Model Possible call E.g., Edgewise typed linear semantic redundant data PyTorch Batched MM Segment MM Multiple MMs copy cuBLAS GEMV cuBLAS GEMM multiple calls libtorch **GEMM** kernels **GEMV** kernels small grid, libcublas tuned for large input low occupancy tuned for large input