Performance for Model: 8B, Seq Len: 1024 on A100 4000 3308.29 3672.17 3994.06 4056.47 4067.64 4071.14 4069.93 3500 3310.87 3669.10 3994.15 4056.95 4071.50 4074.85 4072.64 3000 Host Memory (GB) 70 75 3311.49 3668.78 3994.14 4058.36 4074.71 4071.90 4074.58 **Tokens/sec** 2500 3309.53 4060.99 4075.99 3671.15 4001.13 4073.94 4073.68 2000 4058.39 4069.08 4076.06 3670.76 3990.86 4078.16 3330.98 - 1500 3581.73 3659.51 3850.88 4039.21 4065.61 4075.83 4070.62 24 30 40 50 60 70 78 Device Memory (GB)