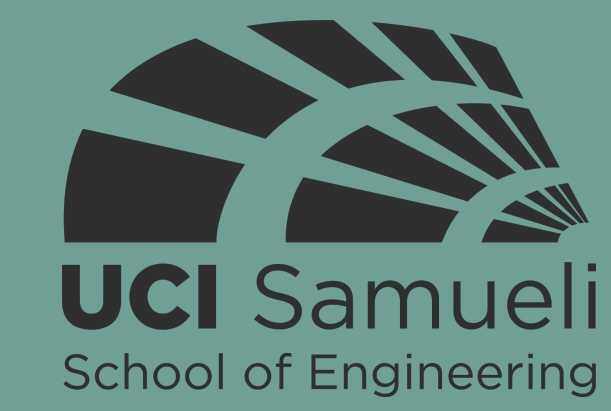


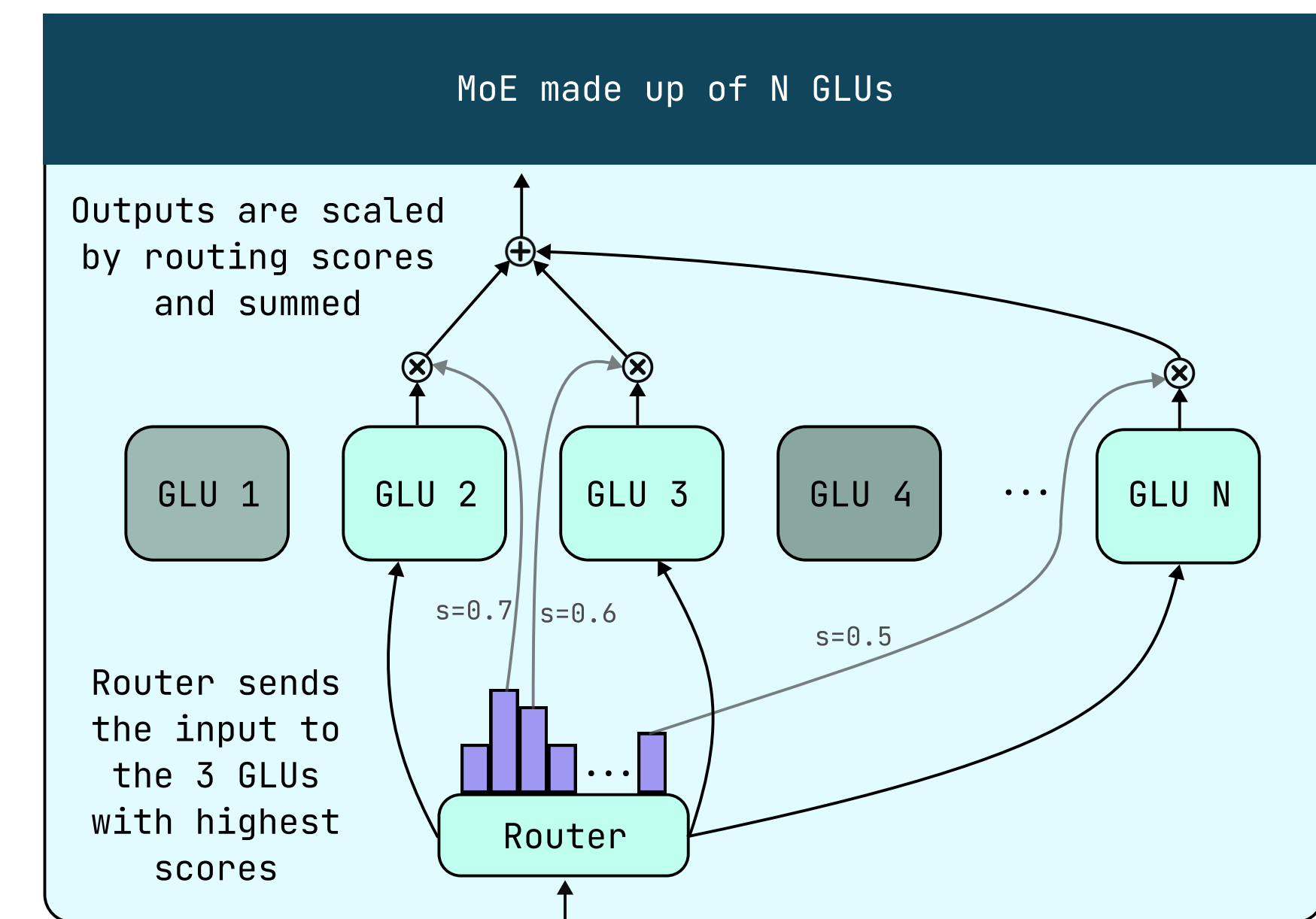
GPU Kernels for Mixture of Experts

Arthur Feeney¹ Ying Wai Li² Aparna Chandramowlishwaran¹

¹University of California Irvine, ²Los Alamos National Laboratory



Background: A Sparse Mixture of Experts (MoE) is a popular module to increase the number of parameters in a model, without a corresponding increase in FLOPs. This is typically implemented by using a router to select a subset of “experts” for different inputs:

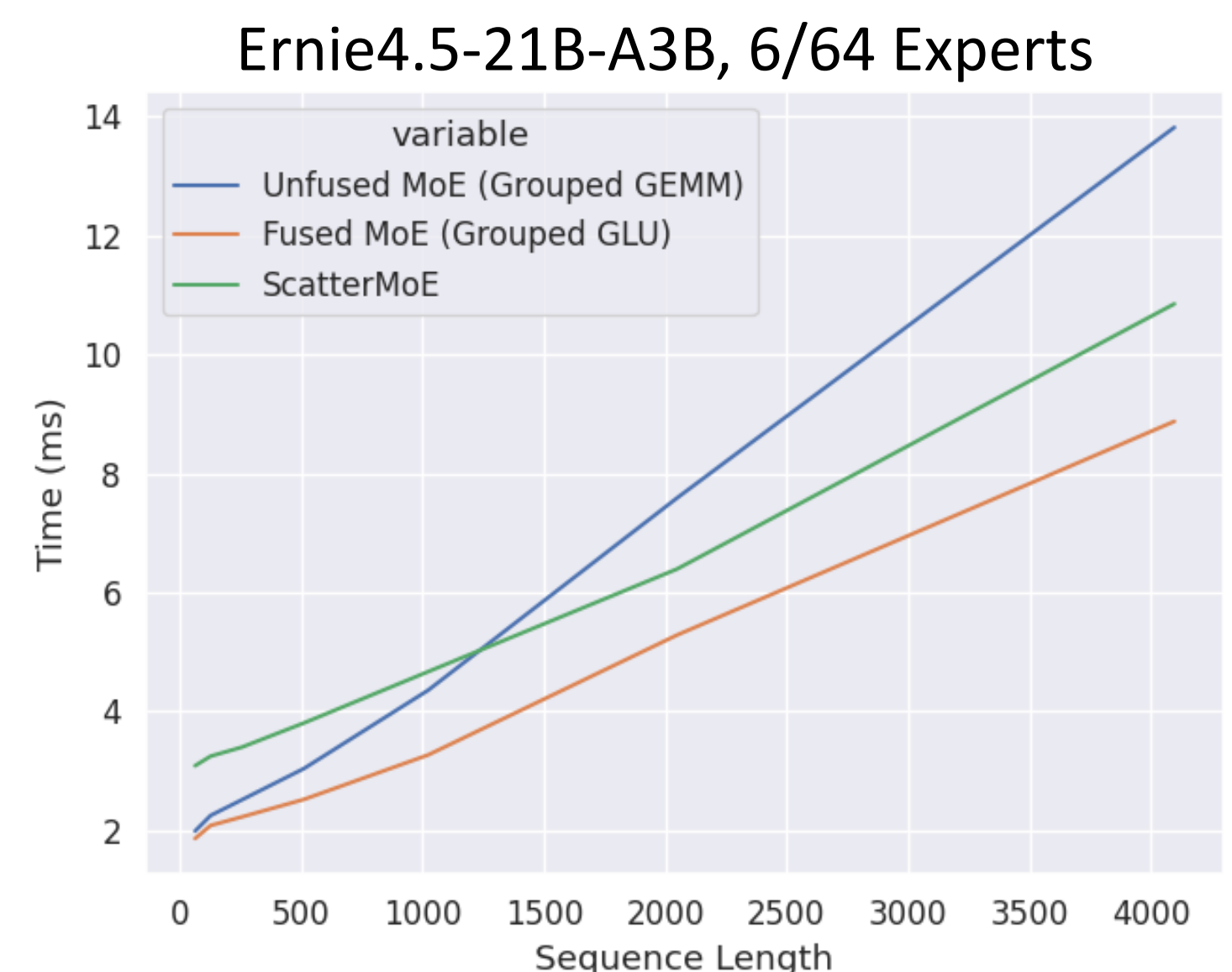
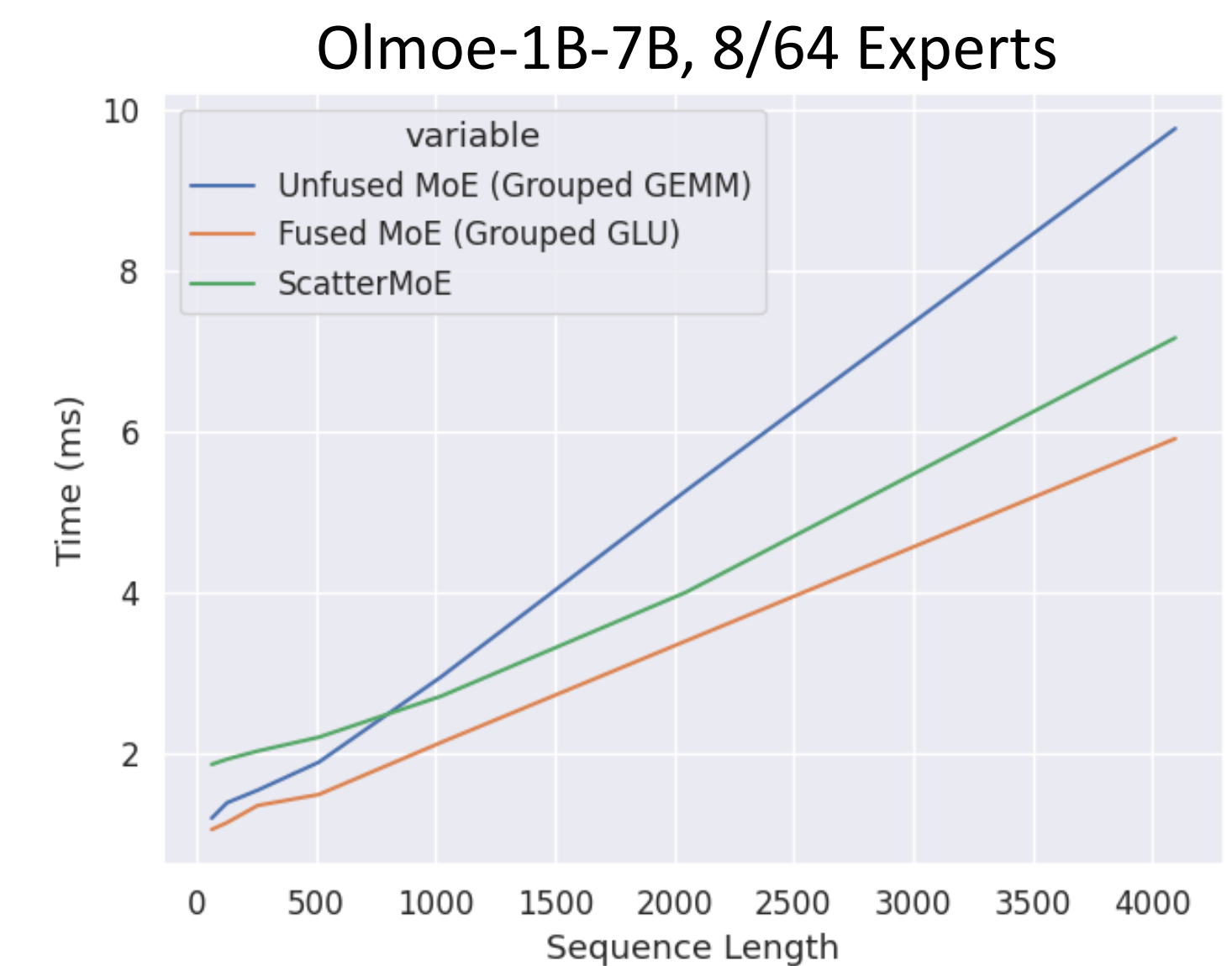
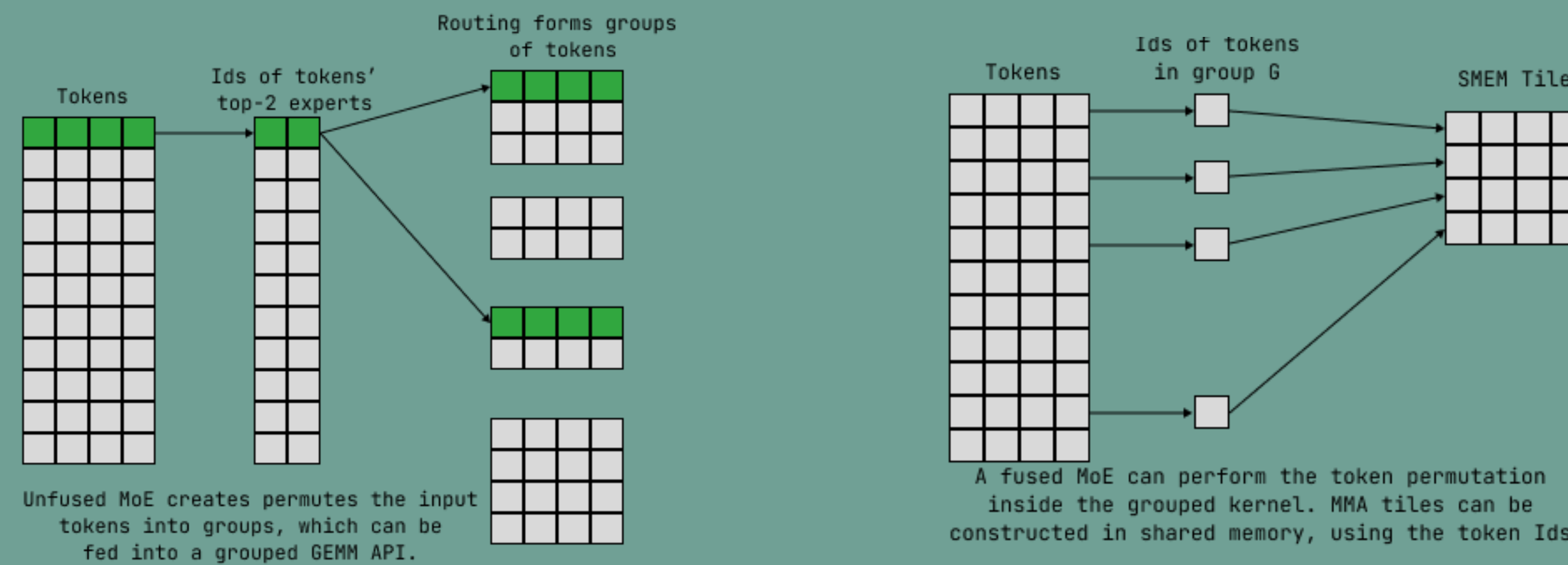


Goal: develop GPU kernels for Mixture of Experts, mainly for inference or fine-tuning.

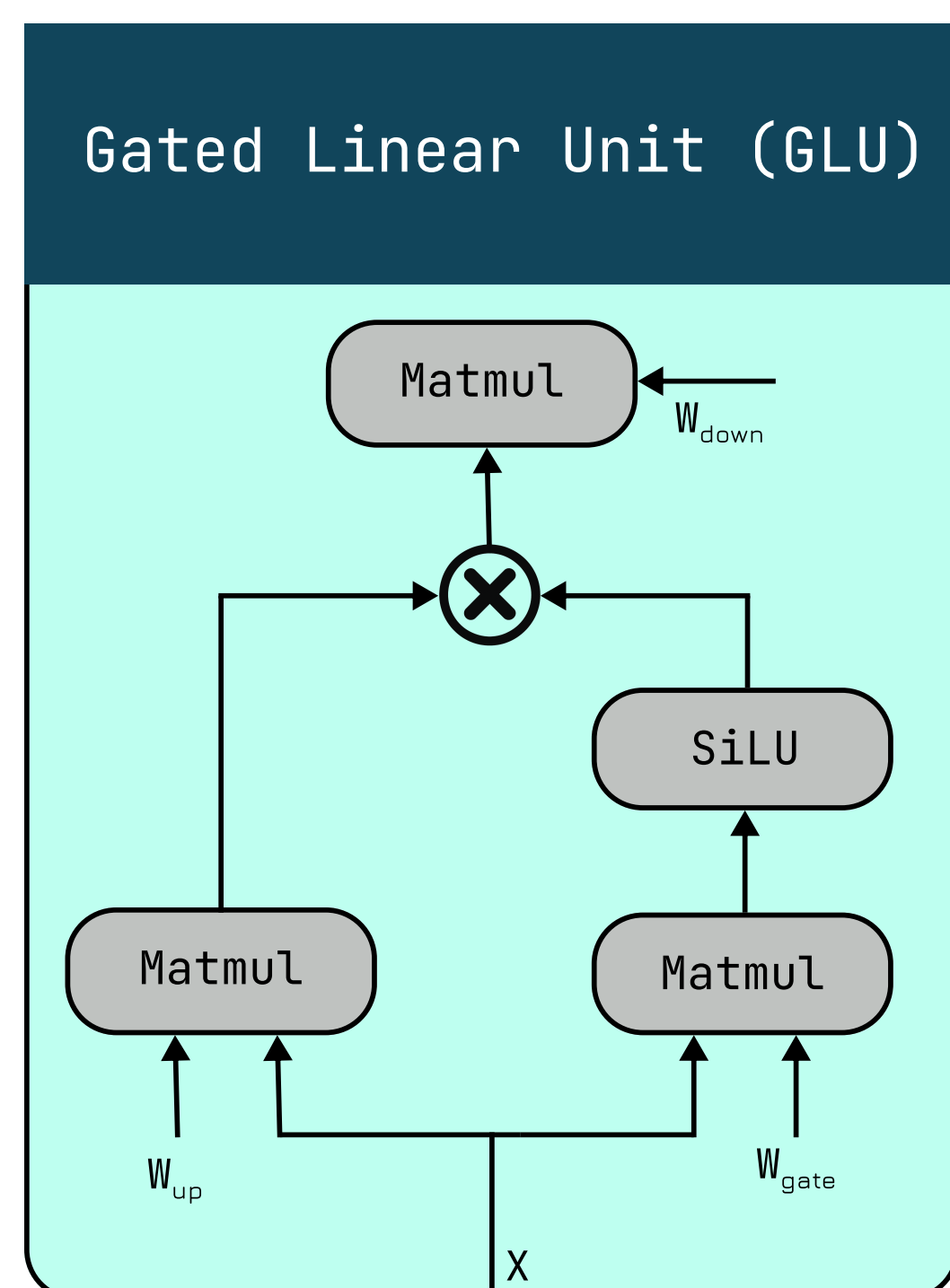
An MoE with N experts performs N matrix multiplications with N weight matrices and N different subsets of tokens. This does not map well to typical batched GEMM APIs. Alternatives APIs to process multiple GEMMs of different sizes include:

1. Magma variable batch gemm,
2. CUTLASS and cuBLAS Grouped GEMM.

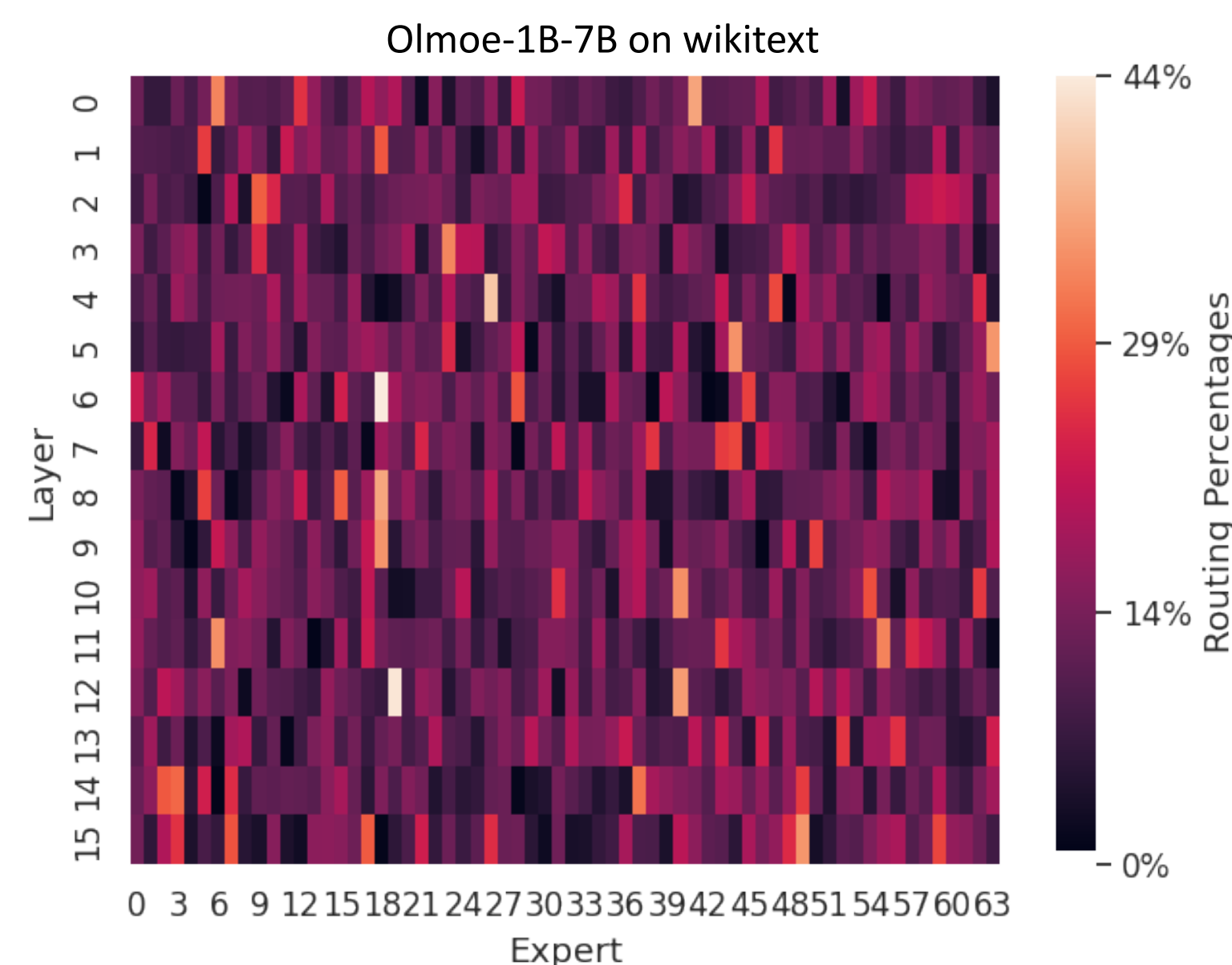
An MoE includes additional operations for routing, which can be fused with the GEMM kernels.



It is common for each experts to be an MLP or GLU:



In Practice, routing in MoEs can be imbalanced so that experts may process different quantities of tokens.



The fused routing does not affect performance relative to unfused grouped GEMM, though the routing distribution does affect performance.

