

Word2Vec

One-hot vector \Rightarrow orthogonal to each other

\Rightarrow no natural notion of "similarity"

High dimensional vector \Rightarrow semantic spaces

\Rightarrow vectors with similar meaning are close by to each other

- Definition of Similarity: Two vectors are similar \Leftrightarrow high dot product

- Turning Similarity into Probability: softmax function $\mathbb{R}^n \rightarrow [0, 1]^n$

② exponentiation makes everything positive

$$P(w|c) = \frac{\exp(U_w^\top V_c)}{\sum_{w \in V} \exp(U_w^\top V_c)}$$

① dot product \Leftrightarrow similarity
③ normalize over entire vocabulary to give probability distribution

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

↑
amplifies probability of largest x_i

Still assigns some probabilities to smaller x_i

* Interesting semantic patterns emerge in the scaled vectors, and they are linear.

Gradient Descent

Update equation in matrix notation:

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla_{\theta} J(\theta)$$

→ step size or learning rate

Update equation for each single parameter θ_j :

$$\theta_j^{\text{new}} = \theta_j^{\text{old}} - \alpha \frac{\partial}{\partial \theta_j^{\text{old}}} J(\theta)$$

→ loss function

- Problem: very expensive to compute the gradient of the entire corpus

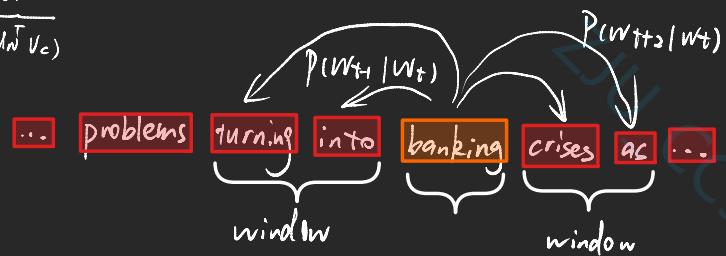
- Solution: Stochastic gradient descent (SGD)

Repeatedly sample windows, only compute the gradient of the window and use it to update the entire corpus.

Main Idea of Word2Vec :

- Start with random word vectors
- Iterate through each word position in the whole corpus.
- Try to predict surrounding words using word vectors:

$$\text{vectors: } P(w_t | c) = \frac{\exp(u_c^T v_t)}{\sum_{w \in V} \exp(u_w^T v_t)}$$



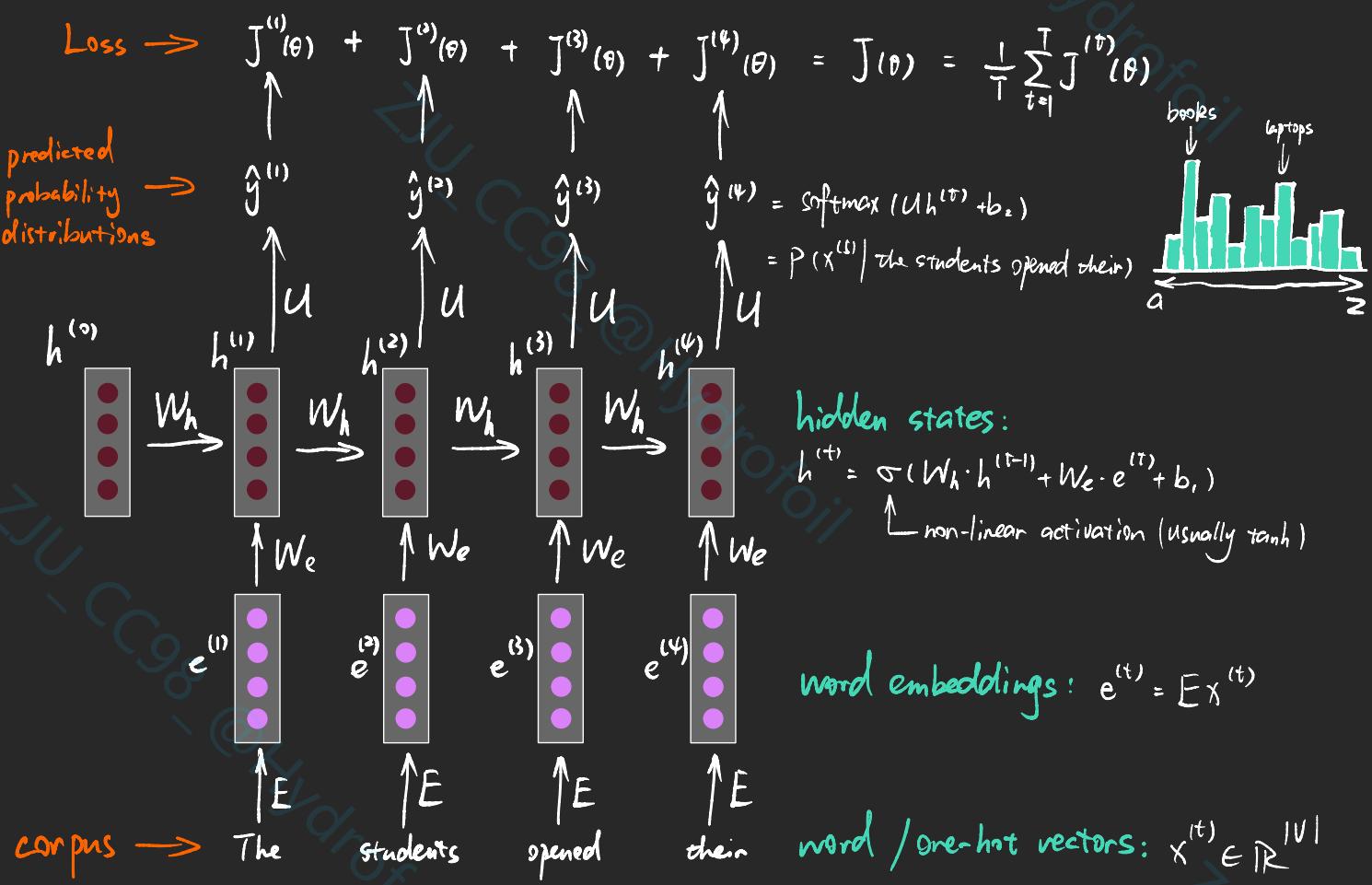
- Learning: Update vectors so each word is at more appropriate position in the semantic space.

\star Polysemy: different senses of a word reside in a linear weighted sum in standard word embeddings like word2vec

$$V_{\text{pike}} = \alpha_1 V_{\text{pike}_1} + \alpha_2 V_{\text{pike}_2} + \alpha_3 V_{\text{pike}_3}$$

where $\alpha_i = \frac{f_i}{f_1 + f_2 + f_3}$, etc., for frequency f .

RNN:



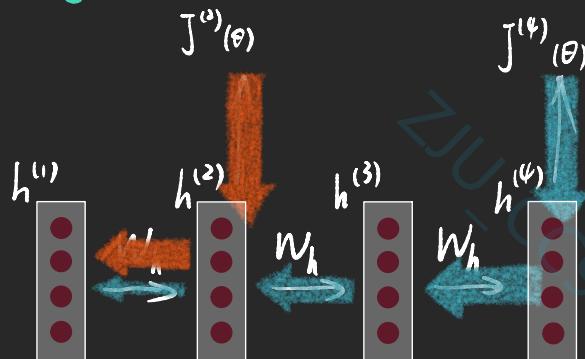
$$J^{(t)}(\theta) = \text{cross_entropy}(y^{(t)}, \hat{y}^{(t)}) = -\sum_{w \in V} y^{(t)}_w \log \hat{y}_w^{(t)} = -\log \hat{y}_{x_{t+1}}^{(t)}$$

RNN Advantages:

- Can process any length input.
- Computation for step t can use information from many steps back.
- Model size doesn't increase for longer input.

RNN Disadvantages:

- Recurrent computation is slow (compute in loop, not parallel)
- In practice, difficult to access information from many steps back due to Vanishing Gradients.



- Exploding Gradients.

LSTM:

Forget gate: controls what's kept vs. forgotten, from previous cell state.

Input gate: controls what parts of the new cell content are written to cell.

Output gate: controls what parts of cell are output to hidden state.

$$\begin{aligned} f^{(t)} &= \sigma(W_f \cdot h^{(t-1)} + U_f \cdot x^{(t)} + b_f) \\ i^{(t)} &= \sigma(W_i \cdot h^{(t-1)} + U_i \cdot x^{(t)} + b_i) \\ o^{(t)} &= \sigma(W_o \cdot h^{(t-1)} + U_o \cdot x^{(t)} + b_o) \end{aligned}$$

sigmoid

New cell content: the new content that is to be written to the cell.

Cell state: erase ("forget") some content from last cell state, and write ("input") some new content.

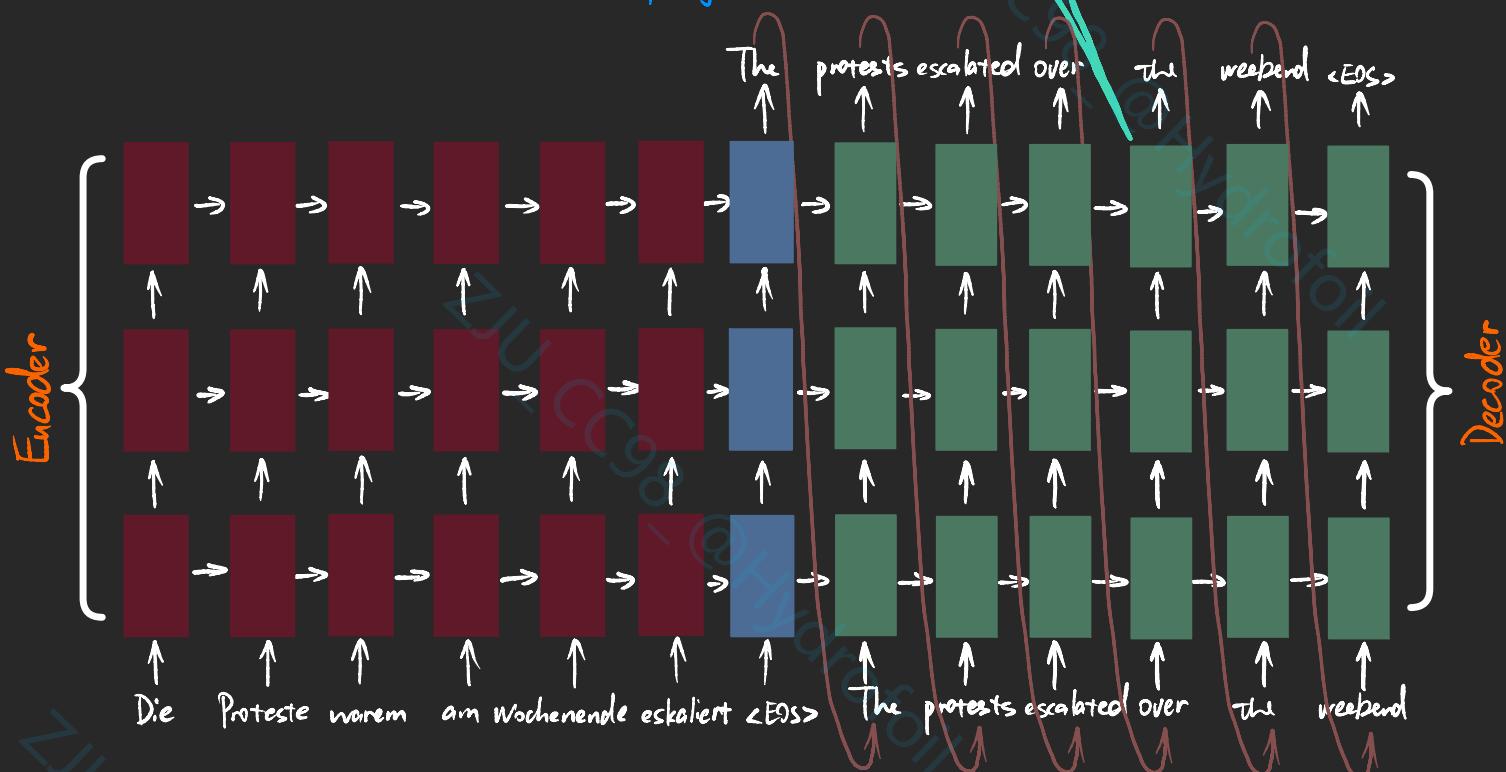
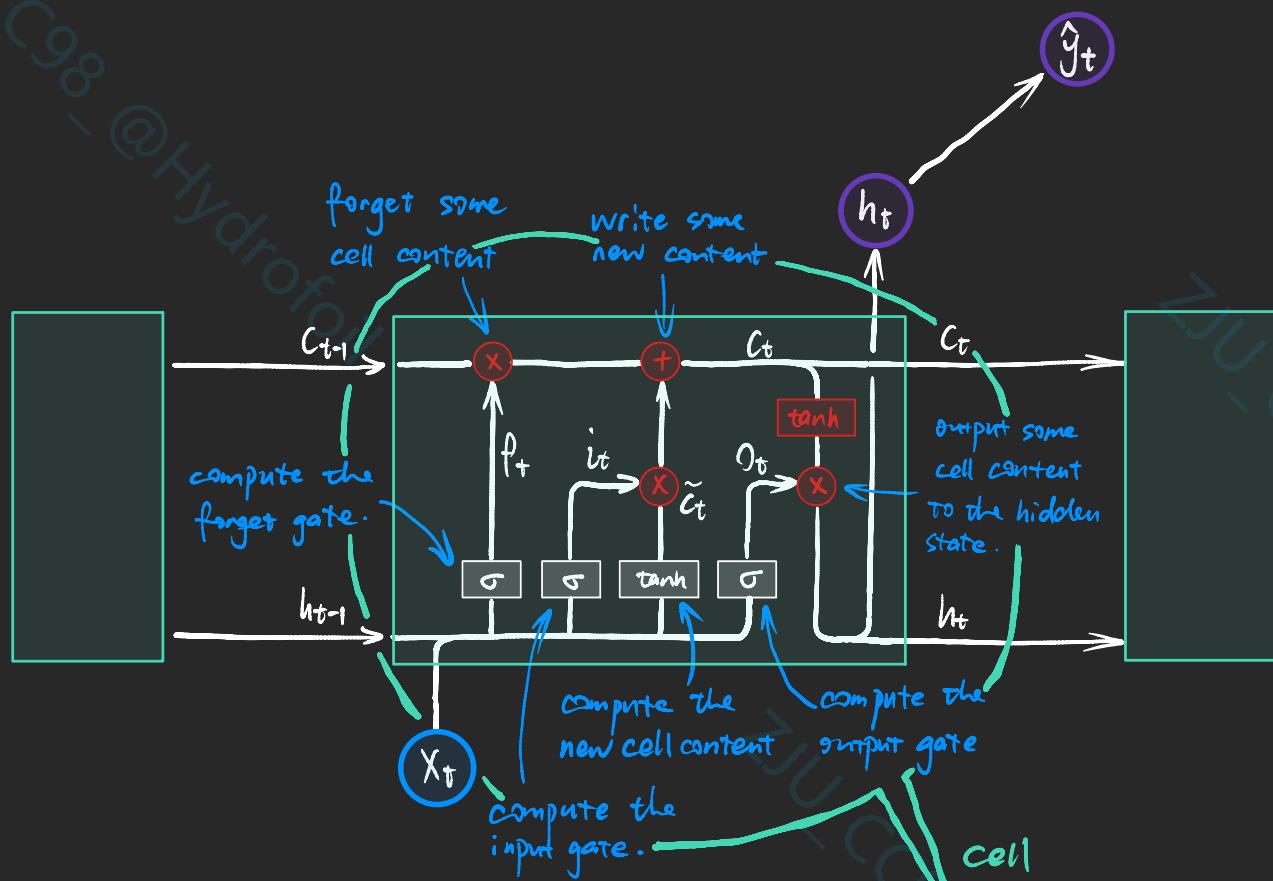
Hidden state: read ("output") some content from the cell.

$$\tilde{c}^{(t)} = \tanh(W_c \cdot h^{(t-1)} + U_c \cdot x^{(t)} + b_c)$$

$$c^{(t)} = f^{(t)} \odot c^{(t-1)} + i^{(t)} \odot \tilde{c}^{(t)}$$

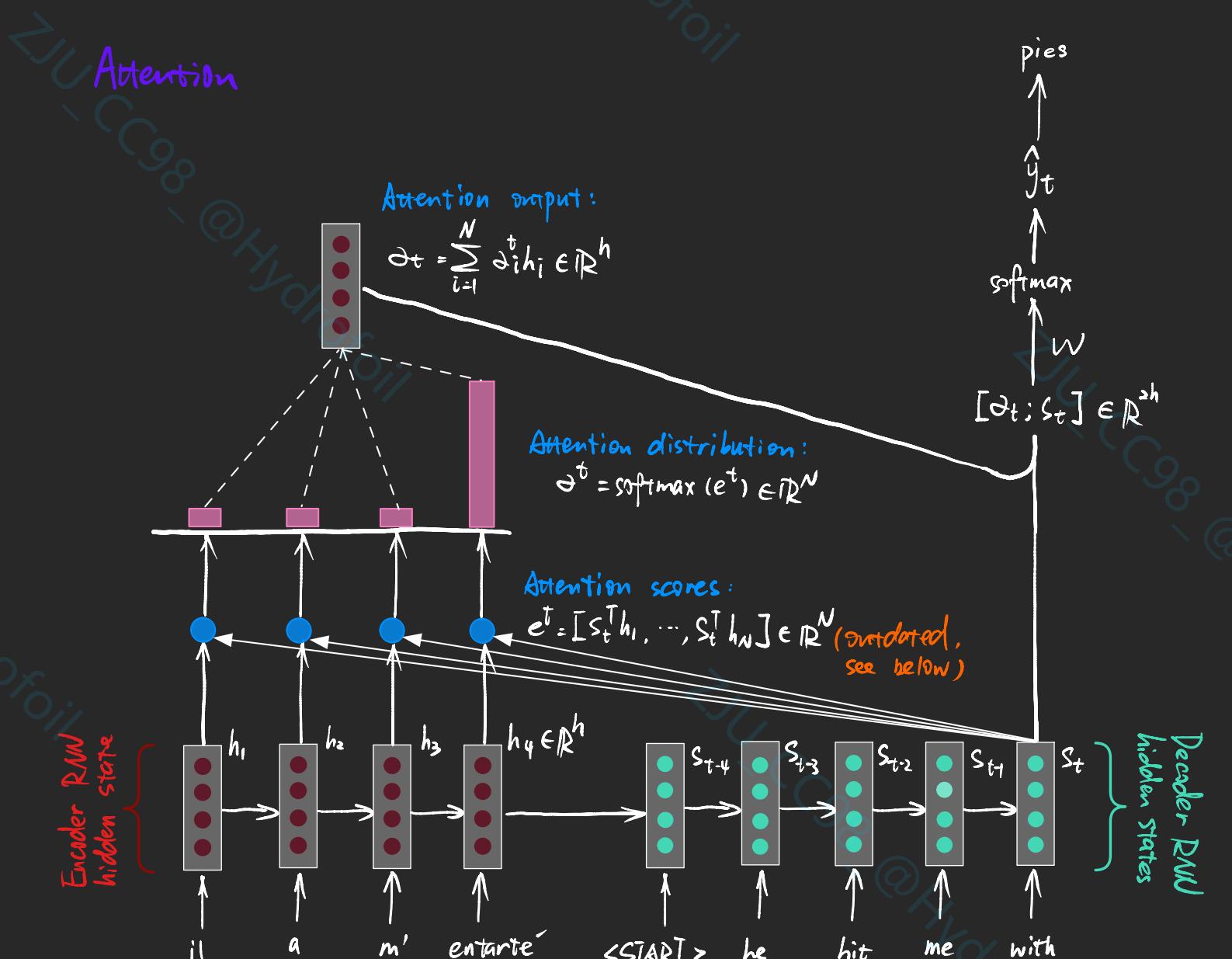
$$h^{(t)} = o^{(t)} \odot \tanh c^{(t)}$$

Multi-layer deep encoder-decoder machine translation net



sequence to sequence model : ① input a sequence \Rightarrow ② output a sequence.

Bottleneck problem of sequence to sequence model : Encoding needs to capture all information about the source sentence, and stuff it into the last vector of the encoder RNN. It's difficult when the source sentence is long.



Computing attention scores: compute $e \in \mathbb{R}^N$ from $h_1, \dots, h_N \in \mathbb{R}^{d_h}$ and $s \in \mathbb{R}^{d_h}$

- Basic dot-product attention: $e_i = s^T h_i \in \mathbb{R}$ (assumes $d_1 = d_2$)
- Multiplicative attention: $e_i = s^T W h_i \in \mathbb{R}$
where $W \in \mathbb{R}^{d_h \times d_h}$ is a weight matrix
problem: too many parameters in W when $d_1 \times d_2$ is big.
- Reduced-rank multiplicative attention: $e_i = s^T (U^T V) h_i = (Us)^T (Vh_i)$
For low rank matrices $U \in \mathbb{R}^{k \times d_h}$, $V \in \mathbb{R}^{k \times d_h}$, $k \ll d_1, d_2$

Attention improves NMT performance

provides a more "human-like" model (look back to the source sentence, rather than remember it all)
solves the bottleneck problem

helps with the vanishing gradient problem

provides some interpretability (we see what the decoder was focusing on)

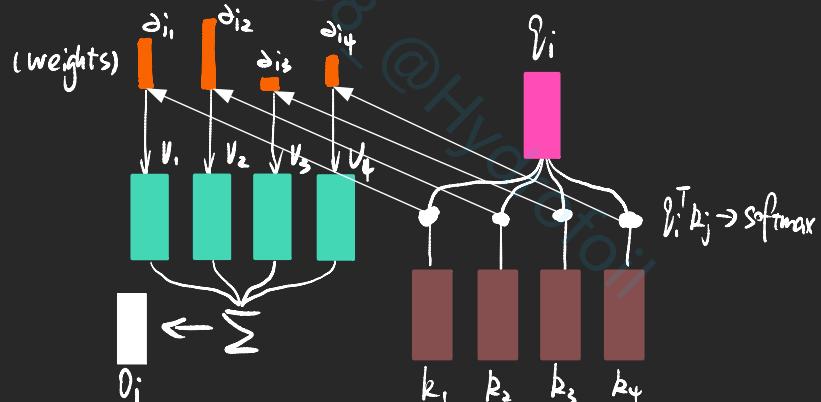
Self - Attention

- word sequence: $W_{1:n} = [w_1, \dots, w_n] \in \mathbb{R}^{|\mathcal{V}| \times n}$
where $|\mathcal{V}|$ is the vocabulary size.
- embedding: for each w_i , let $x_i = E w_i \in \mathbb{R}^{d \times 1}$
where $E \in \mathbb{R}^{d \times |\mathcal{V}|}$ is the embedding matrix.
- position embedding: $\tilde{x}_i = x_i + p_i$ where p_i is the position vector.
- transform each word embedding with weight matrices Q, K, V in $\mathbb{R}^{d \times d}$:
 $q_i = Q x_i$ (queries) $k_i = K x_i$ (keys) $v_i = V x_i$ (values) $q_i, k_i, v_i \in \mathbb{R}^{d \times 1}$
- compute pairwise similarities with keys and queries, normalize with softmax:

$$e_{ij} = q_i^T k_j \in \mathbb{R} \quad \alpha_{ij} = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})} \in \mathbb{R}$$

- compute output for each word as weighted sum of values:

$$o_i = \sum_j \alpha_{ij} v_j \in \mathbb{R}^{d \times 1}$$



- future mask: to use self-attention in the decoders, we need to ensure we can't peek at the future. Otherwise the model can't learn anything during the training.

$$e_{ij} = \begin{cases} q_i^T k_j & j \leq i \\ -\infty & j > i \end{cases}$$

For encoding
these words

We can look at these
(not shadowed out) words

| | | | | |
|---------|--|--|--|--|
| <START> | | | | |
| The | | | | |
| chef | | | | |
| who | | | | |

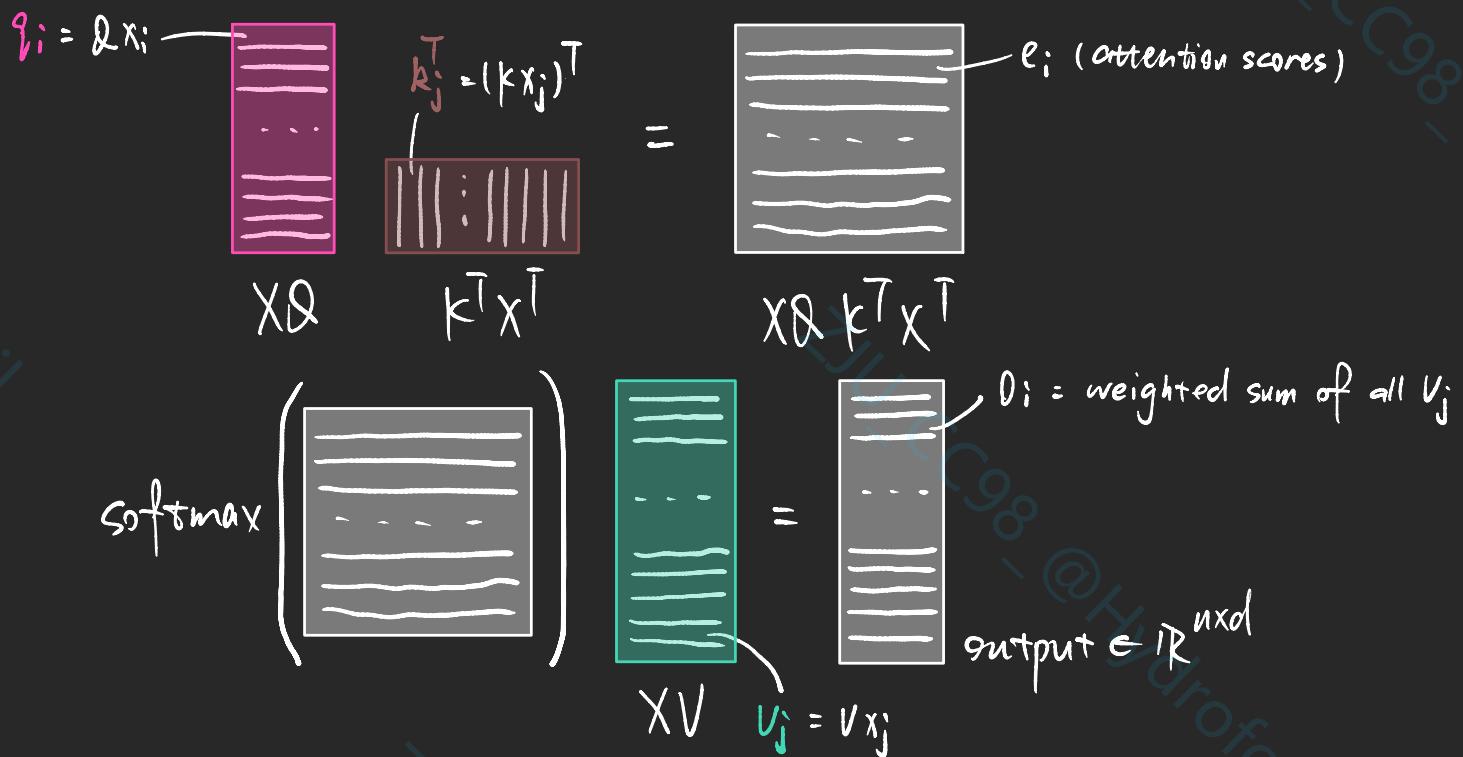
Sequence - Stacked form of Attention:

$$X = [\tilde{x}_1; \tilde{x}_2; \dots; \tilde{x}_n] \in \mathbb{R}^{n \times d}$$

$$XQ \in \mathbb{R}^{n \times d}, XK \in \mathbb{R}^{n \times d}, XU \in \mathbb{R}^{n \times d}$$

$$\text{output} = \text{softmax}(XQ(XK)^T)XU \in \mathbb{R}^{n \times d}$$

$$\begin{array}{c} \tilde{x}_1 \in \mathbb{R}^{1 \times d} \\ \vdots \\ \tilde{x}_n \in \mathbb{R}^{1 \times d} \end{array}$$

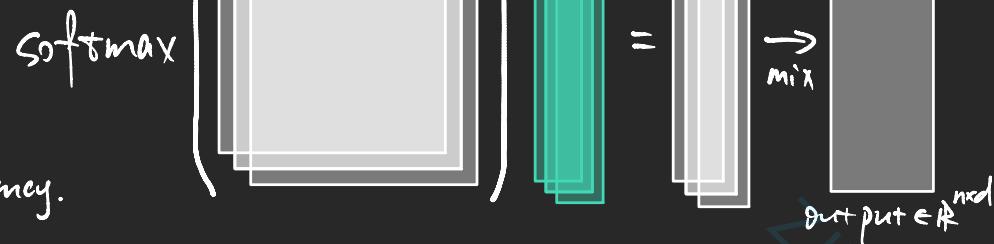


Multi-Head Attention

Let, $Q_l, K_l, V_l \in \mathbb{R}^{d \times h}$,
where h is the number of
attention heads, and l ranges
from 1 to h .

So that $XQ_l \in \mathbb{R}^{n \times \frac{d}{h}}$, likewise
for XK_l and XV_l .

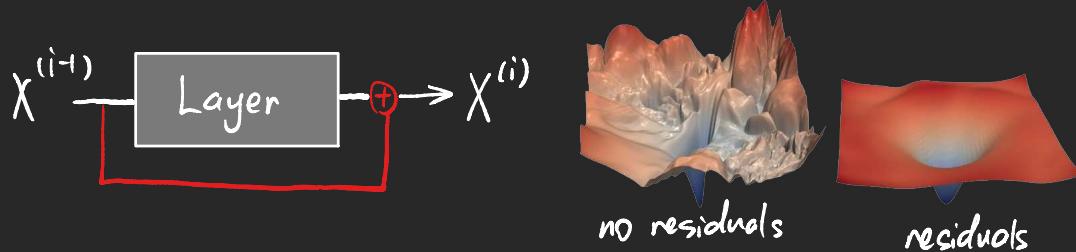
The projection down to $\mathbb{R}^{n \times \frac{d}{h}}$
lower dimensional space $\mathbb{R}^{n \times \frac{d}{h}}$
is for computational efficiency.



Scaled Dot Product: $\text{Output}_l = \text{softmax} \left(\frac{\mathbf{X} \mathbf{Q}_l \mathbf{K}_l^T \mathbf{X}^T}{\sqrt{d/h}} \right) \cdot \mathbf{X} \mathbf{V}_l$

- This is to prevent the case when dimensionality of becomes large
 ⇒ dot products between vectors tend to become large
 ⇒ inputs to the softmax function can be large
 ⇒ the gradient is small.

Residual connections:



- when the gradient of the layer is small, it's at least 1 due to the additive gate.
- make the loss surface smoother, and better to escape local optimum.

Layer normalization:

- problem: Difficult to train the parameters of a given layer because its input from the layer beneath keeps shifting.

- solution: mean: $\mu = \frac{1}{H} \sum_{i=1}^H a_i$

standard deviation: $\sigma = \sqrt{\frac{1}{H} \sum_{i=1}^H (a_i - \mu)^2}$

$$\Rightarrow \text{output} = \frac{x - \mu}{\sigma + \epsilon} \cdot \gamma + \beta$$

Transformer

