

KernelView of Attention

- Suppose ρ is such that $\text{sim}(q, k) = \langle \rho(q), \rho(k) \rangle$

• If $Q' \equiv \varphi(Q)$ and $K' \equiv \varphi(K)$, output is

• Why write this way?

- **Linear time algorithm for computing $LT(A \cdot B^T) \cdot C$**

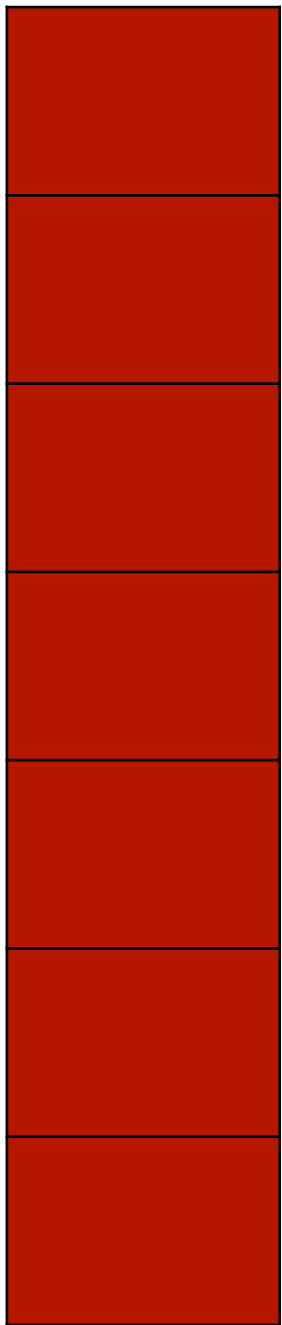
- Runtime depends on dimension of $\varphi(\cdot)$

- What about ρ for softmax?

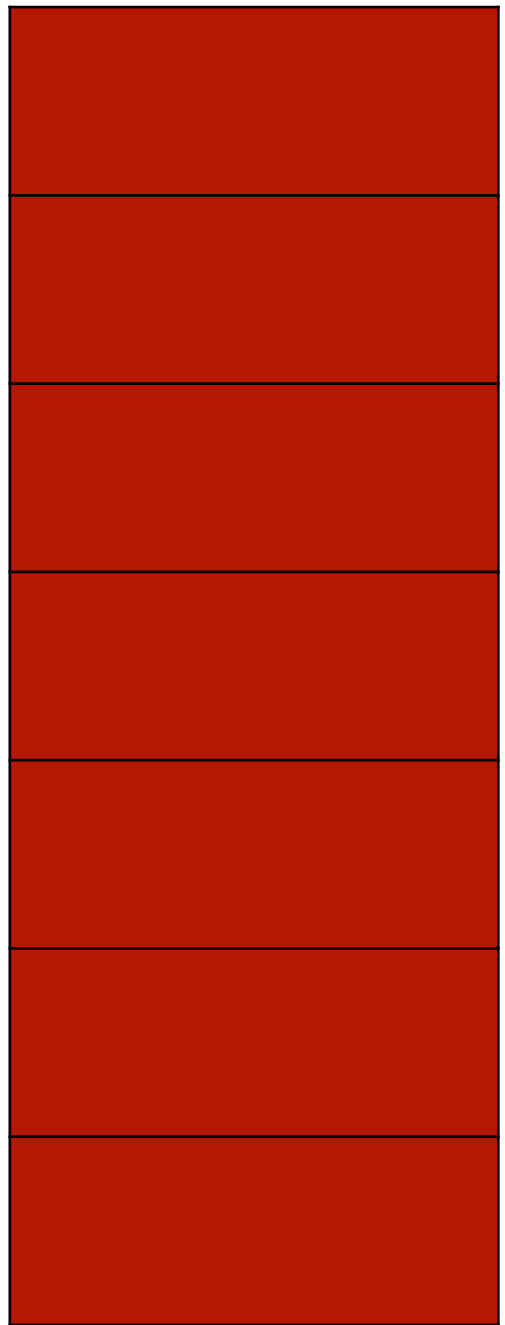
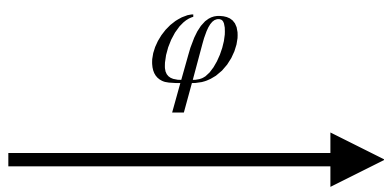
• **No finite dimensional feature maps**



$$D^{-1} \cdot \text{LT}(Q') \cdot (K')^T \cdot V$$



Q



Q'

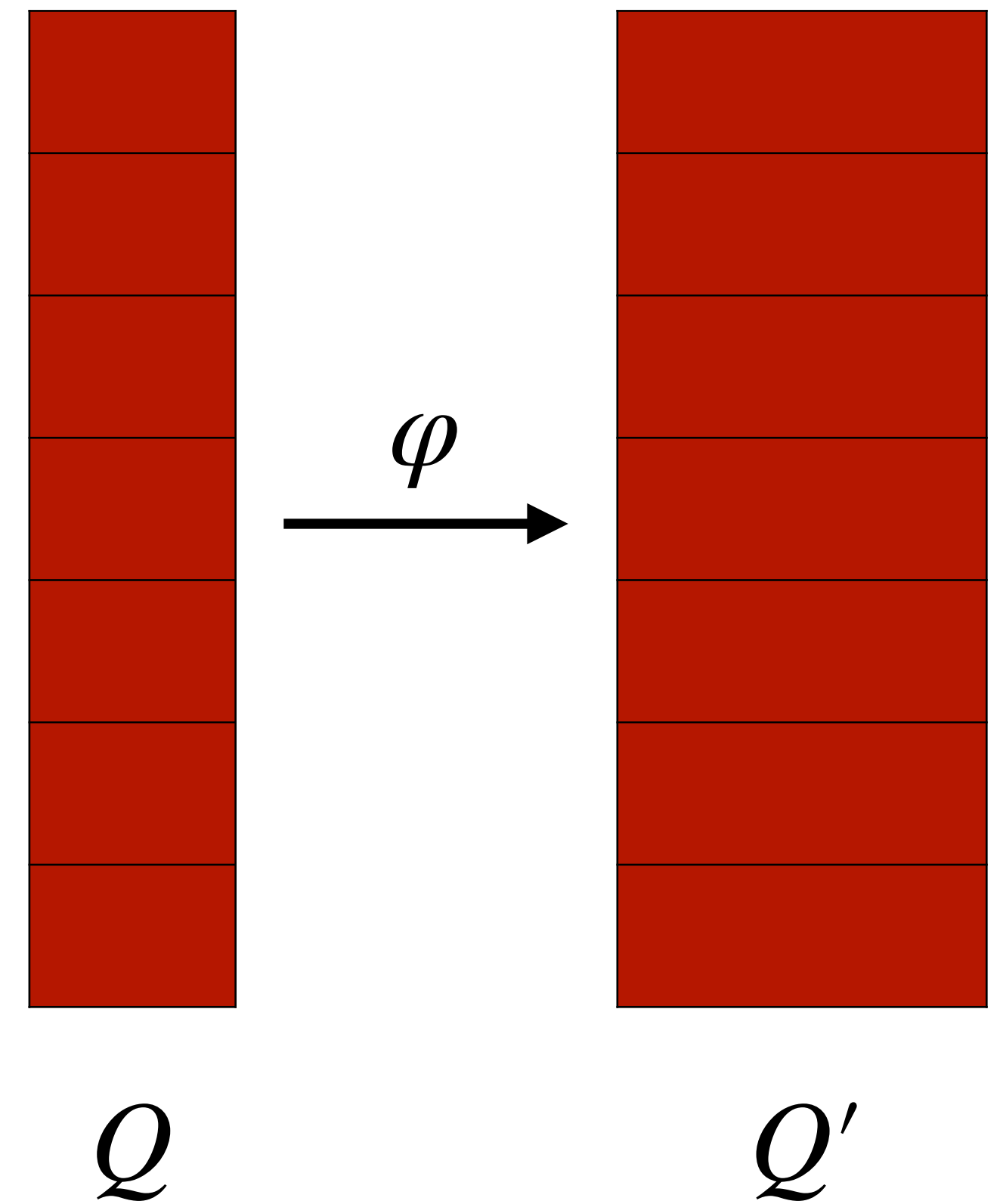
Kernel View of Attention

- Suppose φ is such that $\text{sim}(q, k) = \langle \varphi(q), \varphi(k) \rangle$

- If $Q' = \varphi(Q)$ and $K' = \varphi(K)$, output is

$$D^{-1} \cdot \text{LT}(Q' \cdot (K')^{\top}) \cdot V$$

- Why write this way?
 - **Linear time algorithm** for computing $\text{LT}(A \cdot B^{\top}) \cdot C$
 - Runtime depends on dimension of $\varphi(\cdot)$
- What about φ for softmax?
 - No finite dimensional feature maps



Previous Works