Performance for Model: 8B, Seq Len: 8192 on H100 176 13000 7469.48 8118.46 9190.44 10808.64 10990.16 11001.29 160 7473.79 12000 Host Memory (GB) 112 128 144 11000 **Tokens/sec** 10000 9000 96 - 8000 80 - 7000 40 78 30 50 60 70 Device Memory (GB)