Performance for Model: 1B, Seq Len: 1024 on H100 56 420.12 436.26 476.89 478.92 478.95 479.07 402.13 650 48 420.39 436.41 402.19 477.79 478.58 478.79 479.41 600 Host Memory (GB) 32 40 420.63 - 550 436.24 402.18 477.98 479.32 478.72 480.04 o TFLOPS/ 420.27 436.27 402.09 479.15 479.87 478.87 479.43 450 436.31 402.36 479.50 481.40 480.43 421.83 479.82 400 16 420.54 436.03 402.32 476.49 480.21 479.26 480.39 - 350 16 32 20 28 24 64 78

Device Memory (GB)