

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ

Transformer Language models

Lecture 2 - Decoder-only models

Oct. 1st 2023



Artificial Intelligence Group
Computer Engineering Department, SUT

–Language definition



Language Definition



Chomsky (1959: 137) “A language is a collection of **sentences of finite length** all constructed from a **finite alphabet** (or, where our concern is limited to syntax, a finite vocabulary) of symbols.”



DNA Language

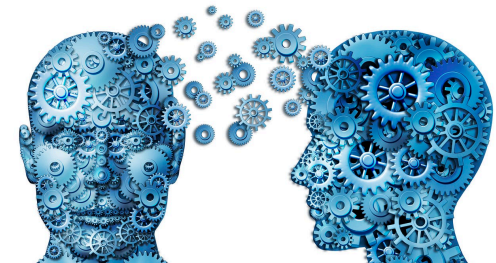
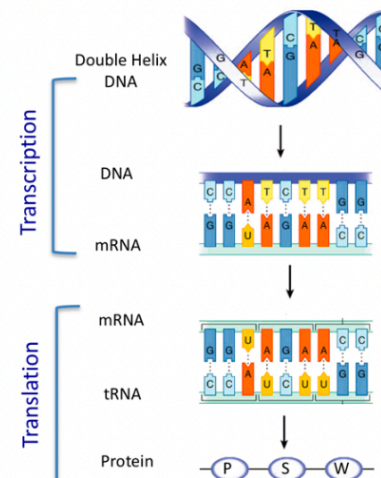
Sentences out of {A,T,C,G}

RNA Language

Sentences out of {A,U,C,G}

Protein Language

Sentences out of
{A,B,C,D,E,F,G,H,I,J,K,L,M,N,O,P,Q,R,S,T,V,W,X,Y}



-Distributional Hypothesis



Distributional hypothesis



Firth (1950) "a word is characterized by the company it keeps"

زیادتی مطلب، کار بر خود آسان کن
صراحی می لعل و بتی چوماهت بس



A picture of a good friendship circle in Persian culture
Made with Bing Image Creator. Powered by DALL-E

Language modeling

$$p(\text{start}, w_1, w_2, \dots, w_n, \text{stop})$$

$$p(\text{start}, w_1, w_2, \dots, w_n, \text{stop}) = \prod_{i=1}^{n+1} \gamma(w_i \mid w_1, w_2, \dots, w_{i-1})$$

N-gram Language modeling

$$p(\text{start}, w_1, w_2, \dots, w_n, \text{stop}) = \prod_{i=1}^{n+1} \gamma(w_i \mid w_1, w_2, \dots, w_{i-1})$$

Bi-gram

$$p(\text{start}, w_1, w_2, \dots, w_n, \text{stop}) = \prod_{i=1}^{n+1} \gamma(w_i \mid w_{i-1})$$

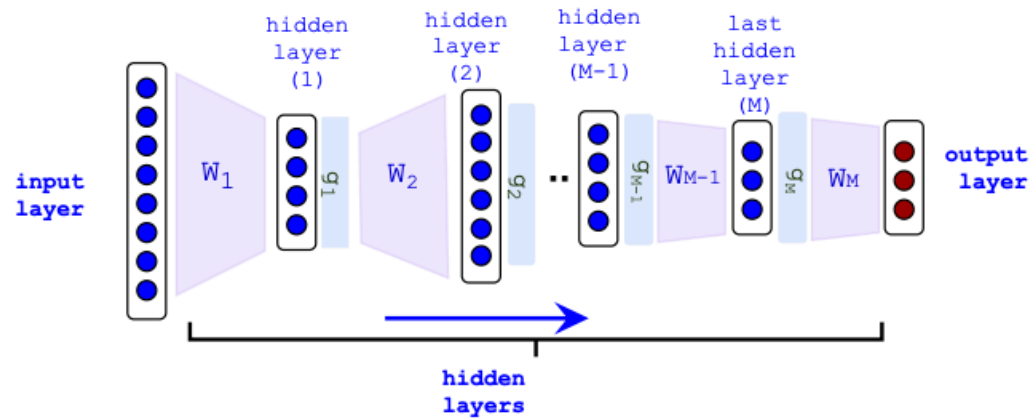
Markov (m^{th} order)

$$p(\text{start}, w_1, w_2, \dots, w_n, \text{stop}) = \prod_{i=1}^{n+1} \gamma(w_i \mid w_{i-m}, \dots, w_{i-1})$$

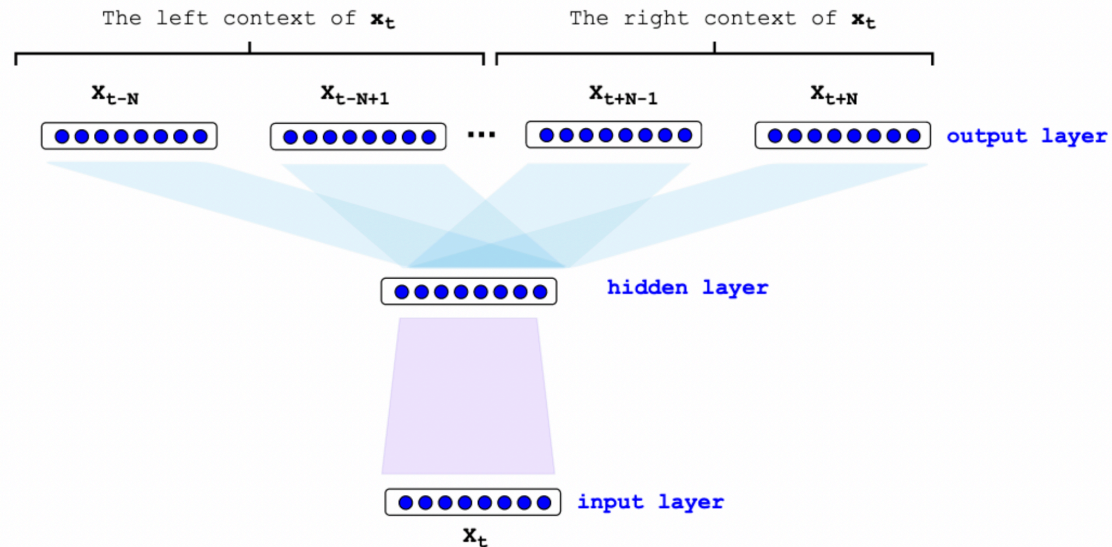
REVIEW

Neural Language modeling

$$p(\text{start}, w_1, w_2, \dots, w_n, \text{stop})$$



Skip-gram – Similar to Language Models



Maximizing the following likelihood:

$$\sum_{t=1}^M \sum_{c \in [t-N, t+N]} \log p(w_c | w_t) \longrightarrow \sum_{t=1}^T \left[\sum_{c \in [t-N, t+N]} \log (1 + e^{-s(w_t, w_c)}) + \sum_{w_r \in \mathcal{N}_{t,c}} \log (1 + e^{s(w_t, w_r)}) \right]$$

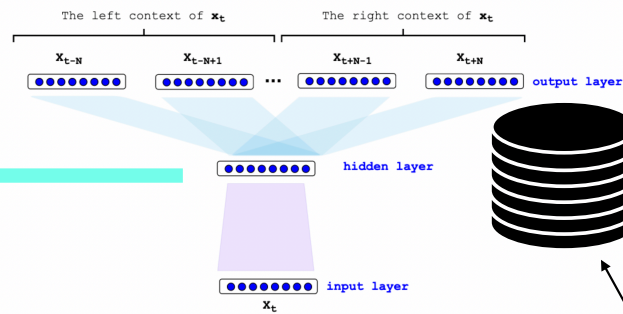
$$p(w_c | w_t; \theta) = \frac{e^{v_c \cdot v_t}}{\sum_{c' \in \mathcal{C}} e^{v_{c'} \cdot v_t}}$$

$$s(w_t, w_c) = v_t^\top \cdot v_c$$

REVIEW



Fixed embeddings – Skip-gram



... که دگر نه عشق خورشیدونه **مهر** ماه دارم

فروزنده ماه و ناپید و **مهر** ...

... فروشت از نگار و نقش ماه **مهر** و آبانش

Fixed embeddings – Skip-gram



... که دگر نه عشق خورشیدونه **مهر** ماه دارم



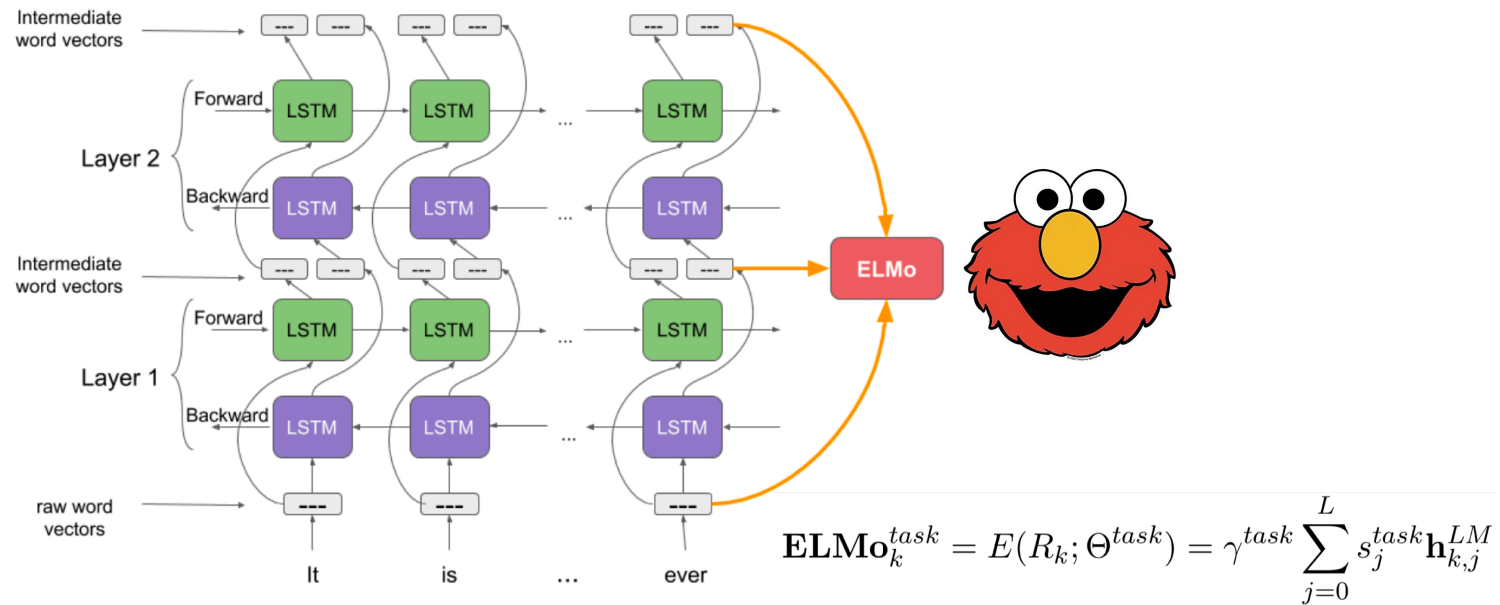
فروزنده ماه و ناپید و **مهر** ...



... فروشت از نگار و نقش ماه **مهر** و آبانش

آبان
ماه
مهر
خورشید
ناپید

ELMO: Deep contextualized word representations



ELMo (Peters et al., 2018; NAACL 2018 best paper)

- Train two separate unidirectional LMs (left-to-right and right-to-left) based on LSTMs
- Feature-based approach: pre-trained representations used as input to task-specific models

– Self-attention



How to contextualize the fixed embeddings?



... که دگر نه عشق خورشیدونه **مهر** ماه دارم



فروزنده ماه و ناپید و **مهر** ...



... فروشت از نگار و نقش ماه **مهر** و آبانش

آبان
ماه
مهر
خورشید
ناپید

Attention

Self-Attention Idea

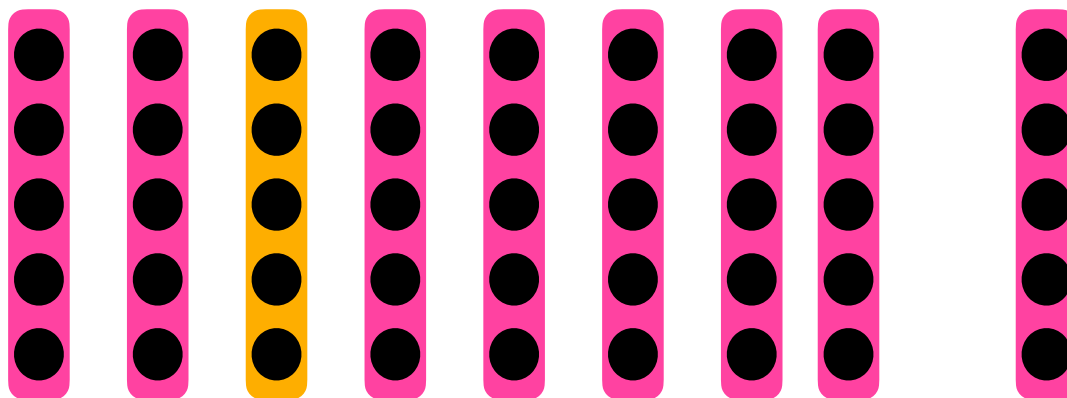
Input embeddings

x_1, x_2, \dots, x_n

Output embeddings

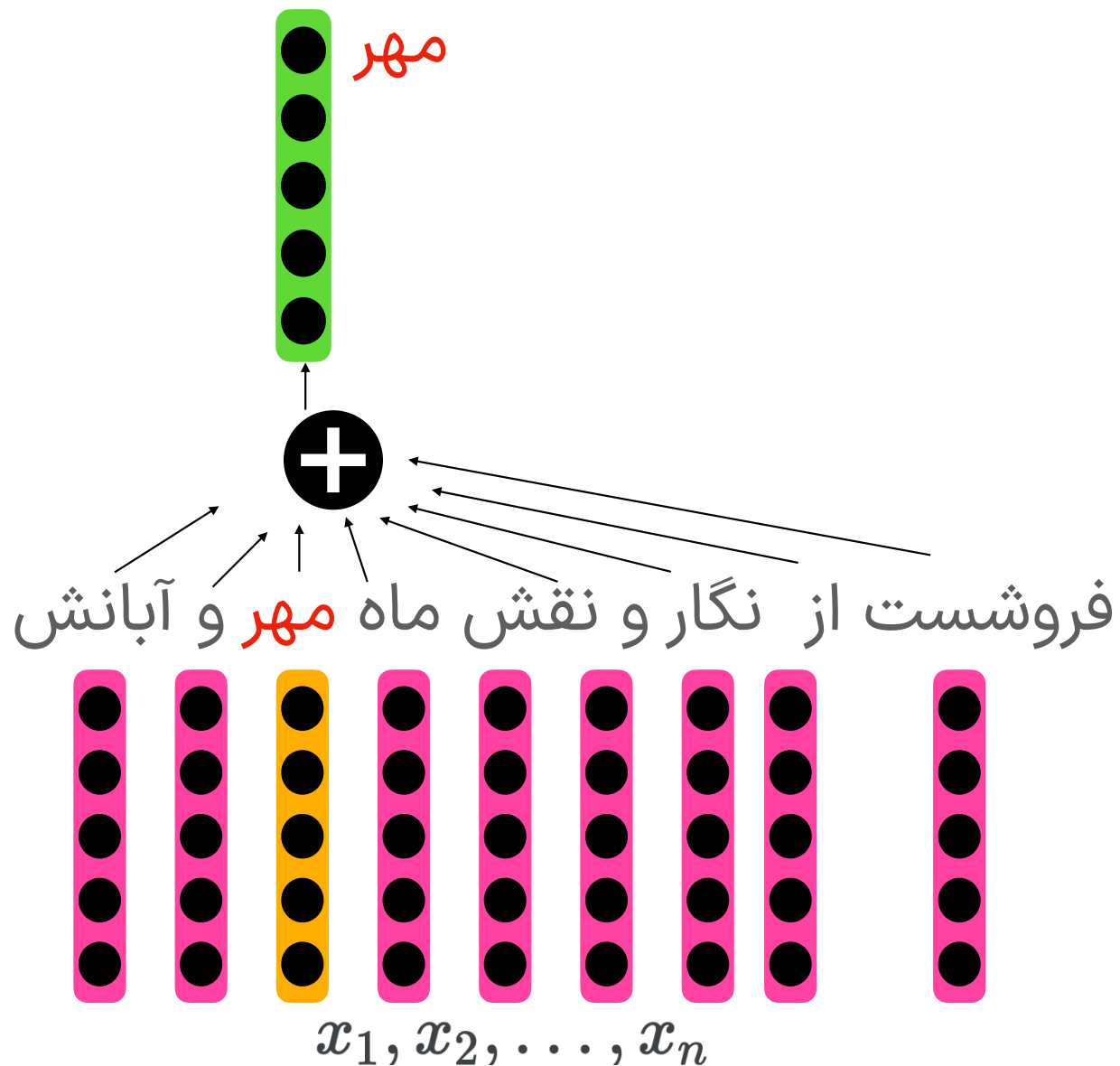
y_1, y_2, \dots, y_n

فرو شست از نگار و نقش ماه مهر و آبانش



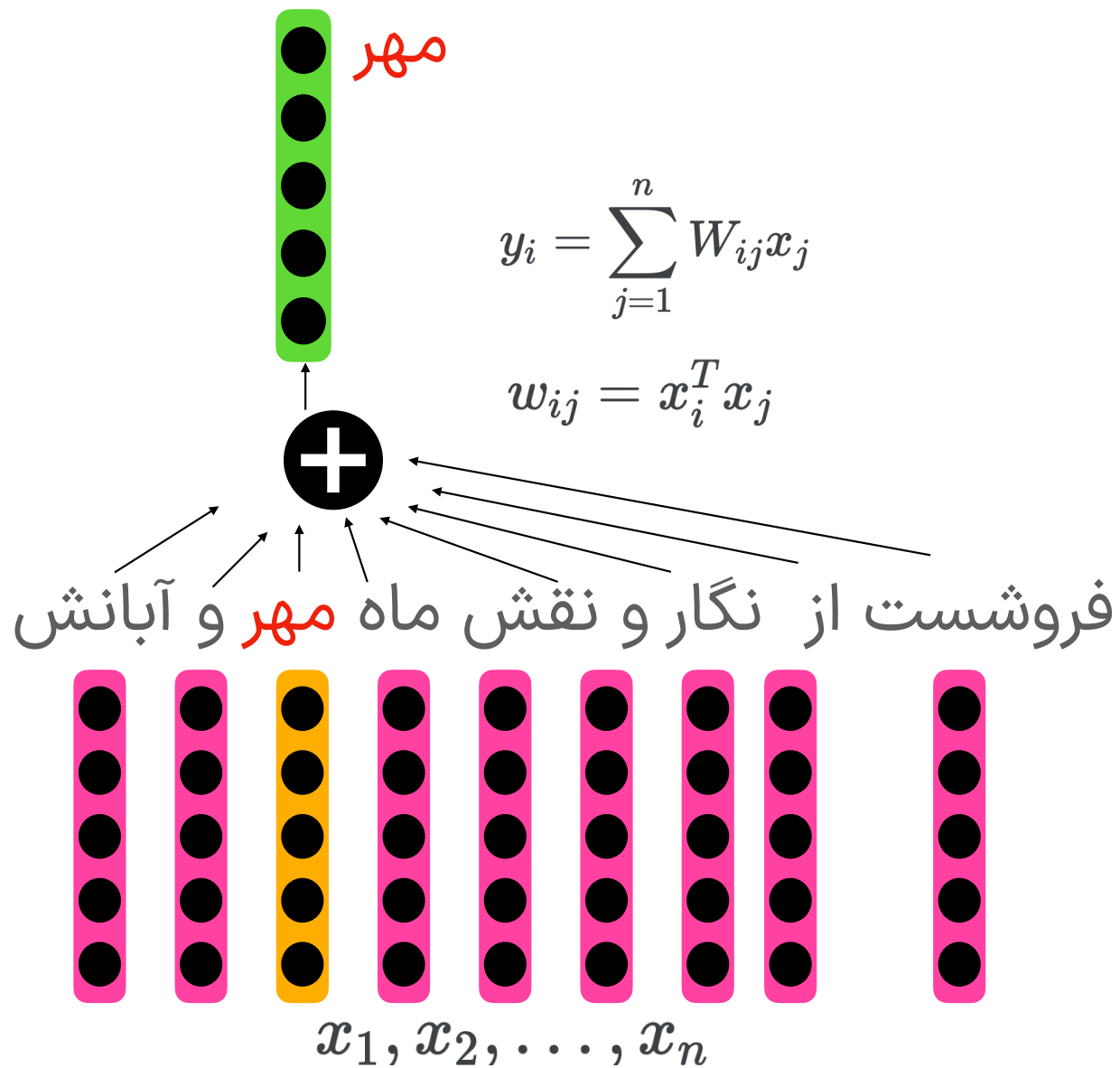
Self-Attention Idea

Attention



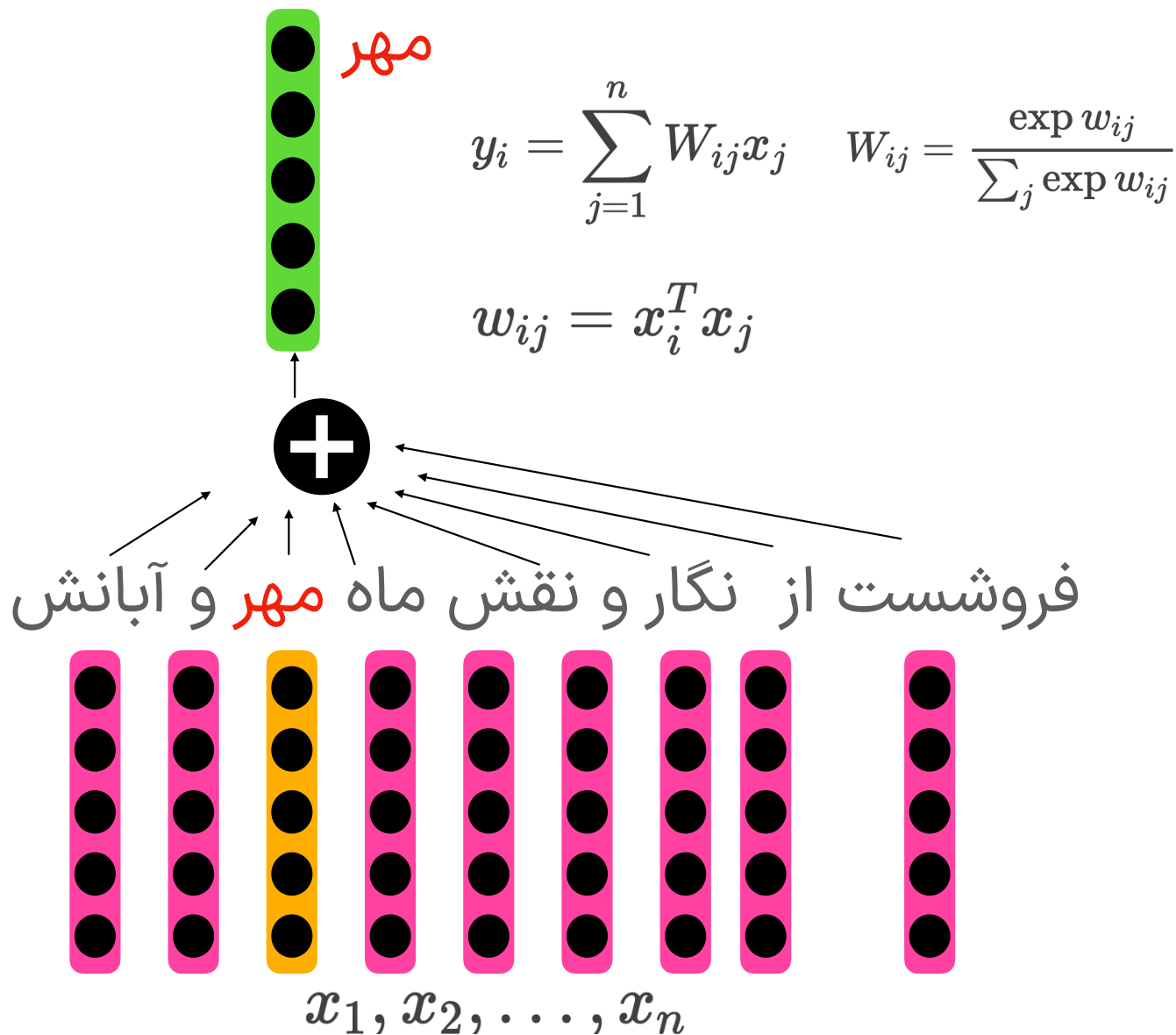
Self-Attention

Attention



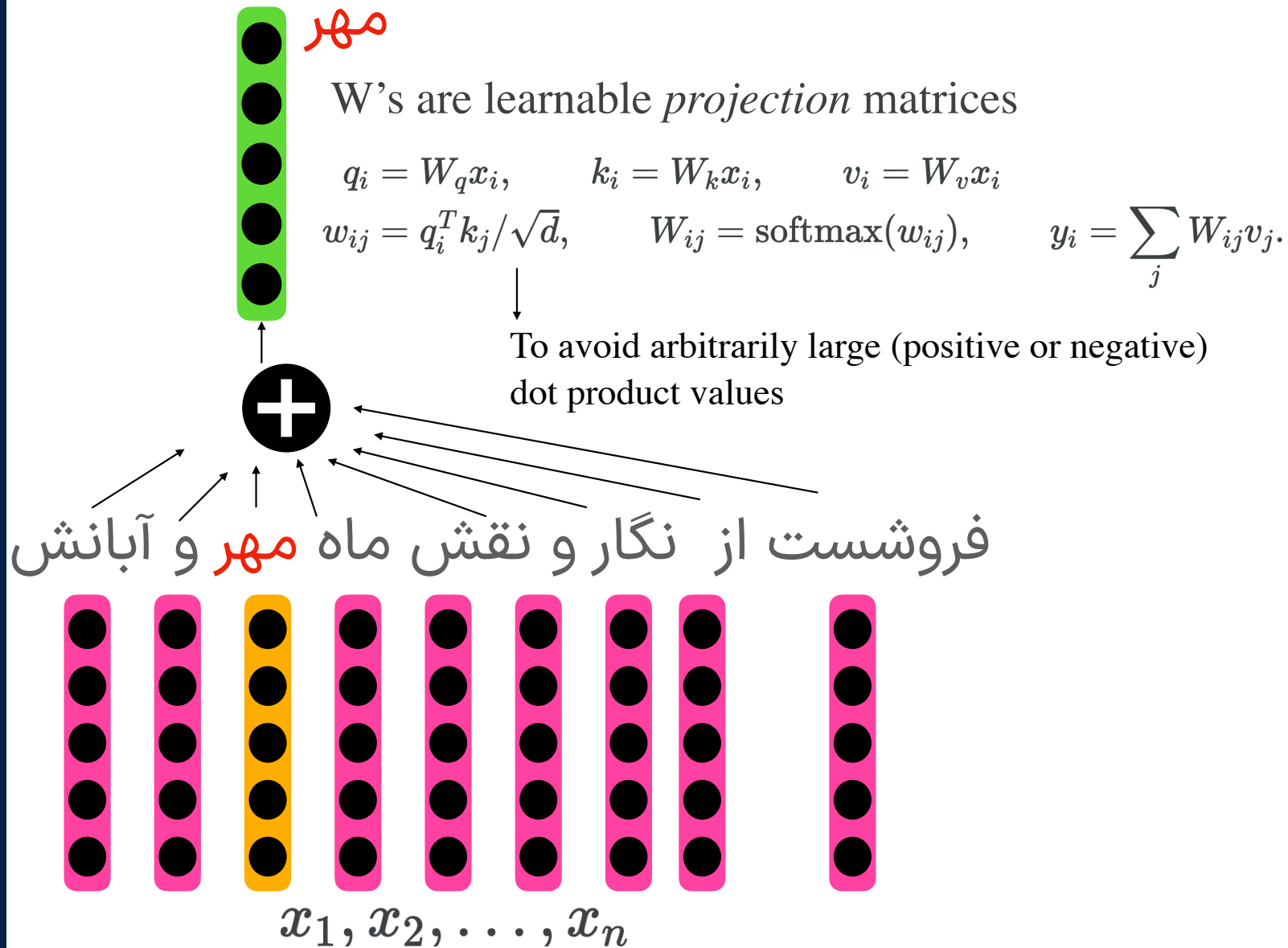
Self-Attention

Attention



Attention

Attention



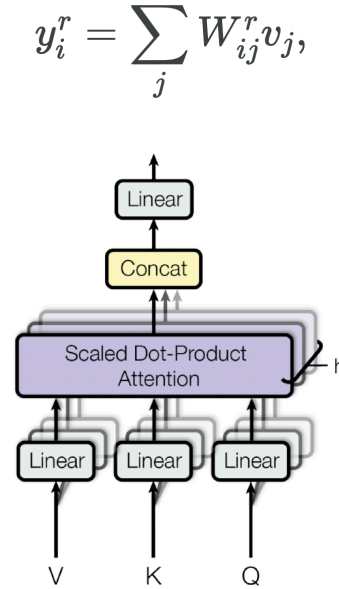
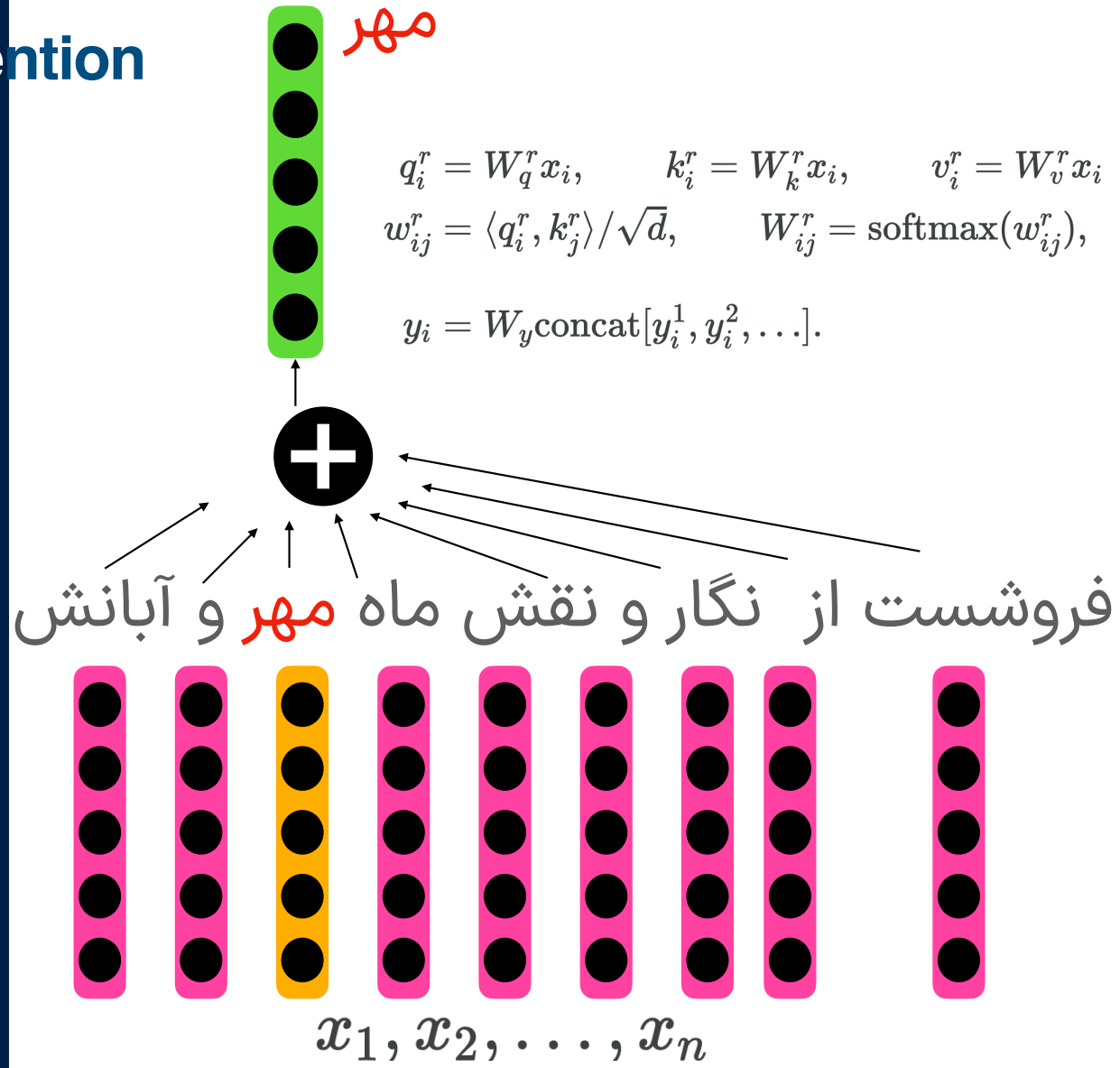
Attention

Attention

- Inputs: a query q and a set of key-value (k-v) pairs to an output
 - All presented as vectors
- Output is weighted sum of values
- Weight of each value: inner product of query and corresponding key

$$A(q, K, V) = \sum_i \frac{e^{q \cdot k_i}}{\sum_j e^{q \cdot k_j}} v_i$$

Multihead Attention



Attention

Matrix Attention

$$X \times W^Q = Q$$

$$X \times W^K = K$$

$$X \times W^V = V$$

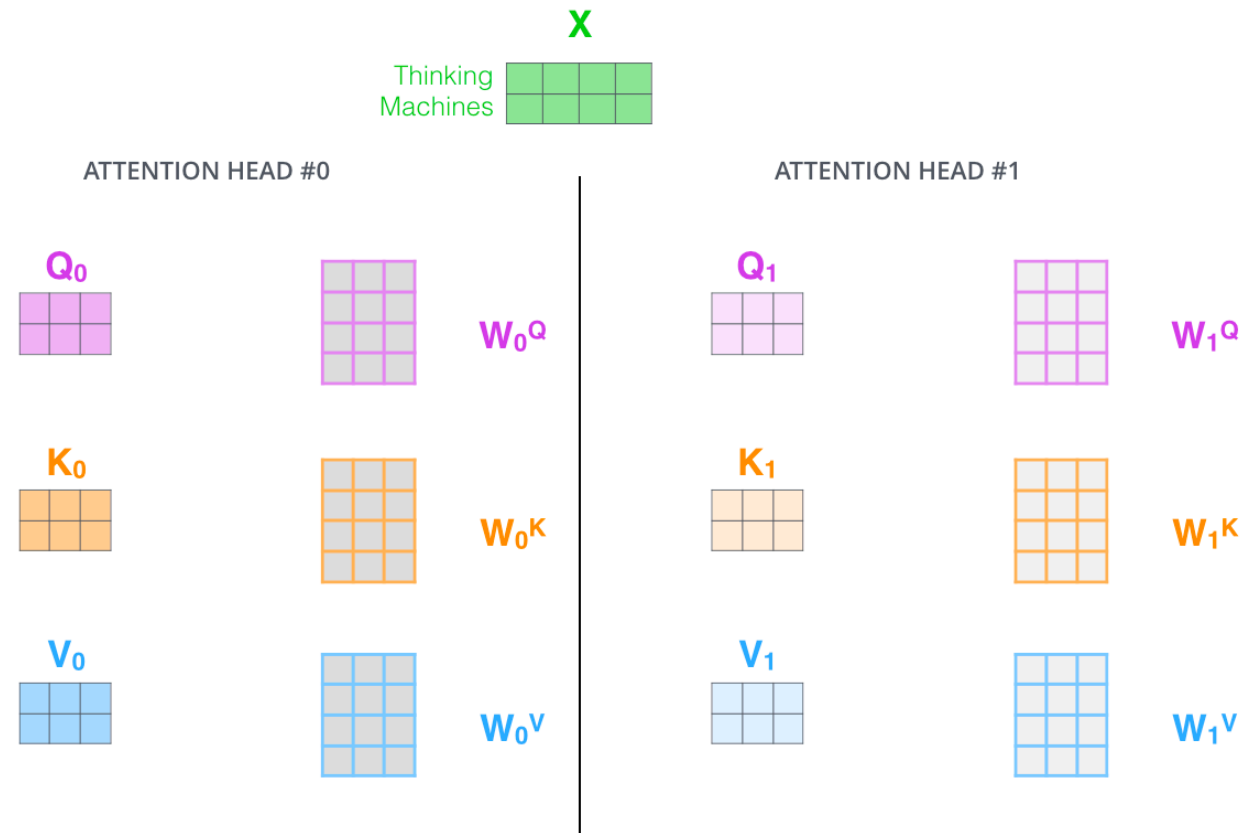
Attention

Matrix Attention

$$\text{softmax}\left(\frac{\begin{matrix} \text{Q} \\ \begin{matrix} \square & \square & \square \end{matrix} \end{matrix} \times \begin{matrix} \text{K}^T \\ \begin{matrix} \square & \square \\ \square & \square \end{matrix} \end{matrix}}{\sqrt{d_k}}\right) \begin{matrix} \text{V} \\ \begin{matrix} \square & \square & \square \end{matrix} \end{matrix}$$
$$= \begin{matrix} \text{Z} \\ \begin{matrix} \square & \square & \square \end{matrix} \end{matrix}$$

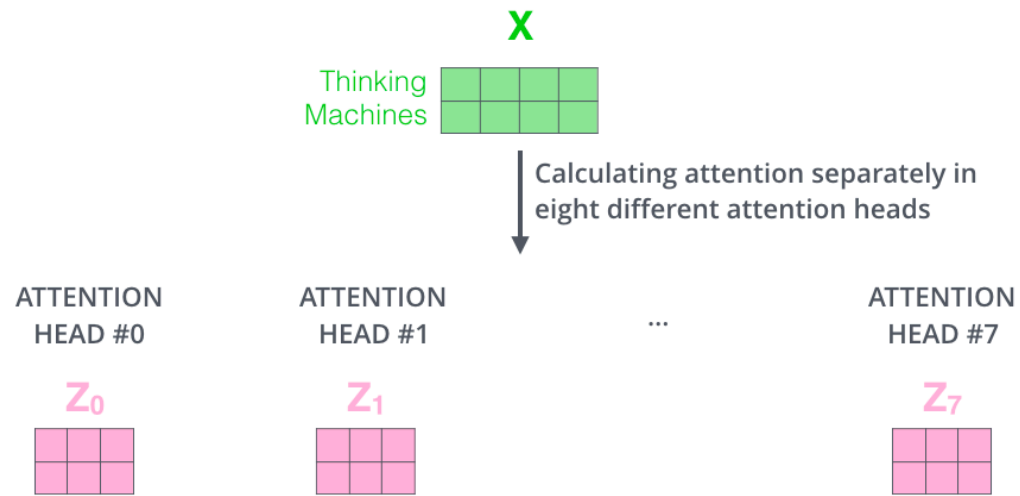
Attention

Matrix Attention



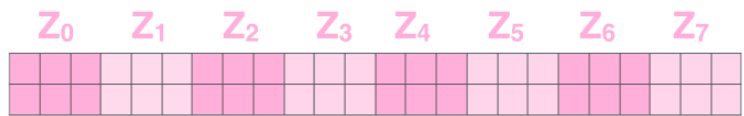
Attention

Matrix Attention



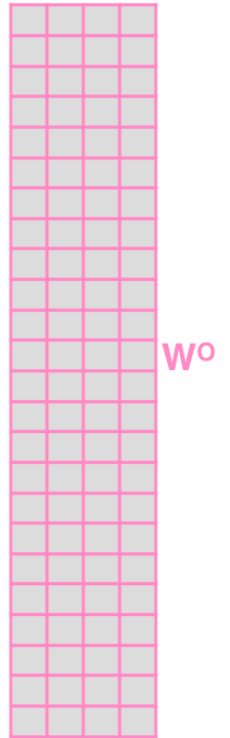
Attention

1) Concatenate all the attention heads

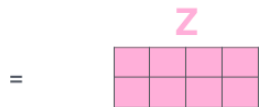


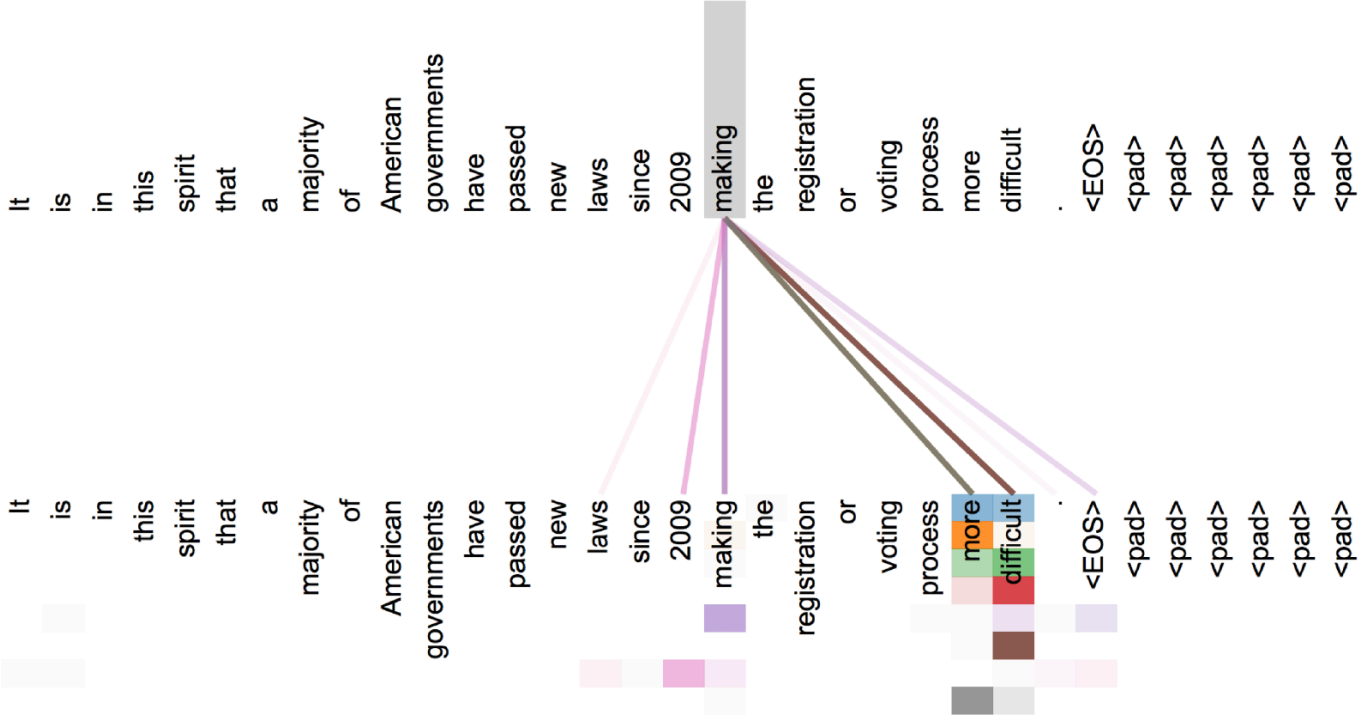
2) Multiply with a weight matrix W^O that was trained jointly with the model

x



3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN





– Transformers



Attention Is All You Need

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Łukasz Kaiser*
Google Brain
lukaszkaizer@google.com

Illia Polosukhin* ‡
illia.polosukhin@gmail.com

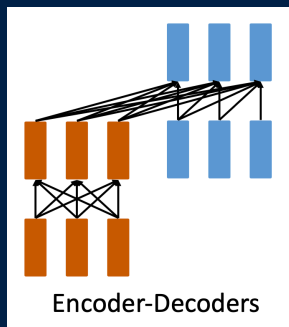
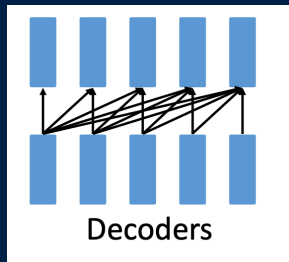
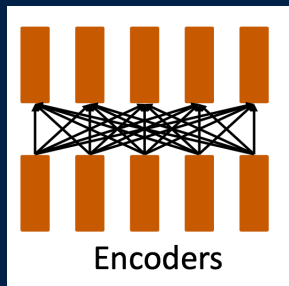
Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... & Polosukhin, I. (2017). Attention is all you need. In *Advances in neural information processing systems* (pp. 5998-6008).

<https://arxiv.org/abs/1706.03762>

–Next lecture



Transformer Architectures



- Encoder-only (e.g., BERT): bidirectional contextual embeddings
- Decoder-only (e.g., GPT-x): unidirectional contextual embeddings, generate one token at a time
- Encoder-decoder (e.g., T5): encode input, decode output