

## CUDA

**M40** → nella machine

NVIDIA  
↳ CUDA

AMD

↳ openCL  
↳ g++ -lopenCL

kernel code executed by GPU has to be provided as a string

different from openCL because openCL is a library and inside the program you can make calls and the library does all the things.

Cuda is a language (extension of C++) because introduces some features.

**NVCC** → compiler

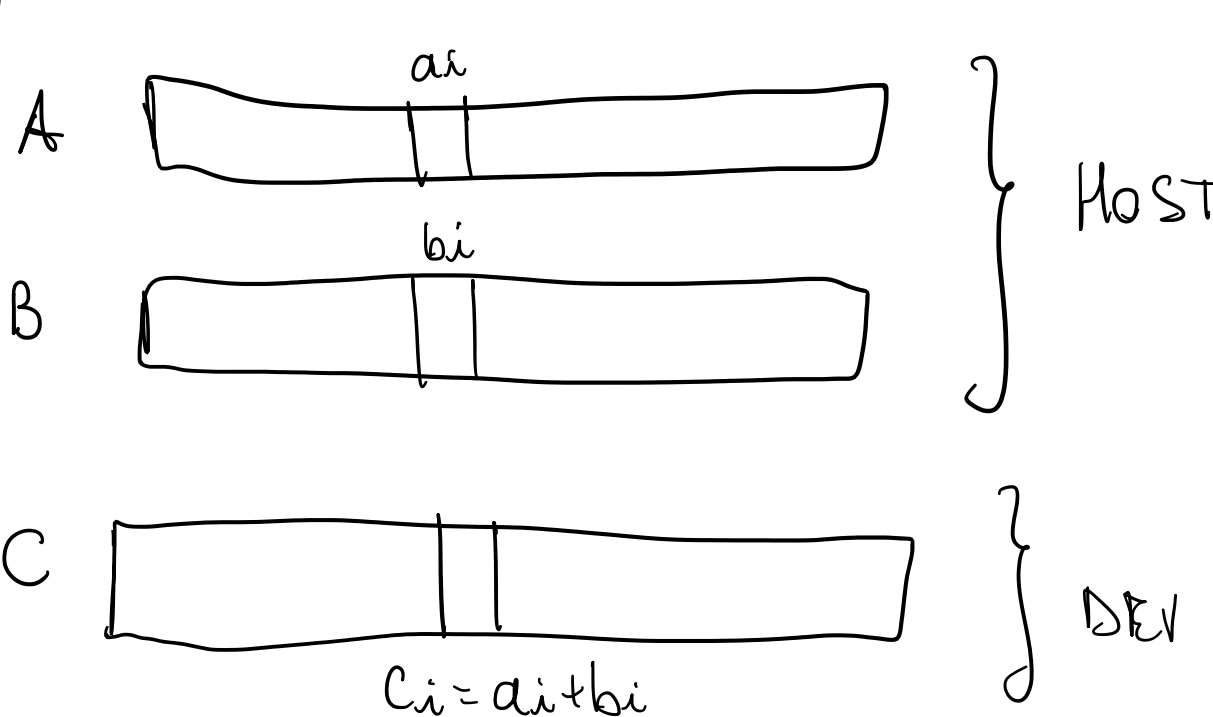
-- global --

<<< ... >>> → is used to tell the compiler to execute the kernel code on GPU.

defined used the kernel keyword in front

name of function

Example of cuda code:



la funzione che somma i due vettori viene definita mettendoli davanti alla definizione  
la keyword -- global --

→ Per allocare la memoria nella GPU  
Cuda Malloc

→ You need to manage events that happens in the GPU.

Example: we can create 2 events, start and stop, before starting the thread and after the end we register the time on this events.

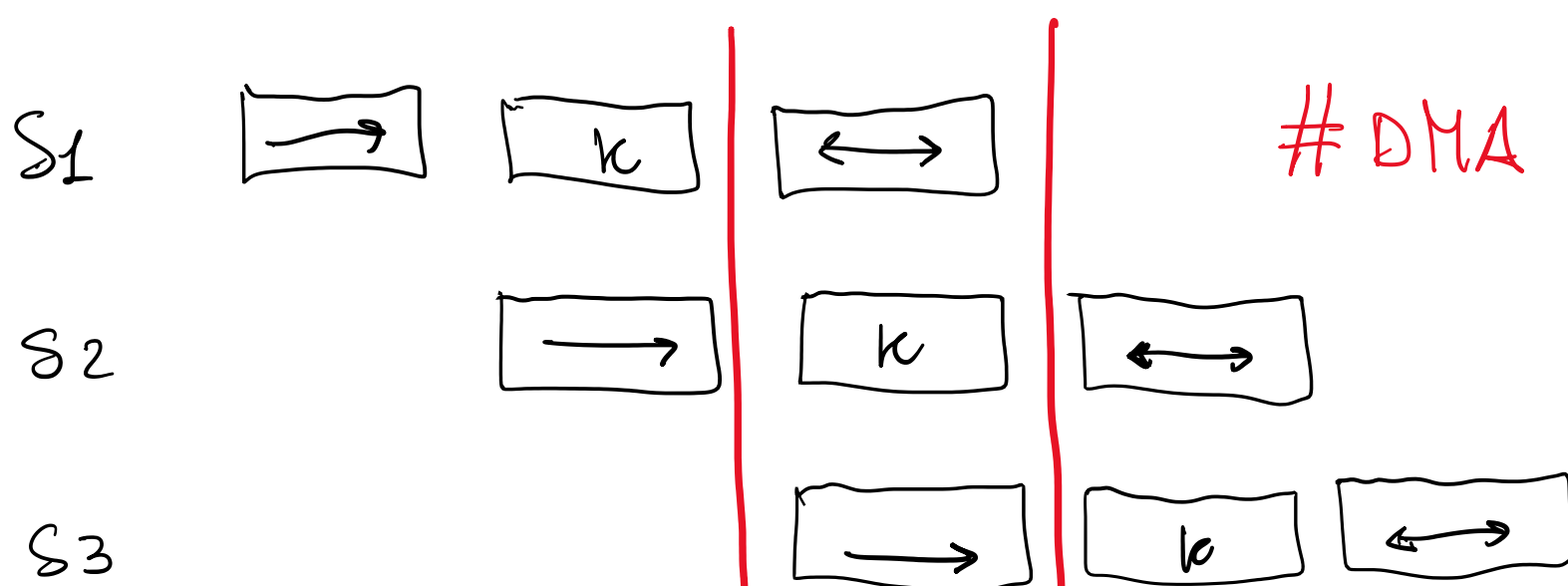
these two events measures the execution time of the kernel.

the events are recorded on the GPU and sometimes we need to synchronize with CPU. **[wait]**

CudaDeviceSynchronize() → wait for GPU to finish before accessing on host

## MACCHINA PER CUDA

131. 112. 136. 61



## Soluzione assignment 3

ooo You need to consider  
o•o the neighbours of the  
ooo node.

Consider two matrix <sup>we write</sup> in one matrix and we read from the other one.

## Assignment 4

Implement a mapReduce object: → tips google

Std::function < Std::pair < tout, tkey, (Tin) >

f(x) → Map function

↓

< y, k >

Std::function < tout (tout, tout) >

function g given < y', k' > → y', k'

reduce function  
Compute  
< y'', k' >

Std::vector < Tin >

↳ return

Std::vector < Std::pair < tout, tkey >

Implement this in parallel. You need to take into account that you are in shared memory.

V1: C++

V2: Fastflow