

BRGM at a glance

The France's leading public institution in Earth Science applications for the management of resources and surface and subsurface risks

> Key objectives

- Understanding geological processes, developing new methodologies and techniques, producing and disseminating relevant high-quality data.
- **Providing** necessary tools:
 - For surface, subsurface and resource management
 - For risk and pollution prevention
 - To support climate change policies



Research activities

> 11 main lines and 42 research programs.



- Geology
- Geothermal Energy
- Post-Mining
- Water
- Lab and experiments
- Mineral resources
- Environement and Ecotechnology
- Geological storage
 - CO2 storage and energy vectors (CH4, compressed air, hydrogen, etc.)
- Risks and impacts
 - Seismic, volcanic and tsunami risks
- Information Systems
 - Scientific Computing, software engineering and 3D.







The H-Cube project

Hydrodynamics, Heterogeneity and Homogeneization in CO2 storage modeling

http://anr-h-cube.brgm.fr



- > Financed by the French Agency of Research (ANR)
- > Call SEED 2012: Systèmes Energétiques Efficaces et Décarbonnées
- > Budget 636k€,
- > Duration: 48 months (Started in January 2013- 2016)
- > Partners:



LSCE: Laboratoire des Sciences du Climat et de l'Environnement (CEA-CNRS-UVQS)



CEREGE: Centre Européen de Recherches Et d'Enseignement (**Université Aix Marseille**)





SPE10 model

- SPE-10 model, a highly heterogeneous reservoir model (Christie et.al 2001)
- irregular nature of sand/shale distribution
- The model size is 6.4km x 8.8 km × 170m (1.2 millions of cells).



CO2 in supercritical state is injected

- with the rate of 390 k tons / year,
- with producing brine groundwater with the rate of 580 k tons / year

TOUGH2

Transport Of Unsaturated Groundwater and Heat

Developed by K. Pruess, C. Oldenburg, G. Moridis from the Earth Science Division (ESD) of Lawrence Berkeley National Laboratory (LBNL)

http://www-esd.lbl.gov/TOUGH2/

- Geothermal reservoir engineering
- Nuclear waste disposal
- Oil and gas
- Carbon storage (sequestration)



Earth Simulator 2

Upgrade of Earth Simulator system (2009) – Jamstec

- NEC SX-9/E system
- 160 nodes with 8 vector processors and 128 GB of memory
- Total of 1280 processors and 131 Tflops

Challenges on the Earth Simulator

 Opportunity based on a collaboration with research teams from Japan (Univ Tokyo, Taisei)

> Challenges

• Code initially assumed for scalar processor (Linear solver \rightarrow 70% of elapsed time) \Rightarrow Increasing of the vector operation ratio (VOR)

Local operation and no global dependency

- Continuous memory access
- Long innermost loop for vectorization
- ⇒Move from the DVBR matrix storage format to the DJDS (Saad 89, Nakajima 05)

Tough2-MP													
AZTE	C (San	dia	natio	nal La	ab.)	•	DVB High <mark>Shor</mark>	R stora ly mem t innern	ge ory eff <mark>nost lo</mark>	ficient			
HPC-	·MW	•	DJDS st <mark>Efficient I</mark> Contiguo	orage <mark>ength o</mark> f us mem	f innerr	<mark>nost l</mark> æss	оор	- In - M	creas ore th	ing ve an 30	ector o times	peratio faster	on ratio
*	NALYSIS LIS	* T											
*	NALYSIS LIS	* T *	12 (10000 7	12 505)									
* FTRACE A * Total CPU	ANALYSIS LIS Time : 3:01	* T * '40"7:	L3 (10900.7	13 sec.)									
* FTRACE A * Total CPU FREQUENCY	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](* T * '40"7: %)	L3 (10900.7 AVER.TIME [msec]	13 sec.) MOPS	MFLOPS	V.OP RATIO	AVER. V.LEN	VECTOR TIME	I-CACHE MISS	O-CACHE MISS	BANK CO CPU PORT	ONFLICT NETWORK	PROC.NAME
* FTRACE A * Total CPU FREQUENCY 1404	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(* T * '40"7 %) 90.5)	L3 (10900.7 AVER.TIME [msec] 7027.487	13 sec.) MOPS 32801.3	MFLOPS 10125.3	V.OP RATIO 99.74	AVER. V.LEN 254.4	VECTOR TIME 8650.944	I-CACHE MISS 308.097	O-CACHE MISS 335.221	BANK CC CPU PORT 2243.675	ONFLICT NETWORK 5600.790	PROC.NAME
* FTRACE A * Total CPU FREQUENCY 1404 solver_vbi 1815	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(cgstab33 450 180(* T * '40"7: %) 90.5)	L3 (10900.7 AVER.TIME [msec] 7027.487 248 033	13 sec.) MOPS 32801.3	MFLOPS 10125.3 2239 7	V.OP RATIO 99.74	AVER. V.LEN 254.4	VECTOR TIME 8650.944	I-CACHE MISS 308.097	0-CACHE MISS 335.221	BANK CC CPU PORT 2243.675 42 371	ONFLICT NETWORK 5600.790	PROC.NAME
* FTRACE # * Total CPU FREQUENCY 1404 solver_vbi 1815 1	DNALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(Ccgstab33 450.180(173.400(* T * '40"7 %) 90.5) 4.1) 1.6)	L3 (10900.7 AVER.TIME [msec] 7027.487 248.033	13 sec.) MOPS 32801.3 15795.4 268.3	MFLOPS 10125.3 2239.7 1 1	V.OP RATIO 99.74 99.24 4.75	AVER. V.LEN 254.4 240.9 134.9	VECTOR TIME 8650.944 313.245 0.279	I-CACHE MISS 308.097 0.027 27.704	O-CACHE MISS 335.221 109.621 21.770	BANK CC CPU PORT 2243.675 42.371 0.284	ONFLICT NETWORK 5600.790 165.832 0.002	PROC.NAME multi readmesh
* FTRACE # * Total CPU FREQUENCY 1404 solver_vbi 1815 1 1601	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(cgstab33 450.180(173.400(154.540(* T * '40"7: %) 90.5) 4.1) 1.6): 1.4)	L3 (10900.7 AVER.TIME [msec] 7027.487 248.033 L73399.760 96.527	13 sec.) MOPS 32801.3 15795.4 268.3 9063.7	MFLOPS 10125.3 2239.7 1.1 4299.5	V.OP RATIO 99.74 99.24 4.75 95.33	AVER. V.LEN 254.4 240.9 134.9 255.4	VECTOR TIME 8650.944 313.245 0.279 37.474	I-CACHE MISS 308.097 0.027 27.704 0.223	O-CACHE MISS 335.221 109.621 21.770 9.599	BANK CC CPU PORT 2243.675 42.371 0.284 2.184	ONFLICT NETWORK 5600.790 165.832 0.002 9.855	PROC.NAME multi readmesh eos
* FTRACE # * Total CPU FREQUENCY 1404 solver_vbi 1815 1 1601 2	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(Cgstab33 450.180(173.400(154.540(95.918(* T * '40"7: %) 90.5) 4.1) 1.6): 1.4) 0.9)	L3 (10900.7 AVER.TIME [msec] 7027.487 248.033 L73399.760 96.527 47958.886	13 sec.) MOPS 32801.3 15795.4 268.3 9063.7 203.0	MFLOPS 10125.3 2239.7 1.1 4299.5 5.6	V.OP RATIO 99.74 99.24 4.75 95.33 0.66	AVER. V.LEN 254.4 240.9 134.9 255.4 12.9	VECTOR TIME 8650.944 313.245 0.279 37.474 0.317	I-CACHE MISS 308.097 0.027 27.704 0.223 20.177	O-CACHE MISS 335.221 109.621 21.770 9.599 10.810	BANK CC CPU PORT 2243.675 42.371 0.284 2.184 0.073	ONFLICT NETWORK 5600.790 165.832 0.002 9.855 0.303	PROC.NAME multi readmesh eos wrifi
* FTRACE A * Total CPU FREQUENCY 1404 solver_vbi 1815 1 1601 2 1	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(cgstab33 450.180(173.400(154.540(95.918(62.822(* T * '40"7 %) 90.5) 4.1) 1.6) 1.4) 0.9) 0.6)	L3 (10900.7 AVER.TIME [msec] 7027.487 248.033 L73399.760 96.527 47958.886 62822.241	13 sec.) MOPS 32801.3 15795.4 268.3 9063.7 203.0 240.8	MFLOPS 10125.3 2239.7 1.1 4299.5 5.6 2.0	V.OP RATIO 99.74 99.24 4.75 95.33 0.66 2.14	AVER. V.LEN 254.4 240.9 134.9 255.4 12.9 47.7	VECTOR TIME 8650.944 313.245 0.279 37.474 0.317 0.411	I-CACHE MISS 308.097 0.027 27.704 0.223 20.177 9.985	O-CACHE MISS 335.221 109.621 21.770 9.599 10.810 10.375	BANK CC CPU PORT 2243.675 42.371 0.284 2.184 0.073 0.048	ONFLICT NETWORK 5600.790 165.832 0.002 9.855 0.303 0.314	PROC.NAME multi readmesh eos wrifi rfile
* FTRACE # * Total CPU FREQUENCY 1404 solver_vbi 1815 1 1601 2 1 1	NALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(cgstab33 450.180(173.400(154.540(95.918(62.822(56.596(* T * '40"7 %) 90.5) 4.1) 1.6) 1.4) 0.9) 0.6) 0.5)	AVER.TIME [msec] 7027.487 248.033 173399.760 96.527 47958.886 62822.241 56596.402	13 sec.) MOPS 32801.3 15795.4 268.3 9063.7 203.0 240.8 373.9	MFLOPS 10125.3 2239.7 1.1 4299.5 5.6 2.0 10.2	V.OP RATIO 99.74 99.24 4.75 95.33 0.66 2.14 67.68	AVER. V.LEN 254.4 240.9 134.9 255.4 12.9 47.7 67.5	VECTOR TIME 8650.944 313.245 0.279 37.474 0.317 0.411 22.890	I-CACHE MISS 308.097 0.027 27.704 0.223 20.177 9.985 0.161	O-CACHE MISS 335.221 109.621 21.770 9.599 10.810 10.375 2.993	BANK CC CPU PORT 2243.675 42.371 0.284 2.184 0.073 0.048 0.057	ONFLICT NETWORK 5600.790 165.832 0.002 9.855 0.303 0.314 19.280	PROC.NAME multi readmesh eos wrifi rfile cycit
* FTRACE # * Total CPU FREQUENCY 1404 solver_vbi 1815 1 1601 2 1 1 2	ANALYSIS LIS Time : 3:01 EXCLUSIVE TIME[sec](9866.592(9866.592(173.400(154.540(95.918(62.822(56.596(8.345(* T * '40"7 %) 90.5) 4.1) 1.6) 1.4) 0.9) 0.6) 0.5) 0.1)	L3 (10900.7 AVER.TIME [msec] 7027.487 248.033 L73399.760 96.527 47958.886 62822.241 56596.402 4172.604	13 sec.) MOPS 32801.3 15795.4 268.3 9063.7 203.0 240.8 373.9 190.5	MFLOPS 10125.3 2239.7 1.1 4299.5 5.6 2.0 10.2 3.1	V.OP RATIO 99.74 99.24 4.75 95.33 0.66 2.14 67.68 0.42	AVER. V.LEN 254.4 240.9 134.9 255.4 12.9 47.7 67.5 79.7	VECTOR TIME 8650.944 313.245 0.279 37.474 0.317 0.411 22.890 0.007	I-CACHE MISS 308.097 0.027 27.704 0.223 20.177 9.985 0.161 1.263	O-CACHE MISS 335.221 109.621 21.770 9.599 10.810 10.375 2.993 1.571	BANK CC CPU PORT 2243.675 42.371 0.284 2.184 0.073 0.048 0.057 0.007	ONFLICT NETWORK 5600.790 165.832 0.002 9.855 0.303 0.314 19.280 0.006	PROC.NAME multi readmesh eos wrifi rfile cycit finalout

E.



Dumux

- > Developed by University of Stuttgart
- > Multiphase, multicomponent flow and transport in porous media

> <u>http://www.dumux.org/</u>

- > Incompressible, immiscible flow
 - Isothermal condition \rightarrow TwoP model (2p)
- > Compressible, Miscible flow
 - Isothermal condition → TwoPTwoC model (2p2c)

Dumux

- Speedup = 4.01 (16)
- Speedup = 4.96 (16)
- Impact of linear solver



	Implicit-Ya	nspGrid	Implicit-ALUGrid		
	CPU Time [s]	Time step	CPU Time [s]	Time step	
24 cores	3934.94	76	4587.27	73	
48 cores	2175.51	73	2503.85	73	
96 cores	1668.16	76	1296.53	73	
192 cores	1169.24	73	1116.46	79	
384 cores	981.37	73	923.626	75	







data transfer and load imbalance







Charm++/AMPI overdecomposes the domain into virtual processors (VPs)

- Multiple VPs per core / Each VP implements one MPI task



Motivation \rightarrow manycores architecture

MPPA-256 from Kalray

- 2 MB on compute cluster 4 GB on I/O nodes
- Theoretical 230 Gflops 5 W



Motivation \rightarrow manycores architecture

Intel Xeon Phi

- 57 cores 512 KB of L2 Cache
- 6 GB of shared DDR memory / 2.0 TFLOPS in SP



Seismic stencil - algorithm



Results



Time-to-solution

- Ratio GPU/Phi \rightarrow Peak = 1.75 Measured = 1.20
- Ratio MPPA/Phi \rightarrow Peak = 14 Measured = 3.68

Motivation

- > Reuse existing CPU and GPU numerical kernels
- > Evaluate task programming model.
- > Tackle Heterogeneous architectures (CPU + GPU / CPU + Xeon Phi)



Ondes3D on top of StarPU

StarPU programming model

- > Task programming model
- > Directed Acyclic Graph according to data dependencies
- > StarPU manages memory transfers .

> For each iteration and each block :

- Update source
- Load stress boundaries & Compute velocity from stress
- Save & send velocity boundaries for each direction
- Load velocity boundaries & Compute stress from velocity
- Save & send stress boundaries for each direction
- Record seismograms

Scheduling strategies



Dependencies in the DAG for a kernel on a given block (here in a grid of 3x3 blocks) Colored tasks can not be scheduled before the initial one is over.

- If a task take longer to compute on a slow PU than the rest of the grid on a fast PU, there will be waiting for the fast PU on the next step
- Time taken to compute a task (here a given kernel) on a slow PU should not exceed the time taken to compute the rest of the grid on faster PU for the same kernel



Incore results



- > Cost of data transfer using the PCI bus
- > Hybrid (DMDAR) vs Hybrid (WS)

Out-of-core results



Impact of data transfer for incore and ooc experiments (pure GPU)
 We are far from peak performance level (i.e. speedup with incore GPU)

Impact \rightarrow Scheduling strategies



> Contention on the memory bus.

- Strong impact on the HPC node \rightarrow 8 GPU + 4 CPU
- Reduced on commodity-node \rightarrow 1 GPU + 3 CPU

Impact \rightarrow Size of the block



> Competition between two behaviors with a fixed problem size

- Create enough (but not too much) tasks for parallelism at StarPU level
- Suitable block size with respect to the underlying architecture.

Stencil computation – space blocking > Standard methodology

- > 2D blocking and reuse in the third direction (*Rivera 2000*)
- > Poor reuse opportunity \rightarrow Blas 1 formulation O(1)



Stencil computation – spacetime blocking



- Reuse data accros several time steps
- Wavefront algorithm
- Blas 3 formulation \rightarrow O(n)

Characterization

How do threads communicate.

- Amount of communication
- Heterogeneity (Communication Pattern)

SPCD: shared pages and communication detection



Characterization

Exclusivity:

highest memory accesses from a node

- # memory accesses from all nodes
 - Max: 100% (page fully exclusive)
 - Min: 1/#nodes (page fully shared)
 → with 4 nodes, min: 25%

Treemap



7-points \rightarrow *Impact* of the « reuse » parameter



Tradeoff between optimal size of the tiles and « reuse » opportunity
 Shift in the upper-bound limit with respect to the number of cores
 Better usage of memory BW available

Seismic wave equation							
	Platform	Spaceblocking	Spacetime blocking				
	Intel Xeon E5-2697v2	x1.19	x3.59				
	Intel Xeon E5-2697v3	x1.27	x2.68				
	 Limited speedup from space only blocking Stronger impact of vectorization on Haswell vs complexity of the kernel. 						

Variability of earthquakes ground motions at sedimentary basin scale (1D)

⇒Uncertain sources: wave velocities, thickness of layers, damping ratio
 ⇒Quantities of interest : PGA, acceleration, spectral ratio

> Meta-modeling:

- Polynomial regression
- Neural network
- Kriging (or Gaussian process)
- Chaos Polynomial: convergence in case of smooth process

> Chaos Polynomial (spectral expansion)

- Need to solve some direct models (Design of experiments)
- Low number of evaluations in comparison with Monte-Carlo
- Evaluation of the polynomial coefficents by projection
- > Difficulties: Stochastic stiffness of the quantities of interest
 - Limited number of simulation (in 3D)

Variability of earthquakes ground motions at sedimentary basin scale (1D)







Scientific workflow

Large Data sets

- o LiDAR: Digital Elevation Model
 - \rightarrow French coastal area: 160Go, 25000 files
 - \rightarrow Several Terabytes for global map
- Data on tides and water levels (SHOM National Hydrographic Service)
- Population and land use (Medde, INSEE, EEA)
- Extrem events and coastal erosion (BRGM)







- > Take full advantage of large scale processing facilities in order to speedup gridbased data analysis including multi-resolution strategies (**Rasdaman**)
- Increase the complexity of overflow model including PDE-based simulation tools (Swash, Swan).
- Extend this framework capabilities to tackle large environmental data collections (EPOS, Copernicus ...)
- > Application to Earthquake early warning → exploiting GPU implementation

Conclusion

> Harnessing high performance computing is critical for future challenges in Geosciences

Large scale, 3D, coupled simulations associated with uncertainties analysis

> Effort at several levels at BRGM

- Balance between internal efforts and external collaborations at the software level.
- Maintain high level scientific collaborations

> New challenges \rightarrow mostly related to data !

- Geoprocessing \rightarrow High Performance Data
- Uncertainties (non intrusive) \rightarrow all aspects including visualization.
- Moving code paradigm and complex high-performance workflows.

> Happy with the collaboration with Brazil

- Ongoing work on I/O (Efispec3D software)
- Co-advising the Ph.D of Victor Martinez

References

> Journal

- R.Keller Tesser, L.Pilla, F.Dupros, P.O.Navaux, J-F Mehaut, C.Mendes : Dynamic load balancing for seismic wave propagation models IJHPCA (revision)
- M.Castro, F.Dupros, E.Francesquini, J-F Méhaut, P.O.Navaux : Seismic Wave Propagation Simulations on Lowpower and Performance-centric Manycores. *Parallel Computing (revision)*
- M.Diener, E.Cruz, L.Pilla, F.Dupros, P.O.Navaux : Characterizing Communication and Page Usage of Parallel Applications for Thread and Data Mapping . Performance Evaluation (2015)
- M.Castro, E.Francesquini, P-H Penna, F.Dupros, H.Freitas, J-F Méhaut, P.O.Navaux : On the Energy Efficiency and Performance of Irregular Applications on Multicore, NUMA and Manycore Processors : JPDC (2014)

> Conferences

- > V.Martinez, F.Dupros, D.Michea, O.Aumage, S.Thibaut, P.O.A.Navaux : Towards seismic wave modeling on heterogeneous many-core architectures using task-based runtime system. SBAC-PAD (2015)
- > F. Dupros, F.Boulahya, H.Aochi, P.Thierry : Communication-avoiding seismic numerical kernels on multicore processors: HPCC (2015)
- M.Castro, F.Dupros, E.Francesquini, J-F Méhaut, P.Navaux : Energy Efficient Seismic Wave Propagation Simulation on a Low-Power Manycore Processor. Proceedings of SBAC-PAD (2014) – Best Paper
- R.Keller Tesser, L.Pilla, F.Dupros, P.O.Navaux, J-F Mehaut, C.Mendes : Improving the performance of seismic wave simulations with dynamic load balancing : Proceedings of PDP (2013)
- Hajime Yamamoto et.al. : Numerical Simulation of Long-Term Fate of CO₂ Stored in Deep Reservoir Rocks on Massively Parallel Vector Supercomputer. Proceedings of VECPAR (2012)

