



## A distributed-heterogeneous framework for of explicit hyperbolic solvers of shallow-water equations

Rui M L Ferreira<sup>1</sup>, **Daniel Conde**<sup>2</sup>, Zixiong Zhao<sup>3</sup>, and Peng Hu<sup>4</sup>

<sup>1</sup>Instituto Superior Tecnico, Universidade de Lisboa, CERIS, Lisboa, Portugal (ruimferreira@tecnico.ulisboa.pt)

<sup>2</sup>CERIS - Civil Engineering Research and Innovation for Sustainability, Lisboa, Portugal (daniel.conde@tecnico.ulisboa.pt)

<sup>3</sup>Zhejiang University, Hangzhou, P.R. China, visiting at CERIS, Lisboa, Portugal (187477@zju.edu.cn)

<sup>4</sup>Zhejiang University, Hangzhou, P.R. China (penghu@zju.edu.cn)

This work addresses current performance limitations by introducing new distributed multi-architecture design approaches for massively parallel hyperbolic solvers. Previous successful implementations that couple MPI with either OpenMP or CUDA have been previously reported in the literature. We present novel approaches that remain intuitive and compatible with developer-centered object-oriented (OOP) paradigm but coupled with a cache-conscious data layout, compatible with both structured and unstructured meshes, promoting memory efficiency and quasi-linear scalability.

One of the approaches is based on a unified object-oriented CPU+GPU framework, augmented with an inter-device communication layer, enabling both coarse and finegrain parallelism on hyperbolic solvers. The framework is implemented through a combination of three different programming models, namely OpenMP, CUDA and MPI. The second approach is also based on a unified object-oriented CPU+GPU framework, augmented with an improved local time step algorithm (LTS) on variable updating. This framework is implemented through a combination of parallel technology (CUDA Fortran) and mathematical algorithm (TLS).

The efficiency of these distributed-heterogeneous frameworks are quantified under static and dynamic loads on consumer and professional grade CPUs and GPUs. In both approaches, an asynchronous communications scheme is implemented and described, showing very reduced overheads and a nearly linear scalability for multiple device combinations. For simulations (or systems) with non-homogeneous workloads (or devices) the domain decomposition algorithm incorporates a low-frequency load-to-device fitting function to ensure computational balance. Real-world applications to high-resolution shallow-water problems are presented. The proposed implementations show speedups of up to two orders of magnitude, opening new perspectives for shallow-water solvers with high-demand requirements.