Speedup of GEMV on A100 (Weight Quantize) *FP*16 $cuBLAS-W_{FP16}A_{FP16}$ BitBLAS-W_{INT2}A_{INT8} Marlin-W_{INT4}A_{FP16} BitBLAS-W_{INT4}A_{FP16}

