



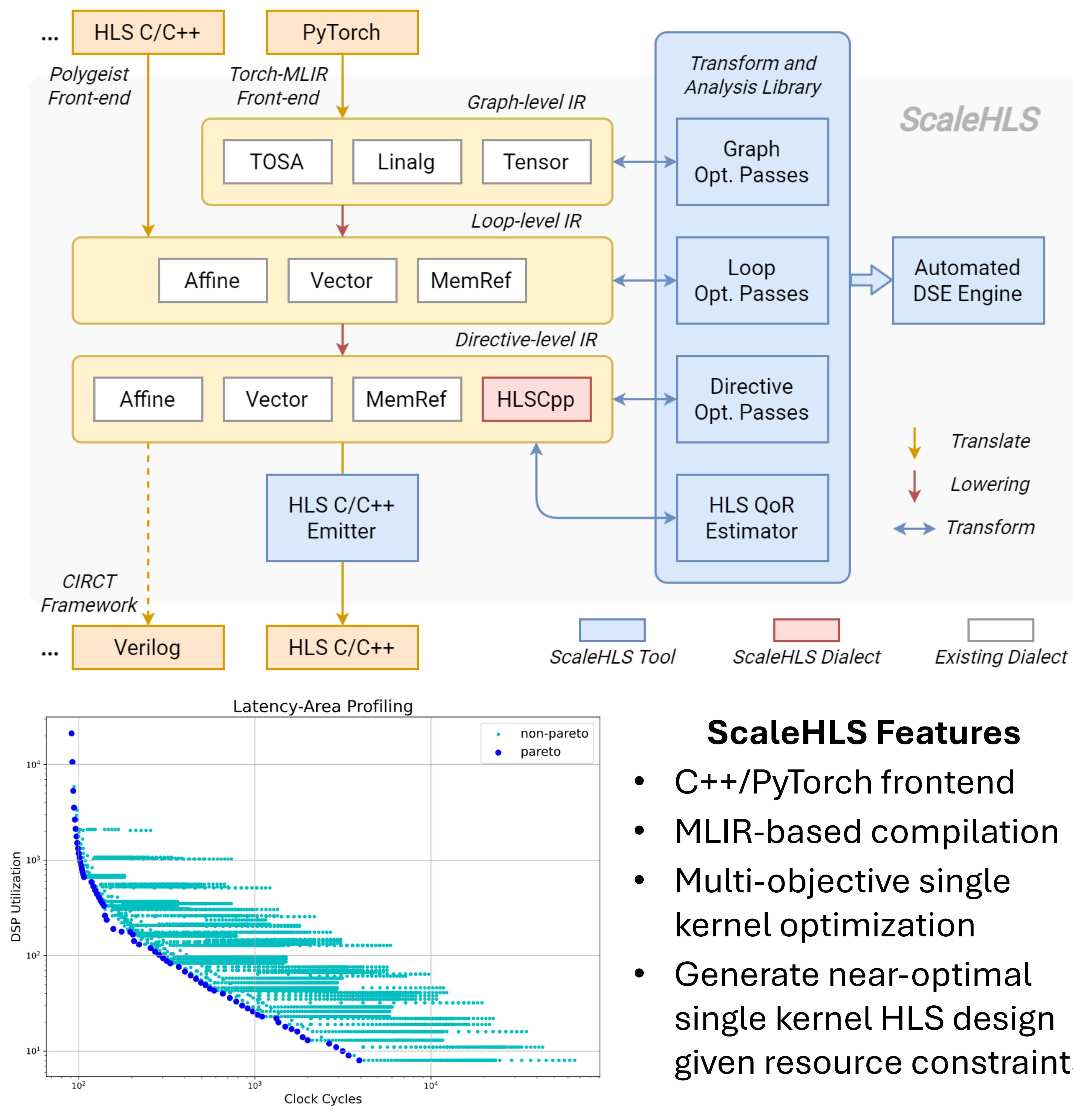
# StreamTensor: A Compiler from PyTorch to FPGA for AI/ML Applications



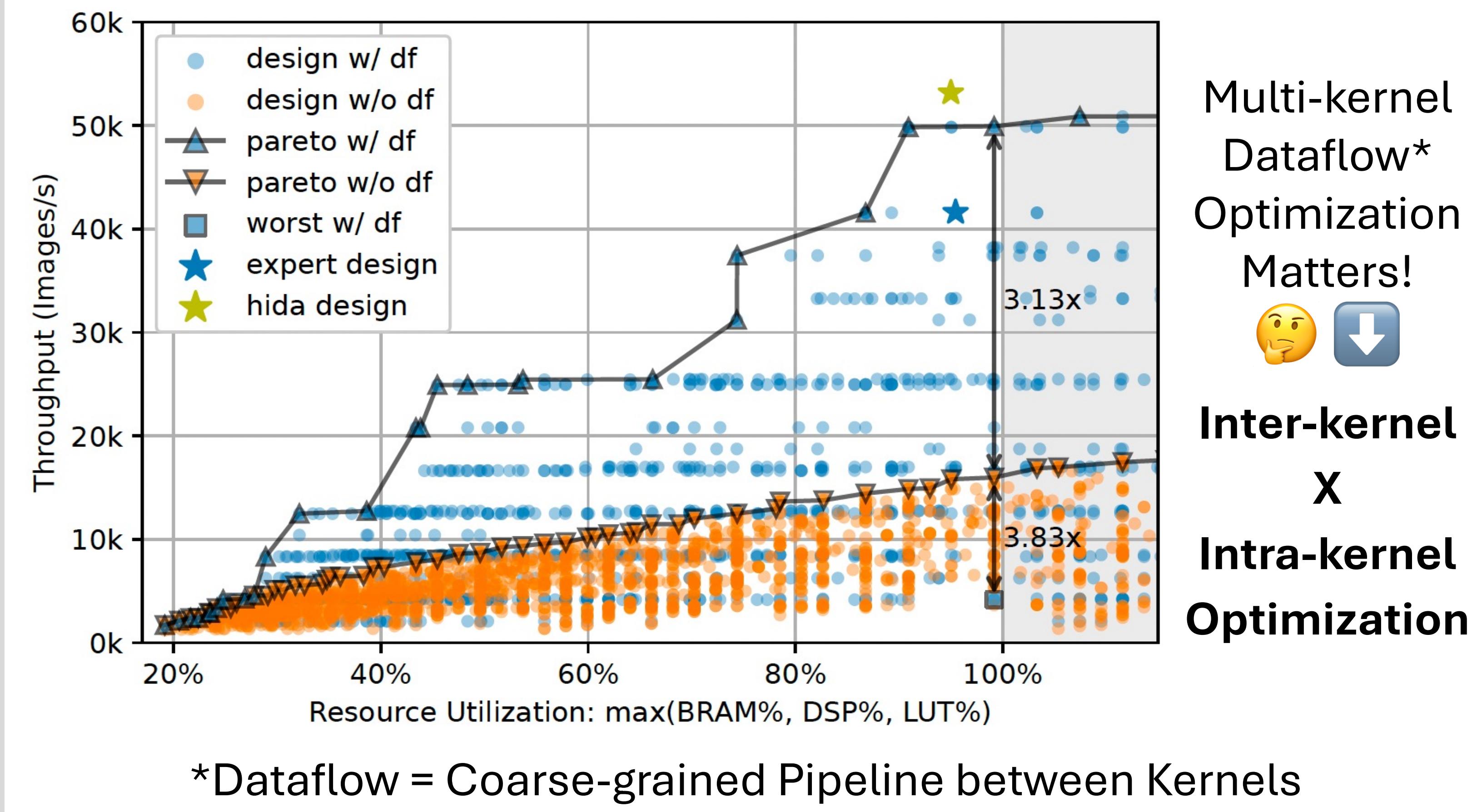
Hanchen Ye<sup>12</sup> (hanchen8@illinois.edu) and Deming Chen<sup>12</sup> (dchen@illinois.edu)

<sup>1</sup>University of Illinois Urbana-Champaign <sup>2</sup>Inspirit IoT

## ScaleHLS: Single-kernel Optimization (2022)



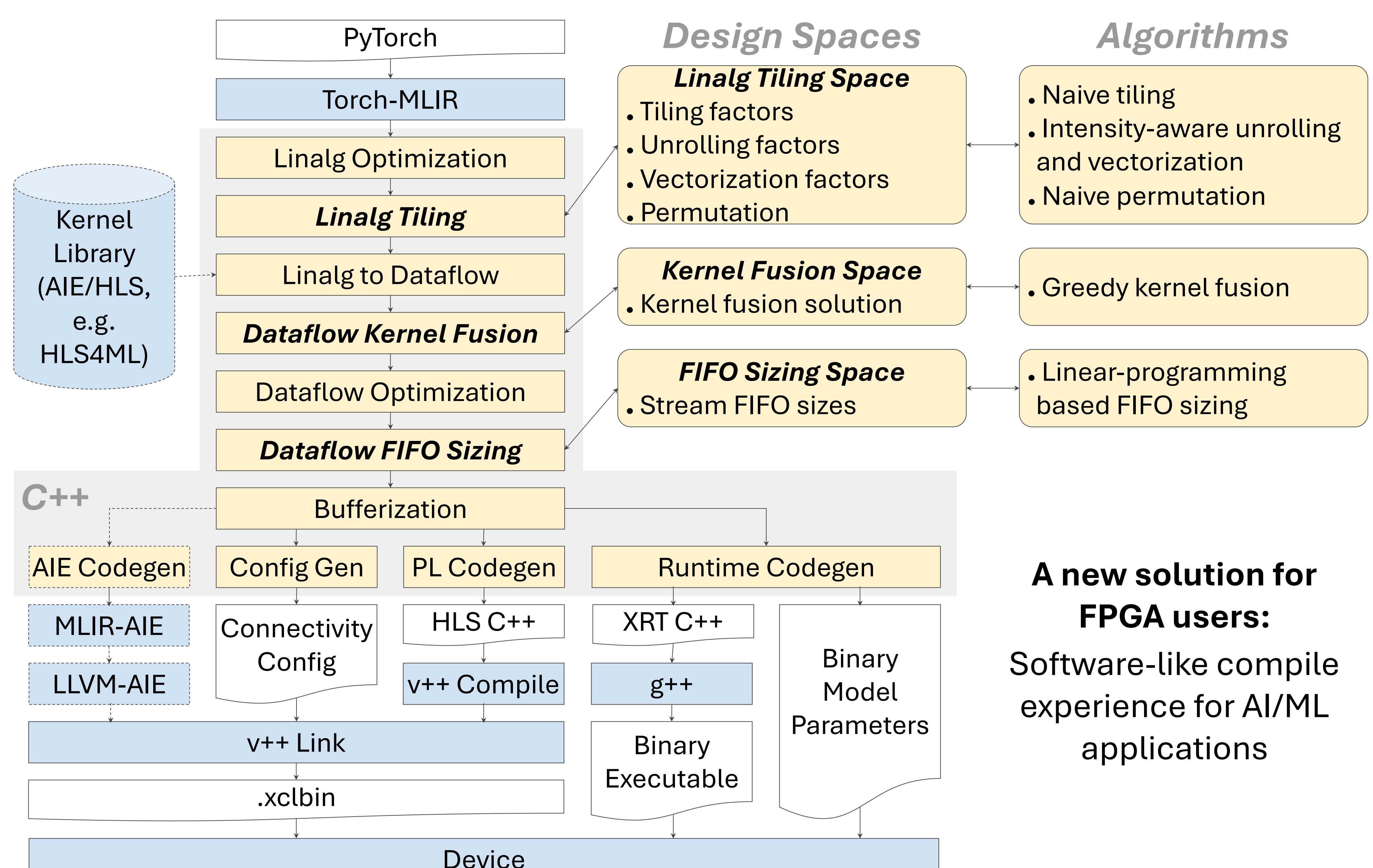
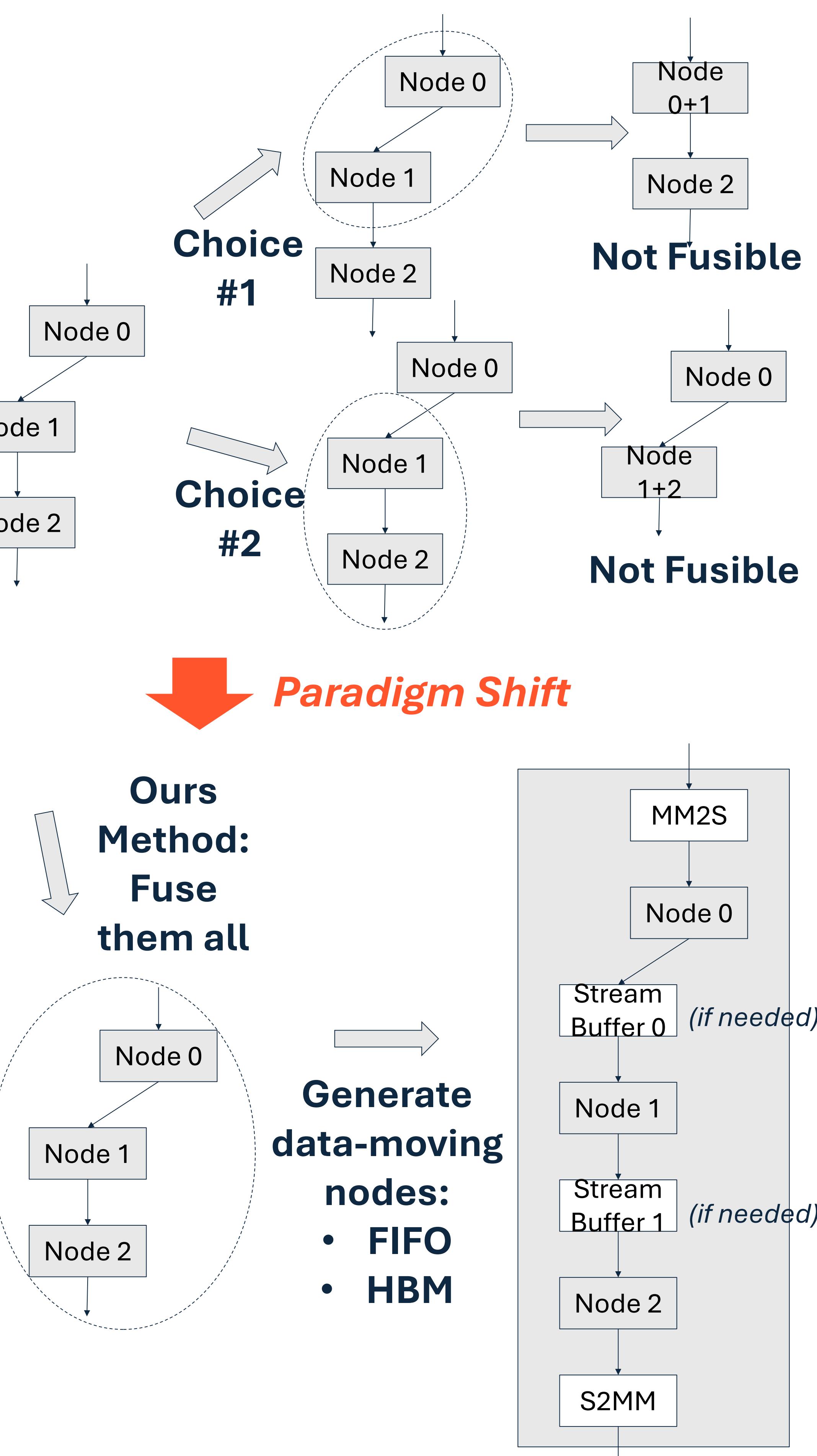
## HIDA: Multi-kernel Optimization (2024)



**Multi-kernel Dataflow\* Optimization Matters!**

**Inter-kernel X Intra-kernel Optimization**

## StreamTensor: Streaming-based Multi-kernel Dataflow (Now)



- End-to-end PyTorch-to-device flow validated on FPGA BOARD
- Aggressive kernel fusion with streaming
  - Single kernel for a single 😊 Mistral-7b or Llama-8b decoder layer on U55C FPGA
- Automated data movement generation
  - DMA, data packing, and interface widening for HBMs
  - Data vectorization, data layout conversion, and FIFO sizing for streaming FIFOs
- Pythonic user-level API
  - Python API for passes and design spaces; auto-tuning-ready