Execution Time (Linear) FP16_CUDA FP16_FP32_MIXED_CUDA → FP16_FP32_MIXED_TENSOR FP16_TENSOR 50000 + **→** FP32 40000 Median Time (µs) 30000 20000 10000 0 128 256 512 1024 2048 8192 4096 Matrix Size (N x N)