Performance for Model: 8B, Seq Len: 4096 on RTX5090 176 162.83 167.36 170.83 171.76 172.98 175.58 176.60 180 160 163.04 168.26 170.77 175.48 171.27 172.91 176.35 -160 Host Memory (GB) 12 128 144 163.86 171.60 167.78 170.89 172.73 175.59 176.71 140 163.79 167.85 170.86 171.56 173.13 175.48 176.43 120 164.34 167.63 170.98 172.18 173.09 176.13 176.83 - 100 96 163.94 167.28 171.37 171.95 172.99 176.07 177.04 163.17 167.46 171.34 171.87 173.14 175.54 176.65 - 80 18 22 24 26 30 20 28 Device Memory (GB)