

Optimization of a discontinuous Galerkin solver with OpenCL and StarPU

Bérenger Bramas, Philippe Helluy, Laura Mendoza, Bruno Weber

► **To cite this version:**

Bérenger Bramas, Philippe Helluy, Laura Mendoza, Bruno Weber. Optimization of a discontinuous Galerkin solver with OpenCL and StarPU. *International Journal on Finite Volumes*, Institut de Mathématiques de Marseille, AMU, 2020, 15 (1), pp.1-19. hal-01942863v3

HAL Id: hal-01942863

<https://hal.archives-ouvertes.fr/hal-01942863v3>

Submitted on 23 Jan 2020

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.

Optimization of a discontinuous Galerkin solver with OpenCL and StarPU

Bérenger Bramas[†]

[†] *Université de Strasbourg, Icube, Inria Camus.*

berenger.bramas@inria.fr

Philippe Helluy, Laura Mendoza^{*}

^{*} *Université de Strasbourg, IRMA, Inria Tonus.*

philippe.helluy@unistra.fr, laura.mendoza@inria.fr

Bruno Weber[‡]

[‡] *AxesSim, Illkirch.*

bruno.weber@axessim.fr

Abstract

Since the recent advance in microprocessor design, the optimization of computing software becomes more and more technical. One of the difficulties is to transform sequential algorithms into parallel ones. A possible solution is the task-based design. In this approach, it is possible to describe the parallelization possibilities of the algorithm automatically. The task-based design is also a good strategy to optimize software in an incremental way. The objective of this paper is to describe a practical experience of a task-based parallelization of a Discontinuous Galerkin method in the context of electromagnetic simulations. The task-based description is managed by the StarPU runtime. Additional acceleration is obtained by OpenCL.

Key words : Discontinuous Galerkin, StarPU, OpenCL, hybrid parallelism.

1 Introduction

The development of engineering simulation software faces two major antagonist issues. On the one hand, the user is interested in handling more and more complex physical phenomena, in devices with arbitrary geometries. These constraints require to adopt a generic and abstract software engineering approach in order to ensure the generality of the code. On the other hand, the user also wants to harness the full computational power of his hybrid computer. This requirement generally imposes to use low-level hardware optimization hacks that are not always compatible with a generic and universal approach. In addition, as the hardware evolves, optimizations of yesterday are not necessarily relevant for the devices of tomorrow.

The Discontinuous Galerkin algorithm is a general method for approximating system of conservation laws coming from physics. For a recent review of this method, we refer to [11]. This method is well adapted to parallel computing, GPU computing and various software optimization. See for instance [14, 10]. However these optimizations are complicated to program with low-level libraries. The task is even more complex when the software is run on a hybrid computer made of different computing accelerators. The aim of this work is to describe a practical optimization of the Discontinuous Galerkin method that tries to keep the software design as simple as possible.

Several software environments have been developed recently for addressing massively multicore accelerators, such as Graphic Processing Units (GPUs). CUDA is one of these tools, it is devoted to NVIDIA GPUs. OpenCL (Open Computing Language) is another tool. Its main advantage over CUDA is that it can be used with hardware of different brands (NVIDIA, AMD, Intel, ARM, *etc.*). OpenCL is a nice software environment for optimization purposes. OpenCL specifies programming languages (based on C99 and C++11) for programming these devices and application programming interfaces (APIs) to control the platform and execute programs on the compute devices. OpenCL provides a standard interface for parallel computing OpenCL has been used successfully in several works for accelerating Finite Volume (FV) or Discontinuous Galerkin (DG) programs. See for instance [4, 14, 9, 17, 10, 5]. It presents the good balance between an abstract view of the computing devices and some important hardware aspects, such as local memory for accelerating data transfers or subgroups for minimizing synchronization barriers. However depending on architecture, optimizations written for one accelerator may be completely irrelevant for another one. A CPU (Central Processing Unit), a discrete GPU (Graphic Processing Unit) or an IGP (Integrated Graphic Processor) have completely different behavior concerning local cache memory optimizations. For instance, cache prefetching is generally efficient for discrete GPUs while inefficient for IGPs.

In this paper, we present our practical approach to this issue, with the help of the StarPU

This work has obtained support from the ANR EXAMAG project and the BPI France HOROCH project.

runtime system. StarPU is developed at Inria Bordeaux since 2006. StarPU is a runtime C library based on the dataflow paradigm. The programmer has to split the computation workload into abstract computational tasks. Each task processes data buffers that can be in `read` (R), `write` (W) or `read/write` (RW) mode.

The tasks are then conveniently implemented into *codelets*. It is possible (and recommended) to write several implementations of the same task into several codelets. For instance, one can imagine writing a standard codelet in C for validation purposes and one or several optimized OpenCL codelets. The user then submits the tasks in a sequential way to StarPU. At runtime, StarPU is able to construct a task graph based on buffer dependencies. The tasks are submitted to the available accelerators, in parallel if the dependencies allow it. StarPU decides at runtime which version of the codelet (C or OpenCL) is executed. The scheduler is in charge of the distribution of the tasks over the processing units. Dozens of schedulers are provided by StarPU and it is possible to create a custom scheduler to plug it inside StarPU. Some schedulers take into account data locality or use performance models by measuring the efficiency of the different codelets' implementations in order to choose the best one. In addition, StarPU automatically handles data transfers between the accelerators: when a task is assigned to a processing unit, StarPU moves the data that are not already on the processing unit's memory and it can also evict old data if the available space is not enough.

With StarPU, implementing a complex dataflow of OpenCL kernels becomes easier. Indeed, the programmer does not have to handle the kernel dependencies with OpenCL events. In addition, it is possible to first write a well validated pure C version of the software, with only C codelets. Then, one can enrich the tasks with OpenCL codelets, which allows an incremental optimization of the code. At each stage, StarPU should be able to use the available codelets in an efficient way.

In this paper, we describe how we applied the StarPU philosophy in order to optimize a discontinuous finite element solver for conservation laws. A similar experience is also presented for instance in [5]. The outlines are as follows: first we will present in its main lines the Discontinuous Galerkin (DG) scheme that is used in the solver. Then, after a short presentation of StarPU, we will explain how we have integrated the OpenCL optimizations into the DG solver. Finally, we will present some numerical experiments.

2 Discontinuous Galerkin method

The Discontinuous Galerkin (DG) method is a general finite element method for approximating systems of conservation laws of the form

$$\partial_t \mathbf{w} + \sum_{k=1}^D \partial_k \mathbf{f}^k(\mathbf{w}) = 0.$$

The unknown is the vector of conservative variables $\mathbf{w}(\mathbf{x}, t) \in \mathbb{R}^m$ depending on the space variable $\mathbf{x} = (x^1, \dots, x^D) \in \mathbb{R}^D$ and on time t . In this paper, the space dimension is $D = 3$. The flux components \mathbf{f}^k , $k = 1 \dots D$, are smooth functions from \mathbb{R}^m to \mathbb{R}^m . We adopt the notations ∂_t for the partial derivative with respect to t and ∂_k for the partial derivative with respect to x^k . Let $\mathbf{n} = (n_1, \dots, n_D) \in \mathbb{R}^D$ be a spatial direction. The flux \mathbf{f} is a smooth function from $\mathbb{R}^m \times \mathbb{R}^D$ to \mathbb{R}^m , defined by

$$\mathbf{f}(\mathbf{w}, \mathbf{n}) = \sum_{k=1}^D n_k \mathbf{f}^k(\mathbf{w}).$$

For instance, in this work we consider the numerical simulation of an electromagnetic wave. In this particular case, the conservative variables are

$$\mathbf{w} = (\mathbf{E}^T, \mathbf{H}^T, \lambda, \mu)^T \in \mathbb{R}^m, \quad m = 8,$$

where $\mathbf{E} \in \mathbb{R}^3$ is the electric field, $\mathbf{H} \in \mathbb{R}^3$ is the magnetic field and λ, μ are divergence cleaning potentials ([16]). The flux is given by

$$\mathbf{f}(\mathbf{w}, \mathbf{n}) = \begin{pmatrix} -\mathbf{n} \times \mathbf{H} + \lambda \mathbf{n} \\ \mathbf{n} \times \mathbf{E} + \mu \mathbf{n} \\ c\mathbf{n} \cdot \mathbf{E} \\ c\mathbf{n} \cdot \mathbf{H} \end{pmatrix}$$

and $c > 1$ is the divergence cleaning parameter.

The computational domain Ω is an open set of \mathbb{R}^D . We consider a mesh \mathcal{M} of Ω made of N_c open sets, called *cells*, $\mathcal{M} = \{L_i, i = 1 \dots N_c\}$. In the most general setting, the cells satisfy

1. $L_i \cap L_j = \emptyset$, if $i \neq j$;
2. $\overline{\cup_i L_i} = \overline{\Omega}$.

In each cell $L \in \mathcal{M}$, we consider a basis of N_d functions $(\varphi_{L,i}(\mathbf{x}))_{i=0 \dots N_d-1}$. The basis functions are constructed by tensor products of polynomials of order d . Therefore $N_d = (d+1)^D$. We denote by h the maximal diameter of the cells. We assume that the local approximation \mathbf{w}_L of \mathbf{w} in cell L is given by

$$\mathbf{w}_L(\mathbf{x}, t) = \sum_{j=0}^{N_d-1} \mathbf{w}_{L,j}(t) \varphi_{L,j}(\mathbf{x}), \quad \mathbf{x} \in L.$$

The DG formulation then reads: find $\mathbf{w}_{L,j}$ such that for all cell L and all test function $\varphi_{L,i}$

$$\int_L \partial_t \mathbf{w}_L \varphi_{L,i} - \int_L \mathbf{f}(\mathbf{w}_L, \nabla \varphi_{L,i}) + \int_{\partial L} \mathbf{f}(\mathbf{w}_L, \mathbf{w}_R, \mathbf{n}) \varphi_{L,i} = 0. \quad (2.1)$$

In this formula (see Figure 2.1):

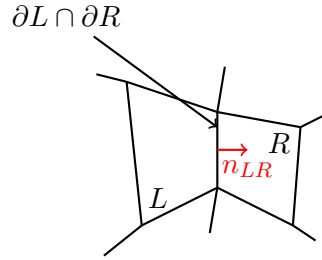


Figure 2.1: Convention for the L and R cells orientation.

- R denotes the neighboring cell to L along its boundary $\partial L \cap \partial R$, or the exterior of Ω on $\partial L \cap \partial \Omega$.
- $\mathbf{n} = \mathbf{n}_{LR}$ is the unit normal vector on ∂L oriented from L to R .
- \mathbf{w}_R denotes the value of \mathbf{w} in the neighboring cell R on $\partial L \cap \partial R$.
- If L is a boundary cell, one may have to use the boundary values instead: $\mathbf{w}_R = \mathbf{w}_b$ on $\partial L \cap \partial \Omega$.
- $\mathbf{f}(\mathbf{w}_L, \mathbf{w}_R, \mathbf{n})$ is the standard upwind numerical flux encountered in many finite volume or DG methods.

In our application, we consider hexahedral cells. We have a reference cell

$$\hat{L} =] - 1, 1[^D$$

and a smooth transformation $\mathbf{x} = \boldsymbol{\tau}_L(\hat{\mathbf{x}})$, $\hat{\mathbf{x}} \in \hat{L}$, that maps \hat{L} on L

$$\boldsymbol{\tau}_L(\hat{L}) = L.$$

We assume that $\boldsymbol{\tau}_L$ is invertible and we denote by $\boldsymbol{\tau}'_L$ its (invertible) Jacobian matrix. We also assume that $\boldsymbol{\tau}_L$ is a direct transformation

$$\det \boldsymbol{\tau}'_L > 0.$$

In our implementation, $\boldsymbol{\tau}_L$ is a quadratic map based on hexahedral curved ‘‘H20’’ finite elements with 20 nodes. The mesh of H20 finite elements is generated by `gmsh` ([6]).

On the reference cell, we consider the Gauss-Lobatto points $(\hat{\mathbf{x}}_i)_{i=0 \dots N_d-1}$, $N_d = (d+1)^D$ and associated weights $(\hat{\omega}_i)_{i=0 \dots N_d-1}$. They are obtained by tensor products of the $(d+1)$ one-dimensional Gauss-Lobatto (GL) points on $] - 1, 1[$. The reference GL points and weights are then mapped to the physical GL points of cell L by

$$\mathbf{x}_{L,i} = \boldsymbol{\tau}_L(\hat{\mathbf{x}}_i), \quad \omega_{L,i} = \hat{\omega}_i \det \boldsymbol{\tau}'_L(\hat{\mathbf{x}}_i) > 0. \quad (2.2)$$

In addition, the six faces of the reference hexahedral cell are denoted by F_ϵ , $\epsilon = 1 \dots 6$ and the corresponding outward normal vectors are denoted by $\hat{\mathbf{n}}_\epsilon$. One advantage of choosing the GL points is that the cells and the faces share the same quadrature points. We use the following notations to define the face quadrature weights:

- if a GL point $\hat{\mathbf{x}}_i \in F_\epsilon$, we denote by μ_i^ϵ the corresponding quadrature weight on face F_ϵ ;
- we also use the convention that $\mu_i^\epsilon = 0$ if $\hat{\mathbf{x}}_i$ does not belong to face F_ϵ .

Let us remark that a given GL point $\hat{\mathbf{x}}_i$ can belong to several faces when it is on an edge or in a corner of \hat{L} . Because of symmetry, we observe that if $\mu_i^\epsilon \neq 0$, then the weight μ_i^ϵ does not depend on ϵ .

We then consider basis functions $\hat{\varphi}_i$ on the reference cell: they are the Lagrange polynomials associated to the Gauss-Lobatto points and thus satisfy the interpolation property

$$\hat{\varphi}_i(\hat{\mathbf{x}}_j) = \delta_{ij}.$$

The basis functions on cell L are then defined according to the formula

$$\varphi_{L,i}(\mathbf{x}) = \hat{\varphi}_i(\boldsymbol{\tau}_L^{-1}(\mathbf{x})).$$

In this way, they also satisfy the interpolation property

$$\varphi_{L,i}(\mathbf{x}_{L,j}) = \delta_{ij}. \quad (2.3)$$

In this paper, we only consider conformal meshes: the GL points on cell L are supposed to match the GL points of cell R on their common face. Dealing with non-matching cells is the object of a forthcoming work. Let L and R be two neighboring cells. Let $\mathbf{x}_{L,j}$ be a GL point in cell L that is also on the common face between L and R . In the case of conformal meshes, it is possible to define the index j' such that

$$\mathbf{x}_{L,j} = \mathbf{x}_{R,j'}.$$

Applying a numerical integration to (2.1), using (2.2) and the interpolation property (2.3), we finally obtain

$$\partial_t \mathbf{w}_{L,i} \omega_{L,i} - \sum_{j=0}^{N_d-1} \mathbf{f}(\mathbf{w}_{L,j}, \nabla \varphi_{L,i}(\mathbf{x}_{L,j})) \omega_{L,j} + \sum_{\epsilon=1}^6 \mu_i^\epsilon \mathbf{f}(\mathbf{w}_{L,i}, \mathbf{w}_{R,i'}, \mathbf{n}_\epsilon(\mathbf{x}_{L,i})) = 0. \quad (2.4)$$

We have to detail how the gradients and normal vectors are computed in the above formula. Let \mathbf{A} be a square matrix. We recall that the cofactor matrix of \mathbf{A} is defined by

$$\text{co}(\mathbf{A}) = \det(\mathbf{A}) (\mathbf{A}^{-1})^T. \quad (2.5)$$

The gradient of the basis function is computed from the gradients on the reference cell using (2.5)

$$\nabla \varphi_{L,i}(\mathbf{x}_{L,j}) = \frac{1}{\det \boldsymbol{\tau}'_L(\hat{\mathbf{x}}_j)} \text{co}(\boldsymbol{\tau}'_L(\hat{\mathbf{x}}_j)) \hat{\nabla} \hat{\varphi}_i(\hat{\mathbf{x}}_j).$$

In the same way, the scaled normal vectors \mathbf{n}_ϵ on the faces are computed by the formula

$$\mathbf{n}_\epsilon(\mathbf{x}_{L,i}) = \text{co}(\boldsymbol{\tau}'_L(\hat{\mathbf{x}}_i)) \hat{\mathbf{n}}_\epsilon.$$

We introduce the following notation for the cofactor matrix

$$\mathbf{c}_{L,i} = \text{co}(\boldsymbol{\tau}'_L(\hat{\mathbf{x}}_i)).$$

The DG scheme then reads

$$\partial_t \mathbf{w}_{L,i} - \frac{1}{\omega_{L,i}} \sum_{j=0}^{N_d-1} \mathbf{f}(\mathbf{w}_{L,j}, \mathbf{c}_{L,j} \hat{\nabla} \hat{\varphi}_i(\hat{\mathbf{x}}_j)) \hat{\omega}_j + \frac{1}{\omega_{L,i}} \sum_{\epsilon=1}^6 \mu_i^\epsilon \mathbf{f}(\mathbf{w}_{L,i}, \mathbf{w}_{R,i'}, \mathbf{c}_{L,i} \hat{\mathbf{n}}_\epsilon) = 0. \quad (2.6)$$

REMARK 1 We observe that in our presentation, we have not used at any time the linearity of the Maxwell flux. This means that the same method could be used for non-linear conservation laws. However, for non-linear equations it is necessary to supplement (2.6) with a limiting strategy in order to avoid nonphysical oscillations (see for instance [7] and included references).

REMARK 2 Formulation (2.6) is an approximation of (2.1) because of numerical integration by the Gauss-Lobato method. Gauss-Lobato integration can generate inaccuracies, known in the literature as ‘‘aliasing’’ effects [11, 18]. It is possible to improve precision with a slight modification of the DG method [18]. In this paper we have not yet implemented this modification, because we mainly address performance issues. But it is an objective for us in the future.

In formula (2.6) we identify volume terms that are associated to Gauss-Lobatto points in the volume

$$V_{i,j} = \mathbf{f}(\mathbf{w}_{L,j}, \mathbf{c}_{L,j} \hat{\nabla} \hat{\varphi}_i(\hat{\mathbf{x}}_j)) \hat{\omega}_j \quad (2.7)$$

and surface terms that are associated to cell boundaries

$$S_{i,i'}^\epsilon = \mu_i^\epsilon \mathbf{f}(\mathbf{w}_{L,i}, \mathbf{w}_{R,i'}, \mathbf{c}_{L,i} \hat{\mathbf{n}}_\epsilon). \quad (2.8)$$

Then, once these terms have been accumulated, one has to apply the inverse of the mass matrix. Here, this matrix is diagonal and it corresponds to a simple multiplication of $\sum_\epsilon S_{i,i'}^\epsilon + \sum_j V_{i,j}$ by

$$\frac{1}{\omega_{L,i}}. \quad (2.9)$$

Generally, i and i' are associated to unknowns of internal cells, *i.e.* cells that do not touch the boundaries. On boundary GL points, the value of $\mathbf{w}_{R,i'}$ is given by the boundary condition

$$\mathbf{w}_{R,i'} = \mathbf{w}_b(\mathbf{x}_{L,i}, t), \quad \mathbf{x}_{L,i} = \mathbf{x}_{R,i'}.$$

For practical reasons, it can be interesting to also consider $\mathbf{w}_{R,i'}$ as an artificial unknown in a fictitious cell. The fictitious unknown is then a solution of the differential equation

$$\partial_t \mathbf{w}_{R,i'} = \partial_t \mathbf{w}_b(\mathbf{x}_{L,i}, \cdot). \quad (2.10)$$

In the end, if we put all the unknowns in a large vector $\mathbf{W}(t)$, (2.6) and (2.10) read as a large system of coupled differential equations

$$\partial_t \mathbf{W} = \mathbf{F}(\mathbf{W}). \quad (2.11)$$

This set of differential equations is then numerically solved by a Runge-Kutta numerical method. In practice, we use a second order Runge-Kutta method (RK2).

3 Data-based parallelism and StarPU

3.1 StarPU philosophy

StarPU is a runtime system library developed at Inria Bordeaux ([2, 1]). It relies on the data-based parallelism paradigm. The user has first to split its whole problem into elementary computational tasks. The elementary tasks are then implemented into *codelets*, which are simple C functions. The same task can be implemented differently into several codelets. This allows the user to harness special accelerators, such as vectorial CPU cores or OpenCL devices, for example. In the StarPU terminology these devices are called *workers*. If a codelet contains OpenCL kernel submissions, special utilities are available in order to map the StarPU buffers to OpenCL buffers.

For each task, the user also has to describe precisely what are the input data, in `read` mode, and the output data, in `write` or `read/write` mode. The user then submits the task in a sequential way to the StarPU system. StarPU is able to construct at runtime a task graph from the data dependencies. The task graph is analyzed and the tasks are scheduled automatically to the available workers (a worker could be a single CPU core, a GPUs, *etc.*). If possible, they are executed in parallel into concurrent threads. The data transfer tasks between the threads are automatically generated and managed by StarPU. OpenCL data transfers are also managed by StarPU.

When a StarPU program is executed, it is possible to choose among several schedulers. The simplest **eager** scheduler adopts a very simple strategy: the tasks are executed in the order of submission by the free workers, without optimization. More sophisticated schedulers,

such as the `dmda` or `dmdar` scheduler, are able to measure the efficiency of the different codelets and the data transfer times, in order to apply a more efficient allocation of tasks.

Recently a new data access mode has been added to StarPU: the `commute` mode. In a task, a buffer of data can now be accessed in `commute` mode, in addition to the `write` or `read/write` modes. A `commute` access tells to StarPU that the execution of the corresponding task may be executed before or after other tasks containing commutative access. This can drastically increase the degree of parallelism as StarPU is able to change the order of the tasks that access the same data in `commute` mode.

There also exists a MPI version of StarPU, which allows to address clusters of hybrid nodes. The MPI version of StarPU is not used in this work.

3.2 Macrocell approach

StarPU is quite efficient, but there is an unavoidable overhead due to the task submissions and to the on-the-fly construction and analysis of the task graph. Therefore it is important to ensure that the computational tasks are not too small, in which case the overhead would not be amortized, or not too big, in which case some workers would be idle. To achieve the right balance, we decided not to apply the DG algorithm at the cell level but to groups of cells instead.

The implementation of the DG scheme has been made into the `schnaps` software, which is a C99 software dedicated to the numerical simulation of conservation laws. In `schnaps`, we first construct a *macromesh* of the computational domain. Then each *macrocell* of the macromesh is split into *subcells*. See Figure 3.1. We also arrange the subcells into a regular sub-mesh of the macrocells. In this way, it is possible to apply additional optimizations. For instance, the subcells L of a same macrocell \mathcal{L} share the same geometrical transformation τ_L , which saves memory access.

In `schnaps` we also defined an *interface* structure in order to manage data communications between the macrocells. An interface contains a copy of the data at the Gauss-Lobatto points that are common to the two neighboring macrocells.

To solve one time step of the DG method, here are the most important tasks:

1. `Interface_extraction`: this task simply extracts the values of \mathbf{w} coming from the neighboring macrocells to the interface. In this task, the macrocell buffers are accessed in `read` mode and the interface data in `write` mode. Of course, if the interface is a boundary interface, only one side is extracted.
2. `Interface_fluxes`: this task only computes the numerical fluxes (2.8) that are located at the Gauss-Lobatto points on the interface. The fluxes are then applied to

Solveur Conservatif Hyperbolique Non-linéaire Appliqué aux PlaSmas: <https://gitlab.inria.fr/bramas/schnaps>

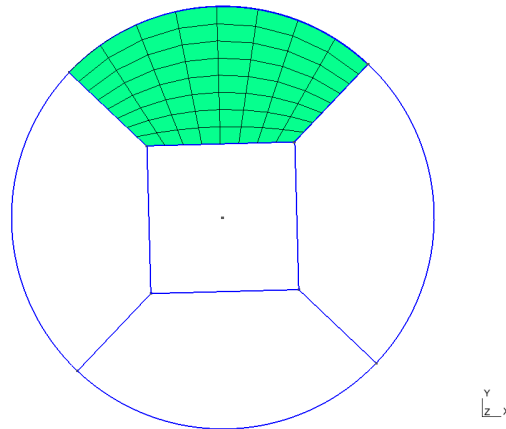


Figure 3.1: Macrocell approach: an example of a mesh made of five macrocells. Each macrocell is then split into several subcells. Only the subcells of the top macrocell are represented here (in green).

the neighboring macrocells. In this task, the interface data are accessed in `read` mode and the macrocell data in `read/write` mode. For a better efficiency, we also assume a `commute` access to the macrocell data. In this way the interface fluxes can be assembled in any order, which can help the StarPU scheduling.

3. **Interface_boundary_fluxes**: this task computes the boundary data and only numerical fluxes (2.8) that are located at the boundary interface. The fluxes are applied to the neighboring macrocell. In this task, the interface data are accessed in `read` mode and the macrocell data in `read/write` and `commute` mode.
4. **Volume_terms**: this task applies the volumic terms (2.7) in a given macrocell. The macrocell data are accessed in `read/write/commute` mode.
5. **Subcell_fluxes**: this task applies the numerical subcell fluxes (2.8) that are internal to a given macrocell. The macrocell data are accessed in `read/write/commute` mode.
6. **Mass_matrix**: this task apply the inverse of the diagonal mass matrix. The macrocell data are accessed in `read/write` mode.

Additional simple tasks are needed to apply the Runge-Kutta algorithm. We do not describe them.

The general sequential DG algorithm is then given in Algorithm 1. Thanks to StarPU, this algorithm can be submitted in a sequential way. StarPU then uses the data dependency in order to distribute the tasks in parallel on the available workers. The wait for task

completion is done at the end of the algorithm and not inside each iteration. This is not an essential point for performance but we have observed in some cases a 15% gain compared to an algorithm with a synchronization point at the end of each time iteration. Let us also observe that this synchronization strategy assumes that the time step is known at the beginning of the computations. This is not always possible especially when solving non-linear systems.

Algorithm 1: Sequential DG algorithm. The parallelism is automatically detected by StarPU, thanks to the data dependency.

```

1 forall interfaces do
  | // task (1)
2   Extract interface (copy the data from the neighboring macrocells to the interface);
3   if the interface is a boundary interface then
  |   // task (3)
4   |   compute the boundary fluxes (Equation 2.8) and apply them to the neighboring
  |   |   macrocell;
5   else
  |   // task (2)
6   |   the interface is an internal one. Compute the numerical fluxes (Equation 2.8)
  |   |   and apply them to the two neighboring macrocells;
7   end
8 end
9 foreach macrocell do
  | // task (5)
10  Compute and apply the numerical fluxes (Equation 2.8) between the subcells ;
  | // task (4)
11  Compute the volume terms (Equation 2.7) inside the subcells;
  | // task(6)
12  Apply the inverse of the mass matrix (Equation 2.9) inside the subcells ;
13 end

```

In the next section we give more details about the implementation of the OpenCL codelets.

4 Hybrid C/OpenCL codelets

In order to attain better performance, we programmed an OpenCL version of the previously described codelets. The CPU workers will continue to use the kernels implemented in C, while the GPU workers will use the ones in OpenCL. When computing on GPUs, it is generally advised to care about memory access in order to improve coalescence [8]. The

values of \mathbf{w} at the Gauss-Lobatto points are stored into what we call a *field* structure. A field is attached to a macrocell. A given component of \mathbf{w} in a field is located thanks to its subcell index \mathbf{ic} , a Gauss-Lobatto index \mathbf{ig} in the cell, and a component index \mathbf{iw} . If the macrocell is cut into n_r subcells in each direction, if the polynomial order is d and for a system of m conservation laws, these indices have the following bounds

$$0 \leq \mathbf{ic} < n_r^3, \quad 0 \leq \mathbf{ig} < (d+1)^3, \quad 0 \leq \mathbf{iw} < m.$$

The values of \mathbf{w} in a given field are stored into a memory buffer `wbuf`. In order to test several memory arrangements, the index in the buffer is computed by a function `varindex(ic,ig,iw)` that we can easily change. For instance, we can consider the following formula, in which the electromagnetic components at a given Gauss-Lobatto point are grouped in memory

$$\text{varindex}(\mathbf{ic}, \mathbf{ig}, \mathbf{iw}) = \mathbf{iw} + m * (\mathbf{ig} + (d+1)^3 * \mathbf{ic}), \quad (4.1)$$

or this formula

$$\text{varindex}(\mathbf{ic}, \mathbf{ig}, \mathbf{iw}) = \mathbf{ig} + (d+1)^3 * (\mathbf{iw} + m * \mathbf{ic}), \quad (4.2)$$

in which the electromagnetic components are separated in memory. A few years ago, we could observe large differences in performance between the two strategies. Here, we have not observed significant differences, maybe because the cache memory of the GPUs used in this work (NVIDIA GTX 1050 or NVIDIA P100) is now larger.

We have programmed an OpenCL codelet for each task described in Section 3.2. The most time-consuming kernels are: (i) the kernel associated to the computations of the volume terms (2.7) inside the macrocells, *volume* kernel, and (ii) the kernel that computes the subcell fluxes, the *surface* kernel. Boundary and interface terms generally involve fewer computations.

A natural distribution of the workload is to associate one subcell to each OpenCL work-group and the computations at a Gauss point to a work-item. With this natural distribution, formula (4.2) ensures optimal coalescent memory access in the volume kernel, but the access is no more optimal in the surface kernel, because neighboring Gauss points on the subcell interfaces are not necessarily close in memory. Therefore, we prefer considering formula (4.1).

5 Numerical results

In this section, we present some practical experiments that we realized with `schnaps`. The first experiment deals with the efficiency of the StarPU C99 codelets on a multicore CPU. The second experiment deals with the efficiency of the OpenCL codelets. Then we test the code in a CPU/GPU configuration.

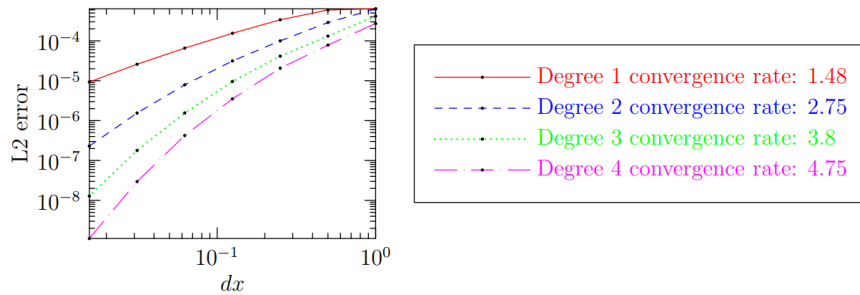


Figure 5.1: Numerical convergence rates for the DG method with Gauss-Lobato numerical integration. Plane wave case in a cube.

We approximate by the DG algorithm an exact plane wave solution of the Maxwell equations

$$\mathbf{E} = c \begin{pmatrix} -v \\ u \\ 0 \end{pmatrix}, \quad \mathbf{H} = c \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \quad \lambda = 0, \quad \mu = 0.$$

with $c = -\cos(\pi/2(ux^1 + vx^2 - t))$ and $(u, v) = (\cos(\pi/4), \sin(\pi/4))$.

For validation purpose, we have first checked the numerical order of the method. For this, we consider a cubic domain

$$\Omega =]0, 1[^3.$$

We apply the exact solution at the boundary $\partial\Omega$ and in Ω at the initial time $t = 0$. We perform numerical simulations for a given final time T . We vary the order d of the polynomial in the DG method and the mesh refinement. We compare the error of the method in the discrete L^2 norm (computed with Gauss-Lobato numerical integration) at that final time $t = T$. In a log scale, we obtain the curves of Figure 5.1. As other authors (see, for instance, [12, 13]), we observe that the order of convergence is between d and $d + 1$. With Gauss-Legendre numerical integration we would have obtained more precise results [15].

We consider the propagation of the previous electromagnetic plane wave into a torus-shaped domain Ω represented on Figure 5.2. In the DG solver, we consider only third order polynomials: $d = 3$. The domain is split into $M = 400$ macrocells. The macrocells are split into n_r subcells in each direction. The number of degrees of freedom is thus

$$n_{\text{d.o.f.}} = Mm(d + 1)^3 n_r^3, \quad M = 400, \quad m = 8, \quad d = 3.$$

StarPU will have to take decisions on how to distribute efficiently the tasks on the available accelerators. We have performed our experiments on two different computers. The first

config.	CPU	# cores	mem. CPU	GPU	mem. GPU
“PC”	AMD FX-8320E 3.2 GHz	8	8 GB	NVIDIA GTX 1050 Ti	4 GB
“WS”	Intel Xeon E5-2683 v4 2.1GHz	2×16	256 GB	NVIDIA P100	16 GB

Table 1: Configurations of the computers used for the tests.

n_r	n_{cpu}	Time (s)	Ideal time	Efficiency
4	2	704	704	1
4	4	365	352	0.96
4	7	209	201	0.96
6	7	689	679	0.99
8	4	2902	2816	0.97
8	7	1650	1609	0.98

 Table 2: Efficiency (with two threads as reference) of the CPU codelets for the “PC” configuration with the `dmdar` scheduler.

configuration “PC” correspond to a standard desktop personal computer. The second configuration “WS” is made of a more powerful two-CPU workstation. The technical details are listed in Table 1.

5.1 Multicore CPU tests

In the first test, we only activate the C99 codelets, and thus the GPU is not activated. We also vary the number of computing CPU cores from 2 to 7 for the “PC” configuration and from 2 to 31 for the “WS” configuration. One CPU core is anyway reserved for the main StarPU thread.

We perform a fixed number of n_t iterations of the RK2 algorithm. We assume that the number of elementary computing operations increases as $O(n_r^3)$. Ideally, with infinitely fast operations at interfaces, and with n_{cpu} identical CPU cores, the computations time would behave like

$$T \sim C n_t \frac{n_r^3}{n_{\text{cpu}}} \quad (5.1)$$

where C is fixed constant, depending on the hardware. In our benchmark the *efficiency* is the ratio of the measured execution time over this ideal time. Normally it should be smaller than one (because of communications). If it is close to one, the algorithm is very efficient. In Table 2 and 3, we observe an excellent efficiency of the StarPU task distribution on the CPU cores even on a NUMA (Non Uniform Memory Access) computer with two CPUs connected to different memory zones.

n_r	n_{cpu}	Time (s)	Ideal time	Efficiency
4	1	646	646	1
4	2	331	323	0.97
4	4	183.4	161.5	0.88
4	8	95.1	80.7	0.85
4	16	48	40.3	0.84
4	30	26	21.5	0.82
6	1	2172.1	2172.1	1
6	2	1106.7	1086	0.98
6	4	612.2	543	0.88
6	8	317	271.5	0.85
6	16	159	135.8	0.85
6	30	85.2	72.4	0.84
8	1	5139.4	5139.4	1
8	2	2613	2569.7	0.98
8	4	1446.6	1284.9	0.88
8	8	749.1	642.4	0.85
8	16	375.2	321.2	0.85
8	30	200.6	171.3	0.85

Table 3: Efficiency of the CPU codelets for the “WS” configuration with the `laheteroprio` scheduler.

n_r	config.	CPU	GPU	2xGPU	CPU/ GPU	CPU/ 2xGPU	CPU+ GPU	CPU+ 2xGPU	CPU/ CPU+GPU	CPU/ CPU+2xGPU
4	PC	209 s	73 s	-	2.86	-	32 s	-	6.53	-
6	PC	689 s	86 s	-	8.01	-	64 s	-	10.77	-
8	PC	1650 s	171 s	-	9.65	-	130 s	-	12.69	-
4	WS	26 s	23.2 s	18.1 s	1.12	1.43	10.7 s	8.9 s	2.4	2.9
6	WS	85.2 s	26.1 s	20.3 s	3.2	4.19	16.7 s	11 s	5.1	7.8
8	WS	200.6 s	29.8 s	24.3 s	6.7	8.23	29.8 s	18.1 s	6.7	11

Table 4: Efficiency of the CPU/GPU implementation. All timings are given in seconds. All ratio are dimensionless. This table shows that StarPU is always able to harness an increase of hybrid computing capacities.

5.2 OpenCL codelets

In this section, we compare the efficiency of the CPU/C99 and GPU/OpenCL codelets on the two configurations “PC” and “WS”. We first perform a CPU-only computation with all the cores activated (7 for the “PC” platform and 31 for the “WS” platform) and a GPU-only computation. The results are given in Table 4. We observe that the OpenCL codelets are faster than the CPU codelets. The highest efficiency is achieved for the finer meshes ($n_r = 8 t_{\text{CPU}+2\text{xGPU}}$).

5.3 Hybrid CPU/GPU computations

Now that we are equipped with verified OpenCL codelets we can try to run the code with StarPU in order to see if it is able to distribute the tasks efficiently on the available accelerators.

One difficulty is to manage the data transfers efficiently between the accelerators. Several scheduling strategies are available in StarPU. With the `eager` scheduler, the tasks are distributed in the order of submission to the inactive device, without taking into account the cost of memory transfers. This strategy is obviously not optimal.

It is better to choose another scheduler. For the “PC” configuration, we use the `dmda` scheduler provided by StarPU and we use the `LAHeteroprio` scheduler ([3]) for the “WS” configuration. With the `LAHeteroprio` scheduler, the different types of tasks are grouped into buckets and the processing unit pick tasks following an order based on their type, which is CPU or GPU in our case, and the location of the data. The idea is to reduce the memory transfer and to promote the consumption of the tasks by the processing units that are more efficient. Therefore, we configure `LAHeteroprio` by sorting the different operations by arithmetic cost and we specify that the CPU workers should follow this

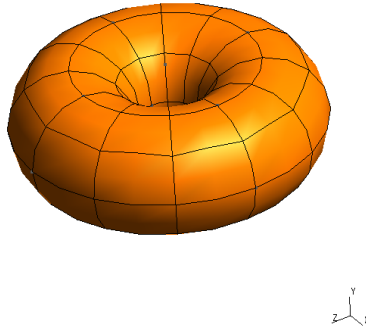


Figure 5.2: Computational domain Ω . Only the macrocells of the mesh are visible. The represented mesh contains 240 macrocells (poloidal refinement set to $n_{pol} = 3$ and toroidal refinement set to $n_{tor} = 3$. In the numerical experiments, we take $n_{pol} = 3$ and $n_{tor} = 5$, which leads to a mesh of 400 macrocells). Each macrocell is then cut regularly into the three directions.

order, while the GPU workers do the opposite: the GPU workers first try to pick tasks with high workloads. Additionally, we ensure that GPU workers privilege data locality over priorities. Therefore, the distribution of the tasks on the available CPU cores and GPU is done dynamically by the scheduler.

The results are given in the last two columns of Table 4. We observe that in most situations, StarPU is able to get an additional gain from the CPU cores, even if the GPU codelets are faster. The fact that for $n_r = 8$ we have the same execution time with one CPU and one GPU ($t_{CPU+GPU}$), or with only one GPU (t_{GPU}) is an illustration of the cost of moving data. More precisely, consider an execution where all the tasks are on the GPU. Using the CPU would mean moving data on it, way and back, and thus can imply an overhead, even if the CPU could be faster for some tasks.

6 Conclusion

We have proposed an optimized implementation of the Discontinuous Galerkin method on hybrid computer made of several GPUs and multicore CPUs. In order to manage the heterogeneous architecture easily and efficiently, we rely on OpenCL for the GPU computations and on the StarPU runtime for distributing the computational tasks on the available devices.

OpenCL programming becomes much easier because the task dependency is computed by

StarPU. We only had to concentrate on the optimization of the individual OpenCL kernels and not on data distribution or memory transfers. We first tested the efficiency of our OpenCL kernels. We verified that the macrocell approach and cache prefetching strategy, while not optimal, give good results.

In addition, with a good choice of scheduler, and for heavy computations, we have shown that StarPU is able to overlap computations and memory transfer in a quite efficient way. It is also able to use the available CPU codelets to achieve even higher acceleration.

References

- [1] Cédric Augonnet, Olivier Aumage, Nathalie Furmento, Raymond Namyst, and Samuel Thibault. StarPU-MPI: Task programming over clusters of machines enhanced with accelerators. In *European MPI Users' Group Meeting*, pages 298–299. Springer, 2012.
- [2] Cédric Augonnet, Samuel Thibault, Raymond Namyst, and Pierre-André Wacrenier. StarPU: a unified platform for task scheduling on heterogeneous multicore architectures. *Concurrency and Computation: Practice and Experience*, 23(2):187–198, 2011.
- [3] Bérenger Bramas. Impact study of data locality on task-based applications through the Heterprio scheduler. *PeerJ Computer Science*, 5:e190, May 2019.
- [4] Anaïs Crestetto and Philippe Helluy. Resolution of the Vlasov-Maxwell system by PIC Discontinuous Galerkin method on GPU with OpenCL. In *ESAIM: Proceedings*, volume 38, pages 257–274. EDP Sciences, 2012.
- [5] Mohamed Essadki, Jonathan Jung, Adam Larat, Milan Pelletier, and Vincent Perrier. A task-driven implementation of a simple numerical solver for hyperbolic conservation laws. *ESAIM: Proceedings and Surveys*, 63:228–247, 2018.
- [6] Christophe Geuzaine and Jean-François Remacle. Gmsh: A 3-D finite element mesh generator with built-in pre-and post-processing facilities. *International journal for numerical methods in engineering*, 79(11):1309–1331, 2009.
- [7] Jean-Luc Guermond, Richard Pasquetti, and Bojan Popov. Entropy viscosity method for nonlinear conservation laws. *Journal of Computational Physics*, 230(11):4248–4267, 2011.
- [8] Philippe Helluy and Jonathan Jung. OpenCL numerical simulations of two-fluid compressible flows with a 2D random choice method. *International Journal on Finite Volumes*, 10:1–38, July 2013.
- [9] Philippe Helluy and Thomas Strub. Multi-GPU numerical simulation of electromagnetic waves. *ESAIM: Proceedings and Surveys*, 45:199–208, 2014.

-
- [10] Philippe Helluy, Thomas Strub, Michel Massaro, and Malcolm Roberts. Asynchronous OpenCL/MPI numerical simulations of conservation laws. In *Software for Exascale Computing-SPPEXA 2013-2015*, pages 547–565. Springer, 2016.
- [11] Jan S Hesthaven and Tim Warburton. *Nodal discontinuous Galerkin methods: algorithms, analysis, and applications*. Springer Science & Business Media, 2007.
- [12] JS Hesthaven and T Warburton. Discontinuous galerkin methods for the time-domain maxwell’s equations. *ACES Newsletter*, 19(ARTICLE):10–29, 2004.
- [13] Florian Hindenlang, Gregor J Gassner, Christoph Altmann, Andrea Beck, Marc Staudenmaier, and Claus-Dieter Munz. Explicit discontinuous galerkin methods for unsteady problems. *Computers & Fluids*, 61:86–93, 2012.
- [14] Andreas Klöckner, Timothy Warburton, and Jan S Hesthaven. High-order Discontinuous Galerkin methods by GPU metaprogramming. In *GPU Solutions to Multi-scale Problems in Science and Engineering*, pages 353–374. Springer, 2013.
- [15] David A Kopriva and Gregor Gassner. On the quadrature and weak form choices in collocation type discontinuous galerkin spectral element methods. *Journal of Scientific Computing*, 44(2):136–155, 2010.
- [16] Claus-Dieter Munz, Pascal Omnes, Rudolf Schneider, Eric Sonnendrücker, and Ursula Voss. Divergence correction techniques for Maxwell solvers based on a hyperbolic model. *Journal of Computational Physics*, 161(2):484–511, 2000.
- [17] Thomas Strub. *Resolution of tridimensional instationary Maxwell’s equations on massively multicore architecture*. Phd thesis, Université de Strasbourg, March 2015.
- [18] Andrew R Winters, Rodrigo C Moura, Gianmarco Mengaldo, Gregor J Gassner, Stefanie Walch, Joaquim Peiro, and Spencer J Sherwin. A comparative study on polynomial dealiasing and split form discontinuous galerkin schemes for under-resolved turbulence computations. *Journal of Computational Physics*, 372:1–21, 2018.