Normalized Speedup vs cuBLAS-W$_{FP16}$A$_{FP16}$ on A100