Expert Parallel (EP) vs Iteration Latency ( ↓ is better)

GPT-3 350M MoE on 4x4 A100 GPUs