Performance for Model: 1B, Seq Len: 65536 on H100 0.43 0.45 0.46 0.46 0.46 0.46 0.46 0.45 8.0 0.43 0.45 0.45 0.46 0.46 0.46 0.46 0.45 0.44 0.43 0.45 0.45 0.45 0.45 0.46 0.46 0.46 0.46 0.7 0.43 0.45 0.45 0.45 0.45 0.46 0.46 0.46 0.46 0.43 0.45 0.45 0.46 0.46 0.46 0.46 0.46 0.46 - 0.6 0.43 0.45 0.45 0.46 0.46 0.46 0.46 0.46 0.46 - 0.5 0.43 0.45 0.45 0.46 0.46 0.46 0.46 0.46 0.46 0.43 0.45 0.46 0.46 0.46 0.46 0.46 0.46 0.46 -0.40.43 0.45 0.45 0.46 0.46 0.46 0.46 0.46 0.46 0.44 0.49 0.46 0.46 0.46 0.46 0.46 0.46 0.46 0.3

0.46

50

0.46

60

0.46

70

0.46

78

176

160

144

128

Host Memory (GB) 80 96 112

64

48

32

16

0.44

20

0.46

24

0.45

28

0.45

30

0.46

40

Device Memory (GB)