

EXASCALE-READY CROSS-DOMAIN WORKFLOWS EXECUTION: THE ACROSS APPROACH

Alberto Scionti, Olivier Terzo, Giacomo Vitali, Paolo Viviani, Chiara Vercellino, Paolo Savio, Klodiana Goga, Giuseppe Caragnano, Fabrizio Bertone.

LINKS Foundation, via P. C. Boggio 61, 10138 Turin (Italy)

Abstract. The boundaries among Machine Learning (ML), Deep Learning (DL) –collectively referred to as AI-domain, Big-Data (BD) and traditional numerical simulation (HPC) computing domains are ever less prominent. The tremendous technology advancement of the last decade is pushing these domains to lay on a common background, as represented by the integration of FPGAs, GPUs, and other AI domain-specific devices into recent supercomputers. Industrial and scientific applications see the usage of these devices as a way to increase the capability of processing larger input data sets, and to boost the overall performance while reducing the energy consumption. This all generated big challenges in the corresponding software stack, now responsible for managing the job submission process while being aware of the specificities of the diverse cross-domain tasks and the availability of different accelerators. Challenges arise by the inadequacy of the traditional *Resource and Job Management Software* (RJMS), which are limited by their heuristic priority functions that, despite their sophistication, are not able to well adapt to the changes in the workloads over time. Also, the current generation of RJMS are poorly-aware of the growing level of heterogeneity in modern HPC clusters. This is obviously exacerbated in the upcoming Exascale machines, where even more pressure on the RJMS and applications in terms of flexibility and performance exists. The ACROSS project (funded by EU under the H2020 and EuroHP-JU frameworks) aims to address such technological challenges on one hand, while reinforcing the use of many types of hardware accelerators to boost application performance and scalability. As such, ACROSS targets the implementation of an execution platform that is based on two main pillars: *i*) the integration and exploitation of a wide range of hardware accelerators complemented by diverse CPU architectures (Intel Xeon and AMD Epyc), which include GPUs, FPGAs and AI-specific (e.g., Habana Goya chips, VPUs); *ii*) the design and integration of a multi-level cross-domain aware RJMS. Moreover the project is exploring the integration of even more innovative architectures, such as neuromorphic architectures, within the devised platform. To this end, the project consortium is investigating the efficient implementation of a spiking neural network (SNNs) accelerator as an overlay of high-performance FPGAs, and to test it in one of the use cases (aeronautic domain). Interestingly, the accompanying software stack will be integrated in the designed RJMS, enabling pilot users to seamlessly integrate SNN models in their existing workflows without significant effort overhead. The second innovation will come from a multi-layer, cross-domain task, workflow-aware RJMS. It will explore advanced mechanisms to acquire (heterogeneous, including Cloud) resources and to plan the execution of workflows, by adapting to the changes over time of their composition. A fine control of the available resources at the lowest level of the sstack will improve the task scheduling and increase the overall energy efficiency. Dedicated modules will ease the integration and management of ML/DL models, as well as getting access to the neuromorphic hardware. The validation of such an ambitious platform will be performed through pilots covering highly resource-demanding domains: *i*) aeronautic; *ii*) climate and weather prediction; and *iii*) energy and carbon sequestration.

Main Author. Alberto Scionti holds a MSc and a PhD (European doctorate degree) in computer science and control engineering, both received from Politecnico di Torino, Italy. He is a senior researcher in the Advanced Computing, Photonics and Electromagnetics group at LINKS Foundation. His main research interests include management of heterogeneous HPC/Cloud infrastructures and quantum computing. Other research interests are in the field of evolutionary algorithms, as well as in the design of FPGA-based accelerators. He was involved in several national and EU-funded projects. He is co-author of more than 50 peer reviewed papers on international conferences and journals. He was editor of a book on heterogeneous computing architectures and is currently leading the orchestrator design in the EU-H2020/EuroHPC-JU ACROSS project.