Performance for Model: 8B, Seq Len: 8192 on A100 3500 96 2837.90 2916.18 3097.67 3135.17 3141.30 3141.82 3139.26 3000 2841.88 2919.48 3094.72 3133.28 3141.97 3144.26 3144.98 Host Memory (GB) 70 75 2840.64 2917.03 3099.75 3134.64 3142.01 3143.21 3145.12 2500 3139.25 2840.16 2918.58 3098.07 3135.97 3142.59 3144.25 2000 3136.01 2832.57 2919.60 3098.15 3143.19 3147.75 3145.91 - 1500 2821.51 2864.55 3012.49 3122.30 3140.48 3146.32 3141.78 - 1000 24 50 78 30 40 60 70

Device Memory (GB)