MatMul Performance: GFLOPS ► FP16_CUDA 120000 FP16_FP32_MIXED_CUDA FP16_FP32_MIXED_TENSOR FP16_TENSOR 100000 -**──** FP32 80000 GFLOPS 60000 40000 20000 0 128 256 512 1024 2048 4096 8192 Matrix Size (N x N)