Performance for Model: 8B, Seq Len: 2048 on H100