Performance for Model: 8B, Seq Len: 65536 on H100 7000 5173.18 5466.06 5516.50 5544.92 5481.96 5524.52 6500 5162.49 5406.81 5424.06 5477.10 5490.86 5455.17 6000 5120.60 5399.09 5380.00 5424.84 5423.74 5459.48 5500 5074.65 5323.10 5354.10 5371.53 5394.67 5403.48 5000 5342.34 5028.26 5283.22 5293.50 5303.65 5356.88 - 4500 5000.85 5236.53 5258.99 5275.75 5295.72 5315.17 - 4000 4505.62 4772.27 4834.62 4948.59 4991.86 5037.36 30 40 50 60 70 78

Device Memory (GB)

176

160

Host Memory (GB) 12 128 144

96