Performance for Model: 8B, Seq Len: 1024 on A100 4000 176 3094.14 3979.34 4018.25 4031.00 4025.72 4032.71 3500 160 3096.21 3983.00 4019.19 4032.82 4027.70 4032.77 Host Memory (GB) 12 128 144 3000 3097.75 3990.24 4034.04 4040.24 4044.72 4042.43 **Tokens/sec** 4035.02 4044.01 4044.40 3088.96 3961.74 4024.36 2500 3977.89 4026.65 3092.49 4011.13 4026.37 4028.32 2000 96 3091.41 4019.97 4002.09 4034.49 3967.97 4028.36 - 1500 80 3093.42 3967.51 4018.49 4026.52 4011.91 4030.30 30 60 50 40 70 78 Device Memory (GB)