**Phase 2: query-only TTT**

$N_{\text{qTTT}} \times$

Cross-entropy Loss

Logits

Scaled Dot-Product Attention

$\nabla_Q \mathcal{L}_{CE}$

Query vectors

Query weight $W_Q$

$x_5$ $x_6$

Sampled span at TTT step $t$

cheap and fast, computed $N_{\text{qTTT}}$ times

**Phase 1: Prefill**

KV vectors

KV weights $W_K, W_V$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$

Long-context tokens

expensive, computed only once