

Per-token latency (ms)

