

Alya - CFD on exascale GPU hardware for the Wind community.

Barcelona Supercomputing Center (BSC) - FAU University

ISC 2022

Herbert OWEN and Dominik ERNST

with help from Filippo Spiga(Nvidia) & the PSC Toolkit





Alya - Wind energy - Energy oriented Center of Excellence



Wind Energy

LES Modelling of wind flow at Bolund hill

Large Eddy Simulation of wind flow at the Bolund hill, one of the best known benchmarks for complex terrain. The visualisation shows three distinct parts of the overall computational domain, where we visualise the turbulent flow structures through a volumetric rendering of the speed (0 to 8 m/s).



Wind is a mature technology but plant-level energy losses —> 20% on flat terrain. energy cheaper, and enable widespread adoption. of scales.

wind)

Today, to Out-Compute is to Out- Compete. Europe needs to react accordingly.

- A better understanding of the turbulent flow is needed to reduce these losses, make wind
- Exascale simulations will enable an unprecedented understanding of a wind farm's wide range

The US is investing 6 million dollars per year in the A2e and Exawind projects. (20X EoCoE-



HPC-based multi-physics simulation code Developed at BSC to run efficiently in parallel supercomputers.

Scalability has been tested up to 100K cores. Part of UEABS - Unified European Bench Suite for CPUs & GPUs..



Numerical Model: Alya - LES

EMA - Energy, momentum and angular momentum conserving convective term. Stabilisation for the p-v interaction coming from Laplacian approximation in Fractional Step Method. Physical based SGS modelling (Vreman in current work). SIMPLE and no user defined numerical parameters.

Lehmkuhl el al. A low-dissipation finite element scheme for scale resolving simulations of turbulent flows J. Comput. Phys., 308:51–65, 2019

- Galerkin discretisation with explicit (RK3/4) treatment of convective and diffusive terms.





Numerical Model: Alya - LES

Test case: Taylor-Green vortex Re = 1600 *



A low-dissipation finite element scheme for scale resolving simulations of turbulent flows, Lehmkuhl et al. submitted to Journal of Computational Physics For reference see: Comparison between several approaches to simulate the Taylor-Green vortex case, Moulinec et al. PARCFD 2016



t = 5



t = 20





FRACTIONAL STEP SCHEME - with Runge Scheme time discretization

For each substep: Obtain F^{J} assembly process — close to 80% on typical run on MN4 $U^{i,*} = U^n + \delta t M^{-1} \sum_{i=1}^{i} a_{ij} F^j,$ Multiply by inverse of Lumped Mass Matrix (cheap) Multiply by Divergence Matrix — D Solve Linear System for Laplacian (2nd most costly step) $\mathsf{D}\mathsf{M}^{-1}\mathsf{G}\left(\Phi^{i}\right) = \frac{1}{c_{i}\delta t}\left(\mathsf{D}\mathsf{U}^{i,*} - \mathsf{R}_{\mathsf{c}}\right),$ $\mathsf{U}^{i} = \mathsf{U}^{i,*} - c_{i}\delta t\mathsf{M}^{-1}\mathsf{G}\left(\Phi^{i}\right),$ Multiply by Gradient Matrix — G Multiply by inverse of Lumped Mass Matrix (cheap)

Other Operations : multiply by scalars; add vectors; apply dirichlet BCs

A low-dissipation finite element scheme for scale resolving simulations of turbulent flows, O. Lehmkuhl, G. Houzeaux, H. Owen, G. Chrysokentis, I. Rodriguez, JCP 2018





GPU implementation and optimization



NAME	time from .log	improvement	Clasification	Work done
GPU2_17	39,61		code reformats	Alya from master
40	33,64	1,18	code reformats	Reincorporate non optimized miniapp
41	19,64	1,71	algorithmic changes	Incorporate optimized miniapp
42-43	18,83	1,04	code reformats	Switch to just P1 - Eliminate call to ke
44-45	4,35	4,33	privatization	all thread private
46-49	3,98	1,09	code reformats	Cleaning- add and use vreman-
51		3,50	code reformats	avoid elmat - obtain directly elrbu - in
			amail E nov 21	



Times for assembly (& properties)

	Initial - 2	Final - 80CPU / 91 GPU	Improvement
1A100 GPU + 9mpi	9.8 s	0.07 s	140 X
MN4 46 mpi Intel Xeon Platinum 8160 - Skylake	1.93 s	0.49 s	3.9 X
	CPU 5 X faster	GPU 7 X faster	
Maximum assembly	time obtained fro	om Alya's .log file	



and ents



Energy estimates from Top 500

The GPU node we have used is very similar to a Perlmutter node From top 500 - 2589KW/1536nodes - 1685W per node We use approx. I/4 node —> 421 W x 0.07s - 29.5

....for I MN-IV node $\longrightarrow 511W \times 0.49s - 250J$

Our assembly is 8.5 times more efficient on the GPU!! 30 (8.5x3.9) times more efficient than our current CPU implementation!!

From Green 500: Perlmutter 27.374 GFlops/watts MN-IV 3.965 GFlops/watts

Power Efficiency Ratio 6.90 - We are getting 8.5!!!!





Optimization

Thread private	e instead of I:VECTOR_SIZE	
real(rp) real(rp)	<pre>:: elvel(VECTOR_DIM,ndime,pnode) :: elcod(VECTOR_DIM,ndime,pnode)</pre>	! u ! x
real(rp) real(rp)	<pre>:: gpvel(VECTOR_DIM,ndime,pgaus) :: gpgve(VECTOR_DIM,ndime,ndime,pgaus)</pre>	! u ! grad(u
gprhs(DEF_VECT,	<pre>idime,igaus) = gprhs(DEF_VECT,idime,igaus) + FACT2X * gravi_ </pre>	_nsi(idime) & _VECT_idime_idaus

- Specialization I: Only tetrahedra: FE Shape functions and its derivatives easier to calculate. Gradient homogeneous for all gauss points -> only one needed
- Specialization 2: Obtain elemental RHS directly only valid for explicit.
- Specialization 3: Turbulent viscosity calculated on the fly nearly for free.
- Pnode, Pgaus, Mnode as Fortran parameters
- Reduction of registries. Better memory use.







NSYS profile — final version 9 mpi

≡ Timeline View 🔹											乓 1x		▲ <u>10</u>
58s •	+950ms 59s	+50ms	+100ms +150ms	+200ms	+250ms	+300ms	+350ms	+400ms	+450ms	+500ms	+550ms	+600ms	+650m
CPU (128)													
Processes (19)													
[1228324] Alya.x													
CUDA HW (0000:41:00.0 -	ter eine an eine eine						in line						
 Threads (9) 													
▼ ✓ [1228324] Alya.x •													
OpenACC	(ığınınığırididididi di≌⊞⊞d+ I-dhdH41111-d	M4MB4M4M4M444M444AM444AM444AM444AM444AM	mandijinaandaand jaaraadaanadaanad				199 9 999999999999999999999999999999999	ակակիվիվի իս(իվիկի՝ իվսփասվափ		
CUDA API													
Profiler overhead													
✓ [1230758] cuda-Evtl-													
7 threads hidden – +													
18 processes hidden – +													
	31ms						20ms						

3 Runge Kutta substesp —> total 71ms matches 0.07s from previous slide

The total GPU compute time from ncu (nsight compute) is 50ms with 9 mpi overhead is only 21 ms





NSYS profile — final version 9 mpi



- Blend between compute and copyin/out
- Using 9 MPIs allows to hide part of the copyin/out
- Copyin/out y mainly velocity and rhsid once we have everything on the GPU those times will disappear!!!

						
+341ms +342ms +3	343ms +344ms	+345ms -	+346ms -	+347ms +348	ms +349ms	+350
						<u> </u>
a : mod_nsi_element_operations_hh)	Compute Construct : mod_nsi_	element_oper		Update :	nsi_elmope_all.f90:38	38
d_nsi_element_operations_hh91.f9_	Wait : mod_nsi_element_opera	ations_hh91.f9_				Wait : ns
cuStreamSynchronize	cuStreamSynchron	nize	(cuEventSynchronize		cuStream.

2nd substep 20ms



NSYS profile — final version serial



- Total 235ms much slower than with 9mpi (71ms) copyin/out become dominant
- In any case all copyin/out will disappear when we port 'the rest' to GPU The rest of the code will be much faster on the GPU

NSYS profile — original version serial



'The rest' looked much smaller because GPU computation was very inefficient

What is being done there?

Multiplication by Gradient and Divergence Matrices Solve Linear System for Laplacian (2nd most costly step)

Mainly:

Nsight Compute

38% of FP64 Peak — before 0%

NVIDIA Nsight Compute

The ratio of peak float (fp32) to double (fp64) performance on this device is 2:1. The kernel achieved 0% of this device's fp32 peak performance. If <u>Compute Workload Analysis</u> determines that this kernel is fp64 bound, consider using 32-bit precision floating point operations to improve its performance.

Original

Roofline Analysis

The ratio of peak float (fp32) to double (fp64) performance on this device is 2:1. The kernel achieved 0% of this device's fp32 peak performance and 0% of its fp64 peak performance.

Explanation, that I nearly believed Low order FE not suited for GPUs (still trying to find papers) High order, Finite Difference and LBM - YES



Alya - linear solver for the pressure



Journal of Computational Science Volume 14, May 2016, Pages 15-27



Alya: Multiphysics engineering simulation toward exascale

-Mariano Vázquez ^{a, b} 🕾 🖾, Guillaume Houzeaux ^a, Seid Koric ^c 🖾, Antoni Artigues ^a, Jazmin Aguado-Sierra ^a, Ruth Arís^a, Daniel Mira^a, Hadrien Calmet^a, Fernando Cucchietti^a, Herbert Owen^a, Ahmed Taha^c, Evan Dering Burness ^c, José María Cela ^a, Mateo Valero ^a



Kiln furnace. Convergence of the DCG solver for the continuity equation





Computational Mechanics solvers have two main steps:

- I) Assembly of a Matrix & RHS
- 2) Solution of a Linear System

For the solution of the linear System Alya has used in house developed code until recently.

Main solvers:

- GMRES
- •CG
- Deflated CG



AMG4PSBLAS

PSCToolkit (Parallel Sparse Computation Toolkit)



A software framework for scalable solution of sparse linear systems by Krylov methods coupled with Algebraic Multigrid (AMG) Preconditioners on hybrid CPU/NVIDIA-GPU architectures

Available at PSCTOOLKIT | Parallel Sparse Computation Toolkit

Main reference: P. D'Ambra, F. Durastante, S. Filippone, AMG Preconditioners for Linear Solvers towards Extreme Scale, to appear. Preprint available at [2006.16147] (arxiv.org)



Bolund case



Original mesh with 5.6 M nodes and 32 M elements

Apply divisor up to 3 times (X 512) 2850M nodes **I6000M** elements



Guillaume Houzeaux, Raúl de la Cruz, Herbert Owen, and Mariano Vázquez. Parallel uniform mesh multiplication applied to a Navier- Stokes solver. Computers and Fluids





Excellent weak and algorithmic scalability from original mesh to 3 divisors

	cores	average iterations
Original mesh	46	6
1 divisor	368	6
2 divisors	2944	4.8
3 divisors	23552	4

time per iteration		a	v. time [m	s] av. time w/ Alya CG
	133		800	1568
	134		802	1198
	150		716	1939
	159		635	4658



Conclusions and next step

- •Low order FE can run efficiently on GPUs !!!! •Algebraic multigrid - weak scalability up to 16000 million elements

- •Put the whole fractional step on the GPU.
- •Generalise (other elements, richer physics) without losing too much efficiency.

•If Europe does not want to be out-computed/out-competed it must react soon.



Thanks for your attention!



The author thankfully acknowledges the computer resources at MareNostrum and the technical support provided by Barcelona Supercomputing Center (RES-AECT-IM-2021-2-0011)



