Performance for Model: 8B, Seq Len: 4096 on RTX5090 176 164.02 167.97 170.27 171.27 172.48 175.39 176.60 180 160 163.77 167.68 170.84 171.39 172.65 175.49 176.13 -160 Host Memory (GB) 12 128 144 163.57 167.66 170.43 171.73 173.28 175.35 176.22 140 164.19 167.52 170.72 171.90 172.57 175.30 176.66 120 163.94 167.65 170.19 171.74 172.83 175.31 176.28 - 100 96 163.84 168.16 171.01 171.73 172.96 175.28 176.68 163.84 168.13 170.89 171.65 172.39 175.45 176.57 - 80 18 20 22 24 26 30 28 Device Memory (GB)