Performance for Model: 1B, Seq Len: 512 on RTX3090 112 7379.97 7425.85 7443.84 7437.25 8000 96 7390.85 7373.68 7317.71 7418.75 7000 Host Memory (GB) 48 64 80 7412.14 7436.62 7387.33 7360.35 Tokens/sec 6000 7396.18 7418.58 7372.51 7398.86 7351.90 7395.00 7353.05 7368.39 5000 32 7343.04 7410.99 7411.10 7366.17 - 4000 16 7419.77 7331.01 7391.28 7384.42 20 16 22 18 Device Memory (GB)