



UNIVERSIDAD DEL BÍO-BÍO  
FACULTAD DE CIENCIAS EMPRESARIALES

# Frameworks de Programación

## Computación Heterogénea

Profesor: Dr. Joel Fuentes - [jfuentes@ubiobio.cl](mailto:jfuentes@ubiobio.cl)

Ayudantes:

- Daniel López - [daniel.lopez1701@alumnos.ubiobio.cl](mailto:daniel.lopez1701@alumnos.ubiobio.cl)
- Sebastián González - [sebastian.gonzalez1801@alumnos.ubiobio.cl](mailto:sebastian.gonzalez1801@alumnos.ubiobio.cl)

Página web del curso: <http://www.face.ubiobio.cl/~jfuentes/classes/ch>

# Contenidos

- CUDA
- SYCL (DPC++)

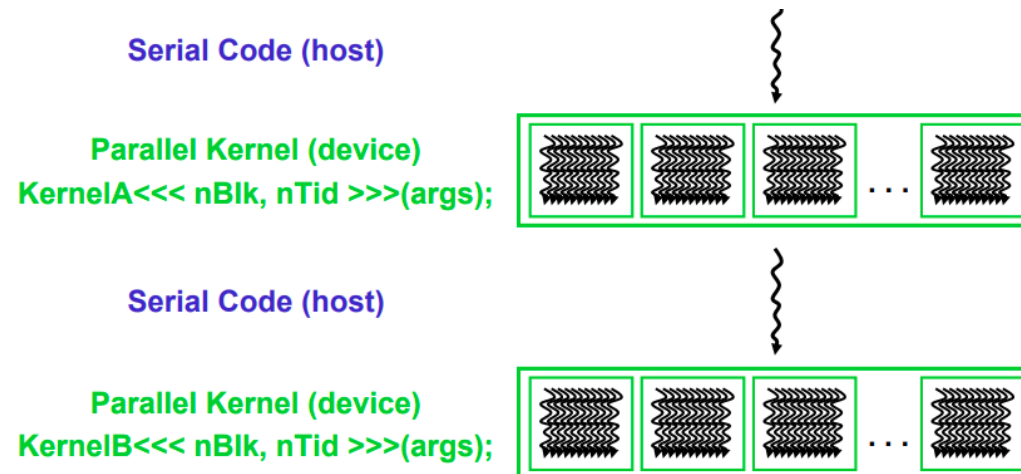
# Lenguajes y Frameworks de Programación

- CPU multi-core
  - C++, Java, OpenMP, DPC++
- GPU
  - OpenCL, CUDA, DPC++, OpenACC
- FPGA
  - OpenCL, DPC++

# CUDA

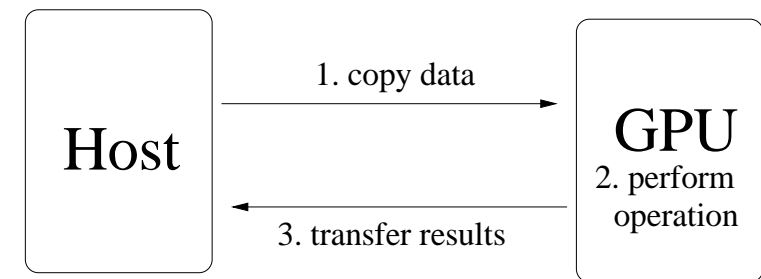
# CUDA

- Compute Unified Device Architecture
- Programación en C extendido
- Programación de código serial en Host (CPU)
- Programación de código paralelo en Device (GPU)



# CUDA

- GPU puede ejecutar muchos hilos simultaneamente, pero no independientemente
  - Hilos en Device conectados en grupos llamados warps
  - Todos los miembros de un warp ejecutan la misma operación
    - SIMT = Single Instruction, Multiple Threads
- Programador escribe función que corre en el dispositivo (kernel)
- Función se invoca con un número de bloques
- Todos los hilos ejecutan la misma función
- Host y GPU tienen espacios de memoria separadas
  - Memoria debe ser transferida explícitamente



# CUDA

- C extendido

## Declaraciones

- `global, device, shared, local, constant`

## Palabras claves

- `threadIdx, blockIdx`

## Intrinsics

- `__syncthreads`

## Runtime API

- `Memory, symbol, execution management`

```
#include <stdio.h>

__global__ void hello() {
    int id = threadIdx.x + blockIdx.x * blockDim.x;
    printf("Hello from thread %d (%d of block %d)\n",
        id, threadIdx.x, blockIdx.x);
}

int main() {
    //launch 3 blocks of 4 threads each
    hello<<<3,4>>>>();

    //make sure kernel completes
    cudaDeviceSynchronize();
}
```

# CUDA

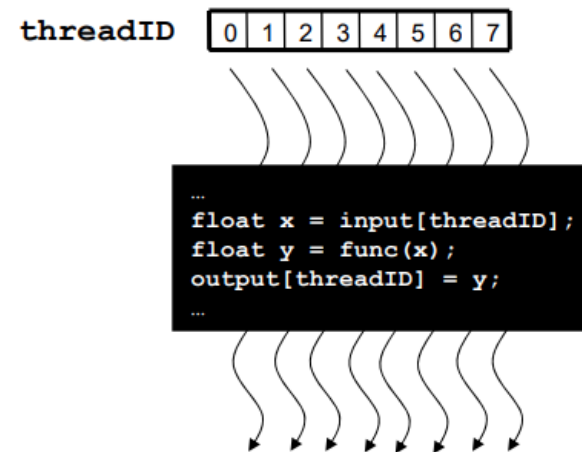
- Posible salida por pantalla
- Hilos y bloques son ejecutados en cualquier orden

Hello from thread 0 (0 of block 0)  
Hello from thread 1 (1 of block 0)  
Hello from thread 2 (2 of block 0)  
Hello from thread 3 (3 of block 0)  
Hello from thread 8 (0 of block 2)  
Hello from thread 9 (1 of block 2)  
Hello from thread 10 (2 of block 2)  
Hello from thread 11 (3 of block 2)  
Hello from thread 4 (0 of block 1)  
Hello from thread 5 (1 of block 1)  
Hello from thread 6 (2 of block 1)  
Hello from thread 7 (3 of block 1)



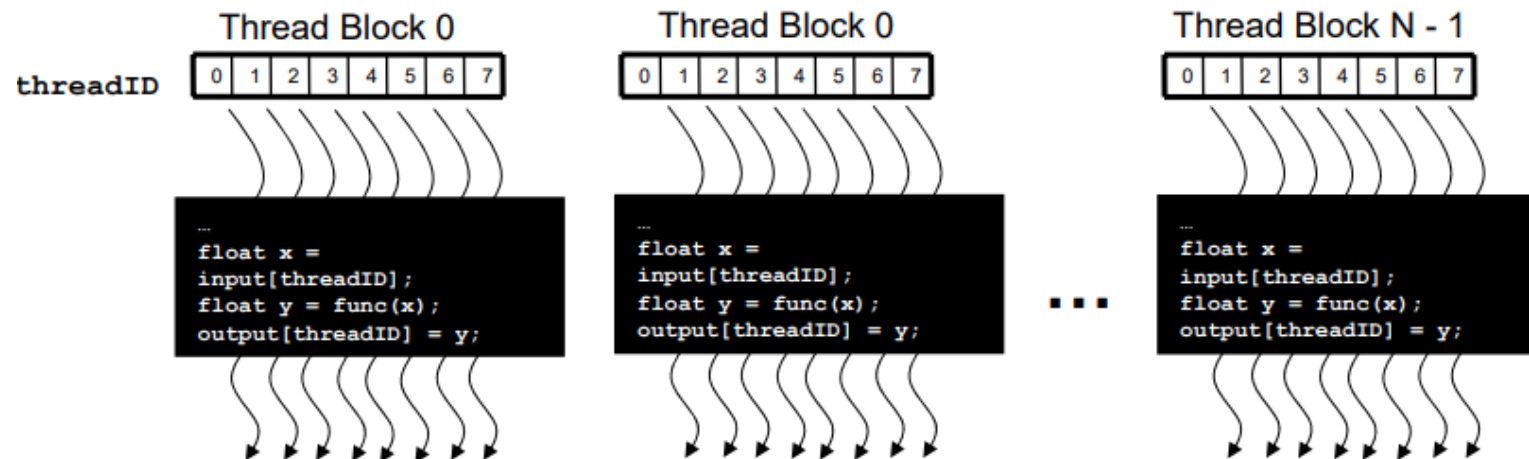
# CUDA

- En CUDA todos los hilos ejecutan el mismo código
- Cada hilo tiene su ID que es usado para calcular accesos a memoria y realizar decisiones de control



# CUDA

- Hilos se organizan en múltiples bloques
- Hilos en un bloque pueden cooperar a través del uso de memoria compartida, operaciones atómicas y barreras de sincronización.
- Hilos en diferentes bloques no pueden cooperar.



# CUDA: Ejemplo suma de vectores

```
int main() {  
    int* a;           //first input array (on host)  
    int* a_dev;        //first input array (on device)  
  
    a = (int*) malloc(N*sizeof(int));  
    cudaMalloc((void**) &a_dev, N*sizeof(int));  
  
    ...               //same for b and res  
  
    free(a);  
    cudaFree(a_dev);  
}
```

- En el host:
  1. Reservar memoria en dispositivo
  2. Copiar datos a dispositivo
  3. Llamar al kernel
  4. Copiar resultados al host
  5. Liberar memoria del dispositivo
- En dispositivo:
  - \_\_global\_\_
  - Determinar thread ID

# CUDA: Ejemplo suma de vectores

```
int main() {  
    ...  
    int threads = 512;          //# threads per block  
    int blocks  = (N+threads-1)/threads;  
                                //# blocks (N/threads rounded up)  
    kernel<<<blocks,threads>>>(res_dev, a_dev, b_dev);  
    ...  
}
```

- En el host:
  1. Reservar memoria en dispositivo
  2. Copiar datos a dispositivo
  3. **Llamar al kernel**
  4. Copiar resultados al host
  5. Liberar memoria del dispositivo
- En dispositivo:
  - `__global__`
  - Determinar thread ID

# CUDA: Ejemplo suma de vectores

```
__global__ void kernel(int* res, int* a, int* b) {  
    //sets res[i] = a[i] + b[i]  
    //each thread is responsible for one value of i  
  
    int thread_id = threadIdx.x + blockIdx.x*blockDim.x;  
  
    if(thread_id < N) {  
        res[thread_id] = a[thread_id] + b[thread_id];  
    }  
}
```

- En el host:
  1. Reservar memoria en dispositivo
  2. Copiar datos a dispositivo
  3. Llamar al kernel
  4. Copiar resultados al host
  5. Liberar memoria del dispositivo
- En dispositivo:
  - \_\_global\_\_
  - Determinar thread ID

# SYCL

# SYCL

- SYCL es una propuesta para XPU (GPU, CPU and FPGA).
- Corresponde a una extensión de C++ para soportar programación heterogénea.
- Construído sobre OpenCL (evolución)
- Desarrollada por el grupo Khronos <https://www.khronos.org/>

# SYCL

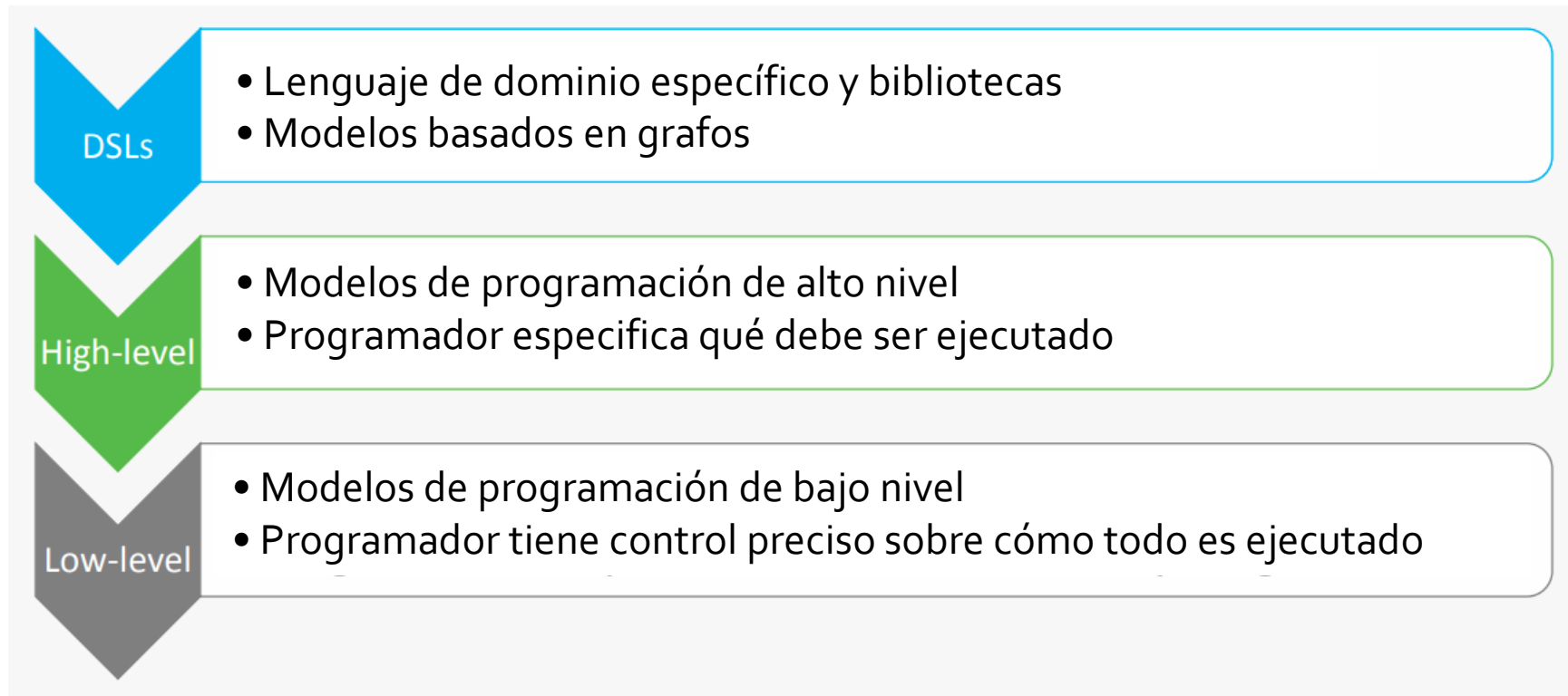
- Plataformas que soportan SYCL:

Implementation	Supported Platforms	Supported Devices	Required Version
ComputeCpp	Windows 10 Visual Studio 2019 (64bit)* Ubuntu 18.04 (64bit)	Intel CPU (OpenCL) Intel GPU (OpenCL)	CE 2.4.0
DPC++	Intel DevCloud Windows 10 Visual Studio 2019 (64bit) Red Hat Enterprise Linux 8, CentOS 8 Ubuntu 18.04 LTS, 20.04 LTS (64bit) Refer to <a href="#">System Requirements</a> for more details	Intel CPU (OpenCL) Intel GPU (OpenCL) Intel FPGA (OpenCL) Nvidia GPU (CUDA)**	2021.2
hipSYCL	Any Linux	CPU (OpenMP) AMD GPU (ROCm)*** Nvidia GPU (CUDA)	Latest develop branch



# SYCL

- Lenguajes heterogéneos



# SYCL

- Problema: Hardware distinto require diferentes enfoques de programación
  - Funcionalidad puede ser portable, pero no el rendimiento
  - Un algoritmo optimizado para una arquitectura en particular puede correr muy mal en otra.
- Solución parcial: Considerar un lenguaje de más alto nivel.

# SYCL

## OpenCL

```
const char *src =
    "__kernel void vecadd(global int *A,\n"
    " global int *B,\n" " global int *C) {\n"
    "  size_t gid = get_global_id(0);\n"
    "  C[gid] = A[gid] + B[gid];\n"
    "}\n"

clSetKernelArg(k, 0, sizeof(cl_mem), &ABuf);
clSetKernelArg(k, 1, sizeof(cl_mem), &BBuf);
clSetKernelArg(k, 2, sizeof(cl_mem), &CBuf);

clEnqueueNDRangeKernel(q, k, 1, NULL, {SIZE}, {32, 1, 1}, 0,
    NULL, NULL);
```

## SYCL

```
auto A = ABuf.get_access(cgh);
auto B = BBuf.get_access(cgh);
auto C = CBuf.get_access(cgh);

cgh.parallel_for(
    cl::sycl::range(CBuf.size()),
    [=](cl::sycl::id idx) {
        C[idx] = A[idx] + B[idx];
    });
```

# Referencias

- ToUCH: Teaching Undergrads Collaborative and Heterogeneous Computing in Consortium for Computing Sciences in Colleges South Central Conference (CCSC19), 2019
- Alastair Murray. Codeplay. Layering Abstractions, Heterogeneous Programming and Performance Portability. 2017.