



HAL
open science

Multi-GPU Implementation of a Hybrid Thermal Lattice Boltzmann Solver using the TheLMA Framework

C. Obrecht, F. Kuznik, Bernard Tourancheau, J.-J. Roux

► **To cite this version:**

C. Obrecht, F. Kuznik, Bernard Tourancheau, J.-J. Roux. Multi-GPU Implementation of a Hybrid Thermal Lattice Boltzmann Solver using the TheLMA Framework. *Computers and Fluids*, 2013, 80, pp.269-275. 10.1016/j.compfluid.2012.02.014 . hal-00731152

HAL Id: hal-00731152

<https://hal.science/hal-00731152>

Submitted on 9 Jun 2014

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.

Multi-GPU Implementation of a Hybrid Thermal Lattice Boltzmann Solver using the TheLMA Framework

Christian Obrecht^{a,b,c,*}, Frédéric Kuznik^{b,c}, Bernard Tourancheau^{b,d,e},
Jean-Jacques Roux^{b,c}

^aEDF R&D, Département EnerBAT, F-77818 Moret-sur-Loing Cedex, France

^bUniversité de Lyon, F-69361 Lyon Cedex 07, France

^cINSA-Lyon, CETHIL, UMR5008, F-69621 Villeurbanne Cedex, France

^dINSA-Lyon, CITI, INRIA, F-69621 Villeurbanne Cedex, France

^eUniversité Lyon 1, F-69622 Villeurbanne Cedex, France

Abstract

In this contribution, a single-node multi-GPU thermal lattice Boltzmann solver is presented. The program is based on the TheLMA framework which was developed for the purpose. The chosen implementation and optimisation strategies are described, both for inter-GPU communication and for coupling with the thermal component of the model. Validation and performance results are provided as well.

Key words: Thermal lattice Boltzmann method, GPU computing, CUDA

1. Introduction

Since its introduction in the late eighties, the lattice Boltzmann method (LBM) has proven to be an effective approach in computational fluid dynamics (CFD). It has been successfully applied to a wide range of engineering issues such as multiphase flows, porous media, or free surface flows. Despite of these achievements, the use of the LBM for thermal flow simulation is not very widespread yet. A possible reason for this situation is the relatively high computational cost of most thermal LBM models.

The use of emerging many-core architectures such as graphics processing units (GPUs) in CFD is fairly promising [2]. Being a regular data-parallel algorithm, the LBM is especially well adapted to such hardware. Nevertheless, implementing the lattice Boltzmann method for the GPU is still a pioneering task. Several important issues, such as multi-physics applications and efficient multi-GPU implementations, remain to be addressed. The present work, presenting a multi-GPU thermal LBM solver, faces both challenges.

The remaining of the paper is organised as follows. In the first section, we briefly present the thermal lattice Boltzmann model we chose to implement.

*christian.obrecht@insa-lyon.fr

Next, we give an overview of the TheLMA framework on which our program is based. In the third section, we describe our implementation and our optimisation strategies. Last, we provide some simulation results for validation purpose and discuss performance issues.

2. Hybrid thermal lattice Boltzmann model

The lattice Boltzmann equation (LBE), i.e. the governing equation of the LBM is interpreted as a discrete version of the Boltzmann equation [6]. In the LBM, as for the Boltzmann equation, a fluid is represented through the distribution of a single particle in phase space (i.e. position \mathbf{x} and particle velocity $\boldsymbol{\xi}$). Space is commonly discretised using a uniform orthogonal lattice and time using constant time steps. A finite set of particle velocities $\boldsymbol{\xi}_\alpha$ is substituted to velocity space. The LBM counterpart of the distribution function $f(\mathbf{x}, \boldsymbol{\xi}, t)$ is a finite set $f_\alpha(\mathbf{x}, t)$ of particle distribution functions associated to the particle velocities $\boldsymbol{\xi}_\alpha$. The LBE writes:

$$|f_\alpha(\mathbf{x} + \delta t \boldsymbol{\xi}_\alpha, t + \delta t)\rangle - |f_\alpha(\mathbf{x}, t)\rangle = \Omega[|f_\alpha(\mathbf{x}, t)\rangle]. \quad (1)$$

where Ω is the collision operator. The mass density ρ and the momentum \mathbf{j} of the fluid are given by:

$$\rho = \sum_\alpha f_\alpha, \quad \mathbf{j} = \sum_\alpha f_\alpha \boldsymbol{\xi}_\alpha. \quad (2)$$

The particle velocity set is usually chosen such as to link the nodes to some of their nearest neighbours, as the three-dimensional D3Q19 stencil illustrated by Fig. 1. It is well-known that such basic models are not energy conserving. To address this issue, several approaches such as multi-speed models [15] or double-population models [3] have been developed. In the former category, a larger set of particle velocities is defined allowing multiple particle speeds along some directions. In the later category, an additional set of energy distribution functions is used. Both approaches suffer from inherent numerical instabilities [4]. Moreover, from a computational standpoint, both methods lead to a markedly higher memory consumption.

Hybrid thermal lattice Boltzmann models constitute an alternative approach in which the flow simulation is decoupled from the solution of the heat equation. These models are free from the aforementioned drawbacks. In the present work, we implemented a simplified version of the hybrid thermal lattice Boltzmann model developed in [4]. Flow simulation is performed by multiple-relaxation-time LBM [1], using the D3Q19 stencil. In the multiple-relaxation-time approach, collision is performed in moment space and the LBE writes:

$$|f_\alpha(\mathbf{x} + \delta t \boldsymbol{\xi}_\alpha, t + \delta t)\rangle - |f_\alpha(\mathbf{x}, t)\rangle = -\mathbf{M}^{-1} \mathbf{S} \left[|m_\alpha(\mathbf{x}, t)\rangle - |m_\alpha^{(\text{eq})}(\mathbf{x}, t)\rangle \right] \quad (3)$$

where \mathbf{M} is an orthogonal matrix mapping the set of particle distributions to a set of moments m_α , and \mathbf{S} is a diagonal matrix containing the relaxation rates.

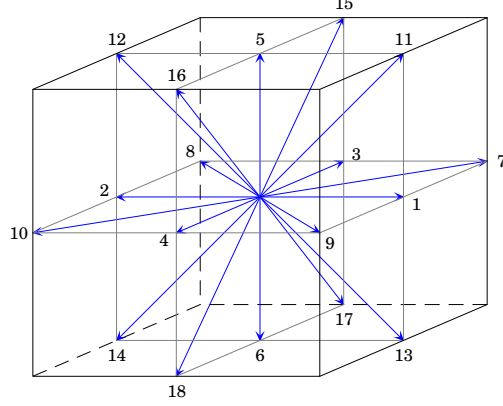


Figure 1: The D3Q19 stencil

Matrices \mathbf{M} and \mathbf{S} as well as the equilibria of the moments for the D3Q19 stencil can be found in appendix A of [1].

In our simulations, we set the ratio of specific heats $\gamma = C_P/C_V$ to $\gamma = 1$. Temperature T is obtained by solving the following finite-difference equation:

$$T(\mathbf{x}, t + \delta t) - T(\mathbf{x}, t) = \kappa \Delta^* T - \mathbf{j} \cdot \nabla^* T \quad (4)$$

where κ denotes the thermal diffusivity and the finite-difference operators are defined as:

$$\begin{aligned} \partial_x^* f(i, j, k) &= f(i + 1, j, k) - f(i - 1, j, k) \\ &- \frac{1}{8} (f(i + 1, j + 1, k) - f(i - 1, j + 1, k) + f(i + 1, j - 1, k) - f(i - 1, j - 1, k) \\ &\quad + f(i + 1, j, k + 1) - f(i - 1, j, k + 1) + f(i + 1, j, k - 1) - f(i - 1, j, k - 1)) \end{aligned}$$

$$\begin{aligned} \Delta^* f(i, j, k) &= 2(f(i + 1, j, k) + f(i - 1, j, k) + f(i, j + 1, k) \\ &\quad + f(i, j - 1, k) + f(i, j, k + 1) + f(i, j, k - 1)) \\ &- \frac{1}{4} (f(i + 1, j + 1, k) + f(i - 1, j + 1, k) + f(i + 1, j - 1, k) \\ &\quad + f(i - 1, j - 1, k) + f(i, j + 1, k + 1) + f(i, j - 1, k + 1) \\ &\quad + f(i, j + 1, k - 1) + f(i, j - 1, k - 1) + f(i + 1, j, k + 1) \\ &\quad + f(i - 1, j, k + 1) + f(i + 1, j, k - 1) + f(i - 1, j, k - 1)) \\ &- 9f(i, j, k) \end{aligned}$$

It should be mentioned that these operators share the same symmetries as the D3Q19 stencil. The coupling of the temperature to the momentum is explicit

in Eq. 4. The coupling of the momentum to the temperature is carried out in the equilibrium of the second order moment m_2 related to internal energy:

$$m_2^{(\text{eq})} = -11\rho + 19\mathbf{j}^2 + T \quad (5)$$

3. The TheLMA framework

Since the introduction of the CUDA technology [9] by the Nvidia company in 2007, several successful attempts to implement the LBM on the GPU were reported [11, 16]. Yet, constraints induced by low-level hardware specificities make GPU programming fairly different from usual software development. Beside other limitations, it is worth mentioning the inability of the compilation tool chain to link several GPU binaries and the inlining¹ of device functions, i.e. functions called by GPU kernels. Library oriented development is therefore not relevant up to now.

To improve code reusability, we designed the TheLMA framework [10], which is outlined in Fig. 2. TheLMA stands for *Thermal LBM on Many-core Architectures*, thermal flow simulations being our main topic of interest. The framework consists in a set of C and CUDA source files. The C files provide a set of utility functions to retrieve simulation parameters, initialise computation devices, extract statistical informations, and export data in various output formats. The CUDA files are included at compile time in the `thelma.cu` file which is mainly a container, additionally providing some general-purpose macros. Implementing a new lattice Boltzmann model within the framework mostly requires to alter the `compute.cu` file.

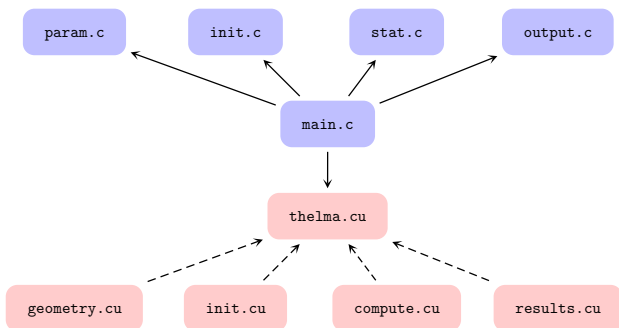


Figure 2: Overall structure of the TheLMA framework

¹Starting with hardware of compute capability 2.0, i.e. the Fermi generation, it is possible to perform actual function calls, yet inlining is still the default behaviour.

Our framework provides native single-node multi-GPU management based on POSIX threads [13, 14]. Each computing device is managed by a specific thread which is responsible for creating the appropriate CUDA context. The computation domain is split along the fastest varying direction. Communication between sub-domains is performed using zero-copy transactions on pinned exchange buffers in CPU memory. The chosen split direction allows the read and store accesses to CPU memory to be coalesced, which leads to an excellent communication and computations overlapping. Figure 3 describes the inter-GPU communication scheme. For the sake of simplicity, only one GPU associated to a single sub-domain interface is presented. Even time step data are displayed in red and odd time step data are displayed in blue. Line G stands for the GPU, lines L0 and L1 stand for the sub-domains stored in GPU memory, R0 and R1 for the buffers containing in-coming data, S0 and S1 for the buffers containing out-going data.

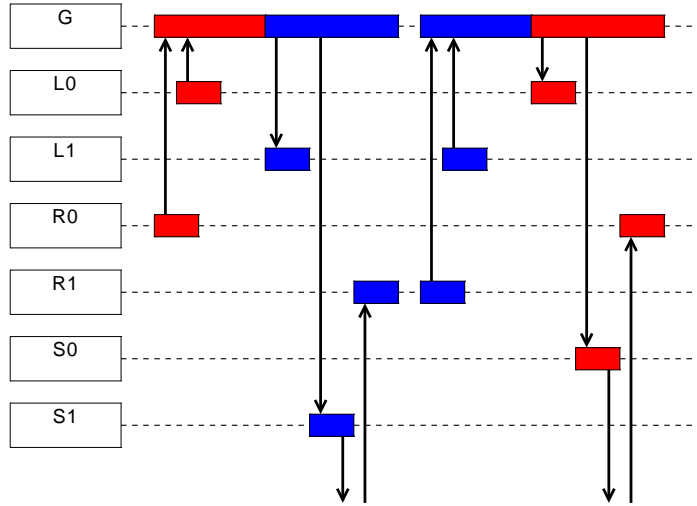


Figure 3: Inter-GPU communication scheme

4. Implementation

For the implementation of a hybrid thermal lattice Boltzmann model on the GPU, there is the alternative of using a single kernel or two distinct kernels for solving the fluid flow and the thermal part. Since Eq. 3 and Eq. 4 are tightly coupled, the two kernels option would increase the communication needs, not mentioning the overhead of kernel switching. We therefore chose to process both parts in the same kernel.

The fluid flow component is derived from the one described in [11]. Beside other optimisations, the kernel uses in-place propagation as illustrated in Fig. 4 instead of the usual out-of-place propagation. This approach allows to minimise the cost of misaligned memory transactions [12]. Misalignment may have a dramatic impact on performance with pre-Fermi hardware, since the device’s main memory is not cached.

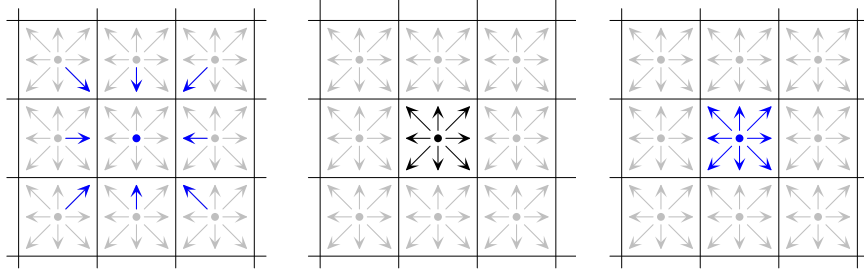


Figure 4: Out-of-place propagation

CUDA implementations of the LBM generally assign a thread to each node in order to take advantage of the massive parallelism of the architecture. This approach often leads to the use of a two-dimensional grid of one-dimensional blocks, which allows a straightforward determination of the coordinates. The grid and block dimensions are therefore identical to the size of the computation sub-domain. The direction of the blocks corresponds to the slowest varying dimension in memory in order to enable coalesced memory accesses.

In our case, these common sense optimisation principles had to be altered. Since the implemented kernel takes care of both the fluid flow part and the thermal part, the register consumption is fairly higher than for usual isothermal LBM kernels. For compute capability 1.3, we could not achieve less than 124 registers per thread. In order to avoid potential register shortage, we use small blocks containing one to four warps, i.e. 32, 64, 96 or 128 threads. Each block is associated to a one-dimensional zone spanning the cavity width, whose nodes are processed in several steps. Figure 5 outlines the chosen configuration (in two dimensions for the sake of simplicity). The blue dots represent the nodes belonging to the zone; the red frame represents the nodes being processed by the block of threads; the white background is used for the nodes whose temperature is required in the finite-difference computations.

The associated grid is two-dimensional, its size corresponding to the remaining dimensions of the sub-domain. It is worth mentioning that we assign the first rather than the second field of the `blockIdx` structure to the fastest varying dimension in memory. This option appears to improve the overlapping of computation and inter-GPU communication.

When implementing stencil computations on the GPU, reducing read redundancy is a key optimisation target [7]. We therefore chose store the temperature of the neighbouring nodes in shared memory. In the case of boundary nodes,

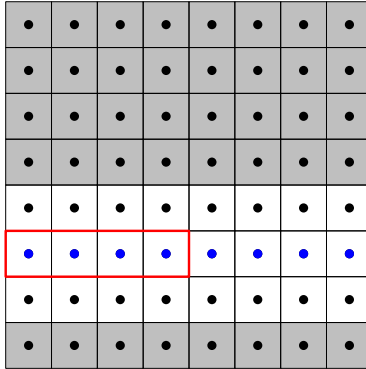


Figure 5: Block configuration

the surplus cells in the temperature array may be used to store shadow values determined by extrapolation. During the read phase, each thread is responsible for gathering the temperatures of all the nodes sharing the same abscissa, as outlined in Fig. 6.

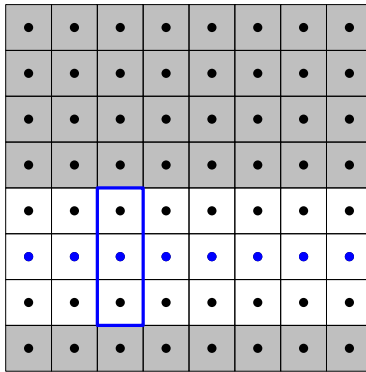


Figure 6: Read access pattern for temperature

Not taking the boundaries into account, the chosen approach reduces the read redundancy in the D3Q19 case from 19 to at most² 9.3125. Moreover, it should be noted that this data access pattern induces no misalignment at all.

²Five additional temperatures for both the first and the last node are read, thus the worst case is for blocks of size 32 and the read redundancy equals $(32 \times 9 + 2 \times 5)/32$.

5. Results and discussion

5.1. Test case

To test our code, we simulated the well-known differentially heated cubic cavity illustrated in Fig.7. In this test case, two vertical opposite walls have imposed temperatures $\pm T_0$ whereas the four remaining walls are adiabatic. The buoyancy force \mathbf{F} is computed using the Boussinesq approximation:

$$\mathbf{F} = -\rho\beta T\mathbf{g} \quad (6)$$

where β is the thermal expansion coefficient, and \mathbf{g} the gravity vector of magnitude g . The Prandtl number is set to $\text{Pr} = 0.71$. The parameters for the simulations are the Rayleigh number Ra and the kinematic viscosity ν , which determine the thermal diffusivity κ and the value of βg .

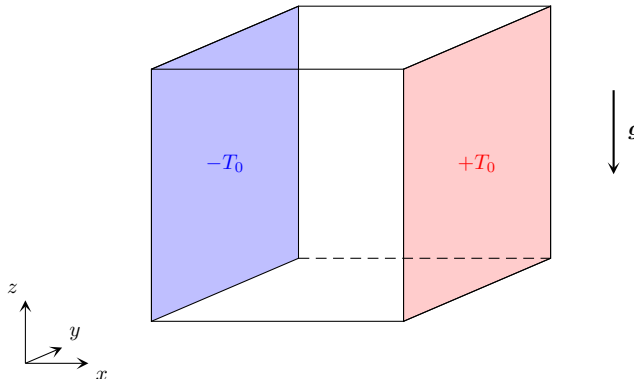


Figure 7: Differentially cubic heated cavity

We ran our program on a Tyan B7015 server fitted with eight Tesla C1060 computing devices. We could therefore perform computations on cavities as large as 512^3 in single precision.

5.2. Simulations

For validation purpose, we performed several simulations of the differentially heated cubic cavity using a 448^3 lattice in single precision. The kinematic viscosity was set to $\nu = 0.05$ and the Rayleigh number ranged from 10^4 to 10^7 . The computations were carried out until convergence to steadiness, which is assumed to be reached when:

$$\max_{\mathbf{x}} |T(\mathbf{x}, t_{n+1}) - T(\mathbf{x}, t_n)| < 10^{-5} \quad (7)$$

where $t_n = n \times 500\delta t$. The obtained Nusselt numbers at the isothermal walls are in good agreement with previously published results [17] as shown in Tab. 1.

Rayleigh number	10^4	10^5	10^6	10^7
Nusselt number	2.050	4.335	8.645	16.432
Time steps	485,000	380,000	266,000	182,000
Computation time (min)	394	309	216	148
Tric <i>et al.</i> [17]	2.054	4.337	8.640	16.342
Relative deviation	0.19%	0.07%	0.06%	0.55%

Table 1: Comparison of Nusselt numbers at the isothermal walls

Using lattices of size 512^3 allowed us to run simulation for Rayleigh numbers up to 10^9 without facing numerical instabilities. From a phenomenological standpoint, although unsteady, the flow rapidly leads to a rather stable vertical stratification. We furthermore observe quasi-symmetric and quasi-periodic flow patterns near the bottom edge of the cold wall and the top edge of the hot wall. Figure 8 shows the temperature field in the symmetry plane after 10^6 iterations. Further investigations on this simulation are required and will be published in a future contribution.

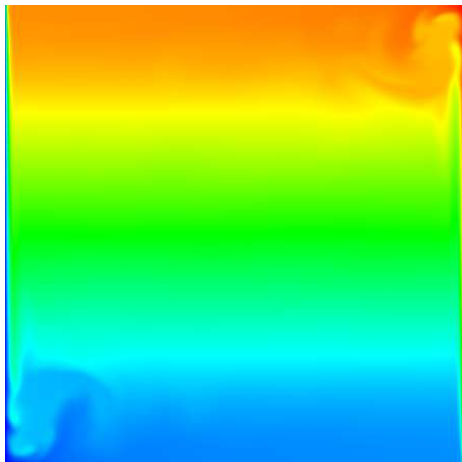


Figure 8: Symmetry plane temperature field at $Ra = 10^9$

5.3. Performance

We recorded performance results of our solver for increasing block size and cavity size (see Fig. 9). The chosen performance metric is the million lattice node updates per second (MLUPS). The cavity size has to be a multiple of the block

size, hence several configurations are not available. Performance obtained with a given block size appears to be correlated to the corresponding occupancy. For compute capability 1.3, global memory is split in eight 256 bytes wide memory banks [8]. Hence, the poor performance obtained for cavity size 256 and 512 is probably caused by partition camping, since the stride between corresponding nodes in distinct blocks is necessarily a multiple of the cavity size.

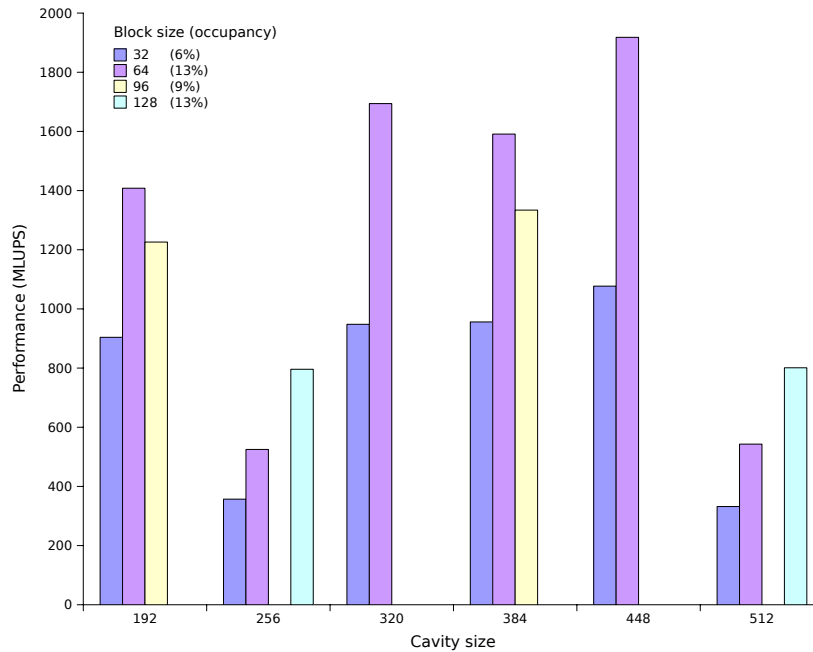


Figure 9: Performance for increasing block size and cavity size

The maximum performance is 1,920 MLUPS, achieved for cavity size 448 and block size 64. The corresponding GPU to device memory data throughput is 46.1 GB/s per GPU, which is about 62.3% of the maximum sustained throughput.³ The multiprocessor occupancy, which is only 13%, appears to be the limiting factor since it is lower than the minimum required to properly hide the global memory latency, i.e. 18.75% for compute capability 1.3.

To evaluate scalability, we also ran our program on a 192^3 lattice using from one to eight GPUs (see Fig. 10). Parallelisation efficiency is very satisfactory with no less than 84% for a fairly small computation domain. As for our isothermal multi-GPU LBM solver [14], our implementation allows excellent overlapping of communication and computations. Moreover, the amount of

³Using the `bandwidthTest` program of the CUDA SDK, we estimate the GPU to device memory maximum sustained throughput to 73.3 GB/s for the Tesla C1060.

data to exchange does not exceed the capacity of the PCI-E links.

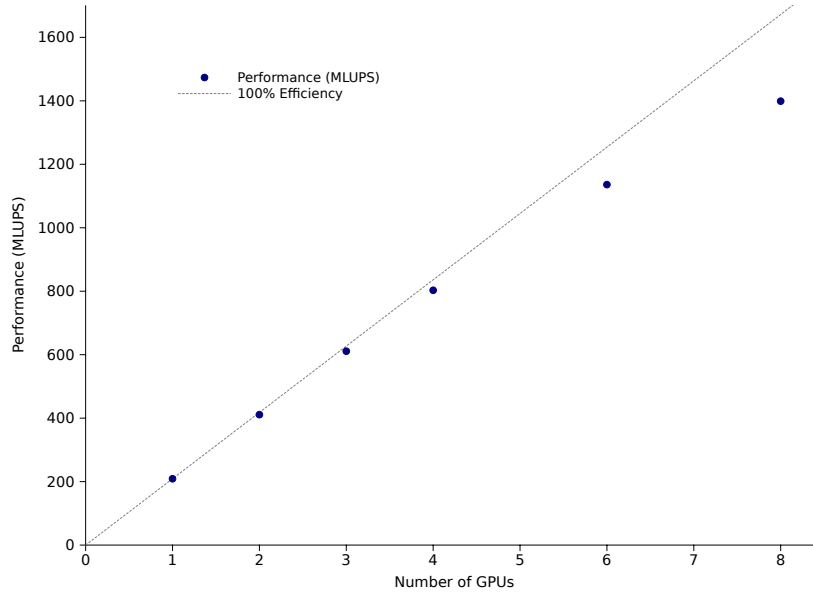


Figure 10: Parallelisation efficiency on a 192^3 cavity

6. Conclusion

In this contribution, we present a multi-GPU implementation of a thermal LBM solver, which to the best of our knowledge was never reported before. Using appropriate hardware, our program is able to run up to eight GPUs in parallel. With the latest generation of Nvidia computing devices, it is therefore possible to perform simulations on lattices containing as much as 3.2×10^8 nodes.

Validation studies have been carried out, showing both the accuracy and the stability of the chosen thermal LBM model and the correctness of our implementation. Although slightly less efficient than the isothermal version of our solver, our program provides unrivaled performance compared to CPU implementations. Recent studies [5] have shown that optimised multi-threaded CPU implementations of isothermal LBM solver running on up-to-date hardware achieve at most 85 MLUPS, which is $22 \times$ less than our maximum performance.

We furthermore study the performance bottlenecks, showing that the limiting factor is the low occupancy. Since the multiprocessor occupancy is bound by the amount of available registers there is little room for improvements using the same hardware. Yet, a more elaborate memory access pattern could avoid the partition camping effects we could observe in some cases.

We believe our work is a significant step towards the use of GPU based LBM solvers in practice. In near future, we intend to add specific optimisations for compute capability 2.0 and 2.1 hardware, i.e. the latest CUDA capable GPUs. We also plan to extend the TheLMA framework on which our program is based to multi-node implementations.

References

- [1] D. d’Humières, I. Ginzburg, M. Krafczyk, P. Lallemand, and L.S. Luo. Multiple-relaxation-time lattice Boltzmann models in three dimensions. *Philosophical Transactions: Mathematical, Physical and Engineering Sciences*, pages 437–451, 2002.
- [2] J. Dongarra, S. Moore, G. Peterson, S. Tomov, J. Allred, V. Natoli, and D. Richie. Exploring new architectures in accelerating CFD for Air Force applications. In *Proceedings of HPCMP Users Group Conference*, pages 14–17. Citeseer, 2008.
- [3] X. He, S. Chen, and G. D. Doolen. A novel thermal model for the lattice boltzmann method in incompressible limit. *Journal of Computational Physics*, 146(1):282–300, 1998.
- [4] P. Lallemand and L. S. Luo. Theory of the lattice Boltzmann method: Acoustic and thermal properties in two and three dimensions. *Physical review E*, 68(3):36706, 2003.
- [5] V.W. Lee, C. Kim, J. Chhugani, M. Deisher, D. Kim, A.D. Nguyen, N. Satish, M. Smelyanskiy, S. Chennupaty, P. Hammarlund, et al. Debunking the 100X GPU vs. CPU myth: an evaluation of throughput computing on CPU and GPU. In *ACM SIGARCH Computer Architecture News*, volume 38, pages 451–460. ACM, 2010.
- [6] G. R. McNamara and G. Zanetti. Use of the Boltzmann Equation to Simulate Lattice-Gas Automata. *Phys. Rev. Lett.*, 61:2332–2335, 1988.
- [7] P. Micikevicius. 3D finite difference computation on GPUs using CUDA. In *Proceedings of 2nd Workshop on General Purpose Processing on Graphics Processing Units*, pages 79–84. ACM, 2009.
- [8] D. Mudigere. Data access optimized applications on the GPU using NVIDIA CUDA. Master’s thesis, Technische Universität München, 2009.
- [9] NVIDIA. *Compute Unified Device Architecture Programming Guide version 3.1.1*, July 2010.
- [10] C. Obrecht, F. Kuznik, B. Tourancheau, and J.-J. Roux. Thermal LBM on Manycore Architectures. www.thelma-project.info, 2010.

- [11] C. Obrecht, F. Kuznik, B. Tourancheau, and J.-J. Roux. A new approach to the lattice Boltzmann method for graphics processing units. *Computers and Mathematics with Applications*, 12(61):3628–3638, June 2011.
- [12] C. Obrecht, F. Kuznik, B. Tourancheau, and J.-J. Roux. Global Memory Access Modelling for Efficient Implementation of the LBM on GPUs. In J.M.L.M. et al. Palma, editor, *High Performance Computing for Computational Science – VECPAR2010, Lecture Notes in Computer Science 6449*, pages 151–161. Springer, 2011.
- [13] C. Obrecht, F. Kuznik, B. Tourancheau, and J.-J. Roux. Multi-GPU implementation of the lattice Boltzmann method. *Computers and Mathematics with Applications*, (doi:10.1016/j.camwa.2011.02.020), 2011.
- [14] C. Obrecht, F. Kuznik, B. Tourancheau, and J.-J. Roux. The TheLMA project: Multi-GPU implementation of the lattice Boltzmann method. *International Journal of High Performance Computing Applications*, 25(3): 295–303, August 2011.
- [15] Y. H. Qian. Simulating thermohydrodynamics with lattice BGK models. *Journal of scientific computing*, 8(3):231–242, 1993.
- [16] J. Tölke. Implementation of a Lattice Boltzmann kernel using the Compute Unified Device Architecture developed by nVIDIA. *Computing and Visualization in Science*, pages 1–11, 2008.
- [17] E. Tric, G. Labrosse, and M. Betrouni. A first incursion into the 3d structure of natural convection of air in a differentially heated cubic cavity, from accurate numerical solutions. *International Journal of Heat and Mass Transfer*, 43(21):4043 – 4056, 2000.