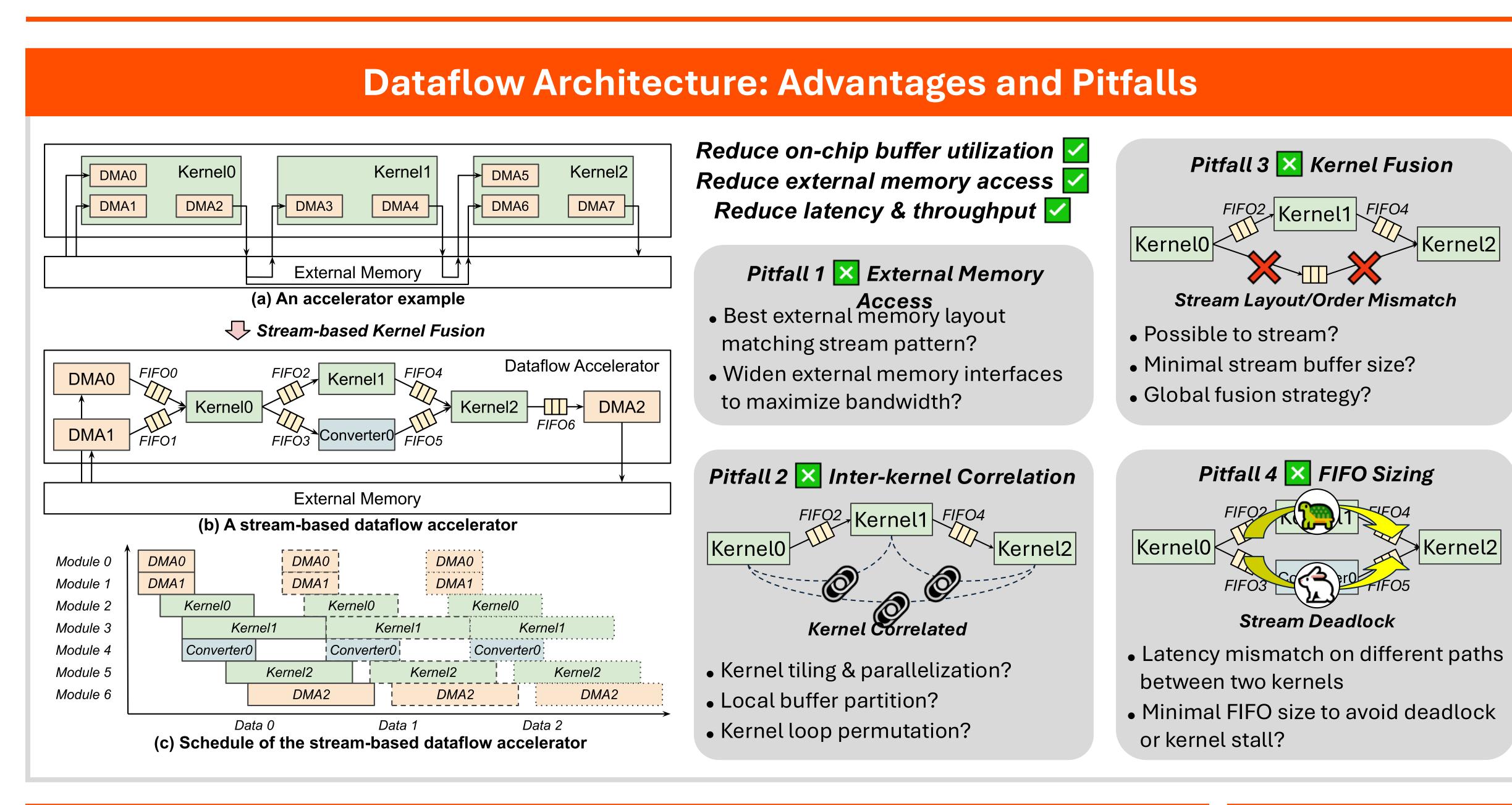
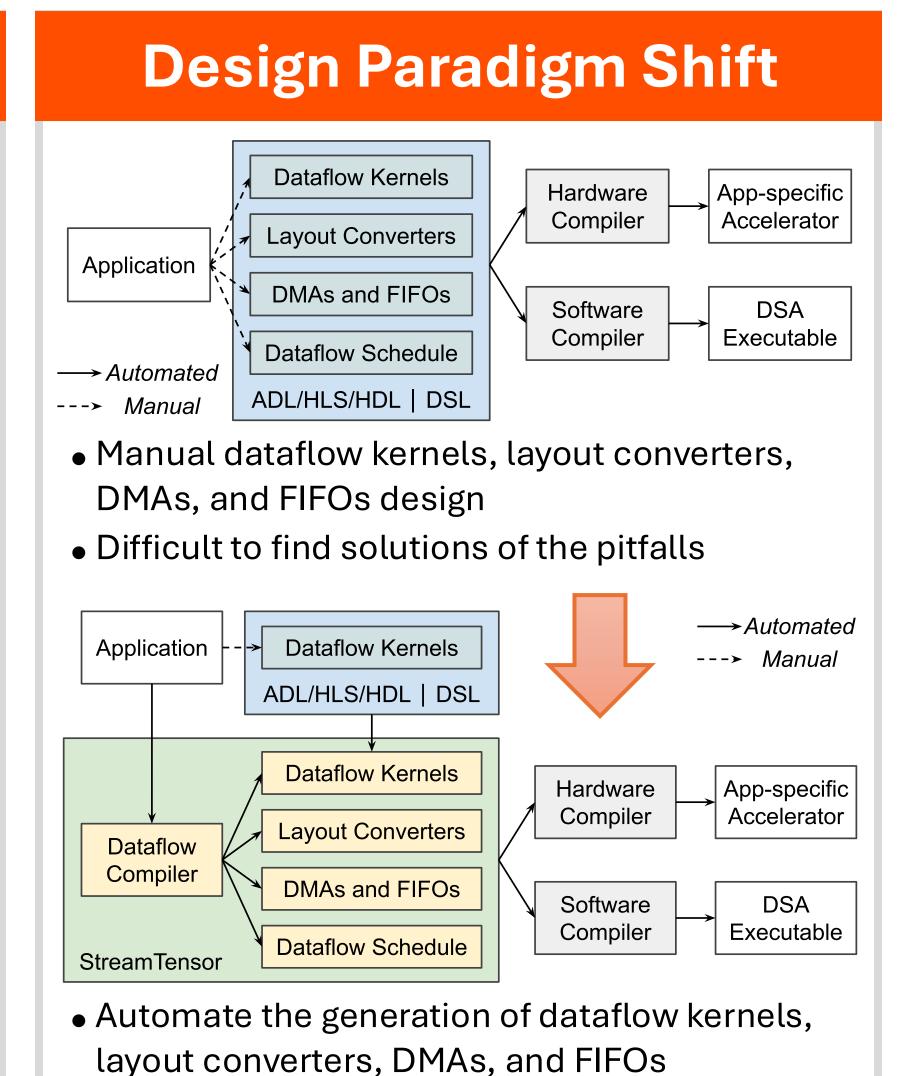


StreamTensor: Make Tensors Stream in Dataflow Accelerators for LLMs



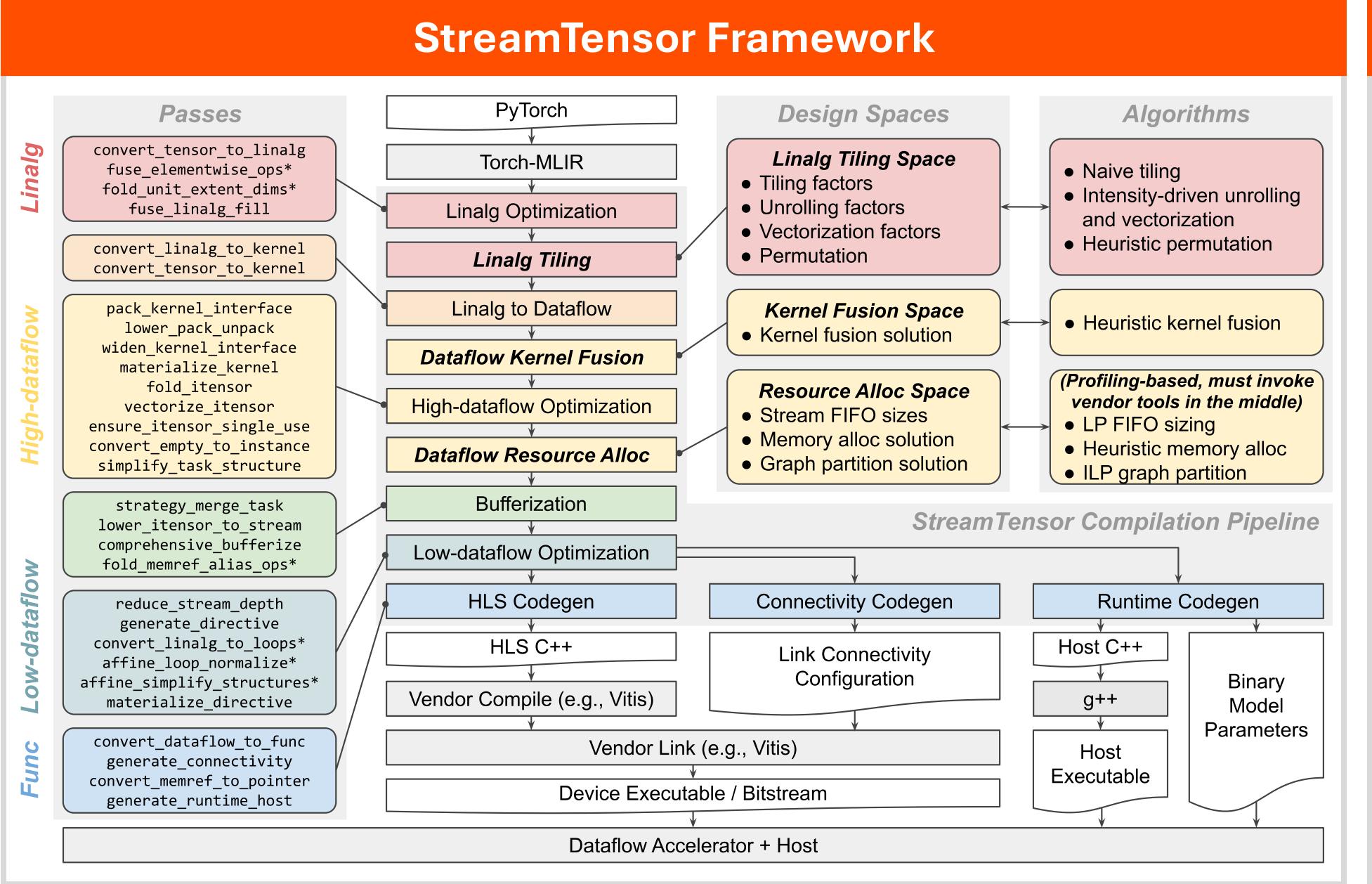
Hanchen Ye¹, Deming Chen^{1,2} ¹University of Illinois at Urbana-Champaign; ²Inspirit IoT, Inc.

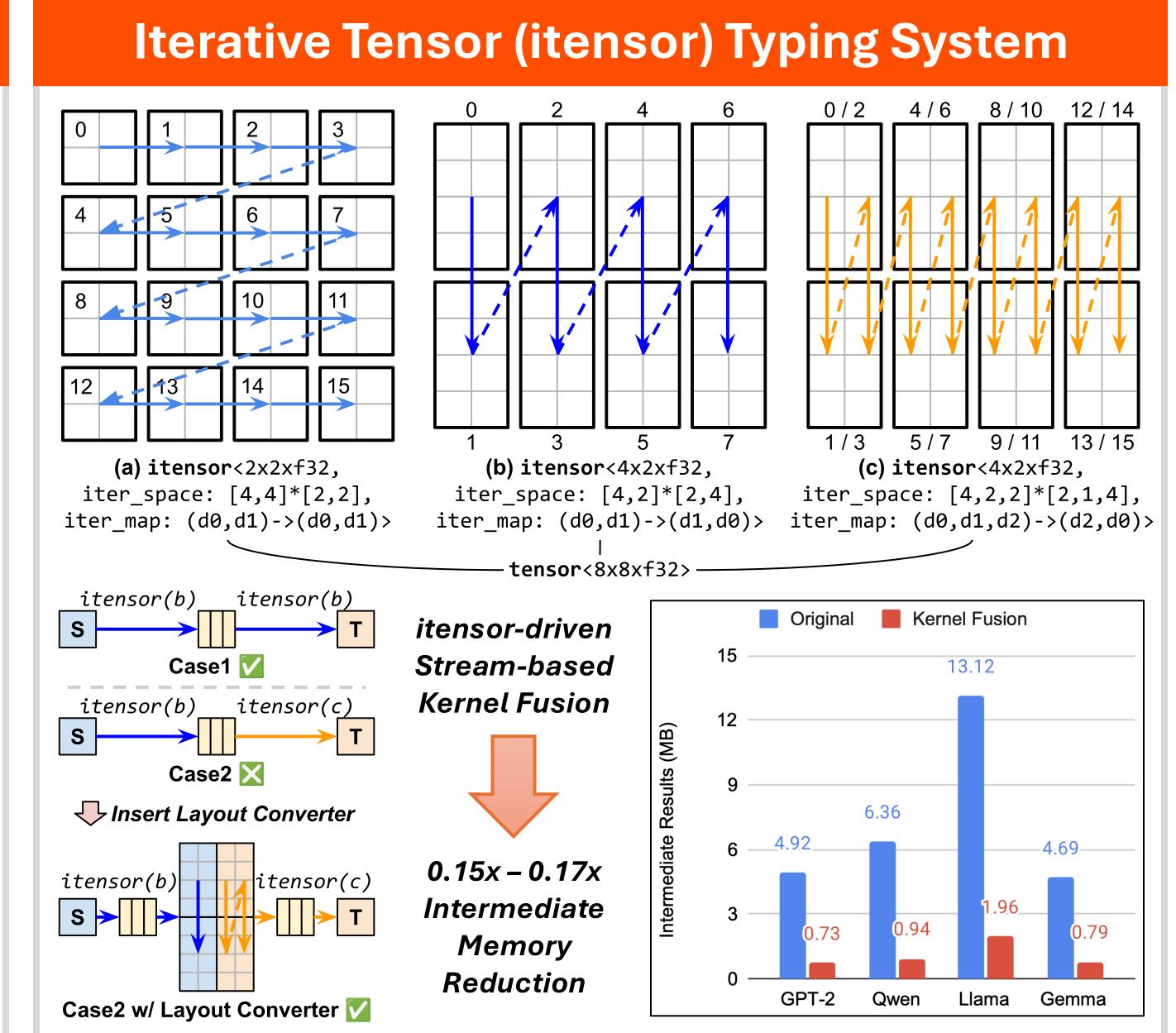




Resolve pitfalls through systematic DSE

Support auto-tuning and external library





StreamTensor Evaluation Results on LLMs

