**Paper Baselines vs Our Batch-Aware Contributions**
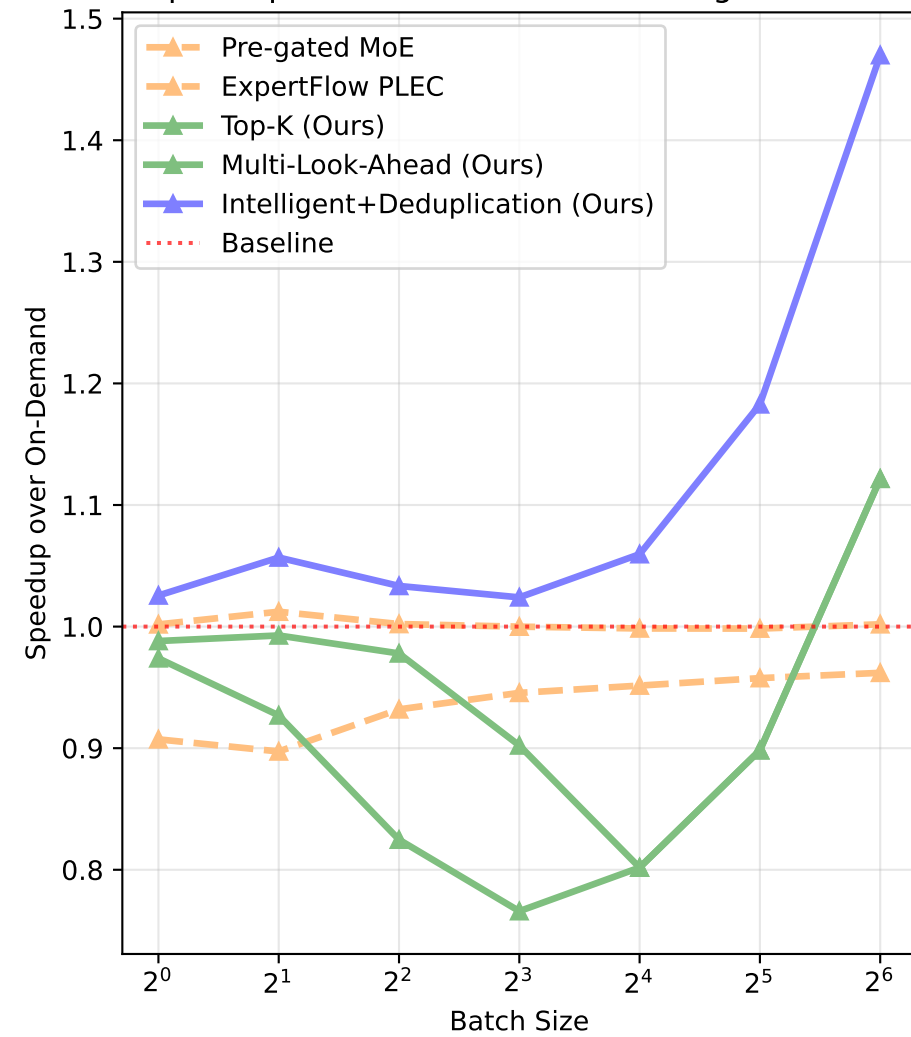
Latency Scaling: Papers Degrade, Ours Improve

Legend: On-Demand, Pre-gated MoE, ExpertFlow PLEC, Top-K (Ours), Multi-Look-Ahead (Ours), Intelligent+Deduplication (Ours)

Hit Rate: Papers Degrade, Ours Maintain

Legend: On-Demand, Pre-gated MoE, ExpertFlow PLEC, Top-K (Ours), Multi-Look-Ahead (Ours), Intelligent+Deduplication (Ours)

Speedup: Our Methods Excel at Large Batches

Legend: Pre-gated MoE, ExpertFlow PLEC, Top-K (Ours), Multi-Look-Ahead (Ours), Intelligent+Deduplication (Ours), Baseline

Our Key Innovation: Expert Deduplication

Memory Savings from Deduplication (%): 0.0%, 0.9%, 2.3%, 5.2%, 11.0%, 20.8%, 36.5%