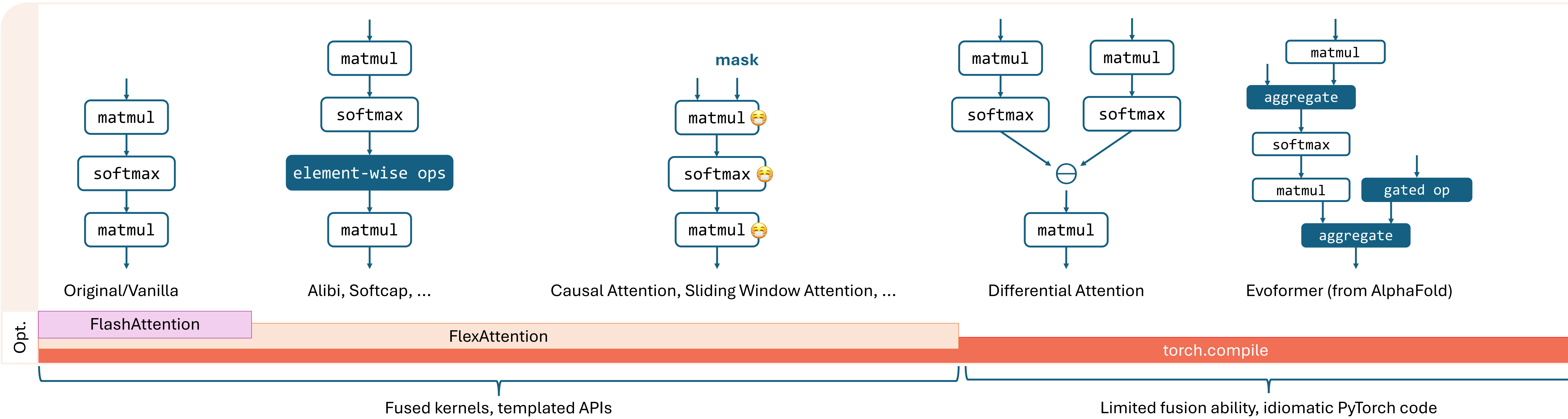
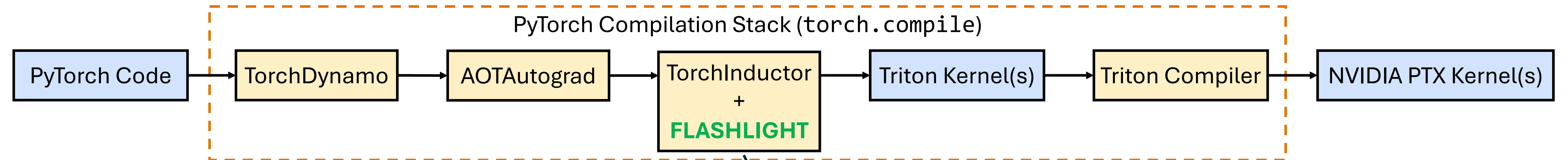


PyTorch Compiler Extensions to Accelerate Attention Variants

Bozhi You, Irene Wang, Zelal Su Mustafaoglu, Abhinav Jangda, Angélica Moreira, Roshan Dathathri, Divya Mahajan, Keshav Pingali

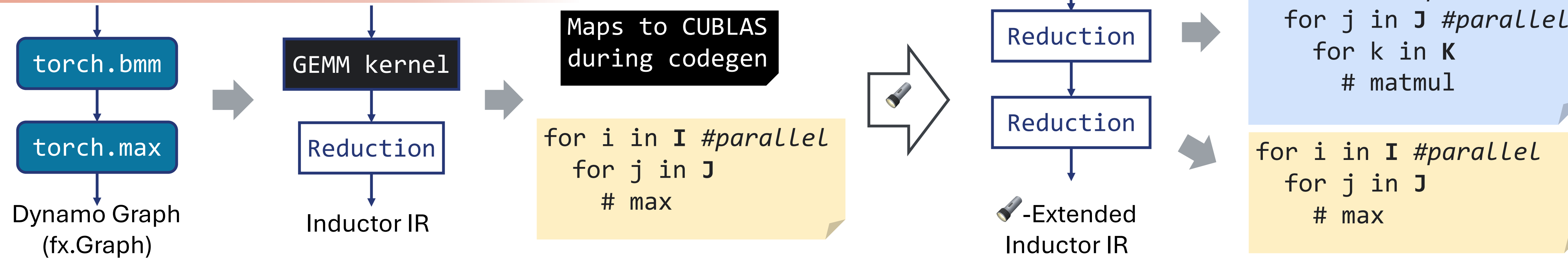


FLASHLIGHT generates optimized, fused kernels from idiomatic PyTorch code



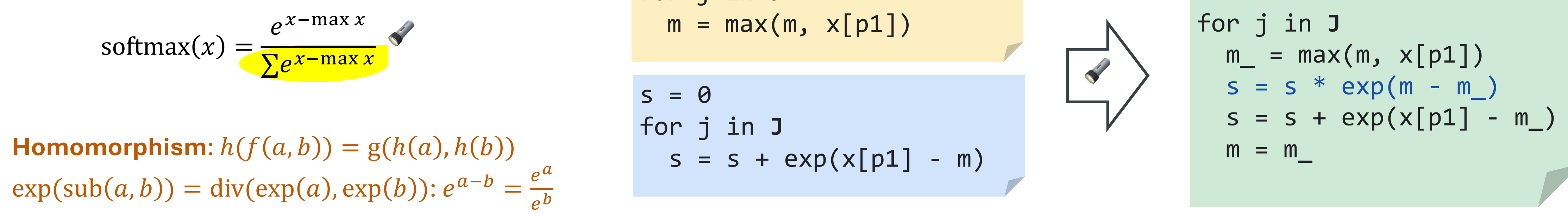
General extensions to TorchInductor, not a standalone compiler

1 Unified, Generalized Reduction IR

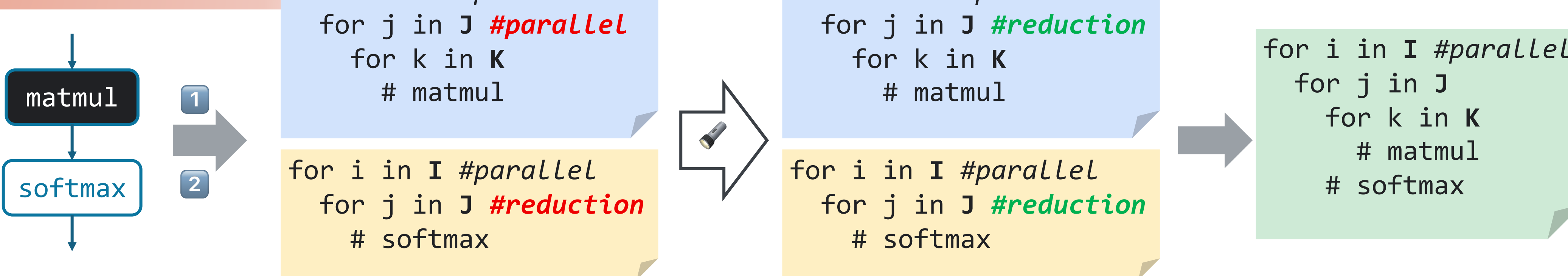


Tradeoff: sacrifices specialized GEMM kernel performance for fusion opportunities

2 Algebraic Transformation

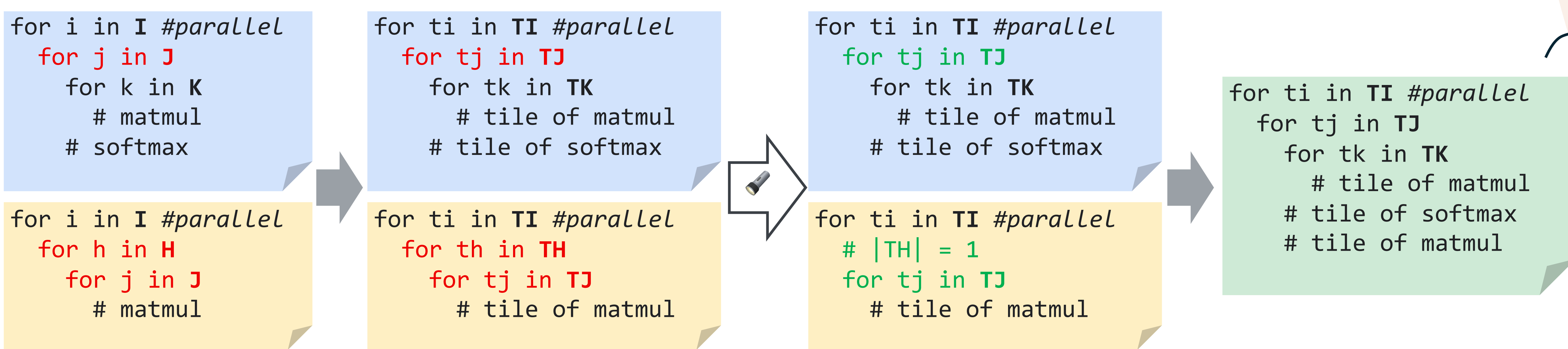


3 Dimension Demotion

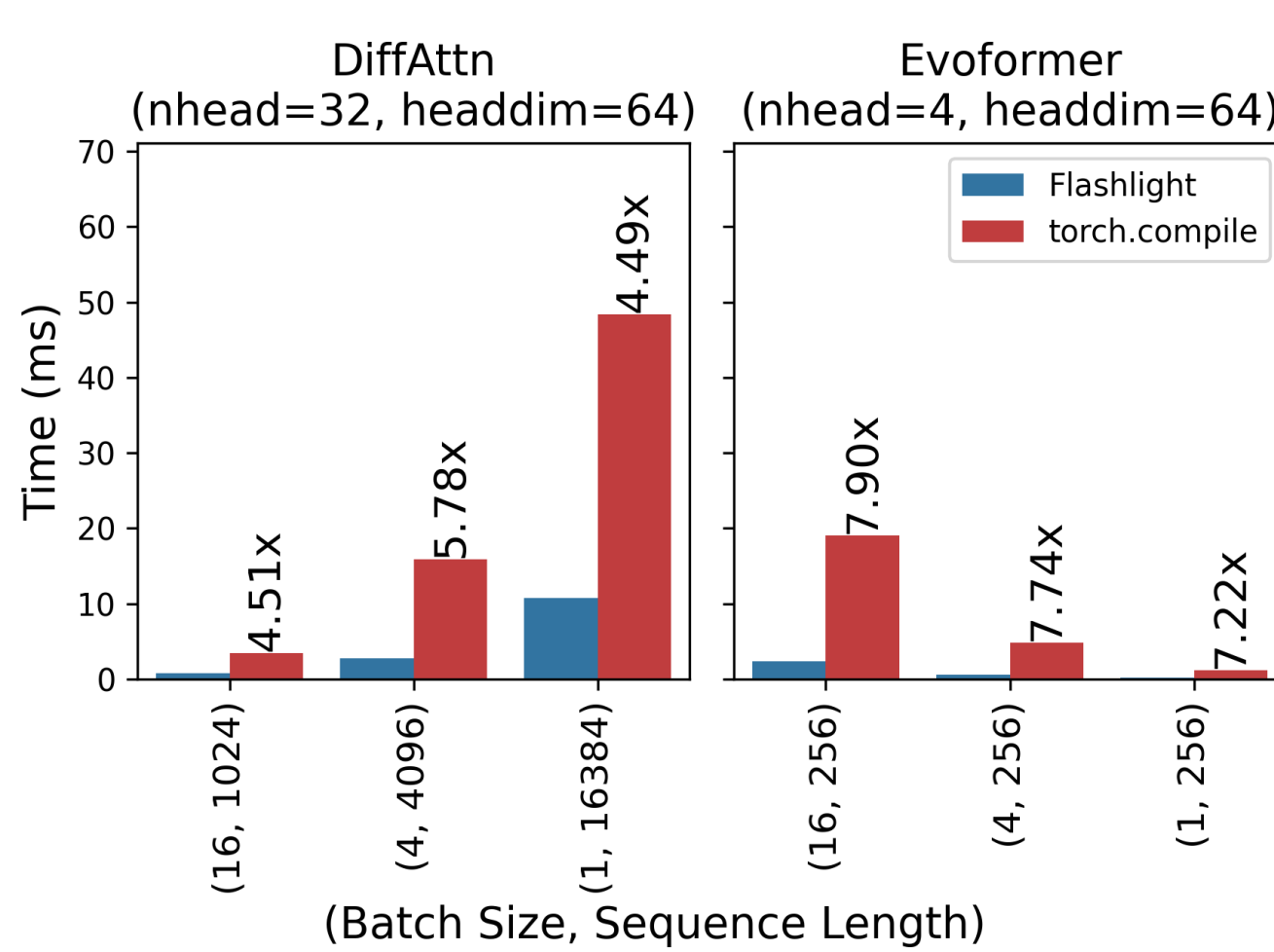
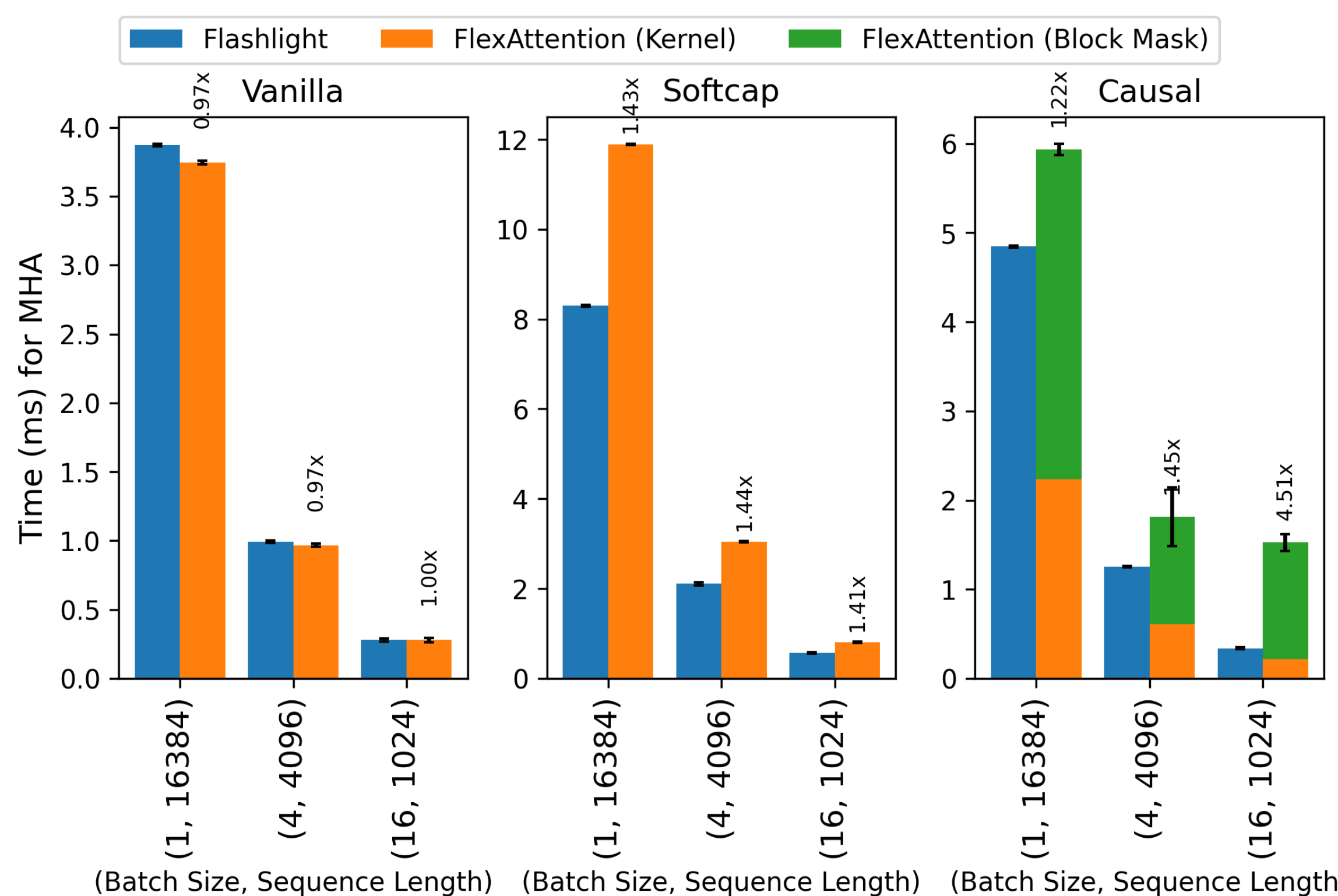


Tradeoff: sacrifices parallelism for I/O efficiency

4 Tiling-Aware Dimension Elimination

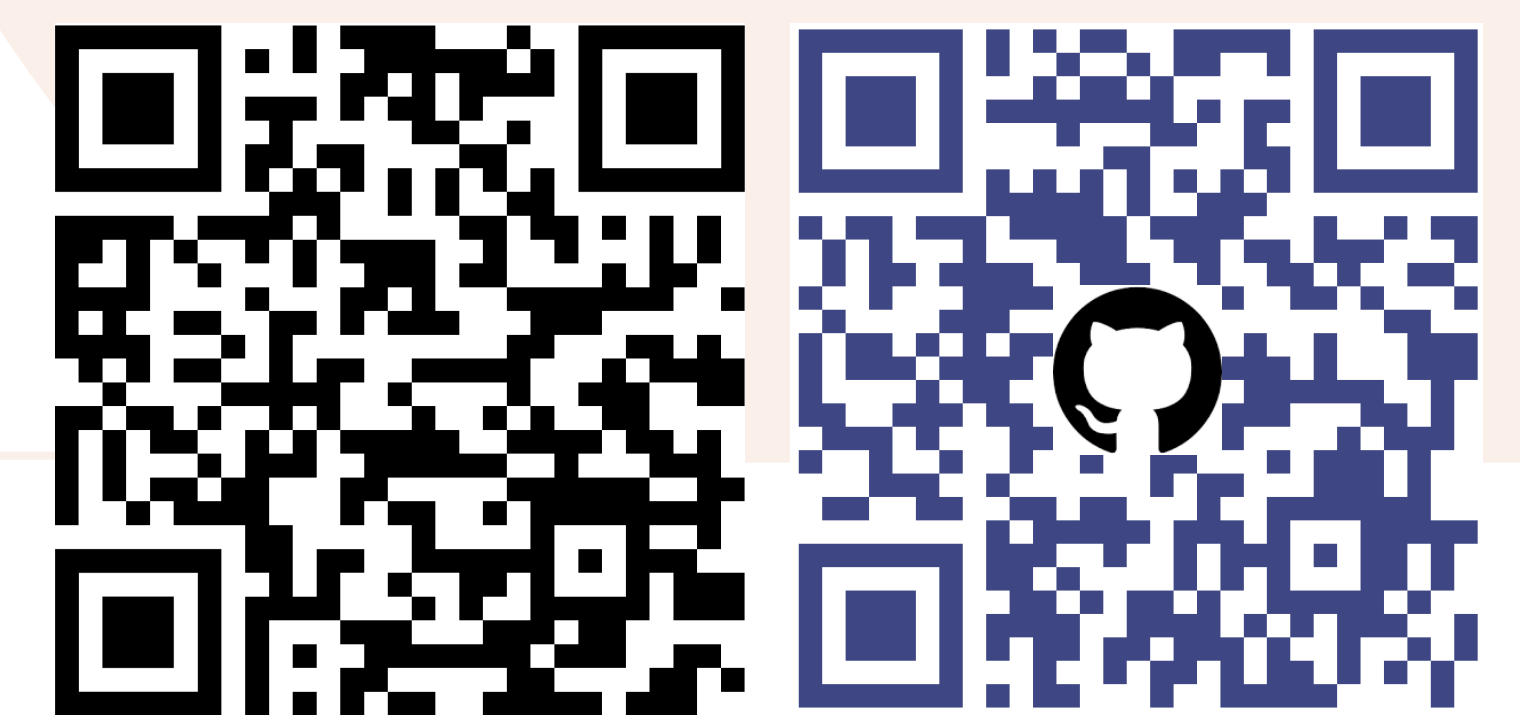


When H (head dim of V) is small, set tile_size_of_H = H so that |TH| = 1 to drop TH loop



Read the paper for more results

- on A100
- for GQA
- for more variants
- for more configurations
- against FlashInfer



Paper

Code

General extensions to PyTorch compiler (torch.compile)

- Generates I/O-efficient Triton kernels from idiomatic PyTorch code
- Moves the optimization burden from the user to the compiler

Superior performance

- Always faster than default PyTorch Compiler
- For FlexAttention-supported variants: similar performance
- Beyond FlexAttention: 5x faster for Evoformer
- Improves end-to-end inference latency for AlphaFold by 6-9%