



**HAL**  
open science

## Some examples of instant computations of fluid dynamics on GPU

Florian de Vuyst, Christophe Labourdette

► **To cite this version:**

Florian de Vuyst, Christophe Labourdette. Some examples of instant computations of fluid dynamics on GPU. CANUM 2012, May 2012, Superbesse, France. [ambp.cedram.org](http://ambp.cedram.org). hal-00732741

**HAL Id: hal-00732741**

**<https://hal.science/hal-00732741>**

Submitted on 16 Sep 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Quelques exemples de calculs instantanés de fluides sur GPU

FLORIAN DE VUYST  
CHRISTOPHE LABOURDETTE

RÉSUMÉ. Dans ce papier, nous partageons notre retour d'expérience sur l'utilisation de GPU et GPGPU pour le calcul de mécanique des fluides bidimensionnels sur grille fine et de problèmes de transport tridimensionnels. Le choix d'une méthode appropriée à l'architecture GPU est critique pour le gain de performance. Pour nos expérimentations numériques, nous testons respectivement une approche Lattice Boltzmann (LBM) pour les équations de Navier-Stokes incompressibles, une méthode de volumes finis de type Flux Vector Splitting (FVS) pour les équations d'Euler compressibles, et une approche particulaire lagrangienne pour le transport cinétique libre.

*Some examples of instant computations of fluid dynamics on GPU*

ABSTRACT. This paper is a summary of our experience feedback on GPU and GPGPU computing for two-dimensional computational fluid dynamics using fine grids and three-dimensional kinetic transport problems. The choice of the computational approach is clearly critical for both performance speedup and efficiency. In our numerical experiments, we used a Lattice Boltzmann approach (LBM) for the incompressible Navier-Stokes equations, a finite volume Flux Vector Splitting (FVS) method for the compressible Euler equations and a lagrangian particle approach for a linear kinetic problem.

---

*Mots-clés:* EDP, GPU, Mécanique des Fluides, interaction, visualisation, calcul instantané, volumes finis, méthode Lattice Boltzmann, méthode particulaire, programmation multicœur.

*Classification math.:* 35L05, 65M08, 76M25, 76N15, 76P05, 97N40.

## 1. GPU computing with instantaneous visualization and interaction in CFD. Experience feedback

### 1.1. Instantaneous computing and interactive visualization

High performance computing (HPC) knows an important growth since recent years. Theoretical peak processing performance and storage capacity in supercomputers gained three or four orders of magnitude in less than a decade. However, scientists and computational engineers would like more flexibility in terms of delay of response and ease of use. Manufacturers develop cluster computers to exceed the petaflop (see the exaflop!), but the cost and planning of very large computations imposes workflow constraints in batch mode.

Recently, the design of manycore processors like graphics processing units GPU, general purpose graphics processing units GPGPU and many-integrated cores MIC allow one to get theoretical teraflop performance into a simple office workstation. This potential flexibility of use with only one user let us imagine new ways of computing and use cases like interactive simulation and instant computations. The applied mathematicians are often little concerned with the very large calculations, they are more interested in the design of methods and algorithms for performing the calculations. That's why we emphasize here on ways of instant computing, in particular on GPU or GPGPU. We especially focus on fluid dynamics problems where time scales of interest allow for interaction with the simulation, and where the models can be controlled by changing parameters and operating conditions with effects viewed instantaneously.

The spin-off effect of such applications is the ease with which a user may “play” with the computational method and the underlying Physics thanks to the visualization. We believe that the coupling between instantaneous computing and interactive visualization really brings an extra dimension to better understand and evaluate methods or schemes. It is a new valuable tool for the applied mathematician. GPU computing today seems to be an inevitable affordable way to build such kinds of applications. Of course there is a price to pay to fully take advantage of GPU resources. We need to reconsider conventional methods and design new innovative computational algorithms to really take advantage of the theoretical peak performance.

## 1.2. Impact on the design of numerical methods

As a simple statement of fact, standard sophisticated discretization methods for partial differential problems or optimization algorithms are not really suited for high-performance parallelization on manycore GPU-like architectures. For example, implicit methods lead to a large (sparse) linear system which is solved either by a direct method involving a sparse factorization and a sequential descent/roll up, or by an iterative algorithm which is also sequential by nature. Of course, one can find BLAS-like libraries on many-cores (like cuBLAS on nVIDIA GPU), but today reported speedups are partially satisfactory. For that reason, explicit methods are certainly much more suitable on GPU architectures. Another aspect is the memory access to neighboring degrees of freedom, which is a common issue for PDE-based problems. It is important to notice that there are strongly optimized data structures and methods for fixed neighbor stencil patterns access. This makes cartesian structured grids very suitable candidates. For complex geometries, one can imagine immersed boundary methods (IBM) into cartesian uniform meshes.

Our belief at CMLA is that we have to reconsider both models and methods, in order to derive efficient single-program multiple-data (SPMD) algorithms with communications rates that do not affect the floating point performance too much. For most of the classical PDEs, there are tracks to achieve manycore-suited computational approaches. For example, Lattice Boltzmann (LB) methods are a particular class of cellular automata able to discretize classical equations like the heat equations, convection-diffusion equations and even the unsteady Navier-Stokes equations or more complicated coupled systems. Because of the explicit nature of the method and the uniform local spatial pattern/stencil of discretization, the LB methods are excellent SPMD computational approaches. In the next section, we will focus and give more details on LB methods for solving the incompressible Navier-Stokes equations.

Another track is the underlying microscopic dynamics behind macroscopic models. Generally PDEs are nothing else but a deterministic macroscopic representation of some microscopic or mesoscopic dynamics with “uncertainty” taken into account (stochastic effects like brownian motion). The Laplace operator for example is the macroscopic diffusion operator of the microscopic brownian random walk. Generally, at the microscopic scale, there is an underlying transport process and a collision/interaction

phenomenon. On the other hand, the possible numerical simulation at microscopic scale will require a large number of “individuals” in order to derive accurate statistics able to return the macroscopic effects accurately. Manycore architectures are excellent hardware candidates to run population-based computational approaches (like particle methods) on a very large number of individuals. In the sequel we shall give some illustrative examples.

### 1.2.1. Incompressible flows

Consider the two-dimensional incompressible Navier-Stokes equations defined on a bounded spatial domain  $\Omega$  of  $\mathbb{R}^2$  :

$$\nabla \cdot \mathbf{u} = 0, \quad t > 0, \quad x \in \Omega, \quad (1.1)$$

$$\rho(\partial_t \mathbf{u} + \mathbf{u} \cdot \nabla \mathbf{u}) + \nabla p - \rho \nu \Delta \mathbf{u} = 0, \quad t > 0, \quad x \in \Omega, \quad (1.2)$$

where  $\rho$  is the (constant) density,  $\mathbf{u}$  the fluid velocity,  $p$  the pressure and  $\nu > 0$  the kinematic

viscosity of the fluid. There is lot of literature on how to discretize this system of equations. Because of the implicit incompressibility constraint  $\nabla \cdot \mathbf{u} = 0$ , generally an implicit solver is used leading to the solution of a large linear system to solve at each time step, what is not very directly suitable for GPU.

Since two decades, a new family of discrete explicit methods, namely the Lattice Boltzmann methods (LBM) (see the book [8] for example or the review paper [7]) are a kind of kinetic-based cellular automata. They are based on a discretization of the Boltzmann equation

$$\partial_t f + \mathbf{e} \cdot \nabla_x f = (\partial_t f)_{coll}$$

governing the distribution  $f = f(\mathbf{x}, \mathbf{e}, t)$  of gas particles having speed  $\mathbf{e}$  at position  $\mathbf{x}$  and time  $t$ . The term  $(\partial_t f)_{coll}$  models all the possible pairwise particle collisions. It is expected that the system fulfils the so-called H-theorem, stating that the entropy functional  $H(f) = \int f \log f d\mathbf{e}$  decreases in time. The equilibrium steady state returns the well-known Maxwellian distribution (see for example [3] for a general theory).

Lattice Boltzmann methods have the advantage to be implemented very easily and even to deal with complex geometries using an immersed boundary approach while being potentially very accurate. Let us consider the

## SOME EXAMPLES OF CFD COMPUTATIONS ON GPU

2D Navier-Stokes case : the basic LB method is the so-called Lattice BGK (LBGK) method that uses a BGK collision operator

$$(\partial_t f)_{coll} = \frac{f^{eq} - f}{\tau}$$

for an equilibrium distribution  $f^{eq}$  and a characteristic collision time scale  $\tau > 0$ . The discretization process first deals with a finite set of discrete velocities  $\{\mathbf{e}_i\}_{i=1,\dots,N}$ ,  $N > 1$ . This leads to a coupled system of spatial transport equations

$$\partial_t f_i + \mathbf{e}_i \cdot \nabla_x f_i = \frac{f_i^{eq} - f_i}{\tau}, \quad i = 1, \dots, N$$

with  $f_i(\mathbf{x}, t) \approx f(\mathbf{x}, \mathbf{e}_i, t)$ . The standard D2Q9 lattice makes use of a uniform spatial grid with constant space step per direction  $\Delta x = 1$ , and (only)  $N = 9$  discrete microscopic velocities  $\{\mathbf{e}_i\}_{i=0,\dots,8}$  (with  $\mathbf{e}_0 = 0$ ) as shown in figure 1. Then we have nine discrete transport-collision equations to solve. Using the characteristic method for integrating transport term and a first order explicit Euler discretization for the collision term, we get

$$f_i(\mathbf{x} + \mathbf{e}_i \Delta t, t + \Delta t) - f_i(\mathbf{x}, t) = \frac{\Delta t}{\tau} [f_i^{eq}(\mathbf{x}, t) - f_i(\mathbf{x}, t)], \quad i = 0, \dots, 8, \quad (1.3)$$

with  $\tau > 0$  the characteristic collision time, for each lattice point  $\mathbf{x}$ . The zeroth-order and first-order moments allow us to retrieve both density and momentum. Denoting by  $\mathcal{S} = \{0, \dots, 8\}$ , we have

$$\sum_{i \in \mathcal{S}} (1, \mathbf{e}_i) f_i(\mathbf{x}, t) = (\rho, \rho \mathbf{u})(\mathbf{x}, t). \quad (1.4)$$

From a formal Chapman-Enskog expansion of the discrete density probability functions  $f_i$ ,

$$f_i = f_i^{(0)} + \varepsilon f_i^{(1)} + \varepsilon^2 f_i^{(2)} + \dots$$

where  $\varepsilon$  is a lattice Knudsen number, for  $\varepsilon \ll 1$  it is possible retrieve the macroscopic Fluid Mechanics equations. In order to reproduce the Navier-Stokes equations, only the first two approximations  $f_i^{(0)}$  and  $f_i^{(1)}$  are required [7]. The zero-th order term  $f_i^{(0)}$  identifies with the equilibrium distribution  $f_i^{eq}$ . Let us consider a dimensionless lattice size  $\Delta x = 1$  and a lattice speed  $c = 1$  ( $\Delta t = 1$ ). By choosing the equilibrium density function

$$f_i^{eq} = w_i \rho \left( 1 + 3\mathbf{e}_i \cdot \mathbf{u} + \frac{9}{2}(\mathbf{e}_i \cdot \mathbf{u})^2 - \frac{3}{2}|\mathbf{u}|^2 \right), \quad (1.5)$$

with weighting factors  $w_0 = 4/9$ ,  $w_k = 1/9$  for  $k = 1, \dots, 4$  and  $w_k = 1/36$  for  $k = 5, \dots, 8$ , then for  $|u| \ll 1$ , we get (up to a scaling) second order accurate approximations of the incompressible Navier-Stokes equations with a kinematic viscosity  $\nu$  equal to

$$\nu = \frac{1}{3} \left( \tau - \frac{1}{2} \right) \quad (1.6)$$

Actually, for a given fluid kinematic viscosity  $\nu$ , we compute  $\tau > \frac{1}{2}$  such that (1.6) holds. It can be shown that the LBGK method is linearly stable as soon as  $\tau > 1/2$ . Practically, it becomes unstable for  $\tau$  close to  $\frac{1}{2}$  (high Reynolds number), but stabilization methods exist in this case (MRT, entropy fix, positivity preserving, etc, see [2]).

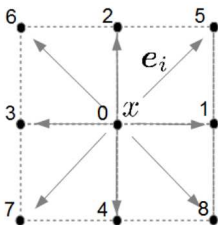


FIGURE 1. The two-dimensional D2Q9 lattice pattern

**LBM code porting on GPU.** It is easy to check that LBM can be rewritten as a two-step fractional step method, with i) a collisionless transport evolution, ii) a pure collision process. The GPU collision step can be perfectly done in parallel because of only pointwise operations. The transport step requires communications with the direct first neighboring lattice points. But, because this communication pattern is uniform over the whole computational domain, there are potentially important ways of improvement and performance gain of memory access. On NVIDIA GPU boards, using CUDA programming, one can use `texture` memory (both structures and access methods) that are optimized for fixed memory patterns.

We implemented the D2Q9 LBGK scheme with a stabilization method proposed by Brownlee et al. in [2]. We used Pixel Buffer Object (PBO) for OpenGL instant visualization and binding between CUDA structures and PBO. On a lattice grid of typical size  $1024 \times 1024$ , we observe speedup

## SOME EXAMPLES OF CFD COMPUTATIONS ON GPU

factors of about 100 compared to a single-thread CPU sequential computation, allowing for interactivity, visual appearance and evolution of von Karman vortex sheddings. Flow interaction is made possible by adding new obstacles during computation with the mouse (see figure 2). This is easily handled by the GPU programming using a solid mask array.

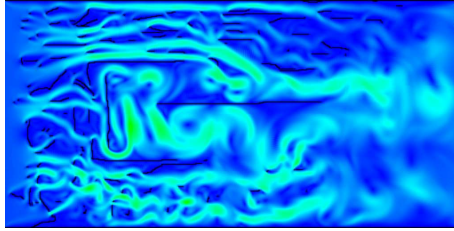


FIGURE 2. Instant Lattice Boltzmann GPU computation of the Navier-Stokes equations on a cartesian grid of typical size  $1024 \times 512$ . Flow interaction is made possible by adding new obstacles during computation with the mouse.

### 1.2.2. Compressible flows

Let us now consider a compressible fluid. The Euler equations govern the dynamics of a perfect fluid. The mass, momentum and total energy conservation equations read

$$\partial_t \rho + \nabla \cdot (\rho \mathbf{u}) = 0, \quad (1.7)$$

$$\partial_t (\rho \mathbf{u}) + \nabla \cdot (\rho \mathbf{u} \otimes \mathbf{u}) + \nabla p = 0, \quad (1.8)$$

$$\partial_t (\rho E) + \nabla \cdot ((\rho E + p) \mathbf{u}) = 0 \quad (1.9)$$

with density  $\rho$ , velocity vector  $\mathbf{u}$ , pressure  $p$  and specific total energy  $E$ . The energy  $E$  is the sum of the kinetic energy  $|\mathbf{u}|^2/2$  and the internal energy  $e$ . For the perfect gas with constant specific heat ratio  $\gamma$ ,  $\gamma \in (1, 3]$ , we have

$$E = e + \frac{1}{2} |\mathbf{u}|^2, \quad e = \frac{1}{\gamma - 1} \frac{p}{\rho}. \quad (1.10)$$

The above system can be written in condensed vector form

$$\partial_t U + \nabla \cdot \mathbf{F}(U) = 0, \quad U = (\rho, \rho \mathbf{u}, E)^T. \quad (1.11)$$

This system is known to be hyperbolic on its admissible state space ([4]).



For discretization, we consider a conservative finite volume scheme built on an unstructured finite volume mesh made of cells  $K$ . We will denote  $A_{KL}$  the edge separating the two volumes  $K$  and  $L$  and  $\nu_{KL}$  the unit exterior vector orthogonal to  $A_{KL}$ . A general explicit first-order finite volume scheme reads

$$U_K^{k+1} = U_K^n - \frac{\Delta t}{|K|} \sum_{L \in \mathcal{V}(K)} |A_{KL}| \Phi(U_K^n, U_L^n, \nu_{KL}), \quad (1.12)$$

with a numerical flux  $\Phi(U_K^n, U_L^n, \nu_{KL})$  having at least Lipschitz-continuous regularity, and being consistent i.e.  $\Phi(U, U, \nu) = \mathbf{F}(U)\nu$ . For stability purposes, numerical fluxes are generally built to fulfil an upwinding property. In this context, two main families of upwind flux are identified (see [6]) : the Flux Difference Splitting (FDS) one, and the Flux Vector Splitting (FVS) one. FDS fluxes (including Godunov, Osher, Roe, HLLC, etc...) are written in the form

$$\begin{aligned} \Phi(U_K^n, U_L^n, \nu_{KL}) = & \quad (1.13) \\ & \frac{1}{2} (F(U_K^n, \nu_{KL}) + F(U_L^n, \nu_{KL})) - \frac{1}{2} \int_{\Gamma_{KL}^n} |A(U(s), \nu_{KL})| U'(s) ds \end{aligned}$$

where  $A(U, \nu)$  denote the (diagonalizable) Jacobian matrix of the flux in the direction  $\nu$ , and  $\Gamma_{KL}^n = \Gamma(U_K^n, U_L^n, s)$  is a Lipschitz continuous path linking the states  $U_K^n$  and  $U_L^n$  with a curvilinear parameter  $s \in [0, 1]$ . The second term of the RHS of (1.13) corresponds to the upwind artificial viscosity term.

From the GPU computational point of view, FDS schemes require at each time step i) the computation of the FDS flux with memory reads of the cell states ; ii) cell updates with memory reads of the FDS fluxes (see figure 3). Memory transfer may be a limiting performance factor for GPUs if the DRAM bandwidth is saturated.

The Flux Vector Splitting (FVS) family [6] has numerical fluxes written in the form

$$\Phi(U_K^n, U_L^n, \nu_{KL}) = F^+(U_K^n, \nu_{KL}) + F^-(U_L^n, \nu_{KL}) \quad (1.14)$$

with  $F^+$  representing the leftward part of the flux and  $F^-$  representing the rightward part. Consistency requirements involve the identity

$$F^+(U, \nu) + F^-(U, \nu) = \mathbf{F}(U)\nu.$$

## SOME EXAMPLES OF CFD COMPUTATIONS ON GPU

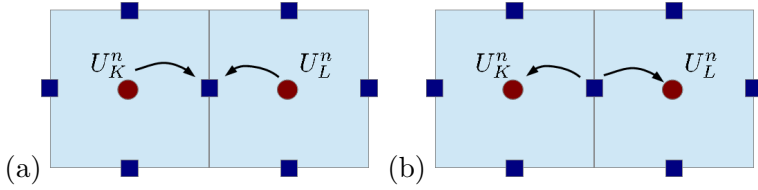


FIGURE 3. FV scheme with Flux Difference Splitting FDS schemes. FDS require two memory transfers : (a) computation of the numerical flux with memory reads of states ; (b) state update into control volumes with memory reads of numerical fluxes.

What is peculiar with FVS is that  $F^+(U_K^n, \nu_{KL})$  can be computed without any knowledge of the neighboring state  $U_L^n$ , and conversely. Then the GPU computation of  $F^+(U_K^n, \nu_{KL})$  may be seen as a pointwise, cell-centered computation, perfectly done in parallel. For that reason, one has only to send  $F^+$  and  $F^-$  to the neighboring cells for state updates, thus reducing memory reads and DRAM transfer (see figure 4). Moreover, FVS fluxes

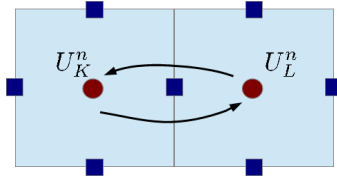


FIGURE 4. FV scheme with Flux Vector Splitting (FVS) fluxes. FVS only require one memory transfers :  $F^+$  ou  $F^-$  are sended to the neighboring cells for state update.

generally do not require neither eigenstructure decomposition nor matrix-vector products, what improves the whole performance. For example, the van Leer's FVS with Hänel-Schwane energy-flux modification [5] leads to

the scripts (written here for  $\nu = (1, 0)$ ) :

$$F_\rho^+ = \frac{\rho c}{4}(M + 1)^2 1_{(|M| \leq 1)} + \max(u, 0) 1_{(|M| > 1)}, \quad (1.15)$$

$$p^+ = \frac{p}{4}(M + 1)^2(2 - M) 1_{(|M| \leq 1)} + p 1_{(M > 1)}, \quad (1.16)$$

$$F_{\rho u}^+ = u F_\rho^+ + p^+, \quad F_{\rho v}^+ = v F_\rho^+, \quad (1.17)$$

$$F_{\rho E}^+ = (E + p/\rho) F_\rho^+ \quad (1.18)$$

where  $c = \sqrt{\gamma p/\rho}$  is the speed of sound and  $M = u/c$  is the normal Mach number. For these reasons, FVS are clearly better candidates for GPU implementation and high-performance computation [9]. On figure 5, we show an instant computation of the well-known 2D Mach 3 forwarding step case on a very fine grid, using (1.15)-(1.18) as FVS scheme. Again, mouse interactivity allows us to add obstacles on-the-fly and observe the fluid response. This is of great interest for training, education and comprehension of Fluid Dynamics.

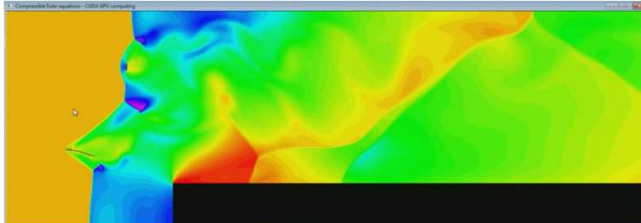


FIGURE 5. Instant GPU computation on the well-known Mach 3 forward step channel 2D problem. Hänel's FVS is here used. The flow can be perturbed by directly adding new wall obstacles with the mouse pointer.

### 1.2.3. Three-dimensional free transport kinetic equations

This case was designed to evaluate GPU performance of particle methods. Let us consider the following homogeneous Vlasov equation in 3D : find  $f = f(\mathbf{x}, \mathbf{v}, t)$ ,  $\mathbf{x} \in \Omega(t) \subset \mathbb{R}^3$ ,  $\mathbf{v} \in \mathbb{R}^3$ ,  $t > 0$ , solution of

$$\partial_t f + \mathbf{v} \cdot \nabla_x f + a(\mathbf{x}) \nabla_v f = 0, \quad \mathbf{x} \in \Omega(t) \subset \mathbb{R}^3, \quad \mathbf{v} \in \mathbb{R}^3, \quad t > 0 \quad (1.19)$$

## SOME EXAMPLES OF CFD COMPUTATIONS ON GPU

with  $f(.,.,0) \in L^1(\Omega \times \mathbb{R}^3)$ . Standard eulerian discretization methods would involve a mesh in a space of dimension 6, what is still not realistic to address at the present time. An alternative approach is to reformulate the transport dynamics behind this equation. Let us consider the following system of motion equations defined on a large set of particles  $\{\mathbf{x}_k\}_k$  :

$$\begin{cases} \dot{\mathbf{x}}_k(t) = \mathbf{v}_k, \\ \dot{\mathbf{v}}_k(t) = \mathbf{a}(\mathbf{x}_k) \end{cases} \quad (1.20)$$

and the discrete measure-valued distribution

$$f = \sum_{k=1}^N \omega_k \delta(\mathbf{x} - \mathbf{x}_k(t)) \delta(\mathbf{v} - \mathbf{v}_k(t)) \quad (1.21)$$

for some given weighting factors  $\{\omega_k\}_k$ ,  $\omega_k > 0$ . We want to evaluate the  $L^1$  norm on  $f$  in the time-dependent domain  $\Omega(t)$  (a pulsating sphere for example, see [10]). For that we have to compute at each time step

$$\|f(t)\|_{L^1(\Omega(t))} = \sum_{k=1}^N \omega_k 1_{(x_k \in \Omega(t))}. \quad (1.22)$$

The summation has to be “optimized” on a GPU architecture. For that we used the sum-reduction algorithm proposed in the GPU code examples of the CUDA SDK. We have observed parallel speedup factors of about 25-30 for particles sets between 1 million and 8 million particles on a GPGPU TESLA C2070.

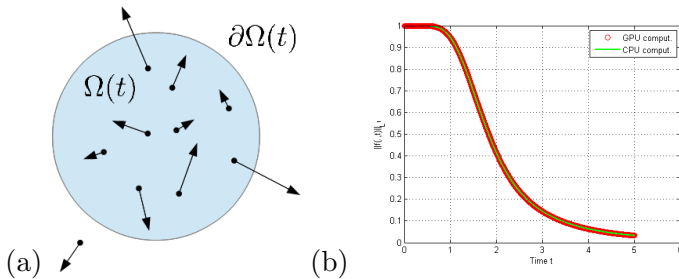


FIGURE 6. Validation of GPU acceleration of free transport equations on a moving domain with parallel reduction at each time step for  $L^1$ -norm computation. (a) Schematic of the particles and moving domain ; (b) Discrete  $L^1(\Omega(t))$ -norm of the distribution during time  $t$ .

### 1.3. Impact on the data structures

In the above sections, we have seen how GPU computing may change the way of thinking both physical models and computational methods. Beyond pure computational aspects, there is of course the programming and code optimization dimensions. The obtained “speedup” and the necessary work of specific GPU programming are subject to a search of thrust performance. Parallel computing strongly alters the balance of power in the classical duality memory-computation. A first simple rule is to try to keep as much as possible data on the parallel “device” and transfer data as less as possible because of limited bandwidth. To illustrate, it is even often cheaper to recompute a result than transmitting it.

Communication performance is also strongly linked to the way to handle complex data. Data coalescence is a critical keypoint for optimal performance. To offset a large latency of global memory a rational way is to read consecutive blocks in memory (coalescing). There are two kinds of

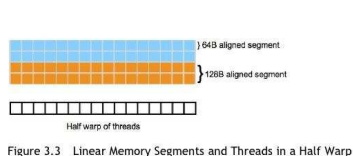


Figure 3.3 Linear Memory Segments and Threads in a Half Warp

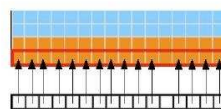


Figure 3.4 Coalesced Access-All Threads but One Access the Corresponding Word in a Segment

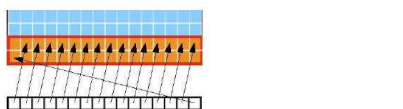


Figure 3.5 Unaligned Sequential Addresses that Fit within a Single 128-Byte Segment

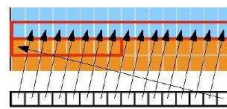


Figure 3.6 Misaligned Sequential Addresses that Fall within Two 128-Byte Segments

FIGURE 7. Schematic of data coalescence, extracted from the “CUDA C Best Practices Guide”, NVIDIA 2012 [1]

non-coalescing memory accesses, described in the figure 7 :

- either the *threads* cannot access to neighboring fields in the right order ;
- or there is an alignment problem, the first *thread* of a *warp* must access one memory multiple of 32, 64 or 128 (depending on the data), see [1].

It is difficult when looking to optimized performance, to work with very sophisticated data models. The notion of array perfectly fits to this

## SOME EXAMPLES OF CFD COMPUTATIONS ON GPU

scheme, data types more developed as structures or classes instead readily scatter in the data. For example it is much more efficient to manipulate Structures of Arrays (SoA) than Arrays of Structures (AoS). To be efficient, it is necessary to stay close to the data and always keep in mind the specific hardware architecture of GPU and the constraints it imposes to the data. This is probably the main reason why we think that today CUDA is one of the most appropriate language to obtain optimal performance on GPU.

### 1.4. General experience feedback and concluding remarks

Let us conclude by some humble advices. Today's GPU for scientific computing is a question a tradeoff between performance and design and/or implementation effort. Reasonable (suboptimal) speedup (say 10 or 20) is easy to reach. GPU parallel programming is far easier for computational methods feined on cartesian grids or meshfree particle methods. For stronger performance, the way is to find a good tradeoff between implementation effort, code readability and performance. GPU computing requires a real reflection on the choice of data structures, especially for the sake of memory coalescence : arrays of structures AoS versus structures of arrays SoA, byte alignment, etc. In the same spirit, the strategy/use of intermediate cache memory level (per streaming multiprocessor for example) is also of great importance for high performance. Texture memory is particularly suited for local partial differential operators involving uniform spatial stencils.

Our belief at CMLA is that GPU/manycore processors will deeply impact the numerical solvers in the next years. We have to think about paradigm shifts for both modeling and discretization for strong better GPU performance.

## 2. Videos of instant GPU computations on youtube

All the interactive computations can be found at the following URL :  
<http://www.youtube.com/user/floriandevuyst/videos>.

## Acknowledgements

This work has been partly supported by the FARMAN Institute of ENS Cachan and by a NVIDIA Equipment Grant in 2011.

## Bibliographie

- [1] *CUDA C Best Practices Guide 4.1*, NVIDIA. 2012.
- [2] R.A. Brownlee, A.N. Gorban, and J. Levesley. Stabilization of the lattice boltzmann method using the Ehrenfests' coarse-graining data. *Physical Review E*, 74 :037703, 2006.
- [3] Carlo Cercignani. *The Boltzmann Equation and Its Applications*, volume 67. Springer-Verlag, 1988.
- [4] E. Godlewski and P.-A. Raviart. *Numerical approximation of hyperbolic conservation laws*, volume 118. Applied Mathematical Sciences, Springer-Verlag, Boston, 1996.
- [5] D. Hanel and R. Schwane. An implicit flux-vector splitting scheme for the computation of viscous hypersonic flow. *AIAA Paper*, 25, 1989. Paper 89-0274.
- [6] A. Harten, P.D. Lax, and B. van Leer. On upstream differencing and godunov-type schemes for hyperbolic conservation laws. *SIAM Review*, 25 :35–61, 1983.
- [7] RR. Nourgaliev, T.N. Dinh, T.G. Theofanous, and D. Joseph. The lattice boltzmann equation method : theoretical interpretation, numerics and implications. *Int. J. of Multiphase Flow*, 29 :117–169, 2003.
- [8] Sauro Succi. *The Lattice Boltzmann Equation for Fluid Dynamics and Beyond*. Oxford, 2001. ISBN :0-19-850398-9.
- [9] F. De Vuyst. A flux vector splitting method that preserves stationary contact discontinuities. *Acta Mathematicae Applicandae*, 2012. submitted.
- [10] F. De Vuyst and F. Salvarani. Gpu-accelerated numerical simulations of the knudsen gas on time-dependent domains. *Computer Physics Communications*, 2012. in press.

## SOME EXAMPLES OF CFD COMPUTATIONS ON GPU

FLORIAN DE VUYST  
Centre de Mathématiques  
et de leurs Applications  
CMLA CNRS UMR 8536  
61, avenue du Président Wilson  
94235 Cachan CEDEX  
FRANCE  
devuyst@cmla.ens-cachan.fr

CHRISTOPHE LABOURDETTE  
Centre de Mathématiques  
et de leurs Applications  
CMLA CNRS UMR 8536  
61, avenue du Président Wilson  
94235 Cachan CEDEX  
FRANCE  
labour@cmla.ens-cachan.fr