Performance for Model: 8B, Seq Len: 4096 on H100 176 361.70 480.78 546.54 415.70 567.69 571.01 650 160 362.19 -600 Host Memory (GB) 12 128 144 - 550 TFLOPS/s - 450 96 -400 80 - 350 60 40 78 30 50 70 Device Memory (GB)