

# Introduction to GPU programming models

**Andrey Alekseenko**

**Science for Life Laboratory, KTH Royal Institute of Technology, Sweden**

Introduction to GPU, Stockholm/Zoom, 15 September, 2023

Presentation based on the **ENCCS lesson "Introduction to GPU programming models"**

# GPU Programming models

GPUs are great! How do we use them?

- C++/Fortran:
  - Standard-based solutions
  - Directive-based programming
  - Non-portable kernel-based models
  - Portable kernel-based models
- Higher-level languages
- Exotic languages

N.B.: The classification above is neither definitive nor exhaustive.

# Standard-based models

Use C++/Fortran to offload algorithms from the standard library to the GPU.

Despite being "standard", the support for them is extremely limited.

Limited flexibility and performance.

# C++ Standard Parallelism

- Introduced in C++17
- Automatic parallelization of algorithms from C++ standard library

```
#include <algorithm>
#include <execution>
// ...
std::transform(std::execution::par_unseq, A, A + n, B, C,
               [](float x, float y) { return x + y; });
```

- Supported on NVIDIA GPUs via [NVIDIA HPC SDK \(nvc++\)](#).
- Experimental support for other GPUs via [hipSYCL/Open SYCL](#).

# Fortran Standard Parallelism

```
subroutine vectoradd(A,B,C,n)
  real :: A(:), B(:), C(:)
  integer :: n, i
  do concurrent (i = 1: n)
    C(i) = A(i) + B(i)
  enddo
end subroutine vectoradd
```

- Supported on NVIDIA GPUs via [NVIDIA HPC SDK \(nvfortran\)](#)

# Directive-based programming

- Annotate existing C++/Fortran *serial* code with *hints* to indicate which loops and regions to execute on the GPU.
- Minimal changes to the source code, but limited performance and expressiveness.
- [OpenACC](#) is more *descriptive*. Developed in 2010 specifically to target accelerators. Supports NVIDIA GPUs with NVIDIA HPC SDK, NVIDIA and AMD GPUs with GCC 12.
- [OpenMP](#) is more *prescriptive*. Developed in 1997 for multicore CPUs, accelerator support added later. Supports most GPUs with GCC and LLVM/Clang.

# C++

```
// OpenMP

#pragma omp target teams distribute parallel for simd
for (i = 0; i < n; i++) {
    C[i] = A[i] + B[i];
}

// OpenACC

#pragma acc parallel loop
for (i = 0; i < n; i++) {
    C[i] = A[i] + B[i];
}
```

# Fortran

```
! OpenMP

!$omp target teams distribute parallel do simd
do i = 1, n
    C(i) = A(i) + B(i)
end do
!$omp end target

! OpenACC

!$acc parallel loop
do i = 1, n
    C(i) = A(i) + B(i)
end do
!$acc end parallel loop
```

# Kernel-based programming

- Requires deep knowledge of GPU architecture.
- The code to execute on the GPU is outlined into a separate function (kernel).
- Low-level code (conceptually), more control, higher performance (or not).



# Non-portable models

## CUDA

- First mainstream GPU computing framework, widely used and supported.
- Developed by NVIDIA, works only on NVIDIA GPUs.
- Extensive ecosystem, libraries, tutorial.
- Stable, feature-rich, well-supported.

## HIP

- Developed by AMD, works on AMD and NVIDIA\* GPUs.
- *Very* similar to CUDA, automatic conversion tools exist.
- Less mature, no official support on consumer GPUs.

# CUDA

```
__global__ void vector_add(float *A, float *B, float *C, int n) {  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n) {  
        C[tid] = A[tid] + B[tid];  
    }  
}
```

```
// Allocate GPU memory  
cudaMalloc((void**)&Ad, n * sizeof(float)); // ...  
// Copy the data from the CPU to the GPU  
cudaMemcpy(Ad, Ah, sizeof(float) * n, cudaMemcpyHostToDevice); // ...  
// Define grid dimensions: how many threads to launch  
dim3 blockDim{256, 1, 1};  
dim3 gridDim{(n/256)*256 + 1, 1, 1};  
vector_add<<<gridDim, blockDim>>>(Ad, Bd, Cd, n);
```

# HIP

```
__global__ void vector_add(float *A, float *B, float *C, int n) {  
    int tid = threadIdx.x + blockIdx.x * blockDim.x;  
    if (tid < n) {  
        C[tid] = A[tid] + B[tid];  
    }  
}
```

```
// Allocate gpu memory  
hipMalloc((void**)&Ad, n * sizeof(float)); // ...  
// Copy the data from the CPU to the GPU  
hipMemcpy(Ad, Ah, sizeof(float) * n, hipMemcpyHostToDevice); // ...  
// Define grid dimensions: how many threads to launch  
dim3 blockDim{256, 1, 1};  
dim3 gridDim{(n/256)*256 + 1, 1, 1};  
vector_add<<<gridDim, blockDim>>>(Ad, Bd, Cd, n);
```

## Portable models: OpenCL

- Cross-platform open standard for accelerator programming, based on C.
- Supports GPUs, FPGAs, embedded devices.
- Supported by all major vendors, but to a varying degree. No latest features.
- Good performance requires vendor extensions, defeating portability aspect.
- Separate-source model, no compiler support required: just headers and a runtime library.

# OpenCL: kernel

```
static const char* kernel_source = R"(
__kernel void vector_add(__global float *A, __global float *B, __global float *C, int n) {
    int tid = get_global_id(0);
    if (tid < n) {
        C[tid] = A[tid] + B[tid];
    }
}
)";
```

# OpenCL: launch code

```
// Boilerplate...
cl_platform_id platform;
clGetPlatformIDs(1, &platform, NULL);
cl_device_id device;
clGetDeviceIDs(platform, CL_DEVICE_TYPE_GPU, 1, &device, NULL);
cl_context context = clCreateContext(NULL, 1, &device, NULL, NULL, NULL);
cl_command_queue queue = clCreateCommandQueue(context, device, 0, NULL);
cl_program program = clCreateProgramWithSource(context, 1, &kernel_source, NULL, NULL);
clBuildProgram(program, 1, &device, NULL, NULL, NULL);
cl_kernel kernel = clCreateKernel(program, "vector_add", NULL);
// Allocate the arrays on GPU
cl_mem Ad = clCreateBuffer(context, CL_MEM_READ_ONLY, n * sizeof(float), NULL, NULL); // ...
// Copy the data from CPU to GPU
clEnqueueWriteBuffer(queue, Ad, CL_TRUE, 0, n * sizeof(float), Ah, 0, NULL, NULL); // ..
// Set arguments and launch the kernel
clSetKernelArg(kernel, 0, sizeof(cl_mem), &Ad);
clSetKernelArg(kernel, 1, sizeof(cl_mem), &Bd);
clSetKernelArg(kernel, 2, sizeof(cl_mem), &Cd);
cl_int n_as_cl_int = n;
clSetKernelArg(kernel, 3, sizeof(cl_int), &n_as_cl_int);
size_t globalSize = n;
clEnqueueNDRangeKernel(queue, kernel, 1, NULL, &globalSize, NULL, 0, NULL, NULL);
```

## Portable models: SYCL

- Cross-platform open standard for accelerator programming, inspired by OpenCL.
- Single-source model based on C++17, requires compiler support and C++ knowledge.
- Two major implementations: Intel oneAPI DPC++ and hipSYCL/Open SYCL.
- Supported officially only by Intel, works via CUDA/HIP for NVIDIA/AMD.
- Interoperability with native features (at the expense of portability).

# SYCL: USM approach

```
// Boilerplate
sycl::queue queue{{sycl::property::queue::in_order()}};
// Allocate GPU memory
float* Ad = sycl::malloc_device<float>(n, queue); // ...
// Copy the data from CPU to GPU
queue.copy<float>(Ah, Ad, n); // ...
// Submit a kernel into a queue; cgh is a helper object
queue.submit([&](sycl::handler &cgh) {
    cgh.parallel_for<class Kernel>(sycl::range<1>{n}, [=](sycl::id<1> i) {
        Cd[i] = Ad[i] + Bd[i];
    });
});
```

There is also support for automatic memory management via SYCL buffers.



# Kokkos

- Open-source programming model for heterogeneous parallel computing, mostly developed at Sandia National Laboratories.
- Single source model, similar to SYCL (but earlier).
- No official vendor support, implemented on top of other frameworks.
- Interoperability with native features (at the expense of portability).

# Higher-level language support: Julia

Julia has first-class support for GPU programming through the following packages that target GPUs from all three major vendors:

- `CUDA.jl` for NVIDIA GPUs.
- `AMDGPU.jl` for AMD GPUs.
- `oneAPI.jl` for Intel GPUs.
- `Metal.jl` for Apple M-series GPUs.

# Higher-level language support: Python

- [CuPy](#): NVIDIA-only.
- [cuDF](#): NVIDIA-only.
- [PyCUDA](#): ... you guessed it.
- [Numba](#): (NVIDIA and deprecated AMD support).

# Exotic GPU-specific languages

- [Chapel](#): developed by Cray/HPE for parallel scientific programming, early AMD/NVIDIA support.

```
on here.gpus[0] {  
  var A, B, C : [0..#n] real(32); // Allocate  
  A = Ah; // Copy from CPU to GPU, ...  
  forall i in 0..#n { C[i] = A[i] + B[i]; };  
}
```

- [Futhark](#): functional data-parallel language for GPU programming with OpenCL and CUDA backends.

```
def vecadd A B =  
  map (\(x,y) -> x + y) (zip A B)
```

# (Opinionated) Summary

There are exceptions, but, in general, prefer:

- Standard-based: if you already use STL algorithms and have supported compiler
- Directive-based: if you have serial code and no time
- CUDA or HIP: if you are targeting specific vendor and want best performance
- OpenCL: if you want portability on *any* hardware
- Kokkos or SYCL: if you want performance and portability on modern GPUs
- Higher-level languages: if you have C, C++, and Fortran
- Exotic languages: if you are feeling adventurous

## More details

- ENCCS GPU Programming course
  - Other ENCCS training materials
- PDC Training