

Alya towards Exascale: efficient finite element assembly on GPUs for LES

Herbert Owen¹, Oriol Lehmkuhl², Guillaume Houzeaux³, Guillermo Oyarzun⁴, Georg Hager⁵, Gerhard Wellein⁶ and Dominik Ernst³

¹ Barcelona Supercomputing Center, herbert.owen@bsc.es

² Barcelona Supercomputing Center, oriol.lehmkuhl@bsc.es

³ Barcelona Supercomputing Center, guillaume.houzeaux@bsc.es

⁴ Barcelona Supercomputing Center, guillermo.oyarzun@bsc.es

⁵ Erlangen National High Performance Computing Center (NHR@FAU), georg.hager@fau.de

⁶ Erlangen National High Performance Computing Center (NHR@FAU), gerhard.wellein@fau.de

⁷ Erlangen National High Performance Computing Center (NHR@FAU), dominik.ernst@fau.de

Key Words: *Finite element, CFD, LES, Exascale, GPU.*

On the path to Exascale, most codes will need significant changes to adapt to new architectures. Currently, GPUs are the most mature among such architectures. This work describes the improvements we have recently introduced in the High-Performance Computing code Alya within the EoCoE-II project. We shall focus on incompressible flow problems for Large Eddy Simulation. The momentum equation is treated explicitly while the pressure is solved implicitly. A fractional step scheme is used to enable elements that do not satisfy the inf-sup condition [1].

Except in cases with mesh deformation, the Laplacian matrix for the pressure remains fixed during the whole simulation. Therefore, the two main kernels are calculating the right-hand side term for the momentum equation and the solution of a linear system for the pressure. In this work, we shall concentrate on the improvements we have obtained in the GPU implementation of the first kernel. Some small comments on the algebraic multigrid strategy will also be presented.

The simple steps that have been followed to obtain more than an order of magnitude reduction in computational time for the right-hand side term assembly compared to our previous GPU implementation [2] will be presented by a finite element CFD specialist. The GPU implementation is based on OpenACC, and it is therefore quite friendly. However, it has allowed us to reach 50% of the maximum floating-point performance on an A100 Nvidia GPU. Finally, while optimizing the GPU version, improvements to the CPU implementation have also been obtained. A comparison between the current CPU and GPU performances will show that we are more energy efficient on the GPU than on the CPU nowadays.

REFERENCES

- [1] O. Lehmkuhl, G. Houzeaux, H. Owen, G. Chrysokentis and I. Rodriguez A low-dissipation finite element scheme for scale resolving simulations of turbulent flows, *Journal of Computational Physics*, 2019.
- [2] R. Borrell et al., Heterogeneous CPU/GPU co-execution of CFD simulations on the POWER9 architecture: Application to airplane aerodynamics, *Future Generation Computer Systems*, 2020.