ESCAPE 824064



GPU Computing with C++20 nvc++ and CUDA

Pierre Aubert











































3



Event



Pierre Aubert, GPU Computing with C++20, nvc++ and CUDA





















Event

Express Parallelism

Independent Events => Independent Computing => Parallelism





Event

Express Parallelism

Independent Events => Independent Computing => Parallelism























































Kernel Development









Kernel Development



















Express Global Computation -> Linear Algebra



Express Global Computation -> Linear Algebra

Calibration -> Broadcast







Express Global Computation -> Linear Algebra

Calibration -> Broadcast Signal Reduction -> SGEMV





Express Global Computation -> Linear Algebra

Calibration -> Broadcast Signal Reduction -> SGEMV







Express Global Computation -> Linear Algebra

Calibration -> Broadcast Signal Reduction -> SGEMV Momenta computation -> SGEMM




Computing on GPU with CUDA

Express Global Computation -> Linear Algebra

Calibration -> Broadcast Signal Reduction -> SGEMV Momenta computation -> SGEMM







Thread

 \bullet







































































•••



•••











































8






















Thread Mapping









Pierre Aubert, GPU Computing with C++20, nvc++ and CUDA



• •



igodol



































••

10



Hadamard product on 1000 elements Z X Y

10





Hadamard product on 1000 elements z 100

10

Scenario #1 10 1d blocks of size 100



Warp of 32 threads

Scenario #1 10 1d blocks of size 100



10







10





















11



1 sector : 32 bytes

Cache line size :

- 128 bytes
- 4 sectors

















0 32 64 96 128 160 192 224 256 288 320 352 384







....

12



12

Read / Write :

- 1 sector

- 32 bytes



Read / Write : - 1 sector - 32 bytes	Array of Structure (AoS)	struct Coeff{ float u, v, w; float x[8], y[8]	80 Bytes , z;
		;	

12



Read / Write : - 1 sector

- 32 bytes

Array of Structure (AoS) Warp of 32 threads struct Coeff{ 80 Bytes
 float u, v, w;
 float x[8], y[8], z;
};











};


Memory Layout



Memory Layout













__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];















13





13





13





13





13





13





13







```
#include <cuda.h>
#include <cuda runtime.b>
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];
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, sizeBvte);
       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:
       lelse/
                dimBlockX = nbElement;
       //Transfer to Device
       cudaMemcpv(vecLeftd, vecLeft, sizeBvte, cudaMemcpvHostToDevice);
       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):
```



13



On A3000 GPU



Δ

On A3000 GPU



14

On A3000 GPU



14

On A3000 GPU



14

On A3000 GPU



14











```
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){</pre>
                h \text{ vec}X[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 * v:
                }):
        // 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;</pre>
        return 0:
```



```
int main(int argc, char** argv){
        size t nbElement(500000000):
        thrust::host vector<float> h vecX(nbElement). h vecY(nbElement)
                                                                          Initialisation
        for(size t i(0lu): i < nbElement: ++i){</pre>
                h \text{ vec}X[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 * v:
                }):
        // 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;</pre>
        return 0:
```















```
#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){</pre>
                h \text{ 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 * v:
        // 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;
```

15



NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, via Spack, and in the Cloud



16



NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, via Spack, and in the Cloud



16



NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, via Spack, and in the Cloud



16



NVIDIA HPC SDK

Available at developer.nvidia.com/hpc-sdk, on NGC, via Spack, and in the Cloud



16



Example : Hadamard Product





Example : Hadamard Product





Example : Hadamard Product

C++

}

for(long unsigned int i(Olu); i < nbElement; ++i){ tabResult[i] = tabX[i]*tabY[i];</pre>




























18





With CMake :

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

18



18





18





18





















19



19



CUDA



19

CUDA





19

CUDA





19

CUDA





19

CUDA



















19







- CUDA HW (0000:01:00.0 - NVIDIA	kernel	-	with a set of some set									_
✓ 100,0% Kernels			oid thrust::THRUST_20	0300_86_NS::cuda_cut	core:_kernel_agent <thru< th=""><th>st::THRUST_200300_86_NS::</th><th>cuda_cub:parallel_for::Par</th><th>allelForAgent<thrust::thru< th=""><th>ST_200300_86_NS::cuda_cub:</th><th>transform:binary_transform</th><th>_f<gnu_cxx::normal_iterator< th=""><th>eflo</th></gnu_cxx::normal_iterator<></th></thrust::thru<></th></thru<>	st::THRUST_200300_86_NS::	cuda_cub:parallel_for::Par	allelForAgent <thrust::thru< th=""><th>ST_200300_86_NS::cuda_cub:</th><th>transform:binary_transform</th><th>_f<gnu_cxx::normal_iterator< th=""><th>eflo</th></gnu_cxx::normal_iterator<></th></thrust::thru<>	ST_200300_86_NS::cuda_cub:	transform:binary_transform	_f <gnu_cxx::normal_iterator< th=""><th>eflo</th></gnu_cxx::normal_iterator<>	eflo
 100,0% _kernel_agent 			oid thrust::THRUST_20	0300_86_NS::cuda_cut	core::_kernel_agent <thru< th=""><th>st::THRUST_200300_86_NS::</th><th>cuda_cub::parallel_for::Par</th><th>allelForAgent<thrust::thru< th=""><th>ST_200300_86_NS::cuda_cub:</th><th>transform:binary_transform</th><th>_f<gnu_coc::normal_iterator< th=""><th>sflo</th></gnu_coc::normal_iterator<></th></thrust::thru<></th></thru<>	st::THRUST_200300_86_NS::	cuda_cub::parallel_for::Par	allelForAgent <thrust::thru< th=""><th>ST_200300_86_NS::cuda_cub:</th><th>transform:binary_transform</th><th>_f<gnu_coc::normal_iterator< th=""><th>sflo</th></gnu_coc::normal_iterator<></th></thrust::thru<>	ST_200300_86_NS::cuda_cub:	transform:binary_transform	_f <gnu_coc::normal_iterator< th=""><th>sflo</th></gnu_coc::normal_iterator<>	sflo
100,0% void thrust::THRU			oid thrust::THRUST_20	0300_86_NS::cuda_cut	core::_kernel_agent <thru< th=""><th>st::THRUST_200300_86_NS::</th><th>cuda_cub::parallel_for::Par</th><th>allelForAgent<thrust::thru< th=""><th>ST_200300_86_NS::cuda_cub:</th><th>transform::binary_transform</th><th>_f<gnu_cxx::normal_iterator< th=""><th>sflo</th></gnu_cxx::normal_iterator<></th></thrust::thru<></th></thru<>	st::THRUST_200300_86_NS::	cuda_cub::parallel_for::Par	allelForAgent <thrust::thru< th=""><th>ST_200300_86_NS::cuda_cub:</th><th>transform::binary_transform</th><th>_f<gnu_cxx::normal_iterator< th=""><th>sflo</th></gnu_cxx::normal_iterator<></th></thrust::thru<>	ST_200300_86_NS::cuda_cub:	transform::binary_transform	_f <gnu_cxx::normal_iterator< th=""><th>sflo</th></gnu_cxx::normal_iterator<>	sflo
✓ 40,5% Unified memory												
 100,0% Memory 												
>99,9% HtoD transfer												
<0,1% DtoH transfer												
 Threads (8) 												
▼ [38155] simple_hadamard												
OS runtime libraries												mmap
CUDA API							cudaStreamSy	nchronize				cuMemFr
 [38164] simple_hadamard 												
OS runtime libraries		poll	poll	poll p	oll poll	poll poll	poll pol	poll	poll poll	poll poll po	li poli poli	poll























21

Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$



Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$

Quadriadic Computation





Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$



Quadriadic Computation





Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$



C++ 17 / 20

Quadriadic Computation





Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$



Quadriadic Computation

Pierre Aubert, GPU Computing with C++20, nvc++ and CUDA

Ζ

x


std::transform : triadic

Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$





std::transform : triadic

Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$



Quadriadic Computation





std::transform : triadic

Ζ

Triadic : $\mathbf{z} = \mathbf{x} + \mathbf{y}$





•••

22





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.

B[j, i] =

/ Multi-dimensional index

<u>j) = id</u>x; <u>A[i, j</u>];



Cartesian Product and iota views







Cartesian Product and iota views



OK for broadcasting but we can do better









OK for broadcasting but we can do better







OK for broadcasting but we can do better

















23



















- 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