



# A parallel Discontinuous Galerkin Time-Domain solver of Maxwell's equations

Stéphane Lanteri, Raphaël Léger<sup>1</sup>

Nachos Project-Team, Inria Sophia-Antipolis Méditerrannée

4th workshop of the HOSCAR project - Gramado - 2014/09/16





<sup>1</sup>raphael.leger@inria.fr

Stéphane LANTERI & Raphaël LÉCER

HPC with a DGTD solver of Maxwell's equations

Context

Traditional approach: MPI

Hybrid MPI/OpenMP

Future Steps

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = 差 - 約900

Stéphane Lanteri & Raphaël Lécer

4th workshop of the HOSCAR project - Gramado - 2014/09/16

HPC with a DGTD solver of Maxwell's equations

### $\operatorname{Context}$

Traditional approach: MPI

Hybrid MPI/OpenMP

Future Steps

・ロト ・回ト ・ヨト ・ヨー うへの

Stéphane Lanteri & Raphaël Lécer

HPC with a DGTD solver of Maxwell's equations

# The DEEP-er european exascale project

DECEP-ER Dynamic Exascale Entry Platform - Extended Reach<sup>2</sup>

- European Union's FP7 Exascale programme
- ▶ 13 partners in 7 countries
- Follow-up to the DEEP project (ends in 12/2014)
- Kick-off meeting in October 2013 at Forschungszentrum Jülich
- ▶ 36 months project



<sup>2</sup>http://www.deep-er.eu/

HPC with a DGTD solver of Maxwell's equations

# The DEEP-er european exascale project

# DEEP-er architecture

- "Cluster/Booster"
- CN = cluster node
   (Xeon Sandy Bridge)
- BN = booster node
   (Xeon Phi KNL)
- ▶ NAM, NVM, ...

### The goals

Stéphane LANTERI & Raphaël LÉCER

- Design a functional prototype
- Assess the potential of this architecture in view of achieving exascale
- Progress on parallel I/O and large scale resiliency
- Co-design with application developers



#### HPC with a DGTD solver of Maxwell's equations

# The Maxwell-Debye model

▶ We are interested in solving the 3D Maxwell-Debye system:

$$\begin{aligned} & \left( \begin{array}{l} \mu \frac{\partial \mathbf{H}}{\partial t} + \nabla \times \mathbf{E} = \mathbf{0} \\ & -\varepsilon_0 \varepsilon_\infty \frac{\partial \mathbf{E}}{\partial t} + \nabla \times \mathbf{H} = \frac{d \mathbf{P}}{dt} + \sigma \mathbf{E} \\ & \left( \begin{array}{l} \frac{d \mathbf{P}}{dt} = \frac{1}{\tau_r} \left[ \varepsilon_0 \left( \varepsilon_s - \varepsilon_\infty \right) \mathbf{E} - \mathbf{P} \right] \end{array} \right) \end{aligned}$$

- ▶ Electromagnetic wave propagation in heterogeneous, dispersive media
- Well-suited dispersion model for human tissues
- $\blacktriangleright$  An application: numerical dosimetry and SAR evaluation





HPC with a DGTD solver of Maxwell's equations

4th workshop of the HOSCAR project - Gramado - 2014/09/16

Stéphane LANTERI & Raphaël LÉGER

# The DGTD method for the Maxwell-Debye system

### Generalities

- ▶ A cell-local finite element formulation [FEZOUI et al. 2005]
- ▶ Unknowns are represented on a set of polynomial basis functions



2D Lagrange basis functions on the unit triangle

- Cells are weakly coupled to their neighbors through a numerical flux likewise finite volume methods
  - Centered
  - Local Lax-Friedrichs (upwind)
  - ► ...

# The DGTD method for the Maxwell-Debye system

### Features

- High-order of accuracy with a compact-stencil
- Unstructured & hybrid cartesian/unstructured (non-conforming) meshes [DUROCHAT et al. 2014]
- Local order refinement hp adaptation
   [FAHS et al. 2009]
- Local time-stepping / Locally-implicit time stepping [MOYA et al. 2012]



Hybrid unstructured/cartesian mesh example

### Stéphane Lanteri & Raphaël Lécer

HPC with a DGTD solver of Maxwell's equations

### The DGTD method for the Maxwell-Debye system

### Semi-discrete formulation

- ▶ For a centered numerical flux and a conforming tetrahedral mesh
- ▶ The semi discrete system reads:

For each cell 
$$i$$
, 
$$\begin{cases} 2\mathbb{M}_{i}^{\mu}\frac{d\overline{\mathbf{H}}_{i}}{dt} - \sum_{k=1}^{3}\mathbb{K}_{i}^{x_{k}}\overline{\mathbf{E}}_{i} - \sum_{a_{ij}\in\mathscr{T}_{d}^{i}}\mathbb{F}_{ij}\overline{\mathbf{E}}_{j} = \mathbf{0}, \\ 2\mathbb{M}_{i}^{\varepsilon}\frac{d\overline{\mathbf{E}}_{i}}{dt} + \sum_{k=1}^{3}\mathbb{K}_{i}^{x_{k}}\overline{\mathbf{H}}_{i} + \sum_{a_{ij}\in\mathscr{T}_{d}^{i}}\mathbb{F}_{ij}\overline{\mathbf{H}}_{j} = \mathbb{M}_{i}\frac{d\overline{\mathbf{P}}_{i}}{dt} + \mathbb{M}_{i}^{\sigma}\overline{\mathbf{E}}_{i}, \\ \frac{d\overline{\mathbf{P}}_{i}}{dt} = \frac{1}{\tau_{r}^{i}}\left(\varepsilon_{0}\left(\varepsilon_{s}^{i} - \varepsilon_{\infty}^{i}\right)\overline{\mathbf{E}}_{i} - \overline{\mathbf{P}}_{i}\right). \end{cases}$$
(1)

- ▶ Matrices are small (e.g. 3-diagonal blocks of size 20 × 20 in P3)
- ▶ Mass and Stiffness: dense / Interface: sparse (for Lagrange polynomials)
- ▶ In practice: 2nd order leap-frog for time-integration explicit

(D) (A) (A) (A) (A)

Context

### Traditional approach: MPI

Hybrid MPI/OpenMP

Future Steps

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ 三回 めんの

Stéphane Lanteri & Raphaël Lécer

HPC with a DGTD solver of Maxwell's equations

# The traditional approach: full MPI parallelization



- ▶ 5 Materials / 256 subdomains
- ▶ 1.8 Million cells = 7.2 Million to 36 Million Lagrange points (depending on the order of precision, P1, P2, P3)
- Full MPI: based on a mesh partitioning and p2p communications at every timestep

### Stéphane Lanteri & Rapha<u>ël Lécer</u>

HPC with a DGTD solver of Maxwell's equations



- ▶ 256 processes
- 16×16 cores (Sandybridge)
- P2-Lagrange basis - 10 nodes per tetrahedron
- 22500 time-iterations
- 9 min. 07s
   walltime

www.youtube.com/embed/PhtRhzzvl94?rel=0

→ B → < B</p>



Stéphane Lanteri & Raphaël Lécer

HPC with a DGTD solver of Maxwell's equations

4th workshop of the HOSCAR project - Gramado - 2014/09/16

3

### MPI parallelization: strong scaling - blocking communications



 $<sup>^3\</sup>mathrm{Bi}\text{-}\mathrm{socket}$  Westmere nodes - 12 cores per node

# MPI-parallelization performance analysis (extrae/Paraver @BSC)

### P2P duration



### Bytes between events



### Blocking communications lead to late senders

# Stéphane LANTERI & Raphaël LÉCER HPC with a DGTD solver of Maxwell's equations 4th workshop of the HOSCAR project - Gramado - 2014/09/16 15/26



Stéphane Lanteri & Raphaël Lécer

HPC with a DGTD solver of Maxwell's equations

4th workshop of the HOSCAR project - Gramado - 2014/09/16

3

### MPI-parallelization performance analysis (extrae/Paraver @BSC)

### Blocking: average portion of time spent within MPI: 32.54%



### Non-Blocking: average portion of time spent within MPI: 3.75 % !!



Stéphane LANTERI & Raphaël LÉCER

HPC with a DGTD solver of Maxwell's equations

### MPI parallelization: strong scaling - non-blocking communications



- Superlinearity up to 32 nodes of Judge (384 processes)!
- Better cache usage wins over the increase of communication complexity ►

Context

Traditional approach: MPI

Hybrid MPI/OpenMP

Future Steps

・ロト ・団ト ・ヨト ・ヨト ・ヨー ぐくの

Stéphane Lanteri & Raphaël Lécer

4th workshop of the HOSCAR project - Gramado - 2014/09/16

HPC with a DGTD solver of Maxwell's equations

# Hybrid MPI/OpenMP: motivations

# $\begin{array}{l} \mathsf{DEEP}\text{-}\mathsf{er} \ \mathrm{cluster}/\mathrm{booster} \\ \mathrm{prototype} \end{array}$

- Located at Jülich
- Cluster: 256 nodes, sandy bridge bi-socket - 16 cores
   2048 cores total
- Booster: 512 Xeon-Phi KNL (due H2 2015)
- (While waiting for the KNL: developing on the KNC)



### Approaching the DEEP-er architecture: time loop on the booster

- ▶ The Booster: 512 nodes  $\times$  72 cores  $\times$  4 threads = 147456 threads total.
- ▶ Full MPI  $\rightarrow$  147456 processes (i.e. subdomains)!!!
- ▶ Hybrid MPI/OpenMP  $\rightarrow$  512 processes × 288 OMP threads  $\checkmark$ .

# OpenMP in the main loops (the example of updating $\overline{\mathbf{E}}_i$ and $\overline{\mathbf{P}}_i$ )

```
Step 1: compute the flux
$!OMP PARALLEL
$!OMP DO SCHEDULE(STATIC) PRIVATE(...)
DO it=1.nt ! loop on cells
  FLUX(:,:,it) = 0.0d0
  DO if=1,4 ! loop on neighbors (4 faces per tetrahedron)
     it2 = elmvoi (if.it)
     FLUX(:,:,jt) = FLUX(:,:,jt) + F(ua(:,:,jt)+ua(:,:,jt2))*0.5d0
  ENDDO
  FLUX(:,:,it) = FLUX(:,:,it) + K(ua(:,:,it))
ENDDO
$!OMP END DO
Step 2: update the fields
$!OMP DO SCHEDULE(STATIC) PRIVATE(...)
DO it=1,nt ! loop on cells
```

### Stéphane Lanteri & Raphaël Lécer

HPC with a DGTD solver of Maxwell's equations

### OpenMP: strong scaling at the level of one bi-socket node



- Compact affinity
- $\blacktriangleright$  Satisfying until 8 cores Tmpi/Tomp = 0.92 for P3
- ▶ Loss when populating the second socket: a NUMA effect

### Stéphane Lanteri & Raphaël Léger

HPC with a DGTD solver of Maxwell's equations

# OpenMP: strong scaling on 32 sockets (32 mpi-processes)



- ▶ 32 processes (8 threads per process here!) compact affinity
- ▶ Communications are outside parallel regions: Amdahl effect
- ▶ 'Best MPI': 256 processes on 32 sockets Tmpi/Tomp = 0.89 for P3

# OpenMP: strong scaling at the level of one Xeon Phi KNC coprocessor



- ▶ Taking 4 threads one one core as the reference
- Satisfying behavior Tmpi/Tomp = 0.92 for P3
- ▶  $T_{Xeon Node} \sim T_{Xeon Phi Node}$  about 10% faster on Xeon

### Stéphane Lanteri & Raphaël Lécer

Context

Traditional approach: MPI

Hybrid MPI/OpenMP

Future Steps

・ロト・西ト・ヨト・ヨト・ヨー りゅつ

Stéphane Lanteri & Raphaël Lécer

4th workshop of the HOSCAR project - Gramado - 2014/09/16

HPC with a DGTD solver of Maxwell's equations

25 / 26

### Future steps

### In DEEP-er: The Hybrid MPI/OpenMP solver

- MPI/OpenMP scalability improvements
- Analysis complements with larger test-cases (a few  $10^7$  cells)
- Fine tuning on the Xeon Phi
  - Linear algebra optimization
  - Study and general improvement of vectorization performance
- Exploiting DEEP-er partners' tools
  - Task-Based parallelism and resiliency with OmpSS @ BSC
  - Parallel I/O with SionLib @ JSC and/or Exascale10 @ Xyratex-Seagate
  - Implementation of the cluster/booster application division

### MHM-DGTD for Maxwell - with F.Valentin and D.Paredes

- ▶ Incorporating this DGTD solver in the MHM framework
  - A new, flexible, numerical methodology,
  - Exposing a different structure of parallelism (master/slaves connections vs. a family of p2p connections)
- $\blacktriangleright$  A 2D demonstrator with parallelism in FORTRAN to start in October
- ▶ Provided satisfying results in 2D: a parallel 3D MHM-DGTD solver in 2015?