Performance for Model: 1B, Seq Len: 4096 on H100 176 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 8.0 160 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 144 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 -0.7128 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 Host Memory (GB) 80 96 112 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 - 0.6 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 - 0.5 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 64 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 -0.448 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 32 0.31 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 - 0.3 16 0.35 0.42 0.47 0.47 0.47 0.47 0.47 0.47 0.47 50 28 30 40 60 70 78 20 24 Device Memory (GB)