Performance for Model: 1B, Seq Len: 512 on RTX5090 176 162.67 169.43 169.45 169.49 169.08 - 180 160 160.94 169.05 169.18 169.12 169.36 144 162.11 169.51 169.24 169.54 169.10 160 128 160.66 169.38 169.31 168.61 169.11 Host Memory (GB) 169.60 161.06 169.74 169.06 168.86 140 160.59 169.07 169.41 168.92 168.83 160.61 168.98 169.14 169.28 168.75 - 120 64 161.21 169.15 168.73 169.01 168.82 48 - 100 161.85 169.40 168.54 168.96 168.23 32 161.05 169.12 168.61 169.12 168.92 - 80 16 160.17 168.64 168.59 168.92 168.95 16 20 24 28 30 Device Memory (GB)