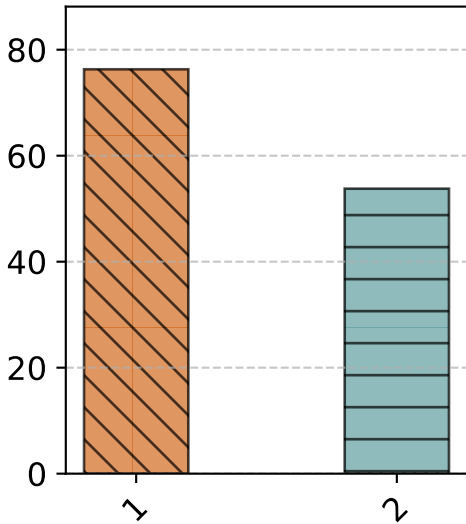


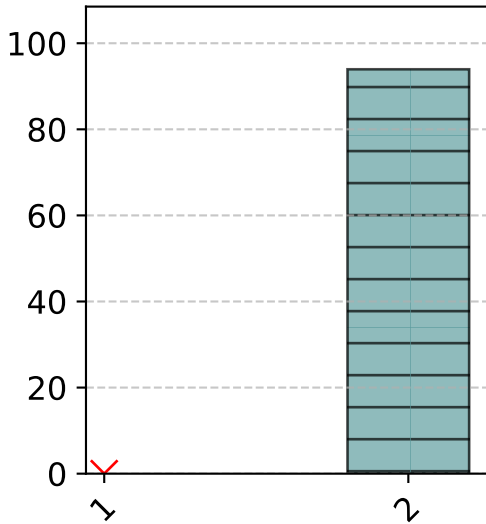
Sequence Length: 4M

Decode Latency (ms)



KV Parallel Degree

Sequence Length: 10M



KV Parallel Degree