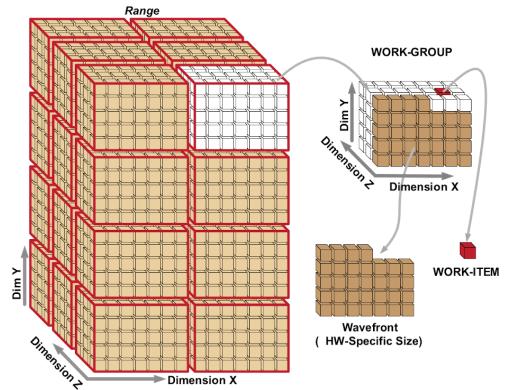


OpenCL for High-Performance Computing

Created by Shromm
Gaind



OPEN MPI



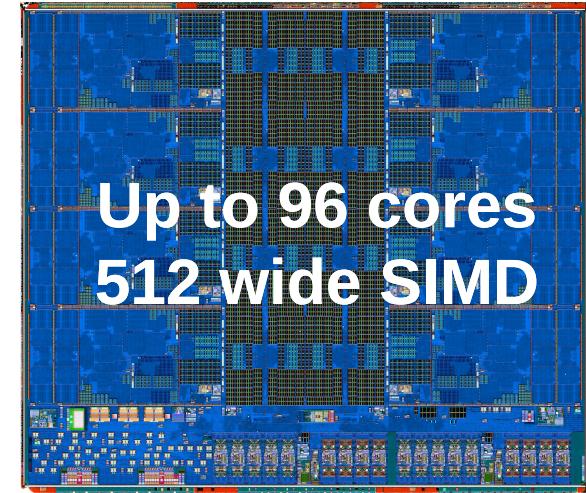
We Can't Just Wait for Faster CPUs

For decades, our code got faster for free. Those days are over. The industry's answer isn't faster cores, it's more cores.

- For performance, concurrency is no longer optional. It's mandatory.



Intel® Granite Rapids
microarchitecture

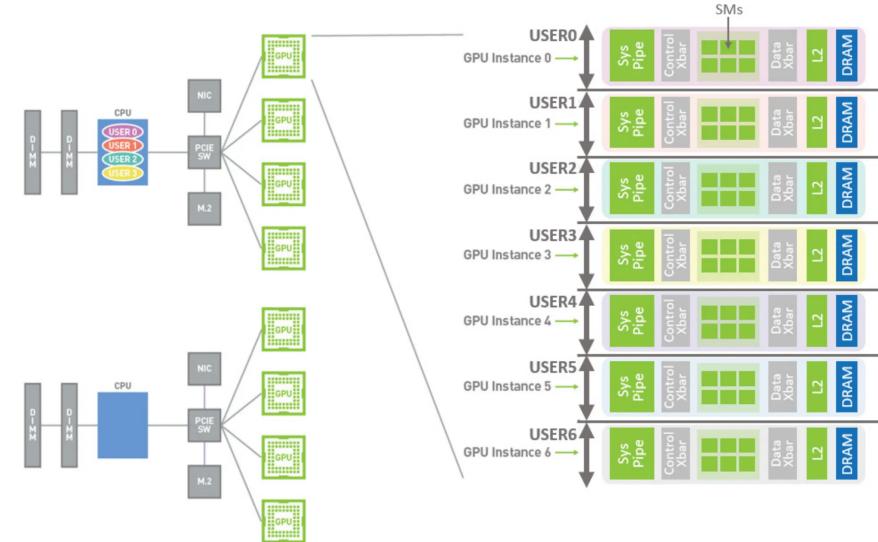


AMD® Zen 5
microarchitecture

Heterogeneous Computing

HPC is about mapping latency-sensitive control flow to the CPU and data-parallel, high-throughput work to the GPU.

- Modern HPC systems are built on this principle.



Why OpenCL for HPC



Powerful, Mature Ecosystem

Proprietary: NVIDIA hardware only



Khronos Open standard for heterogeneous processing

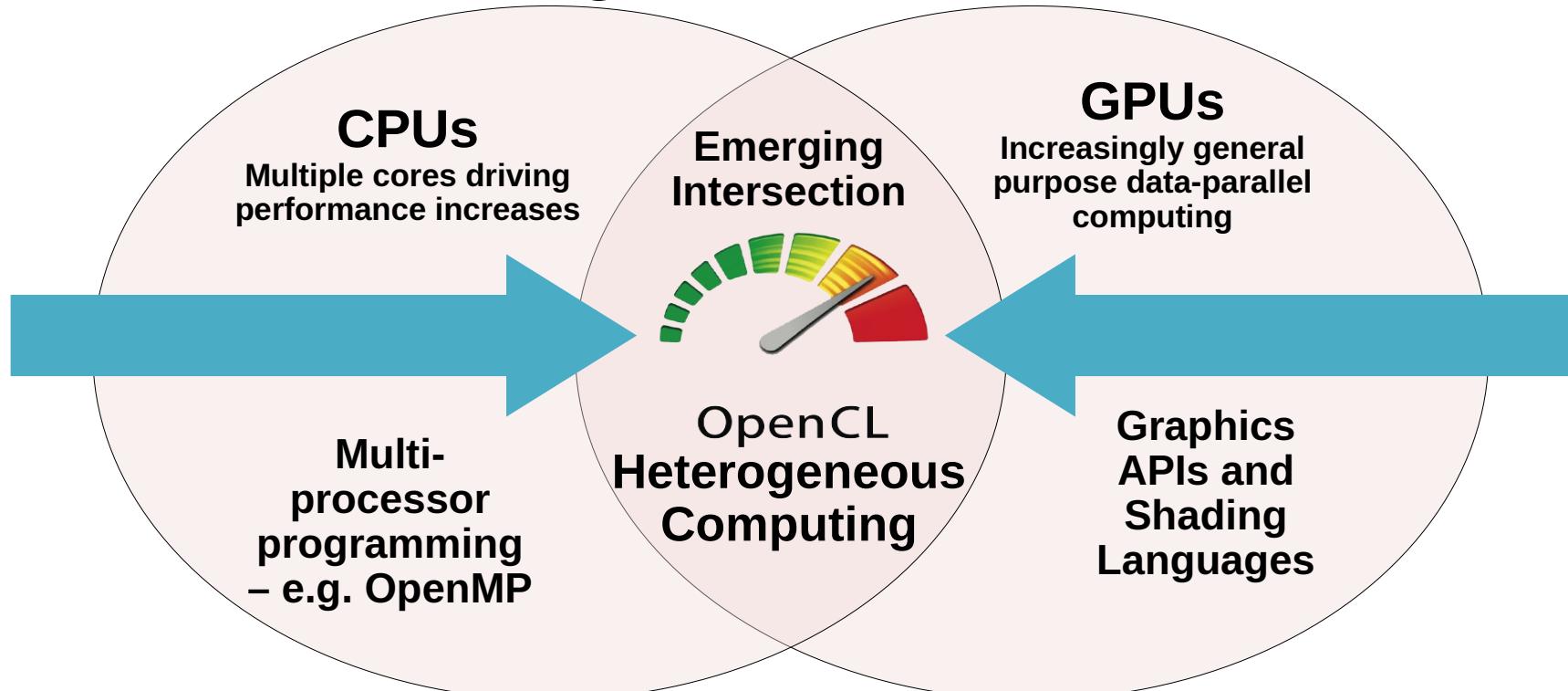
For host-side control (C/C++).

OpenCL C (based on C99 with extensions).

Modern Implementations such as SYCL, DPC++



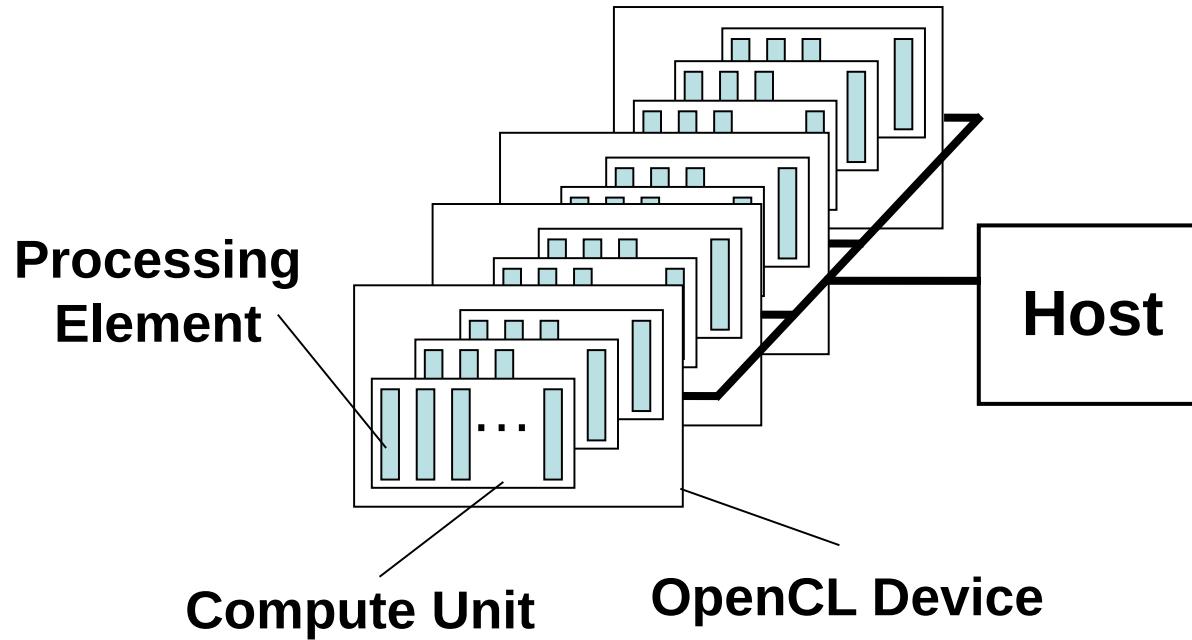
Industry Standards for Programming Heterogeneous Platforms



OpenCL – Open Computing Language

Open, royalty-free standard for portable, parallel programming of heterogeneous parallel computing CPUs, GPUs, and other processors

OpenCL Platform Model



Memory divided into ***host memory (CPU)*** and ***device memory (GPU, FPGA ...)***

Kernels in GPU Programming

- Replace loops with functions (a **kernel**) executing at each point in a problem domain

E.g., process a 1024x1024 image with one kernel invocation per pixel or $1024 \times 1024 = 1,048,576$ kernel executions

Serial loop

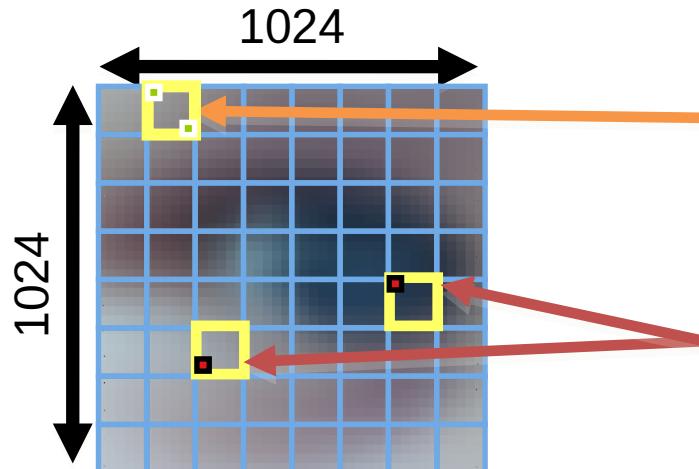
```
for (int y = 0; y < H; ++y) {  
    for (int x = 0; x < W; ++x)  
    {  
        process_pixel(x, y);  
    }  
}
```

Data Parallel OpenCL

```
// Kernel: runs once per (x,y) work-item  
__kernel void process_pixel(__global uchar4*  
image,  
                           const uint width)  
{  
    // each work-item computes its own coords  
    uint x    = get_global_id(0);  
    uint y    = get_global_id(1);  
    uint idx = y * width + x;  
  
    uchar4 px = image[idx];  
    // ...do your per-pixel work...  
    image[idx] = px;  
}
```

An N-dimensional domain of work-items

- **Global Dimensions :**
 - e.g 1024x1024 (whole problem space)
 - **Choose the dimensions that are specific to your algorithm**
- **Local Dimensions:**
 - e.g 64x64 (**work-group**, executes together)



Synchronization between **work-items** possible only within **work-groups**:
barriers and **memory fences**

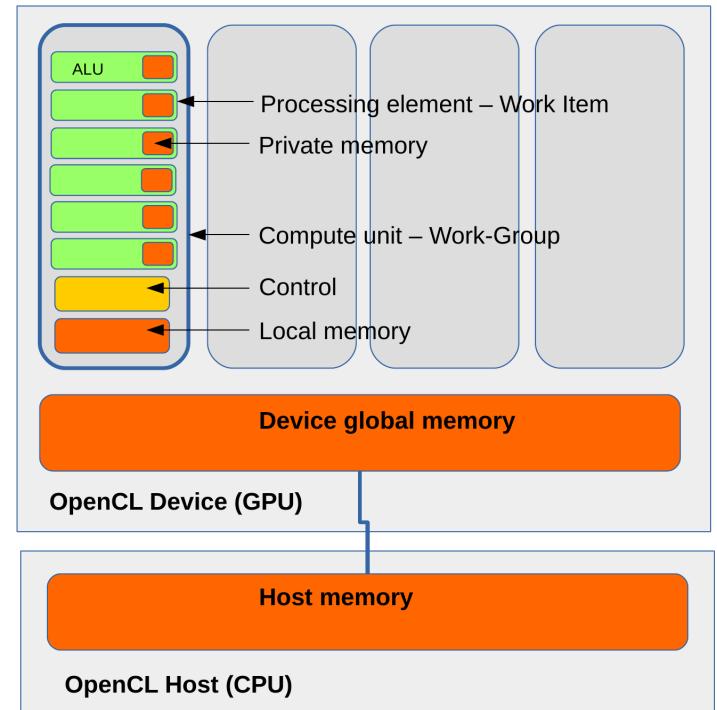
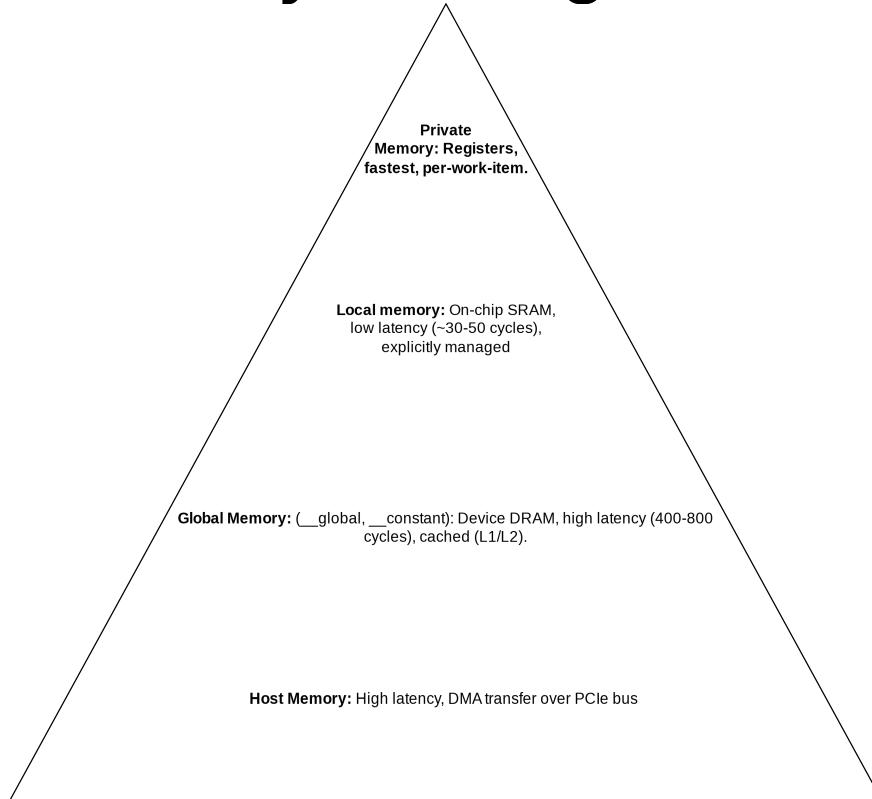
Cannot synchronize
between **work-groups**
within a kernel

Organising the Work: The N-Dimensional Range

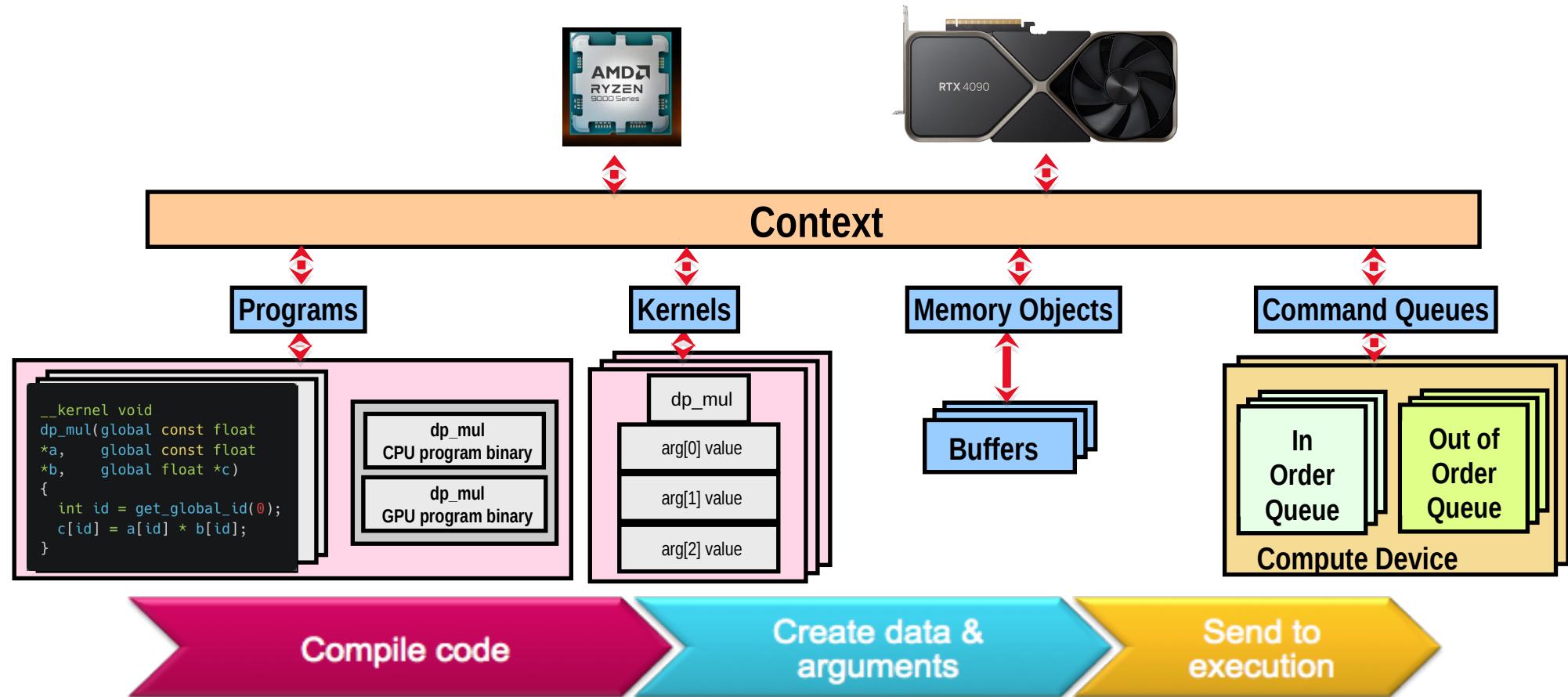
- The **NDRRange** is an N-dimensional index space (where N is 1, 2, or 3) that defines the total domain of computation. Each unique point in this space is a **Work-Item**.
 - This global index space is partitioned into uniform, contiguous blocks called **Work-Groups**.
 - Global Work Size (`global_work_size`)
 - Local Work Size (`local_work_size`)
- All work-items within a single work-group are guaranteed to execute concurrently on a single Compute Unit (CU).
 - This concurrency enables
 - Fast Data Sharing via on-chip local memory.
 - Low-latency Synchronization using the `barrier()` primitive.
- You specify `local_work_size` when enqueueing a kernel. Passing `NULL` lets the runtime decide, but this is **almost never optimal**

OpenCL Memory Model

- Memory management is explicit

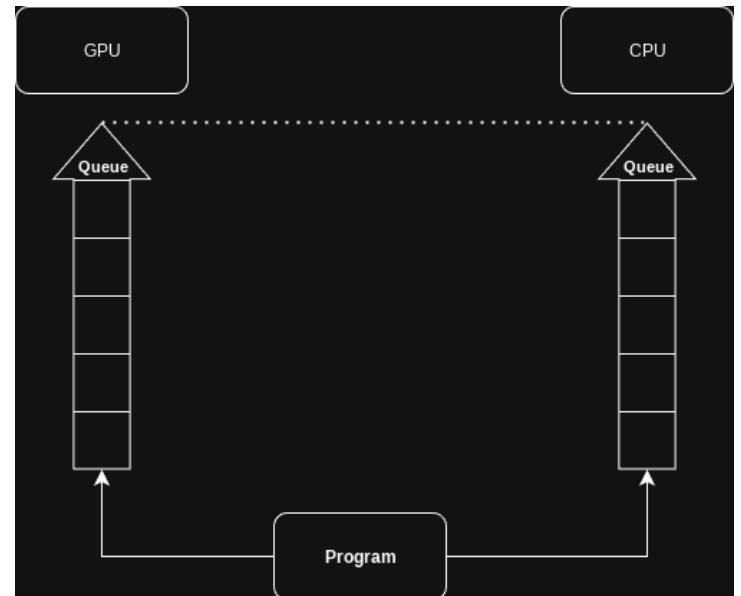


The basic platform and runtime APIs in OpenCL (using C)



Command-Queues

- A **command queue** is a host-side object that accepts commands and submits them in order to a specific device.
- **Asynchronous by Default:** Command submission is non-blocking. The host can queue up a sequence of operations and continue with its own work.
- **Execution Models**
 - *In-Order Queue*: They execute in order they are enqueued
 - *Out-of-Order Queue*: The device driver is free to reorder commands



Buffers and Memory objects

- **Abstracting Device Memory:** A `cl_mem` object is a handle to a reference-counted region of memory on the device (`__global` memory).
- Primary Type: **Buffer Objects**
 - Represents a simple, linear block of **bytes**—essentially a "C array on the device."
 - Naming Convention: A common and highly recommended practice is to prefix host arrays with `h_` and device buffers with `d_`.

```
// On the Host (CPU RAM)
float* h_input = (float*)malloc(bytes);

// On the Device (GPU VRAM) - a handle
cl_mem d_input = clCreateBuffer(context, flags, bytes, NULL,
&err);
```

Creating Buffers: Defining Access Intent

- The `flags` argument in `clCreateBuffer` is critical. It tells the driver how the memory will be used, allowing for significant optimization.
- **Kernel Access:**
 - `CL_MEM_READ_ONLY`: Kernel will only read from this buffer.
 - `CL_MEM_WRITE_ONLY`: Kernel will only write to this buffer.
 - `CL_MEM_READ_WRITE`: Kernel can read and write.
- **Host Initialization:**
 - `CL_MEM_COPY_HOST_PTR`: Allocate device memory and copy data from the provided host pointer during creation.

Data Transfer

- After creation, you move data using explicit commands:

```
// Copy FROM host TO device (blocking call)
clEnqueueWriteBuffer(queue, d_input, CL_TRUE, 0, bytes, h_input, 0, NULL, NULL);

// Copy FROM device TO host (blocking call)
clEnqueueReadBuffer(queue, d_output, CL_TRUE, 0, bytes, h_output, 0, NULL,
NULL);
```

Kernel Programming

- **C syntax (C99)**
- **`__kernel`**: The Entry Point
- **`__global / __constant`**: Pointers passed as arguments must have an address space qualifier. This explicitly defines them as living in device memory.
- Work-Item
Functions: `get_global_id(dim)`, `get_local_id(dim)`, `get_group_id(dim)`, `get_global_size(dim)`.
- These are the intrinsic functions that connect your code to its position in the N-dimensional problem grid, enabling data-parallel computation.
- Core Restriction: **No Recursion**
 - GPU hardware lacks the stack management for deep recursion
 - Algorithms must be re-formulated iteratively, often using a "work queue" pattern.
- **No Function Pointers** (as kernel arguments)

Explicitly Managing Data Locality

- **Address Spaces: YOU EXPLICITLY CONTROL DATA PLACEMENT**
 - **`__global`**: Main device memory (DRAM). High-latency, high-bandwidth.
 - **`__local`**: On-chip SRAM. Low-latency, shared within a work-group.
 - **`__constant`**: Read-only region of global memory, often with a dedicated cache.
 - **`__private`**: Per-work-item registers. Fastest access.

Exploiting Hardware SIMD/SIMT

- Vector Data Types in OpenCL
 - Types like `float4`, `char16` map directly to wide hardware execution units.
 - **Performance Rule:** Always prefer vector operations (`float4 a = b * c;`) over scalar loops (`for(i=0; i<4) a[i]=b[i]*c[i];`).
 - Use `vloadN` and `vstoreN` for efficient, aligned reads/writes of vector types from global memory.
 - `convert_type()`: Efficient, hardware-accelerated type conversions (e.g., `convert_float4(my_int4)`).

LU Decomposition Example

- Transform a matrix into the product of a lower triangular and upper triangular matrix.
Used to solve a linear system of equations

1. The Problem: A System of Linear Equations

We start with a system of equations:

$$x_1 + 2x_2 + 3x_3 = 6$$

$$2x_1 + 5x_2 + 8x_3 = 15$$

$$3x_1 + 8x_2 + 14x_3 = 25$$

Matrix form:

$$\underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 2 & 5 & 8 \\ 3 & 8 & 14 \end{bmatrix}}_{\mathbf{A}} \underbrace{\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}}_{\mathbf{x}} = \underbrace{\begin{bmatrix} 6 \\ 15 \\ 25 \end{bmatrix}}_{\mathbf{b}}$$

2. Goal: Decompose A into L and U

$$A = \begin{bmatrix} 1 & 0 & 0 \\ l_{21} & 1 & 0 \\ l_{31} & l_{32} & 1 \end{bmatrix} \begin{bmatrix} u_{11} & u_{12} & u_{13} \\ 0 & u_{22} & u_{23} \\ 0 & 0 & u_{33} \end{bmatrix}$$

Elimination steps to find U:

$$\begin{bmatrix} 1 & 2 & 3 \\ 2 & 5 & 8 \\ 3 & 8 & 14 \end{bmatrix} \xrightarrow{R_2 \rightarrow R_2 - 2R_1} \xrightarrow{R_3 \rightarrow R_3 - 3R_1} \begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix}$$

Multipliers become L:

$$L = \begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 2 & 5 \end{bmatrix} \xrightarrow{R_3 \rightarrow R_3 - 2R_2} \underbrace{\begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix}}_{U}$$

Solving the System

We replace A with LU:

$$L(Ux) = b, \quad y = Ux$$

Step 3a: Solve Ly = b

$$\begin{bmatrix} 1 & 0 & 0 \\ 2 & 1 & 0 \\ 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = \begin{bmatrix} 6 \\ 15 \\ 25 \end{bmatrix}$$

$$y_1 = 6, \quad y_2 = 3, \quad y_3 = 1$$

Step 3b: Solve Ux = y

$$\begin{bmatrix} 1 & 2 & 3 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 6 \\ 3 \\ 1 \end{bmatrix}$$

$$x_3 = 1, \quad x_2 = 1, \quad x_1 = 1$$

Solution: $x = [1, 1, 1]^T$

Sequential Code in C

- The C code focuses exclusively on the $\mathbf{A} \rightarrow \mathbf{LU}$, decomposition part

- LU Decomposition: $O(N^3)$
- Substitution: $O(N^2)$
- $O(N^3) > O(N^2)$
 - For $N = 1024$:
 - $N^2 = 1,048,576$
 - $N^3 = 1,073,741,824$

Any time spent in
the $O(N^2)$ solvers is just "noise"
in the measurement for large
matrices.

```
int lu_decompose(float A[N][N]) {
    for (int k = 0; k < N; k++) {
        for (int j = k; j < N; j++) {
            float sum = 0.0f;
            for (int p = 0; p < k; p++) {
                sum += A[k][p] * A[p][j];
            }
            A[k][j] = A[k][j] - sum;
        }

        if (fabs(A[k][k]) < 1e-9) { // Use a small epsilon for floating
point
            fprintf(stderr, "Matrix is singular!\n");
            return -1;
        }

        float pivot_inv = 1.0f / A[k][k];
        for (int i = k + 1; i < N; i++) {
            float sum = 0.0f;
            for (int p = 0; p < k; p++) {
                sum += A[i][p] * A[p][k];
            }
            A[i][k] = (A[i][k] - sum) * pivot_inv;
        }
    }

    return 0;
}
```

Benchmark Single Core

Case	MFLOPS	
	*CPU	GPU
Sequential C (not OpenCL)	714.8	N/A

$$\text{MFLOPS} = \left(\frac{2}{3} * N^3 \right) / \left(\text{Time_in_seconds} * 10^6 \right)$$

*Ryzen 9 3900X (12 core 4.2GHz)

Decomposing the Sequential Loops

```
int lu_decompose(float A[N][N]) {
    for (int k = 0; k < N; k++) {
        for (int j = k; j < N; j++) {
            float sum = 0.0f;
            for (int p = 0; p < k; p++) {
                sum += A[k][p] * A[p][j];
            }
            A[k][j] = A[k][j] - sum;
        }
    }

    if (fabs(A[k][k]) < 1e-9) { // Use a small epsilon for floating
        point
        fprintf(stderr, "Matrix is singular!\n");
        return -1;
    }

    float pivot_inv = 1.0f / A[k][k];
    for (int i = k + 1; i < N; i++) {
        float sum = 0.0f;
        for (int p = 0; p < k; p++) {
            sum += A[i][p] * A[p][k];
        }
        A[i][k] = (A[i][k] - sum) * pivot_inv;
    }
}

return 0;
}
```

Parallelizable Work

Remove outer loops and set work-item co-ordinates

OpenCL Kernel (1/2)



```
__kernel void lu_row_kernel(__global float* A, const int k, const int n)
{
    // Each work-item computes one element 'j' in the k-th row of U
    int j = k + get_global_id(0); ← Maps work-item ID to a column index
    float sum = 0.0f;
    for (int p = 0; p < k; p++) {
        sum += A[k * n + p] * A[p * n + j];
    }
    A[k * n + j] = A[k * n + j] - sum; ← Each work-item writes one result
}
```

OpenCL Kernel (2/2)



```
__kernel void lu_column_kernel(__global float* A, const int k, const int n)
{ // Each work-item computes one element 'i' in the k-th column of L
    int i = k + 1 + get_global_id(0);

    // This kernel can ONLY run AFTER the k-th row kernel is complete
    // because it needs the updated pivot value A[k*n + k].
    float pivot_inv = 1.0f / A[k * n + k];
EVERY work-item reads this!

    float sum = 0.0f;
    for (int p = 0; p < k; p++) {
        sum += A[i * n + p] * A[p * n + k];
    }
    A[i * n + k] = (A[i * n + k] - sum) * pivot_inv;
}
```

Host Side (C++) 1: Setup the program

```
// 1. Discover available platforms (e.g., NVIDIA, Intel)
clGetPlatformIDs(...);
cl_platform_id platform = ...; // User selects a platform

// 2. Discover devices on that platform (e.g., a specific GPU)
clGetDeviceIDs(platform, CL_DEVICE_TYPE_ALL, ...);
cl_device_id device = ...; // User selects a device

// 3. Create a context to manage all OpenCL objects
cl_context context = clCreateContext(NULL, 1, &device, ...);

// 4. Create a command queue to send work to the device
cl_command_queue queue = clCreateCommandQueueWithProperties(context, device,
...);
```

Host Side Part 2: Preparing Kernels and Data

```
// 1. Load the kernel source code from the .cl file
std::string kernel_source = read_kernel_file("lu_kernels.cl");
const char* source_ptr = kernel_source.c_str();

// 2. Create and build the OpenCL program
cl_program program = clCreateProgramWithSource(context, 1, &source_ptr,
err); clBuildProgram(program, ...); // JIT compilation for the target GPU!
// ... (Error checking clGetProgramBuildInfo is critical here) ...

// 3. Extract the individual kernel functions from the program
cl_kernel row_kernel = clCreateKernel(program, "lu_row_kernel", ...);
cl_kernel col_kernel = clCreateKernel(program, "lu_column_kernel", ...);

// 4. Allocate memory buffer on the GPU device
cl_mem d_A = clCreateBuffer(context, CL_MEM_READ_WRITE,
                           sizeof(float) * N * N, NULL, ...);

// 5. Copy our input matrix from host RAM to GPU VRAM
clEnqueueWriteBuffer(queue, d_A, CL_TRUE, 0,
                     sizeof(float) * N * N, h_A_working, ...);
```

Host Side Part 3: The Execution Loop



```
// Set the constant arguments for the kernels (device buffer and
clSetKernelArg(row_kernel, 0, sizeof(cl_mem), &d_A);
clSetKernelArg(row_kernel, 2, sizeof(int), &matrix_size);
// (repeat for col_kernel...)
```



```
for (int k = 0; k < N; k++) {
    cl_event row_event; // Used to chain commands

    // Set the changing argument 'k' for this iteration
    clSetKernelArg(row_kernel, 1, sizeof(int), &k);
    size_t global_work_size_row = N - k;
    if (global_work_size_row > 0) {
        clEnqueueNDRangeKernel(queue, row_kernel, 1, NULL, &global_work_size_row,
                               NULL, 0, NULL, &row_event); // Launch kernel
    }

    // Set 'k' for the column kernel and define the dependency
    clSetKernelArg(col_kernel, 1, sizeof(int), &k);
    size_t global_work_size_col = N - k - 1;
    if (global_work_size_col > 0) {
        clEnqueueNDRangeKernel(queue, col_kernel, 1, NULL, &global_work_size_col,
                               NULL, 1, &row_event, NULL); // Launch AFTER
    }
}
```

Host Side Part 4: Retrieving Results and Cleaning

```
// 1. Block the host thread until all queued commands on the GPU have
// finished including all 2*N kernel launches and their dependencies.
clFinish(queue);
clock_t end_ocl = clock(); // Stop the timer

// 2. Create a buffer in host memory for the result
float* h_Result = new float[N * N];

// 3. Copy the final decomposed matrix from GPU VRAM back to host RAM
err = clEnqueueReadBuffer(queue, d_A, CL_TRUE, 0,
                           sizeof(float) * N * N, h_Result, ...);

// 4. Now that the data is on the host, we can verify its correctness
if (h_Result) {
    verify_lu_decomposition(h_A_original, h_Result);
}

// 5. Clean up all OpenCL objects (clReleaseMemObject, clReleaseKernel, etc.)
// ..
```

Host programs can be “ugly”

```
1 int main() {
2     srand(time(NULL));
3
4     std::cout << "Setting up host data for a " << N << "x" << N << " matrix..." << std::endl;
5     float* h_A_original = new float[N * N];
6     initialize_matrix(h_A_original);
7     float* h_A_working = new float[N * N];
8     memcpy(h_A_working, h_A_original, sizeof(float) * N * N);
9
10    cl_int err;
11
12    cl_uint num_platforms;
13    clGetPlatformIDs(0, NULL, &num_platforms);
14    std::vector<cl_platform_id> platforms(num_platforms);
15    clGetPlatformIDs(num_platforms, platforms.data(), NULL);
16
17    std::cout << "Available Platforms:" << std::endl;
18    for (cl_uint l = 0; l < num_platforms; ++l) {
19        char platform_name[128];
20        clGetPlatformInfo(platforms[l], CL_PLATFORM_NAME, sizeof(platform_name), platform_name, NULL);
21        std::cout << " " << l << ":" << platform_name << std::endl;
22    }
23    std::cout << "Select a platform: ";
24    unsigned int plat_idx;
25    std::cin >> plat_idx;
26    cl_platform_id platform = platforms[plat_idx];
27
28    cl_uint num_devices;
29    clGetDeviceIDs(platform, CL_DEVICE_TYPE_ALL, 0, NULL, &num_devices);
30    std::vector<cl_device_id> devices(num_devices);
31    clGetDeviceIDs(platform, CL_DEVICE_TYPE_ALL, num_devices, devices.data(), NULL);
32
33    std::cout << "Available Devices:" << std::endl;
34    for (cl_uint l = 0; l < num_devices; ++l) {
35        char device_name[128];
36        clGetDeviceInfo(devices[l], CL_DEVICE_NAME, sizeof(device_name), device_name, NULL);
37        std::cout << " " << l << ":" << device_name << std::endl;
38    }
39    std::cout << "Select a device: ";
40    unsigned int dev_idx;
41    std::cin >> dev_idx;
42    cl_device_id device = devices[dev_idx];
43
44    cl_context context = clCreateContext(NULL, 1, &device, NULL, NULL, &err); CL_CHECK(err);
45
46    cl_command_queue queue = clCreateCommandQueueWithProperties(context, device, 0, &err); CL_CHECK(err);
47
48    std::string kernel_source = read_kernel_file("lu_kernels.cl");
49    const char* source_ptr = kernel_source.c_str();
50    cl_program program = clCreateProgramWithSource(context, 1, &source_ptr, NULL, &err); CL_CHECK(err);
51    err = clBuildProgram(program, 1, &device, NULL, NULL, NULL);
52    if (err != CL_SUCCESS) {
53        size_t log_size;
54        clGetProgramBuildInfo(program, device, CL_PROGRAM_BUILD_LOG, 0, NULL, &log_size);
55        std::vector<char> log(log_size);
56        clGetProgramBuildInfo(program, device, CL_PROGRAM_BUILD_LOG, log_size, log.data(), NULL);
57        std::cerr << "Build Log:\n" << log.data() << std::endl;
58        exit(1);
59    }
60
61    cl_kernel row_kernel = clCreateKernel(program, "lu_row_kernel", &err); CL_CHECK(err);
62    cl_kernel col_kernel = clCreateKernel(program, "lu_column_kernel", &err); CL_CHECK(err);
63    cl_mem d_A = clCreateBuffer(context, CL_MEM_READ_WRITE, sizeof(float) * N * N, NULL, &err);
64    CL_CHECK(err);
65
```

```
std::cout << "\n--- Running Native OpenCL Implementation ---" << std::endl;
3
4    err = clEnqueueWriteBuffer(queue, d_A, CL_TRUE, 0, sizeof(float) * N * N, h_A_working, 0, NULL, NULL);
5    CL_CHECK(err);
6    int matrix_size = N;
7    clSetKernelArg(row_kernel, 0, sizeof(cl_mem), &d_A);
8    clSetKernelArg(row_kernel, 2, sizeof(int), &matrix_size);
9    clSetKernelArg(col_kernel, 0, sizeof(cl_mem), &d_A);
10   clSetKernelArg(col_kernel, 2, sizeof(int), &matrix_size);
11
12   clock_t start_ocl = clock();
13
14   for (int k = 0; k < N; k++) {
15       cl_event row_event; // Event to track completion of the row kernel
16
17       // Set arguments for the row kernel
18       clSetKernelArg(row_kernel, 1, sizeof(int), &k);
19       size_t global_work_size_row = N - k;
20       if (global_work_size_row > 0) {
21           // Enqueue the row kernel and associate it with row event
22           err = clEnqueueNDRangeKernel(queue, row_kernel, 1, NULL, &global_work_size_row, NULL, 0, &row_event);
23           CL_CHECK(err);
24
25       }
26
27       // Set arguments for the column kernel
28       clSetKernelArg(col_kernel, 1, sizeof(int), &k);
29       size_t global_work_size_col = N - k - 1;
30
31       if (global_work_size_col > 0) {
32           // Enqueue the column kernel.
33           // The '1' argument tells this kernel to wait for row_event to be signaled.
34           err = clEnqueueNDRangeKernel(queue, col_kernel, 1, NULL, &global_work_size_col, NULL, 1, &row_event, NULL);
35           CL_CHECK(err);
36
37       }
38
39       // It's good practice to release events when they are no longer needed in the wait list
40       // to avoid resource leaks, especially in more complex scenarios.
41       if (global_work_size_row > 0) {
42           clReleaseEvent(row_event);
43       }
44
45   }
46
47   // Wait for all enqueued commands to finish before stopping the timer.
48   clFinish(queue);
49
50   clock_t end_ocl = clock();
51
52   double time_ocl = ((double)(end_ocl - start_ocl)) / CLOCKS_PER_SEC;
53   double mflops_ocl = (2.0 * (3.0) * pow(N, 3)) / (time_ocl * 1e6);
54   printf("Native OpenCL Time: %f seconds\n", time_ocl);
55   printf("Native OpenCL Performance: %.2f MFLOPS\n", mflops_ocl);
56
57   // READ RESULT BACK
58   std::cout << "Reading result from GPU for verification..." << std::endl;
59   float* h_Result = new float[N * N];
60   err = clEnqueueReadBuffer(queue, d_A, CL_TRUE, 0, sizeof(float) * N * N, h_Result, 0, NULL, NULL); CL_CHECK(err);
61
62   // Verify Results
63   if (h_Result) {
64       verify_lu_decomposition(h_A_original, h_Result);
65   }
66
67   // MUST DO CLEANUP REMEMBER....!!!!
68   clReleaseMemObject(d_A);
69   clReleaseKernel(row_kernel);
70   clReleaseKernel(col_kernel);
71   clReleaseProgram(program);
72   clReleaseCommandQueue(queue);
73   clReleaseContext(context);
74   delete[] h_A_original;
75   delete[] h_A_working;
76   if (h_Result) delete[] h_Result;
77
78   return 0;
79
```

A LOT OF
BOILERPLATE CODE!

Benchmark Naive OpenCL Kernels

Case	MFLOPS	
	CPU	*GPU
Sequential C (not OpenCL)	714.8	N/A
Naive Kernels (OpenCL)	376.62	20,991.25

Maximum Absolute Error: 1.19e-07

* Nvidia 4090

The Next Step: Optimizing the Kernels

Our first OpenCL implementation was 28.5x faster than the CPU, which is a great start. But it's far from the theoretical peak performance of the GPU.

- **Problem 1: Kernel Launch Overhead**
 - For an N=2048 matrix, the host launches $2 * 2048 = 4096$ kernels.
 - Each kernel launch has a small but significant CPU and driver overhead. This adds up.
- **Problem 2: Memory Bandwidth Saturation**
 - The kernels are memory-bound. They spend most of their time reading from and writing to slow global VRAM.
 - There is almost no data reuse. A value is read from global memory, used in one calculation, and then discarded. The key to GPU performance is to read data once into fast memory and reuse it many times.

Algorithm Change: Block LU Decomposition

- ▶ Partition a large matrix M into four blocks.
- ▶ Goal: factor M into block lower L and block upper U .

$$M = \begin{pmatrix} A & B \\ C & D \end{pmatrix}.$$

Main Advantage: For large matrices, computers can work on the different blocks at the same time

Target factorization:

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} L_{11} & 0 \\ L_{21} & L_{22} \end{pmatrix} \begin{pmatrix} U_{11} & U_{12} \\ 0 & U_{22} \end{pmatrix}.$$

Assumption

Top-left block A is invertible.

Step 2: Solving for the Matrix Blocks

Multiply the block factors and match blocks:

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} L_{11}U_{11} & L_{11}U_{12} \\ L_{21}U_{11} & L_{21}U_{12} + L_{22}U_{22} \end{pmatrix}.$$

1. From $A = L_{11}U_{11}$: LU-decompose A to get L_{11}, U_{11} .
2. From $B = L_{11}U_{12}$: $U_{12} = L_{11}^{-1}B$.
3. From $C = L_{21}U_{11}$: $L_{21} = C U_{11}^{-1}$.

Step 3: The Schur Complement

From the bottom-right block,

$$D = L_{21} U_{12} + L_{22} U_{22} \quad \Rightarrow \quad L_{22} U_{22} = D - L_{21} U_{12}.$$

Substitute L_{21} and U_{12} :

$$L_{22} U_{22} = D - (C U_{11}^{-1})(L_{11}^{-1} B) = D - C(L_{11} U_{11})^{-1} B.$$

Schur Complement

Define $S := D - C A^{-1} B$ (since $A = L_{11} U_{11}$ invertible).

LU-decompose S to obtain L_{22}, U_{22} .

A Hybrid CPU+GPU Block Algorithm

- Step 1: Decompose Diagonal Panel ($A \rightarrow LU$)
 - This is a small LU decomposition (32x32). It's serial and complex.
 - Best suited for the CPU. We'll *copy the panel from GPU to CPU*, decompose it, and *copy it back*.
- Step 2: Update Row/Column Panels
 - These are triangular matrix solves (**TRSM**). They are highly parallel.
- Step 3: Update Trailing Matrix ($D - CA^{-1}B$)
 - This is a large General Matrix-Matrix Multiplication (**GEMM**). This is the most computationally expensive part (**O(N³)**).
- The Big Idea: This approach trades thousands of small, inefficient operations for a few, large, highly optimized ones (TRSM and GEMM) that can saturate the GPU's computational resources.

Synchronization

On a GPU, we have two levels of synchronization:

- **Within a Work-Group (Easy & Fast):**
 - We can use a `barrier()` to create a rendezvous point.
 - When a work-item hits a `barrier()`, it pauses until every other work-item in its group also reaches that barrier.
 - This is essential for safely sharing data through `__local` memory. A common pattern is:
 - All work-items cooperatively load data into `__local` memory.
 - `barrier(CLK_LOCAL_MEM_FENCE); // Wait for all loads to complete.`
 - All work-items can now safely read the shared data for computation.
- **Across Work-Groups (Hard & Explicit):**
 - There is **NO** direct way for different work-groups to communicate or synchronize within a single kernel launch.
 - The hardware scheduler can execute work-groups in any order, concurrently or serially.
- Solution: If you need global synchronization, you must end the current kernel and launch a new one. The host API (using events) becomes your global synchronization tool.

Be careful with barriers

- A `barrier()` must be encountered by all work-items in a group, or none at all. Placing a barrier inside a branch is a common and dangerous source of bugs.

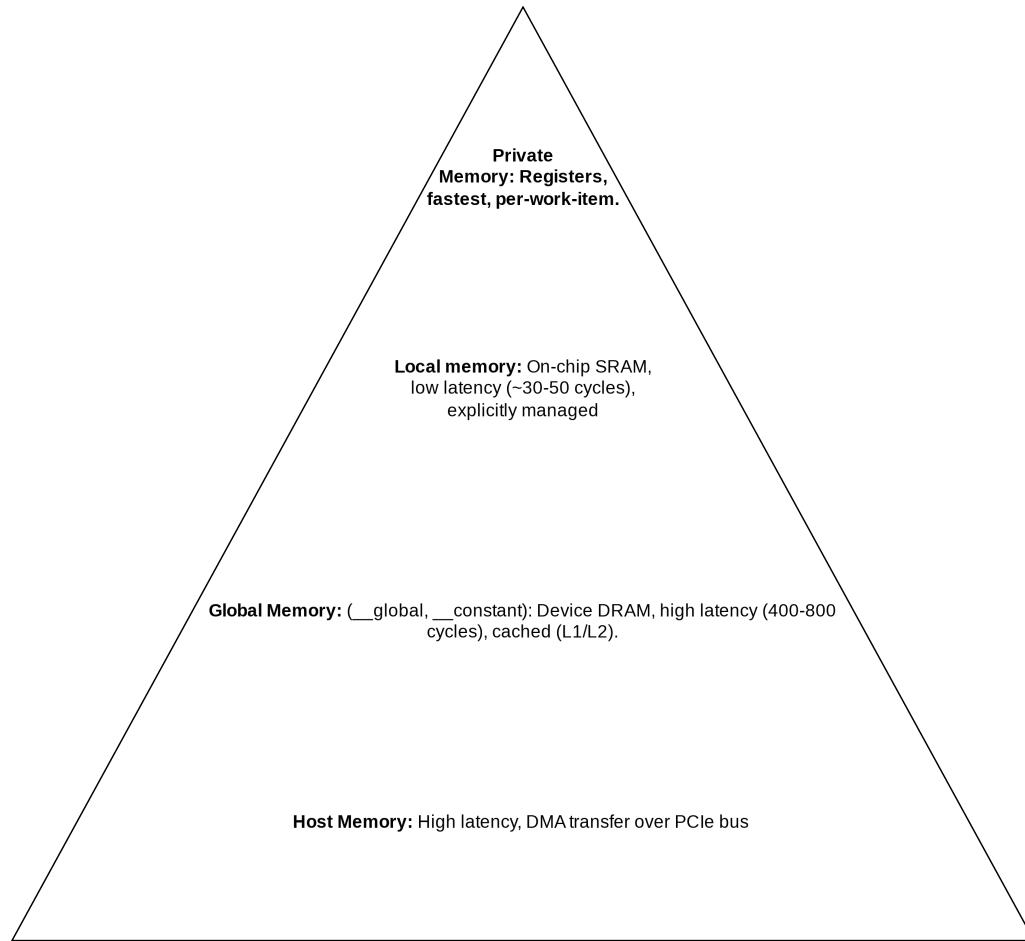
This will cause a deadlock and hang your GPU:

```
1 // WORK-GROUP HANGS!
2 if (get_local_id(0) == 0) {
3     // Only ONE work-item enters this branch...
4     // ...but it will wait FOREVER for the others to
5     // arrive.
6     barrier(CLK_LOCAL_MEM_FENCE);
7 }
```

This is correct:

```
1 int some_value = ...;
2 // ALL work-items check the condition.
3 if (condition_is_uniform_for_the_group)
4 {   // ... do work ...
5     barrier(CLK_LOCAL_MEM_FENCE);
6     // ... do more work ...
7 }
```

Memory Model Again :)

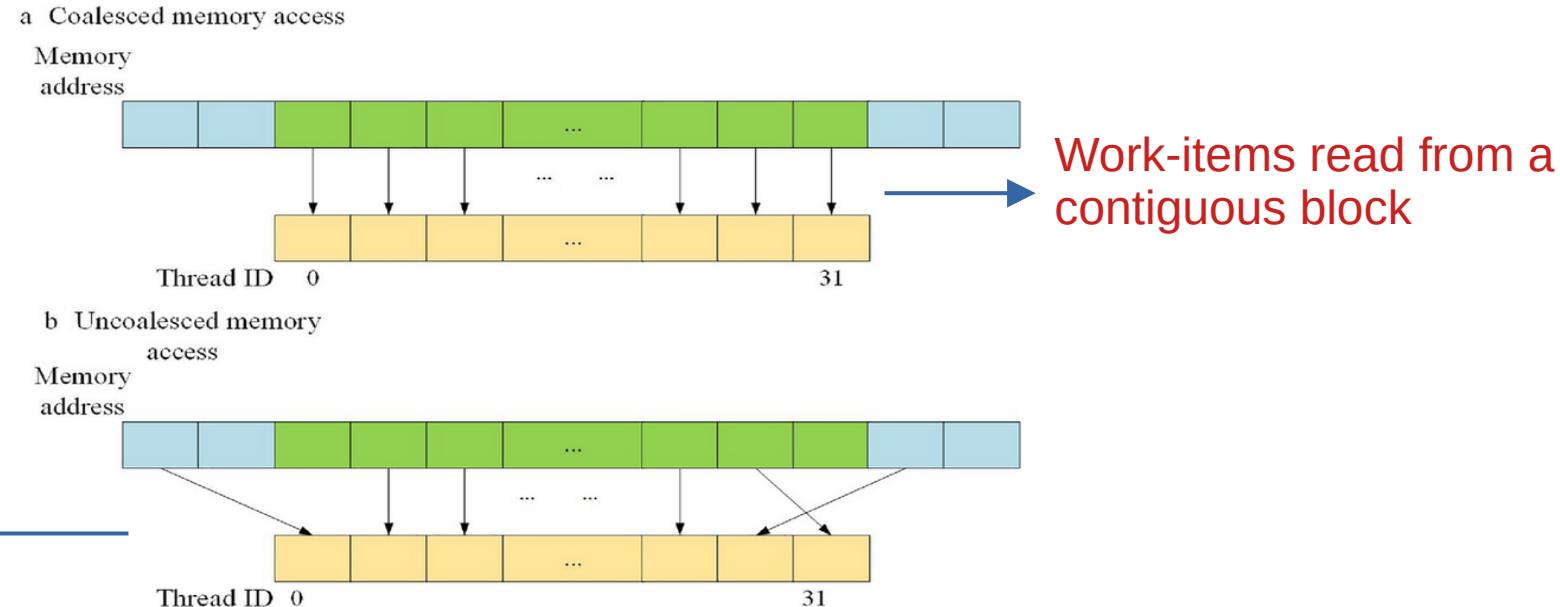


Caching Panels in __local Memory

- **The Problem:** The naive kernels were memory-bound. They spent all their time waiting for data to arrive from slow global VRAM.
- USE the __local memory. The strategy is:
 - Work-items in a group **cooperatively load** a block of data from global to local memory.
 - They **synchronize** to ensure the load is complete.
 - They perform all calculations using the **fast** local memory copy.

Coalesced Loads for Maximum Bandwidth

- **The Problem:** Just *using __local* memory isn't enough. We must fill it from global memory as fast as possible. The GPU's memory bus is extremely wide, and we need to use all of it.
- **The Solution:** **Memory Coalescing**. When adjacent work-items (e.g., IDs 0, 1, 2, 3...) access adjacent memory locations, the hardware combines these individual reads into a single, wide, and highly efficient transaction.



Registers & The Strategic Uncoalesced Read

- **The Problem:** Even local memory has some latency. For data that is used extremely heavily by a single work-item, we want it in the absolute fastest storage: **private registers**.
- **The Solution:** The OpenCL compiler is smart. A small, compile-time-sized array declared inside a kernel will often be mapped directly to registers.

```
1 // 1. This small array will be stored in ultra-fast private registers
2 float x[BLOCK_SIZE];
3
4 // 2. Load the column data from global memory into registers
5 for(int i = 0; i < BLOCK_SIZE; i++) {
6     // This access is strided by N. IT IS NOT COALESCED and looks very slow!
7     x[i] = A[(k + i) * n + global_j];
8 }
9
10 // 3. The Payoff: Perform the entire substitution using only registers and
11 // local memory.
12 for (int i = 0; i < BLOCK_SIZE; i++) {
13     float sum = 0.0f;
14     for (int p = 0; p < i; p++) {
15         // x[p] is read from a register. L_panel is from local memory. ALL
16 FAST!     sum += L_panel[i][p] * x[p];
17     }
18     x[i] = x[i] - sum; // Write to a register
19 }
```

BLOCK_SIZE is typically something like 32

Sizing Work for High Occupancy

- The goal is to keep the Compute Units Busy!!!
 - Have more active **Subgroups** than the GPU can execute at once
 - If a singular subgroup stalls (while waiting for memory)
 - Scheduler instantly swaps in another that is ready to compute.
- **Occupancy:** the ratio of active Subgroups on a CU to the maximum Subgroups that CU can support.
- The Trade-off: There is a finite amount of resources (registers, __local memory) per CU.
 - Large Work-Groups:
 - Pro: Can enable more data sharing in __local memory.
 - Con: Consume more resources
 - Small Work-Groups:
 - Pro: Consume fewer resources,
 - Con: May not be able to cache enough data in __local memory to be effective.

TEST AND PROFILE !!!

Work-Item Divergence

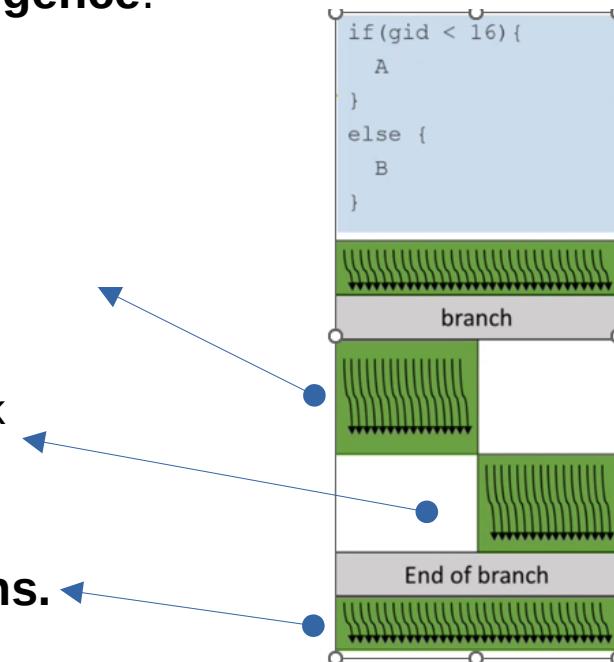
The Concept: All work-items in a Subgroup must execute the same instruction at the same time. What happens when the code has an if/else statement?

- If some work-items in a Subgroup need to enter the if block and others need to enter the else block, this is called **Branch Divergence**.

First execute the if path, disabling all the work items that didn't take it

Then, do the else path, disabling the work-items that took the if path.

The total time taken is the sum of both paths.



TSRM Kernels

Using `__local` Memory as a User-
Managed Cache

Synchronize

Using Private Memory

```
1 __kernel void lu_row_update_kernel(__global float* A, const int k, const int n)
2 {
3     __local float L_panel[BLOCK_SIZE][BLOCK_SIZE];
4
5     int local_id = get_local_id(0);
6
7     // Cooperative load of the L_panel. Every work-item in a group helps.
8     if (local_id < BLOCK_SIZE) {
9         for (int j = 0; j < BLOCK_SIZE; j++) {
10            L_panel[local_id][j] = A[(k + local_id) * n + (k + j)];
11        }
12    }
13    barrier(CLK_LOCAL_MEM_FENCE);
14
15    // Each work-item is responsible for one column of the output panel.
16    int global_j = k + BLOCK_SIZE + get_global_id(0);
17    if (global_j >= n) return;
18
19    // Load column into private memory
20    float x[BLOCK_SIZE];
21    for(int i = 0; i < BLOCK_SIZE; i++) {
22        x[i] = A[(k + i) * n + global_j];
23    }
24
25    // Perform forward substitution on private data
26    for (int i = 0; i < BLOCK_SIZE; i++) {
27        float sum = 0.0f;
28        for (int p = 0; p < i; p++) {
29            sum += L_panel[i][p] * x[p];
30        }
31        // Unit diagonal for L is assumed
32        x[i] = x[i] - sum;
33    }
34
35    // Write result back to global memory
36    for(int i = 0; i < BLOCK_SIZE; i++) {
37        A[(k + i) * n + global_j] = x[i];
38    }
39 }
```

GEMM Kernel

Filling the Cache Efficiently

Checks for matrices whose size isn't a perfect multiple of BLOCK_SIZE.
This is a necessary source of divergence.

```
1 __kernel void gemm_update_kernel(__global float* A, const int k, const int n) {
2     int local_row = get_local_id(1);
3     int local_col = get_local_id(0);
4     int group_row = get_group_id(1);
5     int group_col = get_group_id(0);
6
7     // Top-left corner of the destination tile this work-group is responsible for
8     int dest_tile_row_start = k + BLOCK_SIZE + group_row * BLOCK_SIZE;
9     int dest_tile_col_start = k + BLOCK_SIZE + group_col * BLOCK_SIZE;
10
11    __local float L_tile[BLOCK_SIZE][BLOCK_SIZE];
12    __local float U_tile[BLOCK_SIZE][BLOCK_SIZE];
13
14    // Load the tile from the L panel (A_ik)
15    L_tile[local_row][local_col] = A[(dest_tile_row_start + local_row) * n + (k +
16                                         local_col)];
16
17    // Load the tile from the U panel (A_kj)
18    U_tile[local_row][local_col] = A[(k + local_row) * n + (dest_tile_col_start +
19                                         local_col)];
20
21    barrier(CLK_LOCAL_MEM_FENCE);
22
23    // Standard local memory GEMM computation
24    float accumulator = 0.0f;
25    for (int p = 0; p < BLOCK_SIZE; ++p) {
26        accumulator += L_tile[local_row][p] * U_tile[p][local_col];
27    }
28
29    // Calculate the final global memory address to write to
30    int final_global_row = dest_tile_row_start + local_row;
31    int final_global_col = dest_tile_col_start + local_col;
32
33    // Boundary check and write result
34    if (final_global_row < n && final_global_col < n) {
35        A[final_global_row * n + final_global_col] -= accumulator;
36    }
```

Host Code 1: The Hybrid Main Loop

The main loop now advances in large BLOCK_SIZE steps.

We use clEnqueueRead/WriteBufferRect to transfer only the small diagonal panel, minimizing slow PCIe traffic.



```
1 // The loop now iterates over BLOCKS, not elements
2 for (int k = 0; k < N; k += BLOCK_SIZE) {
3     // === CPU WORK PHASE ===
4
5     // 1. Fetch the diagonal panel from the GPU where it currently
6     // lives
7     clEnqueueReadBufferRect(queue, d_A, CL_TRUE, /* src_origin */,
8                             /* dst_origin */, ..., h_panel, ...);
9
10    // 2. Perform the small, serial LU decomposition on the CPU
11    cpu_panel_decompose(h_panel, BLOCK_SIZE);
12
13    // 3. Send the factorized panel back to the GPU
14    clEnqueueWriteBufferRect(queue, d_A, CL_TRUE, /* dst_origin */,
15                            /* src_origin */, ..., h_panel, ...);
16
17    // === GPU WORK PHASE (details on next slide) ===
18
19    // 4. Enqueue all the parallel GPU work for this step
20    enqueue_gpu_kernels_for_step(k);
21 }
22
23 // 5. After the loop queues EVERYTHING, block the host and wait for the
24 // GPU to finish.
25 clFinish(queue);
26
```

Host Code 2: Asynchronous Kernels with Events

- Defining a Dependency Graph
- Maximizing Concurrency
- Non-Blocking Host

```
1 // Inside the main loop, for a given 'k':
2 // An array to hold our synchronization objects
3 cl_event update_events[2];
4
5 // 1. Launch the row panel solver. It can start immediately.
6 //     The completion of this kernel will signal update_events[0].
7 err = clEnqueueNDRangeKernel(queue, row_update_kernel, ...,
8 &update_events[0]);
9 // 2. Launch the column panel solver. It can also start immediately
10 //    and run CONCURRENTLY with the row kernel.
11 //    Its completion will signal update_events[1].
12 err = clEnqueueNDRangeKernel(queue, col_update_kernel, ...,
13 &update_events[1]);
14 // 3. Launch the massive GEMM update kernel.
15 //     This is the crucial dependency: we tell this kernel it must WAIT for
16 //     the TWO events in our update_events array before it can start.
17 err = clEnqueueNDRangeKernel(queue, gemm_kernel, ..., 2, update_events, NULL);
18
19
```

Benchmark Optimised OpenCL Kernels

Case	MFLOPS	
	CPU	*GPU
Sequential C (not OpenCL)	714.8	N/A
Naive Kernels (OpenCL)	376.62	20,991.25
Block Hybrid (Synchronous)	-	1,339,874.37

Maximum Absolute Error ($|A - L^*U|$): 1.708984e-03

Approaching Peak Performance

Our block-hybrid model is already impressive, achieving over 1.3 TFLOPS. But there are still two remaining bottlenecks:

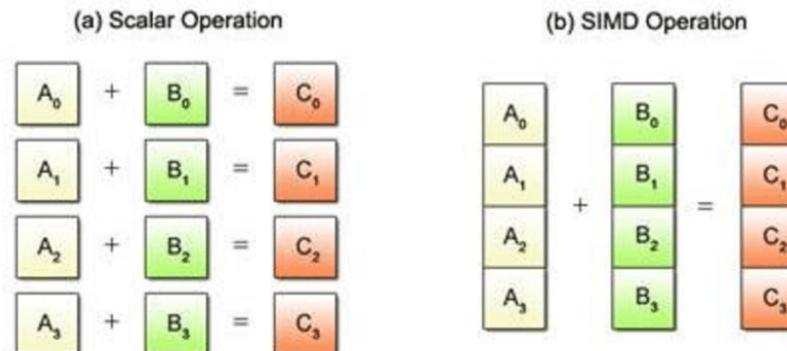
- Inside the Kernel: Our gemm_kernel still has a relatively low arithmetic intensity. Each work-item reads a lot from __local memory relative to the amount of math it does.
- On the Host: Our main loop is synchronous. The GPU waits for the CPU to process a panel, and the CPU waits for the GPU to finish its work. There is still "dead time" where one processor is idle.

To achieve the next level of performance, we must attack both of these issues simultaneously.

GPU: Vectorization (SIMD)

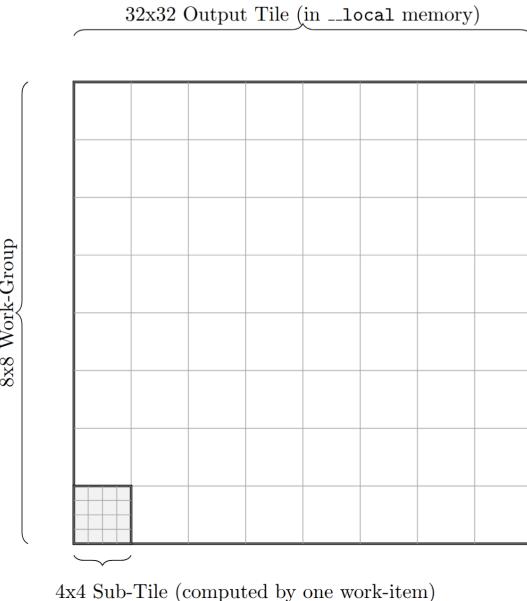
At their core, GPUs are SIMD (Single Instruction, Multiple Data) processors. Their hardware is not designed to operate on single numbers (scalars) but on short vectors of numbers all at once.

- Scalar Operation (Inefficient):
 - $\text{float } a = b * c;$ -> 1 Instruction, 1 Multiplication
- Vector Operation (Efficient):
 - In OpenCL, we can use built-in vector types like float4, float8, or float16.
 - $\text{float4 } a = b * c;$ -> 1 Instruction, 4 Multiplications!



Optimisation 2: Register Blocking

- **The Old Way:** In our previous gemm_kernel, each work-item in a 32x32 work-group computed one (i, j) output element.
- **The New Idea:** What if we make each thread responsible for a small 2D tile of the output? For a 32x32 output tile, we can use an 8x8 work-group where each of the 64 threads computes its own private 4x4 sub-tile.



GEMM Kernel 2

Unroll the loop, compiler
does this



```
1 #define TILE_DIM 4           // Each work-item handles a 4x4 tile
2 #define WGS 8                 // The Work-Group is 8x8 = 64 items
3
4 __kernel void gemm_update_kernel_register_blocked(...) {
5     // 1. POINTER VECTORIZATION: Treat global memory as a grid of
6     // float4's
7     __global float4* A_vec = (__global float4*)A;
8
9     // 2. REGISTER BLOCKING using a vector type.
10    // 'accum' is an array of FOUR float4's, holding all 16 results in
11    // registers.
12    float4 accum[TILE_DIM];
13    for(int i = 0; i < TILE_DIM; ++i) { accum[i] = (float4)(0.0f); }
14
15    // 3. VECTORIZED LOAD into __local memory for maximum bandwidth.
16    // Each work-item loads a float4 vector.
17    *(__local float4*)&L_tile[...] = A_vec[...];
18    barrier(CLK_LOCAL_MEM_FENCE);
19
20    // 4. VECTORIZED COMPUTATION in the main loop
21 #pragma unroll
22    for (int p = 0; p < BLOCK_SIZE; ++p) {
23        // Fetch a full float4 vector from the U tile
24        float4 u_vec = *(__local float4*)&U_tile[p][local_col *
25 TILE_DIM];
26        // Fetch 4 L values for the 4 rows this thread is responsible for
27        float l_vals[TILE_DIM];
28        for(int i=0; i<TILE_DIM; ++i) { l_vals[i] = L_tile[...][p]; }
29
30        // The PAYOFF: Fused-Multiply-Add on 4 floats at once.
31        // The scalar l_vals[i] is broadcast to a vector before the FMA.
32        for(int i=0; i<TILE_DIM; ++i) {
33            accum[i] = fma((float4)l_vals[i], u_vec, accum[i]);
34        }
35    }
36    // 5. VECTORIZED WRITE from registers back to global memory
37    A_vec[...] -= accum[i];
38 }
39
```

Host Code 1: Enabling True Asynchronicity with Pinned Memory

The Problem: Standard C++ new float[...] memory is "pageable." For the GPU to access it, the driver must first copy it to a "pinned," non-pageable staging buffer, and then DMA it to the GPU. This is an implicit, blocking copy.

- The Solution: Pinned (or Page-Locked) Memory.

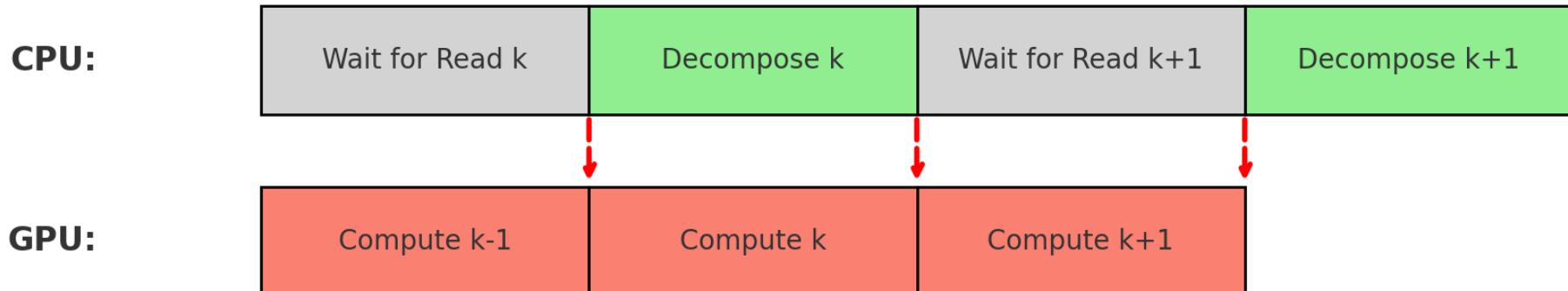
We can ask the OpenCL driver to allocate host memory that the GPU's DMA engine can access directly. This eliminates the staging copy and is the key to enabling non-blocking data transfers.

```
1 // Allocate a buffer on the host, but flag it so the GPU can see it directly
2 cl_mem pinned_buffer = clCreateBuffer(context, CL_MEM_ALLOC_HOST_PTR, size,
3 ...);
4 // Get a normal C++ pointer to this pinned memory region
5 float* h_panel = (float*)clEnqueueMapBuffer(queue, pinned_buffer, ...);
6
7 // Now, transfers using this h_panel pointer can be truly asynchronous.
8 clEnqueueReadBufferRect(..., h_panel, ...); // Can return immediately!
9
```

Host Code 2: The Pipelined Main Loop

To hide all latency, we use two command queues (one for transfers, one for compute) and completely restructure the main loop into a pipeline. While the CPU processes panel k, the GPU can be busy computing the kernels for panel k-1.

Asynchronous Pipelined Execution: CPU prepares k while GPU computes k-1



Host Code 2: The Pipelined Main Loop (Continued)

```
1 // Create TWO queues: one for transfers, one for compute
2 cl_command_queue compute_queue = ...;
3 cl_command_queue transfer_queue = ...;
4
5 for (int k = 0; k < N; k += BLOCK_SIZE) {
6     // 1. Asynchronously start reading panel 'k' from GPU to pinned memory.
7     // This is non-blocking and returns instantly.
8     clEnqueueReadBufferRect(transfer_queue, ..., CL_FALSE, ..., &read_event);
9
10    // 2. If this isn't the first loop, enqueue all the compute kernels for the
11    // PREVIOUS iteration (k-1). They will run while the CPU waits for step 3.
12    if (k > 0) {
13        // All these kernels are submitted to the COMPUTE queue and
14        // depend on the WRITE event from the previous iteration.
15        clEnqueueNDRangeKernel(compute_queue, ...);
16    }
17
18    // 3. This is the only host-side wait. The CPU pauses here until panel 'k'
19    // has finished transferring from the GPU.
20    clWaitForEvents(1, &read_event);
21
22    // 4. While the GPU is now busy on step 2, the CPU does its work on panel 'k'.
23    cpu_panel_decompose(h_panels[...], BLOCK_SIZE);
24
25    // 5. Asynchronously start writing the decomposed panel 'k' back to the GPU.
26    // This is also non-blocking. Save the 'write_event' for the next
27    // iteration.
28    clEnqueueWriteBufferRect(transfer_queue, ..., CL_FALSE, ..., &write_event);
29    last_write_event = write_event;
30 }
31 clFinish(compute_queue); // Final sync
32
```

The Results of Full Optimization

Case	MFLOPS	
	CPU	*GPU
Sequential C (not OpenCL)	714.8	N/A
Naive Kernels (OpenCL)	376.62	20,991.25
Block Hybrid (Synchronous)	-	1,339,874.37
Pipelined & Register-Blocked		3,044,024

Maximum Absolute Error ($|A - L^*U|$):1.870495e-01

Comparing GEMM Kernels

Case	TFLOPS	
	CPU	*GPU
Pipelined & Register-Blocked	N/A	18.8
CLBlast GEMM Kernel		41.3

Learning Resources

- [SimpleOpenCLSamples](#)
- C++ Concurrency in Action [Anthony Williams]
- <https://github.com/ProjectPhysX/FluidX3D>
- **Data Parallel C++**
- <https://sycl.tech/>
- [hands-on-opencl](#)