Execution Time (Log Scale) → FP16_CUDA FP16_FP32_MIXED_CUDA FP16_FP32_MIXED_TENSOR → FP16_TENSOR **─** FP32 10^{3} Median Time (µs) 103 10^{1} 128 256 512 1024 2048 8192 4096

Matrix Size (N x N)