Performance for Model: 8B, Seq Len: 8192 on A100

| Host Memory (GB) | Device Memory (GB) | | | | | |
|---|---|---|---|---|---|---|
| | 30 | 40 | 50 | 60 | 70 | 78 |
| 176 | 2861.97 | 3075.10 | 3106.82 | 3108.61 | 3110.41 | 3124.09 |
| 160 | 2866.05 | 3086.07 | 3118.52 | 3123.73 | 3127.39 | 3103.72 |
| 144 | 2866.68 | 3084.02 | 3122.10 | 3127.53 | 3126.71 | 3129.65 |
| 128 | 2863.14 | 3078.96 | 3111.11 | 3110.18 | 3109.35 | 3129.11 |
| 112 | 2860.30 | 3075.73 | 3111.40 | 3117.01 | 3124.26 | 3119.39 |
| 96 | 2858.93 | 3074.02 | 3114.92 | 3119.51 | 3118.16 | 3121.10 |
| 80 | 2859.71 | 3073.46 | 3105.40 | 3117.44 | 3116.34 | 3118.48 |