# CS371N Lecture 11
Transformers, Transformer Language Modeling

## Announcements
- A3 out

Recap Attention: places a probability distribution over a sequence of $n$ tokens with embeddings $e_1 \ldots e_n$

Simplified version:

① Form keys $\quad K_i = W^K e_i$

$$K = E(W^K)^T$$

query $q$ $\qquad\qquad\qquad$ S $[0\ 0\ 1\ 0]$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \downarrow$

② Compute scores $\quad s_i = K_i^T q \quad \alpha \ [\frac{1}{6}\ \frac{1}{6}\ \frac{1}{2}\ \frac{1}{6}]$

③ Compute attn weights $\quad \alpha = \text{softmax}(s)$

④ Result (output) $= \sum \alpha_i e_i$

## Self-attention

$E$: seq len $\times d$

$E$ now gives rise to $q_i$ and $k_j$ for each word

$W^K$: $d \times d$ matrix $\quad K = E(W^K)^T$

$W^Q$: $d \times d$ matrix $\quad Q = E(W^Q)^T$

$K, Q$: seq len $\times d$

$$S = QK^T \qquad S_{ij} = q_i \cdot k_j$$

Suppose $E = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 1 & 0 \end{bmatrix}$ $\qquad W^Q = W^K = I$

$$S = \begin{bmatrix} 1 & 1 & 0 & 1 \\ 1 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 \end{bmatrix}$$

matrix of similarities

$$A = \text{softmax}(S) = \begin{bmatrix} 3/10 & 3/10 & 1/10 & 3/10 \\ 1/6 & 1/6 & 1/2 & 1/6 \end{bmatrix}$$

Last step:

$$\text{Output} = A \left( E (W^v)^T \right)$$

$$\left( \text{seq len} \times \text{seq len} \right) \cdot \left( \text{seq len} \times d \right) \cdot \left( d \times d \right)$$

$$\text{Out} = \text{seq len} \times d$$

A takes a weighted sum of values according to attention weights at each position

First row of output $= \frac{3}{10} \cdot v_1 + \frac{3}{10} \cdot v_2$
$$+ \frac{1}{10} \cdot v_3 + \frac{3}{10} \cdot v_4$$

Third row $= \frac{1}{6} \cdot v_1 + \frac{1}{6} v_2 + \frac{1}{2} v_3 + \frac{1}{6} v_4$

$e_1 \quad e_2 \quad e_3 \quad e_4$