Performance for Model: 1B, Seq Len: 512 on RTX3090 112 55.46 55.81 55.89 55.94 -60 96 54.99 55.41 55.75 55.54 - 55 Host Memory (GB) 48 64 80 55.70 55.89 55.31 55.52 - 50 55.58 55.75 55.41 55.60 - 40 55.25 55.57 55.26 55.37 - 35 32 55.18 55.36 55.69 55.70 - 30 16 55.76 55.09 55.55 55.49 18 20 22 16 Device Memory (GB)