Performance for Model: 8B, Seq Len: 16384 on H100 96 491.62 509.32 533.66 558.80 559.85 560.60 522.86 650 80 491.67 508.96 533.92 557.61 559.95 522.64 560.88 600 Host Memory (GB) 70 75 491.72 - 550 509.22 522.89 534.30 559.35 560.48 561.44 o 00 TFLOPS/ 491.22 495.57 516.50 527.72 548.82 560.06 560.65 450 65 491.29 495.08 512.09 523.53 550.19 560.40 537.77 400 09 466.01 473.83 494.18 513.12 529.45 540.90 552.95 - 350 24 40 60 30 50 70 78 Device Memory (GB)