

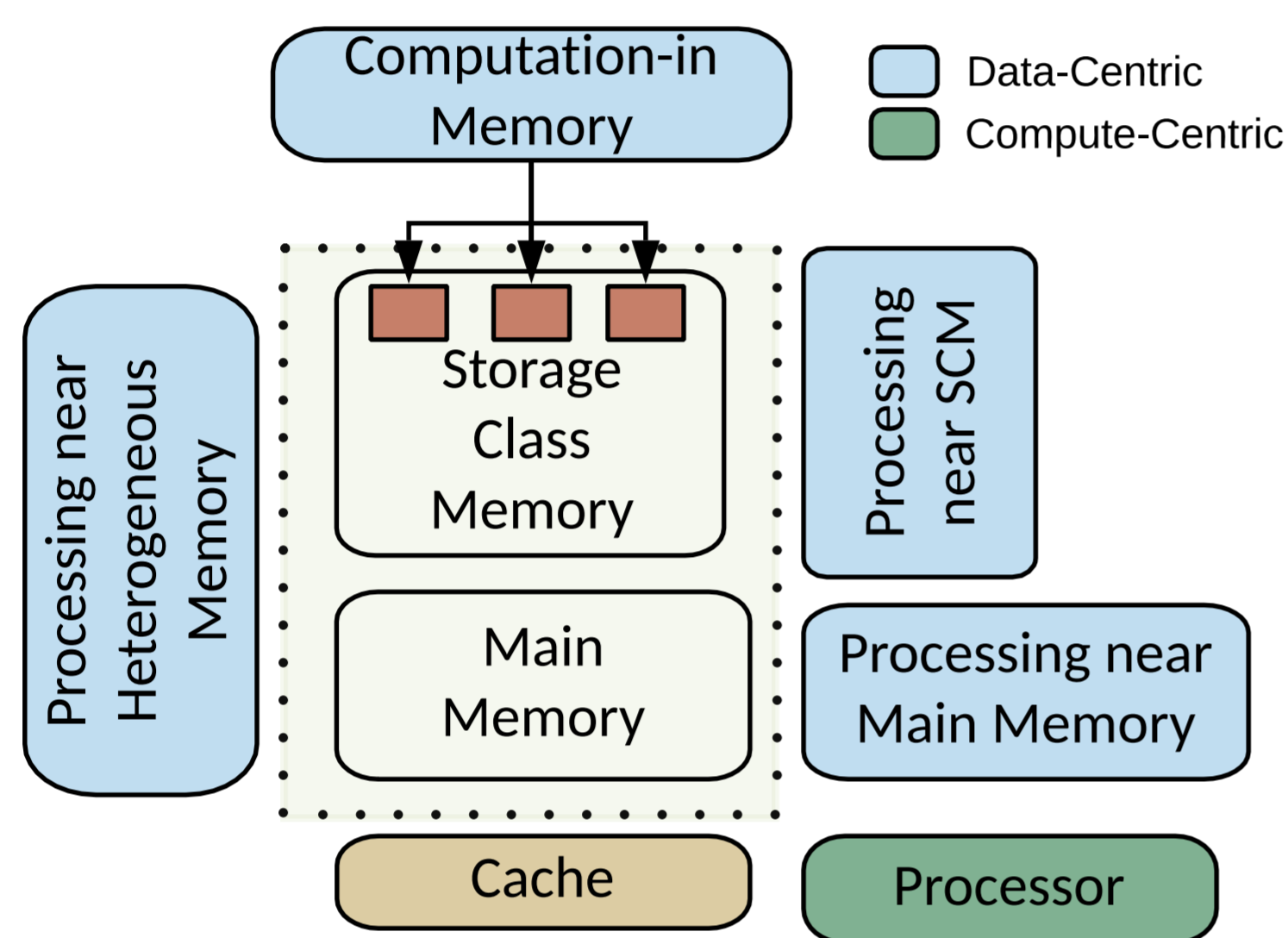
# Scaling Stencil Computation on OpenPOWER Near-Memory Architecture

Gagandeep Singh, Dionysios Diamantopoulos, Sander Stuijk, Henk Corporaal, Christoph Hagleitner

Funded by the Horizon 2020 Framework Programme of the European Union MSCA-ITN-EID

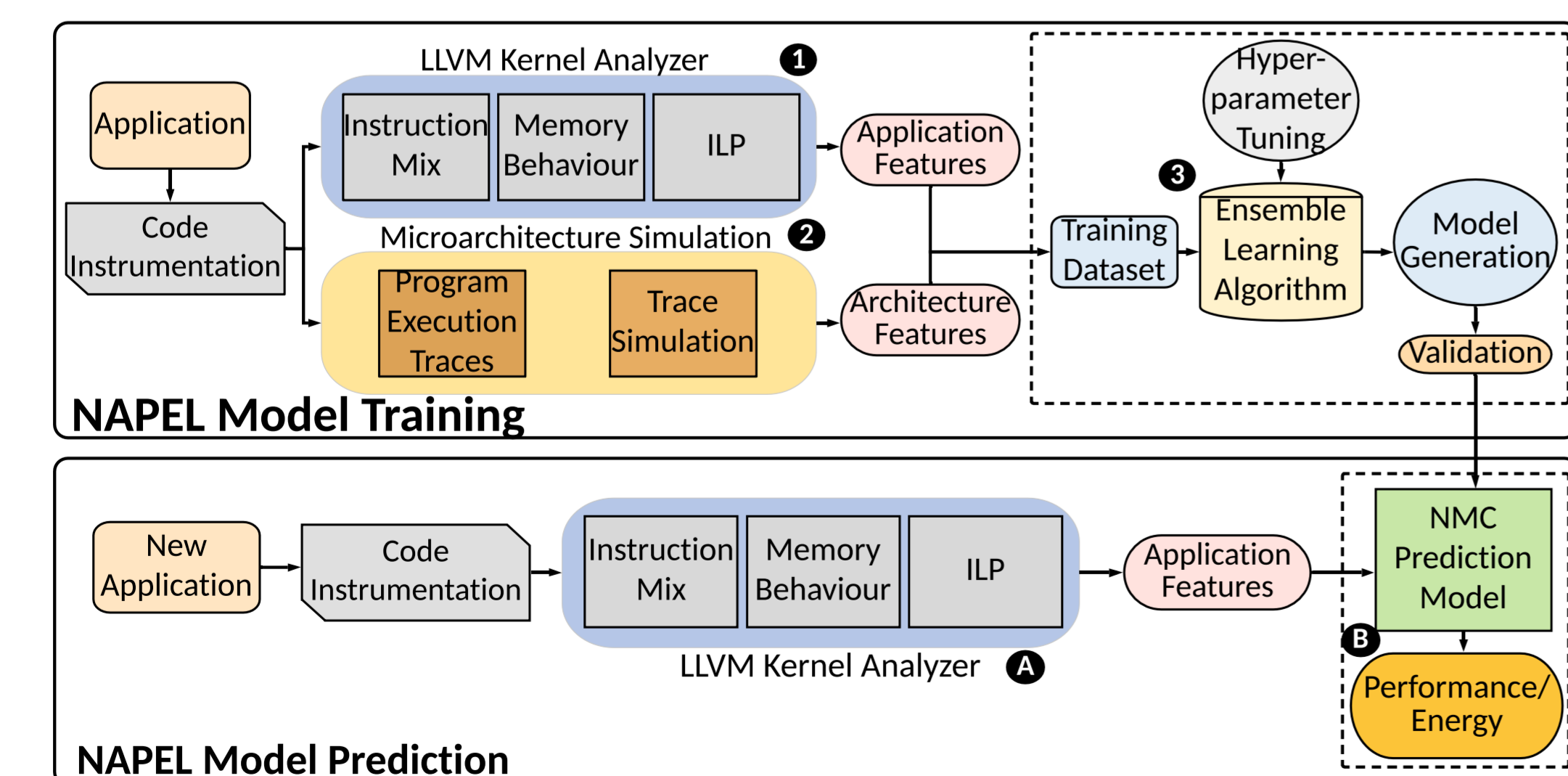
## I. Motivation

- An exorbitant amount of data <sup>1</sup>
- The high cost of energy for data movement <sup>1</sup>
- A paradigm shift towards processing close to the data, i.e., near-memory computing (NMC) <sup>2</sup>
- In early design-stage, simulations are extremely slow, imposing long run-time <sup>3</sup>



## II. Performance Prediction Framework

- NAPEL Framework <sup>3</sup> can provide fast and accurate performance and energy prediction for a previously-unseen application
- Microarchitecture-independent characterization of an application with architectural simulation responses to train an ensemble algorithm
- Intelligent statistical techniques <sup>4</sup> to extract meaningful data with minimum experimental runs

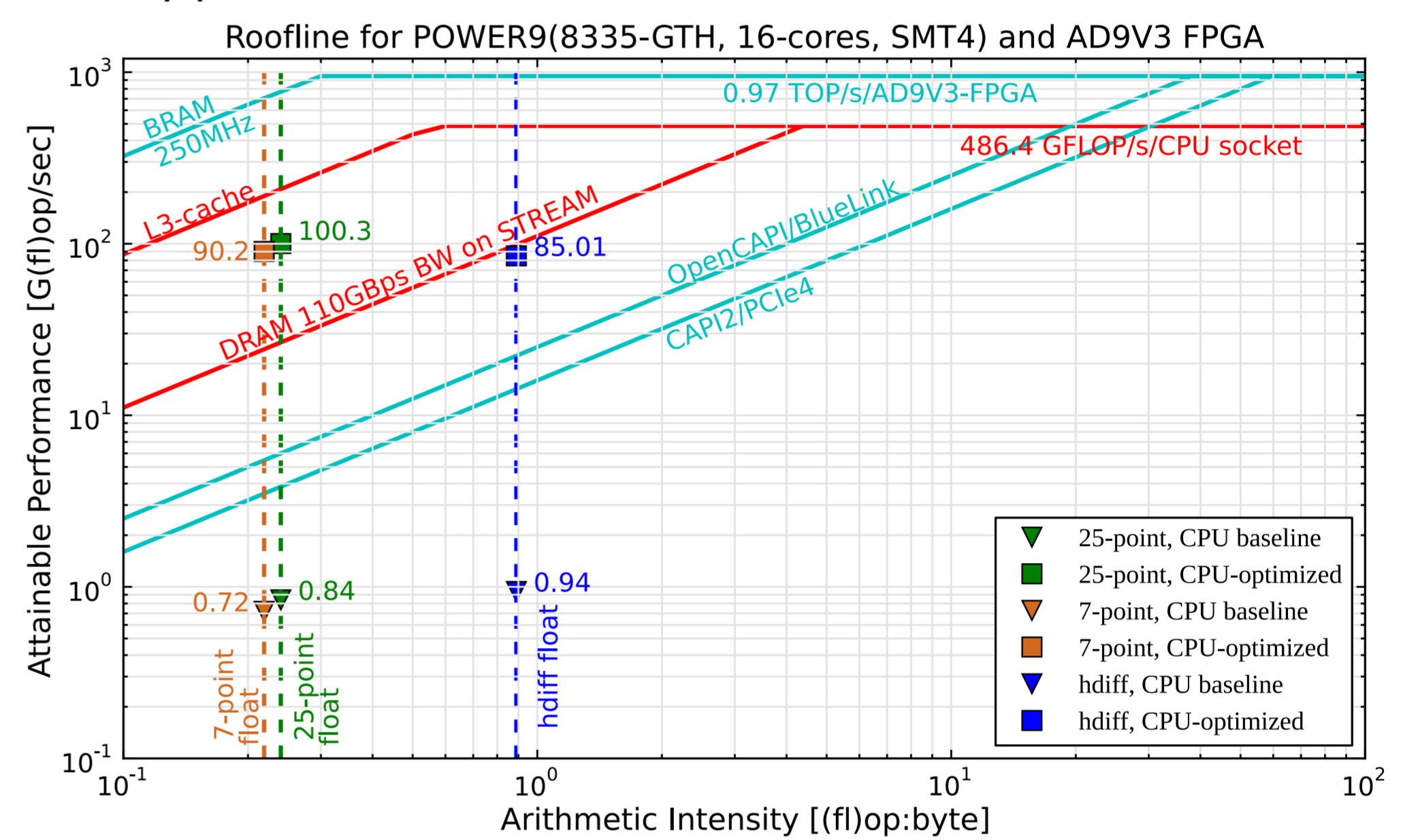


## References

<sup>1</sup> Mutlu, O., et al. "Enabling Practical Processing in and near Memory for Data-Intensive Computing" *Proceedings of the 56th Design Automation Conference (DAC), 2019*  
<sup>2</sup> Singh, G., et al. "A Review of Near-Memory Computing Architectures: Opportunities and Challenges." *21st Euromicro Conference on Digital System Design (DSD), IEEE, 2018*  
<sup>3</sup> Singh, G., et al. "NAPEL: Near-Memory Computing Application Performance Prediction via Ensemble Learning." *Proceedings of the 56th Design Automation Conference (DAC), 2019*  
<sup>4</sup> D. C. Montgomery, Design and analysis of experiments, 2017

## III. Stencil Computation on a CPU-Centric Architecture

- Implementing 3D stencils from weather application on IBM POWER9 CPU
- Performance gap due to a combination of memory hierarchy under-utilization and application's complex access patterns
- FPGA-aware roofline, which consists of AD9V3 FPGA board with CAPI2 and OpenCAPI links, indicates that FPGA is a promising device for near-memory platform



## IV. Scaling Near-Memory Computing

- Ongoing work on the feasibility of multiple NMC devices inside a node
- Followed by inter-node scaling for weather forecasting application
- System integration of IBM POWER9 and near-memory computing units (NMC) via high bandwidth OpenCAPI interface
- Leveraging OpenCAPI lowest point of coherence (LPC) extension to attach an accelerator functional unit (AFU)
- LPC allows AFU memory to act as system memory. Hence, the AFUs can be used to for acceleration near system memory.

