Performance for Model: 1B, Seq Len: 2048 on RTX5090 176 163.36 167.59 167.73 167.79 167.43 180 160 161.21 167.66 167.61 167.95 167.81 144 163.16 167.49 167.68 167.36 167.56 -160 128 161.50 167.71 167.74 167.29 167.36 Host Memory (GB) 162.14 167.10 167.34 167.15 167.30 140 162.30 167.20 167.26 167.19 167.32 161.35 167.19 167.28 166.97 167.22 - 120 64 161.86 167.02 167.27 166.75 167.13 48 - 100 162.09 166.73 167.44 167.01 167.02 32 161.70 166.99 166.78 167.24 167.04 - 80 16 166.89 166.78 166.90 161.92 167.11 16 20 24 28 30 Device Memory (GB)