Performance for Model: 1B, Seq Len: 8192 on H100 56 412.73 413.93 438.60 456.96 457.69 458.58 458.51 650 48 411.85 413.71 437.52 457.37 458.68 458.04 458.02 600 Host Memory (GB) 32 40 411.77 - 550 414.55 438.60 456.99 459.35 458.28 458.20 00 TFLOPS/ 411.59 459.09 414.94 438.30 458.87 458.54 458.30 450 412.94 458.87 460.81 458.31 460.09 414.17 438.74 400 16 460.59 412.19 414.38 438.45 455.49 459.95 457.89 - 350 16 20 28 32 24 64 78

Device Memory (GB)