Performance for Model: 1B, Seq Len: 4096 on RTX5090 176 140.48 159.77 159.89 159.70 159.67 160 160.07 159.89 140.45 159.79 159.70 - 160 144 139.62 159.74 159.89 159.84 160.00 128 159.69 139.60 159.60 159.65 159.61 - 150 Host Memory (GB) 140.50 159.73 159.76 159.64 159.87 140.00 159.69 159.54 159.63 159.74 140 139.18 159.55 159.72 159.50 159.72 64 140.95 159.71 159.64 159.68 159.57 48 - 130 159.79 139.85 159.57 159.39 159.63 32 140.21 159.42 159.10 159.56 159.44 16 140.10 159.20 159.05 158.93 159.31 - 120 20 24 28 16 30 Device Memory (GB)