**Figure 6a: Per-Token Latency Comparison**