

MatMul Speedup vs. Sparsity Level
on GPT-3 Layer (12k*12k MatMul)

