# CS371N Lecture 11
## Transformers, Transformer Language Modeling

Recap   Attention over a sequence of
         $n$ tokens with embeddings
         $e_1 \cdots e_n$

① Form keys $= W^k e_i$
   query $= q$

② Scores $s_i = k_i^T q$



A A B A

③ Attention weights (probs) $\alpha = \text{softmax}(s)$

④ Result (output) $= \sum \alpha_i e_i$
                        $\wedge$
         we'll throw a matrix
                   here later

# Today

- Self-attention recap
- Exercises
- Multi-head self-attention
- Transformers
- Language modeling

# Self-attention

Idea: all words are now keys and queries simultaneously

$E$: seq len $\times d$ matrix

$W^K$: $d \times d$     $K = E(W^K)^T$ $\left.\vphantom{\begin{array}{c}a\\b\\c\end{array}}\right\}$ same as before

$Q$: seq len $\times d$     $\left( Q = E(W^Q)^T \right)$

scores

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

len x len

$$A = softmax(S) \text{ by rows}$$

distribution $A_i$ for each word's query
$$q_i$$

<u>Ex</u>   $A = [1 \ 0]$    $B = [0 \ 1]$

A B ← sequence                    "boosted"
                                      identity

$$E = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{matrix} \leftarrow A \\ \leftarrow B \end{matrix} \qquad W^K = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$$

$$Q = \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{matrix} \leftarrow (\text{Greg says}) \\ \text{"find Bs"} \end{matrix}$$

$$K = E(W^K)^T = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix} = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$$

$$S = Q K^T$$

$$= \begin{bmatrix} 0 & 1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix} = \begin{bmatrix} 0 & 10 \\ 0 & 10 \end{bmatrix}$$

$$A = \text{softmax}(S) = \begin{bmatrix} 0 & 10 \\ 0 & 10 \end{bmatrix} \begin{array}{l} \to sm \to \\ \to sm \to \end{array} \overset{A \quad B}{\begin{bmatrix} 0 & 0.999 \\ 0 & 0.999 \end{bmatrix}}$$

— Big K made our probs. peaked

— Q had B for each row $\Rightarrow$ prob on B

$$\overset{A \quad A \quad B \quad A}{\quad}$$

$$W^K = \begin{bmatrix} 10 & 0 \\ 0 & 10 \end{bmatrix}$$

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

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

$$K = E(W^K)^T = \begin{bmatrix} 10 & 0 \\ 10 & 0 \\ 0 & 10 \\ 10 & 0 \end{bmatrix}$$

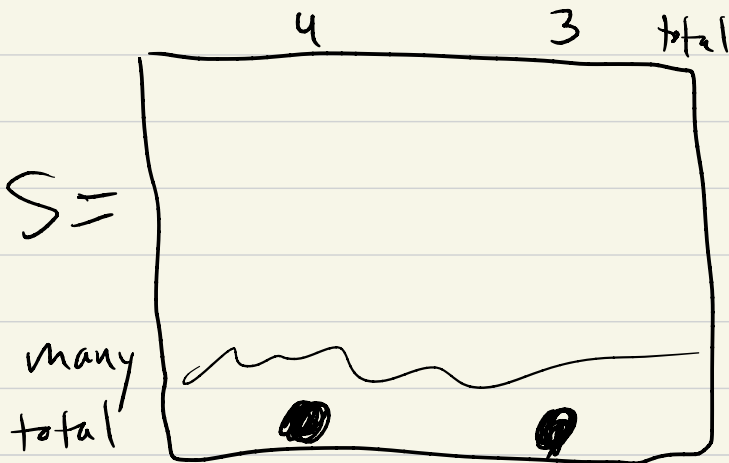$$S = \begin{bmatrix} 0 & 0 & 10 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 10 & 0 \\ 0 & 0 & 10 & 0 \end{bmatrix}$$

$4 \times 4$

Softmax
$\to$ high prob on B

S: for word i, how much does it "attend" to word j

Mary had 4 apples. Jane had 3. How many total? ___

12 words

12×12 matrix

S =

|   | 4 | 3 | total |
|---|---|---|-------|

many total

result vector

$$\sum_{i=1}^{seq\ len} A_{total,\ i} \cdot E_i$$