

# GPU Architecture

Pierre Aubert



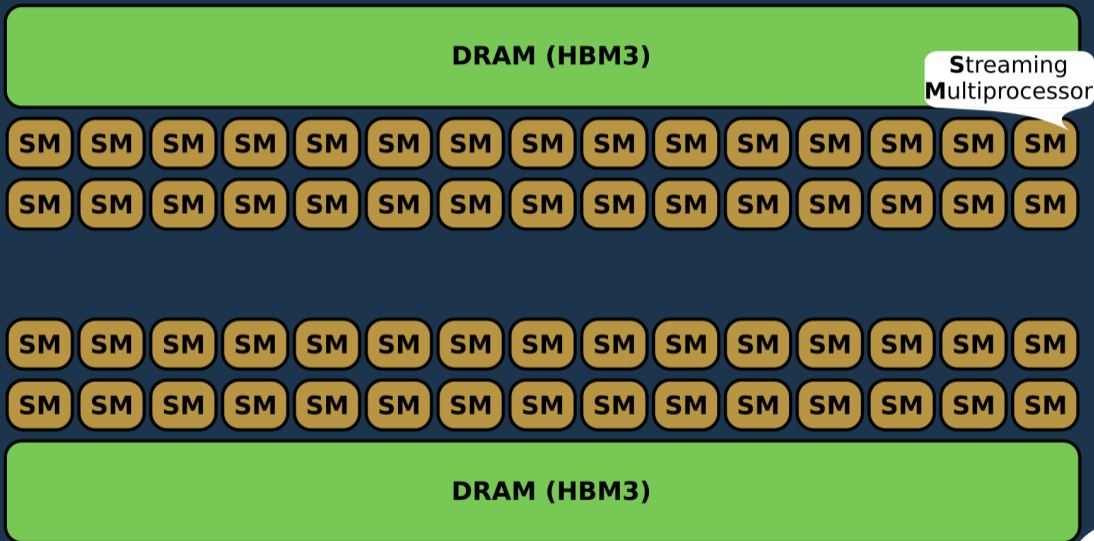




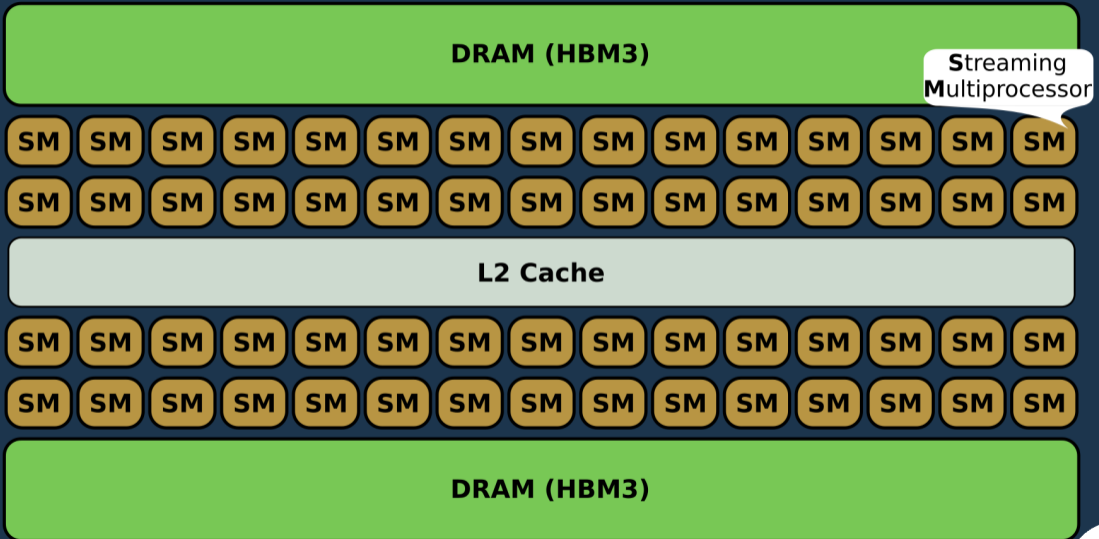
**DRAM (HBM3)**

**DRAM (HBM3)**

# GPU Architecture

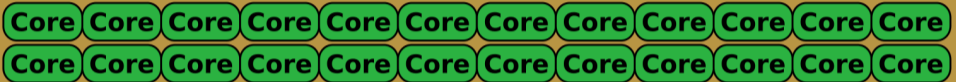


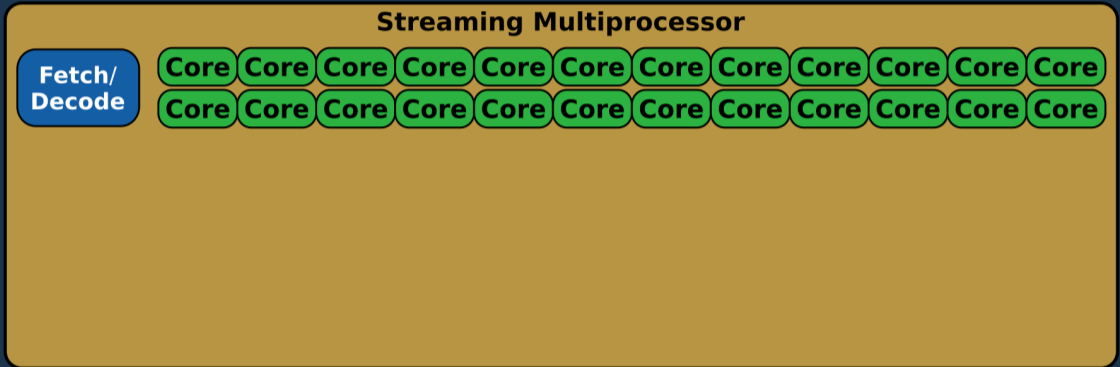
# GPU Architecture

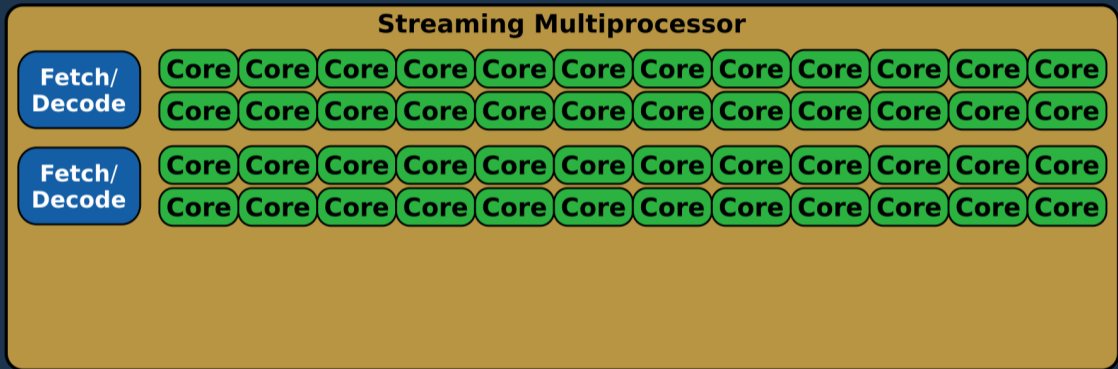


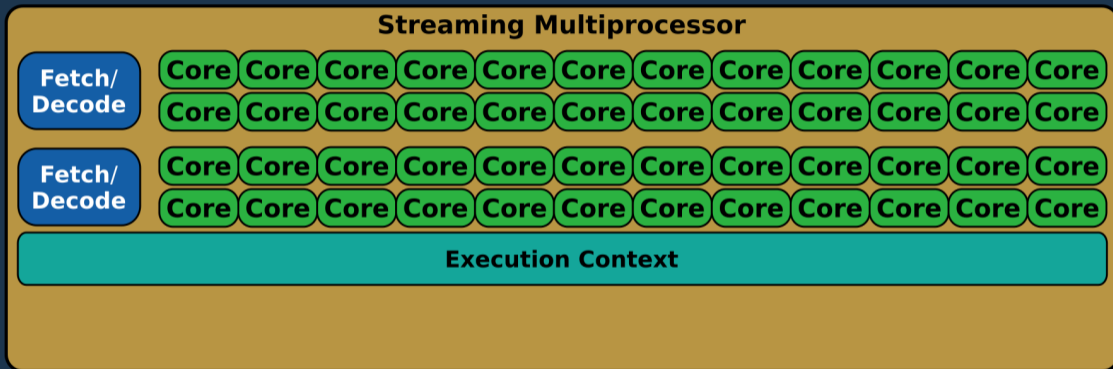
## Streaming Multiprocessor

## Streaming Multiprocessor

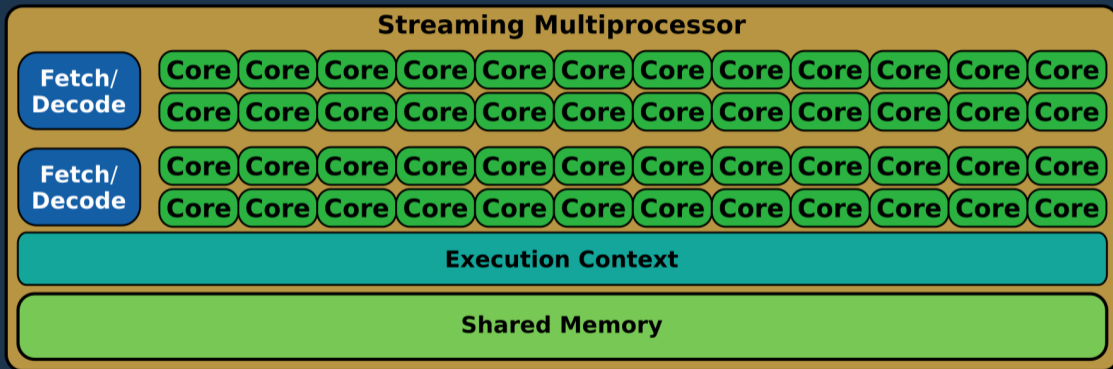












# Streaming Multiprocessor H100

## L1 Instruction Cache



## Tensor Memory Accelerator

## L1 Data Cache / Shared Memory

Tex

Tex

Tex

Tex

# Streaming Multiprocessor H100

## L1 Instruction Cache



## Tensor Memory Accelerator

## L1 Data Cache / Shared Memory

Tex

Tex

Tex

Tex

# Streaming Multiprocessor H100

## L1 Instruction Cache



## Tensor Memory Accelerator

## L1 Data Cache / Shared Memory

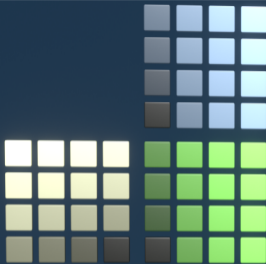
Tex

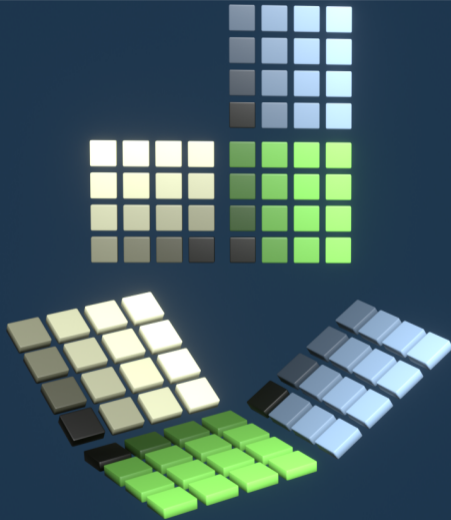
Tex

Tex

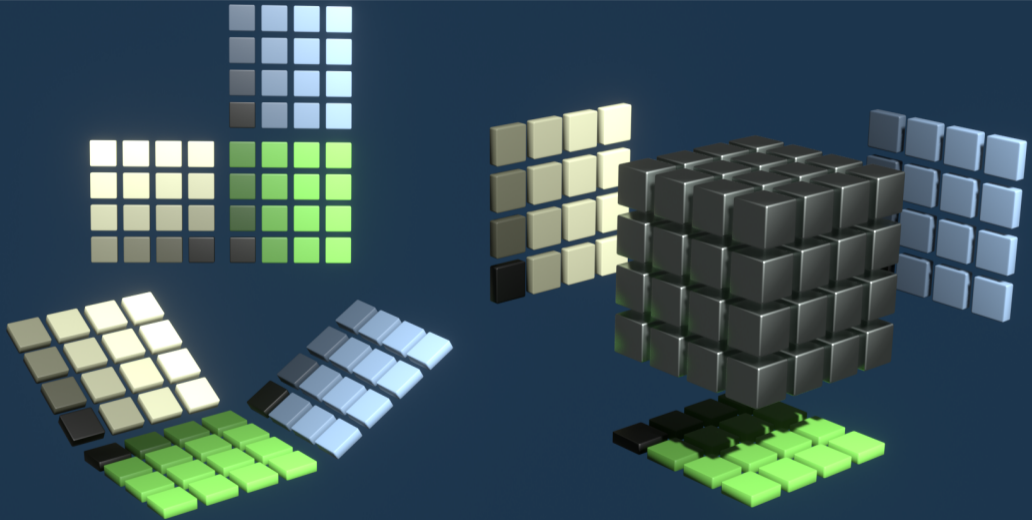
Tex





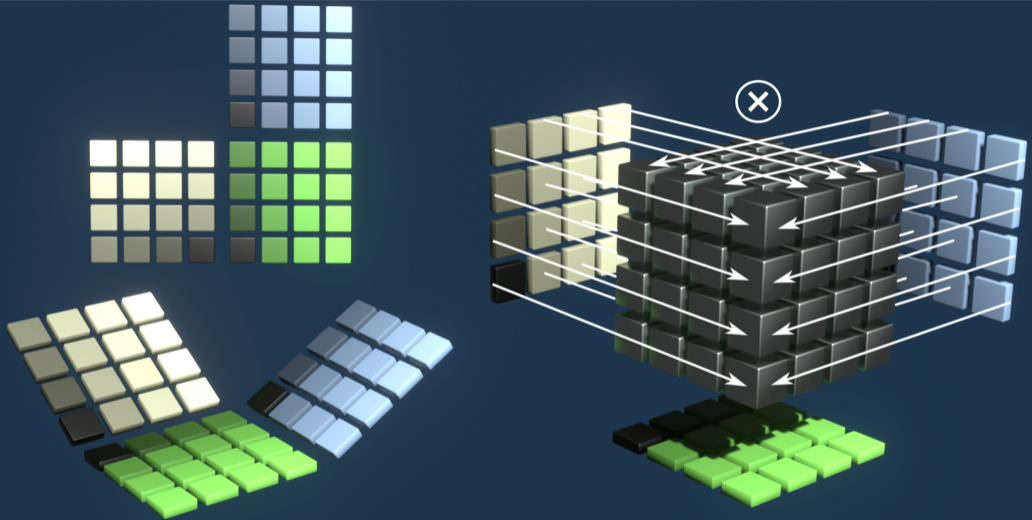


# Tensor Core

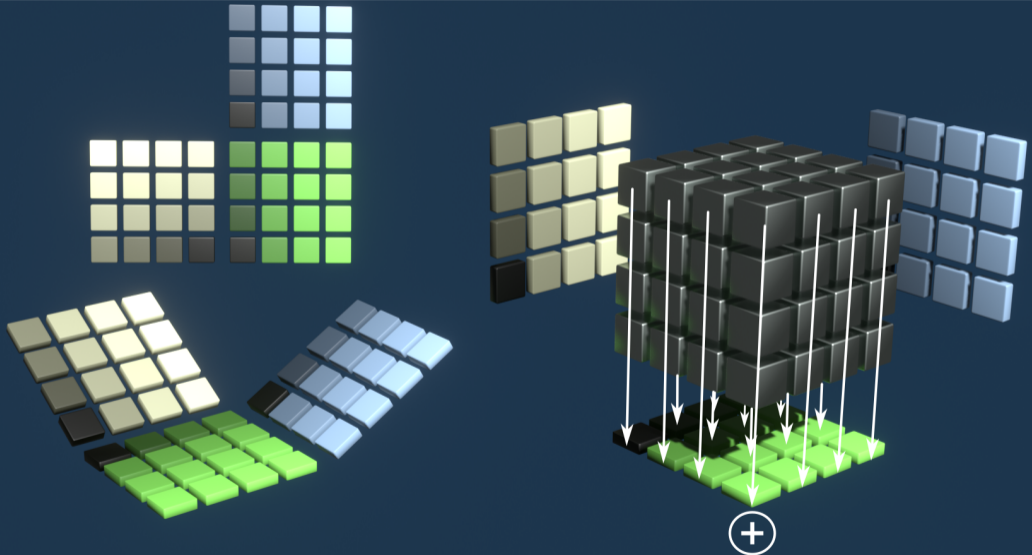




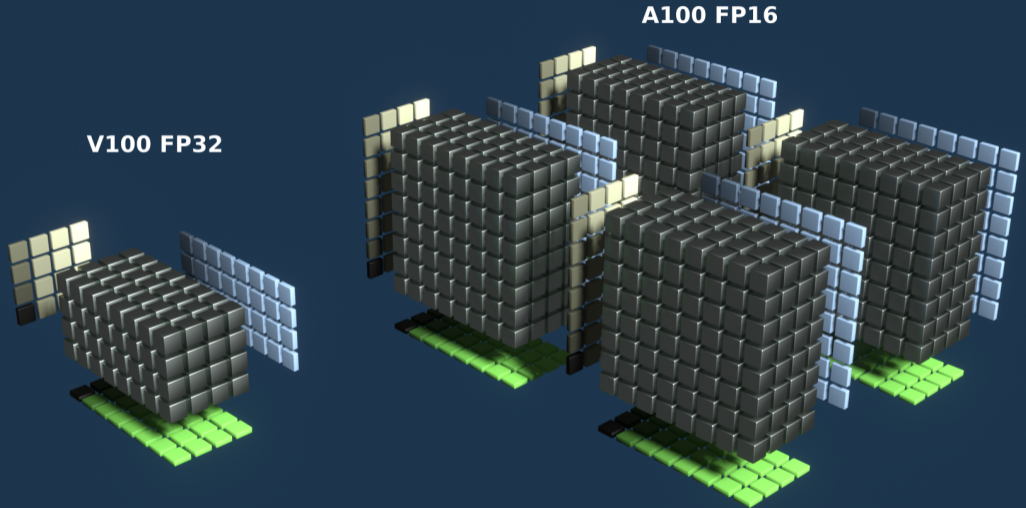
# Tensor Core



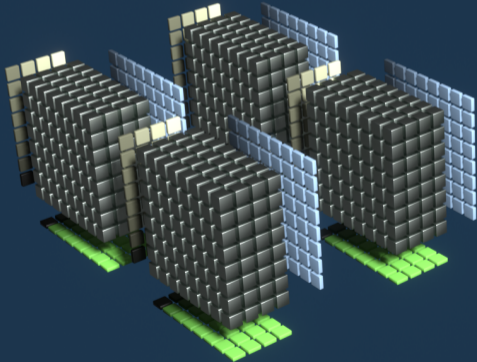
# Tensor Core



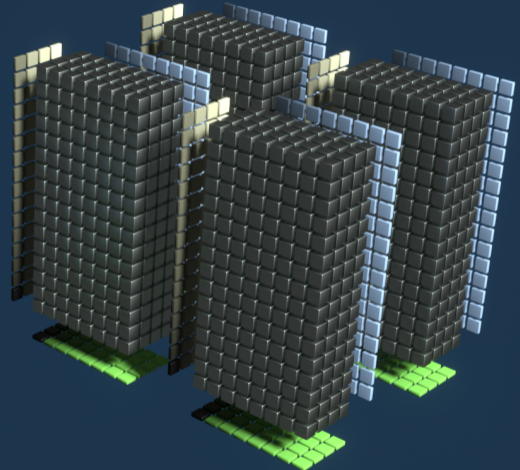
# Tensor Core



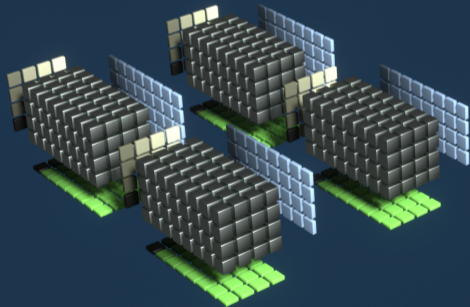
A100 FP16



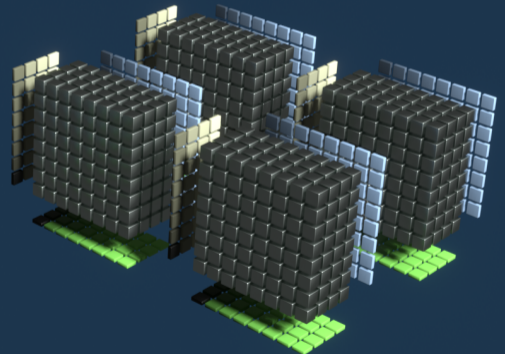
H100 FP16



**A100 TF32**



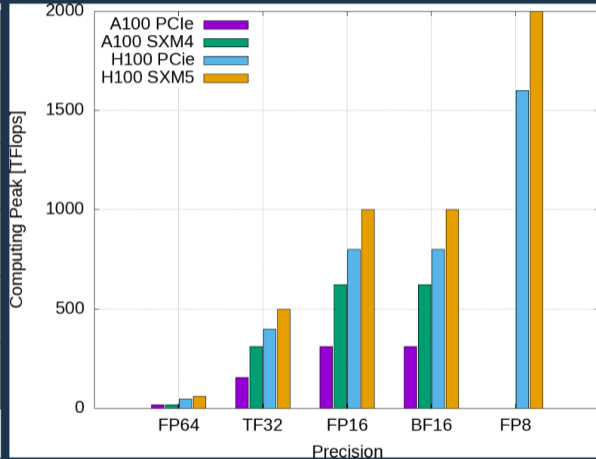
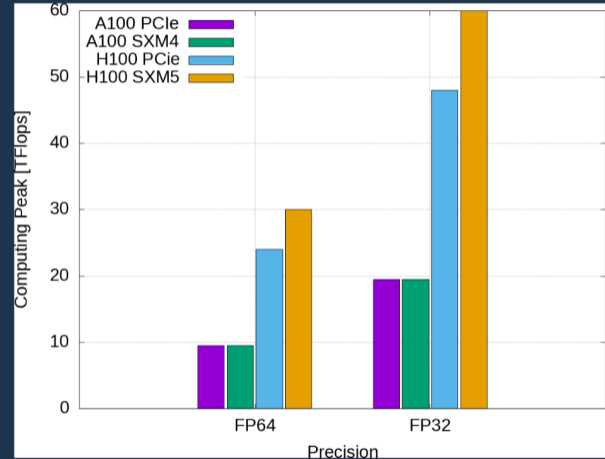
**H100 TF32**



# Why using Tensor Cores ?

## CUDA Cores

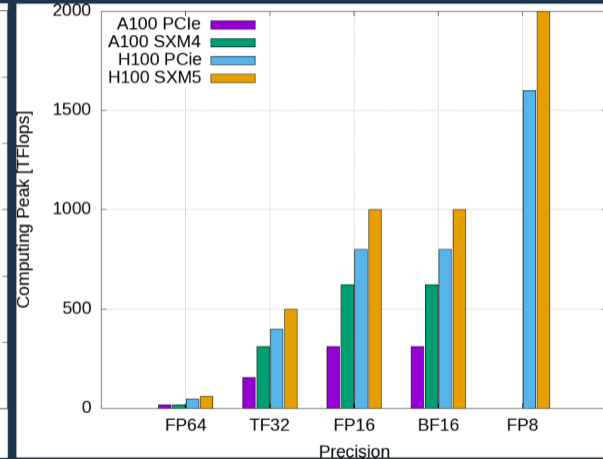
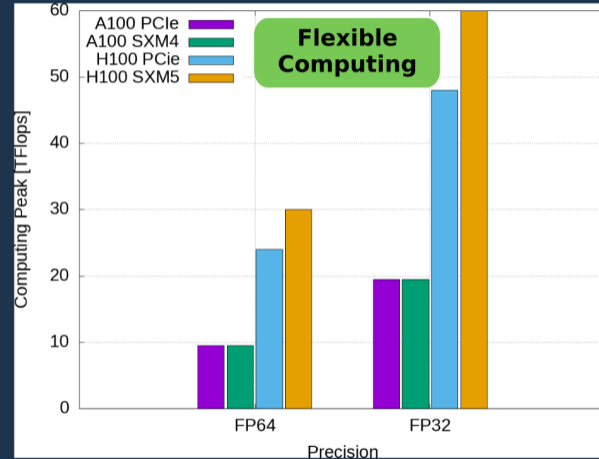
## TensorCores



# Why using Tensor Cores ?

## CUDA Cores

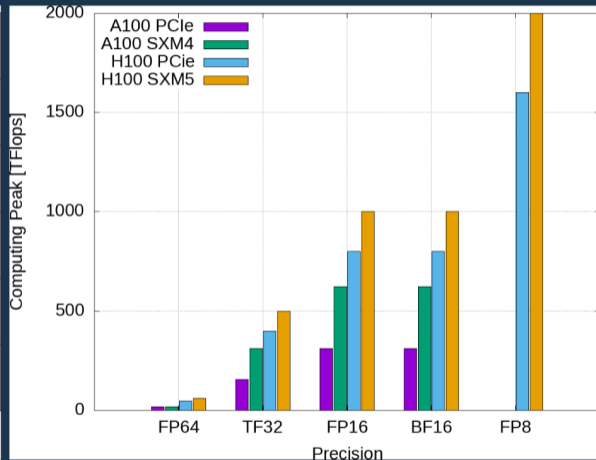
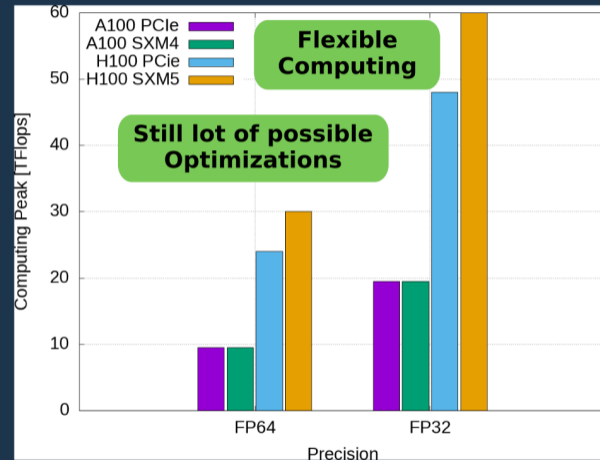
## TensorCores



# Why using Tensor Cores ?

## CUDA Cores

## TensorCores

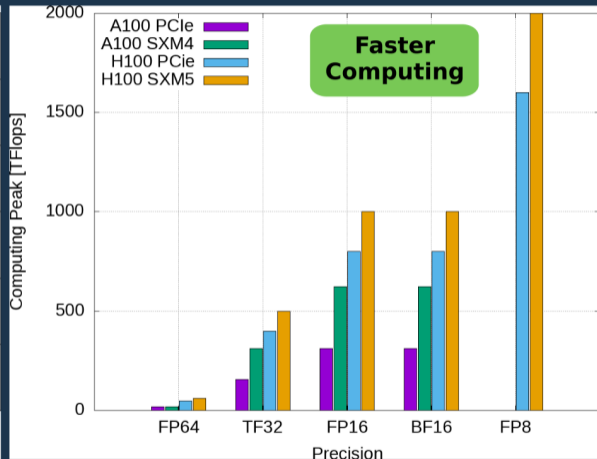
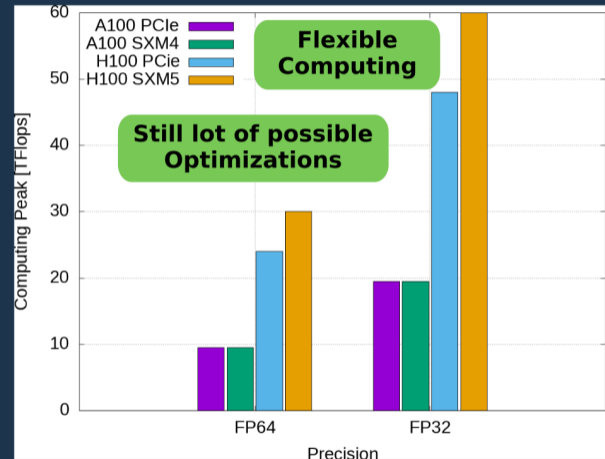




# Why using Tensor Cores ?

## CUDA Cores

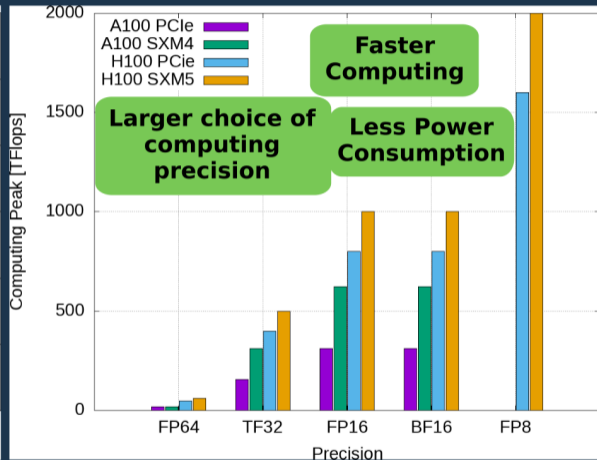
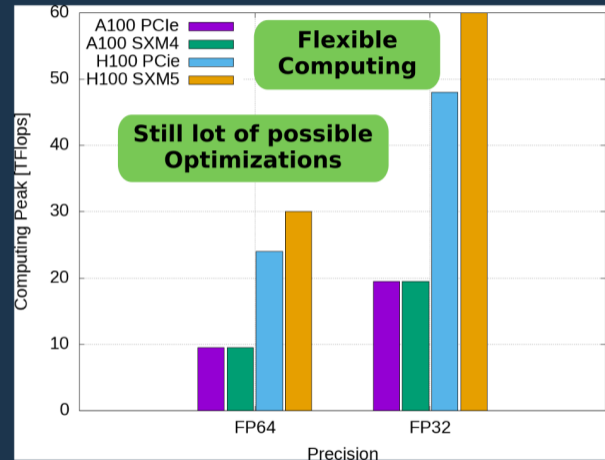
## TensorCores



# Why using Tensor Cores ?

## CUDA Cores

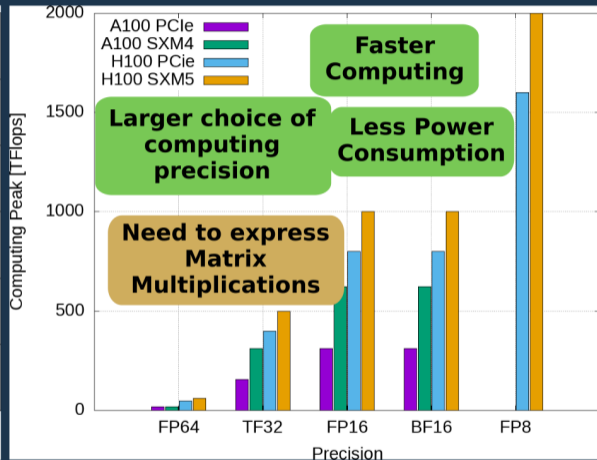
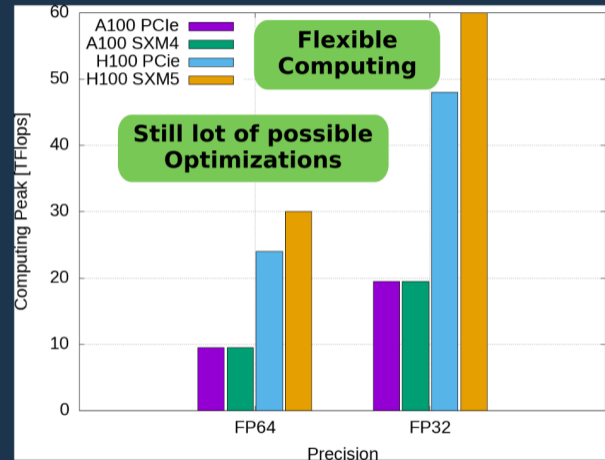
## TensorCores



# Why using Tensor Cores ?

## CUDA Cores

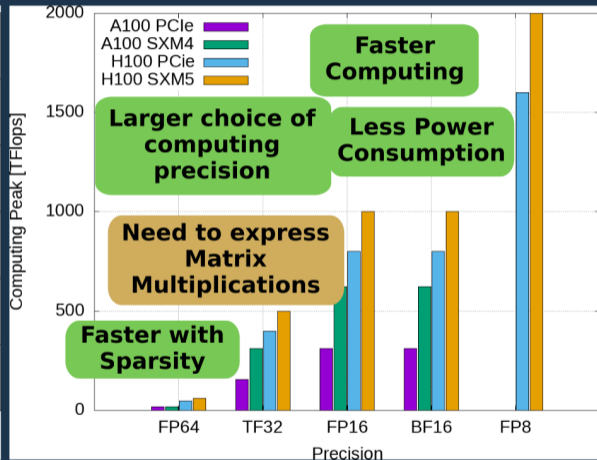
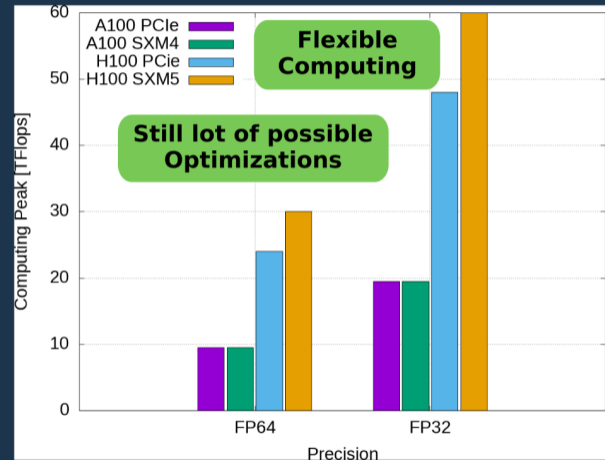
## TensorCores



# Why using Tensor Cores ?

## CUDA Cores

## TensorCores



# The CUDA ecosystem



## Kernel Development

CUDA

cu++

CUTLASS



# The CUDA ecosystem

## Computing Libraries

BLAS

LAPACK

PBLAS

SCALAPACK

TENSOR

SPARSE

RAND

FFT

## Kernel Development

CUDA

cu++

CUTLASS



# The CUDA ecosystem

## Computing Libraries

BLAS

LAPACK

PBLAS

SCALAPACK



TENSOR

SPARSE

RAND

FFT

## Kernel Development

CUDA

cu++

CUTLASS





# The CUDA ecosystem

## Computing Libraries

BLAS

LAPACK

PBLAS

SCALAPACK



TENSOR

SPARSE

RAND

FFT

## Computing + I/O

### RAPIDS

cunumeric

legate

cuDF

## Kernel Development

CUDA

cu++

CUTLASS



# The CUDA ecosystem

## Computing Libraries

BLAS

LAPACK

PBLAS

SCALAPACK



TENSOR

SPARSE

RAND

FFT

## Computing + I/O

## Kernel Development

CUDA

cu++

CUTLASS

## RAPIDS

cunumeric

legate

cuDF



## I/O and communication

## MagnumIO

NVShmem

GPUDirect

NCCL

# The CUDA ecosystem

## Computing Libraries

BLAS

LAPACK

PBLAS

SCALAPACK



TENSOR

SPARSE

RAND

FFT

## Kernel Development

CUDA

cu++

CUTLASS

## Computing + I/O

### RAPIDS

cunumeric

legate

cuDF



## I/O and communication

### MagnumIO

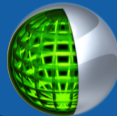
NVShmem

GPUDirect

NCCL

## Profiler / Debugger

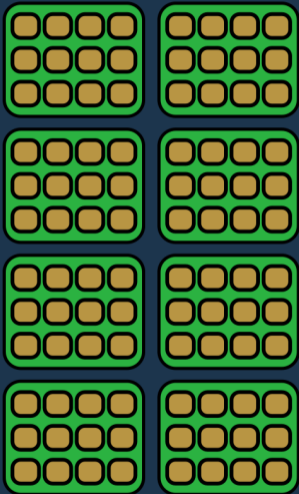
### NSight



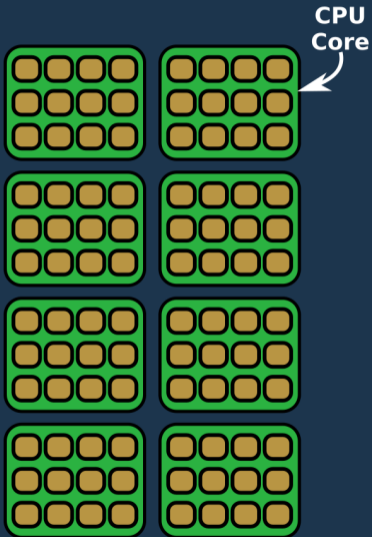
# Threads ?

---

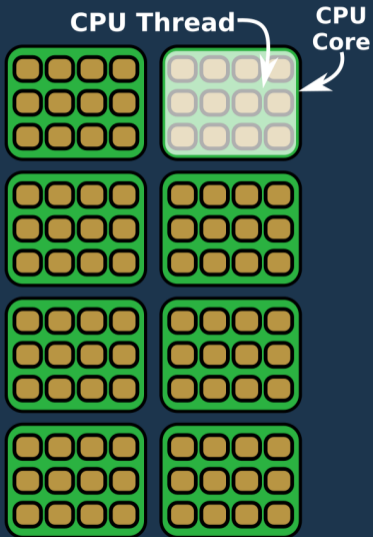
# Threads ?



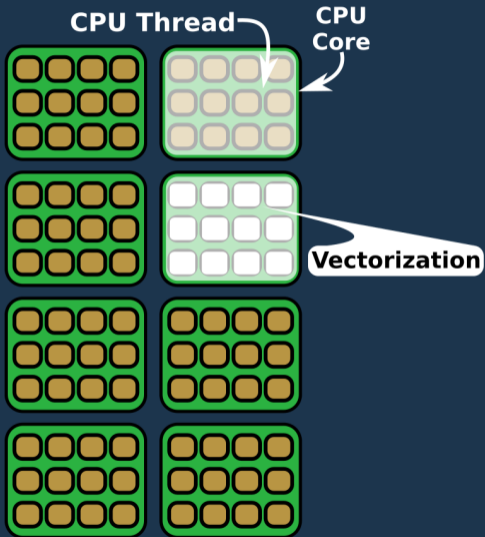
# Threads ?



# Threads ?

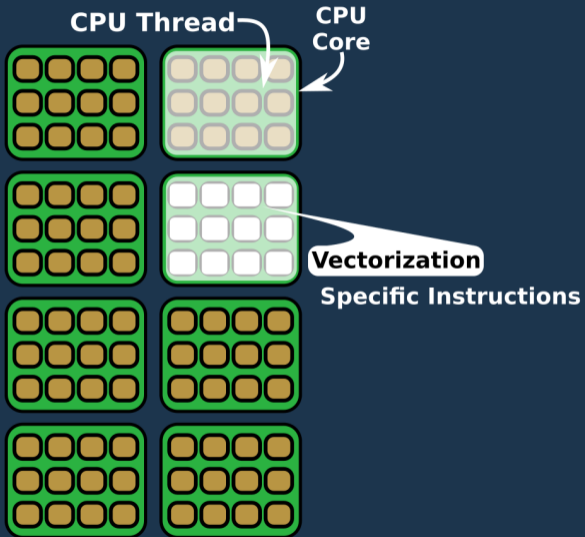


# Threads ?

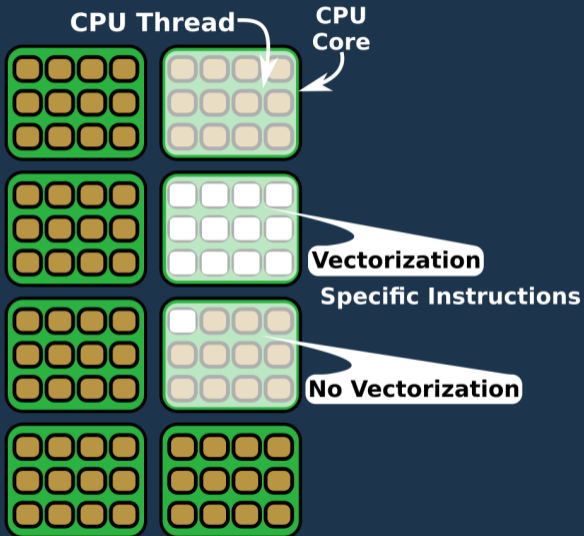




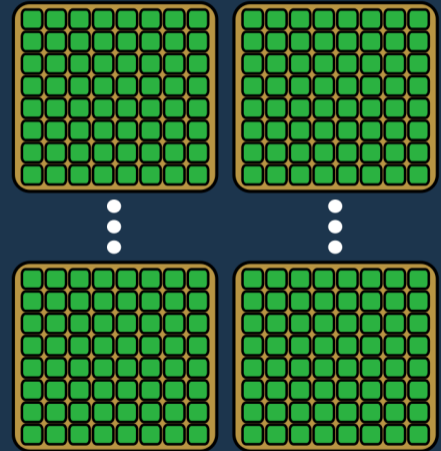
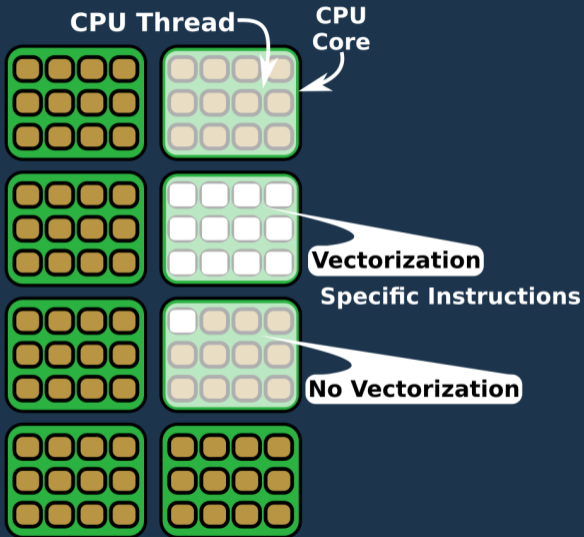
# Threads ?



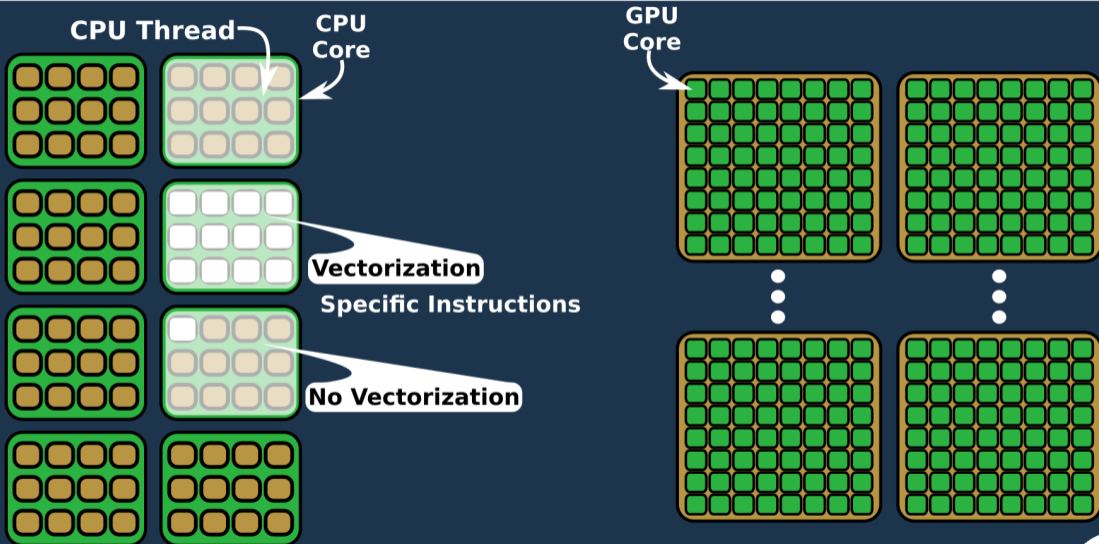
# Threads ?



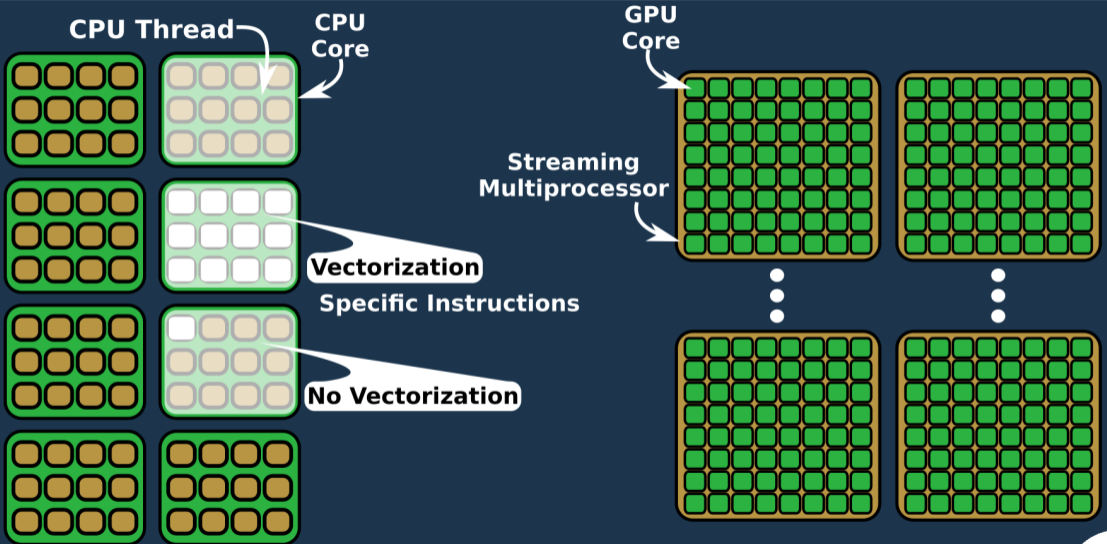
# Threads ?



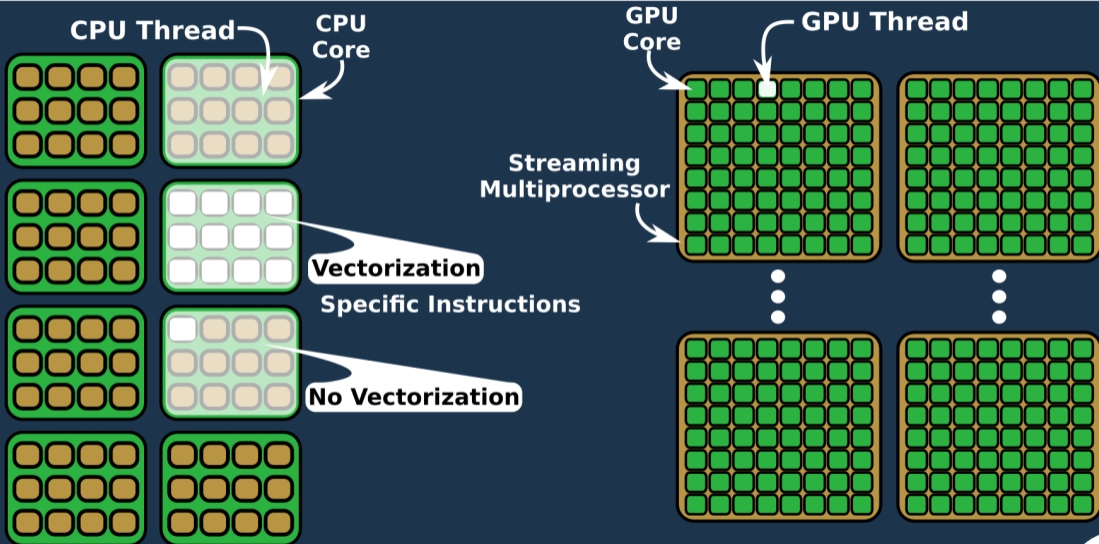
# Threads ?



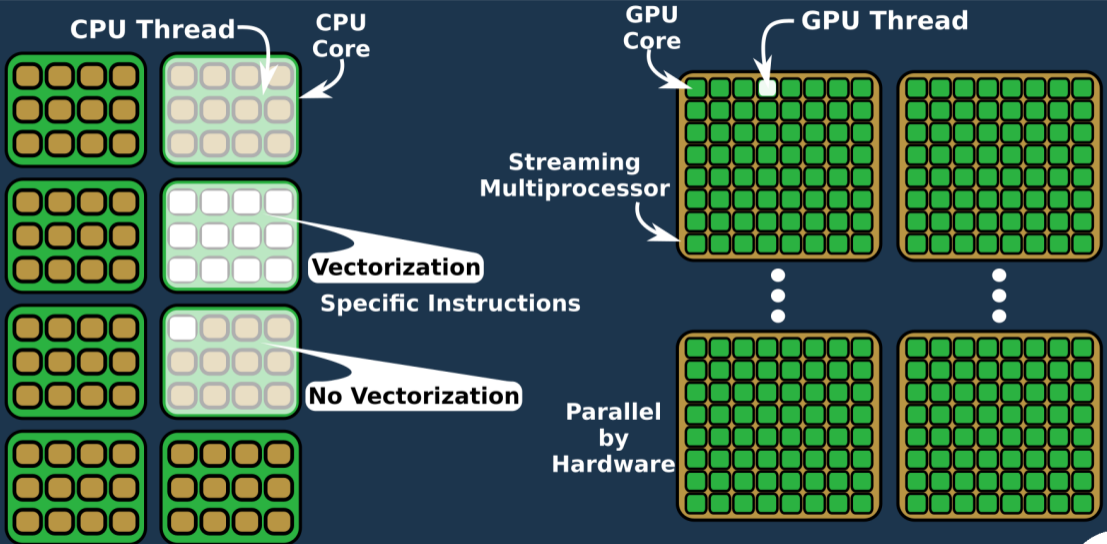
# Threads ?



# Threads ?



# Threads ?



# Threads Hierarchy

---



# Threads Hierarchy

---



**Thread**

# Threads Hierarchy



**Thread**

Core

# Threads Hierarchy



**Thread**

Core

Executes Computing Kernel

# Threads Hierarchy



Warp



Thread

Core

Executes Computing Kernel

# Threads Hierarchy



**Warp**

Hardware  
Scheduler

**Thread**

Core

Executes Computing Kernel

# Threads Hierarchy



**Warp**

Hardware  
Scheduler

32 **Threads**

**Thread**

Core

Executes Computing Kernel

# Threads Hierarchy



**Warp**

Hardware Scheduler

32 **Threads**

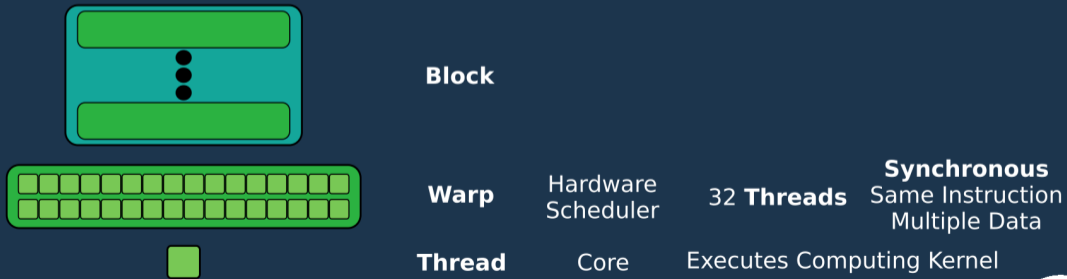
**Synchronous**  
Same Instruction  
Multiple Data

**Thread**

Core

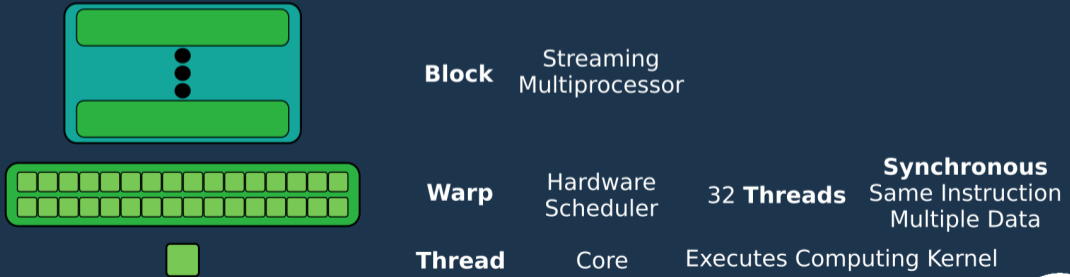
Executes Computing Kernel

# Threads Hierarchy





# Threads Hierarchy



# Threads Hierarchy



**Block** Streaming Multiprocessor 2 - 4 **Warps**  
64 - 128 **Cores**

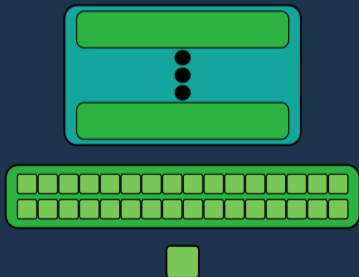


**Warp** Hardware Scheduler 32 **Threads** **Synchronous**  
Same Instruction  
Multiple Data



**Thread** Core Executes Computing Kernel

# Threads Hierarchy

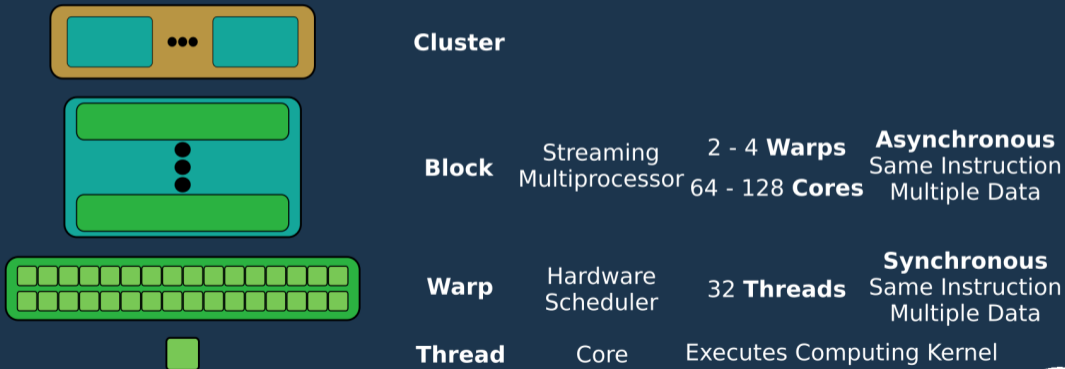


**Block** Streaming Multiprocessor 2 - 4 **Warps** 64 - 128 **Cores** **Asynchronous** Same Instruction Multiple Data

**Warp** Hardware Scheduler 32 **Threads** **Synchronous** Same Instruction Multiple Data

**Thread** Core Executes Computing Kernel

# Threads Hierarchy



# Threads Hierarchy



**Cluster** Blocks  
Communication



**Block** Streaming Multiprocessor 2 - 4 **Warps**  
64 - 128 **Cores** **Asynchronous**  
Same Instruction  
Multiple Data



**Warp** Hardware Scheduler 32 **Threads** **Synchronous**  
Same Instruction  
Multiple Data



**Thread** Core Executes Computing Kernel

# Threads Hierarchy



**Cluster**

Blocks  
Communication

Since :  
- CUDA 12  
- **Hopper**



**Block**

Streaming  
Multiprocessor

2 - 4 **Warps**  
64 - 128 **Cores**

**Asynchronous**  
Same Instruction  
Multiple Data



**Warp**

Hardware  
Scheduler

32 **Threads**

**Synchronous**  
Same Instruction  
Multiple Data



**Thread**

Core

Executes Computing Kernel

# Threads Hierarchy



**Cluster**

Blocks  
Communication

Since :  
- CUDA 12  
- **Hopper**

**Hopper** : 132 SMs  
**Ampere** : 108 SMs



**Block**

Streaming  
Multiprocessor

2 - 4 **Warps**  
64 - 128 **Cores**

**Asynchronous**  
Same Instruction  
Multiple Data



**Warp**

Hardware  
Scheduler

32 **Threads**

**Synchronous**  
Same Instruction  
Multiple Data



**Thread**

Core

Executes Computing Kernel

# Threads Hierarchy



**Grid**



**Cluster**

Blocks  
Communication

Since :  
- CUDA 12  
- **Hopper**

**Hopper** : 132 SMs  
**Ampere** : 108 SMs



**Block**

Streaming  
Multiprocessor

2 - 4 **Warps**  
64 - 128 **Cores**

**Asynchronous**  
Same Instruction  
Multiple Data



**Warp**

Hardware  
Scheduler

32 **Threads**

**Synchronous**  
Same Instruction  
Multiple Data



**Thread**

Core

Executes Computing Kernel



# Threads Hierarchy



**Grid**

Full **GPU**



**Cluster**

Blocks  
Communication

Since :  
- **CUDA 12**  
- **Hopper**

**Hopper** : 132 SMs  
**Ampere** : 108 SMs



**Block**

Streaming  
Multiprocessor

2 - 4 **Warps**  
64 - 128 **Cores**

**Asynchronous**  
Same Instruction  
Multiple Data



**Warp**

Hardware  
Scheduler

32 **Threads**

**Synchronous**  
Same Instruction  
Multiple Data



**Thread**

Core

Executes Computing Kernel

# Threads Hierarchy



**Grid**

Full **GPU**

**Hopper** : 8448 Cores  
**Ampere** : 6912 Cores



**Cluster**

Blocks  
Communication

Since :  
- **CUDA 12**  
- **Hopper**

**Hopper** : 132 SMs  
**Ampere** : 108 SMs



**Block**

Streaming  
Multiprocessor

2 - 4 **Warps**  
64 - 128 **Cores**

**Asynchronous**  
Same Instruction  
Multiple Data



**Warp**

Hardware  
Scheduler

32 **Threads**

**Synchronous**  
Same Instruction  
Multiple Data



**Thread**

Core

Executes Computing Kernel



 Thread

# Threads Execution

■ Thread



# Threads Execution

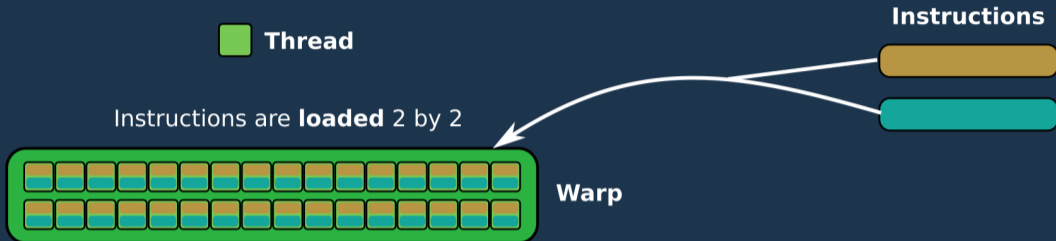
■ Thread

Instructions

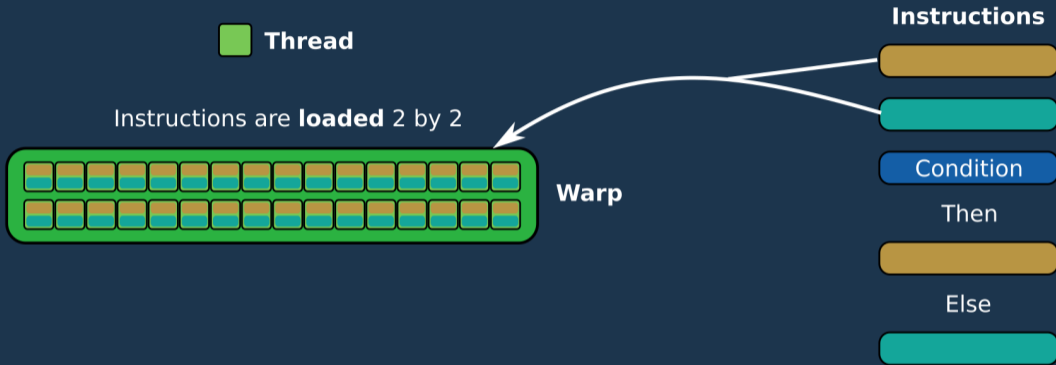


Warp

# Threads Execution

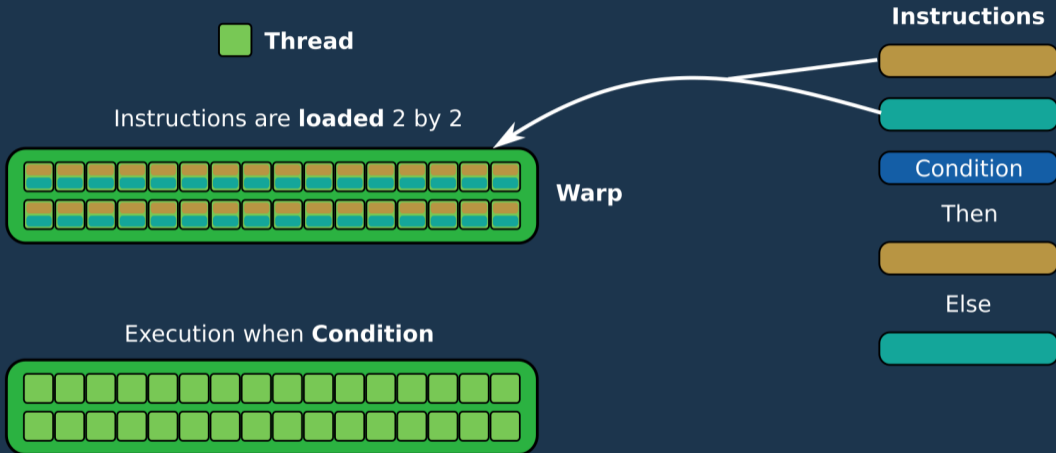


# Threads Execution

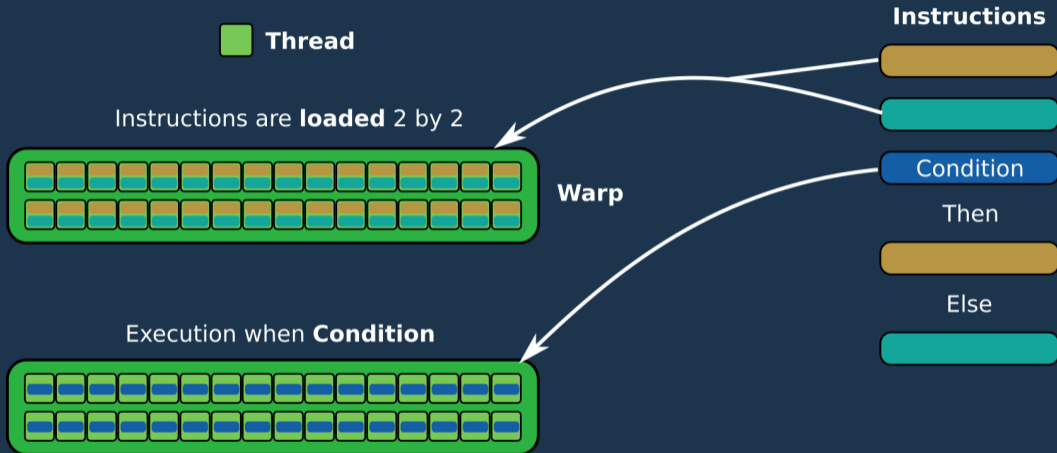




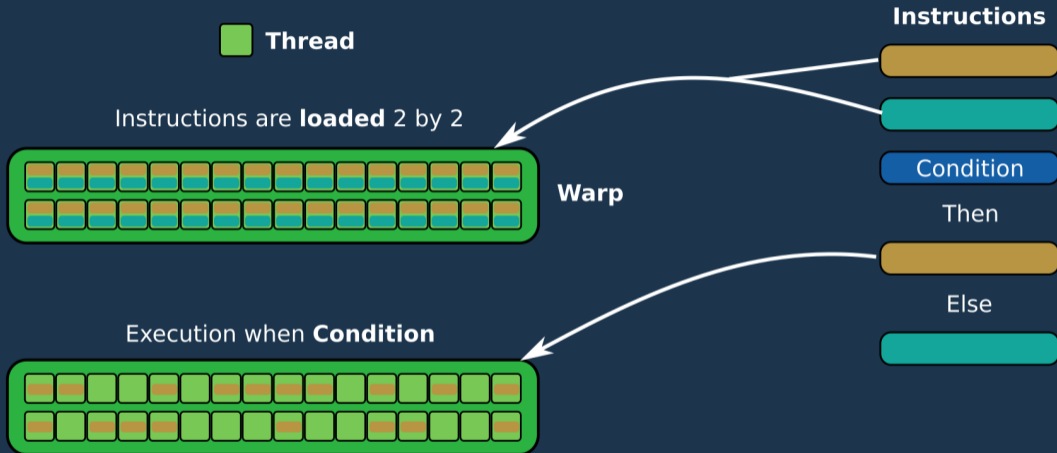
# Threads Execution



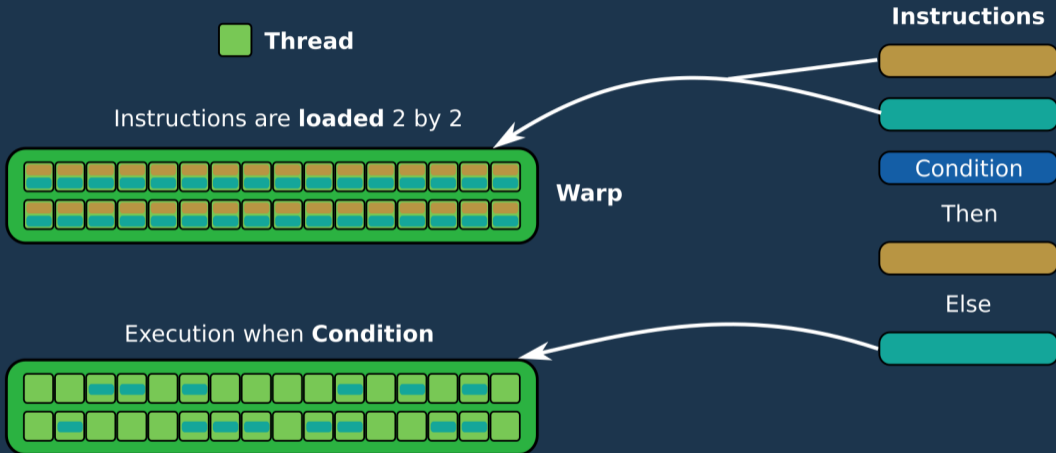
# Threads Execution



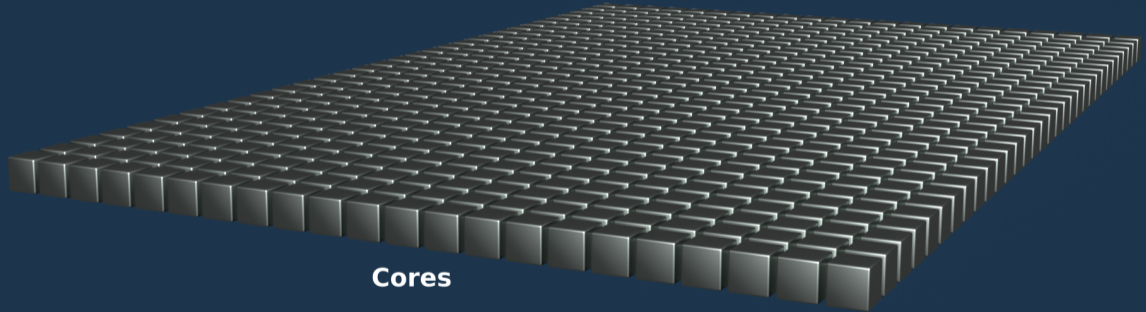
# Threads Execution



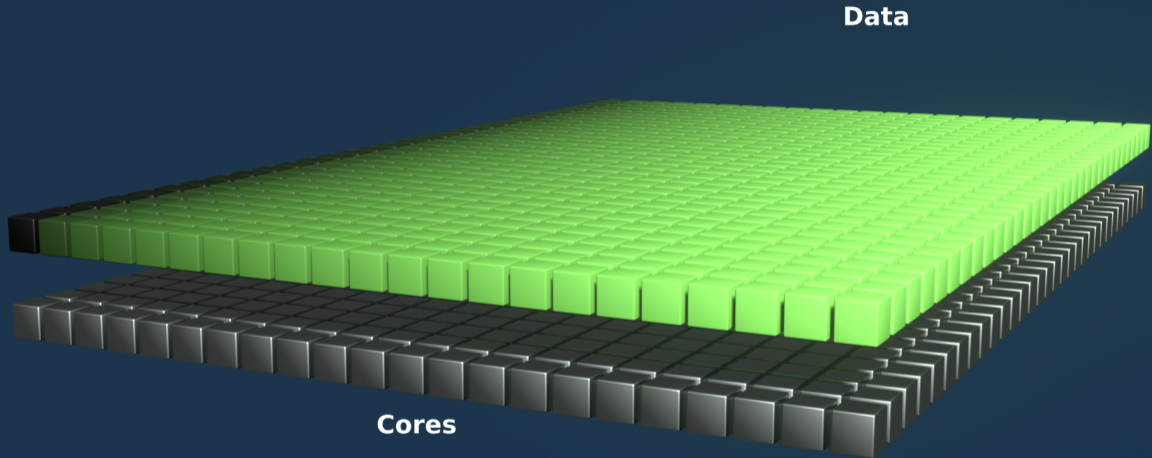
# Threads Execution



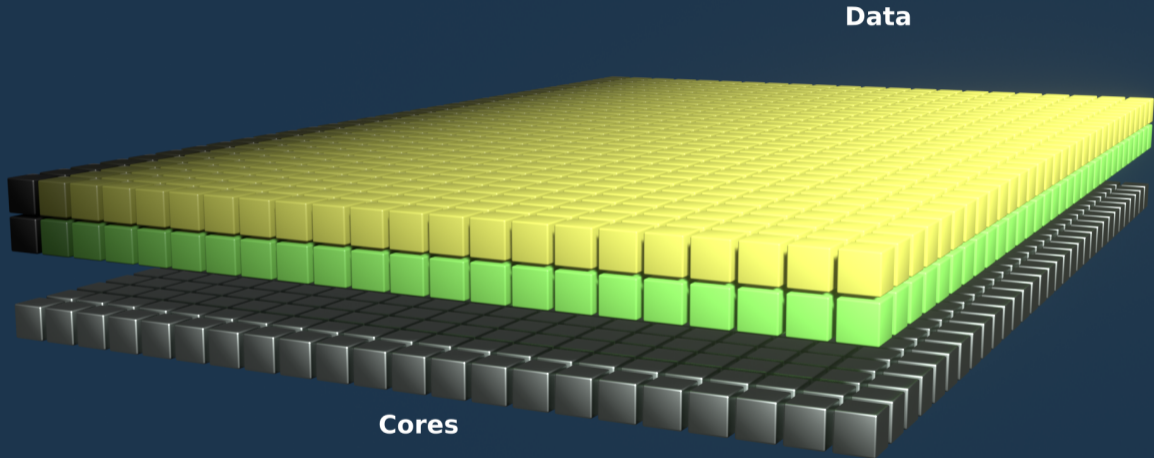
# Thread Mapping



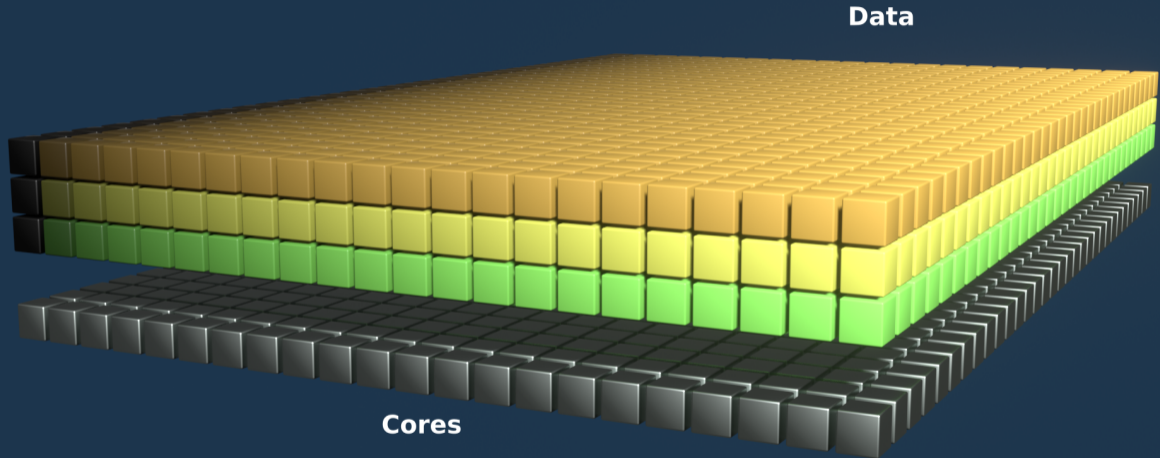
# Thread Mapping



# Thread Mapping

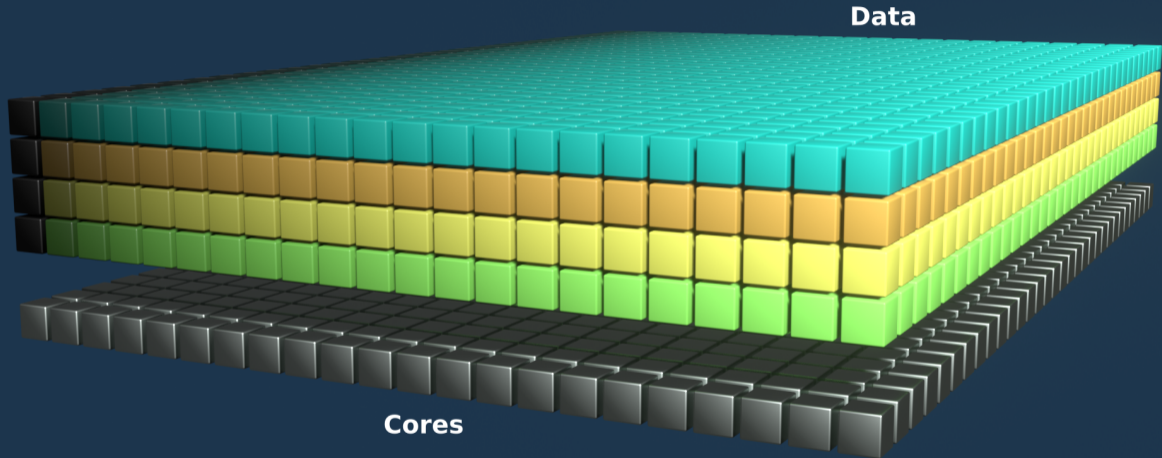


# Thread Mapping

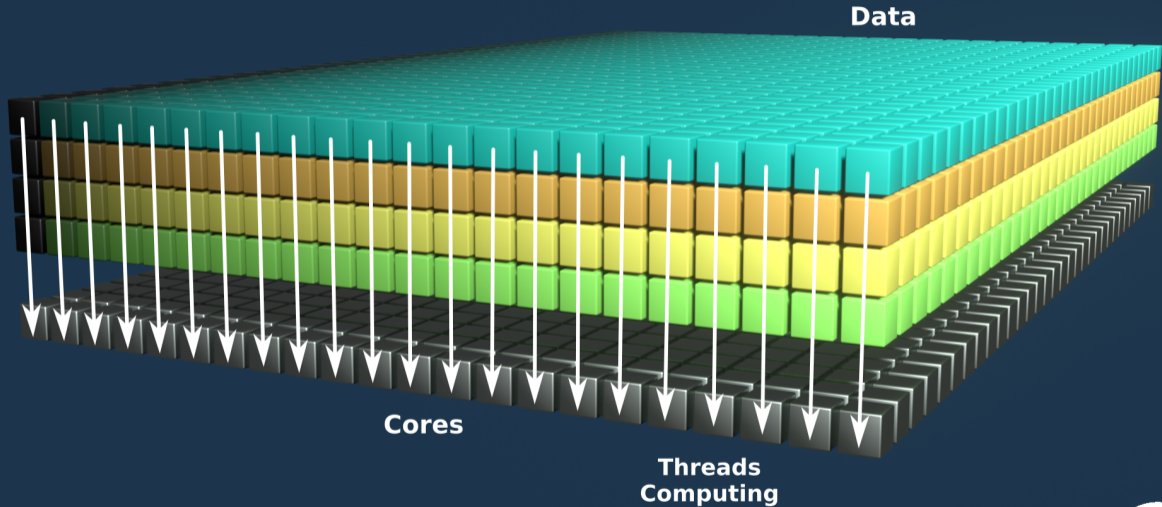




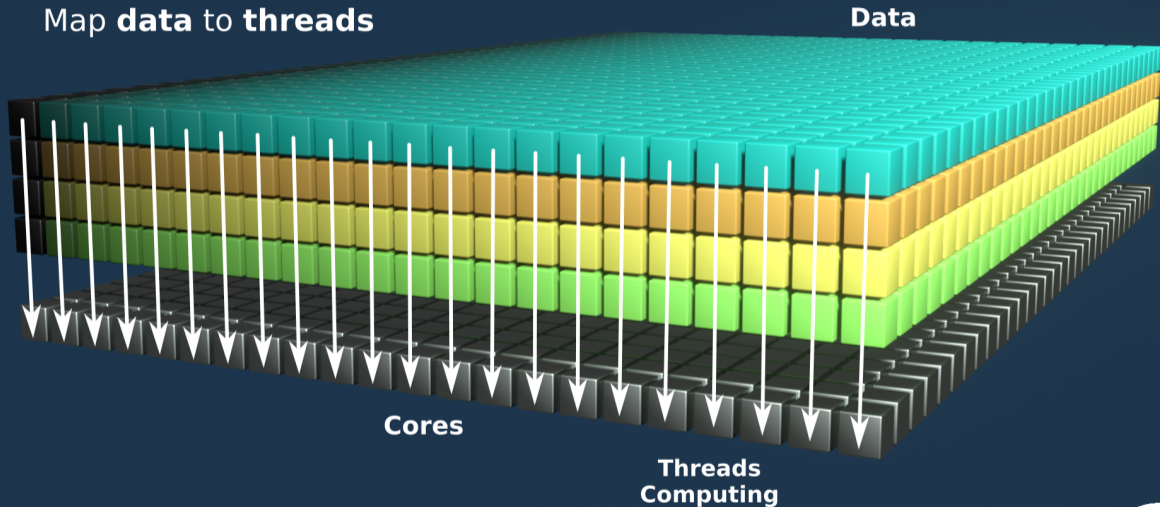
# Thread Mapping



# Thread Mapping



# Thread Mapping

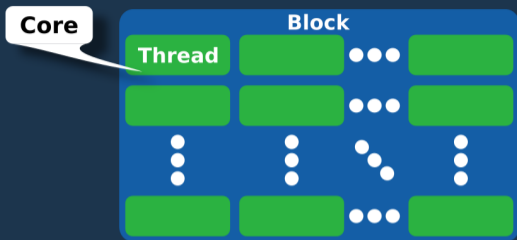




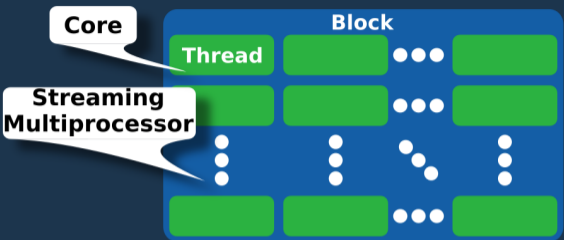
Thread



# Data / Thread mapping



# Data / Thread mapping





# Data / Thread mapping

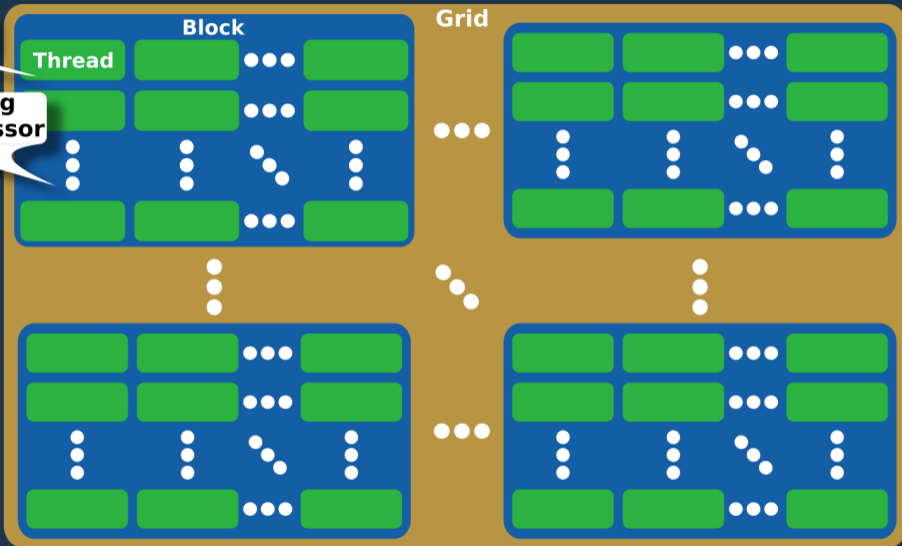
Core

Block

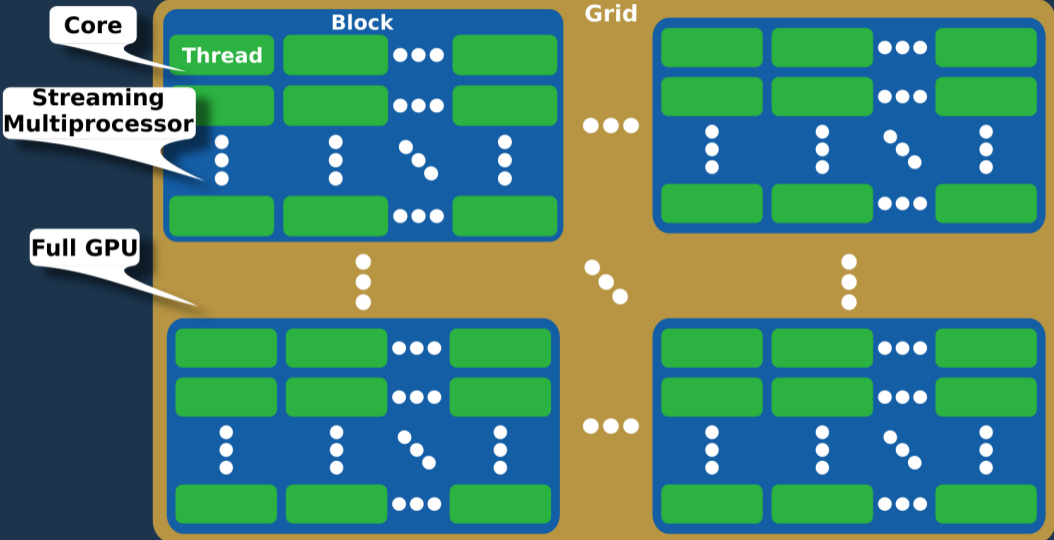
Grid

Thread

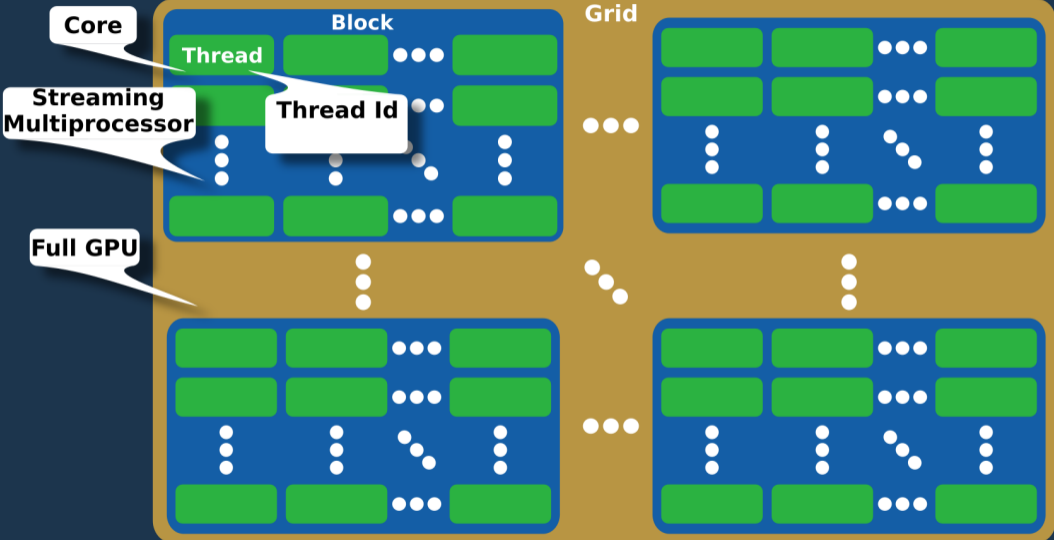
Streaming  
Multiprocessor



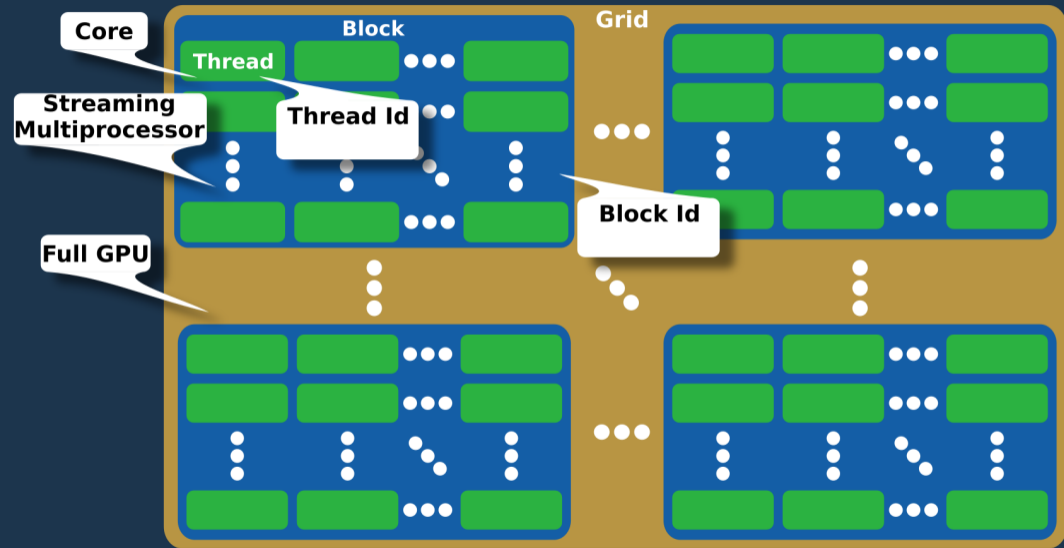
# Data / Thread mapping



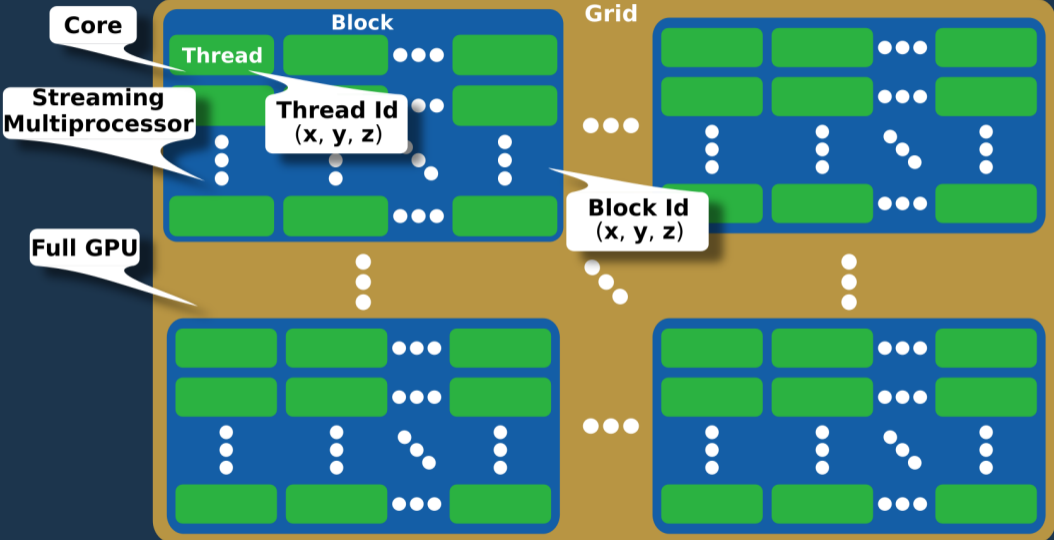
# Data / Thread mapping



# Data / Thread mapping



# Data / Thread mapping



# Block Size ?

---

# Block Size ?

Hadamard product on 1000 elements



# Block Size ?

Hadamard product on 1000 elements

Scenario #1

10 1d blocks of size 100





# Block Size ?

Warp of 32 threads



Scenario #1

10 1d blocks of size 100

Hadamard product on 1000 elements



# Block Size ?

Warp of 32 threads



Scenario #1

10 1d blocks of size 100

Hadamard product on 1000 elements



# Block Size ?

Warp of 32 threads



Hadamard product on 1000 elements



Use: 78%

Scenario #1

10 1d blocks of size 100



# Block Size ?

Warp of 32 threads



Scenario #1

10 1d blocks of size 100

Use: 78%

Hadamard product on 1000 elements



Scenario #2

8 1d blocks of size 128



# Block Size ?

Warp of 32 threads

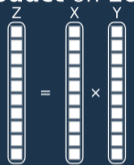


Scenario #1

10 1d blocks of size 100

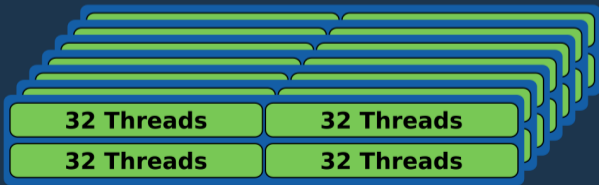
Use: 78%

Hadamard product on 1000 elements



Scenario #2

8 1d blocks of size 128



# Block Size ?

Warp of 32 threads



Scenario #1

10 1d blocks of size 100

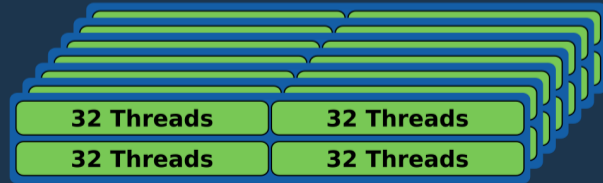
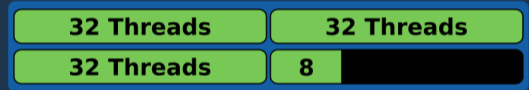
Use: 78%

Hadamard product on 1000 elements



Scenario #2

8 1d blocks of size 128



# Block Size ?

Warp of 32 threads



Scenario #1

10 1d blocks of size 100

Use: 78%

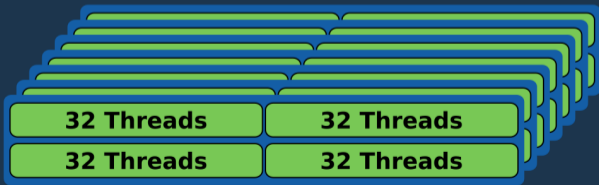
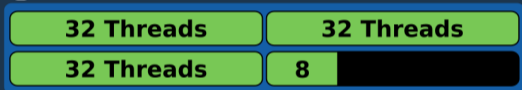
Hadamard product on 1000 elements



Scenario #2

8 1d blocks of size 128

Use: 98%



# Warp Read / Write

---



**1 sector : 32 bytes**

**Cache line size :**

- **128 bytes**
- **4 sectors**

# Warp Read / Write

1 sector : 32 bytes

4-bytes elements access -> 4 sectors

Cache line size :

- 128 bytes

- 4 sectors

Warp of 32 threads



# Warp Read / Write

1 sector : 32 bytes

Cache line size :

- 128 bytes

- 4 sectors

4-bytes elements access -> 4 sectors

Warp of 32 threads



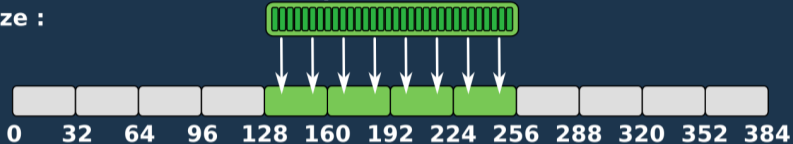
# Warp Read / Write

1 sector : 32 bytes

Cache line size :  
- 128 bytes  
- 4 sectors

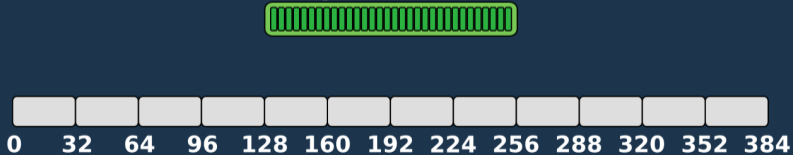
4-bytes elements access -> 4 sectors

Warp of 32 threads



8-bytes elements access -> 8 sectors

Warp of 32 threads



# Warp Read / Write

1 sector : 32 bytes

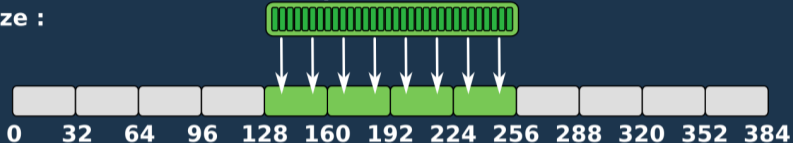
Cache line size :

- 128 bytes

- 4 sectors

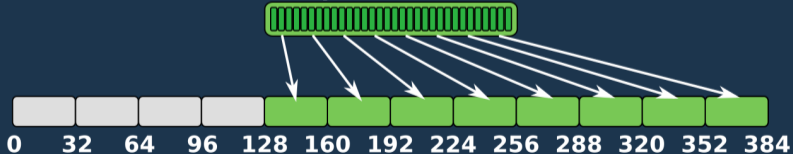
4-bytes elements access -> 4 sectors

Warp of 32 threads



8-bytes elements access -> 8 sectors

Warp of 32 threads



# Warp Read / Write

1 sector : 32 bytes

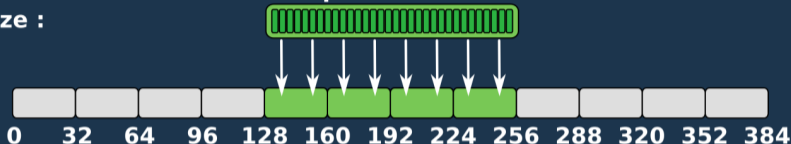
Cache line size :

- 128 bytes

- 4 sectors

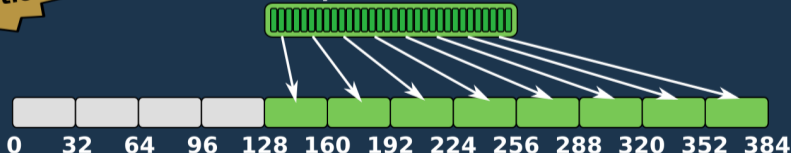
4-bytes elements access -> 4 sectors

Warp of 32 threads



8-bytes elements access -> 8 sectors

Warp of 32 threads



Really close  
to CPU  
Vectorization



**Read / Write :**

- **1 sector**
- **32 bytes**



# Memory Layout

**Read / Write :**

- **1 sector**
- **32 bytes**

**Array of Structure (AoS)**

```
struct Coeff{  
    float u, v, w;  
    float x[8], y[8], z;  
};
```

**80 Bytes**

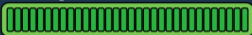
# Memory Layout

Read / Write :

- 1 sector
- 32 bytes

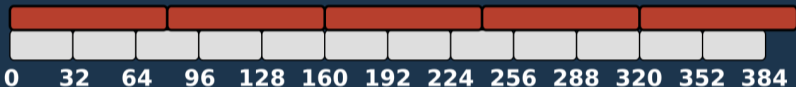
Array of Structure (AoS)

Warp of 32 threads



```
struct Coeff{
    float u, v, w;
    float x[8], y[8], z;
};
```

80 Bytes



# Memory Layout

Read / Write :

- 1 sector
- 32 bytes

Array of Structure (AoS)

Warp of 32 threads

struct Coeff{

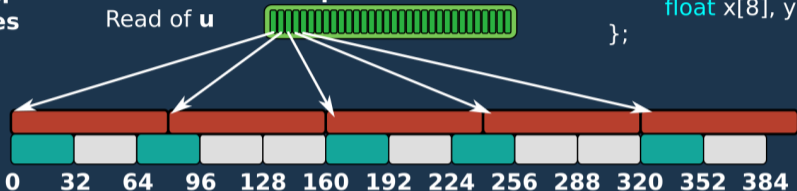
float u, v, w;

float x[8], y[8], z;

};

80 Bytes

Read of u



# Memory Layout

Read / Write :

- 1 sector
- 32 bytes

Array of Structure (AoS)

Warp of 32 threads

struct Coeff{

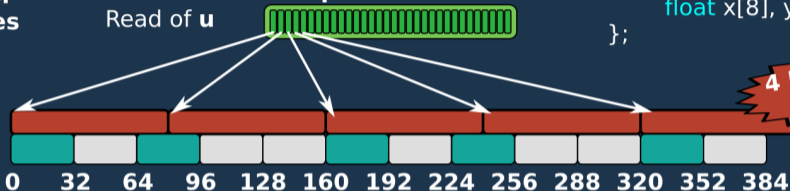
float u, v, w;

float x[8], y[8], z;

};

80 Bytes

Read of u



4 bytes usefull  
out of 32

# Memory Layout

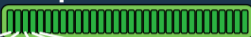
Read / Write :

- 1 sector
- 32 bytes

Array of Structure (AoS)

Warp of 32 threads

Read of u

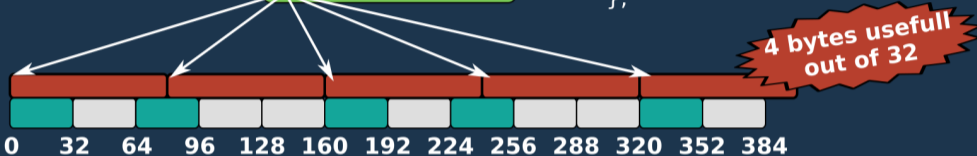


struct Coeff{

```
float u, v, w;
float x[8], y[8], z;
```

80 Bytes

};



Strucutre of Array (SoA)

struct Coeff{

```
float *u, *v, *w;
float *x0, ... *x7, *y0, ... *y7, *z;
```

};

4 bytes usefull out of 32

# Memory Layout

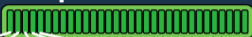
Read / Write :

- 1 sector
- 32 bytes

Array of Structure (AoS)

Warp of 32 threads

Read of u



struct Coeff{

```
float u, v, w;
float x[8], y[8], z;
```

80 Bytes

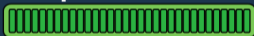
};

4 bytes usefull  
out of 32



Structure of Array (SoA)

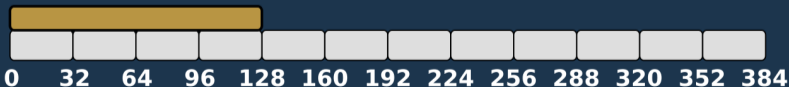
Warp of 32 threads



struct Coeff{

```
float *u, *v, *w;
float *x0, ... *x7, *y0, ... *y7, *z;
```

};



# Memory Layout

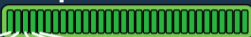
Read / Write :

- 1 sector
- 32 bytes

Array of Structure (AoS)

Warp of 32 threads

Read of u



struct Coeff{

```
float u, v, w;
float x[8], y[8], z;
```

80 Bytes

};

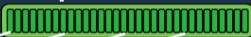
4 bytes usefull  
out of 32



Structure of Array (SoA)

Warp of 32 threads

Read of u



struct Coeff{

```
float *u, *v, *w;
float *x0, ... *x7, *y0, ... *y7, *z;
```

};



# Memory Layout

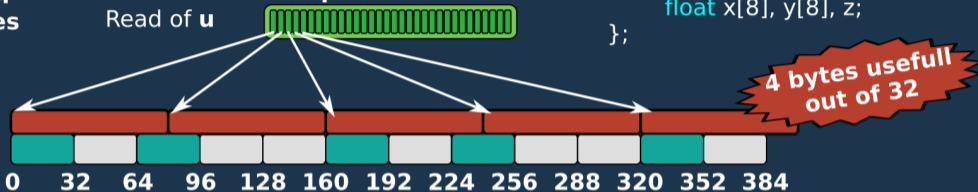
Read / Write :  
- 1 sector  
- 32 bytes

Array of Structure (AoS)

Warp of 32 threads

```
struct Coeff{
    float u, v, w;
    float x[8], y[8], z;
};
```

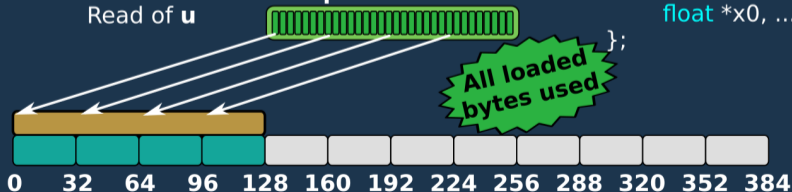
80 Bytes



Structure of Array (SoA)

Warp of 32 threads

```
struct Coeff{
    float *u, *v, *w;
    float *x0, ... *x7, *y0, ... *y7, *z;
};
```



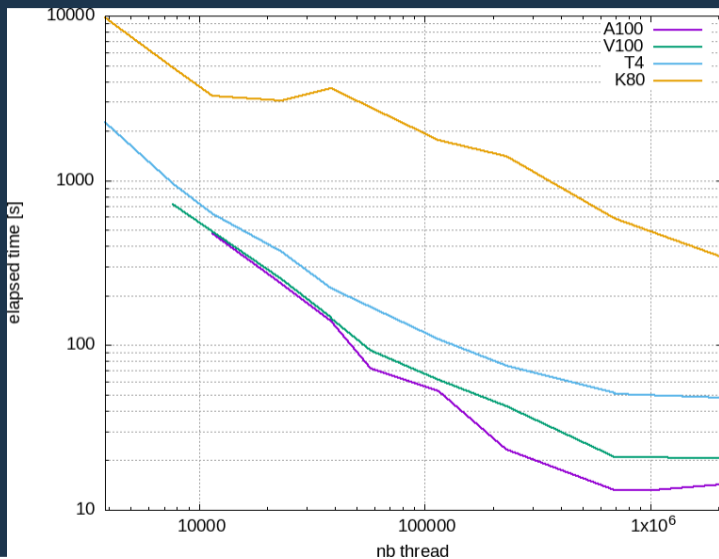


# How many threads ?

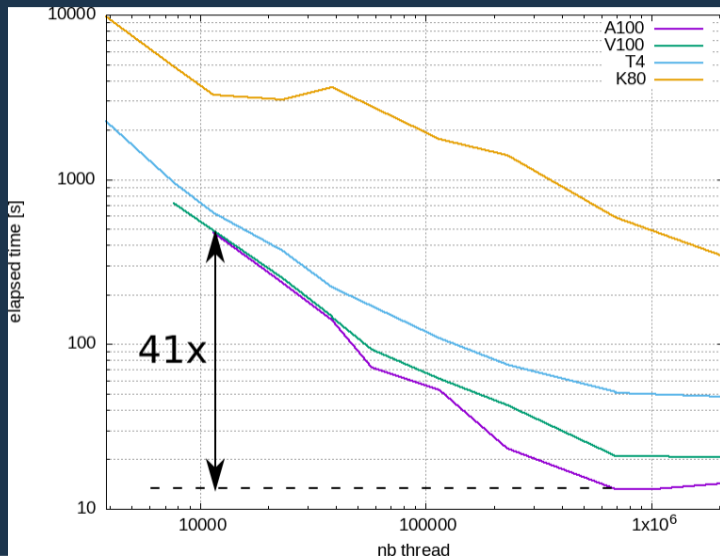
---



# How many threads ?



# How many threads ?





Element



# Express Parallelism

Element

Independent Elements => Independent Computing => Parallelism



# Express Parallelism

Element

Independent Elements => Independent Computing => Parallelism



Particle

# Express Parallelism

Element

Independent Elements => Independent Computing => Parallelism



Particle

$p_x$  

$p_y$  

$p_z$  



# Express Parallelism

Element

Independent Elements => Independent Computing => Parallelism



Particle

Mix Components :  
- Dot Product  
- Cross Product

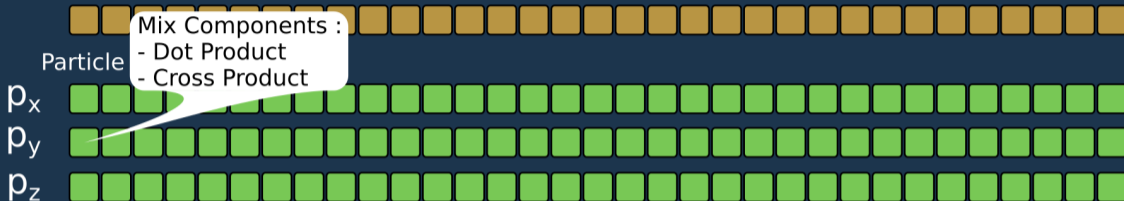
$p_x$  

$p_y$  

$p_z$  

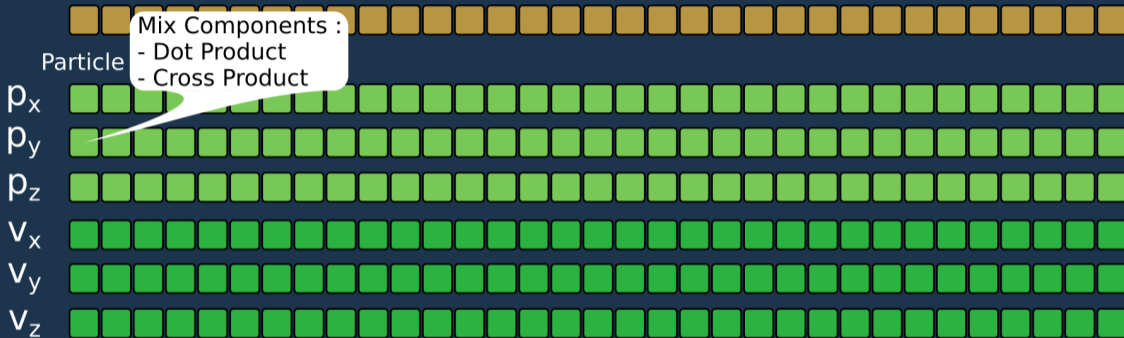
# Express Parallelism

Element Independent Elements => Independent Computing => Parallelism



# Express Parallelism

Element Independent Elements => Independent Computing => Parallelism



# Express Parallelism

Element                      Independent Elements => Independent Computing => Parallelism

Mix Components :  
- Dot Product  
- Cross Product

Particle

$p_x$

$p_y$

$p_z$

$v_x$

$v_y$

$v_z$

$a_x$

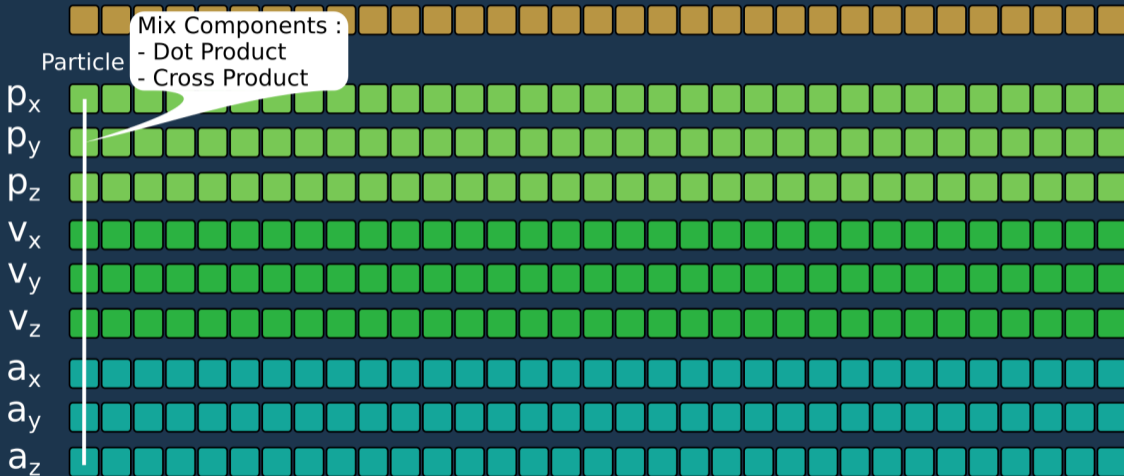
$a_y$

$a_z$



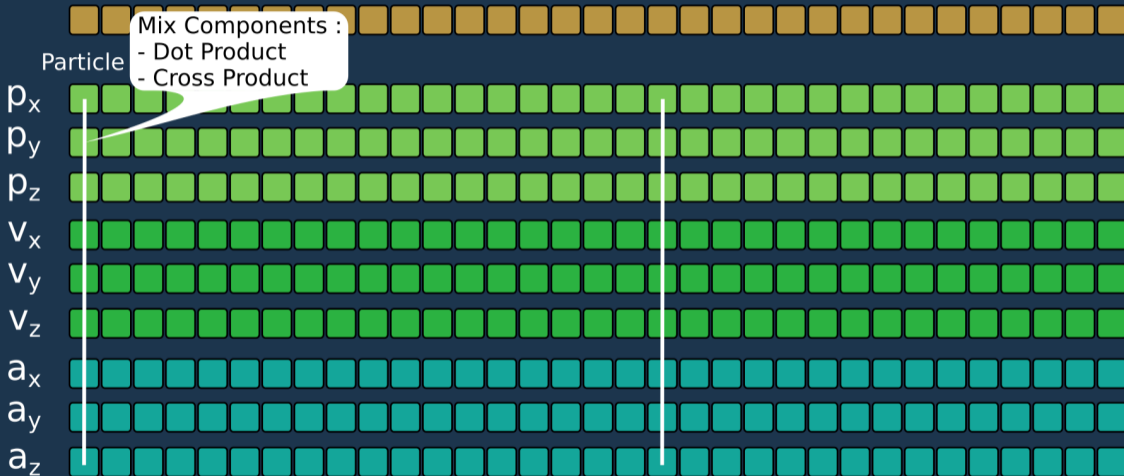
# Express Parallelism

Element Independent Elements => Independent Computing => Parallelism



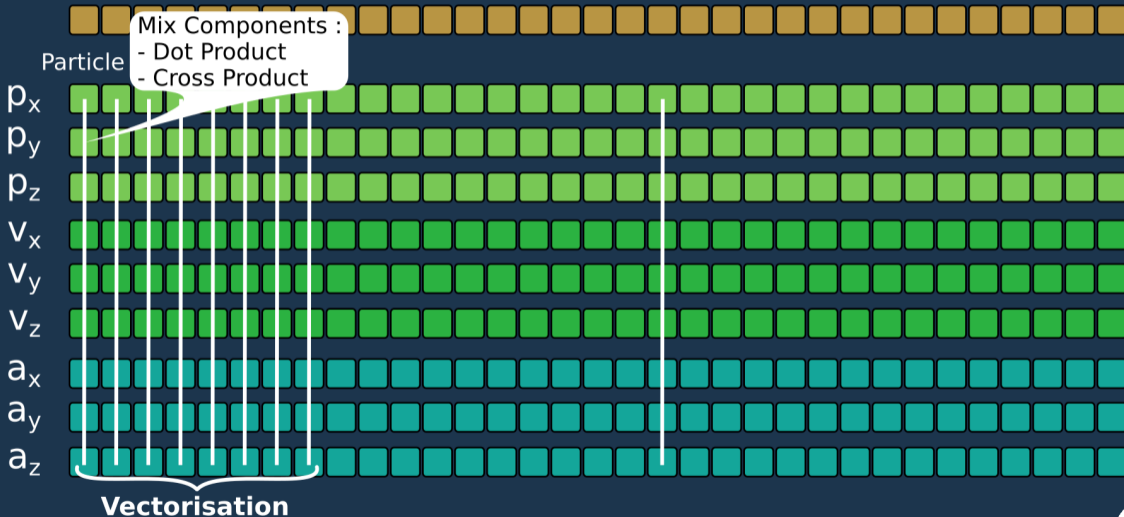
# Express Parallelism

Element Independent Elements => Independent Computing => Parallelism



# Express Parallelism

Element Independent Elements => Independent Computing => Parallelism

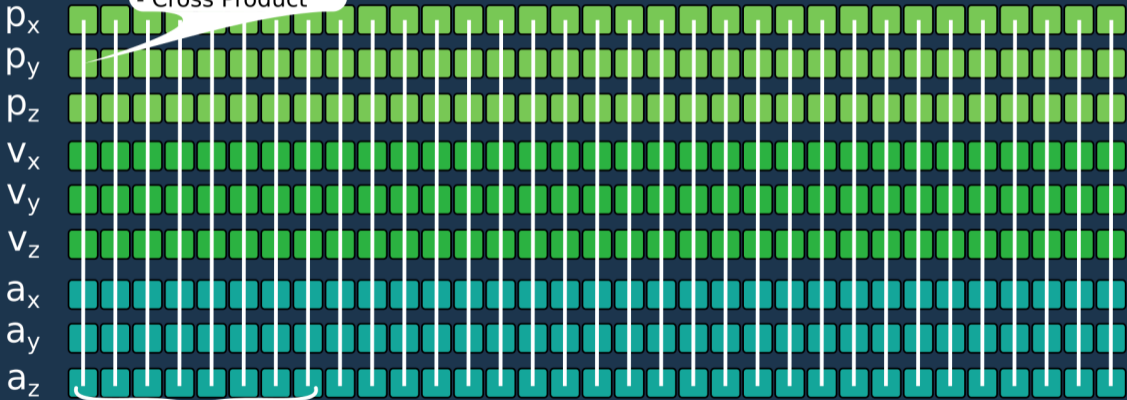


# Express Parallelism

Element Independent Elements => Independent Computing => Parallelism

Mix Components :  
- Dot Product  
- Cross Product

Particle



**Vectorisation**







Event



# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



# Express Parallelism

Event

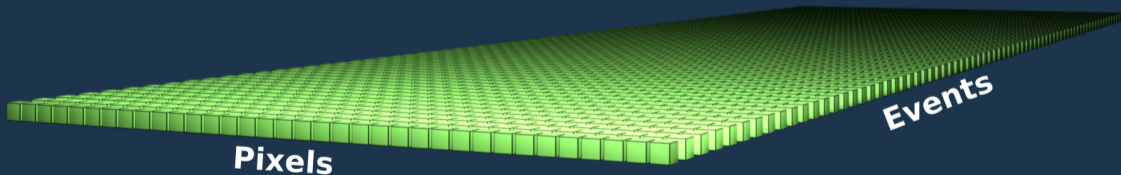
Independent Events => Independent Computing => Parallelism



# Express Parallelism

Event

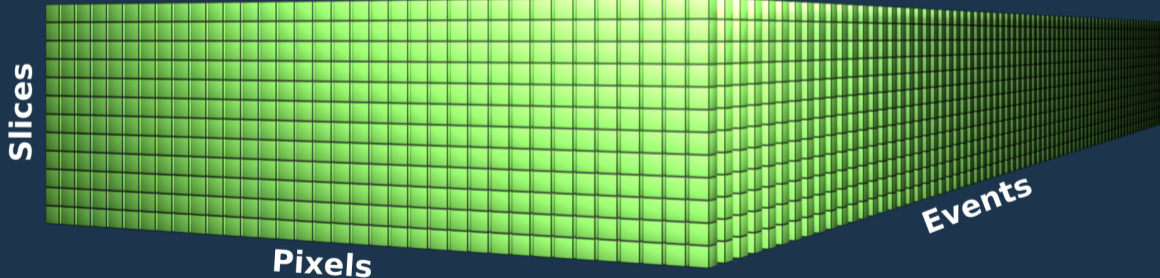
Independent Events => Independent Computing => Parallelism



# Express Parallelism

Event

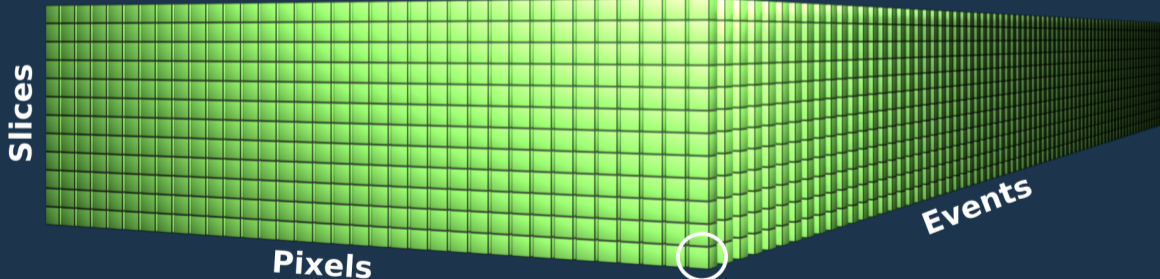
Independent Events => Independent Computing => Parallelism



# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism

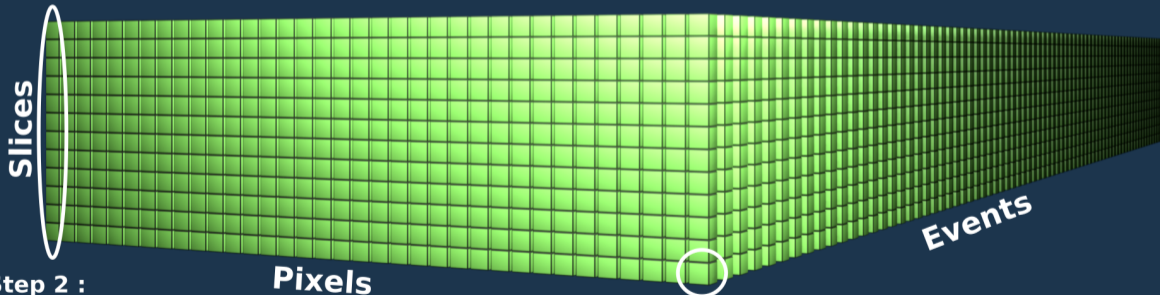


Step 1 : Independent Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Step 2 : **Independent** Computing in Integration

**Pixels**

Step 1 : **Independent** Computing in Calibration



# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Slices

Events

Step 2 :

**Pixels**

**Independent** Computing in Integration

**Step 1 : Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Contiguous Data

Slices

Events

Step 2 :

**Independent** Computing in Integration

**Pixels**

**Step 1 : Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Allows  
Vectorization

Contiguous Data

Slices

Events

Step 2 :

**Independent** Computing in Integration

**Pixels**

**Step 1 : Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Allows  
Vectorization

Even for  
Integration

Contiguous Data

Slices

Events

Step 2 :

**Independent** Computing in Integration

**Pixels**

**Step 1 : Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Allows  
Vectorization

Even for  
Integration

Contiguous Data

Slices

Events

Step 2 :

**Pixels**

**Independent** Computing in Integration

**Step 1 : Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Allows  
Vectorization

Even for  
Integration

Contiguous Data

Slices

Events

Step 2 :

**Independent** Computing in Integration

**Pixels**

**Step 1 : Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Allows Vectorization

Even for Integration

Contiguous Data

Slices

Parallelisms :  
- Cores  
- Nodes  
- Clusters

Events

Step 2 : **Independent** Computing in Integration

**Pixels**

Step 1 : **Independent** Computing in Calibration

# Express Parallelism

Event

Independent Events => Independent Computing => Parallelism



Reduction

Allows Vectorization

Even for Integration

Computing Drives Data Storage

Contiguous Data

Slices

Parallelisms :  
- Cores  
- Nodes  
- Clusters

Events

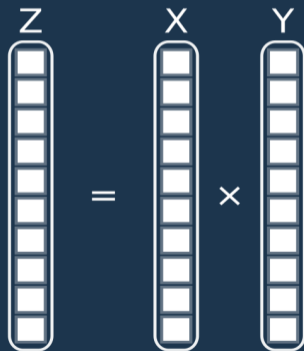
Step 2 : **Independent** Computing in Integration

**Pixels**

Step 1 : **Independent** Computing in Calibration

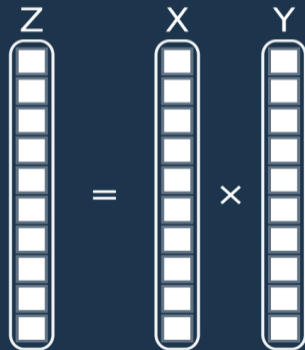


# Hadamard Product CUDA



# Hadamard Product CUDA

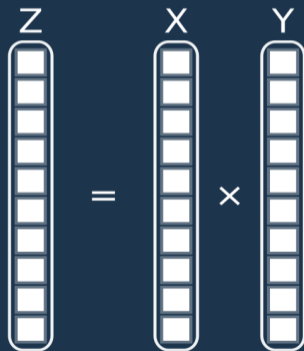
Element Wise  
Operation



# Hadamard Product CUDA

```
__global__ void hadamard_product_kernel(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement){
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    vecResult[i] = vecLeft[i] * vecRight[i];
}
```

Element Wise  
Operation

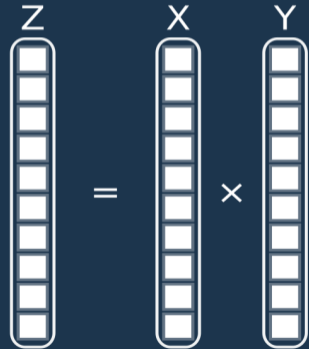


# Hadamard Product CUDA

```
global__ void hadamard_product_kernel(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement){  
    int i = blockIdx.x*blockDim.x + threadIdx.x;  
    vecResult[i] = vecLeft[i] * vecRight[i];  
}
```

Execution on  
**Device or Host**

Element Wise  
Operation



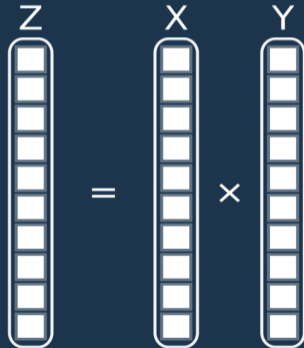
# Hadamard Product CUDA

```
global void hadamard_product_kernel(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement){
    int i = blockIdx.x*blockDim.x + threadIdx.x;
    vecResult[i] = vecLeft[i] * vecRight[i];
}
```

Execution on  
**Device or Host**

Thread position

Element Wise  
Operation



# Hadamard Product CUDA

```

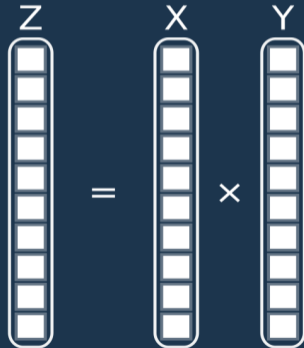
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                           int maxNbThreadPerBlockX)
{
    //Allocations
    float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
    int sizeByte = nbElement*sizeof(float);
    cudaMalloc((void**)&vecResultd, sizeByte);
    cudaMalloc((void**)&vecLeftd, sizeByte);
    cudaMalloc((void**)&vecRightd, sizeByte);

    //Block and grid size
    int dimGridX = 1;
    int dimBlockX = 1;
    if(nbElement > maxNbThreadPerBlockX){
        dimBlockX = maxNbThreadPerBlockX;
        dimGridX = nbElement/dimBlockX;
    }else{
        dimBlockX = nbElement;
    }

    //Transfer to Device
    cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
    PLIB_CUDA_CHECK_FILE
    cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
    PLIB_CUDA_CHECK_FILE

    dim3 dimGrid(dimGridX, 1, 1);
    dim3 dimBlock(dimBlockX, 1, 1);
    hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
    PLIB_CUDA_CHECK_FILE
    //Transfer to Host
    cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
    PLIB_CUDA_CHECK_FILE
    //Free memory
    cudaFree(vecRightd);
    cudaFree(vecLeftd);
    cudaFree(vecResultd);
}
    
```

Element Wise  
Operation



# Hadamard Product CUDA

```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

```
//Allocations
```

```
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void*)&vecResultd, sizeByte);
cudaMalloc((void*)&vecLeftd, sizeByte);
cudaMalloc((void*)&vecRightd, sizeByte);
```

Allocation  
on GPU

```
//Block and grid size
```

```
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
}
```

```
//Transfer to Device
```

```
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
```

```
dim3 dimGrid(dimGridX, 1, 1);
```

```
dim3 dimBlock(dimBlockX, 1, 1);
```

```
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
```

```
PLIB_CUDA_CHECK_FILE
```

```
//Transfer to Host
```

```
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
```

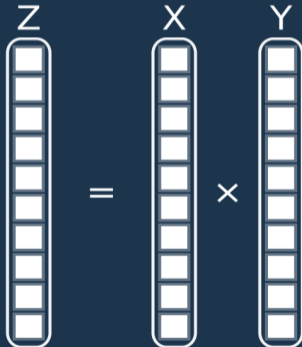
```
PLIB_CUDA_CHECK_FILE
```

```
//Free memory
```

```
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

```
}
```

Element Wise  
Operation



# Hadamard Product CUDA

```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

```
//Allocations
```

```
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void**)&vecResultd, sizeByte);
cudaMalloc((void**)&vecLeftd, sizeByte);
cudaMalloc((void**)&vecRightd, sizeByte);
```

Allocation  
on **GPU**

```
//Block and grid size
```

```
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
```

```
//Transfer to Device
```

```
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
```

```
dim3 dimGrid(dimGridX, 1, 1);
```

```
dim3 dimBlock(dimBlockX, 1, 1);
```

```
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
```

```
PLIB_CUDA_CHECK_FILE
```

```
//Transfer to Host
```

```
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
```

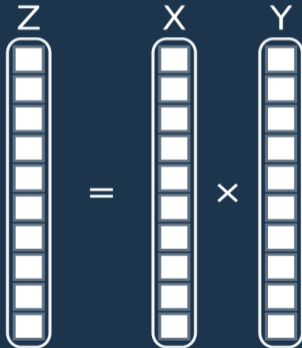
```
PLIB_CUDA_CHECK_FILE
```

```
//Free memory
```

```
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

Free on **GPU**

Element Wise  
Operation





```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

```
//Allocations
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void**)&vecResultd, sizeByte);
cudaMalloc((void**)&vecLeftd, sizeByte);
cudaMalloc((void**)&vecRightd, sizeByte);
```

Allocation  
on **GPU**

```
//Block and grid size
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
```

```
//Transfer to Device
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
```

```
dim3 dimGrid(dimGridX, 1, 1);
dim3 dimBlock(dimBlockX, 1, 1);
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
PLIB_CUDA_CHECK_FILE
```

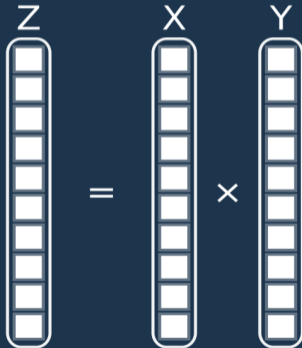
```
//Transfer to Host
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
PLIB_CUDA_CHECK_FILE
```

```
//Free memory
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

Free on **GPU**

Out of this function if many calls

Element Wise  
Operation



# Hadamard Product CUDA

```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

Out of this function if many calls

```
//Allocations
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void**)&vecResultd, sizeByte);
cudaMalloc((void**)&vecLeftd, sizeByte);
cudaMalloc((void**)&vecRightd, sizeByte);
```

Allocation  
on GPU

```
//Block and grid size
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
```

```
//Transfer to Device
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB_CUDA_CHECK_FILE
```

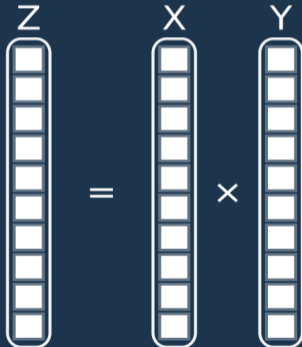
Transfer data  
to GPU

```
dim3 dimGrid(dimGridX, 1, 1);
dim3 dimBlock(dimBlockX, 1, 1);
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
PLIB_CUDA_CHECK_FILE
//Transfer to Host
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
PLIB_CUDA_CHECK_FILE
```

```
//Free memory
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

Free on GPU

Element Wise  
Operation



# Hadamard Product CUDA

```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

```
//Allocations
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void**)&vecResultd, sizeByte);
cudaMalloc((void**)&vecLeftd, sizeByte);
cudaMalloc((void**)&vecRightd, sizeByte);
```

Allocation  
on **GPU**

```
//Block and grid size
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
```

```
//Transfer to Device
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB CUDA CHECK FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB CUDA CHECK FILE
```

Transfer data  
to **GPU**

```
dim3 dimGrid(dimGridX, 1, 1);
dim3 dimBlock(dimBlockX, 1, 1);
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
PLIB CUDA CHECK FILE
```

```
//Transfer to Host
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
PLIB CUDA CHECK FILE
```

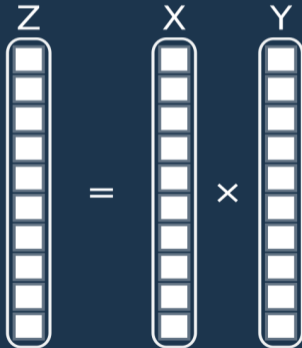
Transfer result  
back to **host**

```
//Free memory
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

Free on **GPU**

Out of this function if many calls

Element Wise  
Operation



# Hadamard Product CUDA

```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

Out of this function if many calls

```
//Allocations
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void**)&vecResultd, sizeByte);
cudaMalloc((void**)&vecLeftd, sizeByte);
cudaMalloc((void**)&vecRightd, sizeByte);
```

Allocation on GPU

```
//Block and grid size
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
```

Determine Grid and Block size

```
//Transfer to Device
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB CUDA CHECK FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB CUDA CHECK FILE
```

Transfer data to GPU

```
dim3 dimGrid(dimGridX, 1, 1);
dim3 dimBlock(dimBlockX, 1, 1);
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
PLIB CUDA CHECK FILE
```

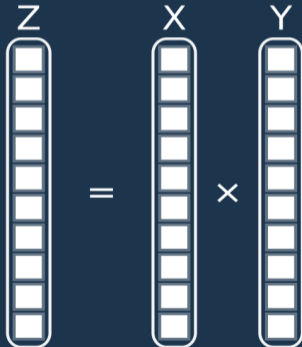
```
//Transfer to Host
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
PLIB CUDA CHECK FILE
```

Transfer result back to host

```
//Free memory
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

Free on GPU

Element Wise Operation



# Hadamard Product CUDA

```
void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
                          int maxNbThreadPerBlockX)
```

Out of this function if many calls

```
//Allocations
float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
int sizeByte = nbElement*sizeof(float);
cudaMalloc((void**)&vecResultd, sizeByte);
cudaMalloc((void**)&vecLeftd, sizeByte);
cudaMalloc((void**)&vecRightd, sizeByte);
```

Allocation  
on GPU

```
//Block and grid size
int dimGridX = 1;
int dimBlockX = 1;
if(nbElement > maxNbThreadPerBlockX){
    dimBlockX = maxNbThreadPerBlockX;
    dimGridX = nbElement/dimBlockX;
}else{
    dimBlockX = nbElement;
}
```

Determine  
Grid and Block size

```
//Transfer to Device
cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
PLIB CUDA CHECK FILE
cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
PLIB CUDA CHECK FILE
```

Transfer data  
to GPU

```
dim3 dimGrid(dimGridX, 1, 1);
dim3 dimBlock(dimBlockX, 1, 1);
hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
PLIB CUDA CHECK FILE
```

Call  
Kernel

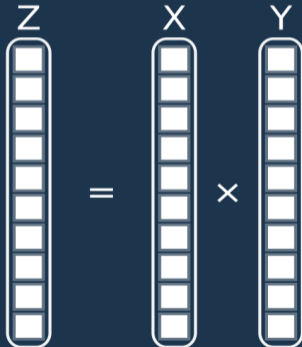
```
//Transfer to Host
cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
PLIB CUDA CHECK FILE
```

Transfer result  
back to host

```
//Free memory
cudaFree(vecRightd);
cudaFree(vecLeftd);
cudaFree(vecResultd);
```

Free on GPU

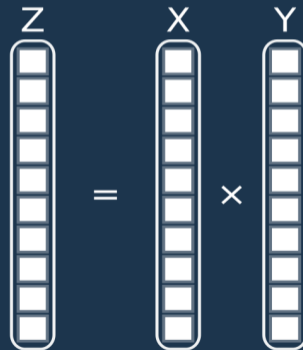
Element Wise  
Operation



```

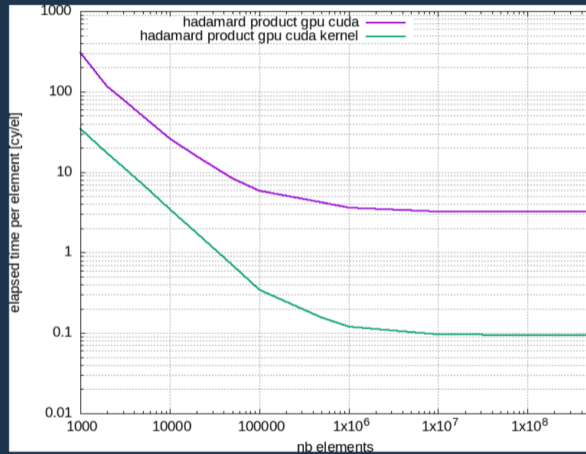
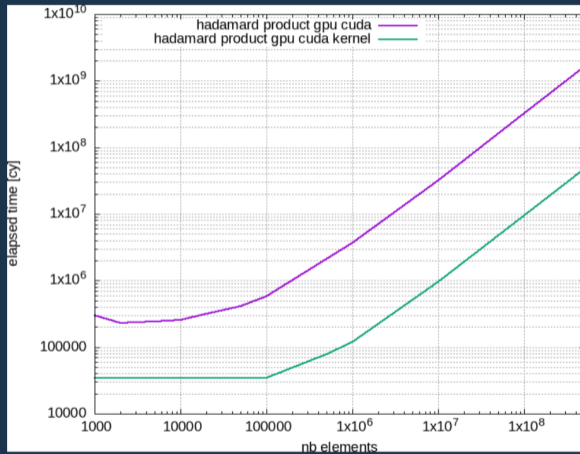
1  #include <cuda.h>
2  #include <cuda_runtime.h>
3
4  __global__ void hadamard_product_kernel(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement){
5      int i = blockIdx.x*blockDim.x + threadIdx.x;
6      vecResult[i] = vecLeft[i] * vecRight[i];
7  }
8
9  void hadamard_product_cuda(float * vecResult, const float * vecLeft, const float * vecRight, int nbElement,
10                             int maxNbThreadPerBlockX)
11  {
12      //Allocations
13      float *vecResultd = NULL, *vecLeftd = NULL, *vecRightd = NULL;
14      int sizeByte = nbElement*sizeof(float);
15      cudaMalloc((void**)&vecResultd, sizeByte);
16      cudaMalloc((void**)&vecLeftd, sizeByte);
17      cudaMalloc((void**)&vecRightd, sizeByte);
18
19      //Block and grid size
20      int dimGridX = 1;
21      int dimBlockX = 1;
22      if(nbElement > maxNbThreadPerBlockX){
23          dimBlockX = maxNbThreadPerBlockX;
24          dimGridX = nbElement/dimBlockX;
25      }else{
26          dimBlockX = nbElement;
27      }
28
29      //Transfer to Device
30      cudaMemcpy(vecLeftd, vecLeft, sizeByte, cudaMemcpyHostToDevice);
31      PLIB_CUDA_CHECK_FILE
32      cudaMemcpy(vecRightd, vecRight, sizeByte, cudaMemcpyHostToDevice);
33      PLIB_CUDA_CHECK_FILE
34
35      dim3 dimGrid(dimGridX, 1, 1);
36      dim3 dimBlock(dimBlockX, 1, 1);
37      hadamard_product_kernel<<<dimGrid, dimBlock>>>(vecResultd, vecLeftd, vecRightd, nbElement);
38      PLIB_CUDA_CHECK_FILE
39      //Transfer to Host
40      cudaMemcpy(vecResult, vecResultd, sizeByte, cudaMemcpyDeviceToHost);
41      PLIB_CUDA_CHECK_FILE
42      //Free memory
43      cudaFree(vecRightd);
44      cudaFree(vecLeftd);
45      cudaFree(vecResultd);
46  }
    
```

Element Wise  
Operation



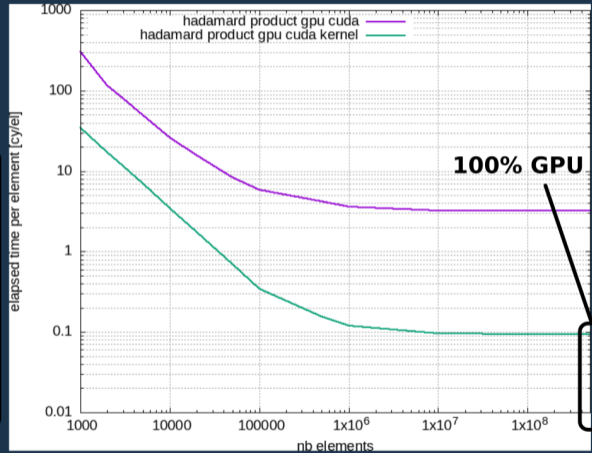
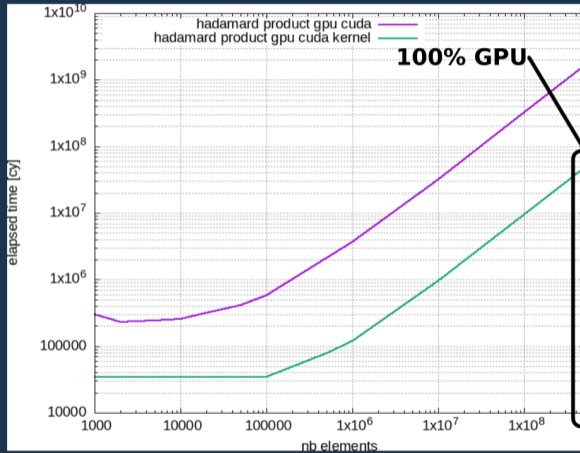
# Perf : Hadamard Product CUDA

On A3000 GPU



# Perf : Hadamard Product CUDA

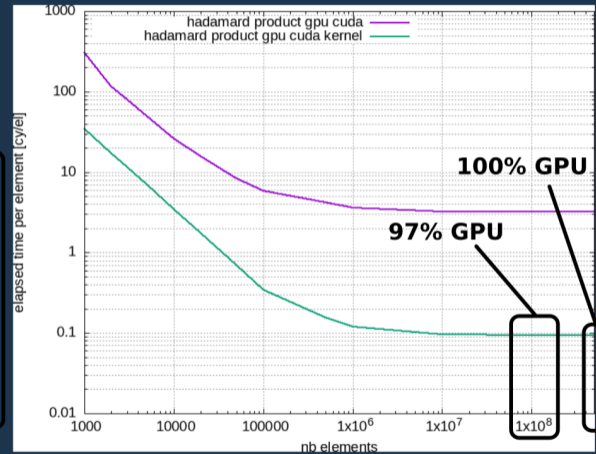
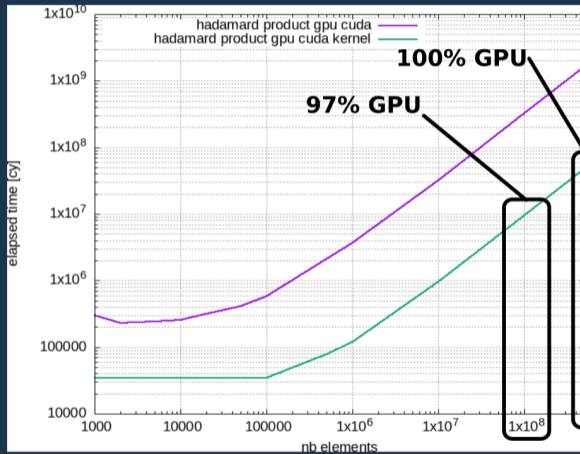
On A3000 GPU





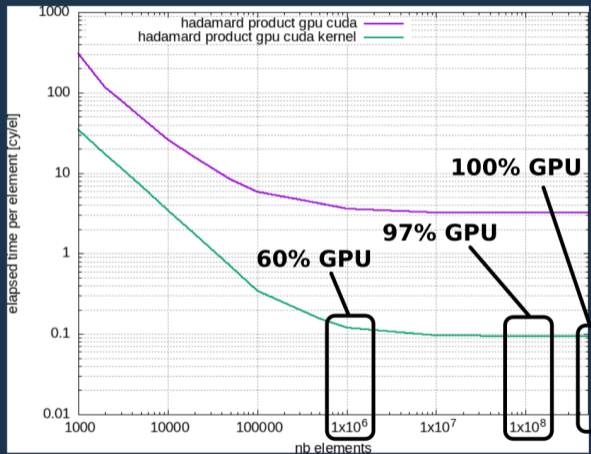
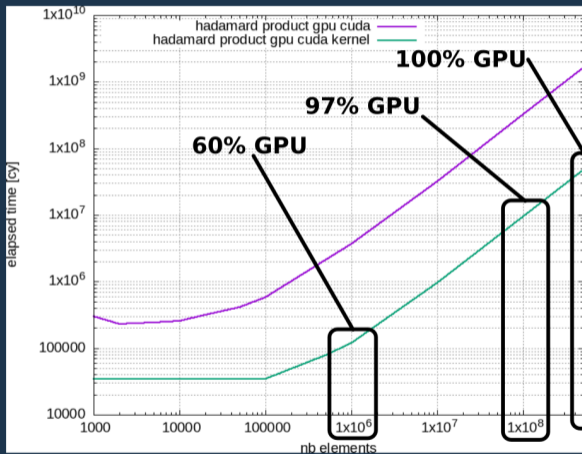
# Perf : Hadamard Product CUDA

On A3000 GPU



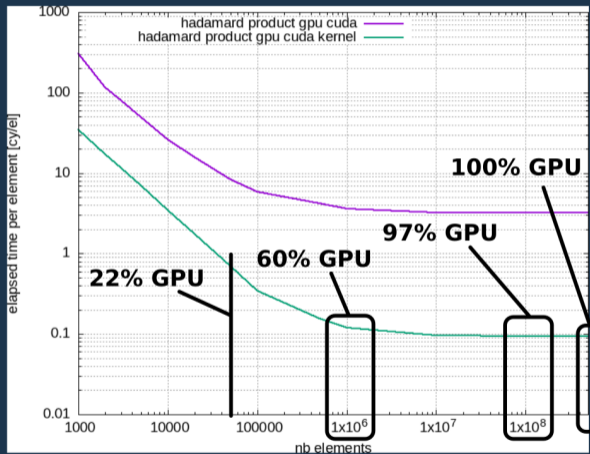
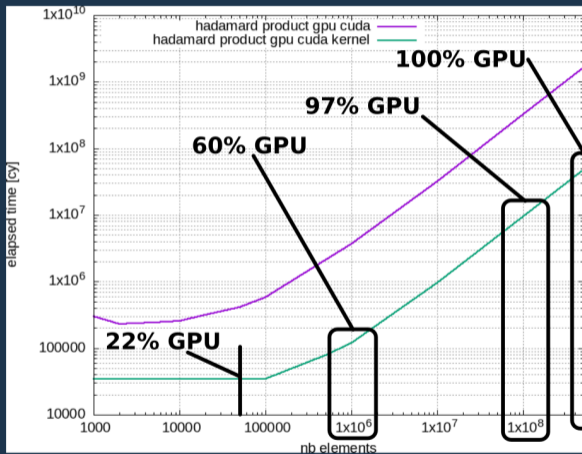
# Perf : Hadamard Product CUDA

On A3000 GPU

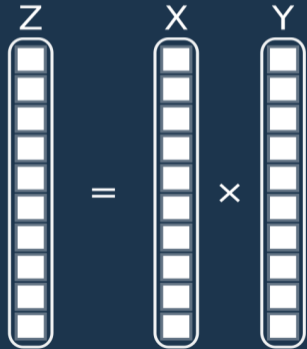


# Perf : Hadamard Product CUDA

On A3000 GPU

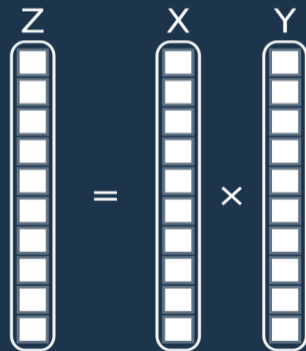


# Hadamard Product Thrust



# Hadamard Product Thrust

Element Wise  
Operation



```

int main(int argc, char** argv){
    size_t nbElement(500000000);

    thrust::host_vector<float> h_vecX(nbElement), h_vecY(nbElement);
    for(size_t i(0lu); i < nbElement; ++i){
        h_vecX[i] = i*19lu%11;
        h_vecY[i] = i*27lu%19;
    }
    // Transfer data to the device.
    thrust::device_vector<float> d_vecX = h_vecX, d_vecY = h_vecY, d_vecRes(nbElement);

    thrust::transform(d_vecX.begin(), d_vecX.end(), d_vecY.begin(), d_vecRes.begin(),
        [=] __host__ __device__ (const float& x, const float& y){
            return x * y;
        });

    // Transfer data to the host
    thrust::host_vector<float> h_vecRes = d_vecRes;

    std::cout << "x = " << h_vecX.front() << ", y = " << h_vecY.front() << ", res = " << h_vecRes.front() << std::endl;
    return 0;
}

```

```
int main(int argc, char** argv){
    size_t nbElement(500000000);

    thrust::host_vector<float> h_vecX(nbElement), h_vecY(nbElement);
    for(size_t i(0lu); i < nbElement; ++i){
        h_vecX[i] = i*19lu%11;
        h_vecY[i] = i*27lu%19;
    }

    // Transfer data to the device.
    thrust::device_vector<float> d_vecX = h_vecX, d_vecY = h_vecY, d_vecRes(nbElement);

    thrust::transform(d_vecX.begin(), d_vecX.end(), d_vecY.begin(), d_vecRes.begin(),
        [=] __host__ __device__ (const float& x, const float& y){
            return x * y;
        });

    // Transfer data to the host
    thrust::host_vector<float> h_vecRes = d_vecRes;

    std::cout << "x = " << h_vecX.front() << ", y = " << h_vecY.front() << ", res = " << h_vecRes.front() << std::endl;
    return 0;
}
```

Initialisation

```

int main(int argc, char** argv){
    size_t nbElement(500000000);

    thrust::host_vector<float> h_vecX(nbElement), h_vecY(nbElement);
    for(size_t i(0lu); i < nbElement; ++i){
        h_vecX[i] = i*19lu%11;
        h_vecY[i] = i*27lu%19;
    }
    // Transfer data to the device.
    thrust::device_vector<float> d_vecX = h_vecX, d_vecY = h_vecY, d_vecRes(nbElement);
    thrust::transform(d_vecX.begin(), d_vecX.end(), d_vecY.begin(), d_vecRes.begin(),
        [=] __host__ __device__ (const float& x, const float& y){
            return x * y;
        });
    // Transfer data to the host
    thrust::host_vector<float> h_vecRes = d_vecRes;

    std::cout << "x = " << h_vecX.front() << ", y = " << h_vecY.front() << ", res = " << h_vecRes.front() << std::endl;
    return 0;
}
    
```

Initialisation

Transfer  
Host to Device



```

int main(int argc, char** argv){
    size_t nbElement(500000000);

    thrust::host_vector<float> h_vecX(nbElement), h_vecY(nbElement);
    for(size_t i(0lu); i < nbElement; ++i){
        h_vecX[i] = i*19lu%11;
        h_vecY[i] = i*27lu%19;
    }

    // Transfer data to the device.
    thrust::device_vector<float> d_vecX = h_vecX, d_vecY = h_vecY, d_vecRes(nbElement);

    thrust::transform(d_vecX.begin(), d_vecX.end(), d_vecY.begin(), d_vecRes.begin(),
        [=] __host__ __device__ (const float& x, const float& y){
            return x * y;
        });

    // Transfer data to the host
    thrust::host_vector<float> h_vecRes = d_vecRes;

    std::cout << "x = " << h_vecX.front() << ", y = " << h_vecY.front() << ", res = " << h_vecRes.front() << std::endl;
    return 0;
}
    
```

**Initialisation**

**Transfer  
Host to Device**

**Compute  
Hadamard Product**

```

int main(int argc, char** argv){
    size_t nbElement(500000000);

    thrust::host_vector<float> h_vecX(nbElement), h_vecY(nbElement);
    for(size_t i(0lu); i < nbElement; ++i){
        h_vecX[i] = i*19lu%11;
        h_vecY[i] = i*27lu%19;
    }
    // Transfer data to the device.
    thrust::device_vector<float> d_vecX = h_vecX, d_vecY = h_vecY, d_vecRes(nbElement);
    thrust::transform(d_vecX.begin(), d_vecX.end(), d_vecY.begin(), d_vecRes.begin(),
        [=] __host__ __device__ (const float& x, const float& y){
            return x * y;
        });
    // Transfer data to the host
    thrust::host_vector<float> h_vecRes = d_vecRes;
    std::cout << "x = " << h_vecX.front() << ", y = " << h_vecY.front() << ", res = " << h_vecRes.front() << std::endl;
    return 0;
}

```

**Initialisation**

**Transfer  
Host to Device**

**Compute  
Hadamard Product**

**Transfer Device to Host**

```

#include <iostream>
#include <vector>

#include <thrust/transform.h>
#include <thrust/functional.h>
#include <thrust/host_vector.h>
#include <thrust/device_vector.h>

int main(int argc, char** argv){
    size_t nbElement(500000000);

    thrust::host_vector<float> h_vecX(nbElement), h_vecY(nbElement);
    for(size_t i(0lu); i < nbElement; ++i){
        h_vecX[i] = i*19lu%11;
        h_vecY[i] = i*27lu%19;
    }
    // Transfer data to the device.
    thrust::device_vector<float> d_vecX = h_vecX, d_vecY = h_vecY, d_vecRes(nbElement);

    thrust::transform(d_vecX.begin(), d_vecX.end(), d_vecY.begin(), d_vecRes.begin(),
        [=] __host__ __device__ (const float& x, const float& y){
            return x * y;
        });

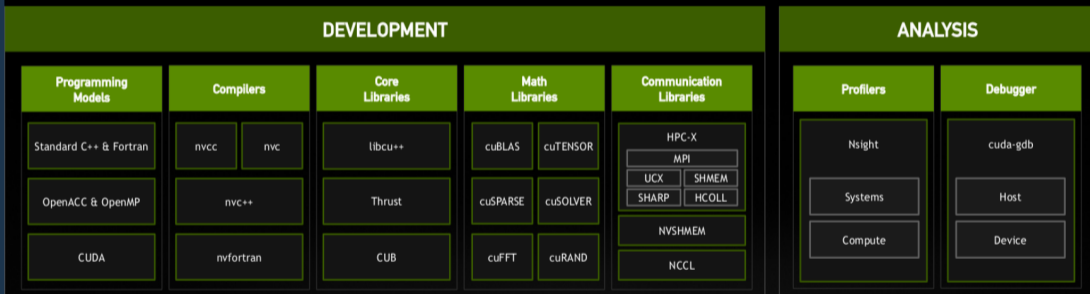
    // Transfer data to the host
    thrust::host_vector<float> h_vecRes = d_vecRes;

    std::cout << "x = " << h_vecX.front() << ", y = " << h_vecY.front() << ", res = " << h_vecRes.front() << std::endl;
    return 0;
}

```

## NVIDIA HPC SDK

Available at [developer.nvidia.com/hpc-sdk](https://developer.nvidia.com/hpc-sdk), on NGC, via Spack, and in the Cloud

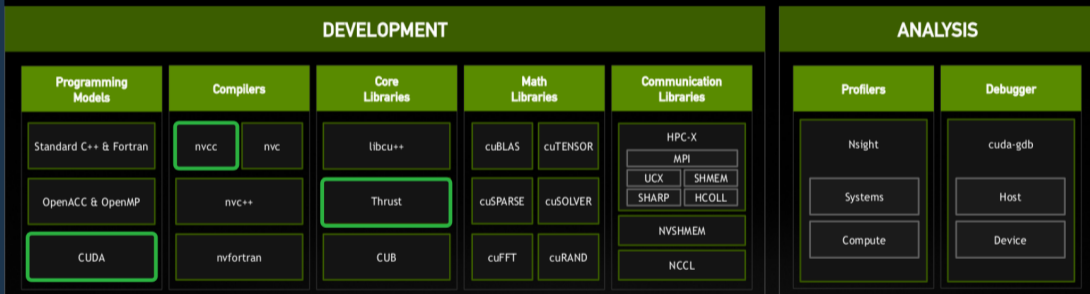


Develop for the NVIDIA Platform: GPU, CPU and Interconnect  
 Libraries | Accelerated C++ and Fortran | Directives | CUDA  
 7-8 Releases Per Year | Freely Available

# NVIDIA HPC SDK Introduction

## NVIDIA HPC SDK

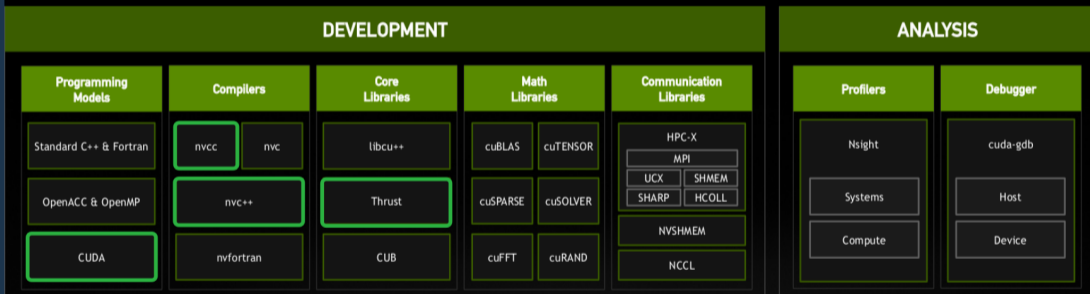
Available at [developer.nvidia.com/hpc-sdk](https://developer.nvidia.com/hpc-sdk), on NGC, via Spack, and in the Cloud



Develop for the NVIDIA Platform: GPU, CPU and Interconnect  
 Libraries | Accelerated C++ and Fortran | Directives | CUDA  
 7-8 Releases Per Year | Freely Available

## NVIDIA HPC SDK

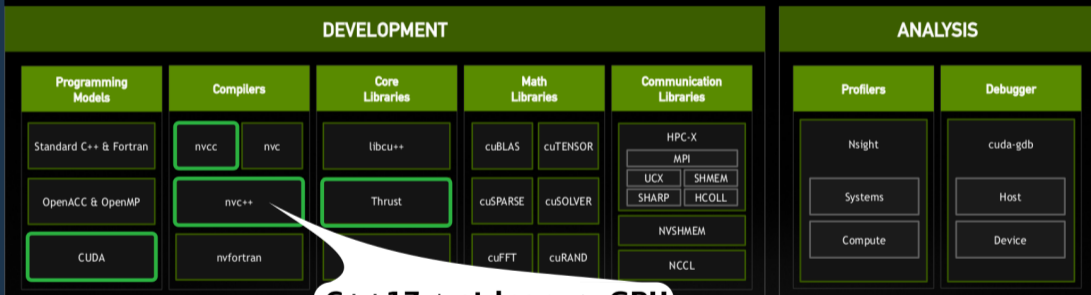
Available at [developer.nvidia.com/hpc-sdk](https://developer.nvidia.com/hpc-sdk), on NGC, via Spack, and in the Cloud



Develop for the NVIDIA Platform: GPU, CPU and Interconnect  
 Libraries | Accelerated C++ and Fortran | Directives | CUDA  
 7-8 Releases Per Year | Freely Available

## NVIDIA HPC SDK

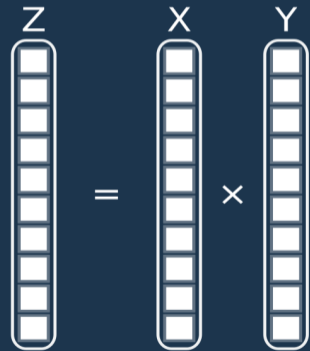
Available at [developer.nvidia.com/hpc-sdk](https://developer.nvidia.com/hpc-sdk), on NGC, via Spack, and in the Cloud



**C++17 + stdpar => GPU**

Develop for the NVIDIA Platform: GPU, CPU and Interconnect  
Libraries | Accelerated C++ and Fortran | Directives | CUDA  
7-8 Releases Per Year | Freely Available

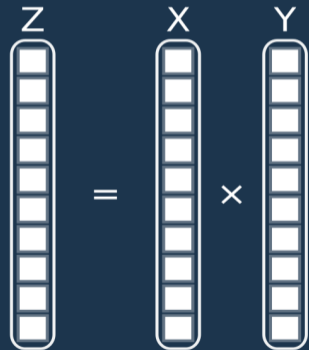
# Example : Hadamard Product





# Example : Hadamard Product

Element Wise  
Operation

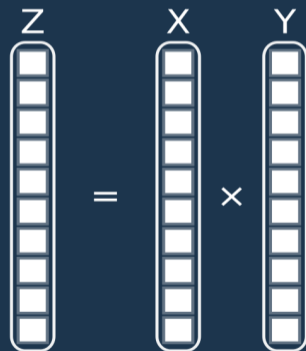


# Example : Hadamard Product

C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    > tabResult[i] = tabX[i]*tabY[i];
}
```

Element Wise  
Operation



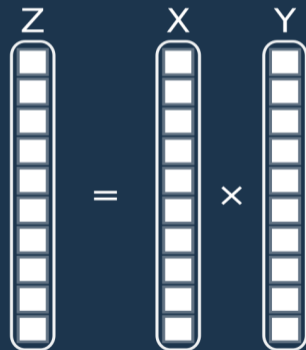
# Example : Hadamard Product

C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    > tabResult[i] = tabX[i]*tabY[i];
}
```

Explicit order  
**not** necessary

Element Wise  
Operation



# Example : Hadamard Product

C++

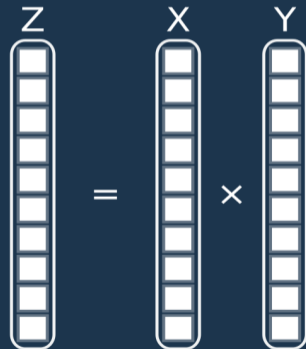
```
for(long unsigned int i(0lu); i < nbElement; ++i){
    > tabResult[i] = tabX[i]*tabY[i];
}
```

Explicit order  
**not** necessary

Element Wise  
Operation

C++ Algorithm : `std::transform`

```
std::transform(std::begin(tabX), std::end(tabX),
    > std::begin(tabY), std::begin(tabRes),
    > [](float xi, float yi){ return xi * yi; });
```



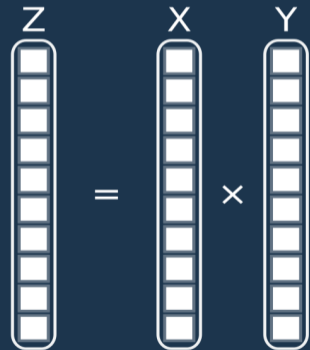
# Example : Hadamard Product

C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >     tabResult[i] = tabX[i]*tabY[i];
}
```

Explicit order  
**not** necessary

Element Wise  
Operation



C++ Algorithm : `std::transform`

```
std::transform(std::begin(tabX), std::end(tabX),
    >     std::begin(tabY), std::begin(tabRes),
    >     [](float xi, float yi){ return xi * yi; });
```

C++ 17 / C++ 20

```
std::transform(std::execution::par_unseq,
    >     std::begin(tabX), std::end(tabX),
    >     std::begin(tabY), std::begin(tabRes),
    >     [](float xi, float yi){ return xi * yi; });
```

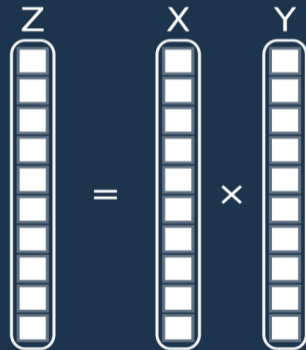
# Example : Hadamard Product

C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >     tabResult[i] = tabX[i]*tabY[i];
}
```

Explicit order  
**not** necessary

Element Wise  
Operation



C++ Algorithm : `std::transform`

```
std::transform(std::begin(tabX), std::end(tabX),
    >     std::begin(tabY), std::begin(tabRes),
    >     [](float xi, float yi){ return xi * yi; });
```

C++ 17 / C++ 20

Execution Policy

```
std::transform(std::execution::par_unseq
    >     std::begin(tabX), std::end(tabX),
    >     std::begin(tabY), std::begin(tabRes),
    >     [](float xi, float yi){ return xi * yi; });
```

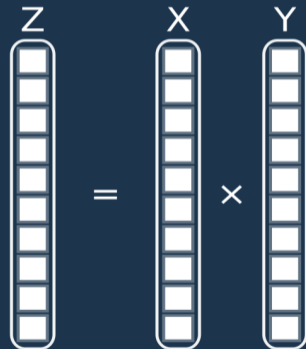
# Example : Hadamard Product

C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    > tabResult[i] = tabX[i]*tabY[i];
}
```

Explicit order  
**not** necessary

Element Wise  
Operation



C++ Algorithm : `std::transform`

```
std::transform(std::begin(tabX), std::end(tabX),
    > std::begin(tabY), std::begin(tabRes),
    > [](float xi, float yi){ return xi * yi; });
```

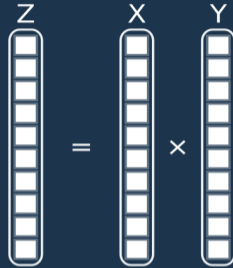
C++ 17 / C++ 20

Execution Policy

```
std::transform(std::execution::par_unseq,
    > std::begin(tabX), std::end(tabX),
    > std::begin(tabY), std::begin(tabRes),
    > [](float xi, float yi){ return xi * yi; });
```

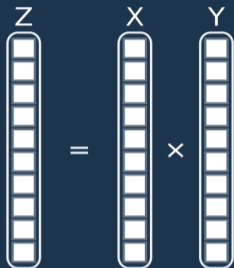
- seq
- unseq
- par
- par\_unseq

# Example : Hadamard Product





# Example : Hadamard Product

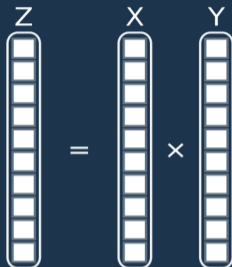


With CMake :

- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

# Example : Hadamard Product

`nvc++ -stdpar=cpu`



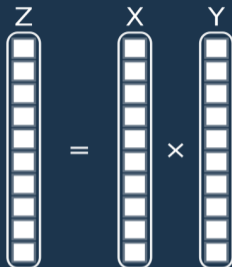
With CMake :

- `export CC=/path/to/gcc`
- `export CXX=/path/to/nvc++`



# Example : Hadamard Product

`nvc++ -stdpar=cpu`



With CMake :

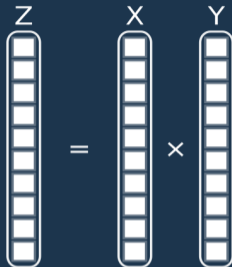
- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

# Example : Hadamard Product

`nvc++ -stdpar=cpu`



**Do not forget to link with TBB**



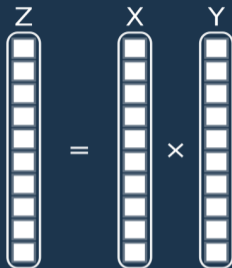
**With CMake :**

- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

# Example : Hadamard Product

`nvc++ -stdpar=cpu`

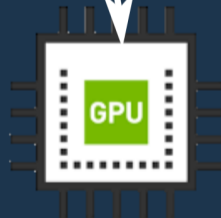
`nvc++ -stdpar=gpu`



**Do not forget to link with TBB**

With CMake :

- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

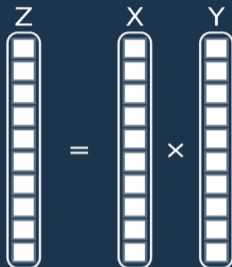


# Example : Hadamard Product

`nvc++ -stdpar=cpu`



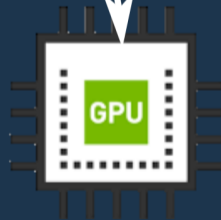
**Do not forget to link with TBB**



With CMake :

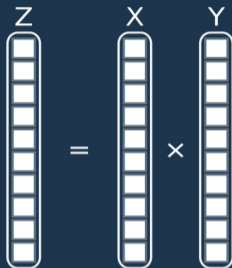
- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

`nvc++ -stdpar=gpu`



# Example : Hadamard Product

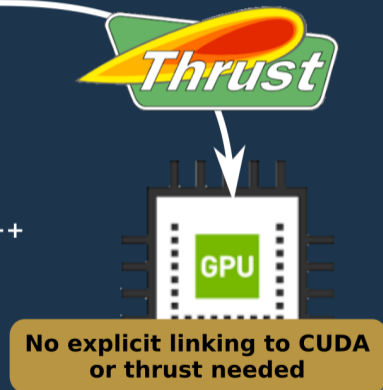
`nvc++ -stdpar=cpu`



With CMake :

- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

`nvc++ -stdpar=gpu`

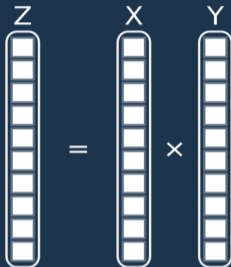


# Example : Hadamard Product

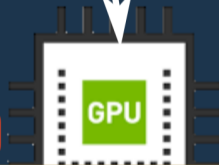
`nvc++ -stdpar=cpu`



Do not forget to link with TBB



`nvc++ -stdpar=gpu`



No explicit linking to CUDA or thrust needed

With CMake :

- export CC=/path/to/gcc
- export CXX=/path/to/nvc++

Avoid static Allocation





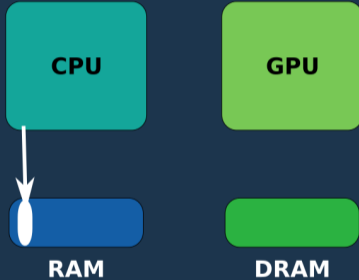
CUDA



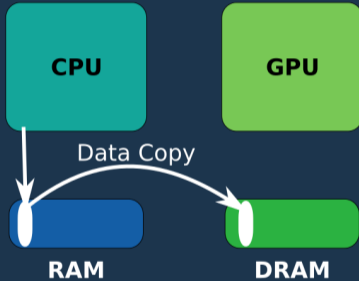
**RAM**

**DRAM**

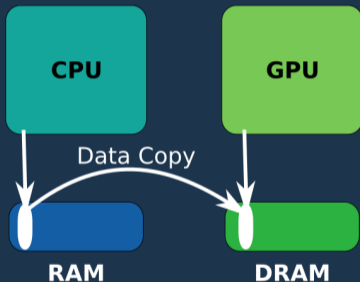
CUDA



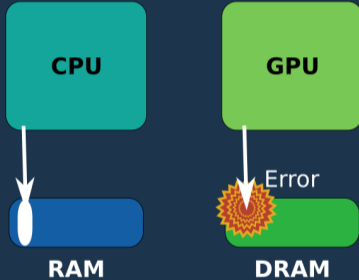
CUDA



CUDA

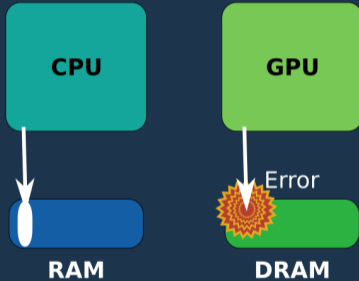


CUDA

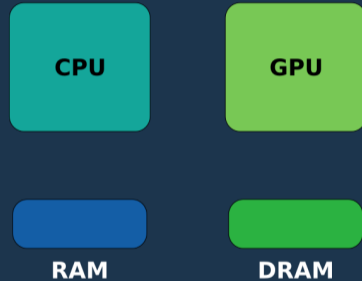


# Memory Management

CUDA

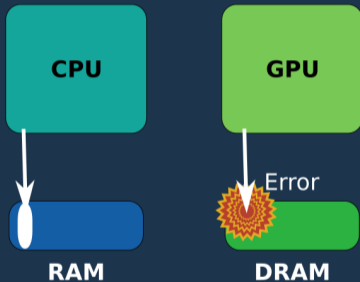


nvc++/nvfortran  
Unified Memory

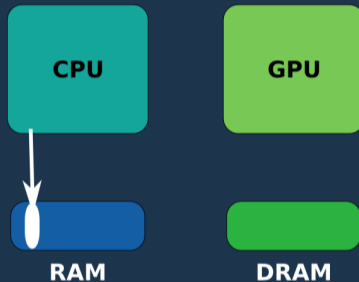


# Memory Management

CUDA



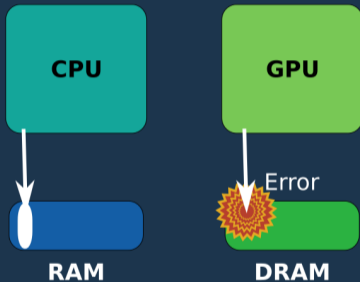
nvc++/nvfortran  
Unified Memory



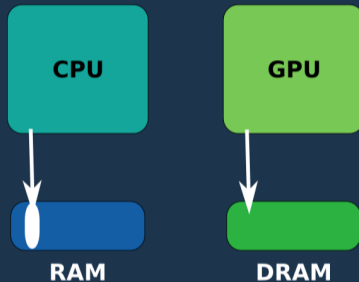


# Memory Management

CUDA

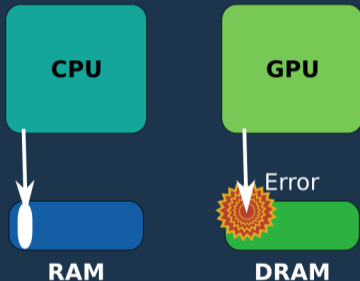


nvc++/nvfortran  
Unified Memory

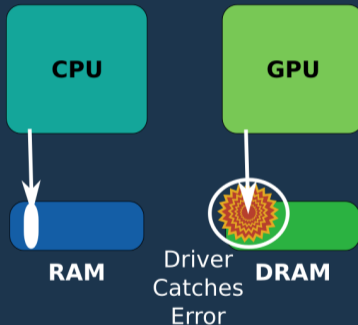


# Memory Management

CUDA

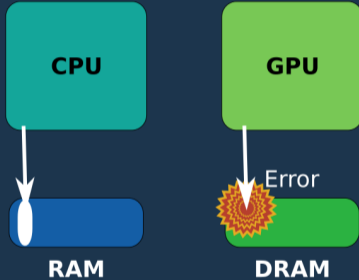


nvc++/nvfortran  
Unified Memory

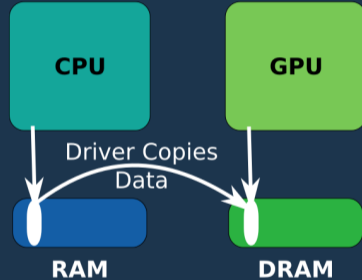


# Memory Management

CUDA



nvc++/nvfortran  
Unified Memory



# NSight : Hadamard Product



# NSight : Hadamard Product



# NSight : Hadamard Product



# NSight : Hadamard Product

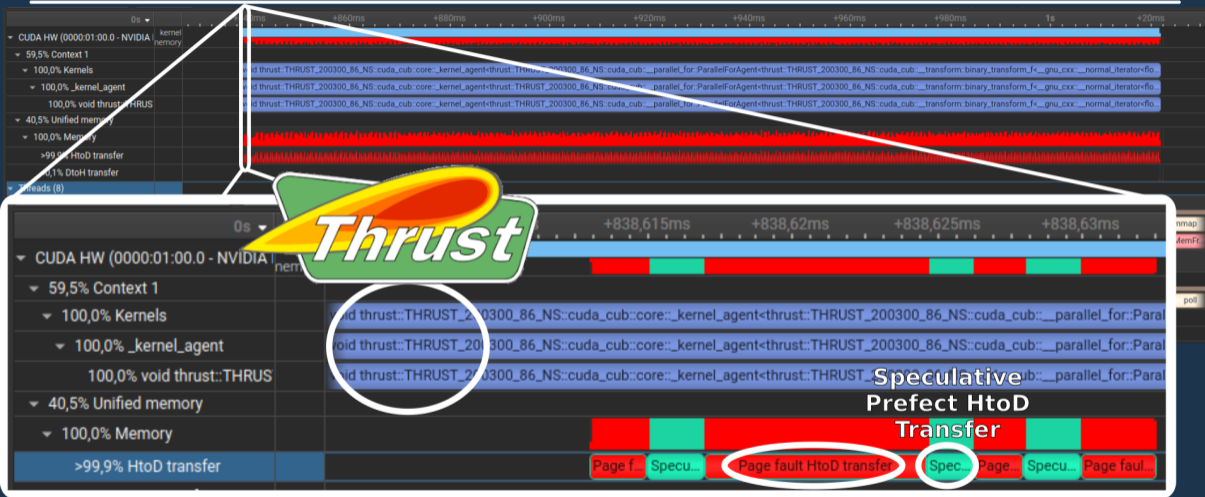


Speculative  
Prefetch HtoD  
Transfer

Page fault HtoD transfer

Specu...

# NSight : Hadamard Product



Speculative  
Prefetch HtoD  
Transfer

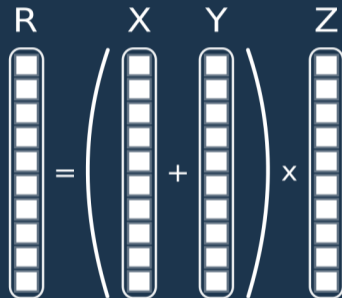




Triadic :  $z = x + y$

Triadic :  $z = x + y$

Quadriadic Computation

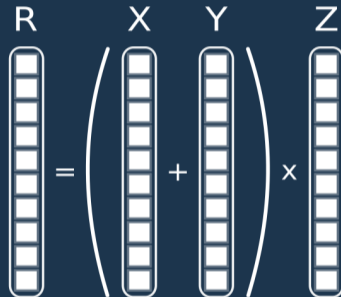


Triadic :  $z = x + y$

## Classic C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >>     tabRes[i] = (tabX[i] + tabY[i])*tabZ[i];
}
```

Quadriadic Computation



Triadic :  $z = x + y$

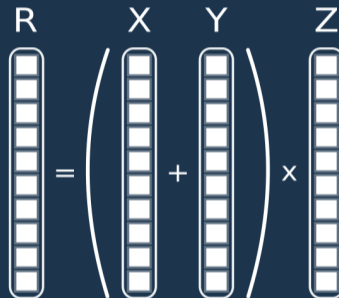
## Classic C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >>     tabRes[i] = (tabX[i] + tabY[i])*tabZ[i];
}
```

## C++ 17 / 20

```
std::transform(std::execution::par_unseq,
    >>     std::begin(vecIndex), std::end(vecIndex),
    >>     std::begin(vecX), std::begin(vecRes),
    >>     [=](int i, float x){
    >>         >>         return (x + vecY[i]) * vecZ[i];
    >>     });
```

Quadriadic Computation



Triadic :  $z = x + y$

## Classic C++

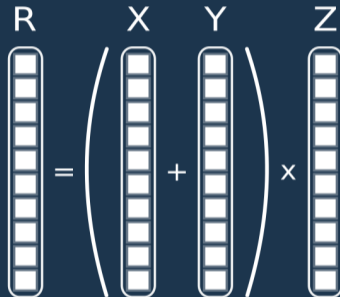
```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >> tabRes[i] = (tabX[i] + tabY[i])*tabZ[i];
}
```

## C++ 17 / 20

```
std::transform(std::execution::par_unseq,
    >> std::begin(vecIndex), std::end(vecIndex),
    >> std::begin(vecX), std::begin(vecRes),
    >> [=](int i, float x){
    >>     >> return (x + vecY[i]) * vecZ[i];
    >> });
```

**vecY, vecZ have to be std::vector**

## Quadriadic Computation



# std::transform : triadic

Triadic :  $z = x + y$

## Classic C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >> tabRes[i] = (tabX[i] + tabY[i])*tabZ[i];
}
```

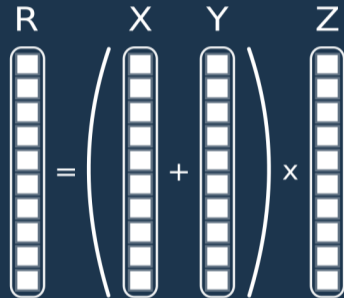
## C++ 17 / 20

```
std::transform(std::execution::par_unseq,
    >> std::begin(vecIndex), std::end(vecIndex),
    >> std::begin(vecX), std::begin(vecRes),
    >> [=](int i, float x){
    >>     >> return (x + vecY[i]) * vecZ[i];
    >> });
```

**vecY, vecZ have to be std::vector**

Works well on GPU

## Quadriadic Computation



# std::transform : triadic

Triadic :  $z = x + y$

## Classic C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >> tabRes[i] = (tabX[i] + tabY[i])*tabZ[i];
}
```

## C++ 17 / 20

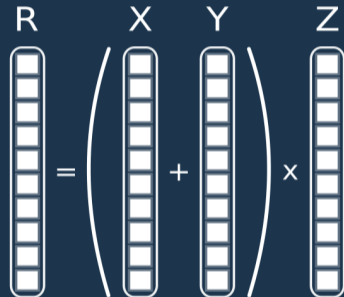
```
std::transform(std::execution::par_unseq,
    >> std::begin(vecIndex), std::end(vecIndex),
    >> std::begin(vecX), std::begin(vecRes),
    >> [=](int i, float x){
    >>     >> return (x + vecY[i]) * vecZ[i];
    >> });
```

**vecY, vecZ have to be std::vector**

Works well on GPU

Needs extra index table

Quadriadic Computation





# std::transform : triadic

Triadic :  $z = x + y$

## Classic C++

```
for(long unsigned int i(0lu); i < nbElement; ++i){
    >> tabRes[i] = (tabX[i] + tabY[i])*tabZ[i];
}
```

## C++ 17 / 20

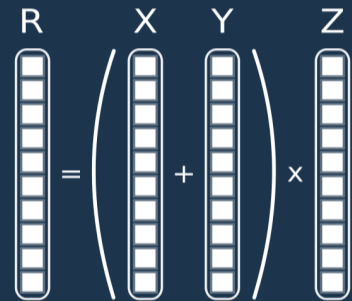
```
std::transform(std::execution::par_unseq,
    >> std::begin(vecIndex), std::end(vecIndex),
    >> std::begin(vecX), std::begin(vecRes),
    >> [=](int i, float x){
    >> >> return (x + vecY[i]) * vecZ[i];
    >> });
```

**vecY, vecZ have to be std::vector**

Works well on GPU

Needs extra index table

## Quadriadic Computation



Close to **CUDA Kernel**



## Cartesian Product and iota views

```

std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M));

std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);

```

## Cartesian Product and iota views

```

std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M));

std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);

```

Multi-dimensional index

## Cartesian Product and iota views

```

std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M));
std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);

```

**No Index Allocation**

**Multi-dimensional index**

## Cartesian Product and iota views

```

std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M));
std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);

```

**No Index Allocation**

Multi-dimensional index

OK for broadcasting but  
we can do better

## Cartesian Product and iota views

```
std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M));
std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);
```

**No Index Allocation**

Multi-dimensional index

OK for broadcasting but  
we can do better

## Broadcast



## Cartesian Product and iota views

```
std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M)); No Index Allocation

std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);
```

Multi-dimensional index

OK for broadcasting but  
we can do better

## Broadcast



```
auto extendedVecB = std::views::join(std::views::repeat(vecB));
std::transform(std::execution::par, vecA.begin(), vecA.end(),
    std::begin(extendedVecB), vecResult.begin(),
    [=](float a, float b){
        return a * b;
    }
);
```



## Cartesian Product and iota views

```
std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M)); No Index Allocation

std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);
```

Multi-dimensional index

OK for broadcasting but  
we can do better

## Broadcast



```
auto extendedVecB = std::views::join(std::views::repeat(vecB));
std::transform(std::execution::par, vecA.begin(), vecA.end(),
    std::begin(extendedVecB), vecResult.begin(),
    [=](float a, float b){
        return a * b;
    }
);
```

Not vectorized yet

## Cartesian Product and iota views

```
std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M)); No Index Allocation

std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);
```

Multi-dimensional index

OK for broadcasting but  
we can do better

## Broadcast



```
auto extendedVecB = std::views::join(std::views::repeat(vecB));
std::transform(std::execution::par, vecA.begin(), vecA.end(),
    std::begin(extendedVecB), vecResult.begin(),
    [=](float a, float b){
        return a * b;
    }
);
```

Not vectorized yet

Not parallelized yet

## Cartesian Product and iota views

```
std::mdspan A{input, N, M};
std::mdspan B{output, M, N};
auto v = std::ranges::views::cartesian_product(
    std::ranges::views::iota(0, N),
    std::ranges::views::iota(0, M)); No Index Allocation

std::for_each(std::execution::par_unseq,
    begin(v), end(v),
    [=] (auto idx) {
        auto [i, j] = idx;
        B[j, i] = A[i, j];
    }
);
```

Multi-dimensional index

OK for broadcasting but we can do better

But only with G++ 13 For now

## Broadcast



```
auto extendedVecB = std::views::join(std::views::repeat(vecB));
std::transform(std::execution::par, vecA.begin(), vecA.end(),
    std::begin(extendedVecB), vecResult.begin(),
    [=](float a, float b){
        return a * b;
    }
);
```

Not vectorized yet

Not parallelized yet

CPU

GPU

CPU

GPU

C++ 17

CPU

GPU

C++ 20

C++ 17

CPU

GPU

C++ 23

C++ 20

C++ 17





- ▶ Good performances on **GPUs**
  - ▶ With **nvcc (CUDA)**
  - ▶ With **nvc++ (C++17 / C++20)**
- ▶ Tests with **HPC SDK 24.3**
- ▶ Compiler **nvc++** powerful and easy to use with **C++17** and **C++20**
  - ▶ **No explicit linking**
  - ▶ Automatic **GPU Targeting** or with **CUDA\_VISIBLE\_DEVICES**
  - ▶ **Avoid static allocation**
- ▶ Warning about industrial software
  - ▶ Will to drive for update
  - ▶ Old **GPUs** become **obsolete** :
    - ▶ **nvc++** : compute capabilities  $\geq 6$  (no **K80**)
    - ▶ **nvcc** : compute capabilities  $\geq 3.5$  (no **K80** in **CUDA 12**)
  - ▶ Need to save binaries to ensure long usability of **GPUs**