Performance for Model: 8B, Seq Len: 4096 on H100 176 361.65 415.34 468.23 538.01 555.13 556.74 800 160 467.60 361.23 556.60 415.13 536.42 555.83 700 Host Memory (GB) 12 128 144 361.26 415.09 467.92 536.54 555.13 557.13 600 361.29 536.31 555.03 555.60 415.26 467.15 500 361.24 536.16 415.06 467.39 554.79 555.85 400 96 467.81 361.65 415.53 536.15 556.44 556.46 300 80 486.29 492.72 510.46 541.82 555.88 557.46 60 50 30 78 40 70

Device Memory (GB)