CEN310 Parallel Programming

Week-5

GPU Programming

Download

- PDF
- DOC
- SLIDE
- PPTX







CEN31 Qutlingar (1)/4) eek-5

- 1. Introduction to GPU Computing
 - GPU Architecture Overview
 - CUDA Programming Model
 - GPU Memory Hierarchy
 - Thread Organization
 - Kernel Functions
- 2. CUDA Programming Basics
 - Memory Management
 - Thread Organization
 - Synchronization
 - Error Handling

CEN310 Week-5A Runtime API

CEN31 Outlinear (21/941) ek-5

3. Memory Management

- Global Memory
- Shared Memory
- Constant Memory
- Texture Memory
- Unified Memory

4. Thread Organization

- Blocks and Grids
- Warps and Scheduling
- Thread Synchronization
- Occupancy
- CEN310 Week-5 Load Balancing

CEN31 Outlinear (3/94) ek-5

5. Performance Optimization

- Memory Coalescing
- Bank Conflicts
- Divergent Branching
- Shared Memory Usage
- Asynchronous Operations

6. Advanced Features

- Streams and Events
- Dynamic Parallelism
- Multi-GPU Programming
- Unified Memory
- CEN310 Week-5 Cooperative Groups

CEN310 Parallel Programming Week-5 Outline (4/4)

7. Best Practices

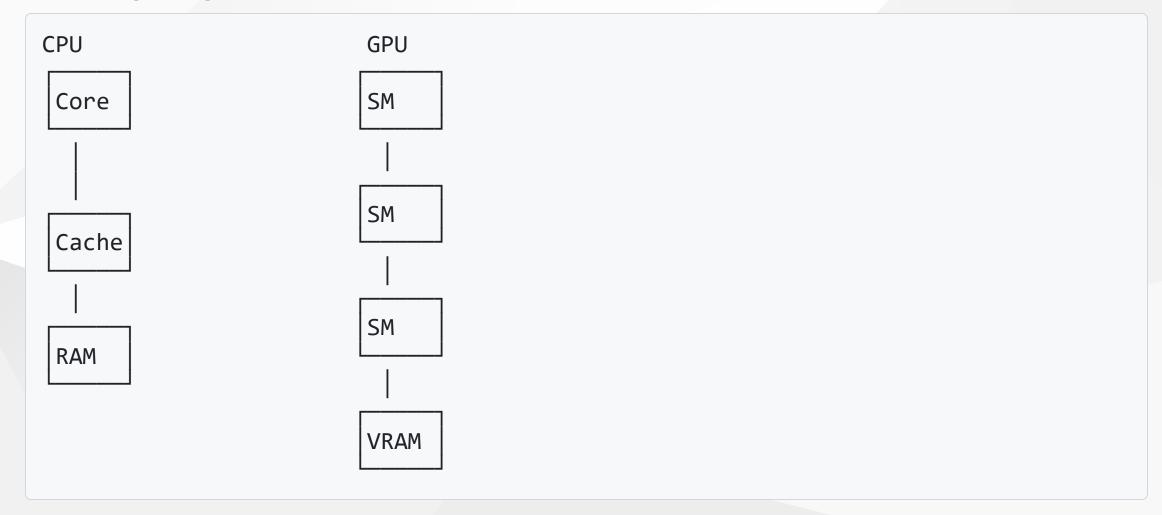
- Code Organization
- Error Handling
- Debugging Techniques
- Profiling Tools
- Common Pitfalls

8. Real-World Applications

- Image Processing
- Scientific Computing
- Machine Learning



CEN3 16 Parallel Architecture (1/4)



Key Components:

Streaming Multiprocessors (SMs)

RTEU CENTUDA Cores

CEN31GP Je Architecture (2/4)

Memory Hierarchy

```
// Example showing different memory types
device constant float device constant[256]; // Constant memory
__shared__ float shared_array[256];
                                       // Shared memory
__global__ void memory_example(float* global_input, // Global memory
                            float* global output) {
    int idx = threadIdx.x + blockIdx.x * blockDim.x;
   // Register (automatic) variables
   float local var = global input[idx];
   // Shared memory usage
    shared array[threadIdx.x] = local var;
    __syncthreads();
   // Constant memory usage
    local var *= device constant[threadIdx.x];
   // Write back to global memory
   global output[idx] = local var;
  CFN310 Week-5
```

CEN31GP Je Architecture (3/4)

Thread Hierarchy

```
__global__ void thread_hierarchy_example() {
    // Thread identification
    int thread idx = threadIdx.x;
    int block idx = blockIdx.x;
    int warp id = threadIdx.x / 32;
    // Block dimensions
    int block size = blockDim.x;
    int grid size = gridDim.x;
    // Global thread ID
    int global idx = thread idx + block idx * block size;
    // Print thread information
    printf("Thread %d in block %d (warp %d)\n",
           thread idx, block idx, warp id);
int main() {
    // Launch configuration
    dim3 block size(256);
    dim3 grid size((N + block size.x - 1) / block size.x);
    thread hierarchy example<<<grid size, block size>>>();
    return 0;
 CEN310 Week-5
```

CEN31GP Je Architecture (4/4)

Basic CUDA Program

```
#include <cuda_runtime.h>
// Kernel definition
__global__ void vector_add(float* a, float* b, float* c, int n) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if(idx < n) {</pre>
        c[idx] = a[idx] + b[idx];
int main() {
    size t size = N * sizeof(float);
    // Allocate host memory
    float *h_a = (float*)malloc(size);
    float *h_b = (float*)malloc(size);
    float *h_c = (float*)malloc(size);
    // Initialize arrays
    for(int i = 0; i < N; i++) {
        h_a[i] = rand() / (float)RAND_MAX;
h_b[i] = rand() / (float)RAND_MAX;
    // Allocate device memory
    float *d_a, *d_b, *d_c;
    cudaMalloc(&d_a, size);
    cudaMalloc(&d_b, size);
    cudaMalloc(&d_c, size);
    // Copy to device
    cudaMemcpy(d_a, h_a, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_b, h_b, size, cudaMemcpyHostToDevice);
    // Launch kernel
    int block_size = 256;
    int num_blocks = (N + block_size - 1) / block_size;
    vector_add<<<num_blocks, block_size>>>(d_a, d_b, d_c, N);
    cudaMemcpy(h_c, d_c, size, cudaMemcpyDeviceToHost);
    // Verify results
    for(int i = 0; i < N; i++) {
        if(fabs(h_c[i] - (h_a[i] + h_b[i])) > 1e-5) {
    fprintf(stderr, "Verification failed at %d\n", i);
            break;
    // Cleanup
    cudaFree(d_a);
    cudaFree(d b);
    cudaFree(d_c);
    free(h_a);
    free(h_b);
    free(h_c);
```



RTEU CEN310 Week-5

CEN312 Rar CUDA Programming Basics

Memory Management (1/4)

```
void memory management example() {
   // Host memory allocation
   float* h data = (float*)malloc(size);
   // Device memory allocation
   float* d data;
   cudaMalloc(&d data, size);
   // Pinned memory allocation
   float* h pinned;
   cudaMallocHost(&h pinned, size);
   // Unified memory allocation
   float* unified;
   cudaMallocManaged(&unified, size);
   // Memory transfers
   cudaMemcpy(d data, h_data, size, cudaMemcpyHostToDevice);
   cudaMemcpy(h data, d data, size, cudaMemcpyDeviceToHost);
   // Asynchronous transfers
   cudaStream t stream;
   cudaStreamCreate(&stream);
   cudaMemcpyAsync(d_data, h_pinned, size,
                    cudaMemcpyHostToDevice, stream);
   // Cleanup
   free(h data);
   cudaFree(d_data);
   cudaFreeHost(h pinned);
    cudaFree(unified);
    cudaStreamDestroy(stream);
```

CEN31MemoryaManagement (2/4)

Shared Memory Usage

```
global void shared memory example(float* input,
                                    float* output,
                                    int n) {
    extern shared float shared[];
    int tid = threadIdx.x;
    int gid = blockIdx.x * blockDim.x + threadIdx.x;
    // Load data into shared memory
    if(gid < n) {</pre>
        shared[tid] = input[gid];
    syncthreads();
    // Process data in shared memory
    if(tid > 0 && tid < blockDim.x-1 && gid < n-1) {</pre>
        float result = 0.25f * (shared[tid-1] +
                               2.0f * shared[tid] +
                                shared[tid+1]);
        output[gid] = result;
// Kernel launch
int block size = 256;
int shared size = block_size * sizeof(float);
shared_memory_example<<<grid_size, block_size, shared_size>>>
FIJ ( [d] zimput e d_output, N);
```



CEN31 Memory a Management (3/4)

Constant Memory

```
__constant__ float const_array[256];
void setup constant memory() {
    float h_const_array[256];
    // Initialize constant data
    for(int i = 0; i < 256; i++) {
        h const array[i] = compute constant(i);
    // Copy to constant memory
    cudaMemcpyToSymbol(const array, h const array,
                       256 * sizeof(float));
 global void use constant memory(float* input,
                                  float* output,
                                  int n) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if(idx < n) {
       // Use constant memory
        output[idx] = input[idx] * const_array[idx % 256];
  CEN310 Week-5
```

CEN31MemoryaManagement (4/4)

Texture Memory

```
texture<float, 2, cudaReadModeElementType> tex ref;
void texture_memory_example() {
    // Allocate and initialize 2D array
   cudaArray* d array;
   cudaChannelFormatDesc channel desc =
       cudaCreateChannelDesc<float>();
   cudaMallocArray(&d_array, &channel_desc,
                   width, height);
   // Copy data to array
   cudaMemcpyToArray(d_array, 0, 0, h_data,
                     width * height * sizeof(float),
                      cudaMemcpyHostToDevice);
    // Bind texture reference
   cudaBindTextureToArray(tex_ref, d_array);
   // Kernel using texture memory
   texture kernel<<<grid size, block size>>>
        (d_output, width, height);
    // Cleanup
   cudaUnbindTexture(tex ref);
   cudaFreeArray(d array);
global void texture kernel(float* output,
                            int width, int height) {
   int x = blockIdx.x * blockDim.x + threadIdx.x;
   int y = blockIdx.y * blockDim.y + threadIdx.y;
   if(x < width && y < height) {</pre>
       // Read from texture
       float value = tex2D(tex_ref, x, y);
       output[y * width + x] = value;
BU CEN310 Week-5
```

CEN313 Par JIhread Organization

Thread Hierarchy (1/4)

```
__global__ void thread organization example() {
   // Thread indices
   int tx = threadIdx.x;
   int ty = threadIdx.y;
   int tz = threadIdx.z;
   // Block indices
   int bx = blockIdx.x;
   int by = blockIdx.y;
   int bz = blockIdx.z;
   // Block dimensions
   int bdx = blockDim.x;
   int bdy = blockDim.y;
   int bdz = blockDim.z;
   // Grid dimensions
   int gdx = gridDim.x;
   int gdy = gridDim.y;
   int gdz = gridDim.z;
   // Calculate global indices
   int global_x = bx * bdx + tx;
   int global_y = by * bdy + ty;
   int global z = bz * bdz + tz;
   // Calculate linear index
   int linear_idx = global_z * gdx * gdy * bdx * bdy +
                    global_y * gdx * bdx +
                    global_x;
```

CEN31Thread Hierarchy (2/4)

Block Configuration

```
void launch configuration example() {
    // 1D configuration
    dim3 block 1d(256);
    dim3 grid 1d((N + block 1d.x - 1) / block 1d.x);
    kernel 1d<<<grid 1d, block 1d>>>();
    // 2D configuration
    dim3 block_2d(16, 16);
    dim3 grid_2d((width + block_2d.x - 1) / block_2d.x,
                 (height + block 2d.y - 1) / block 2d.y);
    kernel 2d<<<grid 2d, block 2d>>>();
    // 3D configuration
    dim3 block_3d(8, 8, 8);
    dim3 grid 3d((width + block 3d.x - 1) / block 3d.x,
                 (height + block 3d.y - 1) / block 3d.y,
                 (depth + block 3d.z - 1) / block 3d.z);
    kernel 3d<<<grid 3d, block 3d>>>();
// Kernel examples
__global__ void kernel_1d() {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
__global__ void kernel_2d() {
    int x = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;
__global__ void kernel_3d() {
    int x = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;
    int z = blockIdx.z * blockDim.z + threadIdx.z;
TBU CEN310 Week-5
```



CEN3 Thread Hierarchy (3/4)

Warp Management

```
__global__ void warp_example() {
   int tid = threadIdx.x;
   int warp id = tid / 32;
   int lane id = tid % 32;
   // Warp-level primitives
   int mask = __ballot_sync(__activemask(), tid < 16);</pre>
   int value = shfl sync( activemask(), tid, 0);
   // Warp-level reduction
   int sum = tid;
   for(int offset = 16; offset > 0; offset /= 2) {
       sum += __shfl_down_sync(__activemask(), sum, offset);
   // Warp-level synchronization
   __syncwarp();
```

CEN31Thread Hierarchy5 (4/4)

Dynamic Parallelism

```
global void child kernel(int* data, int size) {
   int idx = blockIdx.x * blockDim.x + threadIdx.x;
   if(idx < size) {</pre>
       data[idx] *= 2;
__global__ void parent_kernel(int* data,
                            int* sizes,
                            int num_arrays) {
   int array idx = blockIdx.x * blockDim.x + threadIdx.x;
   if(array_idx < num_arrays) {</pre>
       int size = sizes[array idx];
       int* array data = &data[array_idx * MAX_ARRAY_SIZE];
       // Launch child kernel
       int block size = 256;
       int grid size = (size + block size - 1) / block size;
       child_kernel<<<grid_size, block_size>>>
            (array data, size);
```

CEN314 Par Performance Optimization

Memory Coalescing (1/4)

```
// Bad memory access pattern
__global__ void uncoalesced_access(float* input,
                                  float* output,
                                  int width,
                                  int height) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if(idx < width) {</pre>
        for(int y = 0; y < height; y++) {
            output[idx + y * width] =
                input[idx + v * width]; // Strided access
// Good memory access pattern
global void coalesced_access(float* input,
                                float* output,
                                int width,
                                int height) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    int y = blockIdx.y * blockDim.y + threadIdx.y;
    if(idx < width && y < height) {</pre>
        output[y * width + idx] =
            input[y * width + idx]; // Coalesced access
```

Memory Coalescing (2/4)

Bank Conflicts

```
__global__ void bank_conflicts_example(float* data) {
   extern __shared__ float shared[];
   int tid = threadIdx.x;
   // Bad: Bank conflicts
   shared[tid * 32] = data[tid]; // 32-way bank conflict
   // Good: No bank conflicts
   shared[tid] = data[tid];  // Consecutive access
   __syncthreads();
   // Process data
   float result = shared[tid];
   // ...
```

CEN31MemoryaGoalescing (3/4)

Shared Memory Optimization

```
template<int BLOCK SIZE>
__global__ void matrix_multiply(float* A,
                               float* B,
                              float* C,
                               int width) {
    __shared__ float shared_A[BLOCK_SIZE][BLOCK_SIZE];
    __shared__ float shared_B[BLOCK_SIZE][BLOCK_SIZE];
    int tx = threadIdx.x;
    int ty = threadIdx.y;
    int bx = blockIdx.x;
    int by = blockIdx.y;
    int row = by * BLOCK SIZE + ty;
    int col = bx * BLOCK SIZE + tx;
    float sum = 0.0f;
    // Loop over blocks
    for(int block = 0; block < width/BLOCK_SIZE; block++) {</pre>
        // Load data into shared memory
        shared_A[ty][tx] = A[row * width +
                            block * BLOCK_SIZE + tx];
        shared_B[ty][tx] = B[(block * BLOCK_SIZE + ty) * width +
                             col];
        syncthreads();
        // Compute partial dot product
        for(int k = 0; k < BLOCK SIZE; k++) {</pre>
            sum += shared A[ty][k] * shared B[k][tx];
        __syncthreads();
    // Store result
    C[row * width + col] = sum;
TPU CEN310 Week-5
```



CEN31MemoryaGoalescing (4/4)

Memory Access Patterns

```
// Structure of Arrays (SoA)
struct ParticlesSoA {
   float* x;
   float* y;
   float* z;
    float* vx;
    float* vy;
    float* vz;
// Array of Structures (AoS)
struct ParticleAoS {
   float x, y, z;
    float vx, vy, vz;
};
// SoA kernel (better coalescing)
__global__ void update_particles_soa(ParticlesSoA particles,
                                   int n) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if(idx < n) {</pre>
        particles.x[idx] += particles.vx[idx];
        particles.y[idx] += particles.vy[idx];
        particles.z[idx] += particles.vz[idx];
// AoS kernel (worse coalescing)
__global__ void update_particles_aos(ParticleAoS* particles,
                                   int n) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if(idx < n) {
        particles[idx].x += particles[idx].vx;
        particles[idx].y += particles[idx].vy;
        particles[idx].z += particles[idx].vz;
BU CEN310 Week-5
```



CEN315 Randvanced - Freatures

Streams and Events (1/3)

```
void stream_example() {
    const int num streams = 4;
    cudaStream t streams[num streams];
    // Create streams
    for(int i = 0; i < num_streams; i++) {</pre>
        cudaStreamCreate(&streams[i]);
    // Allocate memory
    float *h_input, *d_input, *h_output, *d_output;
    cudaMallocHost(&h_input, size); // Pinned memory
    cudaMallocHost(&h_output, size); // Pinned memory
    cudaMalloc(&d_input, size);
    cudaMalloc(&d_output, size);
    // Launch kernels in different streams
    int chunk size = N / num streams;
    for(int i = 0; i < num_streams; i++) {</pre>
        int offset = i * chunk_size;
        cudaMemcpyAsync(&d input[offset],
                        &h input[offset],
                        chunk_size * sizeof(float),
                        cudaMemcpyHostToDevice,
                        streams[i]);
        process_kernel<<<grid_size, block_size, 0, streams[i]>>>
            (&d_input[offset], &d_output[offset], chunk_size);
        cudaMemcpyAsync(&h_output[offset],
                        &d output[offset],
                        chunk size * sizeof(float),
                        cudaMemcpyDeviceToHost,
                        streams[i]);
    // Synchronize all streams
    cudaDeviceSynchronize();
    for(int i = 0; i < num streams; i++) {</pre>
        cudaStreamDestroy(streams[i]);
    cudaFreeHost(h_input);
    cudaFreeHost(h_output);
    cudaFree(d_input);
```

CEN31StreamsrandgEvents (2/3)

Event Management

```
void event_example() {
    cudaEvent_t start, stop;
    cudaEventCreate(&start);
    cudaEventCreate(&stop);
    // Record start event
    cudaEventRecord(start);
    // Launch kernel
    process kernel<<<grid size, block size>>>(d data, N);
    // Record stop event
    cudaEventRecord(stop);
    // Wait for completion
    cudaEventSynchronize(stop);
    // Calculate elapsed time
    float milliseconds = 0;
    cudaEventElapsedTime(&milliseconds, start, stop);
    printf("Kernel execution time: %f ms\n", milliseconds);
    // Cleanup
    cudaEventDestroy(start);
    cudaEventDestroy(stop);
FU CEN310 Week-5
```

CEN31StreamsrandgEvents (3/3)

Inter-stream Synchronization

```
void stream_synchronization() {
    cudaStream_t stream1, stream2;
    cudaStreamCreate(&stream1);
    cudaStreamCreate(&stream2);
    cudaEvent t event;
    cudaEventCreate(&event);
    // Launch work in stream1
    kernel1<<<grid_size, block_size, 0, stream1>>>
        (d data1, N);
    cudaEventRecord(event, stream1);
    // Make stream2 wait for stream1
    cudaStreamWaitEvent(stream2, event);
    // Launch work in stream2
    kernel2<<<grid_size, block_size, 0, stream2>>>
        (d data2, N);
    // Cleanup
    cudaEventDestroy(event);
    cudaStreamDestroy(stream1);
    cudaStreamDestroy(stream2);
EU CEN310 Week-5
```

CEN316RarBestgrPmractices

Error Handling (1/3)

```
#define CUDA_CHECK(call) do {
    cudaError_t error = call;
    if(error != cudaSuccess) {
        fprintf(stderr, "CUDA error at %s:%d: %s\n",
                __FILE__, __LINE__,
                cudaGetErrorString(error));
        exit(EXIT FAILURE);
} while(0)
void cuda_error_handling() {
    // Allocate memory
    float* d data;
    CUDA CHECK(cudaMalloc(&d data, size));
    // Launch kernel
    process_kernel<<<grid_size, block_size>>>(d_data, N);
    CUDA CHECK(cudaGetLastError());
    // Synchronize and check for errors
    CUDA CHECK(cudaDeviceSynchronize());
    // Cleanup
    CUDA_CHECK(cudaFree(d_data));
```

CEN31Error Handling (25/3)

Debug Tools

```
void debug example() {
    // Enable device synchronization for debugging
    cudaDeviceSetLimit(cudaLimitDevRuntimeSyncDepth, 1);
    // Print device properties
    cudaDeviceProp prop;
    cudaGetDeviceProperties(&prop, 0);
    printf("Device: %s\n", prop.name);
    printf("Compute capability: %d.%d\n",
           prop.major, prop.minor);
    // Launch kernel with debug info
    #ifdef DEBUG
        printf("Launching kernel with grid=%d, block=%d\n",
               grid size.x, block size.x);
   #endif
    process_kernel<<<grid_size, block_size>>>(d_data, N);
    // Check for kernel errors
    cudaError t error = cudaGetLastError();
    if(error != cudaSuccess) {
        printf("Kernel error: %s\n",
               cudaGetErrorString(error));
  CEN310 Week-5
```

CEN31Error Handling (35/3)

Resource Management

```
class CUDAResource {
private:
    void* ptr;
    size_t size;
public:
    CUDAResource(size_t s) : size(s), ptr(nullptr) {
        CUDA CHECK(cudaMalloc(&ptr, size));
    ~CUDAResource() {
        if(ptr) {
            cudaFree(ptr);
    void* get() { return ptr; }
    size_t get_size() { return size; }
    // Prevent copying
    CUDAResource(const CUDAResource&) = delete;
    CUDAResource& operator=(const CUDAResource&) = delete;
};
void resource_management_example() {
    try {
        CUDAResource d_input(1024);
        CUDAResource d_output(1024);
        // Use resources
        process_kernel<<<grid_size, block_size>>>
            (d_input.get(), d_output.get(), N);
    catch(const std::exception& e) {
        fprintf(stderr, "Error: %s\n", e.what());
BU CEN310 Week-5
```



7. Real-World Applications

Image Processing (1/3)

```
// Image convolution kernel
__global__ void convolution_2d(unsigned char* input,
```