# Kelis

**Kelis: performance, power, area and cost modeling of AI datacenters built for large language model training and inference**
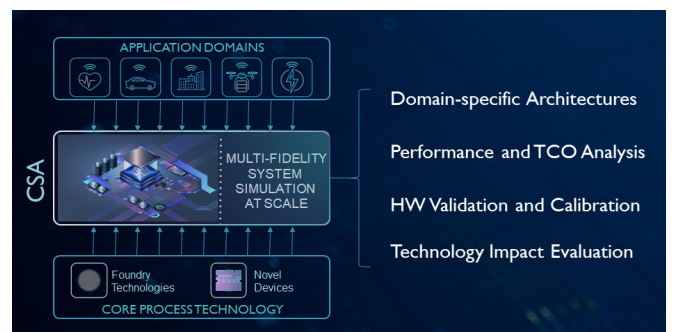
Large language models (LLMs) have proven to be of high value and show an even larger promise for future applications in numerous areas of science and technology. Many innovations stem from scaling the complexity of those models, requiring an **ever-increasing computational need** from the systems that execute them. Uniquely positioned in the R&D landscape across the entire system and technology stack, imec aims to solve scaling issues by creating new approaches to computing, connectivity, and architecture for future high-performance systems.

Compute System Architecture (CSA) is a center of excellence in imec enabling **true hardware-software-technology codesign to architect HPC and AI systems of the future**. The goal is to pathfind energy-conscious, high-performance computing solutions at scale for the chiplet era. The team analyzes emerging usage models and architects compute systems capable of post-exascale performance. Imec's expertise in system-level modelling, performance analysis and hardware validation enables system exploration encompassing the entire spectrum of workload complexity, modelling granularity and technology maturity.

Imec's Kelis tool provides fast and accurate performance, power, area and cost evaluation and design space exploration for AI datacenters running large language models, helping to quickly evaluate and optimize design choices. **Validated within 12% worst case error** for large scale LLM training and inference executions on Nvidia A100 and H100 systems, Kelis returns results within seconds, allowing for interactive exploration. The framework leverages imec's expertise in analytical performance modeling for high-performance computing and artificial intelligence.



Imec, the world-leading independent R&D centre on semiconductor technology, headquartered in Leuven, Belgium, hosts over 6,000 experts from 100 countries.
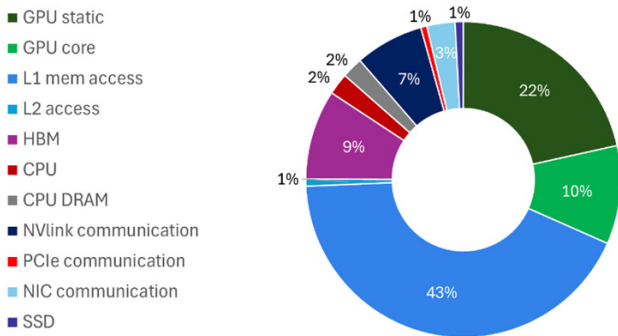


CSA – imec centre of excellence for HW-SW codesign of future compute systems, enables novel system architectures using scalable multi-fidelity simulation framework.

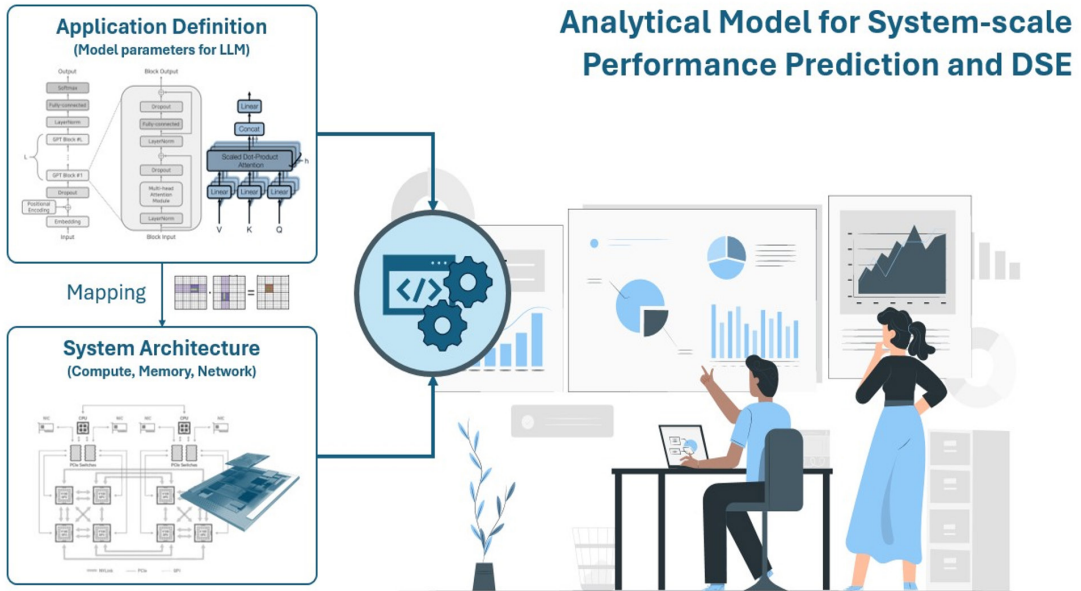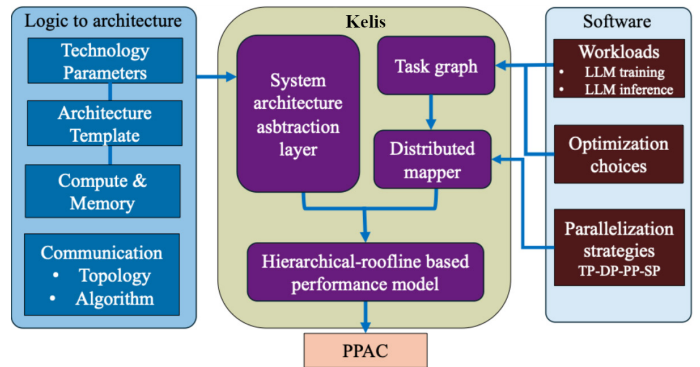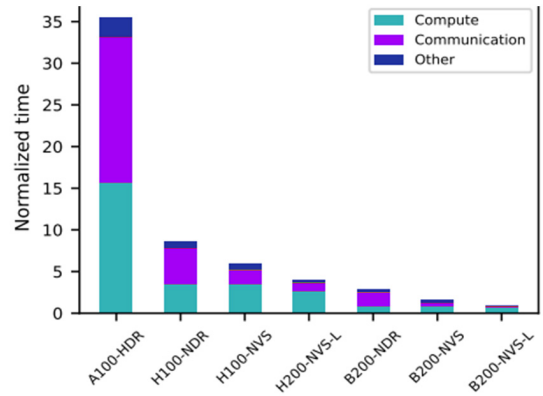# Use Kelis to optimize AI datacenter performance/TCO

The Kelis tool offers modeling of the compute, communication, and memory subsystem of accelerators, as well as the larger systems they are part of, all the way to complete datacenter scale. Kelis consists of an LLM task-graph analyzer, a parallelism mapper, a hierarchical roofline model, a topology-aware collective communication library along with many knobs to play with, that give the user an **end-to-end framework for all the explorations regarding large-scale AI datacenters and relevant workloads**. The Kelis framework is capable of evaluating the impact of new compute, memory, or communication technology for AI workloads.

**Contact us for more information and licensing.**

Normalized training performance scaling across GPU generations and networking systems for GPT3-175B

Energy breakdown of a single node during training of GPT3-175B on a system representing a 16k A100 cluster

- GPU static — 22%
- GPU core — 10%
- L1 mem access — 43%
- L2 access — 1%
- HBM — 9%
- CPU — 2%
- CPU DRAM — 2%
- NVlink communication — 7%
- PCIe communication — 1%
- NIC communication — 3%
- SSD — 1%

**Analytical Model for System-scale Performance Prediction and DSE**

Kelis offers an easy-to-use, interactive interface exposing key parameters, showing results on a dynamic dashboard.

Reference: J. Kundu, W. Guo, A. BanaGozar, U. De Alwis, S. Sengupta, P. Gupta & A. Mallik, (2024). Performance Modeling and Workload Analysis of Distributed Large Language Model Training and Inference. arXiv preprint arXiv:2407.14645 (IISWC'24).

FIND OUT MORE:

CONTACT US
**WWW.CONTACTIMEC.COM**

imec ● Kapeldreef 75 ● 3001 Leuven ● Belgium ● www.imec-int.com