Performance for Model: 8B, Seq Len: 2048 on RTX5090 160.36 166.89 171.96 177.76 181.46 182.59 183.93 180 161.18 167.65 181.48 183.15 172.54 175.73 184.46 -160 161.37 166.28 172.75 176.31 182.01 182.36 184.35 140 160.88 167.34 172.31 177.52 181.40 182.20 184.45 - 120 160.77 167.21 172.25 177.34 181.81 182.43 183.90 - 100 161.44 167.37 172.25 177.59 181.80 182.47 183.94 160.37 166.83 172.16 177.24 181.24 182.92 183.40 - 80

26

30

28

176

160

Host Memory (GB) 12 128 144

96

18

20

22

24

Device Memory (GB)