Performance for Model: 8B, Seq Len: 4096 on H100 14000 176 7495.39 8608.15 9704.49 11150.68 11505.45 11538.81 13000 160 7486.68 8603.96 9691.27 11117.70 11535.92 11519.99 12000 Host Memory (GB) 12 128 144 7487.39 11505.49 8603.05 9698.01 11120.12 11547.01 11000 **Tokens/sec** 11503.40 11515.18 7488.02 8606.59 9682.00 11115.42 10000 11498.45 7486.97 8602.52 9686.97 11112.23 11520.37 9000 - 8000 11532.66 11532.94 7495.45 8612.22 9695.78 11112.19 - 7000 10078.74 10211.89 10579.63 11229.58 11520.94 11553.69 30 40 50 60 70 78 Device Memory (GB)