Performance for Model: 8B, Seq Len: 512 on A100 4000 3097.58 4020.04 4099.00 4105.66 4111.05 4113.98 3500 3103.75 4021.25 4093.25 4103.37 4108.57 4110.01 3107.51 4043.62 4127.45 4139.31 4142.91 4131.98 3000 **Tokens/sec** 4040.02 4121.50 4140.26 4141.76 3099.35 4131.97 2500 4106.90 4115.66 3104.34 4030.14 4118.95 4116.12 2000 3103.88 4123.44 4135.23 4025.22 4133.16 4138.31 - 1500 3102.51 4023.87 4122.98 4130.38 4121.95 4135.00

60

78

70

50

Device Memory (GB)

176

160

Host Memory (GB) 12 128 144

96

80

30

40