Performance for Model: 8B, Seq Len: 2048 on RTX5090 - 180 153.71 159.38 164.19 168.35 175.86 177.70 179.64 153.26 159.87 163.72 167.79 175.18 177.74 179.34 -160 157.22 158.95 164.00 168.11 175.11 178.07 179.82 140 157.14 169.59 171.04 173.13 176.11 177.78 179.10 - 120 159.62 166.29 170.13 171.01 174.04 174.42 176.70 - 100 159.42 161.05 163.42 164.27 166.73 168.52 169.51 - 80

26

30

28

96

80

Host Memory (GB) 70 75

65

09

18

20

22

24

Device Memory (GB)