Performance for Model: 1B, Seq Len: 16384 on H100 33756.3421578.3421654.5451776.1461811.3401765.1481602.9471585.3461798.39 9 33725.5**4**1512.84 Host Memory (GB) 80 96 112 128 000 Tokens per Second

Device Memory (GB)