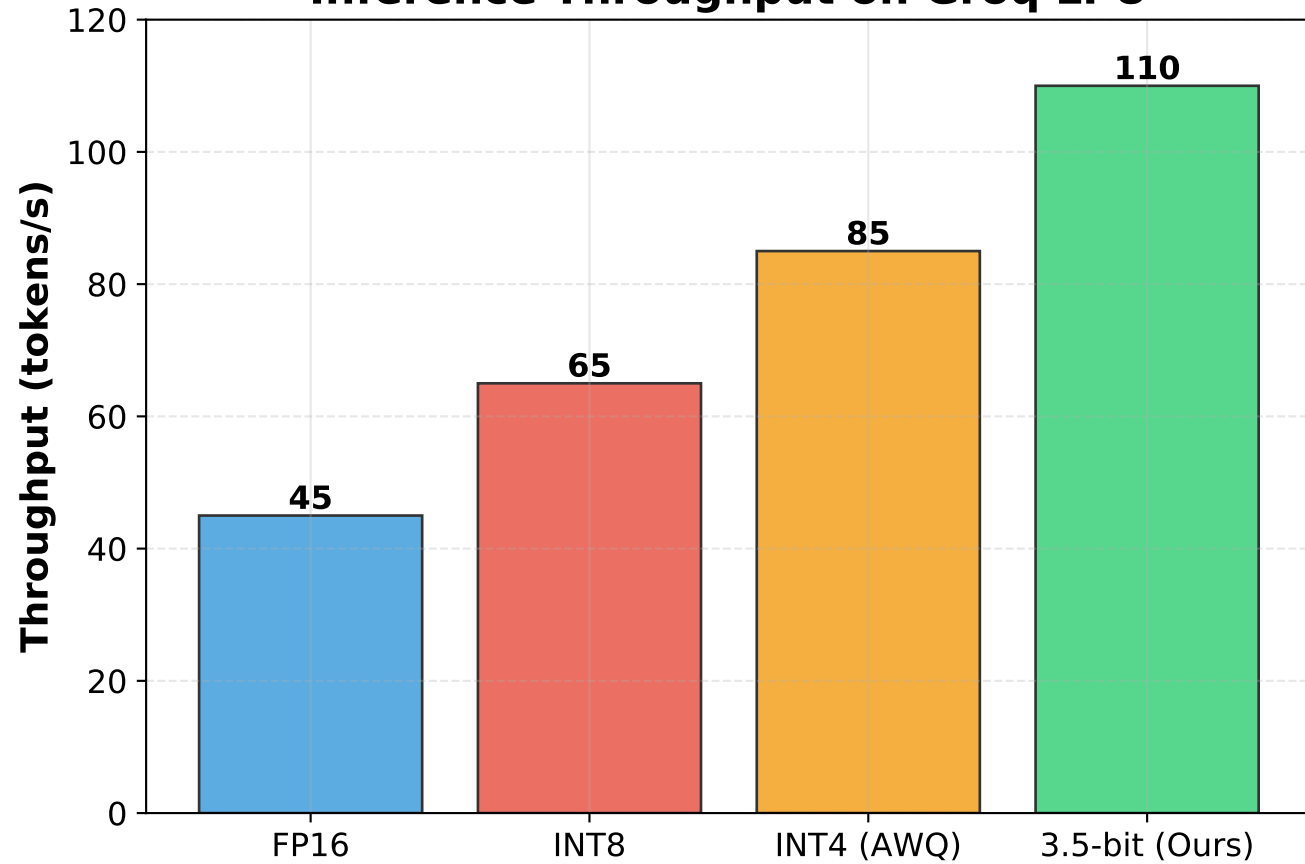


Inference Throughput on Groq LPU



Memory Footprint (LLaMA-70B)

