Natural Language processing
- ⤷ Machine translation
- ⤷ Parts of speech tagging
- ⤷ Generating missing words
- ⤷ Autocomplete

] Understanding language

Rule → $x_{rule}$

$x_n$

$\|x_{rule} - x_{law}\|$

$< \|x_{chows} - x_{law}\|$

vectors
or a sequence
of number

embedding
vector

$x$

$a(it, Law) = 0.8$
attention

$a(it, application)$
$0.1$

The Law will never be perfect , but its application should be just ' this is what we are missing , in my opinion . <EOS> <pad>

$v_1$

The Law will never be perfect , but its application should be just ' this is what we are missing , in my opinion . <EOS> <pad>

$v_n$

$v_{73}$

$h = a(W x + b)$

$h_1$  $h_2$  $h_3$ · · · · $h_n$
O    O O         O

$\omega_{11}$  $\omega_{12}$

O  O O    O · · · · O
$x_1$ $x_2$ $x_3$  · ·    $x_n$

Self attention

$$\underline{v}_1 \quad \square^{\underline{v}_2} \quad \square^{\underline{v}_3} \quad \cdots \quad \square^{\underline{v}_N}$$

$$\underline{v}_n = \Omega_v \, \underline{x}_n + \beta_v$$
$$D \times 1 \qquad D \times D \quad D \times 1 \quad D \times 1$$

$$\underline{x}_1 \quad \underline{x}_2 \quad \underline{x}_3 \qquad \square^{\underline{x}_N}$$

$$sa(\underline{x}_m) = \sum_{n=1}^{N} a(\underline{x}_n, \underline{x}_m) \, \underline{v}_n$$
$$a \in \mathbb{R}$$

$D = 4$

$N = 3$

a) Inputs  Values  Outputs

$a(x_1, x_1)$  sa[$x_1$]  0.1
$a(x_1, x_2)$  0.3
$a(x_1, x_3)$  0.6
sa[$x_2$]
sa[$x_3$]

$\underline{v}_1, \underline{v}_2, \underline{v}_3$

word embedding

b) Inputs  Values  Outputs
sa[$x_1$]  0.5
sa[$x_2$]  0.2
sa[$x_3$]  0.3

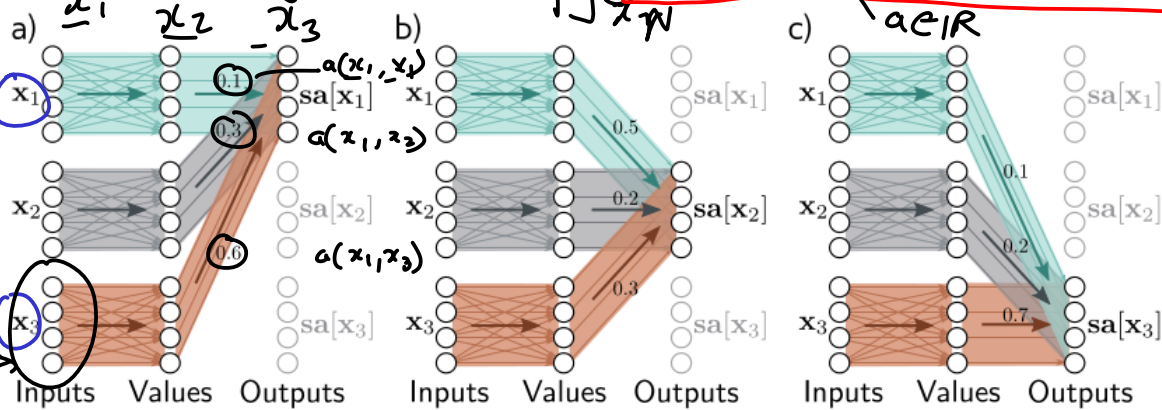c) Inputs  Values  Outputs
sa[$x_1$]  0.1
sa[$x_2$]  0.2
sa[$x_3$]  0.7

$$a(\underline{x}_n, \underline{x}_m) \in [0, 1] \qquad \forall \, n \in [1, N] \text{ and } m \in \{1, N\}$$

$$\sum_{m=1}^{N} a(\underline{x}_n, \underline{x}_m) = 1$$
$$\uparrow$$
$$\text{input}$$

$$a\left(\underline{x}_n, \underline{x}_m\right)$$

value
vector

$$\cdot \underline{v}_n = \underline{\Omega}_v \underline{x}_n + \underline{\beta}_v$$

for each word m the
sentence
$$\forall \; n \in [1, N]$$

key
vector

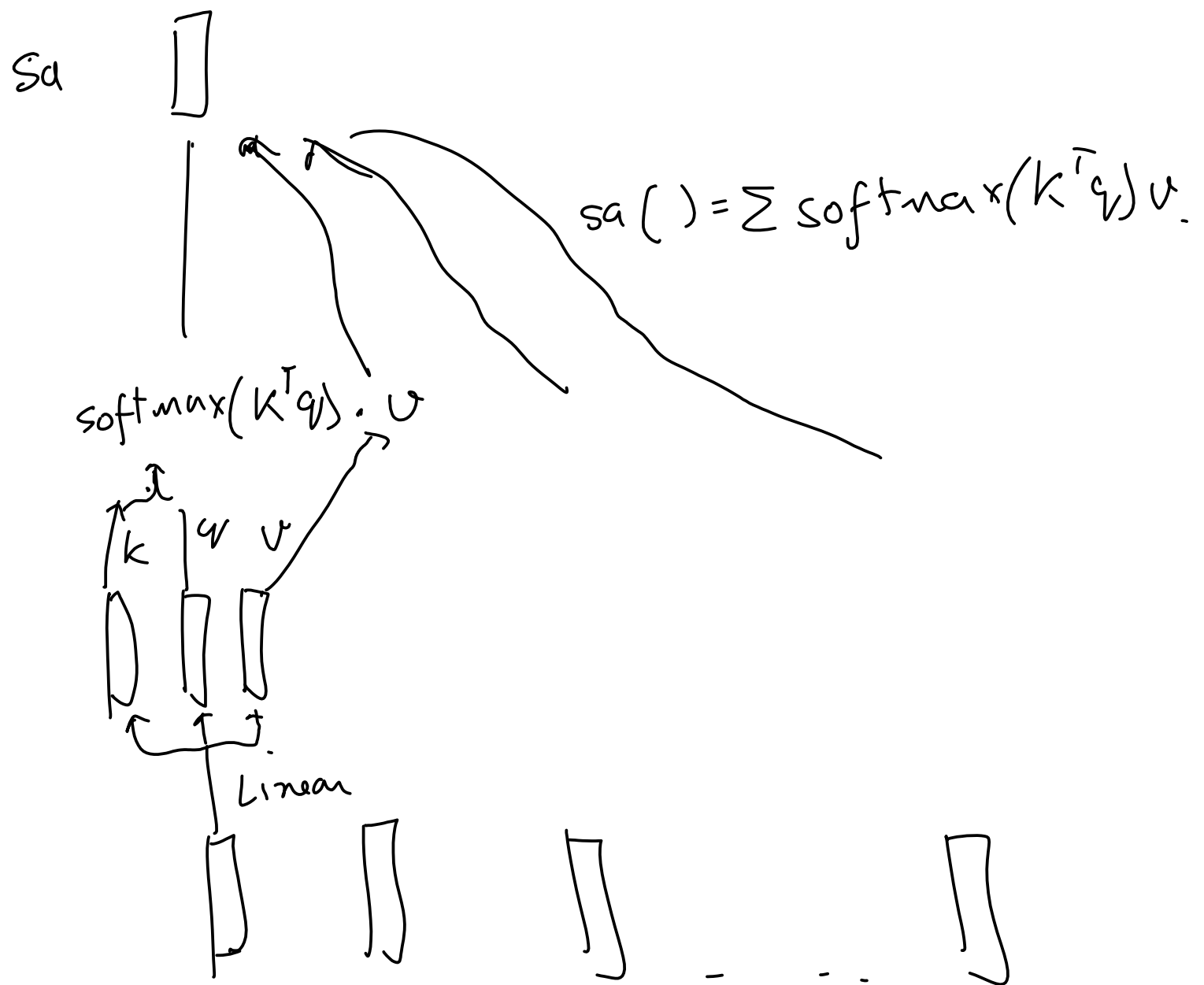$$\underline{k}_n = \underline{\Omega}_k \underline{x}_n + \underline{\beta}_k$$

query
vector

$$\underline{q}_n = \underline{\Omega}_q \underline{x}_n + \underline{\beta}_q$$

$$a\left(\underline{x}_n, \underline{x}_m\right) = \text{softmax}_m\left(\underline{k}_m^T \underline{q}_n\right)$$

output input

$$= \frac{\exp\left(\underline{k}_m^T \underline{q}_n\right)}{\sum_{m'=1}^{N} \exp\left(\underline{k}_m^T \underline{q}_n\right)}$$

$$sa\left(\underline{x}_n\right) = \sum_{m=1}^{N} \text{softmax}_m\left(\underline{k}_m^T \underline{q}_n\right) \underline{v}_m$$

$$\in \mathbb{R}$$

$$\underline{v}_m = \underline{\Omega}_v \underline{x}_m + \underline{\beta}_v$$

sa

$$sa\,(\ ) = \sum \text{softmax}(k^T q)\, v.$$

$$\text{softmax}(k^T q) \cdot v$$

$k \quad q \quad v$

Linear

key
query $\}$ softmax $(\underline{K^T q})$

value

Database

key - value pair

question — answer pair

query is a new key that
you want to search in the database

Self-attention

$q_n$

Queries,
$\mathbf{Q} = \boldsymbol{\beta}_q \mathbf{1}^T + \boldsymbol{\Omega}_q \mathbf{X}$

Keys, $K_n$
$\mathbf{K} = \boldsymbol{\beta}_k \mathbf{1}^T + \boldsymbol{\Omega}_k \mathbf{X}$

Values, $V_n$
$\mathbf{V} = \boldsymbol{\beta}_v \mathbf{1}^T + \boldsymbol{\Omega}_v \mathbf{X}$

Input, $\mathbf{X}$

$x_n$

word sentence

$N = 4096$

sum to 1   $\underline{k}_i^T \underline{q}_j$

Attention,
$\mathbf{Softmax}\left[\mathbf{K}^T\mathbf{Q}\right]$

Output,
$\underset{D \times N}{\mathbf{V}} \cdot \underset{N \times N}{\mathbf{Softmax}\left[\mathbf{K}^T\mathbf{Q}\right]}$

$\in \mathbb{R}_{D \times N}$

$$\underline{\beta}_v \mathbf{1}^T$$
$$= \underline{\beta}_v \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \end{bmatrix}_N$$
$$= \begin{bmatrix} \underline{\beta}_v & \underline{\beta}_v & \underline{\beta}_v & \dots & \underline{\beta}_v \end{bmatrix}^N$$

$$sa(\underline{x}_n) = \sum_{m=1}^{N} softmax_m\left(\underline{k}_m^T \underline{q}_n\right) \underline{v}_m$$

$$\in \mathbb{R}$$

Scaled Dot product self attention

$$Sa[X] = V[X]\, Softmax\left(K(X)^T Q(X)\right)$$

$$Var\left(\underset{D \times N}{K(X)}\right) = 1$$
$$Var\left(\underset{D \times N}{Q(X)}\right) = 1$$

$$Var\left(\underset{N \times D}{K(X)^T}\, \underset{D \times N}{Q(X)}\right)$$
$$= D$$

$$Sa[X] = V[\dot{X}] \text{ softmax} \left( \frac{K(x)^T Q(x)}{\sqrt{D}} \right)$$

# Positional encoding

The sentence) The woman ate the raccoon has a quite different meaning to The raccoon ate the woman.
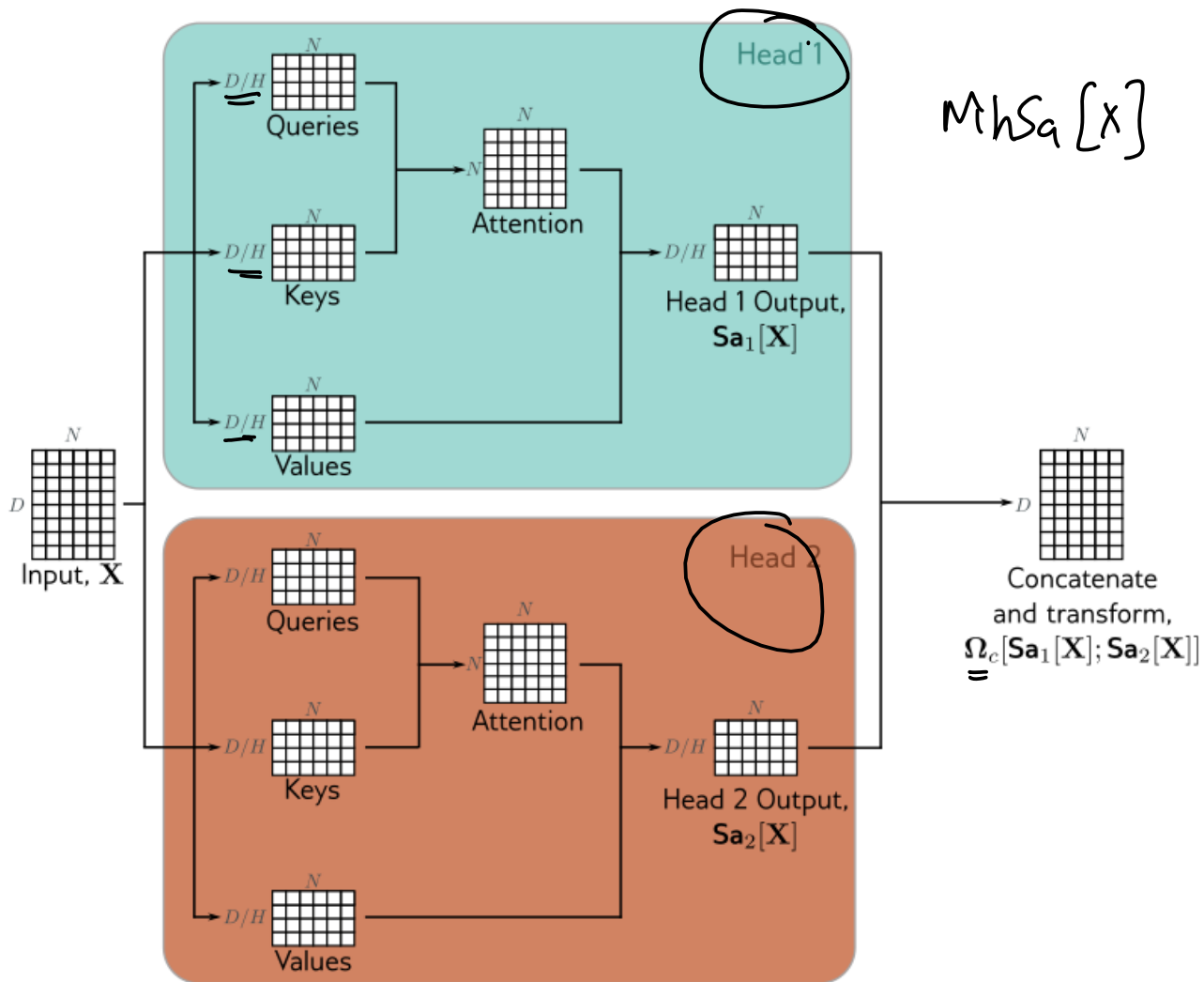
$x_1$ for The

$x_2$ for Raccoon

$$x_i = \begin{bmatrix} & & \\ D & & \\ & & \end{bmatrix} WE \qquad i \uparrow \begin{bmatrix} PE \end{bmatrix} D'$$

$$PE_{(pos, 2i)} = \sin\left(pos/10000^{2i/D'}\right)$$

$$PE_{(pos, 2i+1)} = \cos\left(pos/10000^{2i+1/D'}\right)$$

Scaled dot product self-attention

Input, $\mathbf{X}$

Head 1

Queries

Keys

Attention

Values

Head 1 Output, $\mathbf{Sa}_1[\mathbf{X}]$

Head 2

Queries

Keys

Attention

Values

Head 2 Output, $\mathbf{Sa}_2[\mathbf{X}]$

Concatenate and transform, $\underline{\Omega}_c[\mathbf{Sa}_1[\mathbf{X}]; \mathbf{Sa}_2[\mathbf{X}]]$

$MhSa[x]$

Transformer layer



Residual connection    Residual connection

Multi-head self-attention    LayerNorm    Parallel neural networks $(\times N)$    LayerNorm

Input    Output

# Layer Norm

$$Batch\ Norm\ (x_n) = \frac{x_n - \mu}{\sigma}$$

$$\mu = \frac{1}{B}\sum_{b=1} x_b \qquad \sigma^2 = \frac{1}{B}\sum (x_b - \mu)^2 \qquad x_1, x_2, x_3 \dots x_B$$

Layer Norm

the mean and variance are computed over the "channel" dimension

a) Linear layer

$$x_1 = \begin{bmatrix} \\ \\ \end{bmatrix} D \qquad h_\ell = \begin{bmatrix} \\ \\ \end{bmatrix} \begin{matrix} \# \ of \\ hidden \\ units \end{matrix}$$

b) Conv layer

$$\xrightarrow{Conv2D(3,16)}$$

W    H    $\underset{RGB}{\longleftrightarrow}$ = channels = 3    16 = channels

c) NLP os Sq

Channels $= N =$ number of word in
the sentence

Word embeddings

Transformer block (×K)

layer norm

Linear + softmax

Probability of masked token

&lt;cls&gt;
The
&lt;mask&gt;
pulled
into
&lt;mask&gt;
station

a
aardvark
abacus
...
zeta
zero

MhSA

Self-supervised

parallel network

pre training

GPT = Generative Pre trained Transformers