Performance for Model: 1B, Seq Len: 65536 on RTX5090 176 6364.64 6427.70 6444.76 6460.52 6461.66 9000 160 6372.00 6440.45 6459.93 6465.31 6458.17 144 6367.61 6424.60 6441.65 6466.55 6456.15 8000 128 6363.96 6418.52 6443.81 6451.28 6461.76 Host Memory (GB) 7000 6364.50 6429.53 6436.83 6458.39 6453.63 6423.22 6431.00 6453.49 6449.06 6360.33 6000 6353.96 6418.34 6437.02 6434.91 6424.04 64 6350.25 6394.42 6414.66 6441.82 6438.61 - 5000 48 6433.90 6305.45 6364.07 6413.89 6440.60 6197.32 6325.24 6354.47 6341.82 6348.41 - 4000 5864.63 6126.85 6160.62 6232.40 6046.57 16 30 28 20 24 Device Memory (GB)