Input vs Output Comparison of LLaMA-3-8B for Batch Size = 1 on One A100 using TensorRT-LLM (fp16)

