## Speedup of GEMM on A100 (Weight Quantize) 9 LA 2.5 -BitBLAS-W<sub>INT2</sub>A<sub>INT8</sub> $cuBLAS-W_{FP16}A_{FP16}$ Marlin-W<sub>INT4</sub>A<sub>FP16</sub> BitBLAS-W<sub>INT4</sub>A<sub>FP16</sub> BitBLAS-W<sub>NF4</sub>A<sub>FP16</sub> CUTLASS-W<sub>INT4</sub>A<sub>FP16</sub> BitsAndBytes-W<sub>NF4</sub>A<sub>FP16</sub> BitBLAS-W<sub>INT2</sub>A<sub>FP16</sub>

