Performance for Model: 8B, Seq Len: 4096 on H100 176 361.70 415.70 480.78 546.54 571.01 567.69 650 160 362.19 416.16 480.33 546.48 568.78 571.74 600 Host Memory (GB) 12 128 144 362.00 415.87 481.50 546.53 568.23 570.70 - 550 oo TFLOPS/ 362.22 416.30 481.31 547.13 568.91 571.20 362.09 416.40 480.75 546.80 569.07 571.29 450 96 362.00 415.99 481.29 546.88 570.08 571.43 - 400 80 489.71 494.85 518.36 553.44 568.34 569.63 - 350 60 50 30 78 40 70 Device Memory (GB)