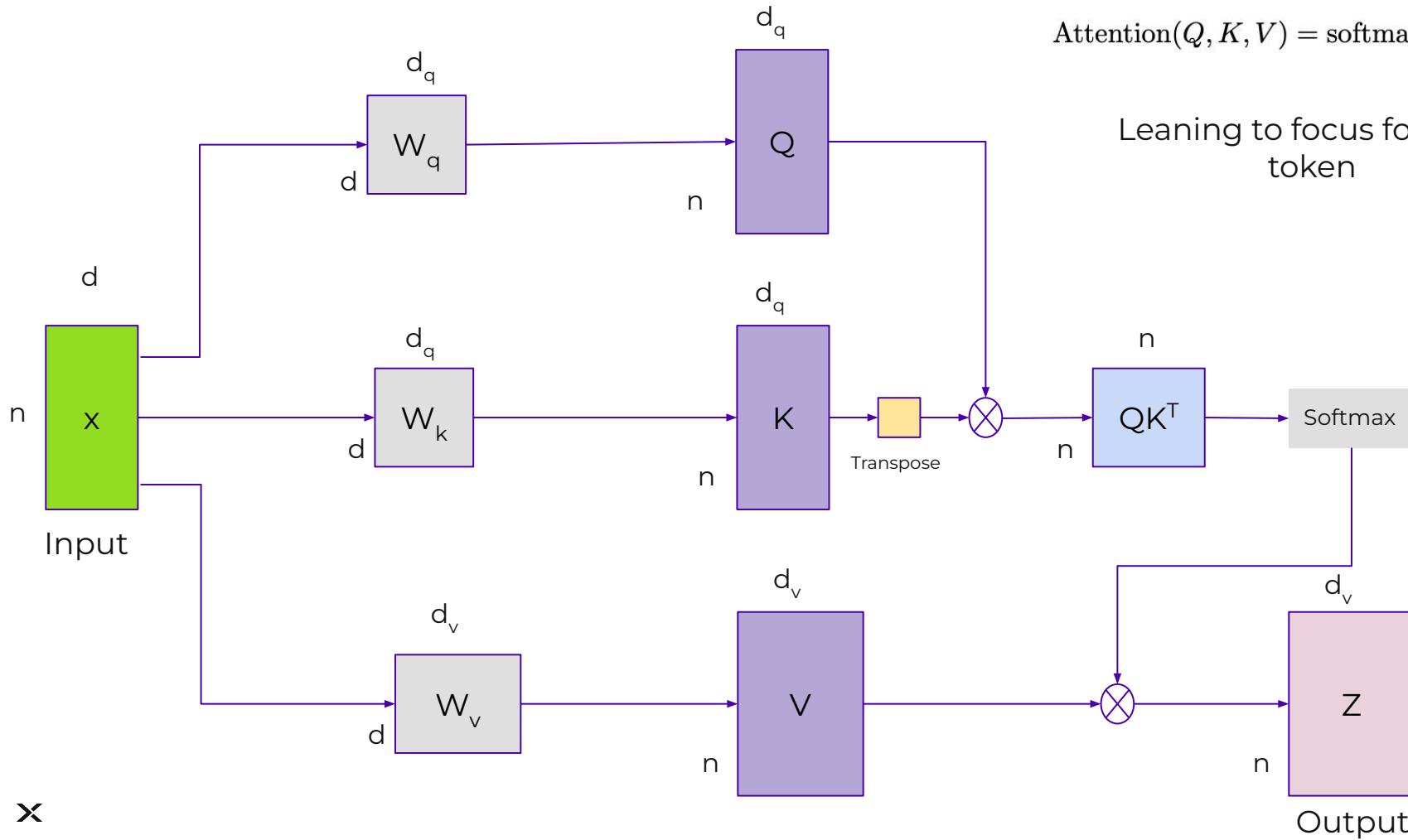
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

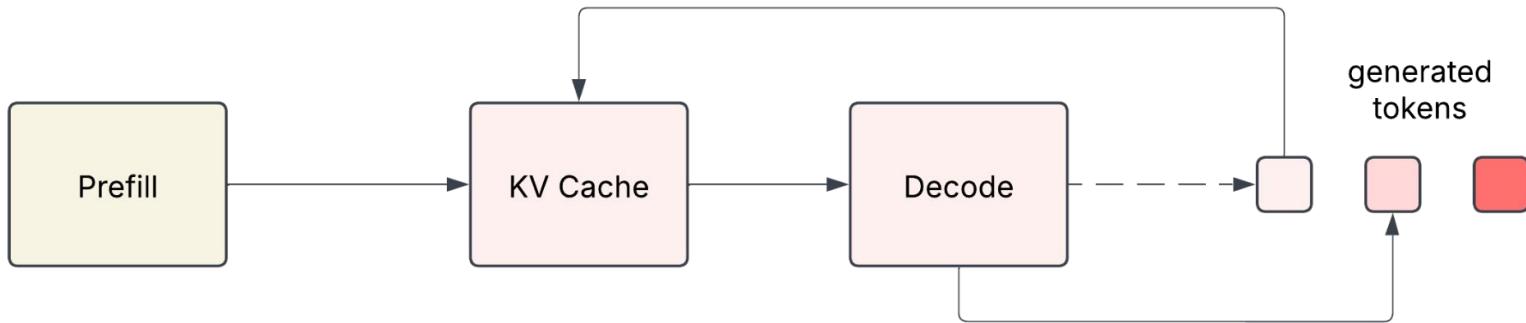
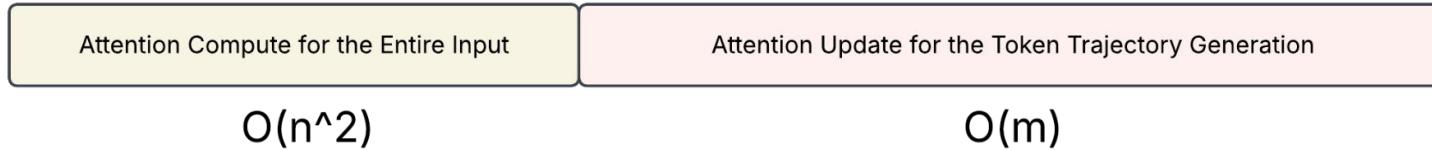
$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O$$

$$\text{where } \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V)$$

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Learning to focus for each token





Paged Attention: KV Cache  
represented as a logical KV blocks

LLM Inference Without KV Cache	LLM Inference With KV Cache
<pre> context = initial_prompt_tokens  for t in range(prompt_length + 1, max_length):     Q = W_q x(context) [ nd<sub>q</sub> ]     K = W_k x (context) [ nd<sub>q</sub> ]     V = W_v x (context) [ nd<sub>v</sub> ]      q_t = Q[-1] # Query for last token [ d<sub>q</sub> ]      attention_scores = softmax(q_t @ K.T / sqrt(d_k)) [ nd<sub>q</sub> ]     o_t = attention_scores @ V [ nd<sub>v</sub> ]      next_token = sample_from(o_t)     context.append(next_token) </pre> <p style="text-align: center;"><b>Overall Complexity: <math>O(nd_q n^2)</math></b></p>	<pre> KV_cache = prefill(prompt_tokens)  for t in range(prompt_length + 1, max_length):     q_t = W_q(x[t-1]) [ dd<sub>q</sub> ]     k_t = W_k(x[t-1]) [ dd<sub>q</sub> ]     v_t = W_v(x[t-1]) [ dd<sub>v</sub> ]      KV_cache.append(k_t, v_t)      attention_scores = softmax(q_t @ KV_cache.K / sqrt(d_k)) [ nd<sub>q</sub> ]     o_t = attention_scores @ KV_cache.V [ nd<sub>q</sub> ]      x[t] = sample_from(o_t) </pre> <p style="text-align: center;"><b>Overall Complexity: <math>O(d<sub>q</sub> n^2)</math></b></p>

```
// paged_attention_config.h
#pragma once

namespace paged_attention {

    struct PagedAttentionConfig {
        int batch_size;
        int num_heads;
        int head_dim;
        int max_seq_len;
        int page_size;
        int num_pages;
        bool recycle_pages;
    };

    // Minimal configuration that will definitely work
    constexpr PagedAttentionConfig DEFAULT_CONFIG = {
        .batch_size = 64, // [ 2, 4, 8, 16, 32, 64]
        .num_heads = 32,
        .head_dim = 64,
        .max_seq_len = 512, // [512, 1024, 2048, 4096, 8192]
        .page_size = 16,
        .num_pages = 2048,
        .recycle_pages = true
    };

} // namespace paged_attention
```

## Experimental Parameters

```
// Simple CUDA kernel for updating KV cache
__global__ void update_kv_cache_kernel(
    float* key_cache,
    float* value_cache,
    const float* keys,
    const float* values,
    const int* page_table,
    int seq_idx,
    int seq_len,
    int max_seq_len,
    int num_heads,
    int head_dim,
    int page_size
) {
    int token_idx = blockIdx.x;
    int head_idx = blockIdx.y;
    int dim_idx = threadIdx.x;

    if (token_idx >= seq_len || head_idx >= num_heads || dim_idx >= head_dim) {
        return;
    }

    int page_idx = page_table[seq_idx * max_seq_len + token_idx];
    if (page_idx < 0) {
        return;
    }

    int offset = token_idx % page_size;

    size_t cache_idx = (page_idx * page_size + offset) * num_heads * head_dim +
                      head_idx * head_dim + dim_idx;

    size_t input_idx = token_idx * num_heads * head_dim +
                      head_idx * head_dim + dim_idx;

    key_cache[cache_idx] = keys[input_idx];
    value_cache[cache_idx] = values[input_idx];
}
```

Divide KV Cache into Contiguous KV Blocks and then Store Key and Value Information there.

## Attention Kernel

```
// Extremely simple attention kernel - focus on correctness
__global__ void simple_attention_kernel(
    const float* key_cache,
    const float* value_cache,
    const float* queries,
    float* outputs,
    const int* page_table,
    const int* seq_lengths,
    int batch_size,
    int num_heads,
    int head_dim,
    int max_seq_len,
    int page_size,
    float scale
) {
    int batch_idx = blockIdx.x;
    int head_idx = blockIdx.y;
    int dim_idx = threadIdx.x;
    if (batch_idx >= batch_size || head_idx >= num_heads || dim_idx >= head_dim) {
        return;
    }
    int seq_len = seq_lengths[batch_idx];
    if (seq_len <= 0) {
        return;
    }
    // Use global memory for scores
    float* scores = (float*)malloc(seq_len * sizeof(float));
    if (!scores) {
        if (dim_idx == 0 && head_idx == 0 && batch_idx == 0) {
            printf("[ERROR] Failed to allocate global memory for scores (seq_len=%d)\n", seq_len);
        }
        return;
    }
```

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

$$\text{softmax}(x) = \frac{e^{x_i}}{\sum_j e^{x_j}}$$

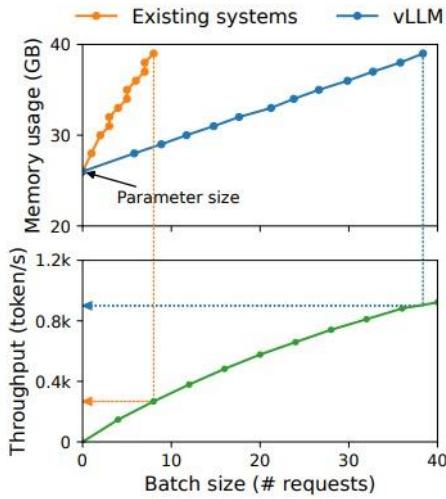
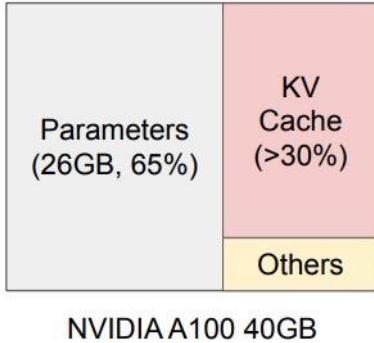
$$\text{softmax}(x) = \text{softmax}(x - \alpha)$$

$$\alpha = \max_i x_i$$

```

int q_idx = (batch_idx * num_heads + head_idx) * head_dim + dim_idx;
int out_idx = q_idx;
float output_val = 0.0f;
float max_score = -INFINITY;
for (int token_idx = 0; token_idx < seq_len; token_idx++) {
    int page_idx = page_table[batch_idx * max_seq_len + token_idx];
    if (page_idx < 0) {
        scores[token_idx] = -INFINITY;
        continue;
    }
    int offset = token_idx % page_size;
    float dot_product = 0.0f;
    for (int d = 0; d < head_dim; d++) {
        int k_idx = (page_idx * page_size + offset) * num_heads * head_dim +
                   head_idx * head_dim + d;
        int q_d_idx = (batch_idx * num_heads + head_idx) * head_dim + d;
        dot_product += queries[q_d_idx] * key_cache[k_idx];
    }
    float score = dot_product * scale;
    scores[token_idx] = score;
    if (score > max_score) {
        max_score = score;
    }
}
float sum_exp = 0.0f;
for (int token_idx = 0; token_idx < seq_len; token_idx++) {
    if (scores[token_idx] == -INFINITY) {
        continue;
    }
    int page_idx = page_table[batch_idx * max_seq_len + token_idx];
    int offset = token_idx % page_size;
    float exp_score = expf(scores[token_idx] - max_score);
    sum_exp += exp_score;
    int v_idx = (page_idx * page_size + offset) * num_heads * head_dim +
               head_idx * head_dim + dim_idx;
    output_val += exp_score * value_cache[v_idx];
}
if (sum_exp > 0.0f) {
    outputs[out_idx] = output_val / sum_exp;
} else {
    outputs[out_idx] = 0.0f;
}

```



**Figure 1.** Left: Memory layout when serving an LLM with 13B parameters on NVIDIA A100. The parameters (gray) persist in GPU memory throughout serving. The memory for the KV cache (red) is (de)allocated per serving request. A small amount of memory (yellow) is used ephemerally for activation. Right: vLLM smooths out the rapid growth curve of KV cache memory seen in existing systems [31, 60], leading to a notable boost in serving throughput.

## Efficient Memory Management for Large Language Model Serving with *PagedAttention*

Woosuk Kwon<sup>1,\*</sup> Zhuohan Li<sup>1,\*</sup> Siyuan Zhuang<sup>1</sup> Ying Sheng<sup>1,2</sup> Lianmin Zheng<sup>1</sup> Cody Hao Yu<sup>3</sup>  
Joseph E. Gonzalez<sup>1</sup> Hao Zhang<sup>4</sup> Ion Stoica<sup>1</sup>

<sup>1</sup>UC Berkeley <sup>2</sup>Stanford University <sup>3</sup>Independent Researcher <sup>4</sup>UC San Diego

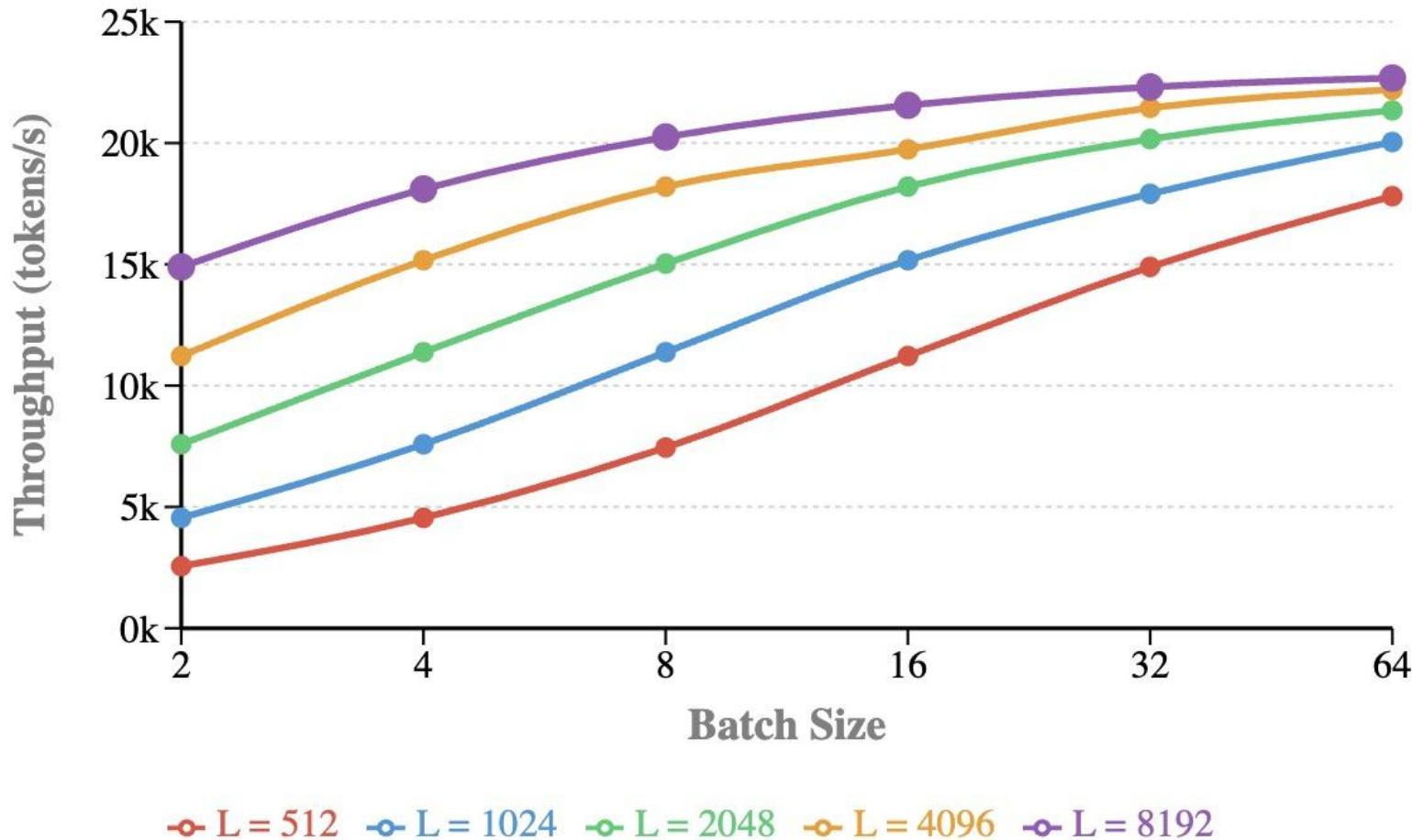
Can we reproduced the patterns shown in the Paged Attention Paper?

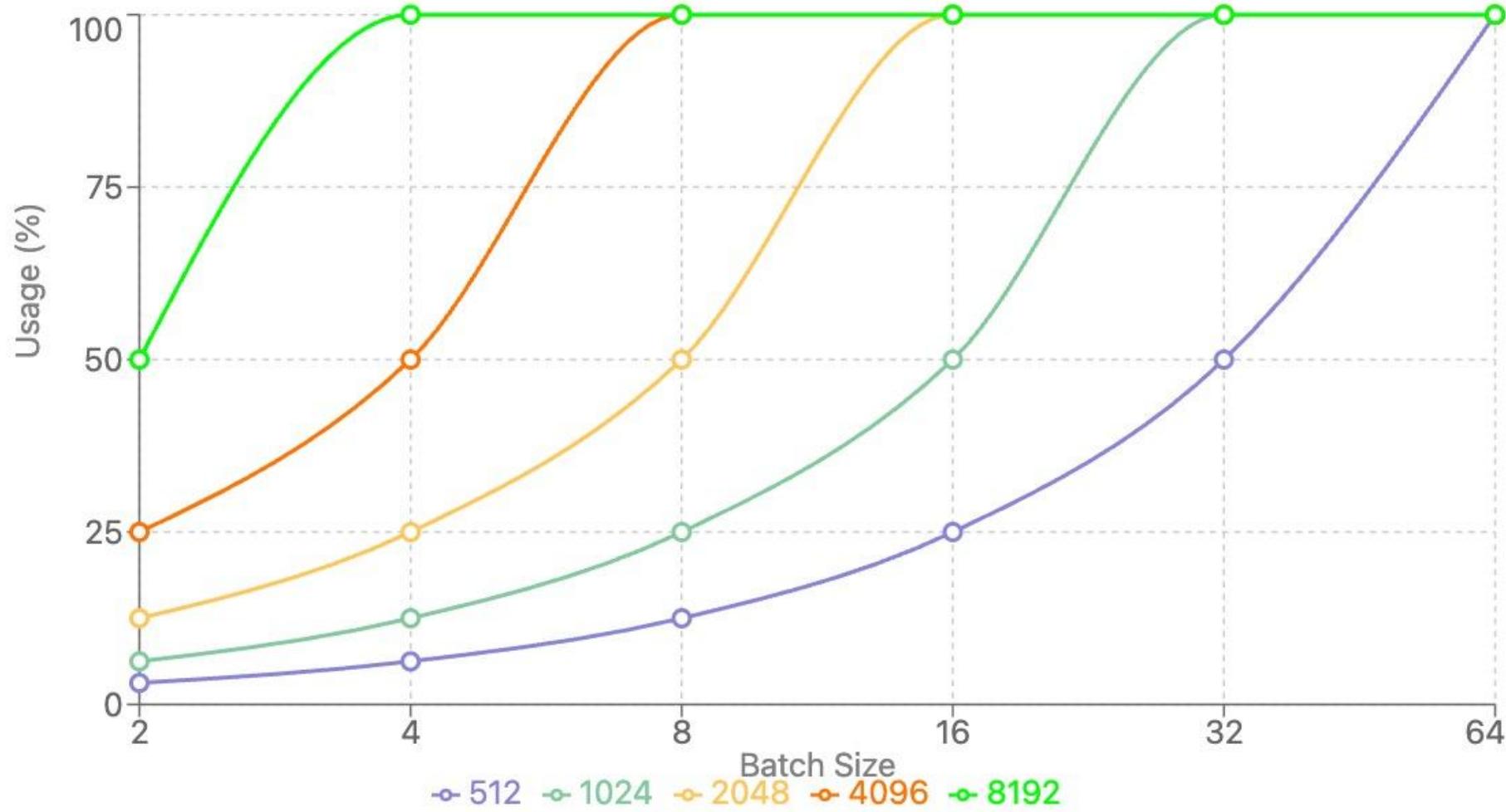
Batch Size	Num of Heads	Head Dimension	Page Size	Num of Pages	Max Seq Length	Response Time (s)	Throughput (tokens/s)
4	32	64	16	2048	512	0.45	4551
4	32	64	16	2048	1024	0.54	7585
4	32	64	16	2048	2048	0.72	11378
4	32	64	16	2048	4096	1.08	15170
4	32	64	16	2048	8192	1.81	18104
2	32	64	16	2048	512	0.40	2560
2	32	64	16	2048	1024	0.45	4551
2	32	64	16	2048	2048	0.54	7585
2	32	64	16	2048	4096	0.73	11222
2	32	64	16	2048	8192	1.10	14895
2	32	64	16	2048	16384	Failed X	Failed X
1	32	32	16	2048	16384	Failed X	Failed X

X

Batch Size	Num of Heads	Head Dimension	Page Size	Num of Pages	Max Seq Length	Response Time (s)	Throughput (tokens/s)
16	32	64	16	2048	512	0.73	11222
16	32	64	16	2048	1024	1.08	15170
16	32	64	16	2048	2048	1.80	18204
16	32	64	16	2048	4096	3.32	19740
16	32	64	16	2048	8192	6.08	21558
8	32	64	16	2048	512	0.55	7447
8	32	64	16	2048	1024	0.72	11378
8	32	64	16	2048	2048	1.09	15031
8	32	64	16	2048	4096	1.8	18204
8	32	64	16	2048	8192	3.237	20246

Batch Size	Num of Heads	Head Dimension	Page Size	Num of Pages	Max Seq Length	Response Time (s)	Throughput (tokens/s)
64	32	64	16	2048	512	1.84	17809
64	32	64	16	2048	1024	3.27	20042
64	32	64	16	2048	2048	6.14	21347
64	32	64	16	2048	4096	11.81	22197
64	32	64	16	2048	8192	23.12	22677
32	32	64	16	2048	512	1.10	14895
32	32	64	16	2048	1024	1.83	17906
32	32	64	16	2048	2048	3.25	20165
32	32	64	16	2048	4096	6.11	21452
32	32	64	16	2048	8192	11.75	22310





# Thank You!

“People who are really serious about software should make their own hardware.” Alan Kay