Performance for Model: 8B, Seq Len: 4096 on H100 14000 7496.51 8615.60 11327.41 11765.84 11834.52 9964.61 13000 160 7506.64 8625.26 11849.76 9955.26 11326.14 11788.38 12000 Host Memory (GB) 12 128 144 7502.73 8619.25 9979.49 11777.07 11828.24 11327.33 11000 11791.1611838.61 7507.30 8628.03 9975.52 11339.60 10000 11840.51 7504.49 8630.10 9963.89 11332.90 11794.37 - 9000 9 - 7502.63 11815.38 8621.72 9975.04 11334.52 11843.23 - 8000 10149.65 10256.18 10743.30 11470.42 11779.25 11805.99 30 40 50 60 70 78 Device Memory (GB)