Performance for Model: 1B, Seq Len: 2048 on RTX3090 112 50.60 51.00 50.84 50.75 -60 96 50.52 50.83 50.76 50.95 - 55 Host Memory (GB) 48 64 80 50.58 50.63 50.79 50.74 - 50 50.32 50.77 50.78 50.85 - 40 50.58 50.35 50.30 50.42 - 35 32 50.38 50.26 50.50 50.55 - 30 16 50.22 50.14 50.45 50.11 18 20 22 16 Device Memory (GB)