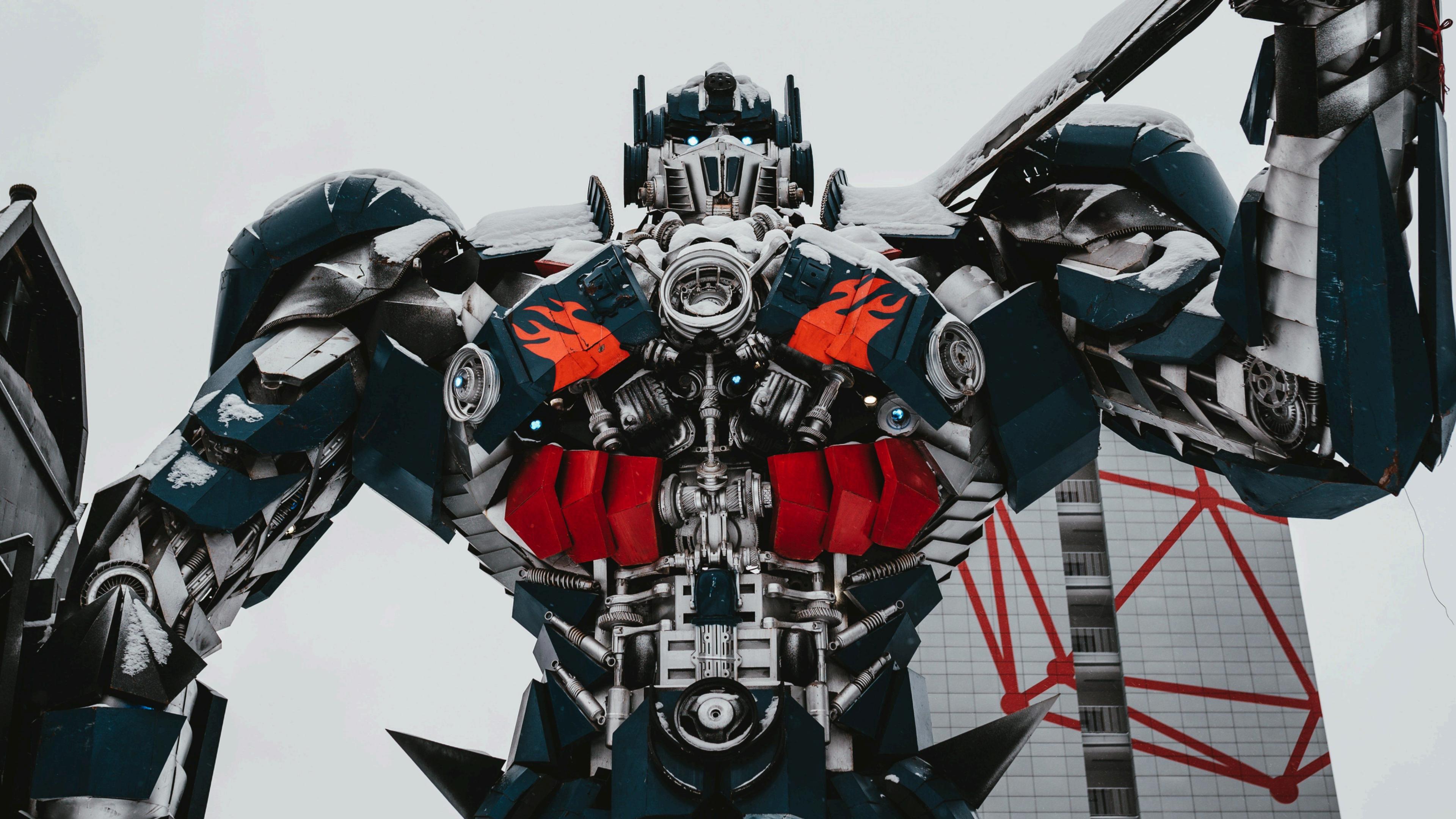
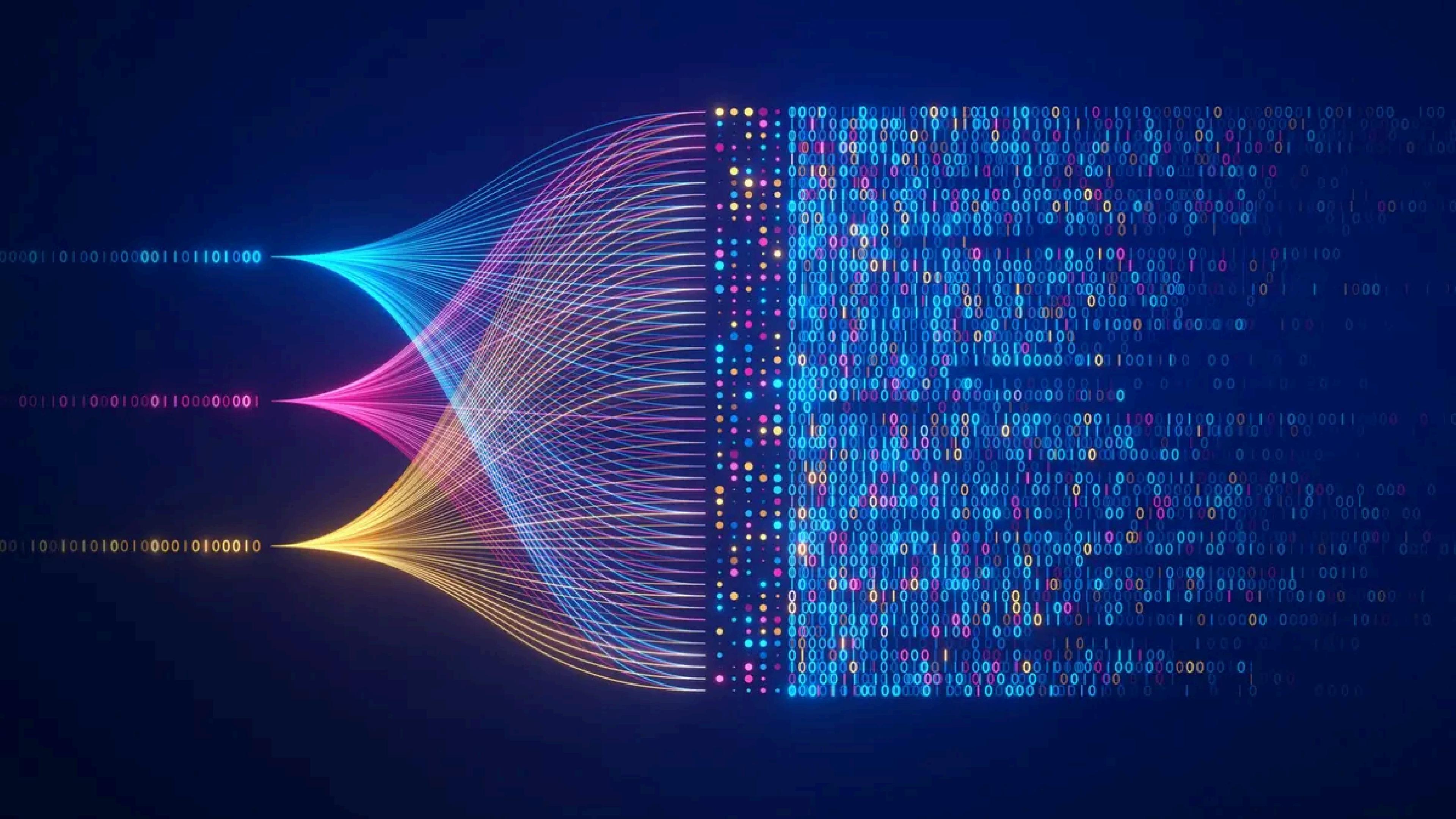
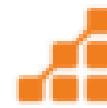


Attention Is All You Need

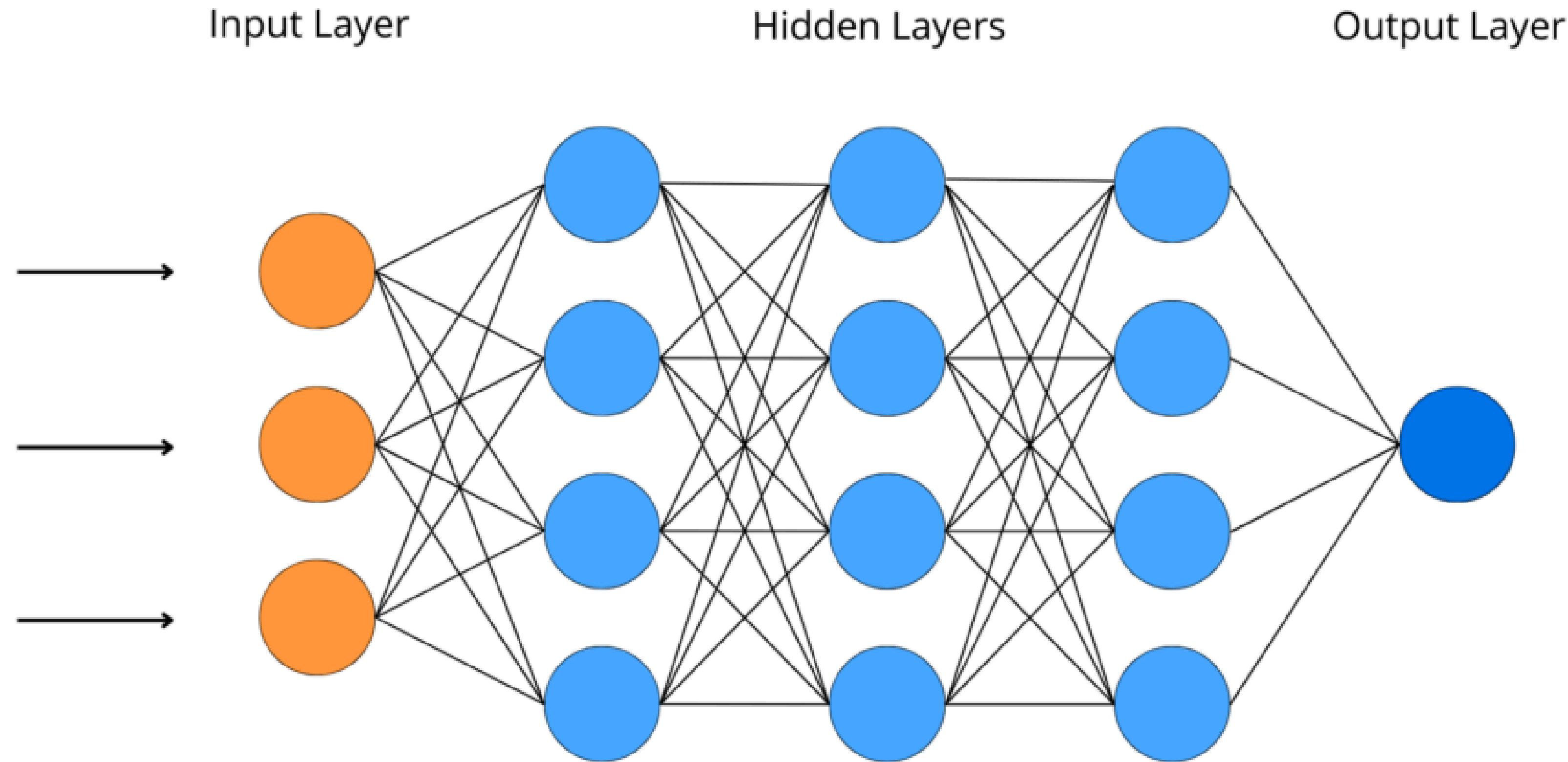
M3hdī Rahmani



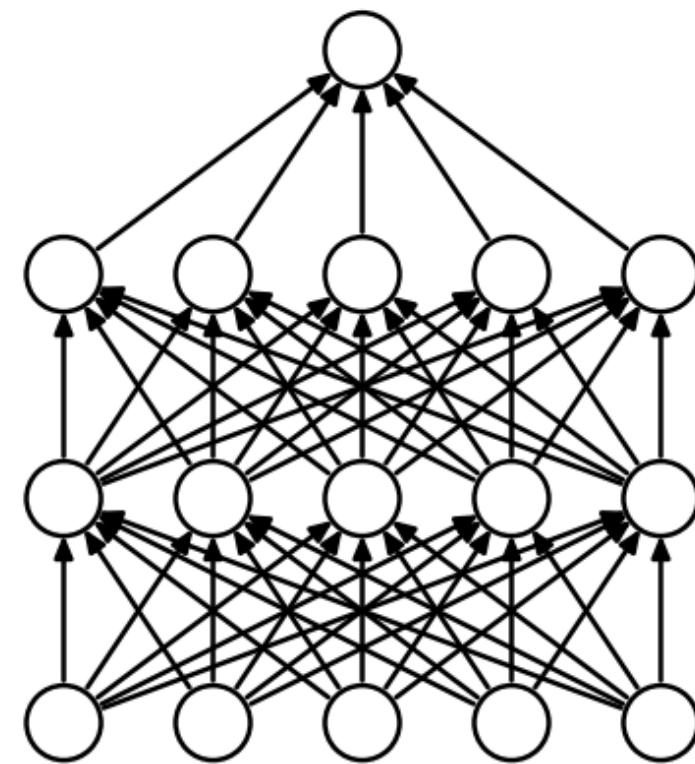




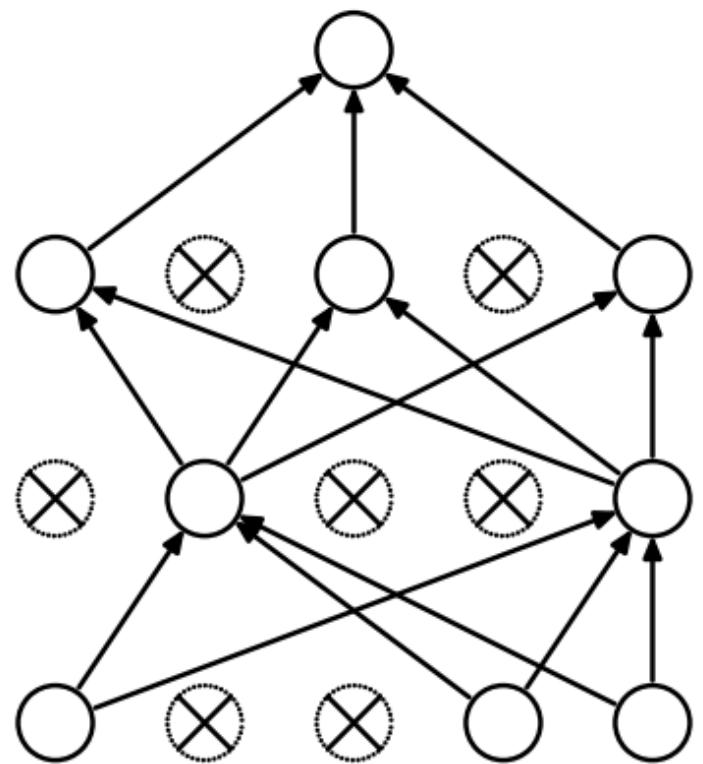
ChatGPT'S Neural Network Architecture



Dropout



(a) Standard Neural Net



(b) After applying dropout.

Left: A standard neural net with **2 hidden layers**.

Right: An example of a thinned net produced by applying dropout to the network on the left. Crossed units have been dropped.

$$P = \bullet / \Lambda$$

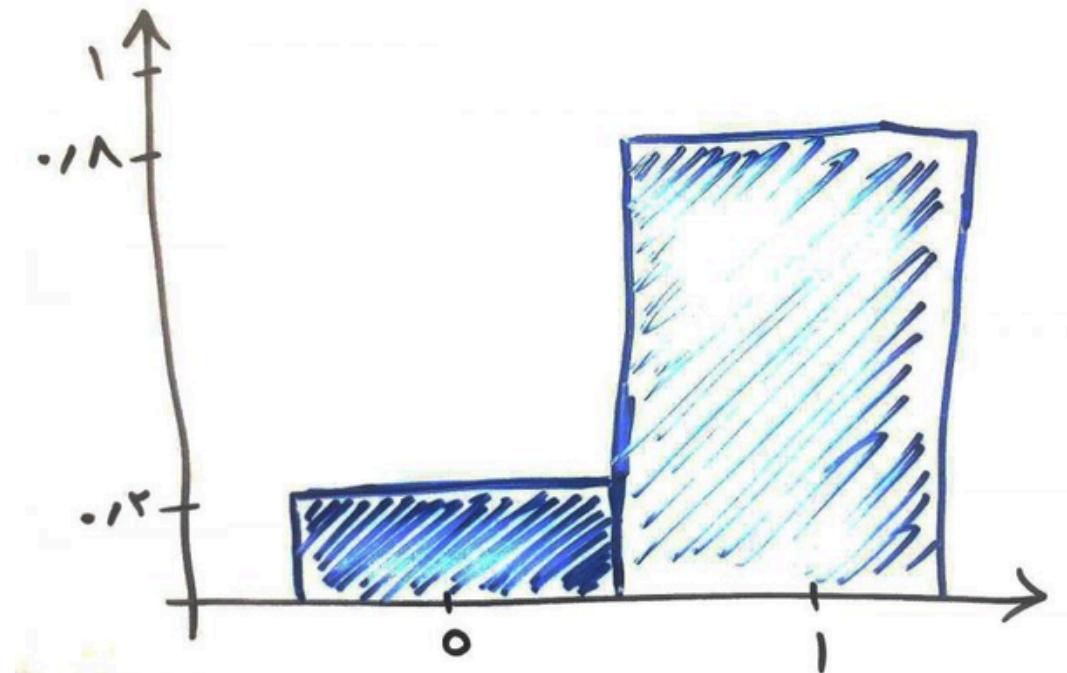
Bernoulli distribution

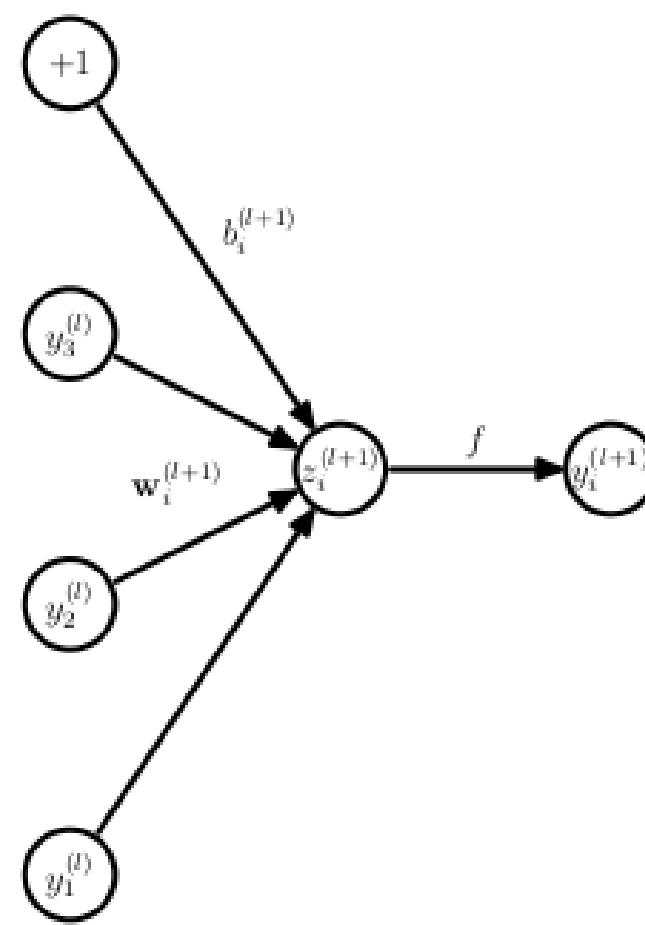
$$P(x^{\text{"success"}}) = \frac{n!}{x!(n-x)!} P^x (1-P)^{(n-x)}$$

تعداد مساحه ادار

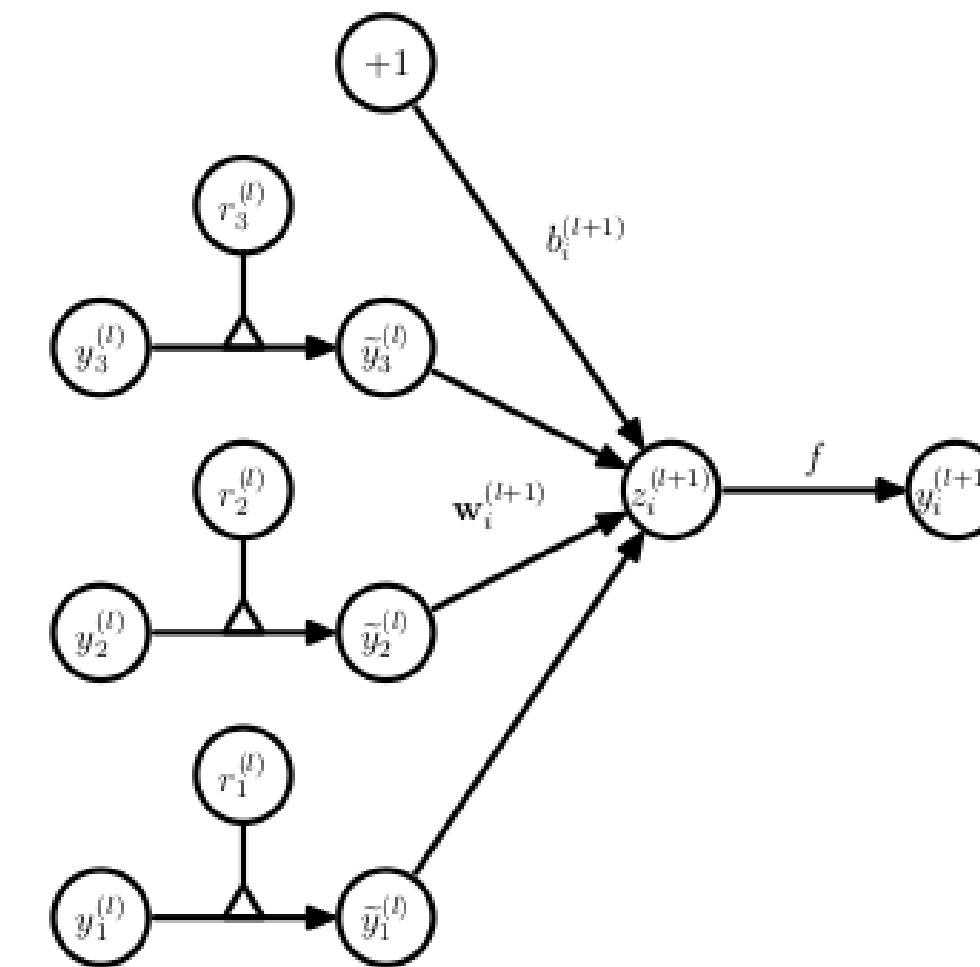
تعداد موفقیت ها و نظر

احتمال موفقیت





(a) Standard network

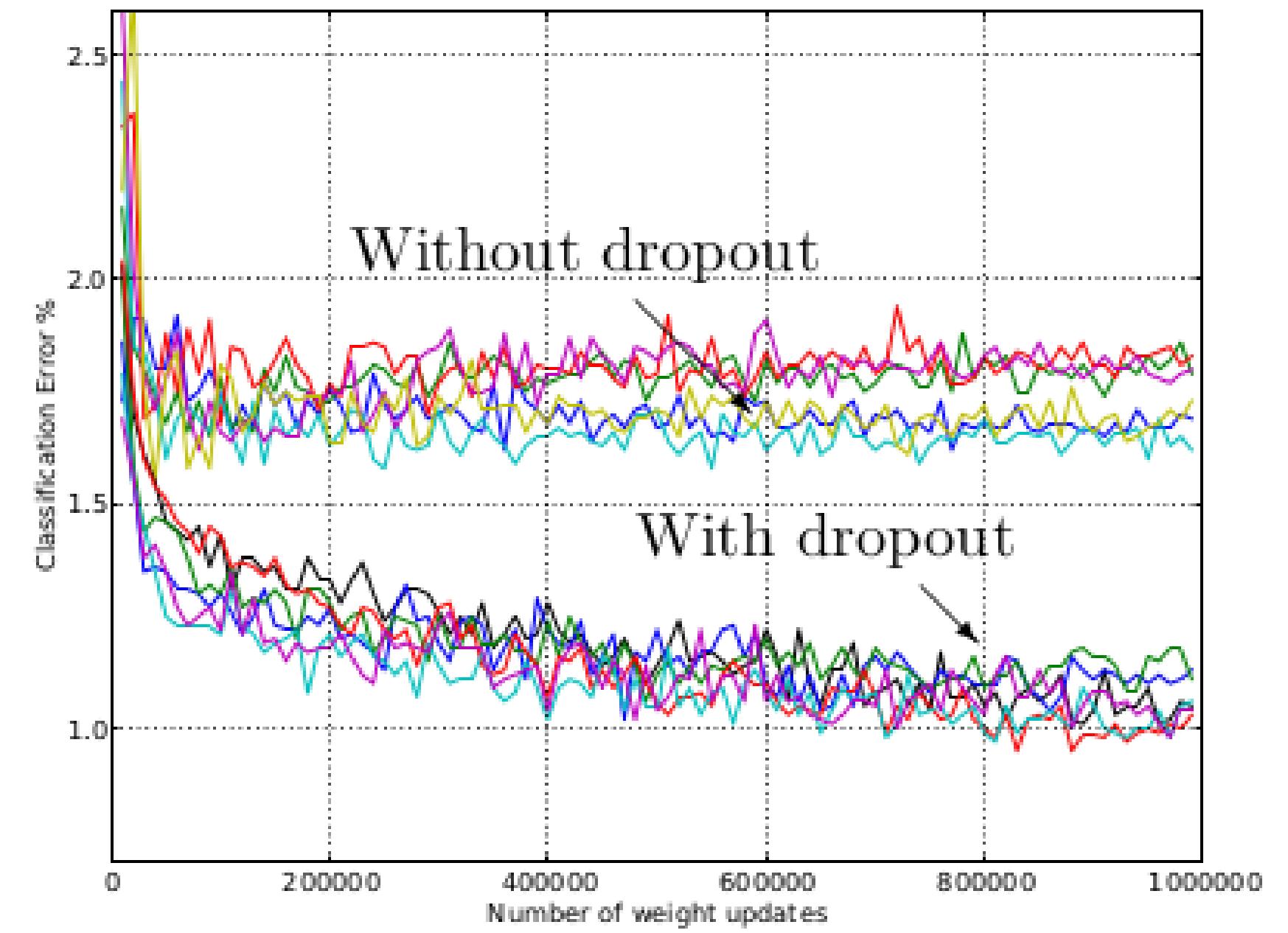


(b) Dropout network

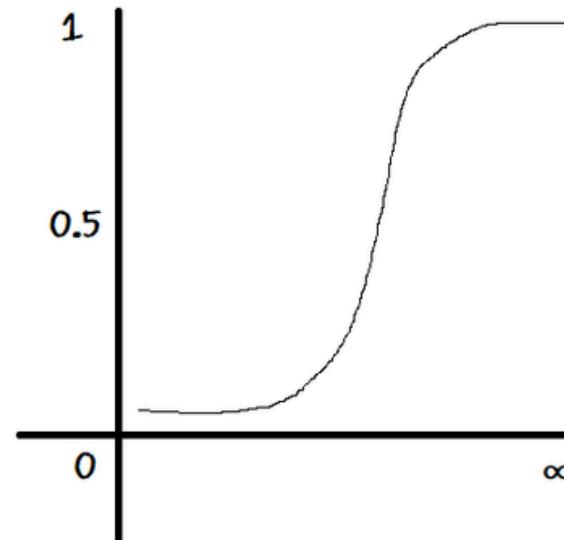
$$\begin{aligned}
 r_j^{(l)} &\sim \text{Bernoulli}(p), \\
 \tilde{\mathbf{y}}^{(l)} &= \mathbf{r}^{(l)} * \mathbf{y}^{(l)}, \\
 z_i^{(l+1)} &= \mathbf{w}_i^{(l+1)} \tilde{\mathbf{y}}^{(l)} + b_i^{(l+1)}, \\
 y_i^{(l+1)} &= f(z_i^{(l+1)}).
 \end{aligned}$$

$$f(x) = 1 / (1 + \exp(-x))$$

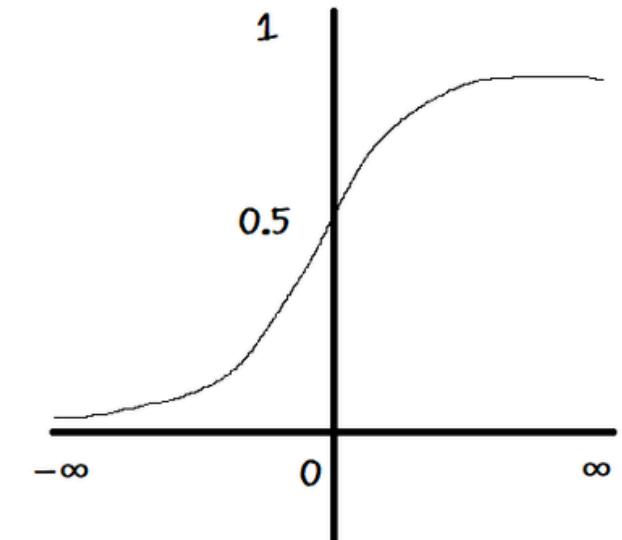
Test error for different architectures



softmax



Sigmoid

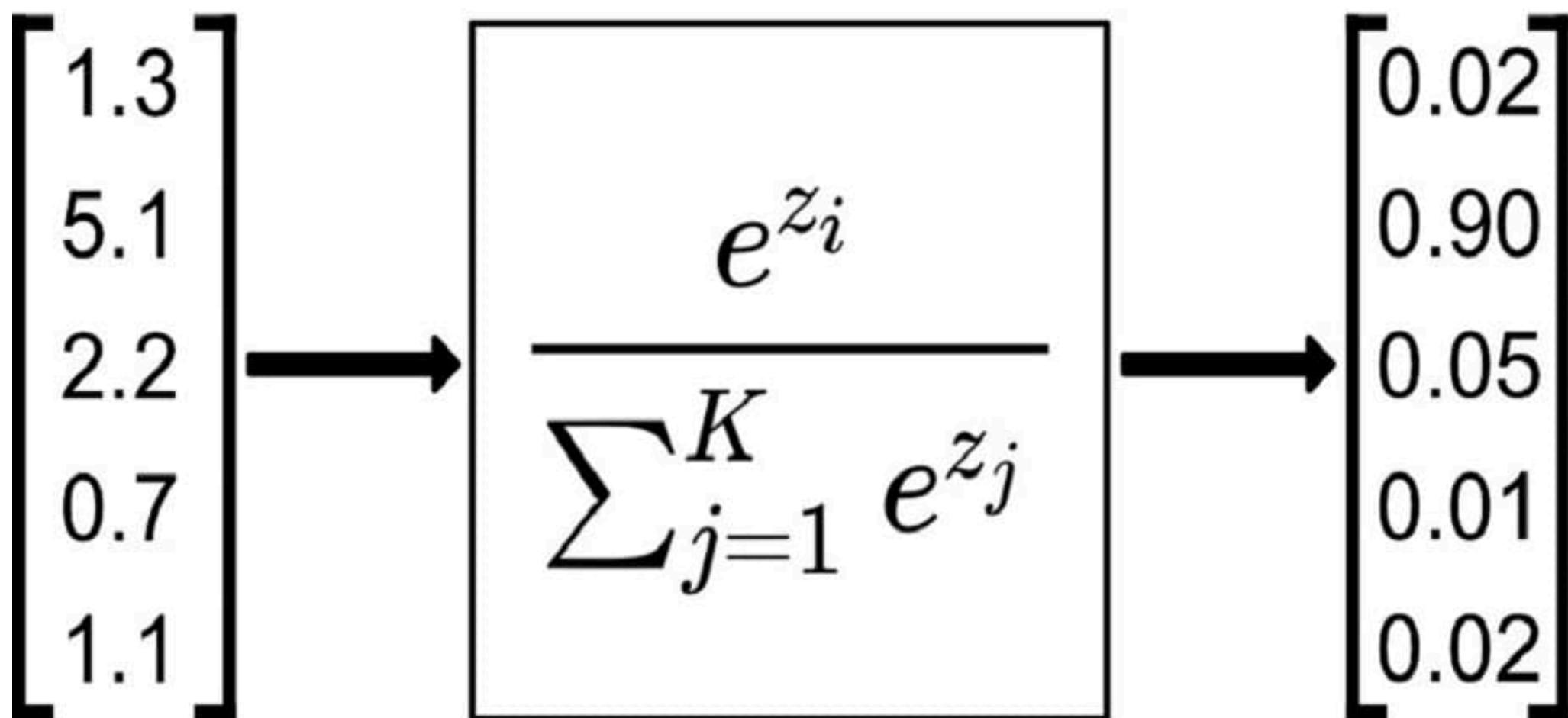


Softmax

Output
layer

Softmax
activation function

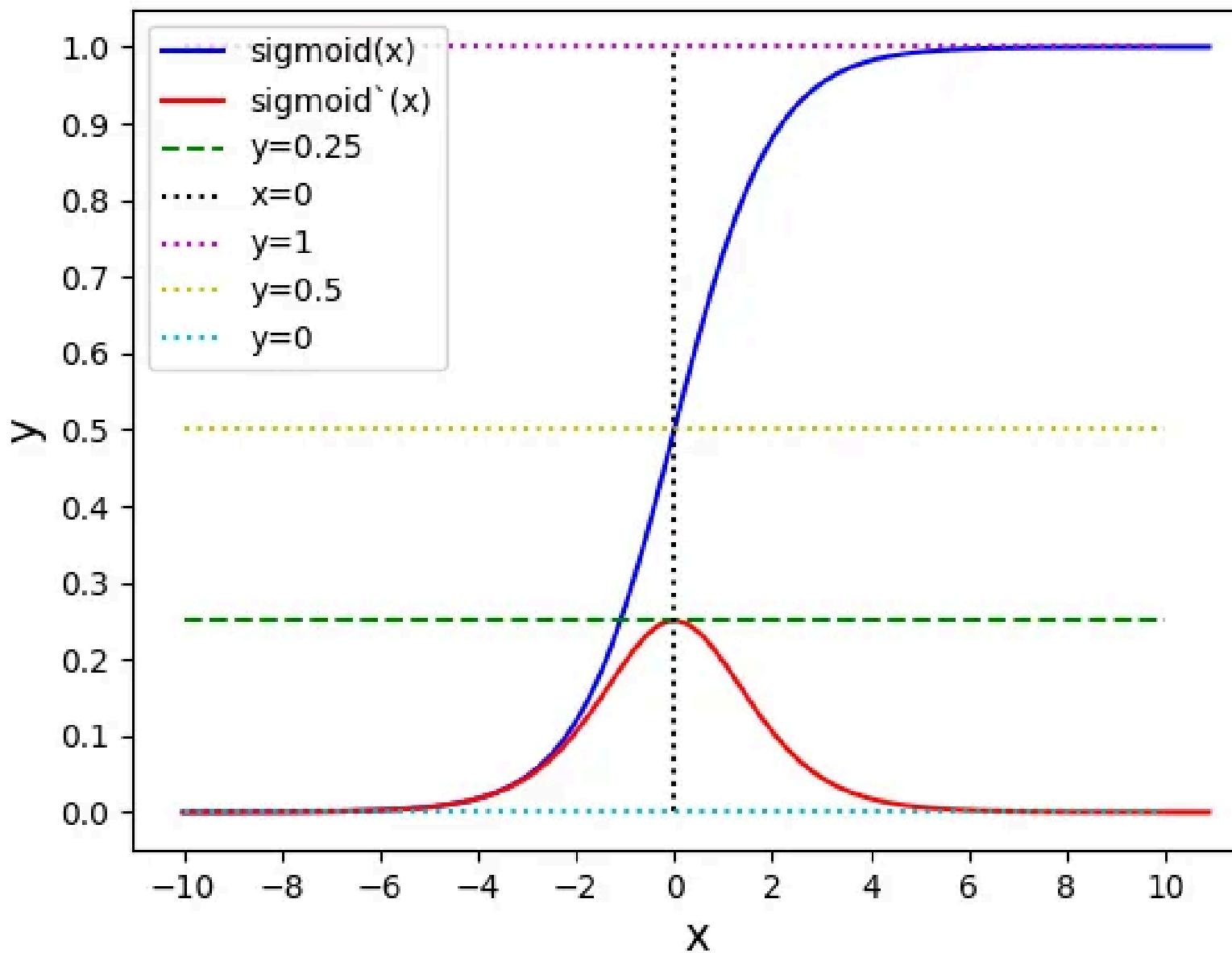
Probabilities



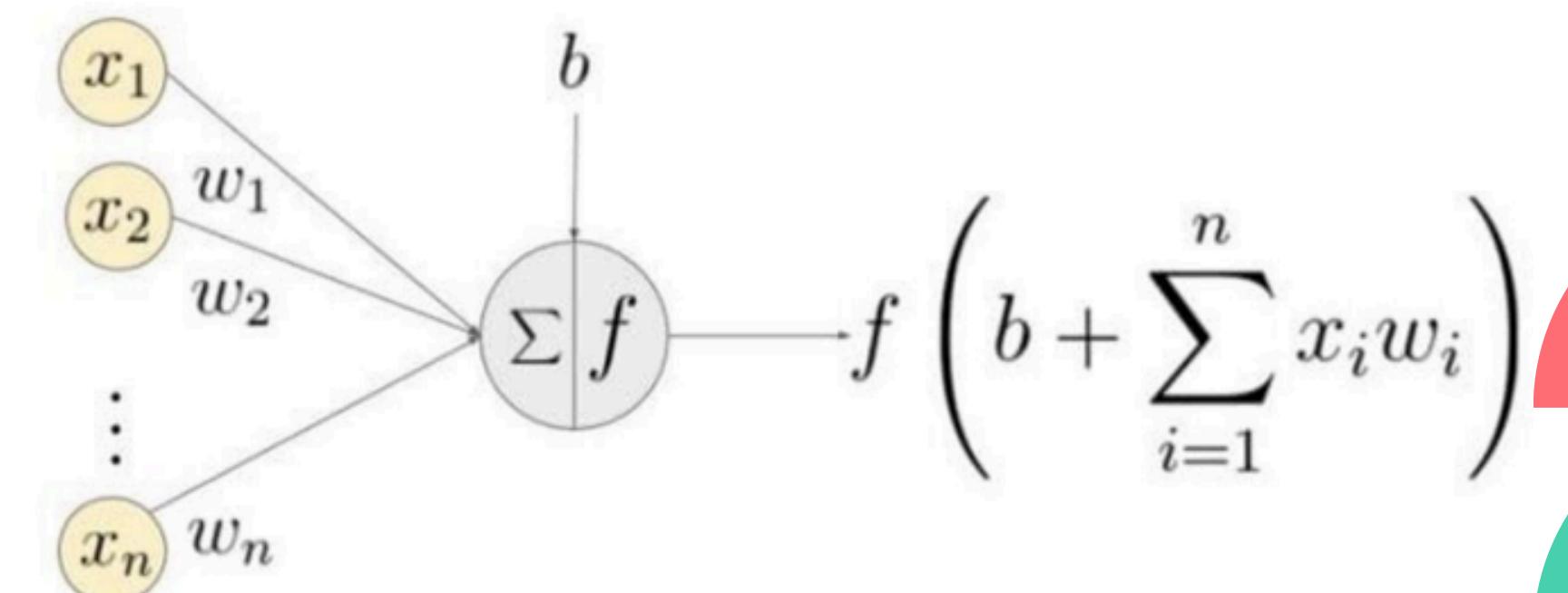
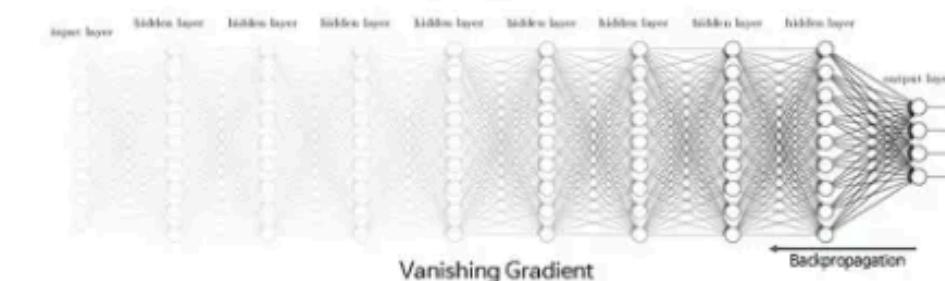
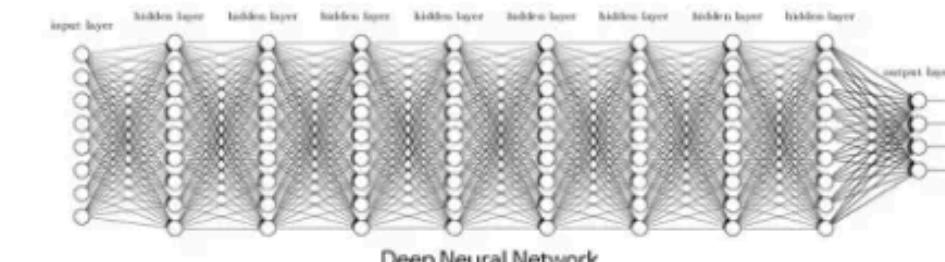
Vanishing Gradient

$$W_{\text{new}} = W_{\text{old}} - \eta * \frac{\partial C}{\partial w}$$

learning rate



Sigmoid: Vanishing Gradient Problem



Method to overcome the problem

01 - Activation function

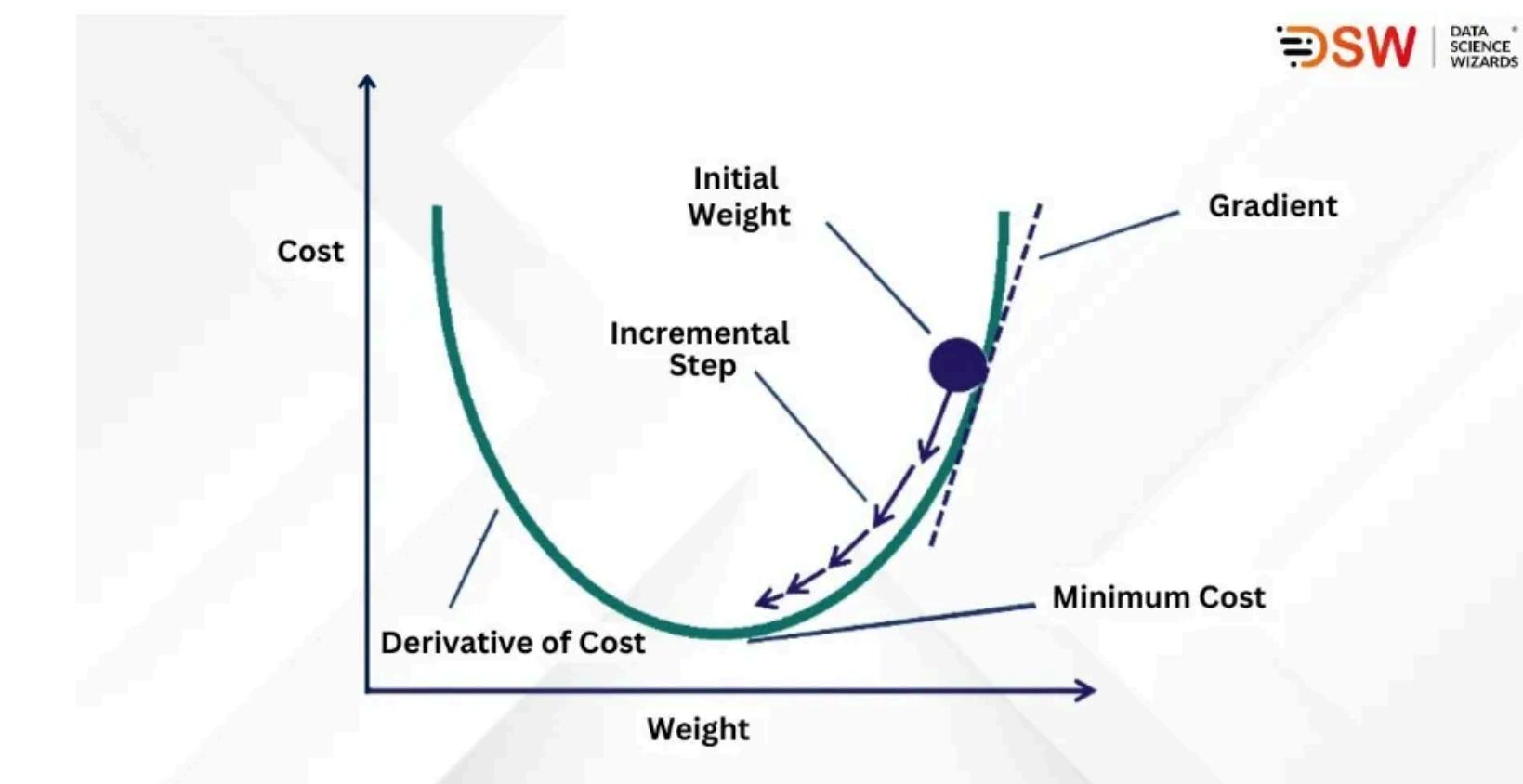
ReLU activation functions keep linearity for regions where sigmoid and TanH are saturated, thus responding better to gradient vanishing / exploding. You can also use different types like Leaky-ReLu, Randomized ReLu, etc.

02 - Cross Entropy Loss Functions

$$L = -\frac{1}{N} \left[\sum_{j=1}^N [t_j \log(p_j) + (1 - t_j) \log(1 - p_j)] \right]$$

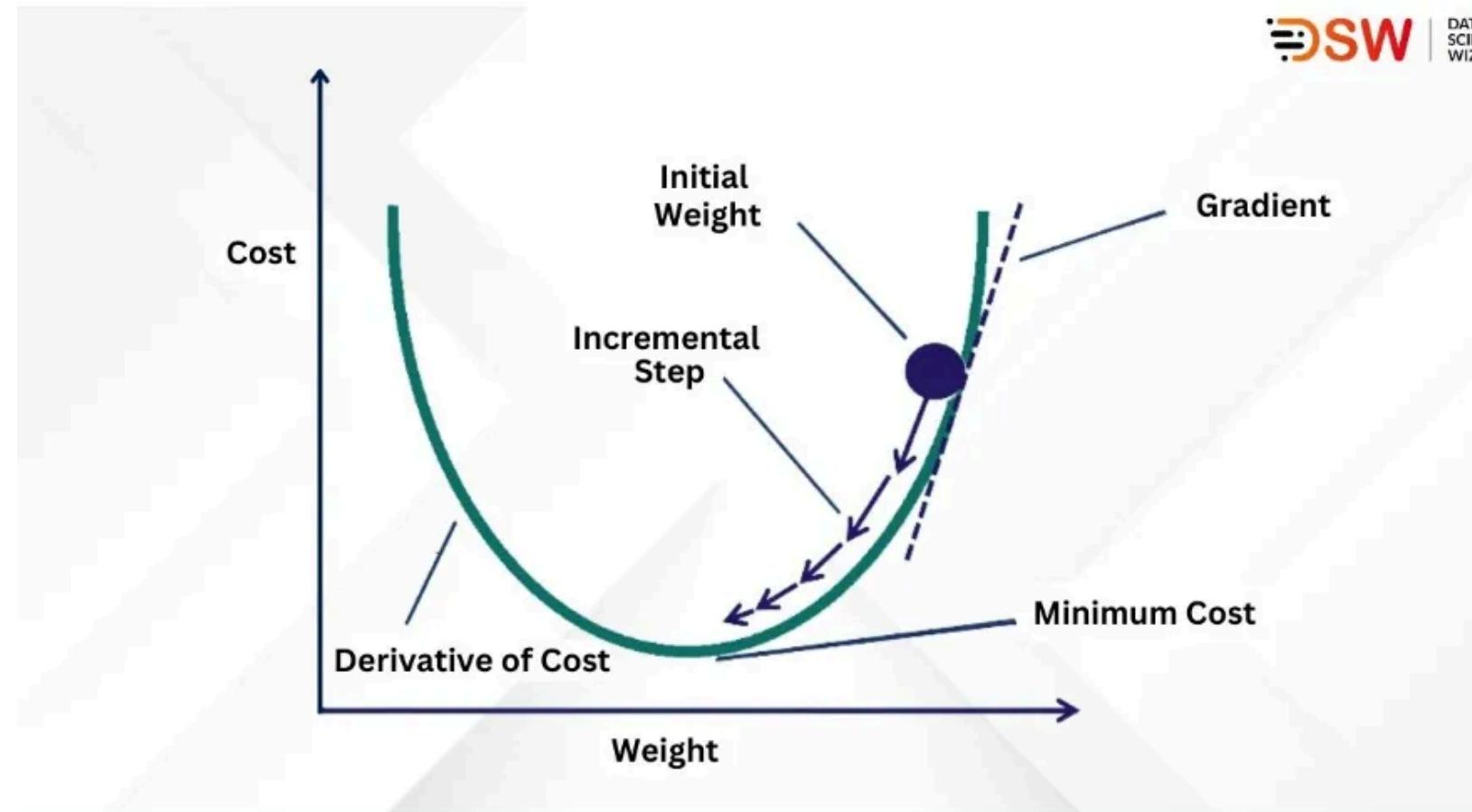
for N data points where t_i is the truth value taking a value 0 or 1 and p_i is the Softmax probability for the i^{th} data point.

03 - Gradient Descent



Gradient Descent Explained with an Example

مقداردهی اولیه



$$X_0 = 3$$

Learning rate = 0.01

$$\frac{dy}{dx} = \frac{d}{dx}(x + 5)^2 = 2 * (x + 5)$$

تکرار اول

$$X_1 = X_0 - (\text{learning rate}) * \left(\frac{dy}{dx}\right)$$

$$X_1 = 3 - (0.01) * (2 * (3 + 5)) = 2.84$$

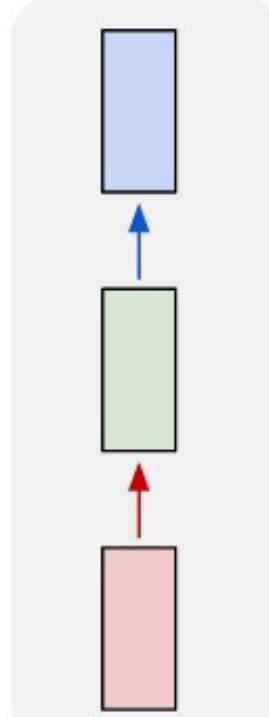
تکرار دوم

$$X_2 = X_1 - (\text{learning rate}) * \left(\frac{dy}{dx}\right)$$

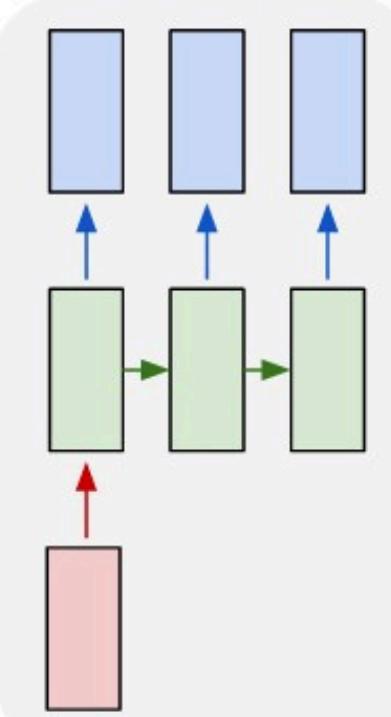
$$X_2 = 2.84 - (0.01) * (2 * (2.84 + 5)) = 2.6832$$

RNN & LSTM

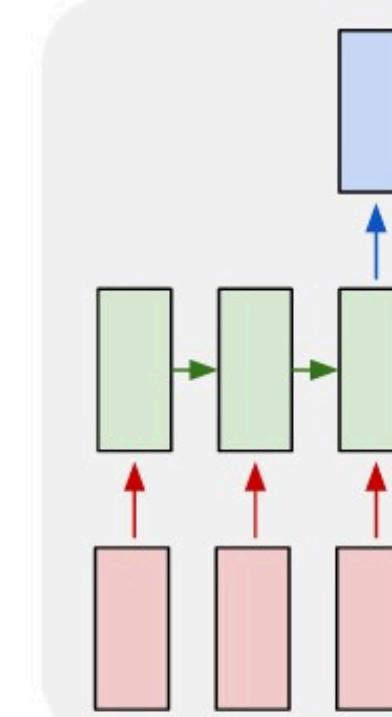
one to one



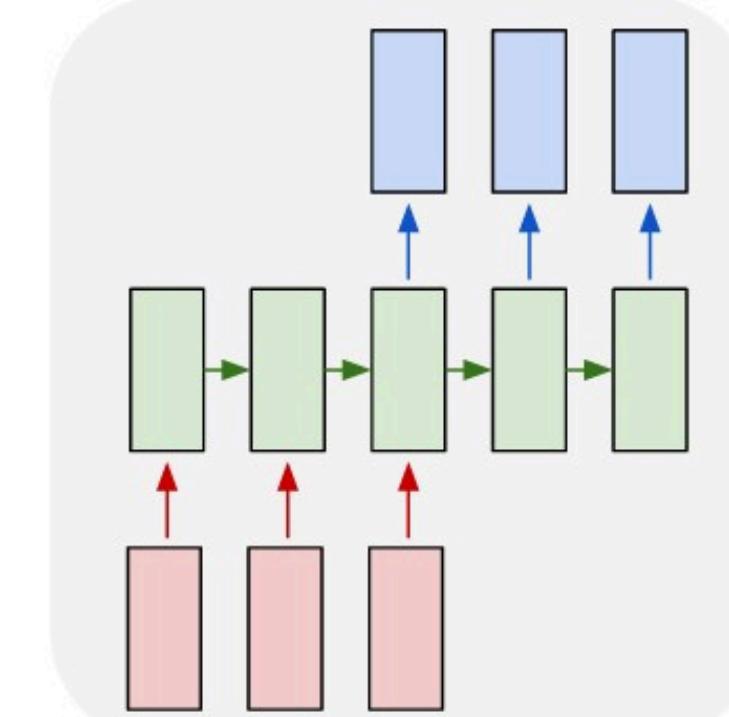
one to many



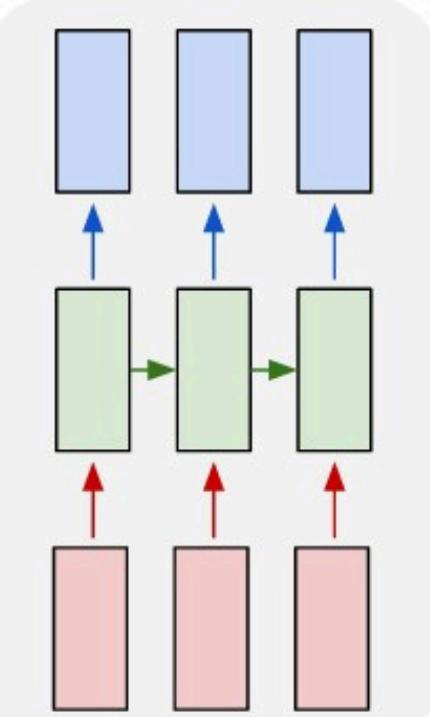
many to one



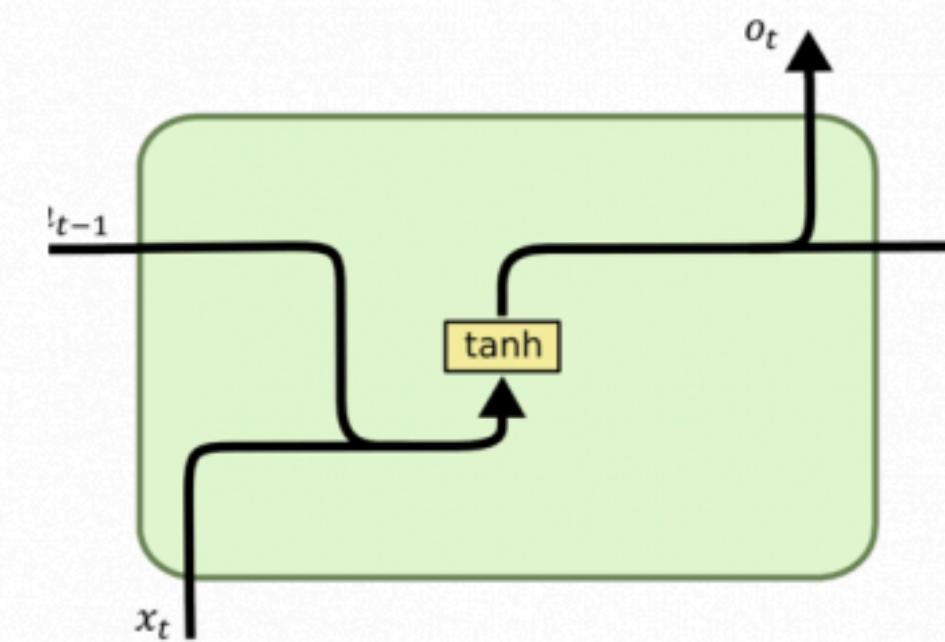
many to many



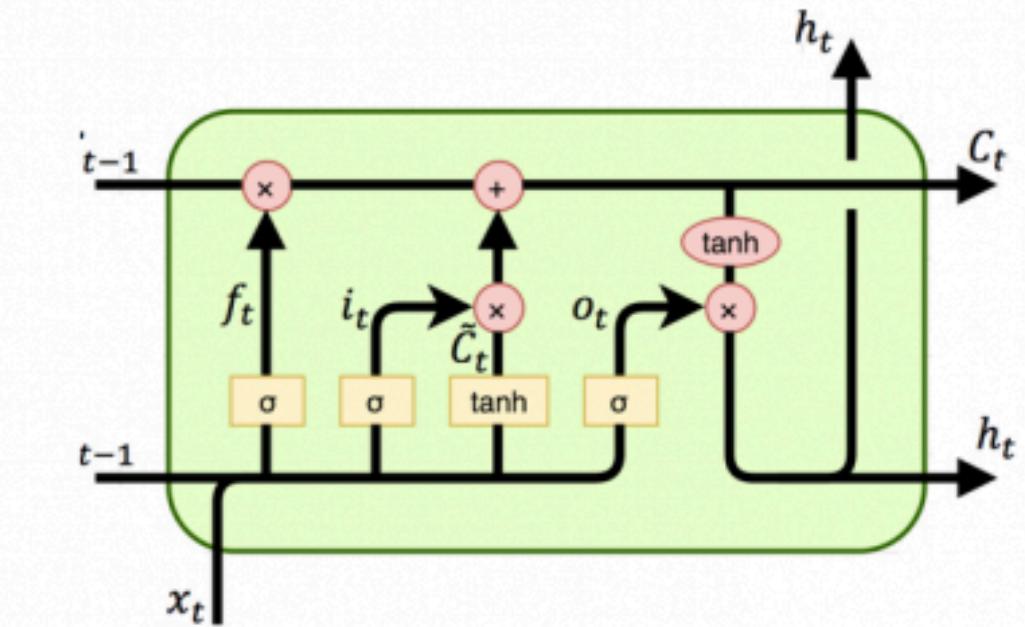
many to many



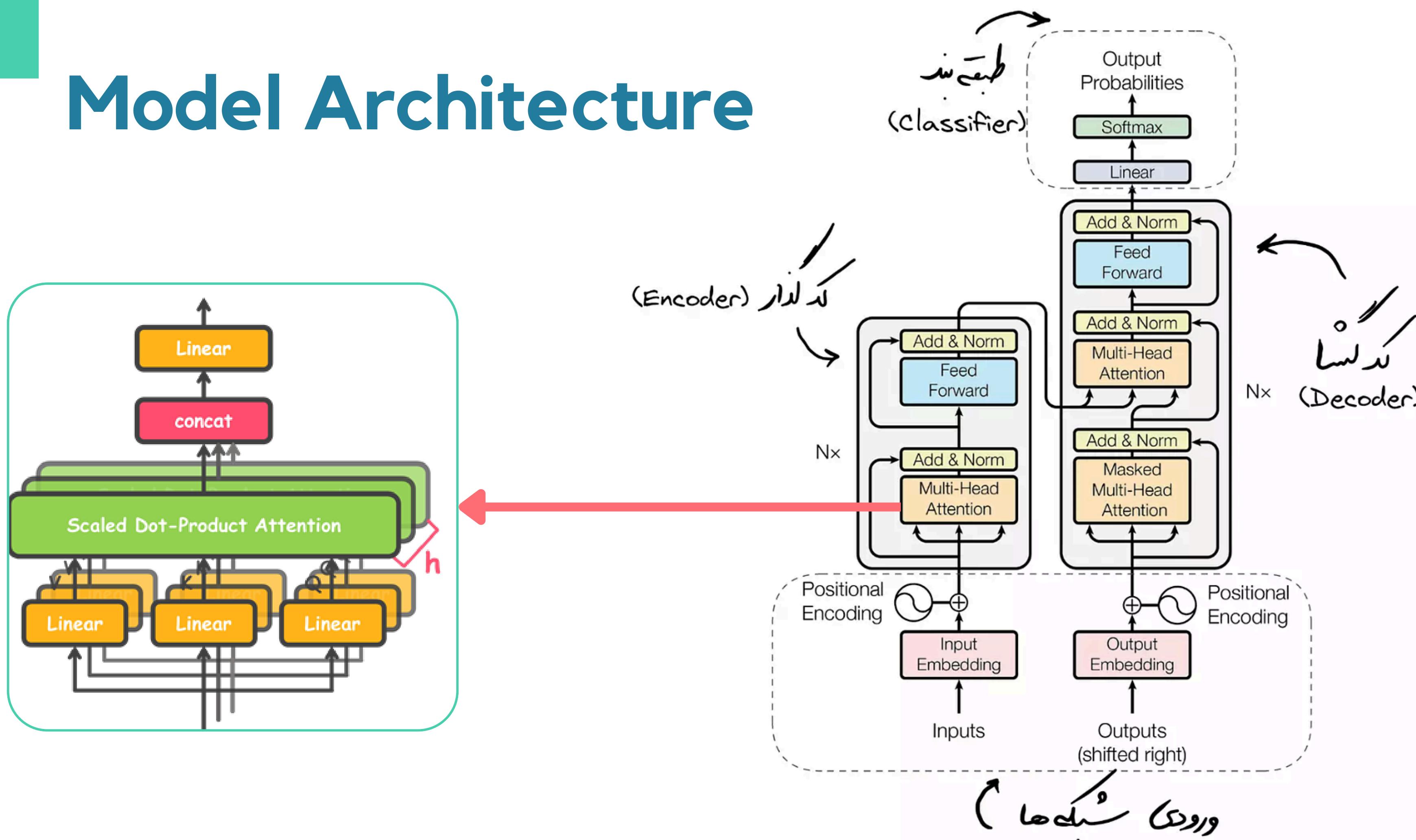
RNN



LSTM



Model Architecture



- “I like deep learning.”
- “I like NLP.”
- “I enjoy flying.”

Tokenization



```
[ "I", "like", "deep", "learning", ".", "NLP",  
  "enjoy", "flying" ]
```

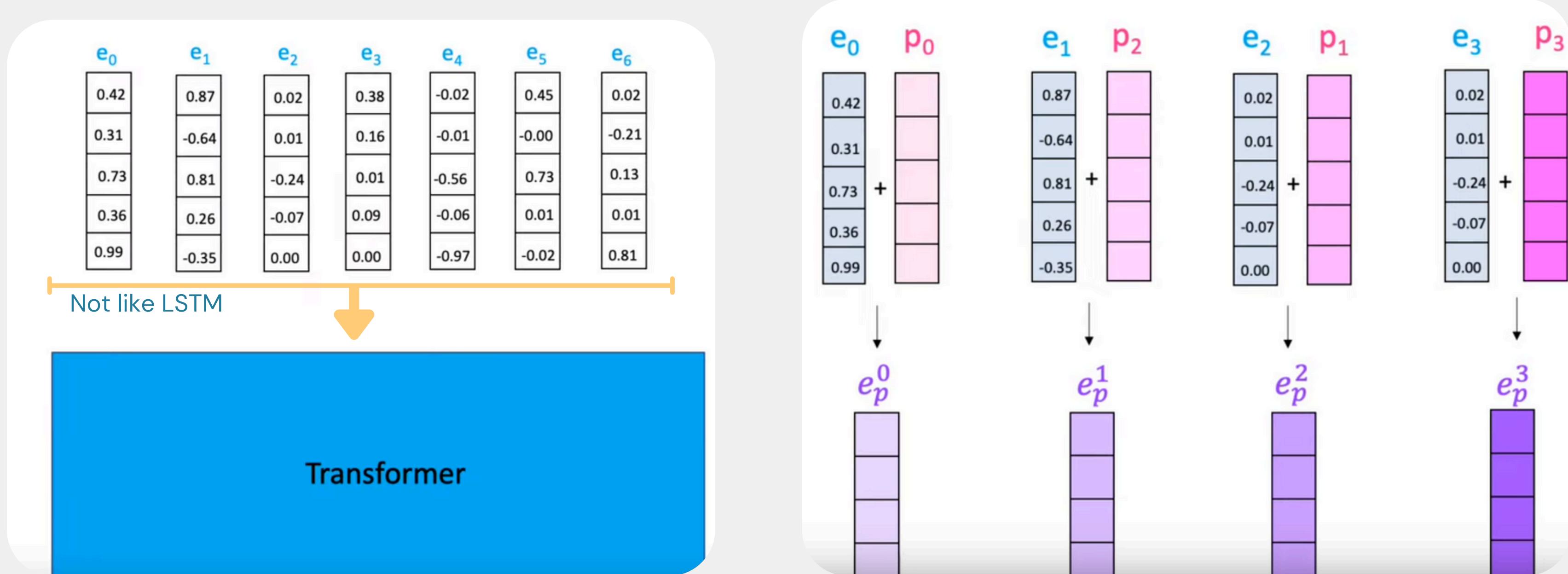
Vectorization (GloVe)

- “I like deep learning.”
- “I like NLP.”
- “I enjoy flying.”

[
 "I", "like", "deep",
 "learning", ".", "NLP",
 "enjoy", "flying"
]

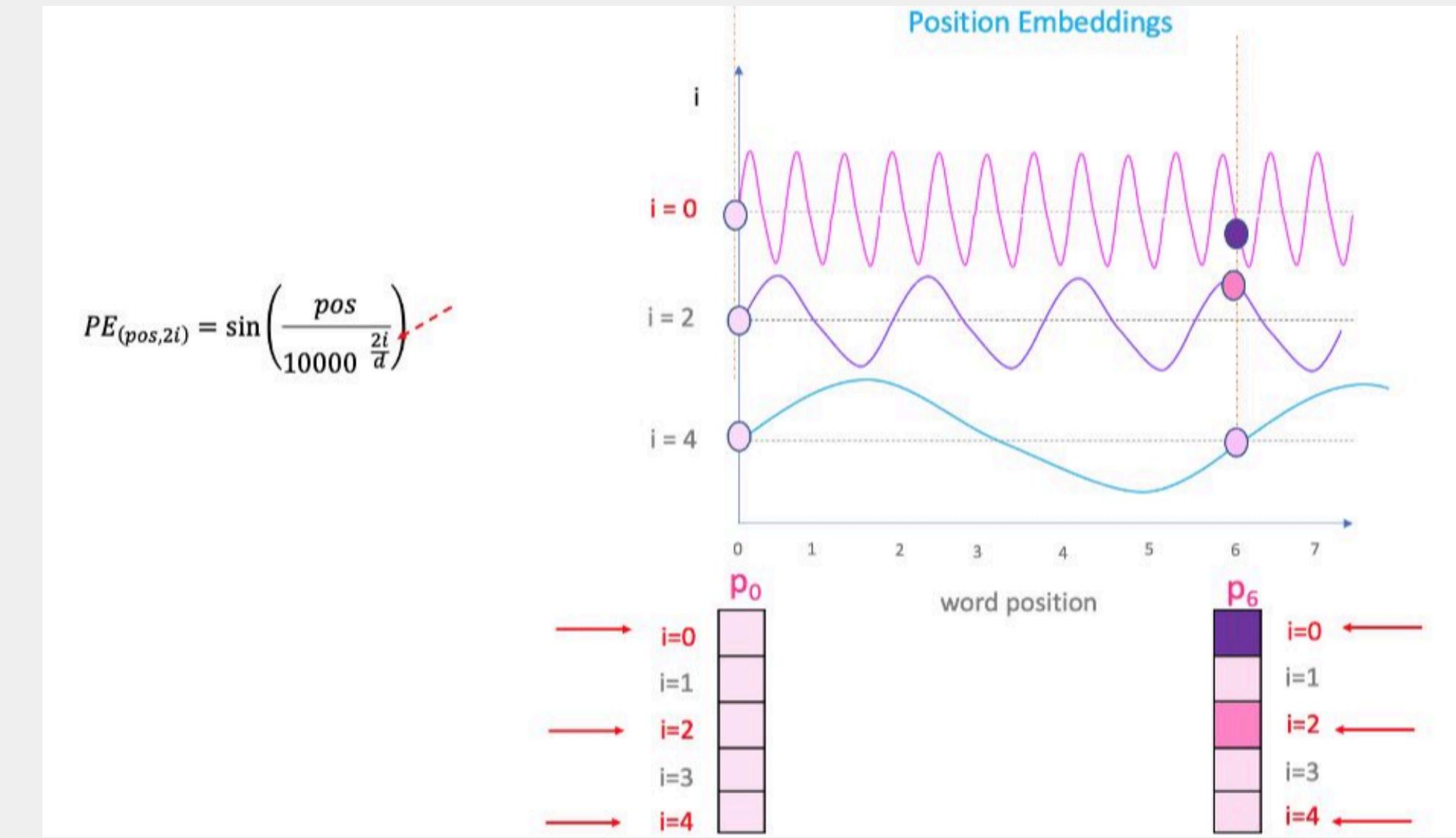
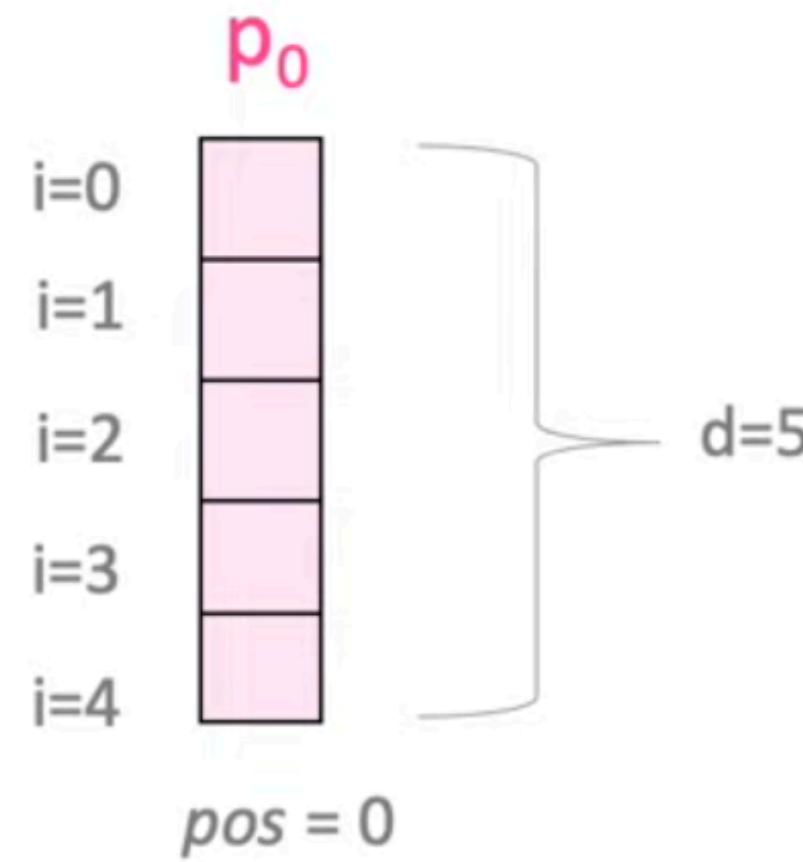
$$X = \begin{matrix} & I & like & enjoy & deep & learning & NLP & flying & . \\ I & 0 & 2 & 1 & 0 & 0 & 0 & 0 & 0 \\ like & 2 & 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ enjoy & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\ deep & 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 \\ learning & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 \\ NLP & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ flying & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 1 \\ . & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \end{matrix}$$

Positional Encoding





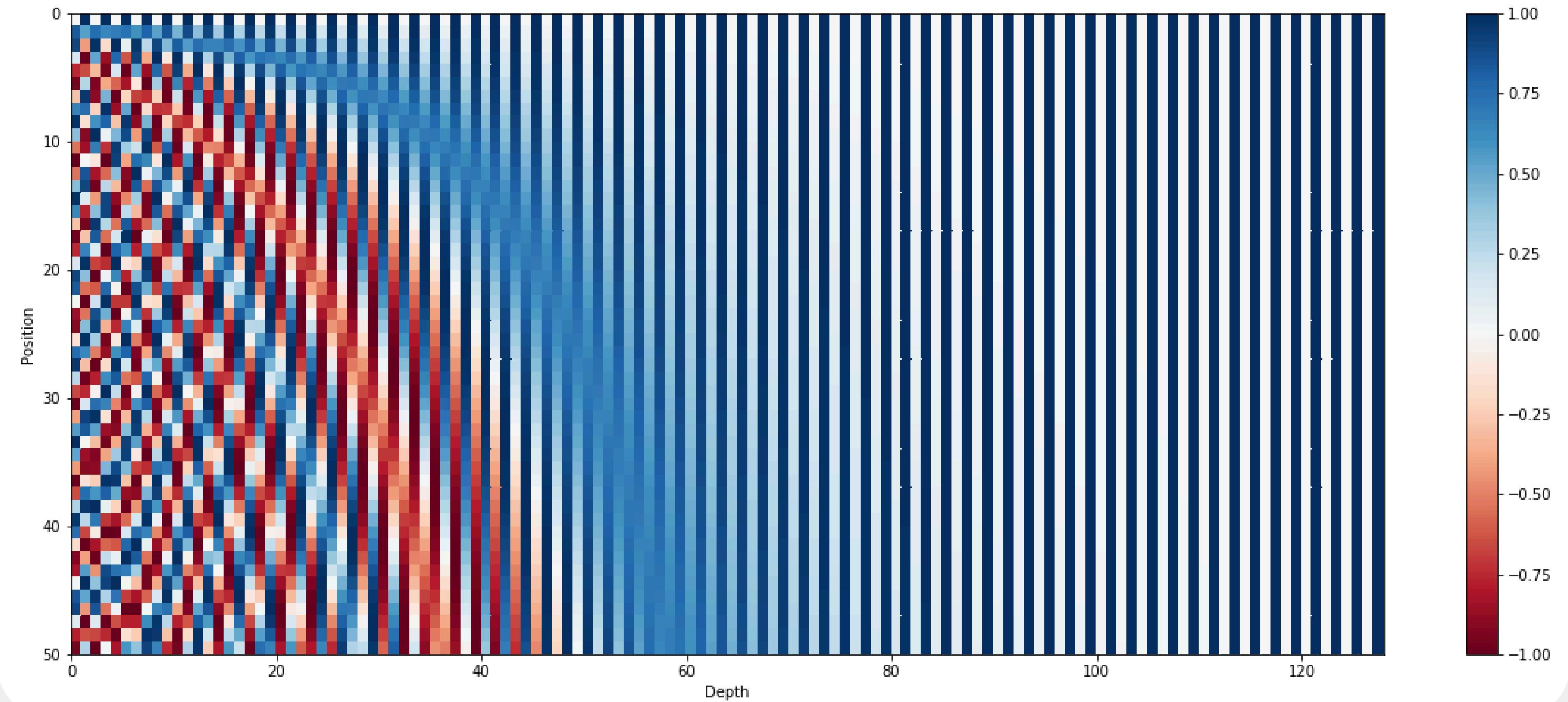
Positional Encoding



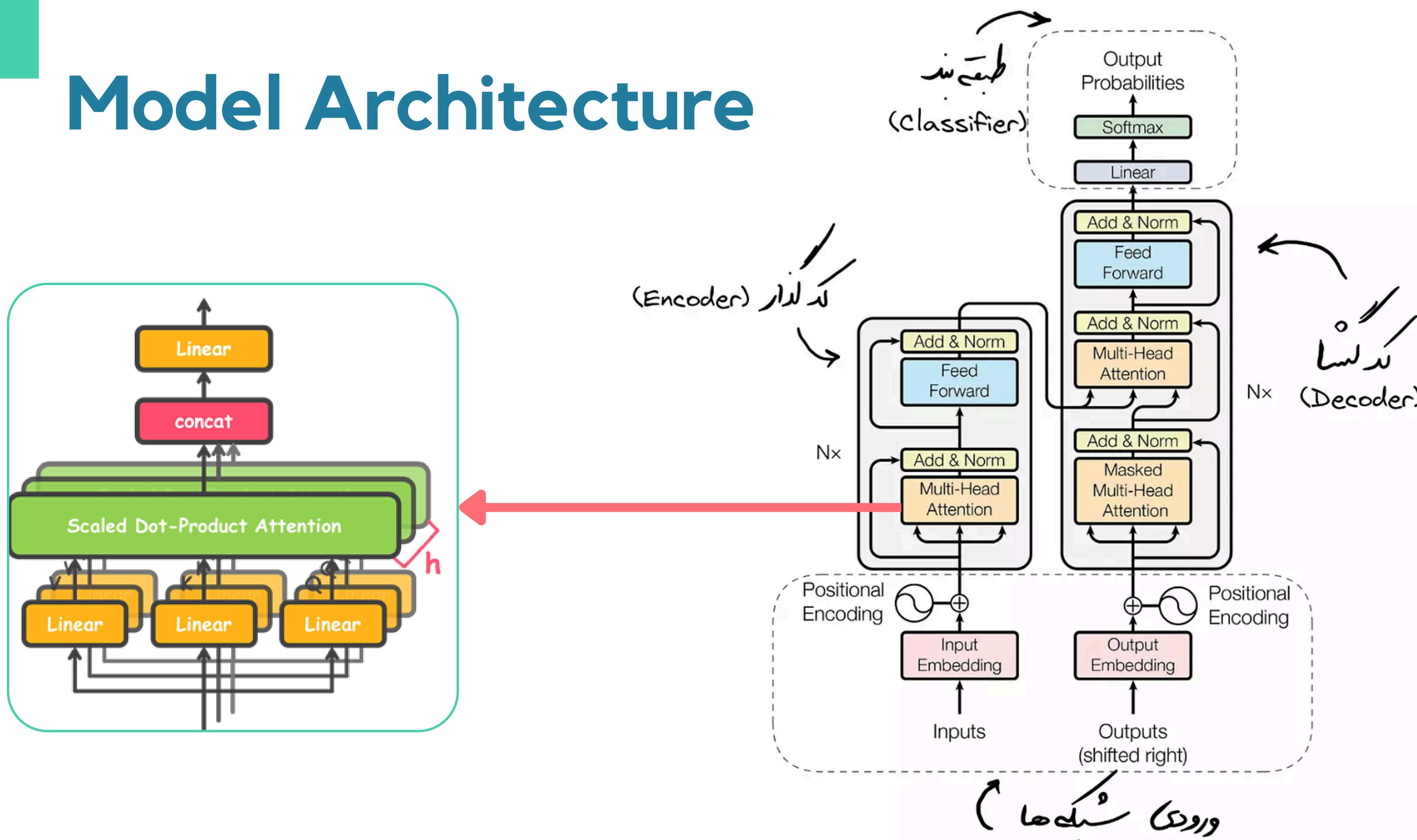
$$PE(pos, 2i) = \sin\left(\frac{pos}{10000} \frac{2i}{d_{model}}\right)$$

$$PE(pos, 2i+1) = \cos\left(\frac{pos}{10000} \frac{2i}{d_{model}}\right)$$

Positional Encoding

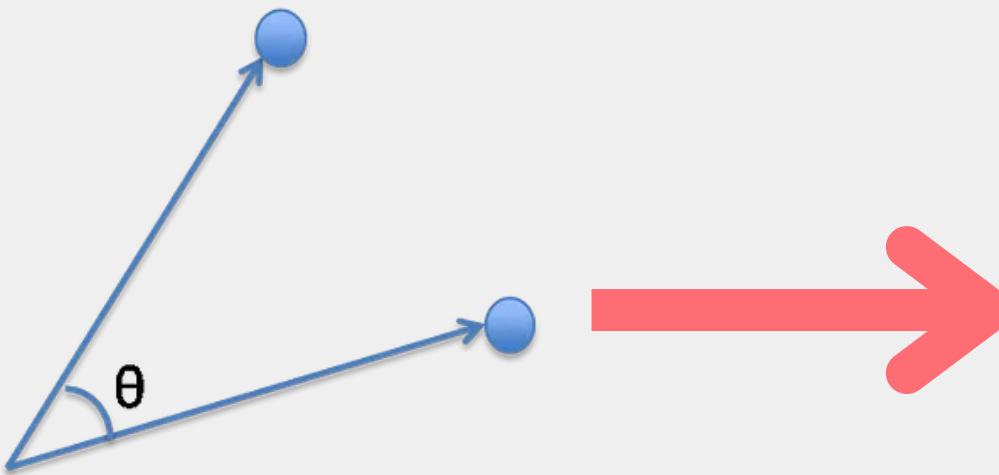


Model Architecture



Cosine Similarity

$$sim(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



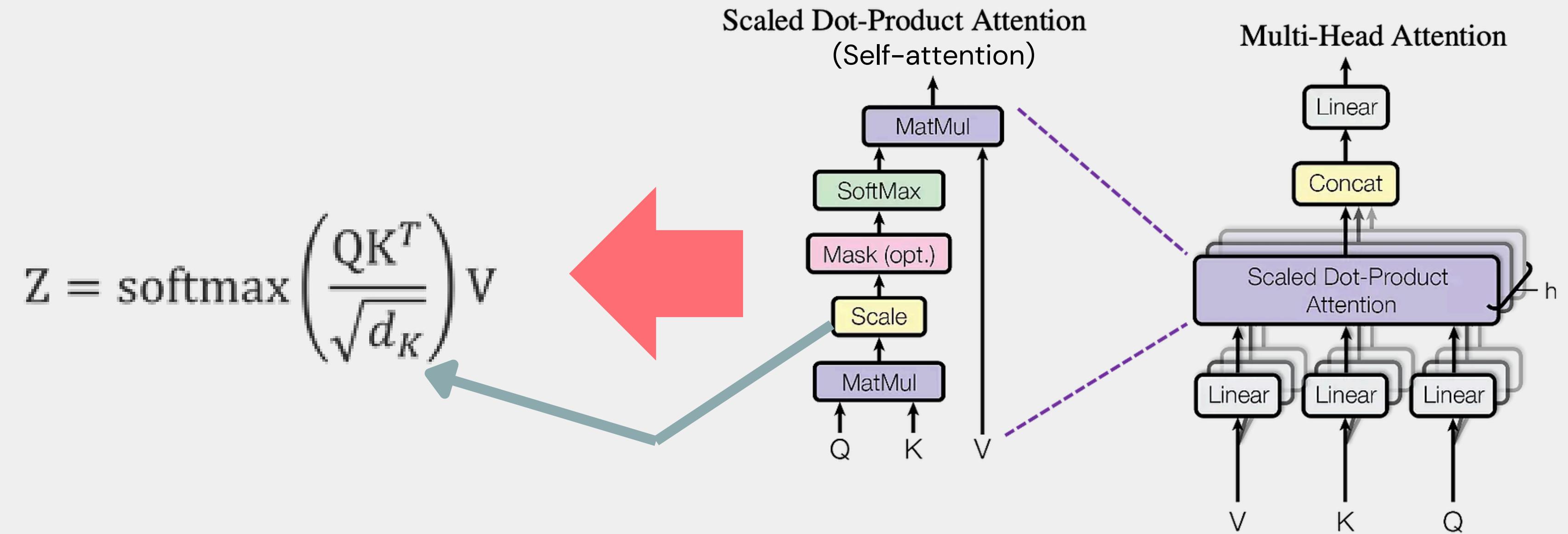
$$Similarity(A, B) = \frac{A \cdot B}{scaling}$$

$$Similarity(A, B) = \frac{A \cdot B^T}{scaling}$$

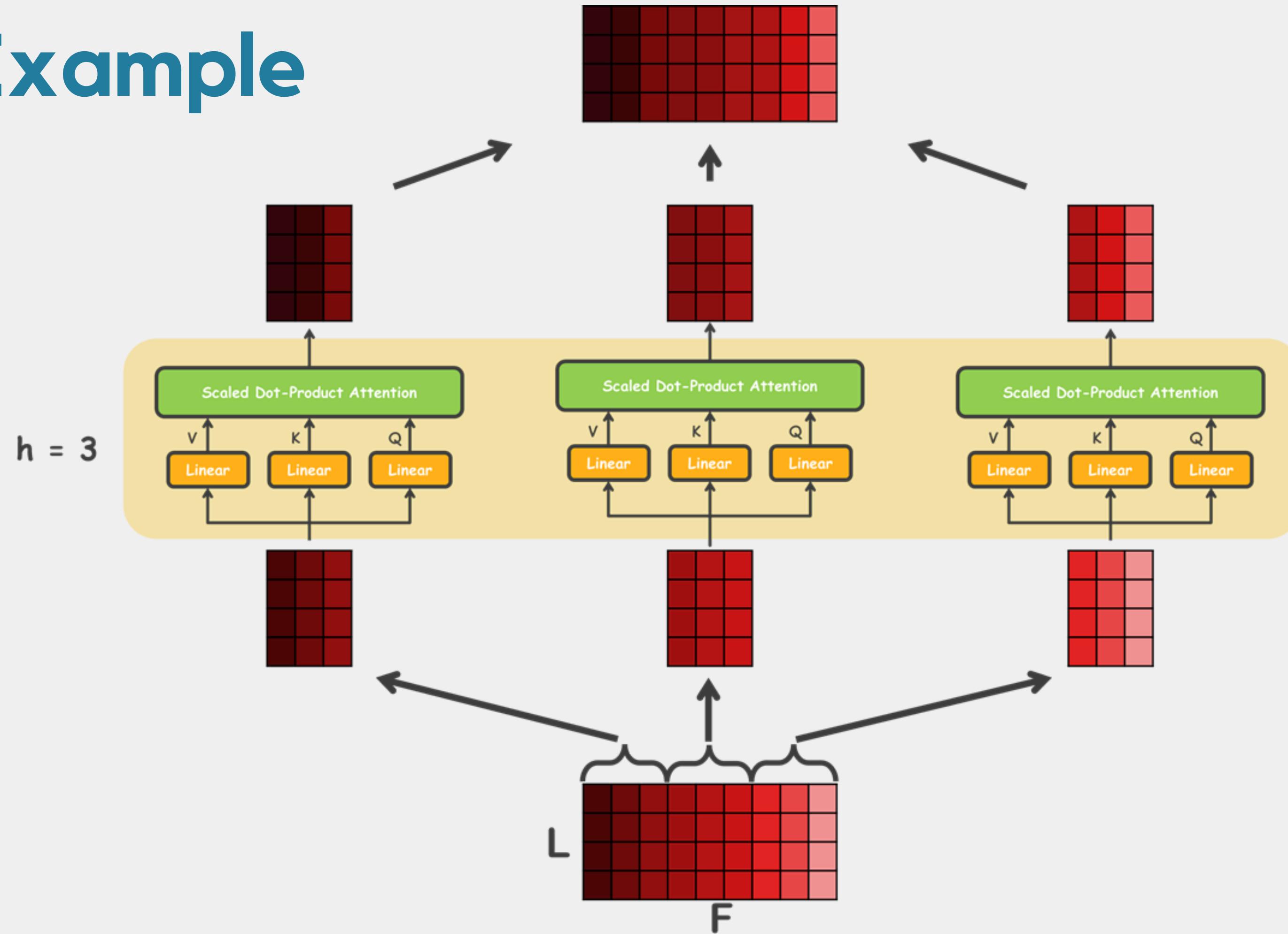


$$Similarity(Q, K) = \frac{Q \cdot K^T}{scaling}$$

Multi-Head Attention

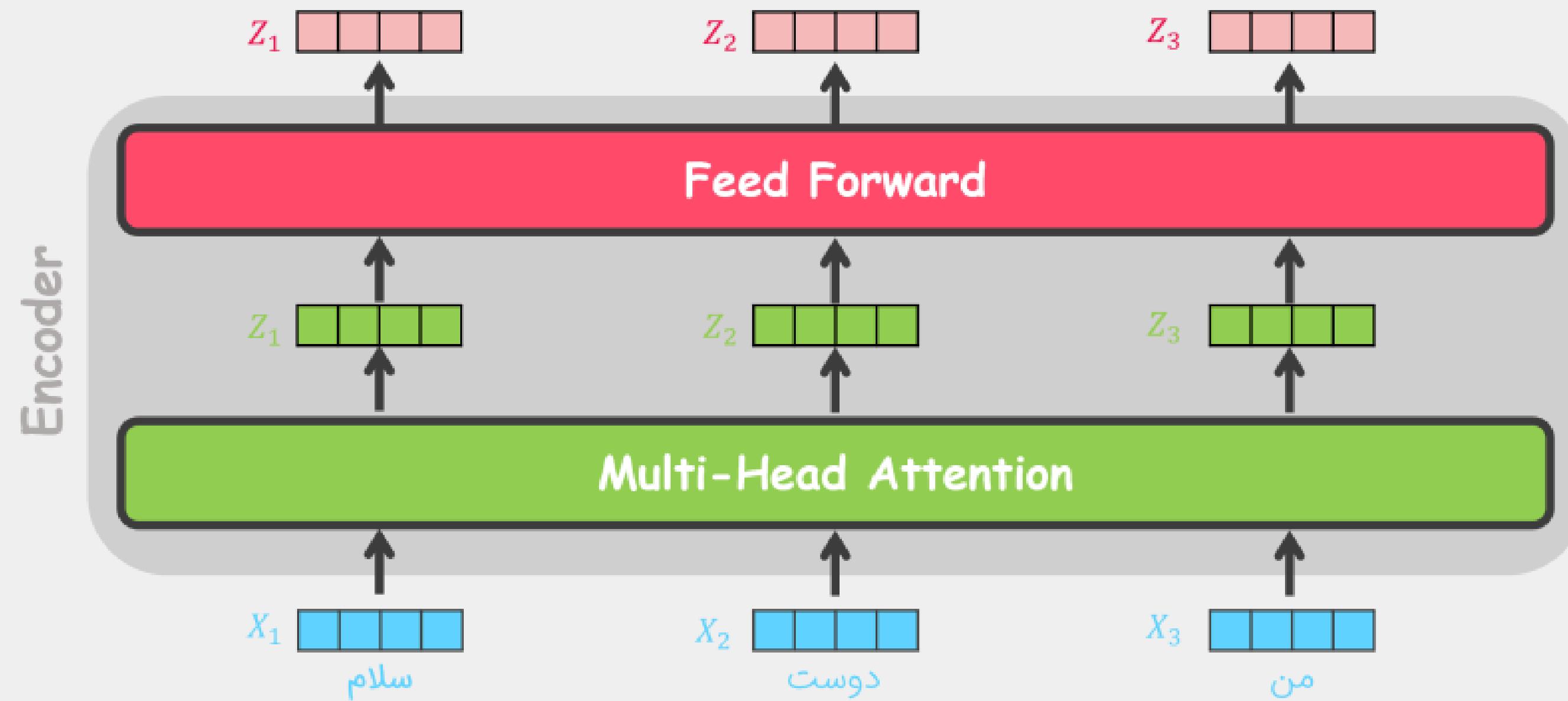


Example



Feed Forward

پس از اتمام مراحل قبل، نتیجه برای پردازش بیشتر به یک شبکه‌ی پیش‌خور داده می‌شود که متشکل از چندین لایه با تابع فعال‌ساز ReLU می‌باشد. هم‌چنین، مراحل استفاده از اتصالات اضافی (residual connection) و نرمال‌سازی نیز مجدد تکرار می‌شوند.



Neural Machine Translation

SEQUENCE TO SEQUENCE MODEL WITH ATTENTION

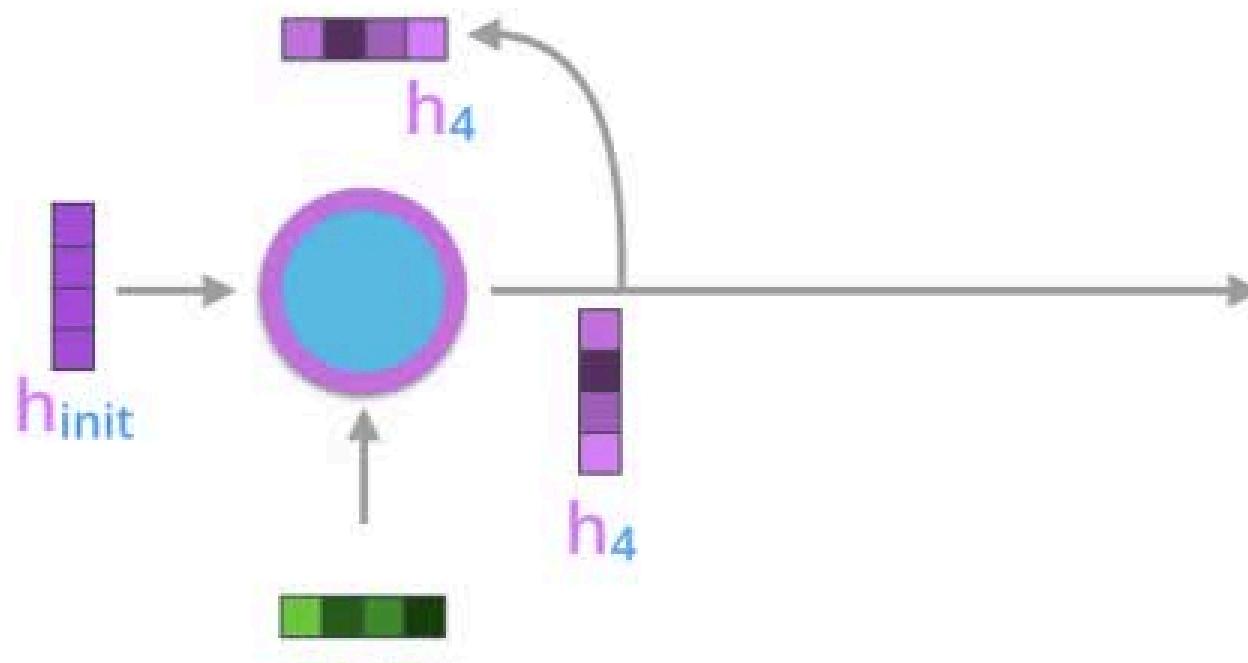
Encoding Stage



$h_1 \ h_2 \ h_3$

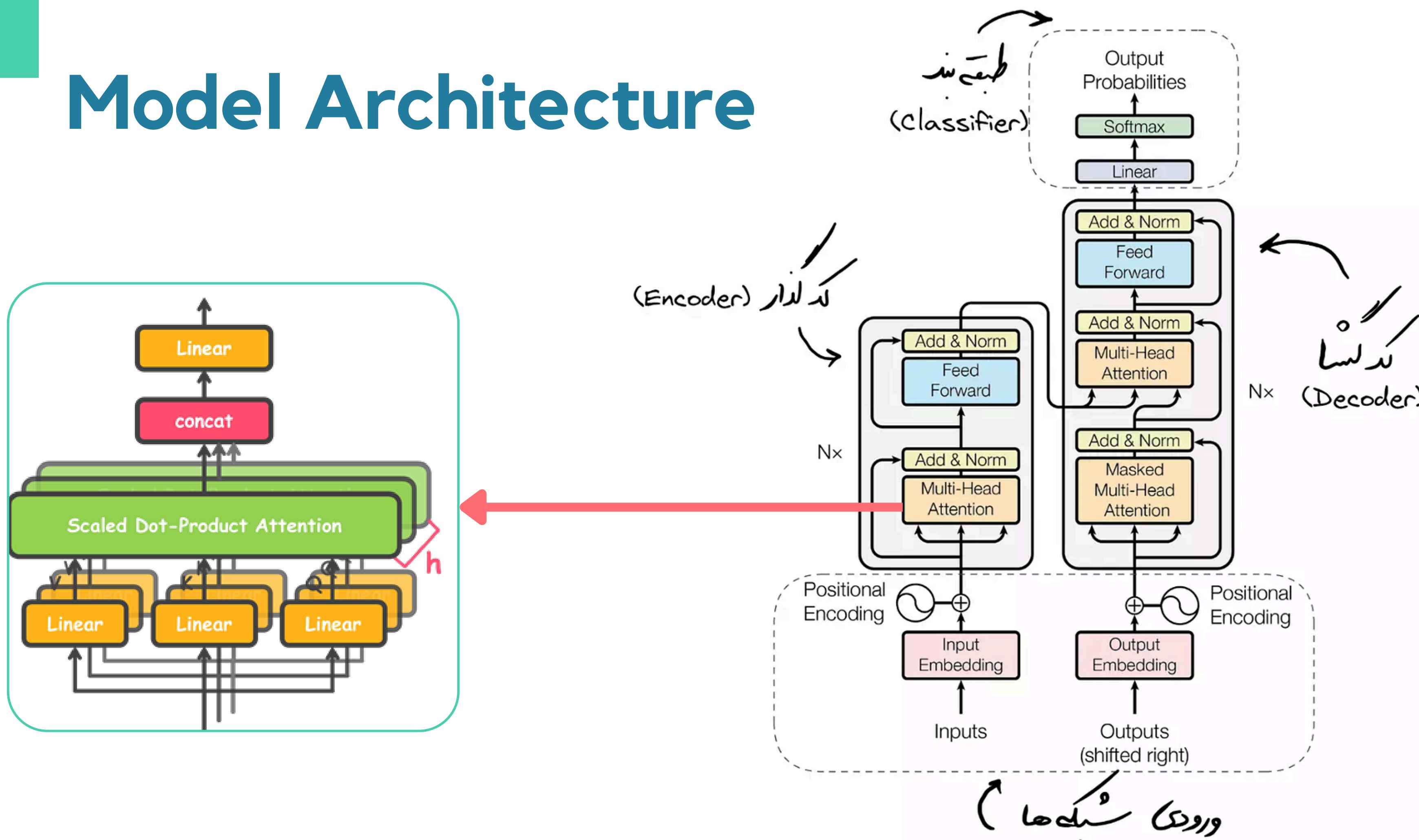
Attention Decoding Stage

Attention₄



4

Model Architecture



References

- <https://medium.com/@geetkal67/attention-networks-a-simple-way-to-understand-self-attention-f5fb363c736d>
- <https://blog.faradars.org/implement-gradient-descent-in-python/>
- <https://blog.faradars.org/implement-gradient-descent-in-python/>
- <https://www.analyticsvidhya.com/blog/2020/10/how-does-the-gradient-descent-algorithm-work-in-machine-learning/>
- <https://paperswithcode.com/paper/squeezebert-what-can-computer-vision-teach>
- <https://biasvariance.net/articles/868ss/softmax-behaviour/>

References

- <https://paperswithcode.com/method/dropout>
- <http://jmlr.org/papers/v15/srivastava14a.html>
- <https://howsam.org/transformer/>
- <https://virgool.io/@shenasa/%D9%85%D9%A9%D8%A7%D9%86%DB%8C%D8%B2%D9%85-%D8%AA%D9%88%D8%AC%D9%87-%D9%88-%D9%85%D8%AF%D9%84-%D9%87%D8%A7%DB%8C-%D8%AA%D8%A8%D8%AF%DB%8C%D9%84-%DA%A9%D9%86%D9%86%D8%AF%D9%87-transformers-zxboffthj9xi>
- <https://www.youtube.com/watch?v=dichlcUZfOw>



THANK YOU

mahdirahmani8@yahoo.com