Performance for Model: 8B, Seq Len: 8192 on RTX5090 - 180 96 104.29 158.61 163.61 157.49 160.68 161.81 162.25 80 104.21 157.48 158.00 160.58 161.09 162.11 163.57 -160 Host Memory (GB) 70 75 106.57 158.66 160.36 161.45 162.63 - 140 ص 157.45 163.69 106.71 155.88 154.75 159.37 159.98 162.02 162.78 - 120 65 109.84 156.22 152.21 151.18 155.31 157.06 157.67 - 100 09 109.70 149.52 150.75 151.52 147.58 147.18 152.68 - 80 18 22 24 26 30 20 28 Device Memory (GB)