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