Performance for Model: 8B, Seq Len: 2048 on H100 7519.27 8632.18 10199.03 11726.39 12285.32 12273.32 14000 160 7527.95 8646.91 10202.83 11722.42 12248.25 12285.74 13000 Host Memory (GB) 12 128 144 7524.94 8640.66 10206.59 12265.83 12296.28 11743.10 12000 Tokens/sec 12283.17 7525.89 8642.06 10202.93 11737.74 12295.91 11000 12268.04 7522.54 8641.40 10210.18 11739.93 12247.91 10000 9 - 7522.39 12243.17 8642.45 10200.06 11730.34 12296.90 - 9000 - 8000 10220.65 10516.59 11105.60 11883.44 12270.10 12287.69 30 40 50 60 70 78 Device Memory (GB)