**Congestion Mitigation (Queue Depth)**

Queue Depth

- Baseline (1 Replica)
- WVA

**Resource Efficiency (KV Cache)**

KV Cache Utilization (%)

- Baseline (12 Replicas)
- WVA

**Scaling & Energy Consumption**

Replica Count

Time (s)

- Baseline (Over-provisioned)
- WVA Replicas
- Energy Savings