Performance for Model: 8B, Seq Len: 2048 on H100 176 350.69 402.59 475.67 546.90 572.97 572.41 650 160 351.09 403.28 475.85 546.72 571.24 572.99 600 Host Memory (GB) 12 128 144 350.95 402.99 476.02 547.68 572.06 573.48 - 550 oo TFLOPS/ 351.00 403.05 475.85 547.43 572.87 573.46 350.84 403.02 476.19 547.53 571.22 572.16 450 96 350.83 547.09 403.07 475.72 571.00 573.51 - 400 80 476.68 490.48 517.95 554.23 572.26 573.08 - 350 60 30 50 78 40 70 Device Memory (GB)