Performance for Model: 1B, Seq Len: 4096 on H100

Heatmap with Device Memory (GB) on x-axis (20, 24, 28, 30, 40, 50, 60, 70, 78) and Host Memory (GB) on y-axis (16, 32, 48, 64, 80, 96, 112, 128, 144, 160, 176). Color scale represents Tokens per Second (from ~50000 to 100000).

Values shown:
Row 176: 36158.94, 49835.25, 55716.25, 54916.55, 55571.15, 55673.95, 55612.95, 55844.15, 55811.46
Row 160: 36179.54, 29801.17