Single-Token Decode Throughput vs Context Length