Performance for Model: 8B, Seq Len: 1024 on H100 176 345.16 396.94 468.82 546.03 572.06 573.03 650 160 345.38 468.56 545.05 396.77 573.01 572.92 600 Host Memory (GB) 12 128 144 345.26 396.85 468.50 546.28 573.88 574.61 - 550 oo TFLOPS/ 396.60 468.51 546.13 345.13 572.70 573.02 345.35 396.60 468.46 546.24 572.74 573.05 450 96 345.20 396.83 468.61 546.33 572.55 573.32 - 400 80 469.26 483.84 514.95 552.91 573.63 574.96 - 350 60 50 30 78 40 70

Device Memory (GB)