

بِسْمِ اللّٰهِ الرَّحْمٰنِ الرَّحِیْمِ

Transformer Language models

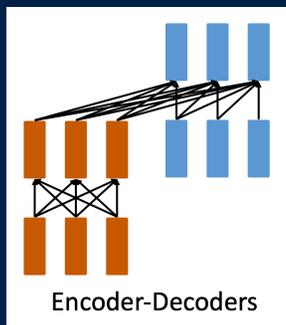
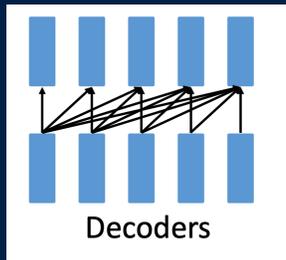
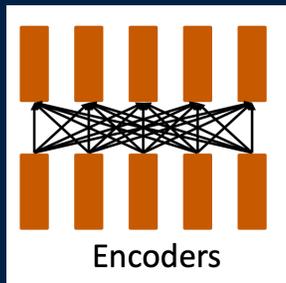
Lecture 3 - Transformers (ii)

Oct. 10th 2023



Artificial Intelligence Group
Computer Engineering Department, SUT

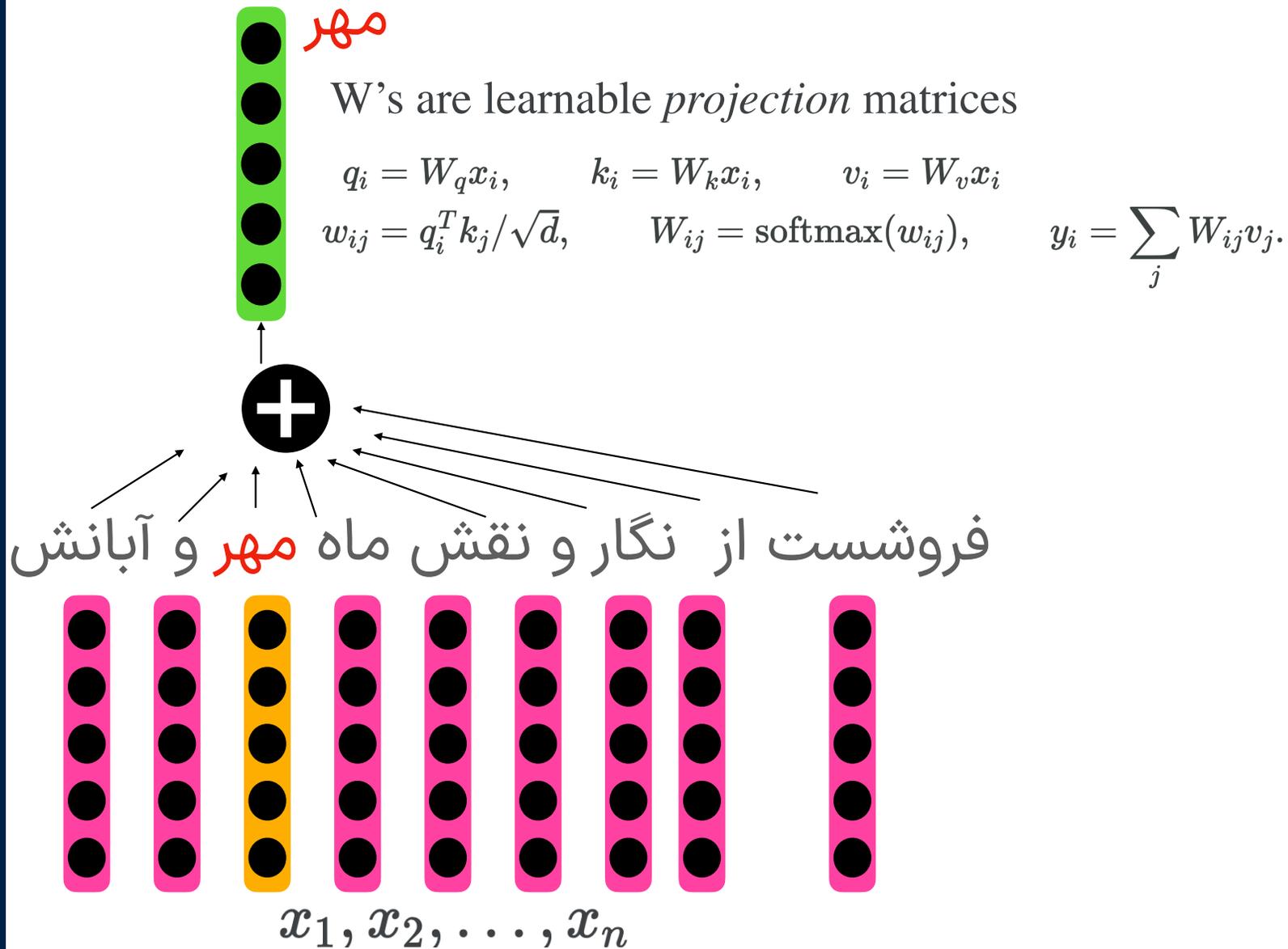
Transformer Architectures



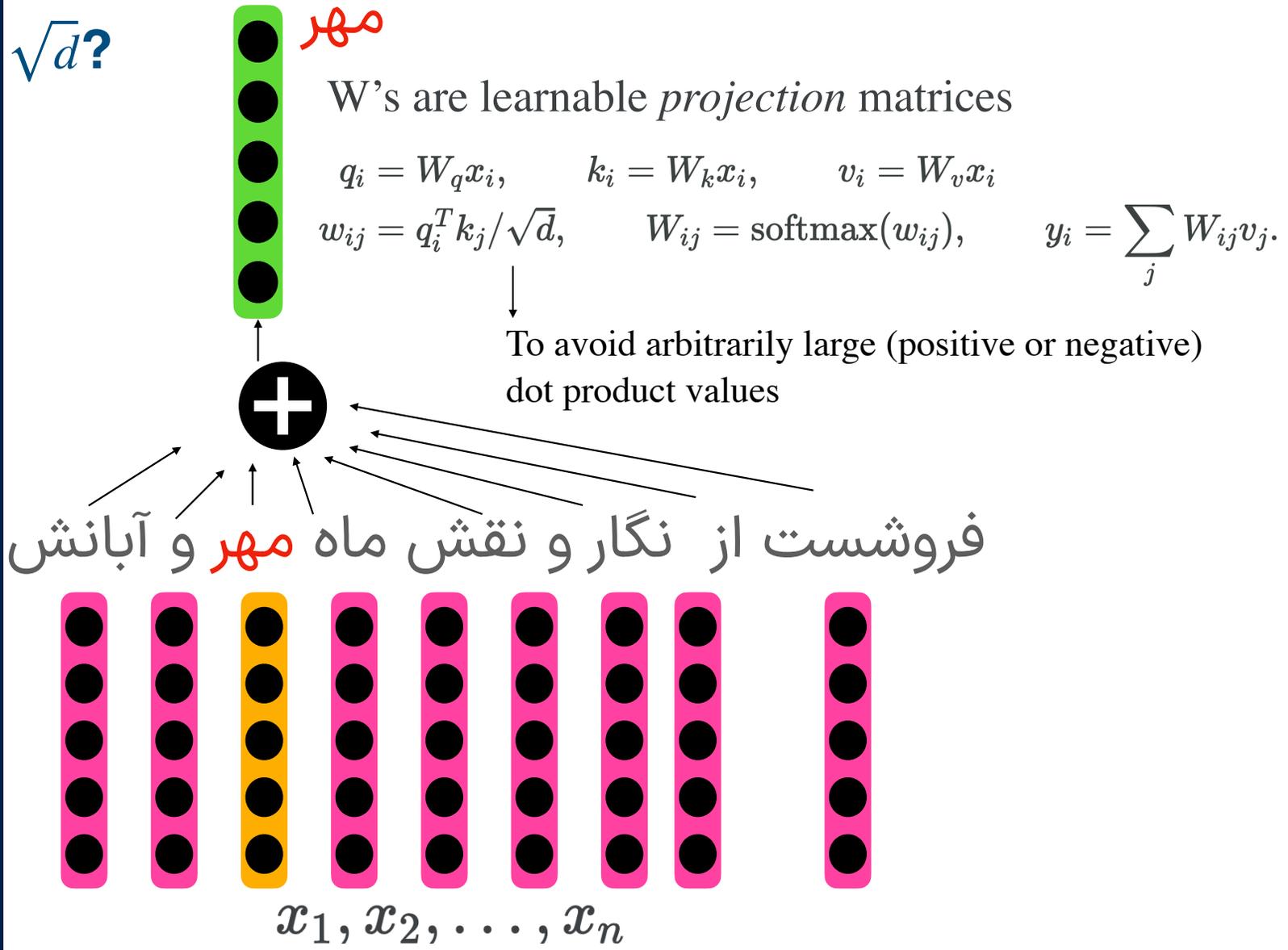
- Encoder-only (e.g., BERT): bidirectional contextual embeddings
- Decoder-only (e.g., GPT-x): unidirectional contextual embeddings, generate one token at a time
- Encoder-decoder (e.g., T5): encode input, decode output

Attention

REVIEW



Why scaling by \sqrt{d} ?



WHY \sqrt{d} ?

$$q_i = W_q x_i, \quad k_i = W_k x_i, \quad v_i = W_v x_i$$
$$w_{ij} = q_i^T k_j / \sqrt{d}, \quad W_{ij} = \text{softmax}(w_{ij}), \quad y_i = \sum_j W_{ij} v_j.$$

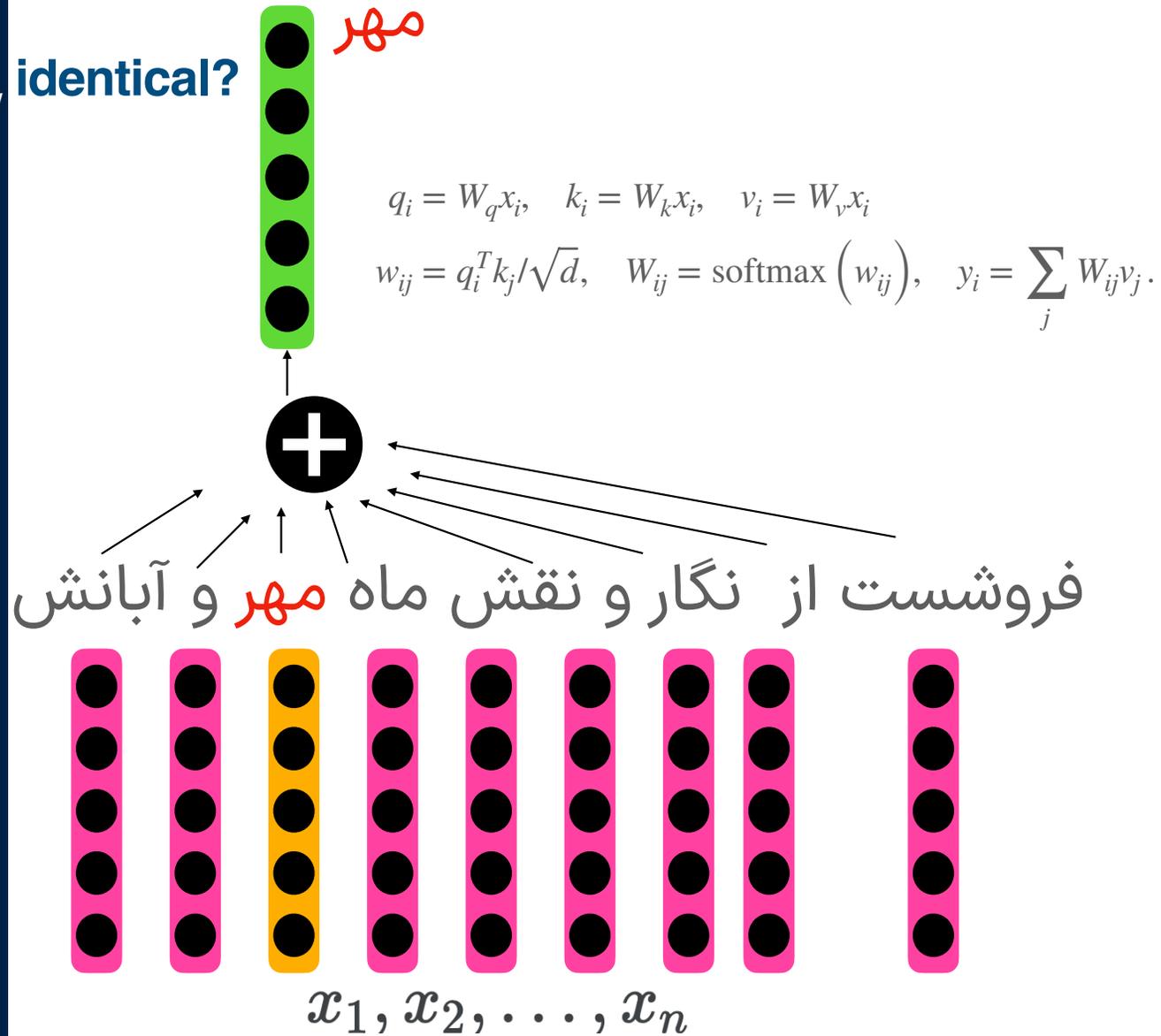
Assume that \mathbf{q} and \mathbf{k} are unit vectors with dimension \mathbf{d} , whose dimensions are independent RV with the following properties:

$$\begin{aligned} \mathbb{E}[q_i] &= \mathbb{E}[k_i] = 0 \\ \text{var}[q_i] &= \text{var}[k_i] = 1 \end{aligned}$$

$$\begin{aligned} \mathbb{E}[q \cdot k] &= \mathbb{E}\left[\sum_{i=1}^d q_i k_i\right] & \text{var}[q \cdot k] &= \text{var}\left[\sum_{i=1}^d q_i k_i\right] \\ &= \sum_{i=1}^d \mathbb{E}[q_i k_i] & &= \sum_{i=1}^d \text{var}[q_i k_i] \\ &= \sum_{i=1}^d \mathbb{E}[q_i] \mathbb{E}[k_i] & &= \sum_{i=1}^d \text{var}[q_i] \text{var}[k_i] \\ &= 0 & &= \sum_{i=1}^d 1 \\ & & &= d \end{aligned} \quad \longrightarrow \quad w_{ij} = \frac{q_i^T k_j - \mu}{\sigma} = \frac{q_i^T k_j}{\sqrt{d}}$$

More detailed proof: <https://github.com/BAI-Yeqi/Statistical-Properties-of-Dot-Product/blob/master/proof.pdf>

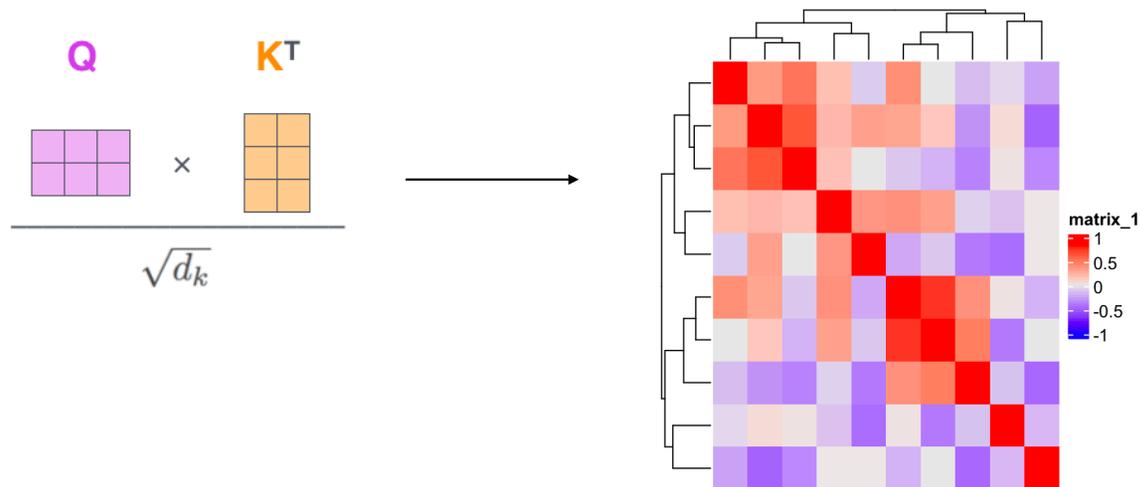
Are W_k and W_q identical?



Are W_k and W_q identical? Better not to be identical!

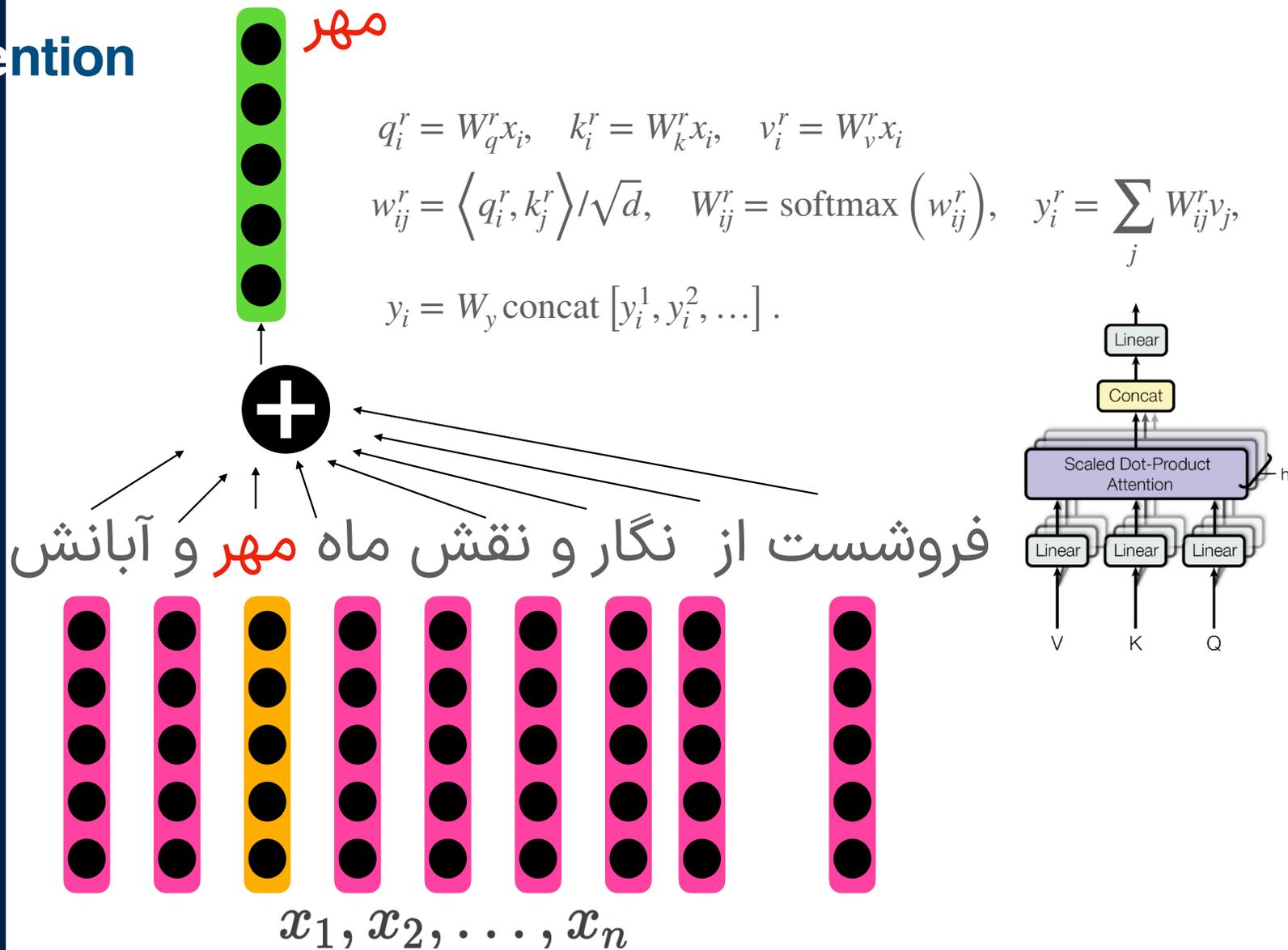
$$q_i = W_q x_i, \quad k_i = W_k x_i, \quad v_i = W_v x_i$$

$$w_{ij} = q_i^T k_j / \sqrt{d}, \quad W_{ij} = \text{softmax}(w_{ij}), \quad y_i = \sum_j W_{ij} v_j.$$

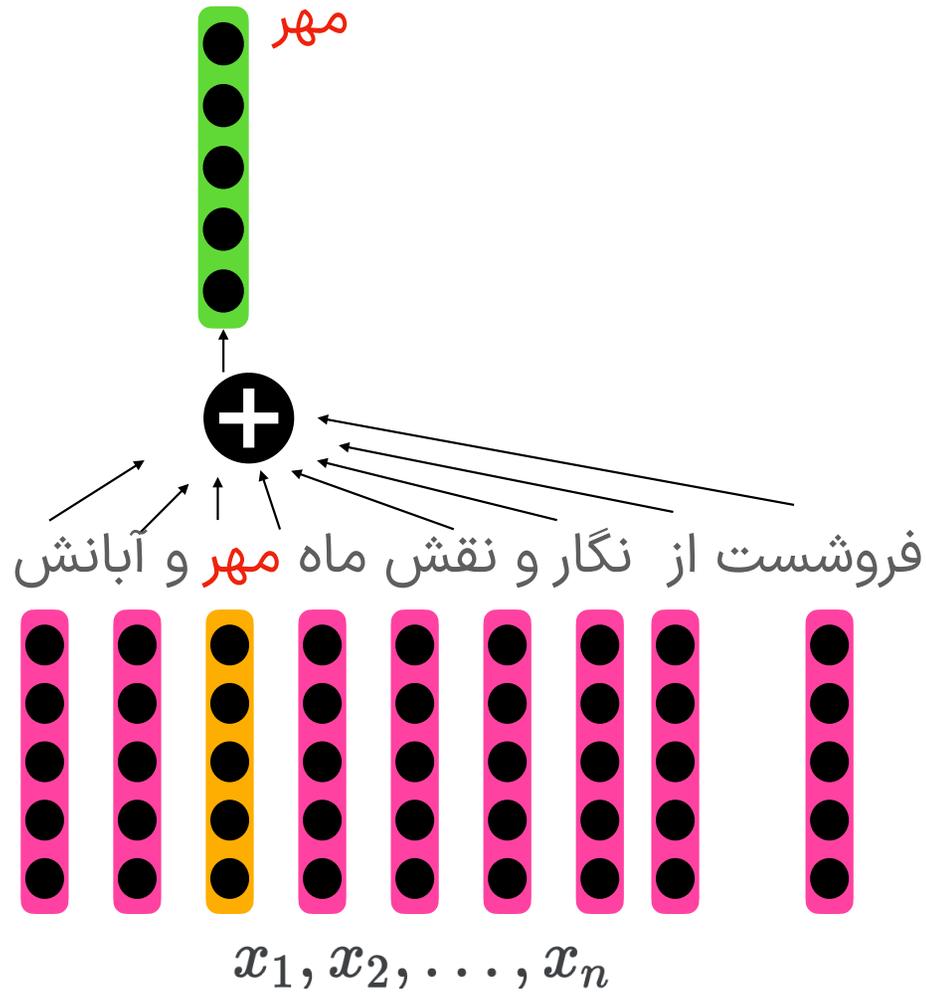


Multihead Attention

REVIEW



Order?



Positional Embedding

- * Assign a number to each time-step within the $[0, 1]$
 - * Time-step differences are not consistent in different sentences.
- * Assign a natural number to each time-step
 - * Long sentences
 - * Differences in the training and the inference

Positional Embedding

- ★ Unique encoding for each time-step.
- ★ Consistent distance between time-steps in varying sentence lengths.
- ★ Easily adapts to longer sentences with bounded values.
- ★ Deterministic output.

Positional Embedding types?

ABSOLUTE VS. RELATIVE POSITION ENCODING

مہر ماہ باران می بارد

معمولاً ہر سال مہر ماہ باران می بارد

Positional Embedding types?

ABSOLUTE VS. RELATIVE POSITION ENCODING

نمبر ماه باران می بارد

معمولاً هر سال **نمبر ماه باران** می بارد

باران ماه مهر

نمبر ماه باران

$$\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \quad \begin{bmatrix} p_{11} & p_{12} & p_{13} \\ p_{21} & p_{22} & p_{23} \\ p_{31} & p_{32} & p_{33} \end{bmatrix} \quad \begin{bmatrix} r_0 & r_1 & r_2 \\ r_{-1} & r_0 & r_1 \\ r_{-2} & r_{-1} & r_0 \end{bmatrix}$$

Attention Matrix *Absolute Position Bias* *Relative Position Bias*

Absolute position embeddings are favorable for classification tasks and relative embeddings perform better for span prediction tasks.

Adding Position Embeddings

Input Embedding $U \in \mathbb{R} \times d$

Position Embedding $P \in \mathbb{R} \times d$

$$\tilde{\mathbf{A}} = \sqrt{\frac{1}{d}}(\mathbf{U} + \mathbf{P})\mathbf{W}^{(q)}\mathbf{W}^{(k)\top}(\mathbf{U} + \mathbf{P})^\top$$

$$\tilde{\mathbf{M}} = \text{SoftMax}(\tilde{\mathbf{A}})(\mathbf{U} + \mathbf{P})\mathbf{W}^{(v)}$$

$$\tilde{\mathbf{O}} = \text{LayerNorm}_2(\tilde{\mathbf{M}} + \mathbf{U} + \mathbf{P})$$

$$\tilde{\mathbf{F}} = \text{ReLU}(\tilde{\mathbf{O}}\mathbf{W}^{(f_1)} + \mathbf{b}^{(f_1)})\mathbf{W}^{(f_2)} + \mathbf{b}^{(f_2)}$$

$$\tilde{\mathbf{Z}} = \text{LayerNorm}_1(\tilde{\mathbf{O}} + \tilde{\mathbf{F}})$$

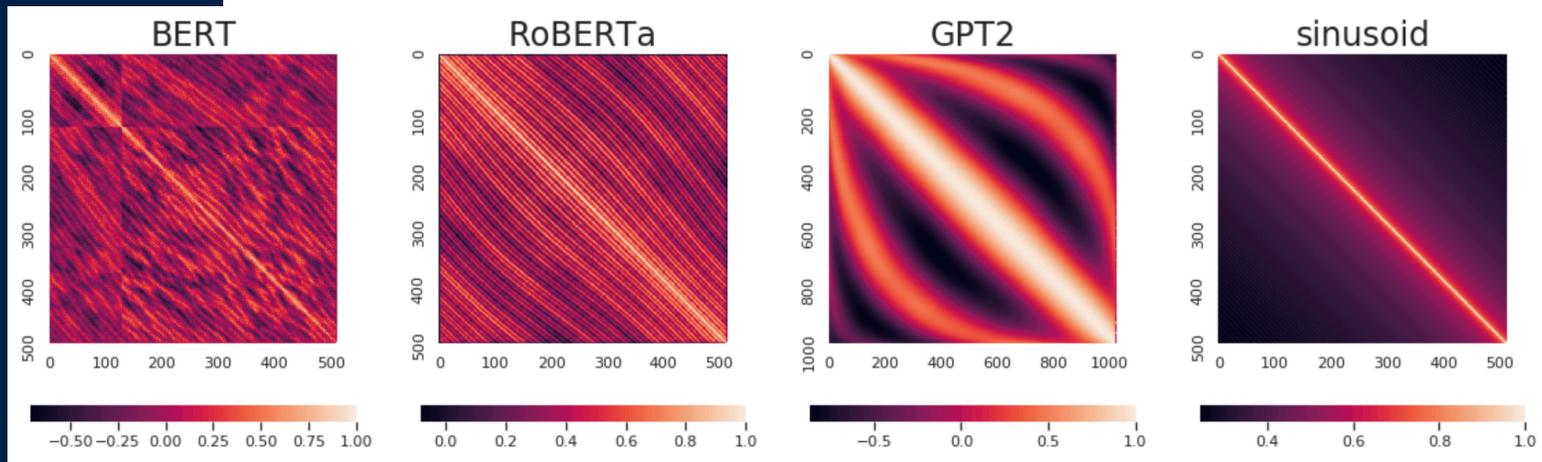
Modifying Attention Matrix

Input Embedding $U \in \mathbb{R} \times d$

Position Embedding $P \in \mathbb{R} \times d$

$$\hat{\mathbf{A}} \sim \underbrace{\mathbf{U}\mathbf{W}^{(q)}\mathbf{W}^{(k)\top}\mathbf{U}^\top}_{\text{unit-unit } \sim \mathbf{A}} + \underbrace{\mathbf{P}\mathbf{W}^{(q)}\mathbf{W}^{(k)\top}\mathbf{U}^\top + \mathbf{U}\mathbf{W}^{(q)}\mathbf{W}^{(k)\top}\mathbf{P}^\top}_{\text{unit-position}} + \underbrace{\mathbf{P}\mathbf{W}^{(q)}\mathbf{W}^{(k)\top}\mathbf{P}^\top}_{\text{position-position}}$$

Positional Embedding types?



Yu-An Wang and Yun-Nung Chen. 2020. [What Do Position Embeddings Learn? An Empirical Study of Pre-Trained Language Model Positional Encoding](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6840–6849, Online. Association for Computational Linguistics.

Sinusoidal Positional Embedding

We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , PE_{pos+k} can be represented as a linear function of PE_{pos} .

$$T^{(k)}E_{t,:} = E_{t+k,:}$$

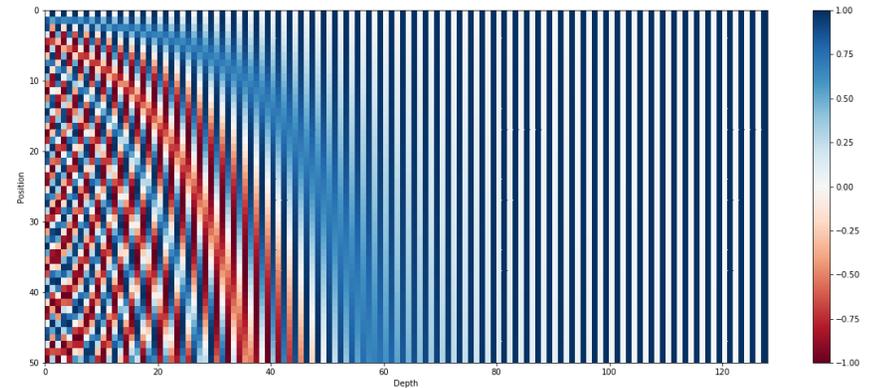
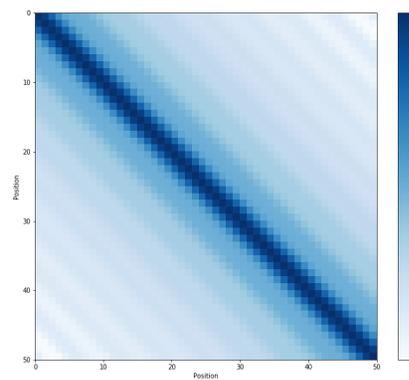
$$T^{(k)} = \begin{bmatrix} \Phi_1^{(k)} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \Phi_2^{(k)} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \dots & \Phi_{\frac{d_{\text{model}}}{2}}^{(k)} \end{bmatrix}$$
$$\Phi_m^{(k)} = \begin{bmatrix} \cos(\lambda_m k) & \sin(\lambda_m k) \\ -\sin(\lambda_m k) & \cos(\lambda_m k) \end{bmatrix}$$
$$\lambda_m = 10000^{\frac{-2m}{d_{\text{model}}}}$$

Sinusoidal Positional Embedding

We chose this function because we hypothesized it would allow the model to easily learn to attend by relative positions, since for any fixed offset k , PE_{pos+k} can be represented as a linear function of PE_{pos} .

$$\vec{p}_t^{(i)} = f(t)^{(i)} := \begin{cases} \sin(\omega_k \cdot t), & \text{if } i = 2k \\ \cos(\omega_k \cdot t), & \text{if } i = 2k + 1 \end{cases} \quad \omega_k = \frac{1}{10000^{2k/d}}$$

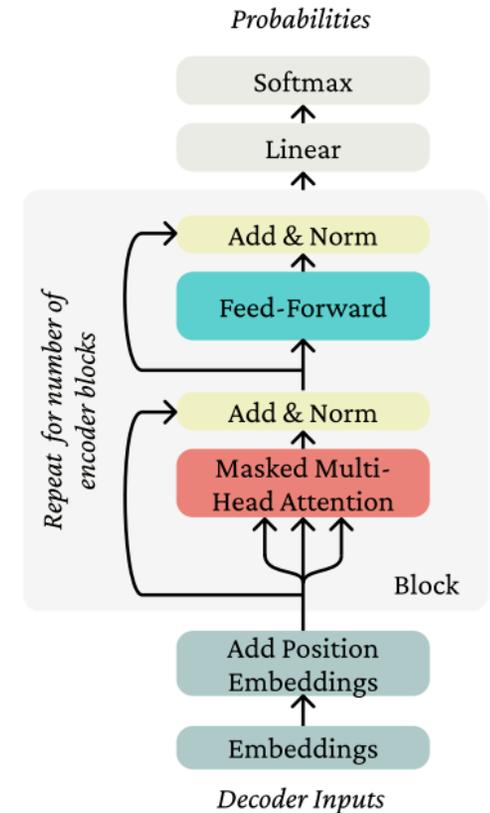
$$\vec{p}_t = \begin{bmatrix} \sin(\omega_1 \cdot t) \\ \cos(\omega_1 \cdot t) \\ \sin(\omega_2 \cdot t) \\ \cos(\omega_2 \cdot t) \\ \vdots \\ \sin(\omega_{d/2} \cdot t) \\ \cos(\omega_{d/2} \cdot t) \end{bmatrix}_{d \times 1}$$



Vaswani, Ashish, et al. "Attention is all you need." Advances in neural information processing systems 30 (2017).

Transformer block

- Each block has two “sublayers”
 1. Multihead attention
 2. Feed-forward NNet (with ReLU)
- Residual: $x + \text{Sublayer}(x)$
- Layernorm changes input to have mean 0 and variance 1



Layer normalization

Main idea: batch normalization is very helpful, but hard with sequences of different lengths

- Resulting more stable input to the next layer
 - Simple solution: “layer normalization” – like batch norm, but not across the batch
- Batch norm Layer norm

a_1, a_2, \dots, a_B ← d -dimensional vectors for each sample in batch

d -dim → $\mu = \frac{1}{B} \sum_{i=1}^B a_i$ $\sigma = \sqrt{\frac{1}{B} \sum_{i=1}^B (a_i - \mu)^2}$

← 1 -dim $\mu = \frac{1}{d} \sum_{j=1}^d a_j$ $\sigma = \sqrt{\frac{1}{d} \sum_{j=1}^d (a_j - \mu)^2}$ ← different *dimensions* of a

$\bar{a}_i = \frac{a_i - \mu}{\sigma}$ $\bar{a} = \frac{a - \mu}{\sigma}$

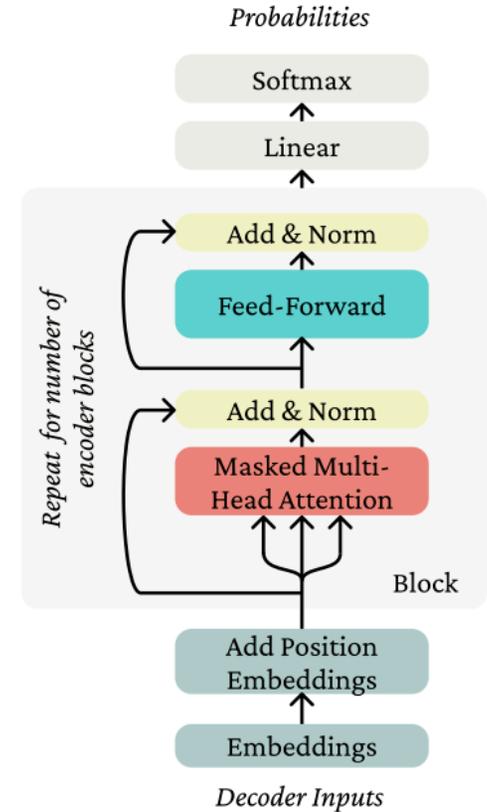
Why transformers?

Pros:

- + Much easier to parallelize
- + Much better long-range connections
- + In practice, can make it much deeper (more layers) than RNN

Cons:

- Attention computations are technically $O(n^2)$
- Somewhat more complex to implement (positional encodings, etc.)



- Encoder Language Model
- BERT LM Architecture





BERT: key contributions

- It is a **fine-tuning approach** based on a deep **Transformer encoder**
- The key: learn representations based on **bidirectional context**

Why? Because both left and right contexts are important to understand the meaning of words.

Example #1: we went to the river bank.

Example #2: I need to go to bank to make a deposit.

- **Pre-training objectives:** masked language modeling + next sentence prediction
- State-of-the-art performance on a large set of **sentence-level** and **token-level** tasks

Masked Language Modeling (MLM)

- Q: Why we can't do language modeling with bidirectional models?



- Solution: Mask out $k\%$ of the input words, and then predict the masked words

store gallon
 ↑ ↑
the man went to [MASK] to buy a [MASK] of milk

MLM:masking rate and strategy

- **Q: What is the value of k ?**
 - They always use $k = 15\%$.
 - Too little masking: computationally expensive
 - Too much masking: not enough context
 - See (Wettig et al., 2022) for more discussion of masking rates
- **Q: How are masked tokens selected?**
 - 15% tokens are uniformly sampled
 - Is it optimal? See span masking (Joshi et al., 2020) and PMI masking (Levine et al., 2021)

Example: He [MASK] from Kuala [MASK] , Malaysia.