

# 1.<sup>st</sup> Task in *HPC-I*

Deadline: Nov 5, 2024, 11:59pm

---

## GPU computing using CUDA

---

### 1. Start with CUDA

Install CUDA on your computer

- manjaro-Linux: `sudo pacman -S nvidia nvidia-utils cuda` (hints<sup>1</sup>)
- ubuntu-Linux: `sudo apt install nvidia-cuda-toolkit cuda` (hints<sup>2</sup>)
- Running Ubuntu in Windows via WSL2 requires appropriate drivers which is described in general by Microsoft<sup>3</sup> and in detail by NVIDIA<sup>4</sup>.

and check

- Compiler: `nvcc --version`
- GPU/CUDA: `nvidia-smi` (or `nvidia-smi -L`)
- example Code<sup>5</sup>

```
cd CUDA/firstSteps
nvcc data_mv_GH.cu
./a.out
```

### 2. Your first CUDA Code

Copy the *CUDA* directory into your directory and download<sup>6</sup>

the Code from [HaSh19] and use `float` as data type for all tasks.

Start with your copy of file *data\_mv\_GH.cu*:

[5 pts]

- The given code realizes on GPU  $\underline{b} := \underline{a}$  followed by  $++\underline{b}$ .  
Extend the code with a kernel function for  $\underline{c} := \underline{a} + \underline{b}$ . Check the result on the host.
- Compare your code with code *vector\_addition\_gpu\_thread\_block.cu*<sup>7</sup> from [HaSh19, §1].
- Extend your code with two kernel functions for  $\underline{b} := \ln(\underline{a})$  and  $\underline{c} := \exp(\underline{b})$ . Check the result  $\underline{c}$  on the host.
- Write a second main function that uses unified memory<sup>8</sup> (intro<sup>9</sup>) instead of the `malloc-cudaMalloc-cudaMemcpy` framework. Have a look at code *unified\_memory.cu*<sup>10</sup> from [HaSh19, p.70f] on how to use `cudaMallocManaged`.
- Add timing for the kernel calls, see general performance metrics<sup>11</sup> on GPU.

---

<sup>1</sup><https://miloserdov.org/?p=4181>

<sup>2</sup><https://linuxconfig.org/how-to-install-cuda-on-ubuntu-20-04-focal-fossa-linux>

<sup>3</sup><https://docs.microsoft.com/en-us/windows/win32/direct3d12/gpu-cuda-in-wsl>

<sup>4</sup><https://docs.nvidia.com/cuda/wsl-user-guide/index.html>

<sup>5</sup>[http://imsc.uni-graz.at/haasegu/Lectures/GPU\\_CUDA/CUDA.zip](http://imsc.uni-graz.at/haasegu/Lectures/GPU_CUDA/CUDA.zip)

<sup>6</sup><https://github.com/PacktPublishing/Learn-CUDA-Programming>

<sup>7</sup>[https://github.com/PacktPublishing/Learn-CUDA-Programming/tree/master/Chapter02/02\\_memory\\_overview/02\\_vector\\_addition/vector\\_addition\\_gpu\\_thread\\_block.cu](https://github.com/PacktPublishing/Learn-CUDA-Programming/tree/master/Chapter02/02_memory_overview/02_vector_addition/vector_addition_gpu_thread_block.cu)

<sup>8</sup><https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html#um-unified-memory-programming-hd>

<sup>9</sup><https://devblogs.nvidia.com/unified-memory-cuda-beginners/>

<sup>10</sup>[https://github.com/PacktPublishing/Learn-CUDA-Programming/blob/master/Chapter02/02\\_memory\\_overview/06\\_unified\\_memory/unified\\_memory.cu](https://github.com/PacktPublishing/Learn-CUDA-Programming/blob/master/Chapter02/02_memory_overview/06_unified_memory/unified_memory.cu)

<sup>11</sup><https://devblogs.nvidia.com/how-implement-performance-metrics-cuda-cc/>

- By using CPU-timing (`clock()` or `std::chrono::system_clock`). Don't forget the call `cudaDeviceSynchronize()` to wait until the kernel functions finish on GPU.
- By using `cudaEvent_t`, see example.

**3. Reduction** We need frequently a reduction function, i.e., to compare two vectors on the GPU in the previous tasks.

Have a look at codes in *CUDA/skalar* computing the inner product of two vectors.

- *skalar\_3\_fast.cu*: The data management is similar to the previous task. The inner product calculation follows Mark Harris' presentation in *CUDA\_intro/reduction\_Mark\_Harris.pdf* with its own kernel functions.

```
nvcc --ptxas-options=-v -restrict skalar_3_fast.cu
```

- *skalar\_4.cu*: Uses cuBLAS for inner product calculation.

```
nvcc -restrict skalar_4.cu -lcublas
```

- *skalar\_5.cu*: Uses the Thrust<sup>12</sup> library for vector management in combination with STL-algorithms. See the documentation<sup>13</sup>.

```
nvcc --ptxas-options=-v -restrict skalar_5.cu
```

Your tasks for realizing  $\underline{c} == \underline{a}$  on GPU:

[8 pts]

- Write a kernel function (e.g.: `equal`) for comparing  $\underline{c} == \underline{a}$ . You need a reduction operation similar to the inner product in *CUDA/skalar/skalar\_3\_fast.cu*.
- See also the reduction kernel in [HaSh19, §3, p.117-126] and the code versions<sup>14</sup> of it.
- Copy and Rewrite your code from Task 2 by using Thrust. You might use only `thrust::reduce()`, see [NLS14] for combining unified memory with Thrust. Even better, your code can use `thrust::inner_product`, see the documentation<sup>15</sup> and predefined function objects<sup>16</sup>.
- ?? Can we call cuBLAS routines as `cublasDdot` with unified memory vectors or/and Thrust vector ??

**4. Using cuBLAS 1-3** The basic idea of these tasks consist in applying cuBLAS<sup>17</sup> (BLAS) calls for all vector and/or matrix operations. Check also the runtime for reasonable matrix sizes.

Start with columnwise stored dense matrices (example *CUDA/densmatrices.libs*. Use `float` as data type. See also [HaSh19, §8, p.320ff] and its code<sup>18</sup>.

[8 pts]

- BLAS1<sup>19</sup>: Realize  $\underline{y} := \alpha \underline{x} + \underline{y}$ ;  $\underline{x} := \alpha \underline{x} + \underline{y}$ ;  $\underline{z} := \alpha \underline{x} + \beta \underline{y}$ ;  $\langle \underline{x}, \underline{y} \rangle$ ;  $\| \underline{x} \|^2$ ;
- BLAS2: Realize  $\underline{r} = M * \underline{x}$ ;  $\underline{r} = M^T * \underline{x}$  with a non-symmetric dense real matrix  $M$ .

<sup>12</sup><https://developer.nvidia.com/thrust>

<sup>13</sup><https://docs.nvidia.com/cuda/thrust/index.html>

<sup>14</sup>[https://github.com/PacktPublishing/Learn-CUDA-Programming/tree/master/Chapter03/03\\_cuda\\_thread\\_programming](https://github.com/PacktPublishing/Learn-CUDA-Programming/tree/master/Chapter03/03_cuda_thread_programming)

<sup>15</sup>[https://nvidia.github.io/cccl/thrust/api/function\\_group\\_\\_transformed\\_\\_reductions\\_1ga7dd3c0d0f64ef48165e5d56ebda94739.html](https://nvidia.github.io/cccl/thrust/api/function_group__transformed__reductions_1ga7dd3c0d0f64ef48165e5d56ebda94739.html)

<sup>16</sup>[https://nvidia.github.io/cccl/thrust/api\\_docs/function\\_objects/predefined.html](https://nvidia.github.io/cccl/thrust/api_docs/function_objects/predefined.html)

<sup>17</sup><https://developer.nvidia.com/cublas>

<sup>18</sup>[https://github.com/PacktPublishing/Learn-CUDA-Programming/tree/master/Chapter08/08\\_cuda\\_libs\\_and\\_other\\_languages/01\\_sgemm](https://github.com/PacktPublishing/Learn-CUDA-Programming/tree/master/Chapter08/08_cuda_libs_and_other_languages/01_sgemm)

<sup>19</sup>[http://www.netlib.org/blas/#\\_level\\_1](http://www.netlib.org/blas/#_level_1)

- (iii) Realize  $\underline{r} = T \cdot \underline{x}$  with a tridiagonal matrix  $T = [-1, 2, -1]$  (assuming non-constant sub-/diagonal entries). Take advantage of the matrix structure.
- (iv) BLAS3: Realize  $A = M * M^T$  and compare the result of  $\underline{z} := A * \underline{x}$  with  $\underline{y} := M * (M^T * \underline{x})$ .
- (v) Extract the diagonal from a dense matrix  $M$  with a BLAS1-function.

## References

- [HaSh19] Jaegeun Han and Bharatkumar Sharma. *Learn CUDA Programming*. Packt> (2019); buy<sup>20</sup>; download code<sup>21</sup>, download figures<sup>22</sup>
- [NSLS14] Dan Negrut, Radu Serban, Ang Li and Andrew Seidl. *Unified Memory in CUDA 6: A Brief Overview and Related Data Access/Transfer Issues*. TR-2014-09<sup>23</sup> (2014)
- [ChAn17] Andrzej Chruszczczyk and Jacob Anders *Matrix computationson the GPUCUBLAS, CUSOLVER and MAGMA by examples*. NVIDIA (2017); download book<sup>24</sup>.

<sup>20</sup><https://www.packtpub.com/eu/application-development/cuda-cookbook>

<sup>21</sup><https://github.com/PacktPublishing/Learn-CUDA-Programming>

<sup>22</sup>[https://static.packt-cdn.com/downloads/9781788996242\\_ColorImages.pdf](https://static.packt-cdn.com/downloads/9781788996242_ColorImages.pdf)

<sup>23</sup>[https://www.researchgate.net/publication/326920953\\_Unified\\_Memory\\_in\\_CUDA\\_6\\_A\\_Brief\\_Overview\\_and\\_Related\\_Data\\_AccessTransfer\\_Issues](https://www.researchgate.net/publication/326920953_Unified_Memory_in_CUDA_6_A_Brief_Overview_and_Related_Data_AccessTransfer_Issues)

<sup>24</sup><https://developer.nvidia.com/sites/default/files/akamai/cuda/files/Misc/mygpu.pdf>