



Université Claude Bernard



Lyon 1

LBM BASED FLOW SIMULATION USING GPU COMPUTING PROCESSOR

Frédéric Kuznik, frederic.kuznik@insa-lyon.fr

Framework

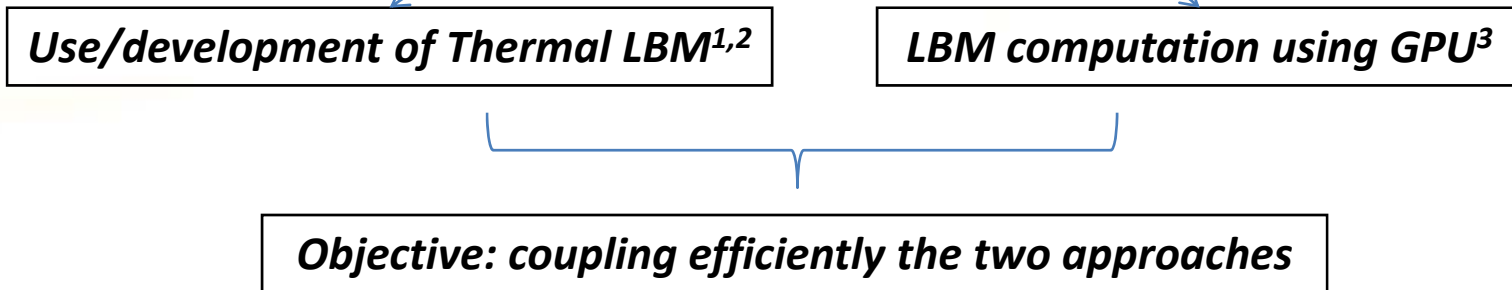
- ④ Introduction
- ④ Hardware architecture
- ④ CUDA overview
- ④ Implementation details
- ④ A simple case: the lid driven cavity problem

Presentation of my activities

🌐 CETHIL: Thermal Sciences Center of Lyon

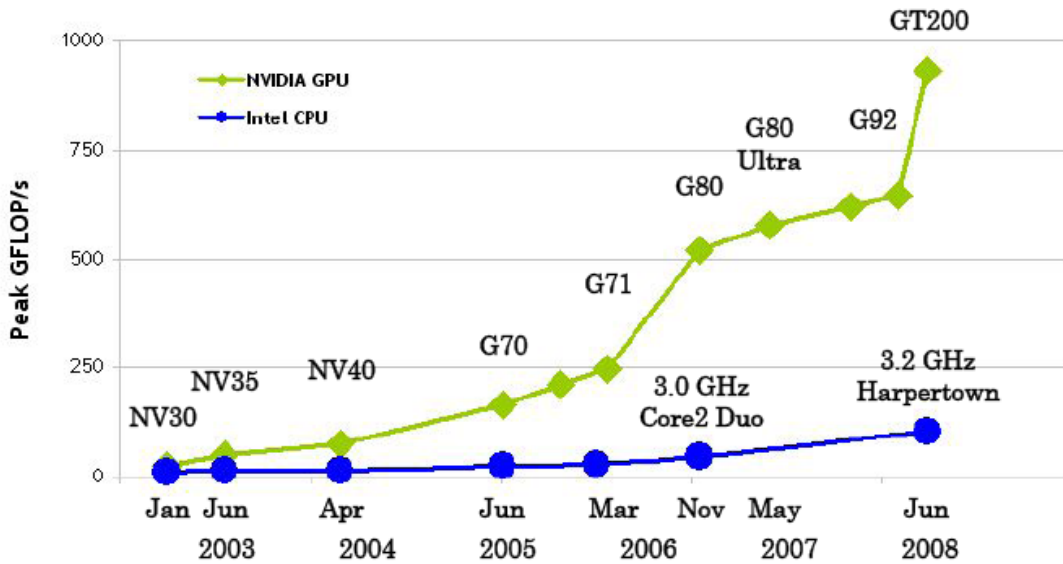


🌐 Current work: simulation of fluid flows with heat transfer



- 1 A double population lattice Boltzman method with non-uniform mesh for the simulation of natural convection in a square cavity, Int. J. Heat and Fluid Flow, vol. 28(5), pp. 862-870, 2007
- 2 Numerical prediction of natural convection occurring in building components: a double population lattice Boltzmann method, Numerical Heat Transfer A, vol. 52(4), pp. 315-335, 2007
- 3 LBM based flow simulation using GPU computing processor, Computers & Mathematics with Applications, under review

Introduction



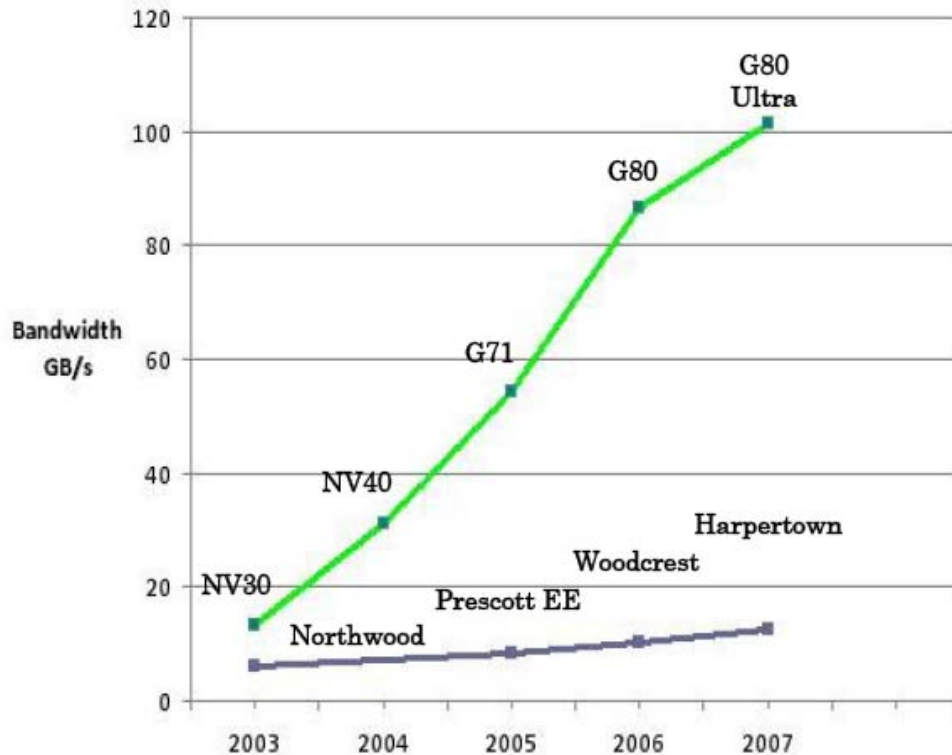
Due to the market (realtime and high definition 3-D graphics), GPU has evolved into highly parallel, multithreaded, multi-core processor.

GT200 = GeForce GTX 280	G71 = GeForce 7900 GTX	NV35 = GeForce FX 5950 Ultra
G92 = GeForce 9800 GTX	G70 = GeForce 7800 GTX	NV30 = GeForce FX 5800
G80 = GeForce 8800 GTX	NV40 = GeForce 6800 Ultra	

(from *CUDA Programming Guide 06/07/2008*)

GFLOPS= 10^9 floating point operations per second

Introduction



The bandwidth of GPU has also evolved for faster graphics purpose.

(from *CUDA Programming Guide 06/07/2008*)

Previous works using GPU & LBM

- ④ GPU Cluster for high performance computing (*Fan et al. 2004*): 32 nodes cluster of GeForce 5800 ultra
- ④ LB-Stream Computing or over 1 Billion Lattice Updates per second on a single PC (*J. Tölke ICMMES 2007*): GeForce 8800 ultra

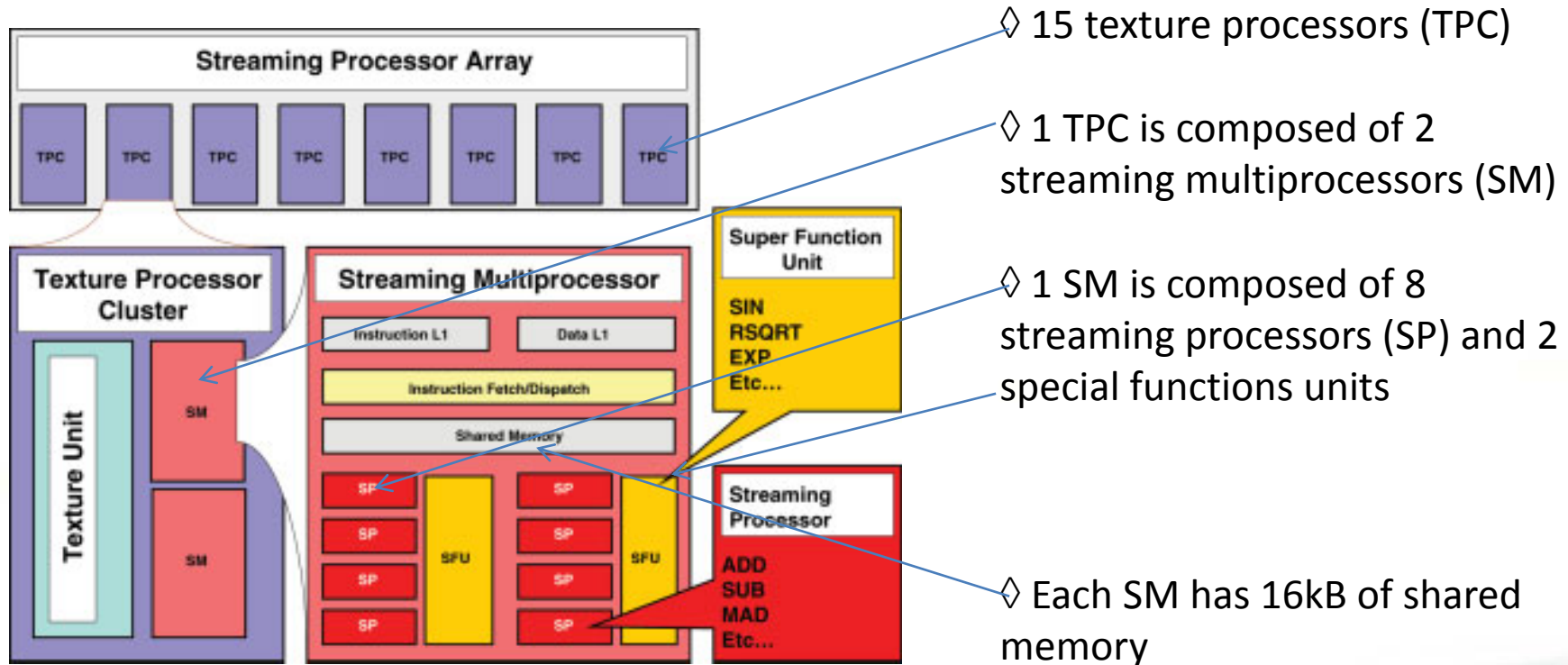
Hardware

- ④ Commercial graphics card (can be included in a desktop PC)
- ④ 240 processors running at 1.35GHz
- ④ 1 Go DDR3 with 141.7GB/s bandwidth
- ④ Theoretical peak at 1000GFLOPS
- ④ Price around 500€



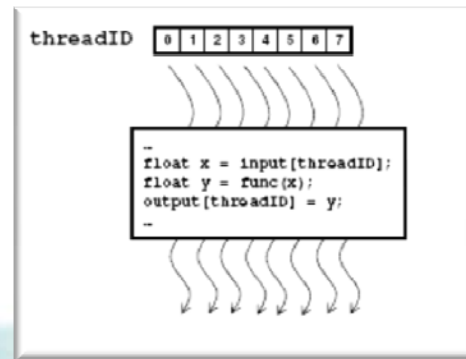
nVIDIA GeForce GTX 280

GPU Hardware Architecture



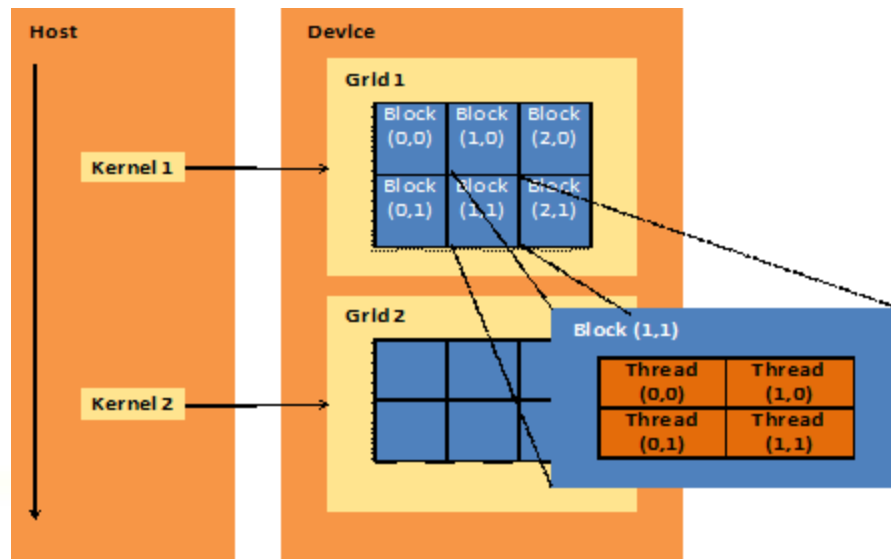
GPU Programming Interface

- ④ NVIDIA® CUDA™ C language programming environment (version 2.0)
- ④ Definitions:
 - ④ Device=GPU
 - ④ Host=CPU
 - ④ Kernel=function that is called from the host and runs on the device
- ④ A cuda kernel is executed by an array of threads

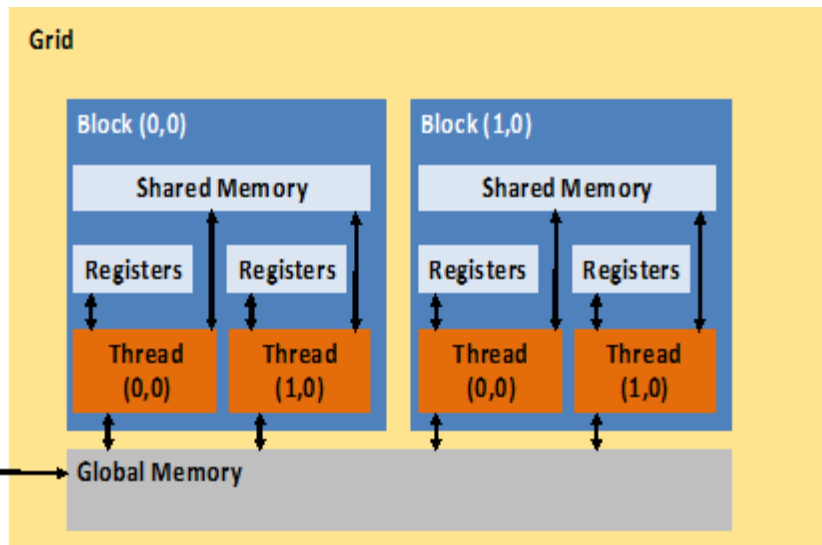


GPU Programming Interface

- ④ A kernel is executed by a grid of thread block
- ④ 1 thread block is executed by 1 multiprocessor
- ④ A thread block contains a maximum of 1024 threads
- ④ Threads are executed by processors within a single multiprocessor
- ④ => Threads from different blocks cannot cooperate



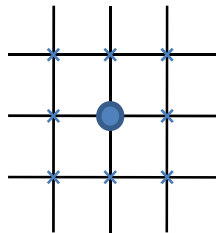
Kernel memory access



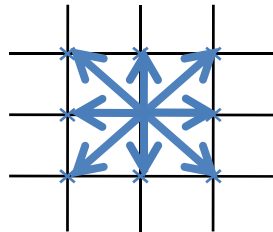
- ④ Registers
- ④ Shared memory: on-chip (fast), small (16kB)
- ④ Global memory: off-chip, large
- ④ The host can read or write global memory only.

LBM ALGORITHM

- Step 1: Collision (Local – 70% computational time)



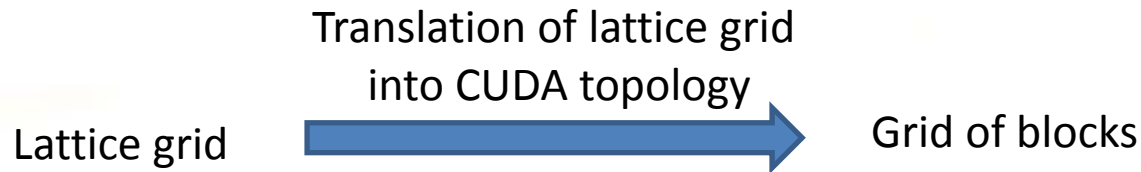
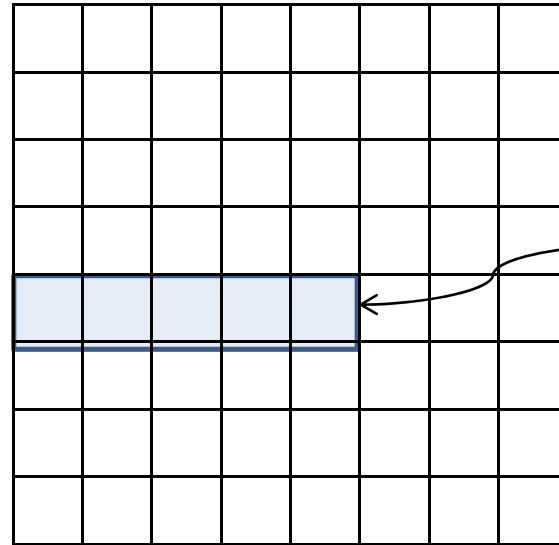
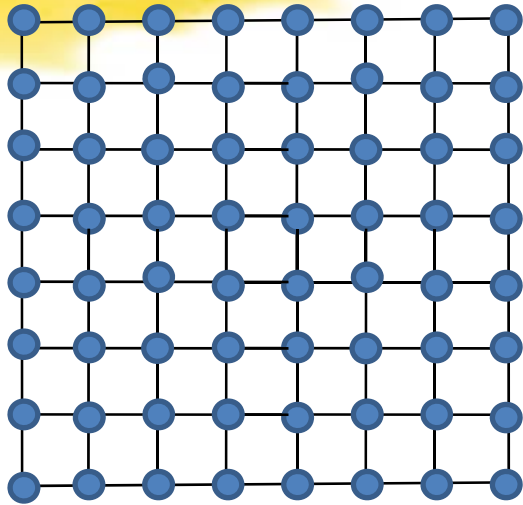
- Step 2: Propagation (Neighbors – 28% computational time)



- + Boundary conditions

(from *BERNSDORF J.*, *How to make my LB-code faster – software planning, implementation and performance tuning, ICMMES'08, Netherlands*)

Implementation details



- 1- Decomposition of the fluid domain into a lattice grid
- 2- Indexing the lattice nodes using thread ID
- 3- Decomposition of the CUDA domain into thread blocks
- 4- Execution of the kernel by the thread blocks

Pseudo-code

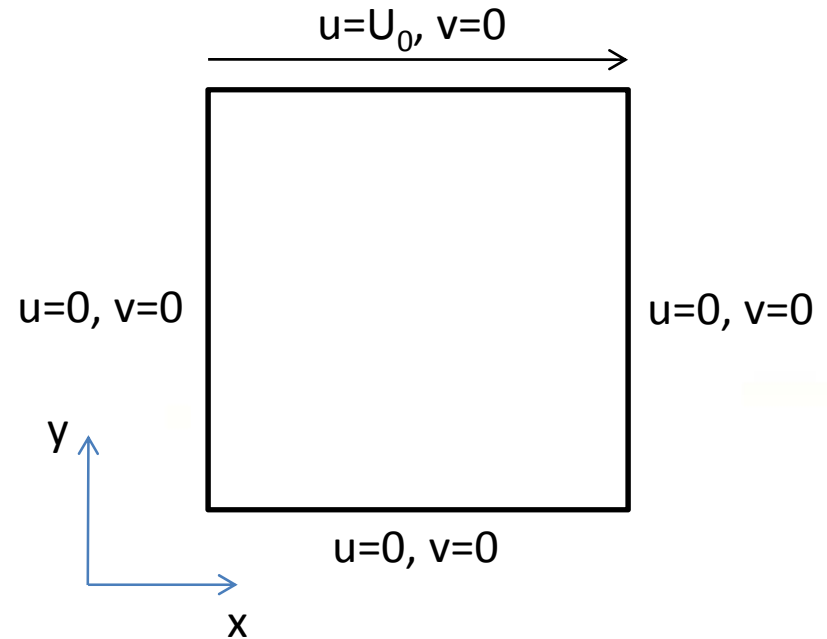
Combine collision and propagation steps:

```
for each thread block
  for each thread
    load  $f_i$  in shared memory
    compute collision step
    do the propagation step
  end
end
```

Exchange informations across boundaries

Application: 2D lid driven cavity

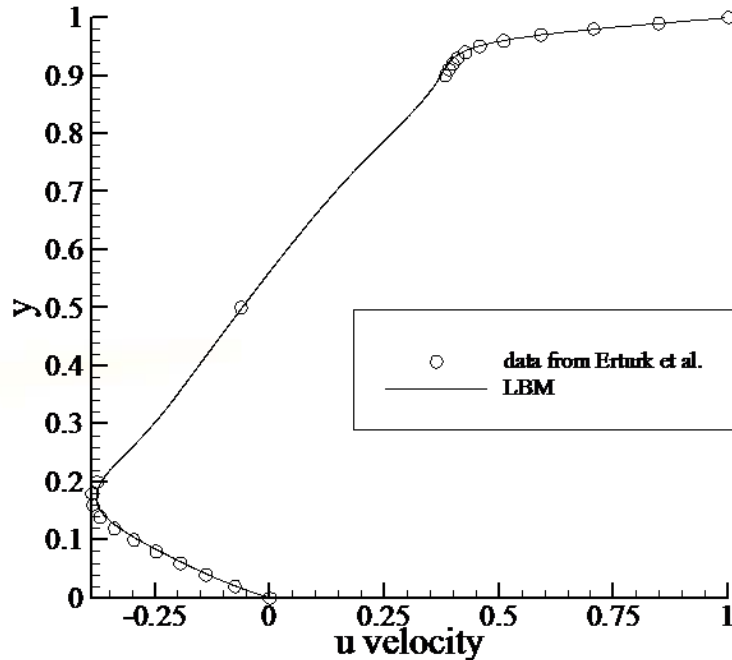
ⓐ Problem description



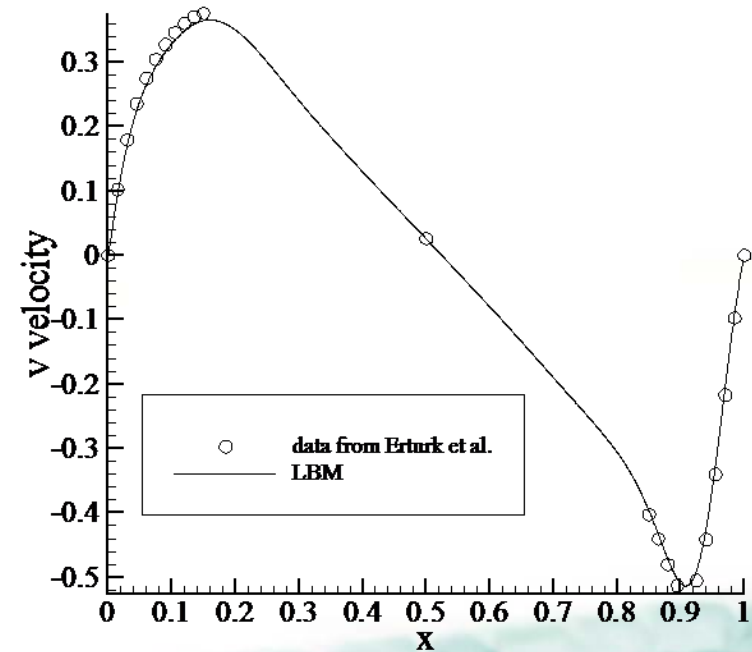
ⓐ D2Q9 MRT model of d'Humières (1992)

Results Re=1000

Reference data from *Erturk et al. 2005*:



Vertical velocity profil for $x=0.5$



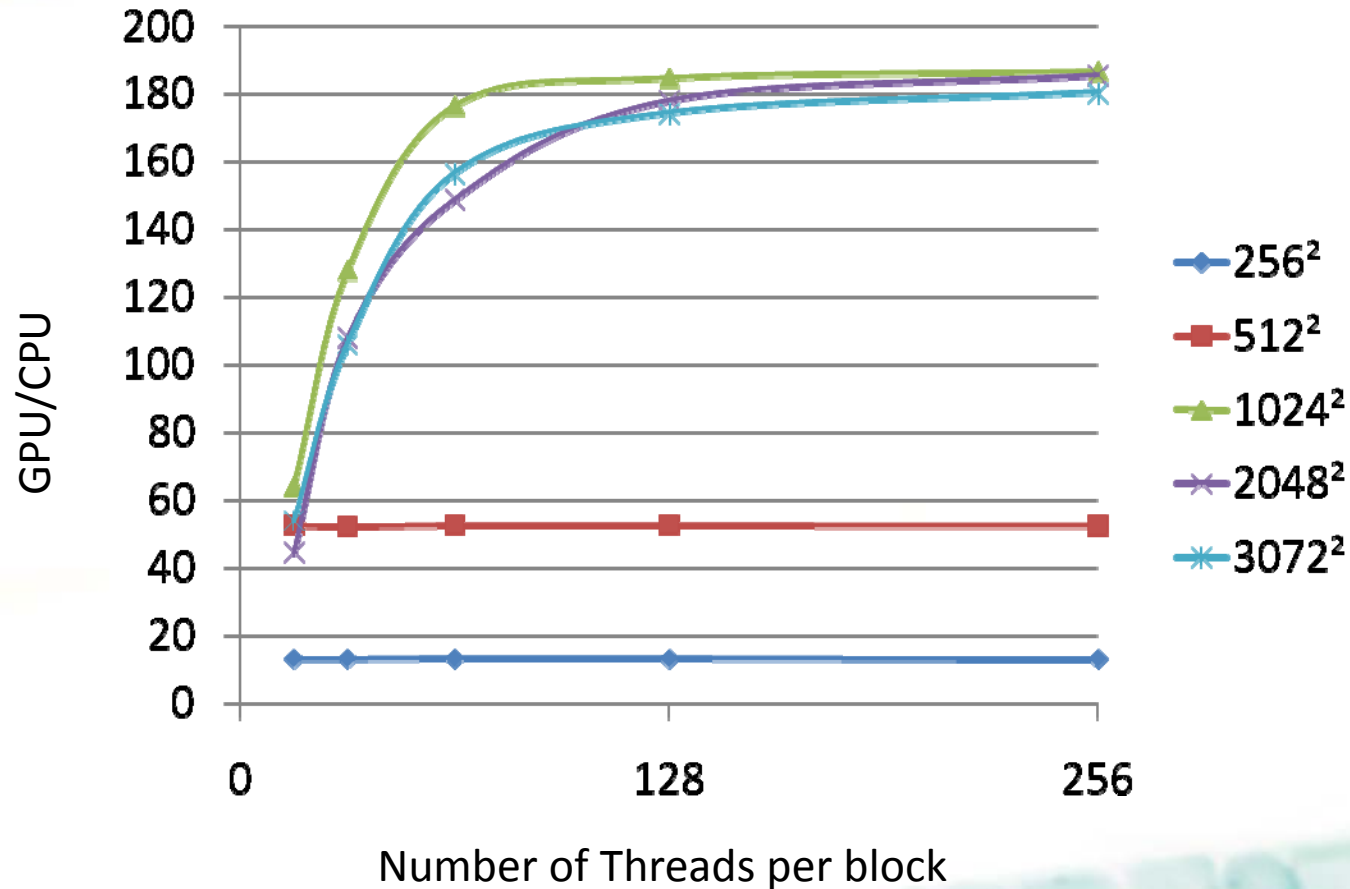
Horizontal velocity profil for $y=0.5$

Performances

Mesh grid size	Number of Threads				
	16	32	64	128	256
256 ²	71.0	71.0	71.1	71.1	70.9
512 ²	284.2	282.7	284.3	284.3	283.7
1024 ²	346.4	690.2	953.0	997.5	1007.9
2048 ²	240.0	583.0	803.3	961.3	1001.5
3072 ²	289.3	572.4	845.2	941.2	974.2

Performances in MLUPS

Performances: CPU comparisons



CPU = Pentium IV, 3.0 GHz

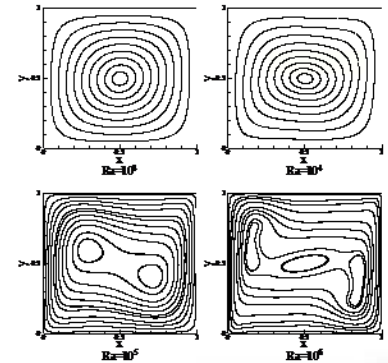
Conclusions

- The GPU allows a gain of about 180 for the calculation time
- Multiple GPU server exists: NVIDIA TESLA S1070 composed of 4 GTX 280 (theoretical gain 720)
- The evolution of GPU is not finished !
- But using double precision floating point, the gain falls to 20 !

Outlooks

- A workgroup concerning LBM and GPU
 - 1 master student with Prof. Bernard TOURANCHEAU (ENS Lyon – INRIA) and a PhD student in 2009
 - 1 master student with Prof. Eric Favier (ENISE – DIPI)

- A workgroup concerning Thermal LBM ...?



Everybody is welcome to join the workgroup !

THANK YOU FOR YOUR ATTENTION

frederic.kuznik@insa-lyon.fr