# *Lecture: Deep Learning and Differential Programming*

# 3.3 GPUs - Software

INSA LYON — Christian Wolf

# APIs

GPUs can be programmed on multiple different levels of abstraction.

We will (very briefly) study two cases:

CUDA/C++: direct low-level GPU programming

```
1  __global__
2  void mult_kernel_simple(int mxWidth, float *mx1, float *mx2, float *output)
3  {
4      int c = blockIdx.x*blockDim.x + threadIdx.x;
```
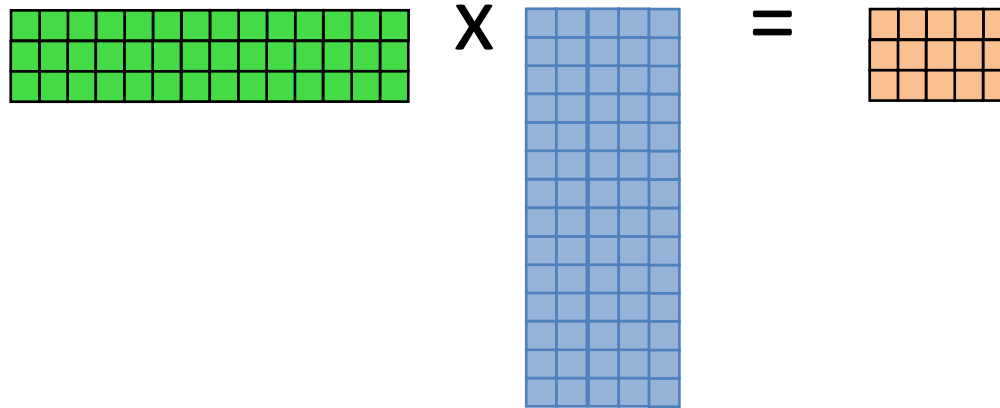
PyTorch/Python: high-level of abstraction

```
31  model = LeNet().to("cuda:0")
32
33  for batch_idx, (data, labels) in enumerate(train_loader):
34      data = data.to("cuda:0")
```

LIRIS cnrs

# 1 CUDA – low level GPU programmng

# 2 GPUs & PyTorch

# Example : matrix multiplication



The classical sequential solution:

```
1    for( int r = 0; r < h; r++)
2    {
3        for( int c = 0; c < w; c++)
4        {
5            m1xm2[r*w + c] = 0;
6
7            for( int k = 0; k < w; k++)
8                m1xm2[r*w + c] += m1[r*w + k] * m2[k*w + c];
9        }
10   }
```

# The parallel solution

Parallel execution of a function called for each individual result value of the result matrix, with arguments being the indices of the value.

```
1  void mult_kernel_simple(int c, int r)
2  {
3      output[r*mxWidth + c] = 0.0f;
4      for( int k = 0; k < mxWidth; k++)
5          output[r*mxWidth + c] += mx1[r*mxWidth + k] * mx2[k*mxWidth + c];
6  }
```

Called in parallel for all c,r

In CUDA, this function is called a *kernel.*

Each tuple (c,r) corresponds to a *thread.*

# Organisation into *blocks*

– Threads are organized into *blocks*
– Faster local memory can be shared by threads of the same block.

```
1  void mult_kernel_simple(int c, int r)
2  {
3      output[r*mxWidth + c] = 0.0f;
4      for( int k = 0; k < mxWidth; k++)
5          output[r*mxWidth + c] += mx1[r*mxWidth + k] * mx2[k*mxWidth + c];
6  }
```
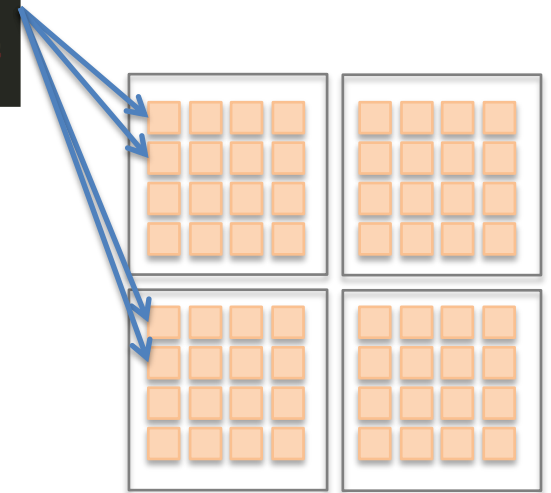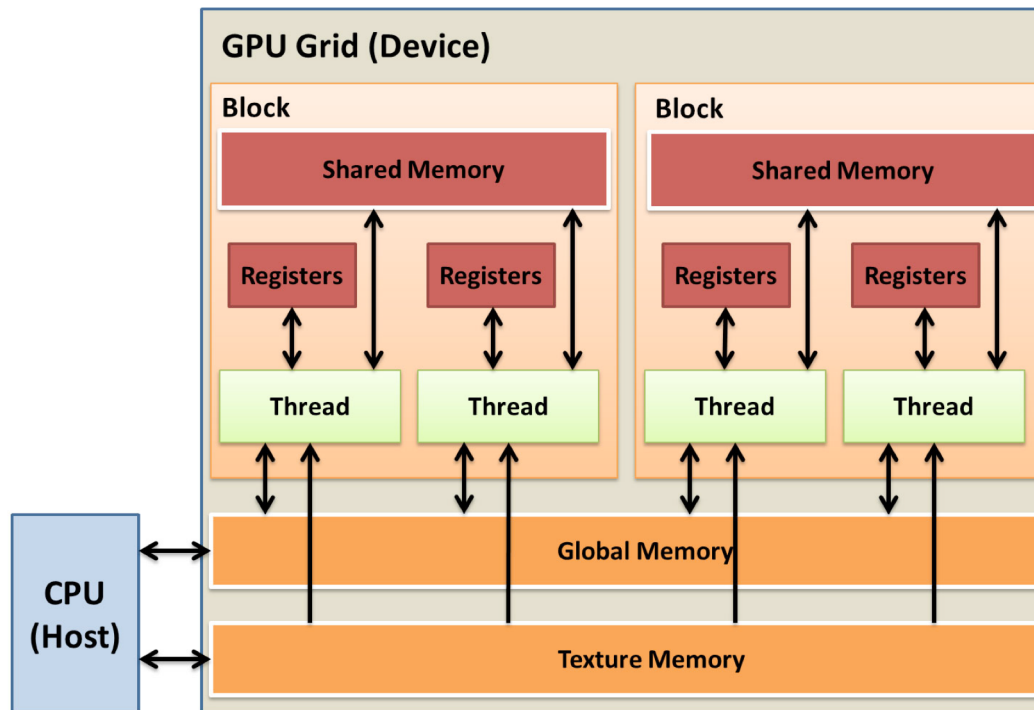
# SM

A block will be sent to a SM (*streaming multi-processor*)

# The classical CUDA sequence

– The CPU allocates memory on the GPU
– The CPU copies the data to GPU global memory
– The CPU launches the kernel on the GPU
– The GPU executes the kernel in parallel
– The CPU copies the result data back to CPU host memory

# The CUDA syntax of the kernel

Key word declares the *kernel*

```
1  __global__
2  void mult_kernel_simple(int mxWidth, float *mx1, float *mx2, float *output)
3  {
4      int c = blockIdx.x*blockDim.x + threadIdx.x;
5      int r = blockIdx.y*blockDim.y + threadIdx.y;
6
7      output[r*mxWidth + c] = 0.0;
8      for( int k = 0; k < mxWidth; k++)
9          output[r*mxWidth + c] += mx1[r*mxWidth + k] * mx2[k*mxWidth + c];
10 }
```

Index of the block in the grid

Index of the thread in the block

# Calling the kernel

```
1   int main(void)
2   {
3       ...
4       // allocate memory on GPU
5       cudaMalloc( (void**) &gpuMatrix1, matrixSizeInBytes);
6       ...
7
8       // copy data from CPU memory to GPU memory
9       cudaMemcpy(gpuMatrix1, matrix1, matrixSizeInBytes, cudaMemcpyHostToDevice);
10      ...
11
12      // Set grid and block size
13      dim3 dimBlock(32, 32);
14      dim3 dimGrid(matrixWidth/dimBlock.x, matrixWidth/dimBlock.y);
15
16      // run kernel
17      mult_kernel_simple<<<dimGrid, dimBlock>>>( matrixWidth, gpuMatrix1, gpuMatrix2, gpuO
18
19      // copy back results from GPU memory to CPU memory
20      cudaMemcpy( outputData, gpuOutput, matrixSizeInBytes, cudaMemcpyDeviceToHost);
21      ...
22  }
```

Call the *kernel* in parallel for a set of threads

# Compilation

CUDA uses a specific compiler, which is based on a generic C++ compiler (gcc)

```
g++ -c  -O3 -Wall file_seqcode.cpp -o file_seqcode.o

nvcc -c -O3 file_parcode.cu -o file_parcode.o

g++ file_seqcode.o file_parcode.o -o executable
```
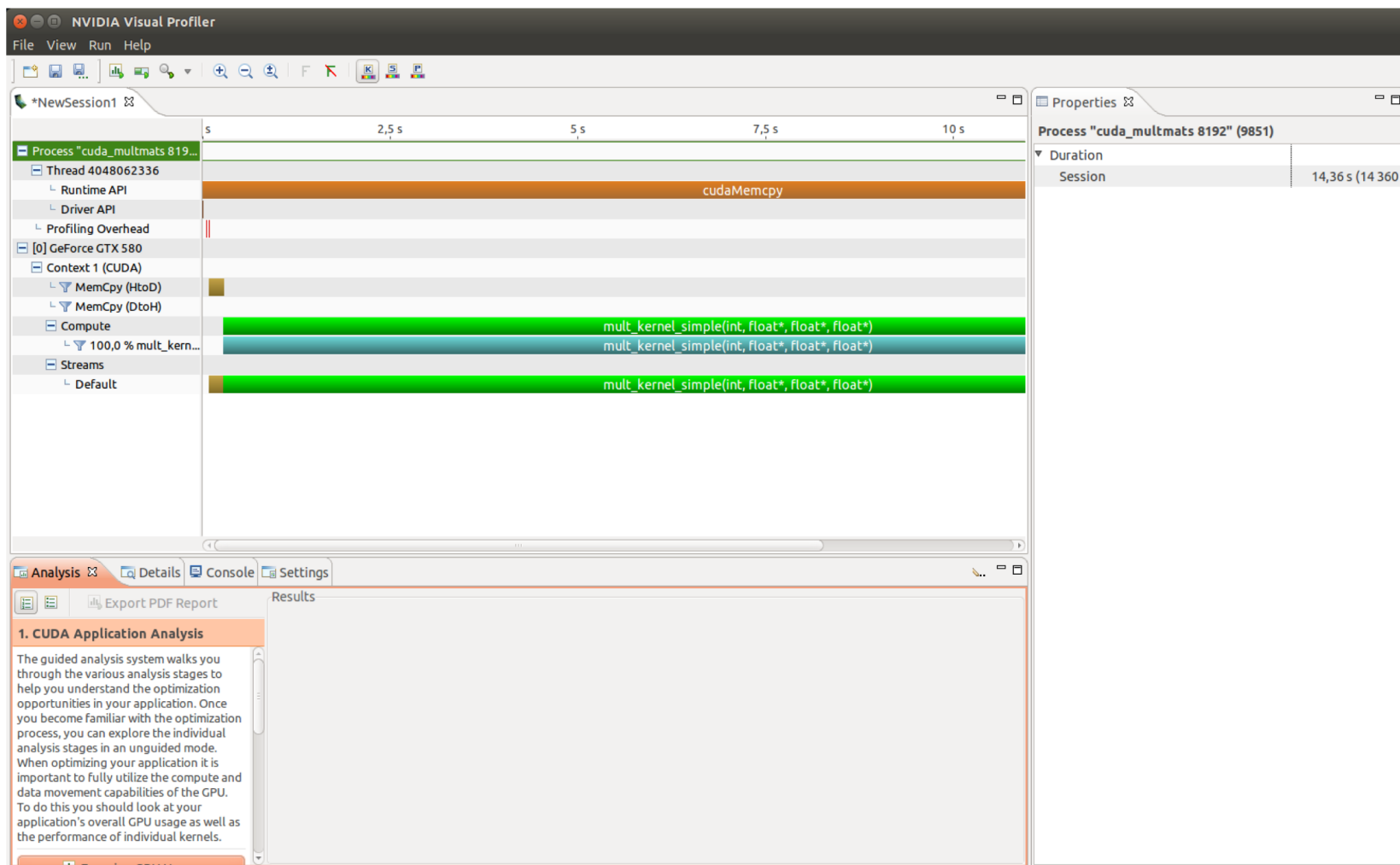
Classical C++ code          CUDA code                    Linking

# Debugging and profiling

The nvvp profiler is part of the CUDA toolkit.

# 1 CUDA – low level GPU programmng

# 2 GPUs & PyTorch

# The PyTorch GPU interface

The model is transferred to GPU memory with `.to(device)`
Device can be "cpu", "cuda:0", "cuda:1" etc.

```python
1 model = LeNet()
2 model = model.to("cuda:0")
```

We also send the data to GPU memory.
We get data back to the CPU with the `.cpu()` method:

```python
1  # Cycle through batches
2  for idx, (data,labels) in enumerate(train_loader):
3    data = data.to("cuda:0")
4    optimizer.zero_grad()
5    y = model(data)
6    loss = crossentropy(y, labels)
7    loss.backward()
8    running_loss += loss.cpu().item()
9    optimizer.step()
10
11   _, predicted = torch.max(y.data.cpu(), 1)
```

# PyTorch vs. Cuda

— For standard functions (Linear, Conv2d, Pooling etc.), PyTorch ships GPU support.

— PyTorch allows to write custom neural network layers (requiring to specify the forward and and the backward pass)

— Custom layers require CUDA programming to run on GPUs.