Performance for Model: 1B, Seq Len: 32768 on H100 0.43 0.46 0.46 0.46 0.46 0.47 0.47 0.47 0.46 8.0 0.43 0.45 0.45 0.45 0.46 0.46 0.46 0.46 0.45 0.43 0.46 0.46 0.46 0.46 0.47 0.47 0.46 0.46 0.70.43 0.45 0.45 0.45 0.46 0.46 0.46 0.46 0.46 0.43 0.46 0.46 0.46 0.46 0.47 0.47 0.47 0.46 - 0.6 0.43 0.45 0.45 0.45 0.46 0.46 0.46 0.46 0.46 - 0.5 0.43 0.46 0.46 0.46 0.46 0.47 0.47 0.47 0.46 0.43 0.45 0.45 0.46 0.47 0.46 0.46 0.46 0.46 -0.40.43 0.46 0.46 0.46 0.46 0.47 0.46 0.47 0.47 0.43 0.45 0.45 0.45 0.46 0.46 0.46 0.46 0.46 0.3 0.46 0.47 0.47 0.47 0.46 0.47 0.47 0.47 0.47 24 28 30 40 50 60 70 78 20

Device Memory (GB)

176

160

144

128

Host Memory (GB) 80 96 112

64

48

32

16