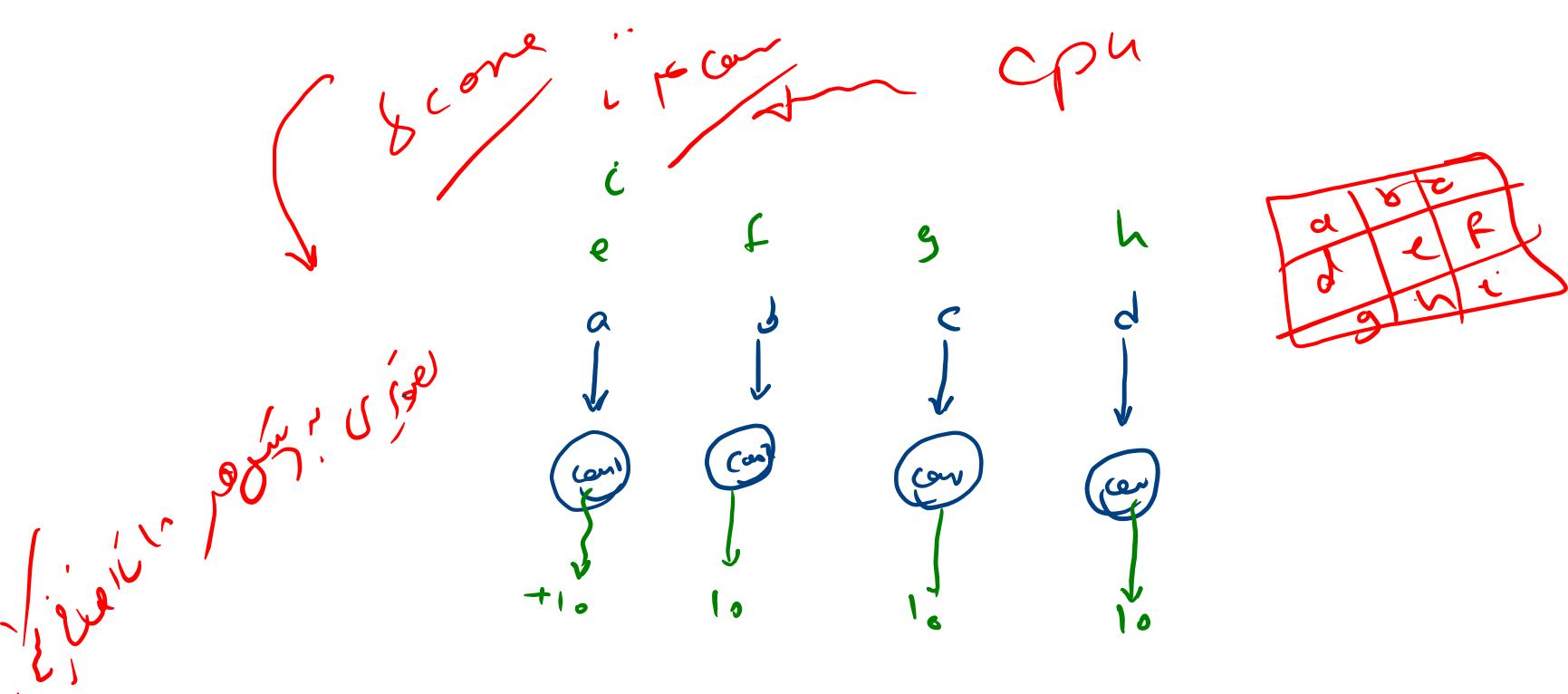


GPU One Step



GPU





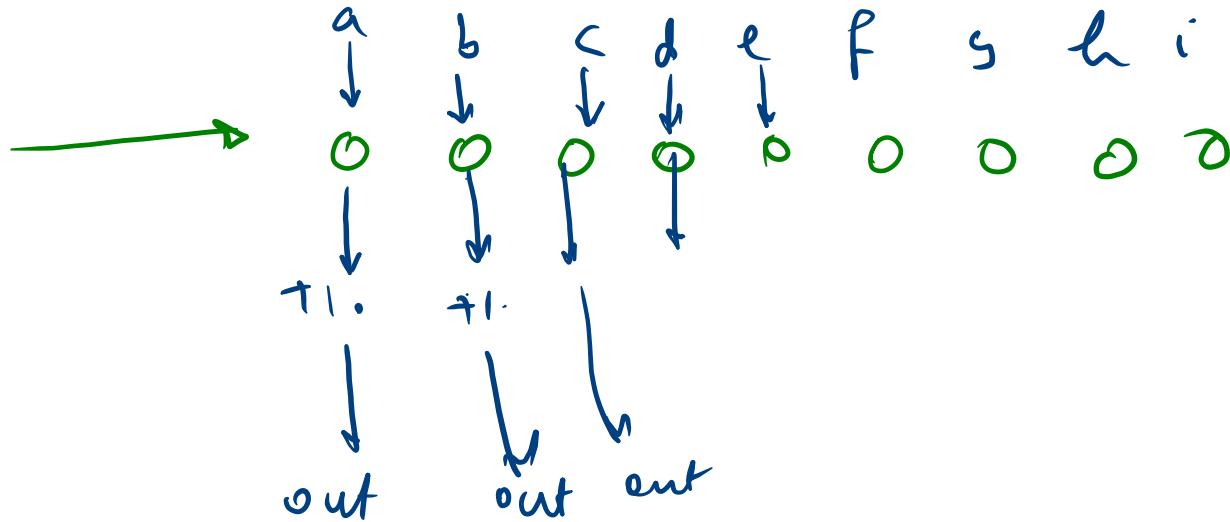
Sequential projects go to jobs on CPU

word

I/O

Word cycle

| | | |
|---|---|---|
| a | b | c |
| d | e | f |
| g | h | i |



nb come ~~using~~ GPU , ~~using~~

~~using~~ GPU ~~not~~ CPU , ~~or~~ this ~~is~~

نحو Gpr نویز کرنل CPU cones



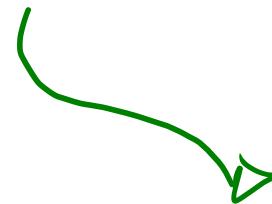
Core

400 MHz

200 MHz

memories

200 GHz
memories



Core

8 GHz



$$a \times b$$

≈ 5 clockcycles on CPU \rightarrow 1μs

≈ 100 clockcycles on GPU

Parallel Segm. mehrere GPUs

Computing

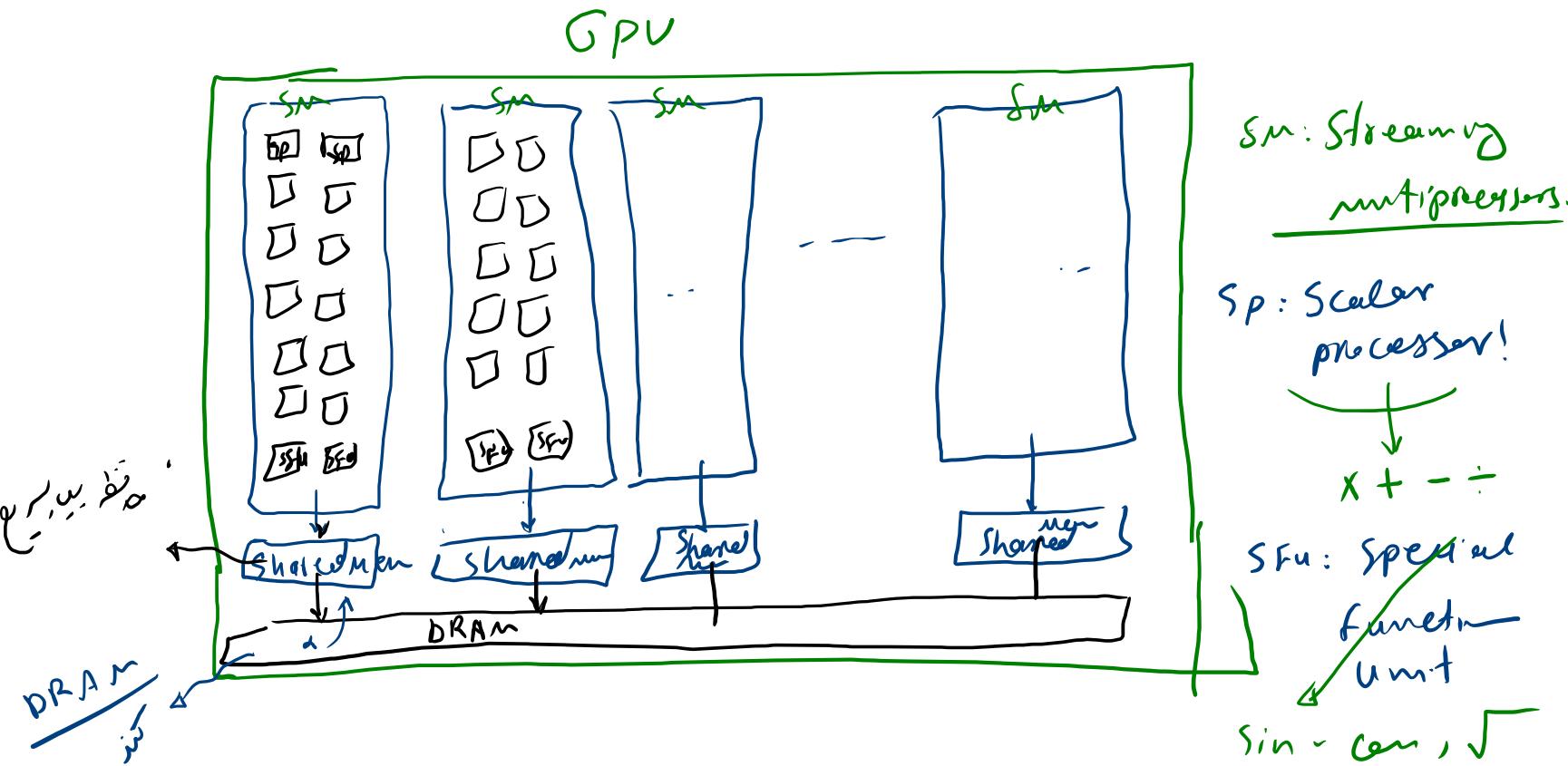
Heterogeneous.

High performance

CPU on seq
+ GPU on parallel

GPU میکرو

میں دنیا سفیر، ۶۸۵
ماہی نو اسٹر مکرم بامن ملک



✓ Shared mem

√8

CPU mem

pixels ← Shared mem

pixels ← CPU mem



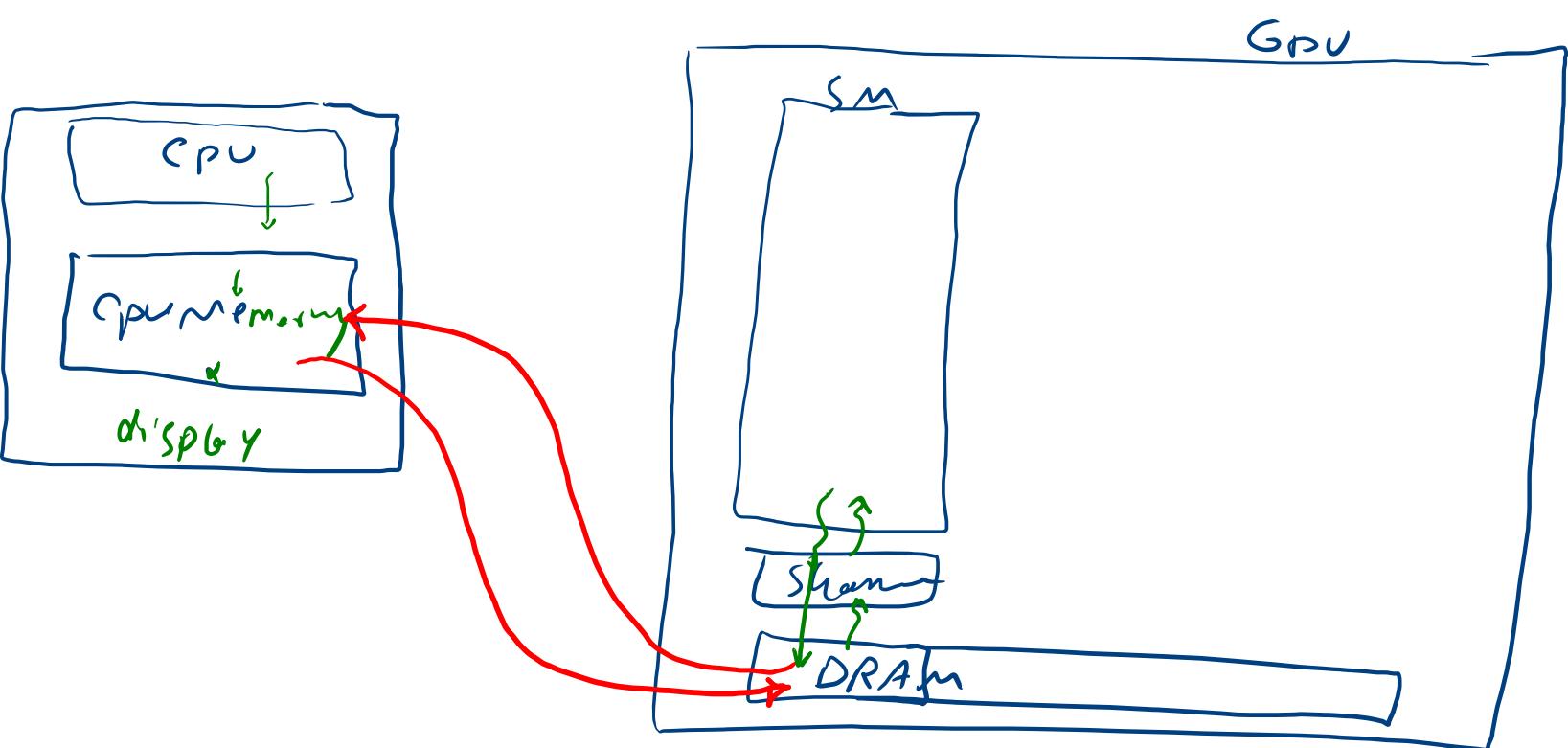
CPU → 4 core

GPU → 4000 core!

NVIDIA

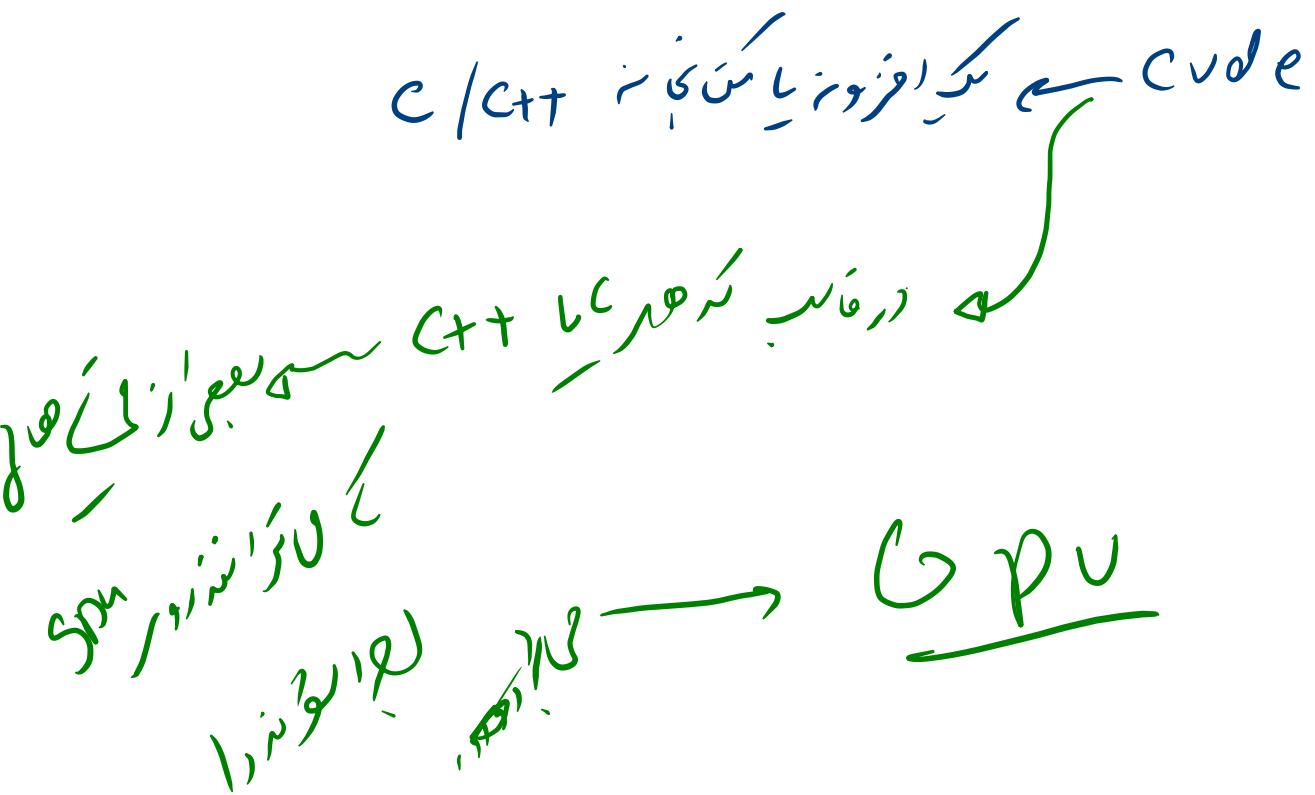
GForce RTX 2080 Ti

{ 4352 SP
64 SP per SM
+



Cuda

Compute
unified
device
Arch.

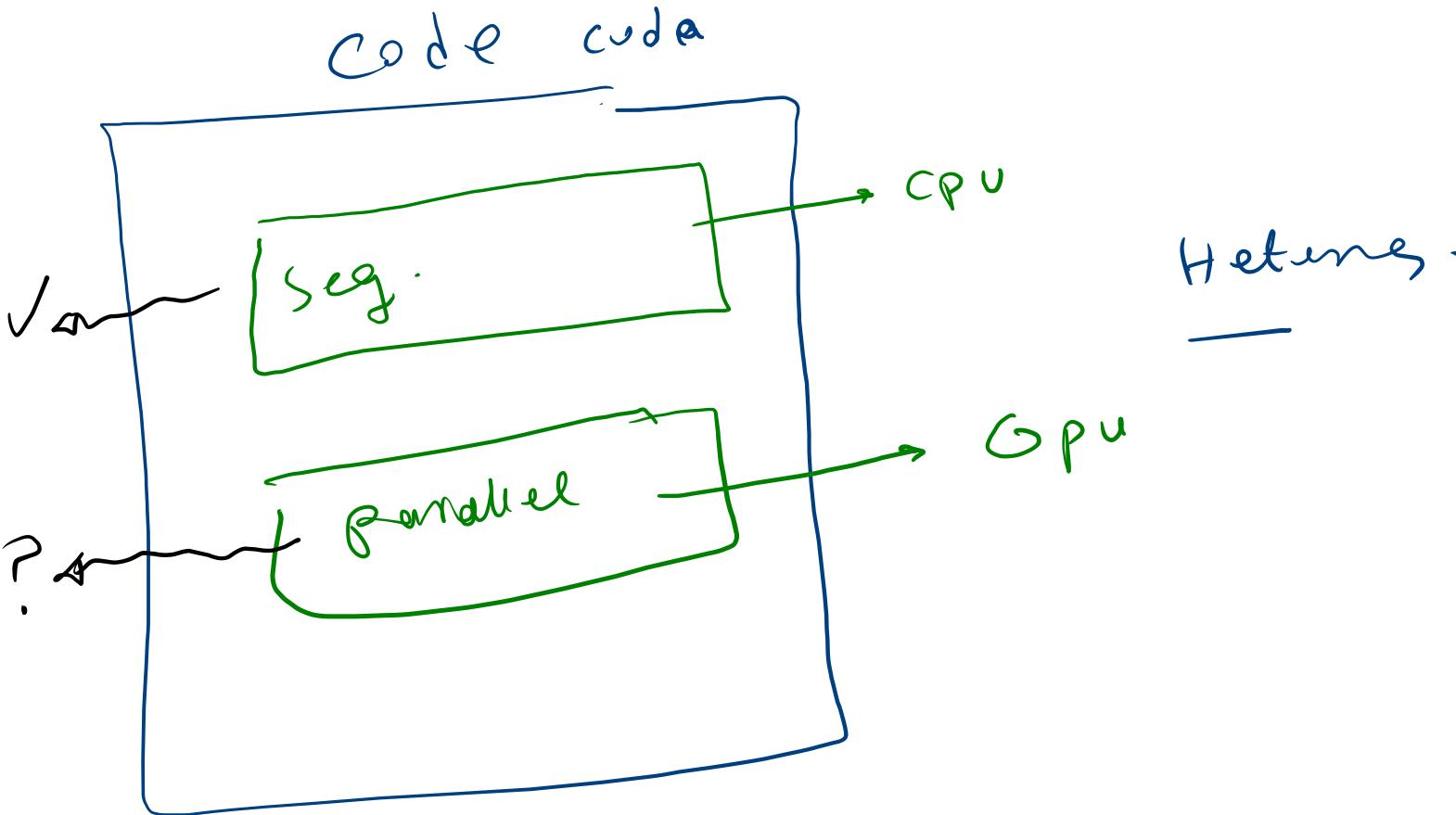


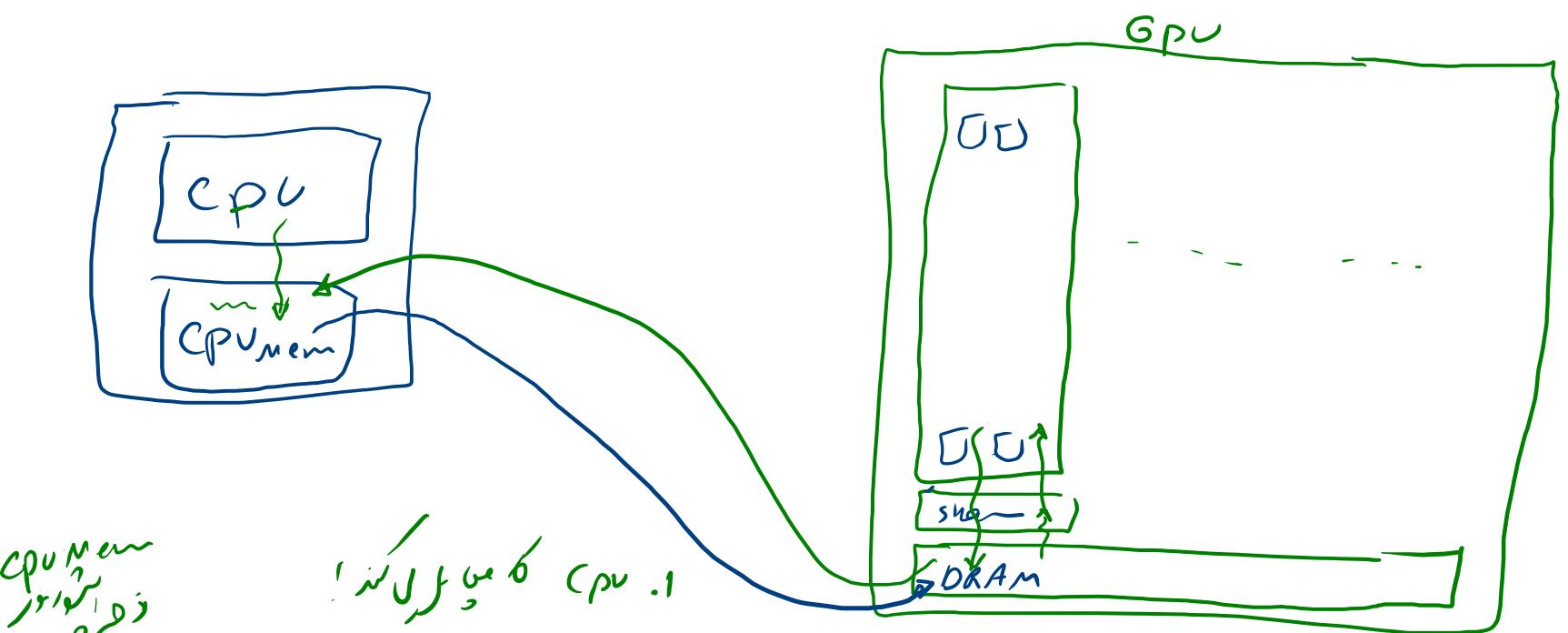
Code

No
CVJ e

(default) acceleration

CP V





CPU میں کام کرنے والے دو حصے ہیں۔
CPU mem اور CPU.
CPU mem کا کام DRAM پر کیا جاتا ہے اور DRAM پر کام کرنے والے دو حصے ہیں۔
CPU mem اور DRAM کو CPU میں کام کرنے والے دو حصے کے مقابلے میں سب سے بڑے ہیں۔

GPU میں کام کرنے والے دو حصے ہیں۔
GPU میں کام کرنے والے دو حصے کے مقابلے میں سب سے بڑے ہیں۔

Cuda basic command

C++ / C

• `cuda malloc()`

• `cuda memcpy()` → cpu to gpu gpu to cpu copy from host to device GPU

• `cudaFree()`

allocate
no name

$a \rightarrow \text{cpu}$
 \downarrow
 gpu
 $\hookrightarrow \text{cuda malloc()}$

$a \text{ copy}$

GPU process

\hookrightarrow a copy to CPU from GPU

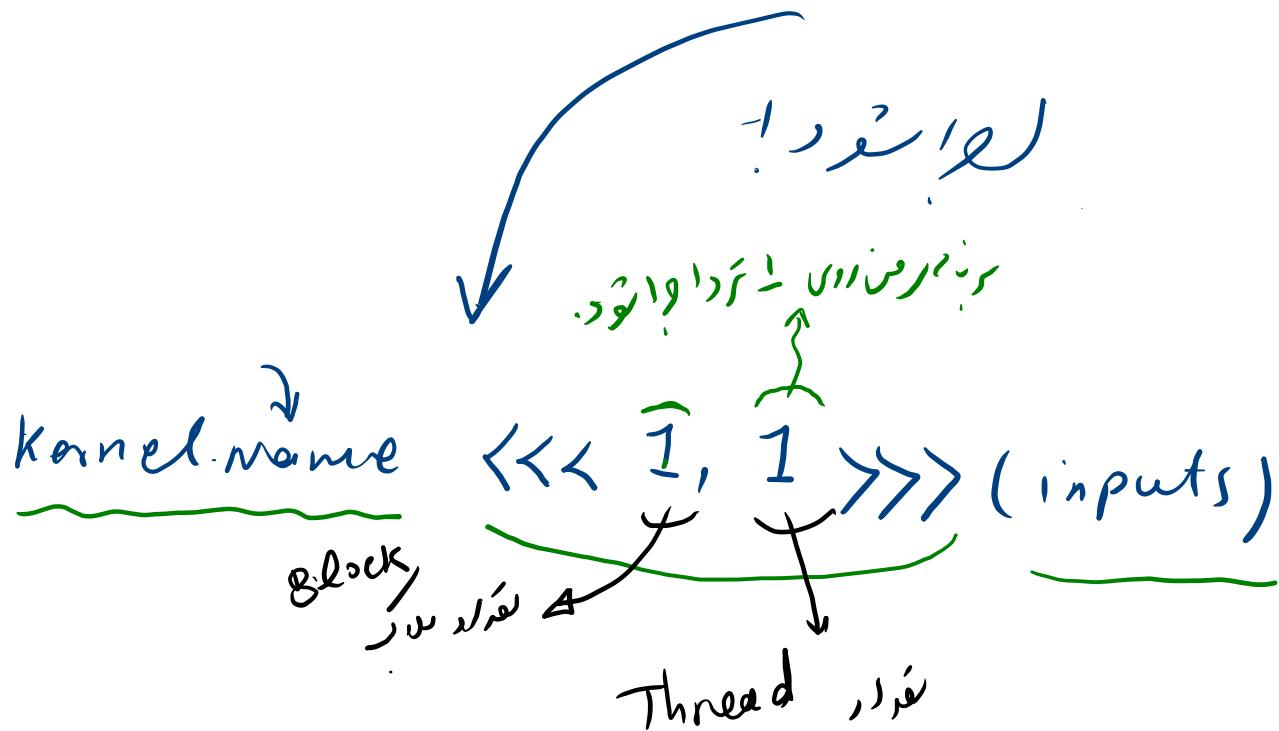
گلوبال کورس کیا ہے

-- global --

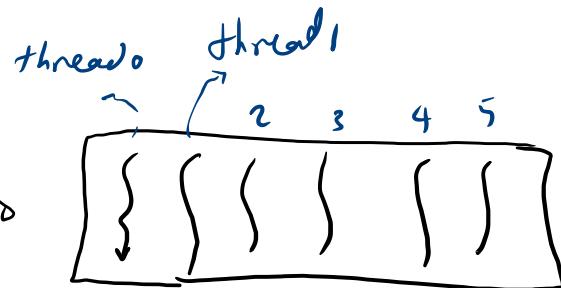
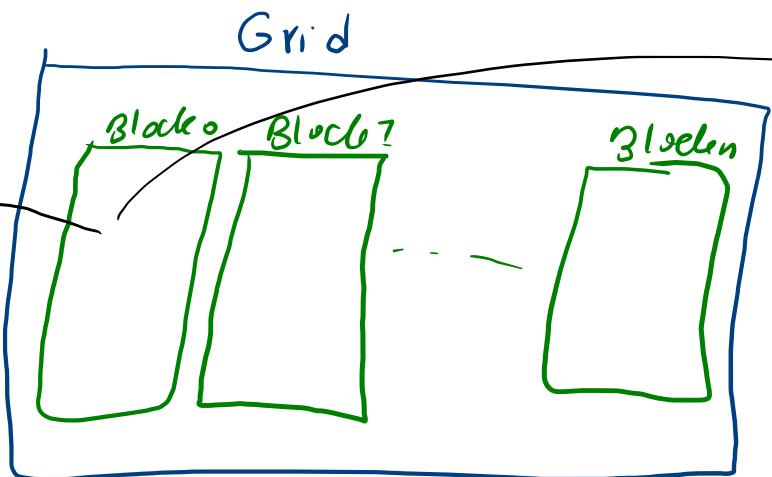
-- global -- void kernel (inputs)
name

کورس کیا ہے
کورس کیا ہے

GPU واردی ها را CPU می خوانند.



لینک چند جزوی را در پایه



$$\rightarrow a[0] + b[0] = c[0] \rightsquigarrow TD$$

$$\rightarrow a[1] + b[1] = c[1] \rightsquigarrow T1$$

$$a[2] + b[2] = c[2] \rightsquigarrow T3$$

$$a[3] + b[3] = c[3] \xrightarrow{T_B}$$

1

$$a[6] + b[6] = c[6] \rightsquigarrow T \text{ true}$$

لقد رأى حارس المخابرات
إسراءيل مراقباً في عليهم

Grid

Block 0

Block 1

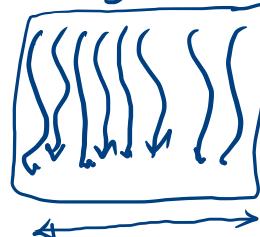
Block 2

Kernel num << 1, 7 >>> ()

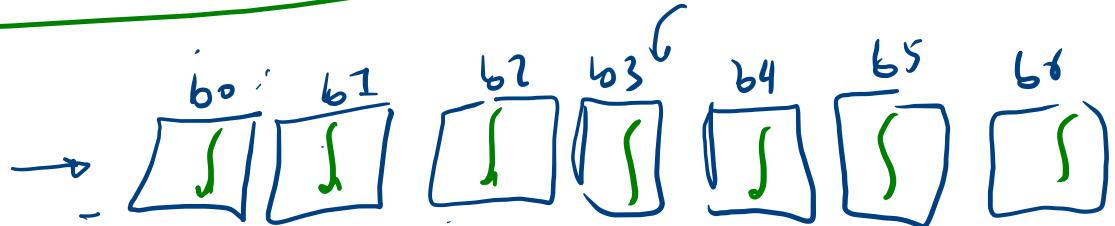
1 2 3 4 5 6 7

grid size is 7

block 0



Kernel num << 7, 1 >>> ()



$\text{KN} \lll 1, 7 \ggg$

$\text{vs} \quad \text{KN} \lll 7, 1 \ggg$



Speed



shared
memory!

local memory

Bleek
Bleek
Bleek
Bleek
Bleek
Bleek
Bleek
Bleek
Bleek

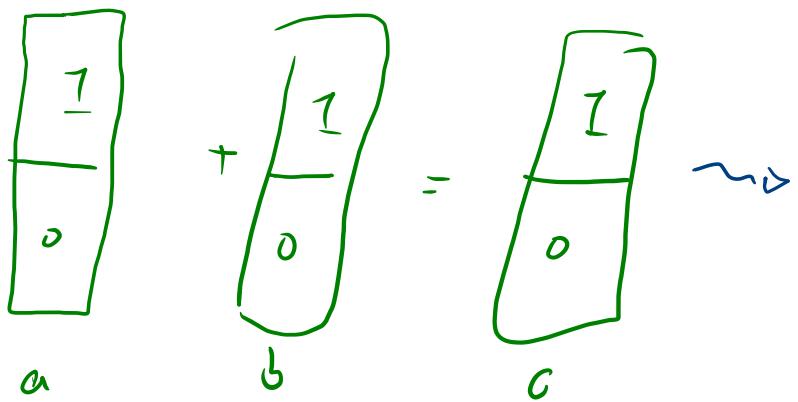


1, R⁴

(M)

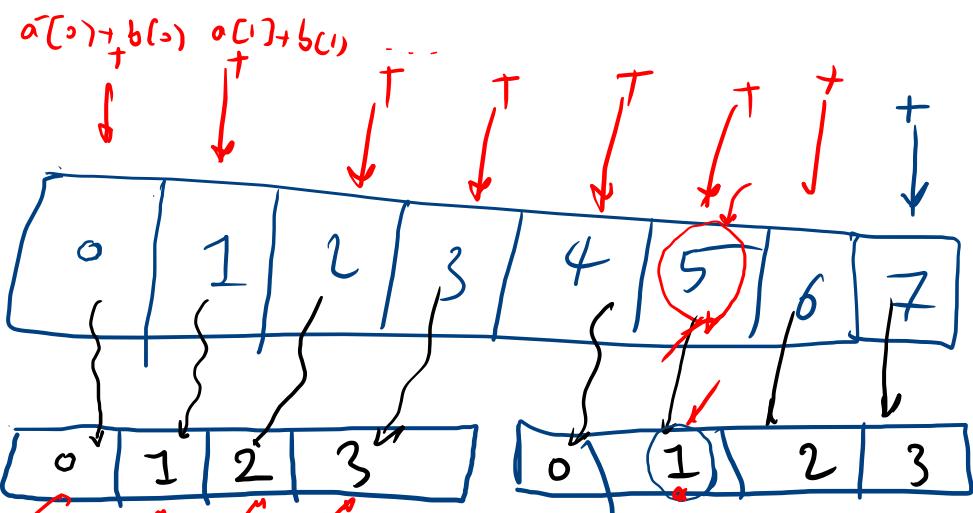
$\frac{512}{8}$

9



$$a[\cdot] + b[\cdot] = c[\cdot] \rightarrow T_0$$

↳ ThreadIdx.x
 ↳ BlockIdx.x.x
 → GridIdx.x.x



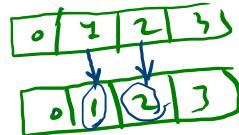
$$\text{BlockIdx.x} = 0$$

$$\text{BlockIdx.x} = 1$$

For $i = 0, j = 0, k = 0 \leftarrow \begin{cases} T: 1 \\ b: 1 \end{cases}$ $\sigma: 0^{15}$

$$\text{Index} = \underline{\text{ThreadIdx.n}} + \underline{\text{BlockIdx.x} * M} = 1 + (1)4 = 5$$

مکالمہ کرنے کا طریقہ



-- global -- void multiply (float * result, float * a, float * b)
{
 int i = threadId * x
 result [i] = a [i] * b [i]
}

int i = blockIdx.x * blockDim.n + ThreadIdx.x

 m

i
 ^
 |
 |

pycuda

size-byte

②

→ `cuda.memalloc(↑)`

→ memalloc GPU

`cuda.memcpy_rectod(gpu-mem, cpu-mem)`

kernel
launched

→ `source model(↓)`

↓
mod. git function (f) Kernel name

Kernel name (inputs, block=())



α