Performance for Model: 1B, Seq Len: 65536 on H100 18999.9240016.8220062.8250098.2250218.1250234.3250217.559738.73 9 18992.020011.73 **Tokens per Second** - 20000

Device Memory (GB)

Host Memory (GB) 80 96 112