Performance for Model: 1B, Seq Len: 65536 on A100 - 180 56.80 56.81 56.66 56.81 56.93 56.91 160 56.72 56.74 56.84 56.77 56.94 56.95 56.94 56.83 56.93 -140 TFLOPS/s 56.36 56.67 56.82 56.50 56.89 56.85 56.85 56.84 56.92 100 56.19 58.11 56.33 56.71 56.40 56.80 56.76 56.81 56.85 - 80 48.73 52.10 54.22 52.92 55.80 56.19 56.59 56.75 56.78 60 20 24 28 30 40 50 60 70 78

Device Memory (GB)

80

64

Host Memory (GB) 48

32