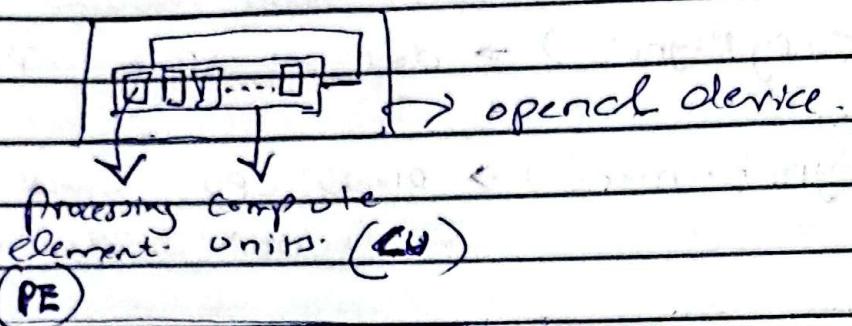


OPENCL (Open Compute language)

- One Host & one or more OpenCL devices
 - Each OpenCL Device is composed of one or more Compute Units.
 - Each Compute unit is divided into one or more Processing elements.



- Host & Kernel.

host Program (CPU Chaining program)

↓
load-Program C (load and initialize problem data on host) threads shared

↓
kernel C (executed on GPU (device))
(Host terminate by) Execute kernel

device parallel per
threads shared

device-
kernel start

Parallel

↓
Kernel end

host function
which waits Sync ()
until
Kernel execution
Complete
Synchronization

get_results () Kernel key results to
host with receiver
body by

Minimum 2 files

- 1 host program file (CPU code) .c or .cpp
- 1 Kernel file (GPU code) .cl

• OpenCL Platform Examples:-

CPU

GPU

- While CPU is single OpenCL - Every GPU is a device.
- Every CPU core is CU
- Every CU has PE
- If GPU uses multiple PE's are there.
- * - CPU app. host can directly use CPU or GPU off concurrently.
- Every GPU is a device.
- GPU has only streaming multiprocessor for CU fetching main.

CPU

GPU

~~Host~~ CPU core → PE ← Streaming multiprocessor

Thread

executed by Core

Thread Block

→ Streaming Multiprocessor.

Kernel code

→ Complete GPU unit.

- Private Memory → Different for every work-item (thread)

• Local Memory → Work-group has own threads main shared by all threads in group. 32 threads can access memory.

• Global / Constant Memory → accessible to all work-groups.

• Host Memory → CPU memory
Data moves from CPU to GPU & vice versa
(come back).

Memory Management is explicit.



Programmer is responsible for.

moving data from

host → global → local & back

Kernel → Code program or function to accelerator device for computation.

- OpenCL main CPU Se kernel execute kernel by host using `clEnqueueNDRangeKernel()`

- Vector addition - Host

Basic host program → Define Platform = devices + context + queues.

• Create + build Program.

• Setup memory objects

• Define Kernel

• Submit Commands.

Kernel void add (global const int^{float} a, global b) {

{ int id = get_global_id(0);
c[id] = a[id] + b[id]; } }.

Kernel code

CPU pV shading gr.

- Set up Host Program in C++ ?

CL_ENABLE_EXCEPTIONS,

#include <iostream>

#include <CL/cl.hpp>

- Context → No environment for kernels computing host.
for memory management and synchronization
define host by.

- Command queue → A queue jobs say send commands device to bhosting host.

only for 1 device.
with 1 context.

Device.

Device memory



Context

— Kernel (GPU पर कैलज जा).

CL:: Context

Context (CL_DEVICE_TYPE_DEFAULT);

" " " CPU

" " " GPU

" " " ACCELERATOR.

CL:: Command Queue

queue (context);

2) Commands include: ① Kernel execution

② Memory object management

③ Synchronization

inorder queues (Jo commands host program jise order mein dayej
out of order queues (device can execute in any order)

CL:: Program program(context, kernel source, true); \rightarrow written
on host.

String or
for function
jo kernel code file
say kernel.h

means openCL to
build the program to
compile our code to de.

- Buffers (no memory object go host (CPU) se declare htag main
CL:: Buffer)
- Host has arrays which has actual data.
- Create device-side buffer which holds data in GPU
or device memory.

CL:: Buffer d_q (Context, h-a.begin(), h-a.end(), true);

Kernel operates on

this device buffer.

host array
ke data kaha say
kaha + K

Read only
buffer
Only device
can read.

3) Setup Memory Objects

- Create input vectors on host.

```
std::vector<float> h-a(LENGTH), h-b(LENGTH), h-c(LENGTH);  
for (i=0; i<LENGTH; i++)  
{  
    h-a[i] = rand();  
    h-b[i] = rand();  
}
```

```
cl::Buffer d-a(Context, CL_MEM_READ_ONLY, a.size() * sizeof(float));  
cl::Buffer d-b(Context, CL_MEM_READ_ONLY, b.size() * sizeof(float));  
cl::Buffer d-c(Context, CL_MEM_READ_ONLY, c.size() * sizeof(float));
```

- Device memory mein 2 taraf ke memory objects hoga hain.

a) Buffer Object.

- It's linear.

- Simple arrays, matrices

jo continuous line main

hota hoga wo buffer mein

hain.

b) Image Object.

- 2D or 3D region of memory.

- Image data can only be accessed through read or write function.

- Device buffer se host memory mein data kaise copy hoga?

```
cl::copy(queue, d-c, hc.begin(), hc.end());
```

↓
device
buffer.

- Host memory, say device buffer.

```
cl::copy(queue, *bac.begin(), hc.end(), d-c);
```

4) Define Kernel. (Create a kernel function for kernel we want to call)
cl:: Kernel kernel (program, "vecAdd");

kernel.setArg (0, bufA);
kernel.setArg (1, bufB);
kernel.setArg (2, bufC);

Kernel void vecAdd(

~~You can call the kernel as a function in your host code~~

5) Enqueue Commands.

For kernel launcher, specify global & local dimensions.

cl::NDRange global (1024)

cl::NDRange local (64) "If we don't specify local dimension it is assumed as cl::NullRange .

• Kernel Enqueue. (non-blocking)

queue.enqueueNDRangeKernel (kernel, cl::NullRange,
cl::NDRange (a.size()), cl::NullRange);

• Read Back Result (blocking)

queue.enqueueReadBuffer (bufC, CL_TRUE, 0,
c.size () * sizeof (float), c.data ());

Do not use this. Just for practice.

Vector addition by OpenCL

#define CL_ENABLE_EXCEPTIONS.

#include <CL/cl.h>

#include <iostream>

#include <vector>

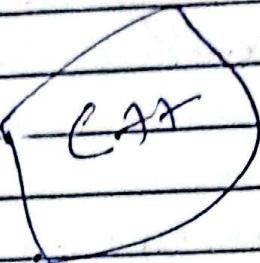
using namespace std;

int main()

1) vector<float> A = { 1, 2, 3, 4 };

vector<float> B = { 5, 6, 7, 8 };

vector<float> C(4);



2) "Setup opencl context and queue.

cl::Context context(CL_DEVICE_TYPE_DEFAULT);

cl::CommandQueue queue(context);

3) "Write kernel source as string.

~~kernel void add(global float* A, global float* B, global float* C)~~

const char KernelSource = R"clc(kernel void vecAdd(

~~global const float* A,~~

~~global const float* B,~~

~~global float* C)~~

{ int i = get_global_id(0);

C[i] = A[i] + B[i]; }) clc;

4) // Build the Program

~~cl::Program::Sources sources;~~

~~sources.push_back({&kernelSource, std::string("KernelSource")});~~

~~cl::Program program(context, sources);~~

~~program.build();~~

5) // Create device buffers, copy A & B from host \rightarrow device.

~~cl::Buffer bufA(context, CL_MEM_READ_ONLY |~~

~~CL_MEM_COPY_HOST_PTR, sizeof(float)*A.size(),~~
~~A.data());~~

~~cl::Buffer bufB(context, CL_MEM_READ_ONLY | CL_MEM_COPY_HOST_PTR,~~
~~sizeof(float)*B.size(), B.data());~~

~~cl::Buffer bufC(context, CL_MEM_WRITE_ONLY, sizeof(float)*~~
~~C.size());~~

6) // Set up Kernel & its arguments

~~cl::Kernel kernel(program, "readd");~~

~~kernel.setArg(0, bufA);~~

~~kernel.setArg(1, bufB);~~

~~kernel.setArg(2, bufC);~~

7) // Enqueue kernel for execution

~~queue.enqueueNDRangeKernel(kernel, cl::NullRange, cl::NDRange~~
~~(A.size()), cl::NullRange);~~

8) // Read Back Results (blocking read)

~~queue.enqueueReadBuffer(bufC, CL_TRUE, 0, sizeof(float)*C.size(),~~
~~C.data());~~

Use this.

```
#include <CL/CL.h>
#include <iostream>
#include <vector>

const char * Kernel Source =
"__kernel void kernel (global const float *a, b... , c)
{ int i = get_global_id(0);
  c[i] = A[i] + B[i]; }";

int main()
{
  float h_a[N] = {1, 2, 3, 4};
  float h_b[N] = {5, 6, 7, 8};
  float h_c[N];
}

cl_platform_id platform;
cl_device_id device;
cl_context context;
cl_command_queue queue;
cl_program program;
cl_kernel kernel;
cl_mem bufA, bufB, bufC;
cl_int status;
```

① // Select platform & device.

```
status = clGetPlatformIDs(1, &platform, NULL);
if (status != CL_SUCCESS) { cout << "Error" << endl; }

status = clGetDeviceIDs(platform, CL_DEVICE_TYPE_DEFAULT,
1, &device, NULL);
```

*

cl CreateProgramWithSource.

(2) // Context Create Kernel handle.

Context = ~~clCreateContext~~ clCreateContext(NULL, 1, &device, NULL, NULL, &status);

*

(3) // Command queue.

queue = clCreateCommandQueue(Context, device, 0, &status);

*

(4) // Program + kernel Source code.

Program = clCreateProgramWithSource(Context, 1, &KernelSource, NULL, &status);

~~clBuildProgram~~

*

status = clBuildProgram(Program, 1, &device, NULL, NULL, NULL);

*

Kernel = clCreateKernel(Program, "vecadd", &status);

*

(5) // Buffer Create KBO.

bufA = clCreateBuffer(Context, CL_MEM_READ_ONLY, c.sizeof(T), &sizeof(float), H-a, &status);

bufB = clCreateBuffer(Context, CL_MEM_READ_ONLY, b.sizeof() * sizeof(float), H-b, &status);

bufC = clCreateBuffer(Context, CL_MEM_WRITE_ONLY, c.sizeof() * sizeof(float), NULL, &status);

(6) // Set Kernel arguments

status = clSetKernelArg(Kernel, 0, sizeof(CLmem), &bufA);

*

status = clSetKernelArg(Kernel, 1, sizeof(CLmem), &bufB);

status = clSetKernelArg(Kernel, 2, sizeof(CLmem), &bufC);

(7) // Enqueue Kernel

size_t global_work_size = N;

status = clEnqueueNDRangeKernel(queue, kernel, 1, NULL,
global_work_size, NULL, 0, NULL, NULL);

(8) // Results wapis host memory mein read tero.

status = clEnqueueReadBuffer(queue, bufC, CL_TRUE,
0, sizeof(float) * N, h -> a, NULL, NULL);

(9) clReleaseMemObject(bufA);

clReleaseMemObject(bufB);

clReleaseMemObject(bufC);

clReleaseKernel(kernel);

clReleaseProgram(program);

clReleaseCommandQueue(queue);

clReleaseContext(context);

return 0;

}

OpenCL Kernel

• Derived from ISO C99. → no recursions

• OpenCL has a lot of threads which work parallel.

• These thread groups (work group) main divided hain. aur har group hardware ke compute unit pr chalta h.

• Har thread apni unique ID se apna data process

krta h.

• Main purpose → loops ki jagah, parallel executions kرنے by AZ performance.

• small and fast memories like private or local should be used more. Global & host memories from use ltra sec

Slow.

- NDRange means ~~poor grid~~ johan. Kernels execute two Valley bin.
- Nav grid chose Chose parts mein divided hoga by size work group kehte bin (local size).
- Work groups main work items hitay bin.
- Synchronization functions.
 - Barrier \rightarrow 1 workgroup koy sab work items ek chay key se synchronize hitay bin.
 - Memory fences \rightarrow memory read/write sahi order mein perform hain.
- Restrictions in OpenCL
 - Kernel key and functions key pointers use nahi kar sktey.
 - Arrays whose size vary on run-time, we cannot use.
 - Recursion not supported
- - Matrix Multiplication cost depends on FLOPs and memory movement.
- FLOPs calculations
$$2 \times n^3 = O(n^3)$$
$$3 \times n^2 = O(n^2)$$
- Whenever memory access hog max FLOPs perform kernel.