Performance for Model: 8B, Seq Len: 1024 on H100 344.72 396.15 536.92 558.29 559.78 800 467.41 344.37 395.59 466.94 535.87 558.68 559.26 700 344.27 558.01 395.62 466.56 536.11 559.81 600 344.34 395.71 467.07 535.99 557.57 558.72 500 344.23 395.95 467.08 535.69 557.98 558.54 400 344.70 396.11 536.65 558.69 559.61 467.37 - 300 468.01 480.32 506.76 542.04 558.50 558.85 60 50 30 78 40 70

Device Memory (GB)

176

160

Host Memory (GB) 12 128 144

96

80