Performance for Model: 1B, Seq Len: 4096 on H100 \mathcal{L} - $\frac{50427.99}{6}$ 48850.42 50409.83 56624.47 56734.30 56770.59 56857.54 80000 75000 $\frac{\infty}{7}$ - 50436.25 48916.52 50390.50 56619.73 56843.54 56784.85 56786.93 70000 Q -<mark>50331.86 48911.77 50376.35</mark> 56660.41 56818.44 56770.04 56889.87 65000 -<mark>50425.41 48929.97 50349.97</mark> 56862.38 56769.61 56851.66 56930.41 60000 - 55000 ₹ -<mark>50535.92 48877.04 50432.27</mark> 56857.31 57036.42 56845.05 57054.52 - 50000 <u>9</u> -<mark>50392.83 48896.47 50375.64</mark> 56590.97 56995.94 56775.45 57126.12 - 45000 16 64 20 28 32 24 78 Device Memory (GB)

Host Memory (GB)