Performance for Model: 1B, Seq Len: 65536 on H100 0.41 0.41 0.41 0.41 0.39 0.42 0.41 0.41 8.0 0.39 0.41 0.41 0.41 0.41 0.41 0.42 0.40 0.40 0.39 0.41 0.41 0.41 0.40 0.41 0.42 0.42 0.41 -0.70.41 0.41 0.40 0.39 0.41 0.41 0.41 0.42 0.42 0.39 0.41 0.41 0.41 0.41 0.41 0.42 0.41 0.42 - 0.6 0.39 0.41 0.41 0.41 0.41 0.41 0.42 0.42 0.42 - 0.5 0.39 0.41 0.41 0.41 0.41 0.41 0.42 0.42 0.42 0.39 0.41 0.41 0.41 0.41 0.42 0.41 0.42 0.42 -0.4 0.39 0.41 0.41 0.41 0.41 0.41 0.41 0.42 0.42 0.39 0.43 0.41 0.41 0.40 0.41 0.41 0.42 0.41 0.3 0.36 0.37 0.38 0.38 0.40 0.41 0.41 0.41 0.42 40 20 24 28 30 50 60 70 78

Device Memory (GB)

176

160

144

128

Host Memory (GB) 80 96 112

64

48

32

16