Performance for Model: 8B, Seq Len: 8192 on H100 176 540.10 383.68 415.71 470.85 553.82 554.71 650 160 383.63 415.64 470.92 553.55 539.41 553.57 600 Host Memory (GB) 12 128 144 383.55 415.34 470.24 538.59 553.80 552.69 - 550 o 00 TFLOPS/ 383.63 415.72 470.37 539.35 552.65 553.15 383.64 470.23 415.40 539.03 552.74 552.82 450 96 471.01 540.29 554.43 383.92 415.97 553.89 - 400 80 494.86 509.84 525.62 542.81 553.59 553.89 - 350 60 50 30 78 40 70 Device Memory (GB)