



香港科技大學
THE HONG KONG
UNIVERSITY OF SCIENCE
AND TECHNOLOGY

COMP 4901B

Large Language Models

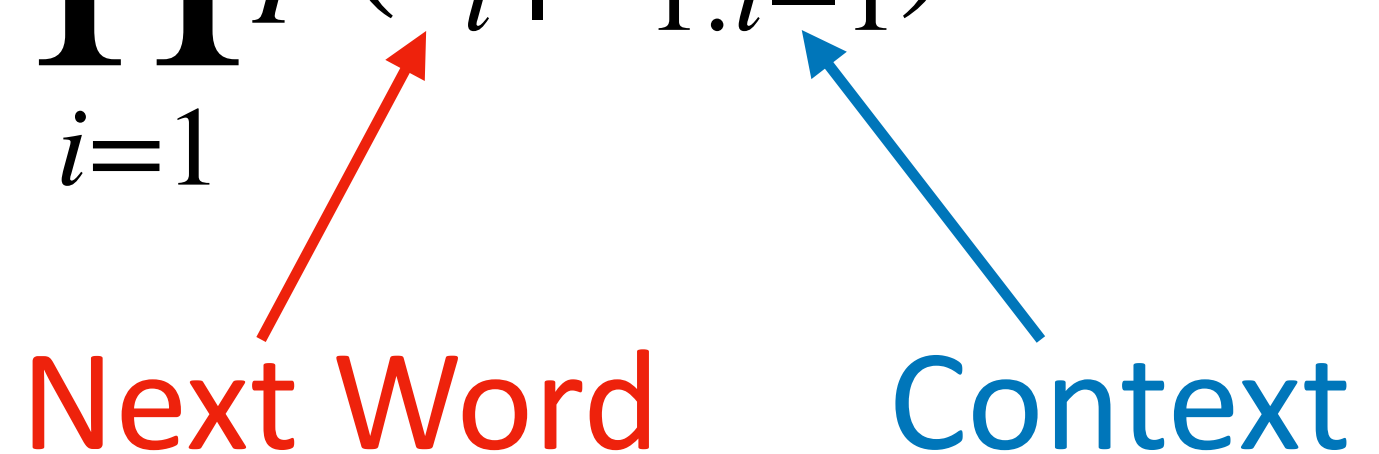
Recurrent Neural Networks, Transformers

Junxian He

Sep 12, 2025

Recap: Autoregressive Language Models

$$\begin{aligned} p(\text{the, mouse, ate, the, cheese}) &= p(\text{the}) \\ &\quad p(\text{mouse} \mid \text{the}) \\ &\quad p(\text{ate} \mid \text{the, mouse}) \\ &\quad p(\text{the} \mid \text{the, mouse, ate}) \\ &\quad p(\text{cheese} \mid \text{the, mouse, ate, the}). \end{aligned}$$

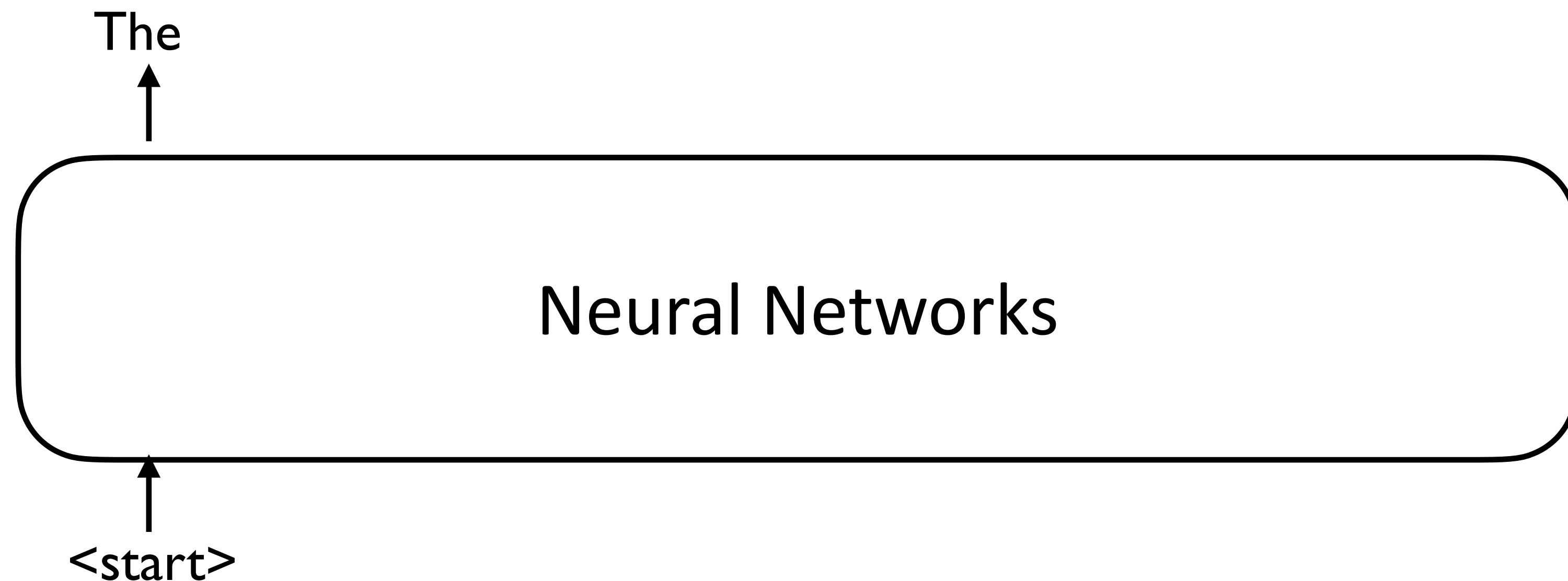
$$p(x_1, x_2, \dots, x_I) = \prod_{i=1}^I p(x_i \mid x_{1:i-1})$$


The diagram illustrates the components of the autoregressive model equation. A red arrow points from the text "Next Word" to the variable x_i in the probability term $p(x_i \mid x_{1:i-1})$. A blue arrow points from the text "Context" to the sequence $x_{1:i-1}$ in the same term.

Recap: Neural Language Models

Neural language models are typically autoregressive

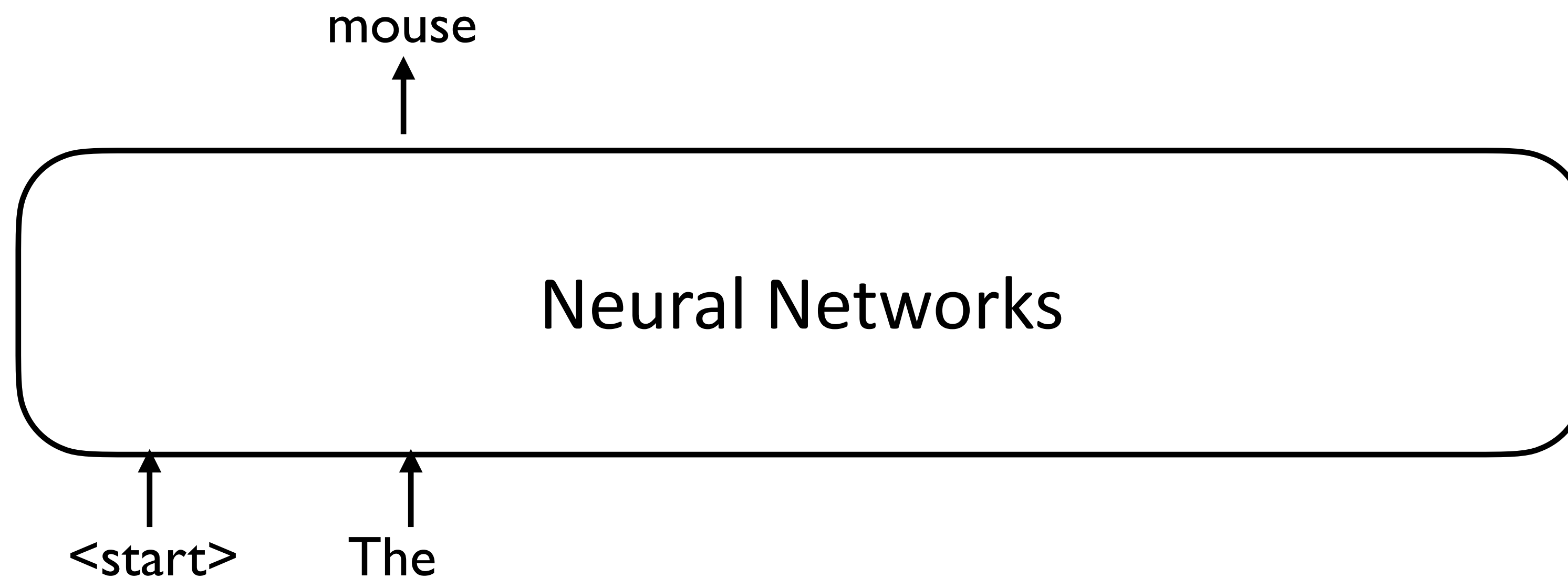
Data: “The mouse ate the cheese .”



Recap: Neural Language Models

Neural language models are typically autoregressive

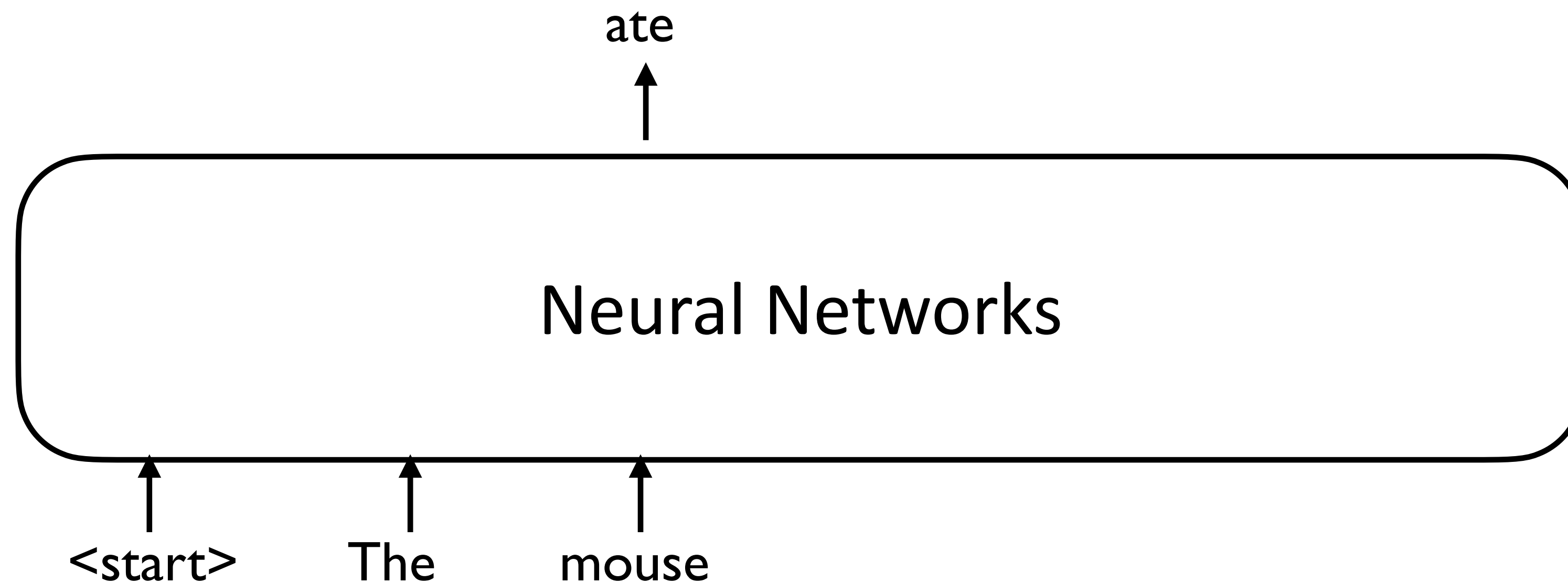
Data: “The mouse ate the cheese .”



Recap: Neural Language Models

Neural language models are typically autoregressive

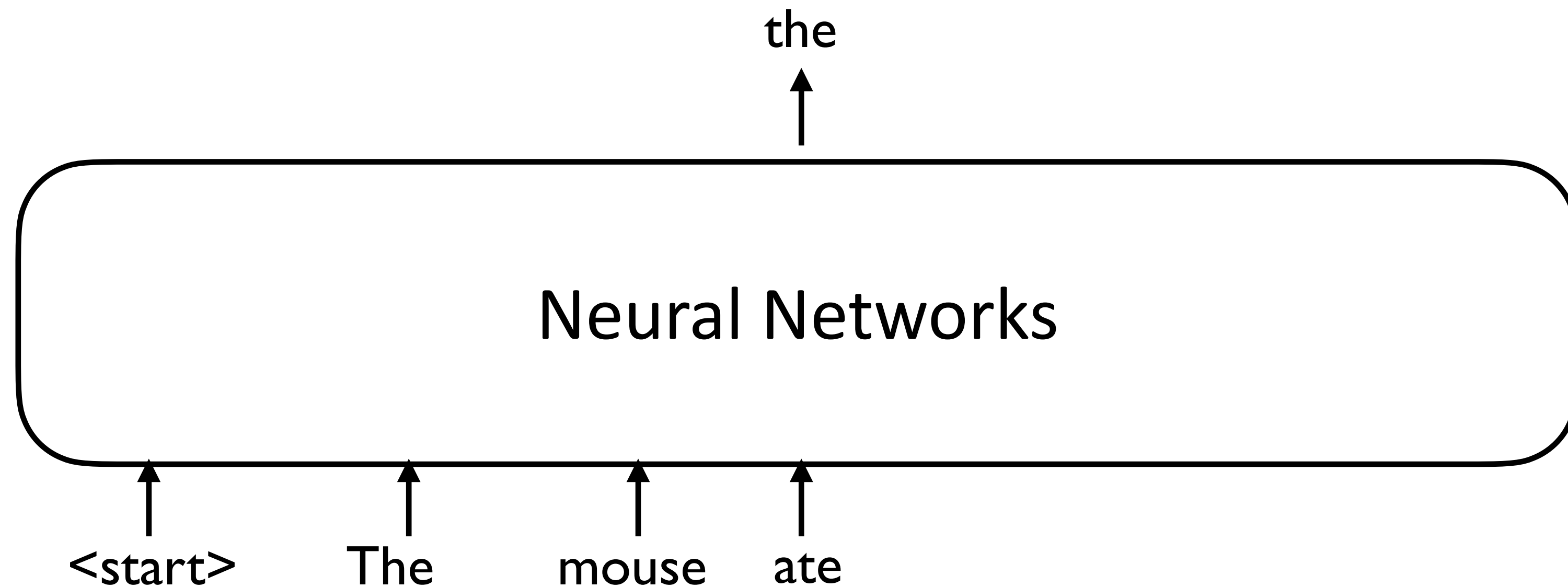
Data: “The mouse ate the cheese .”



Recap: Neural Language Models

Neural language models are typically autoregressive

Data: “The mouse ate the cheese .”

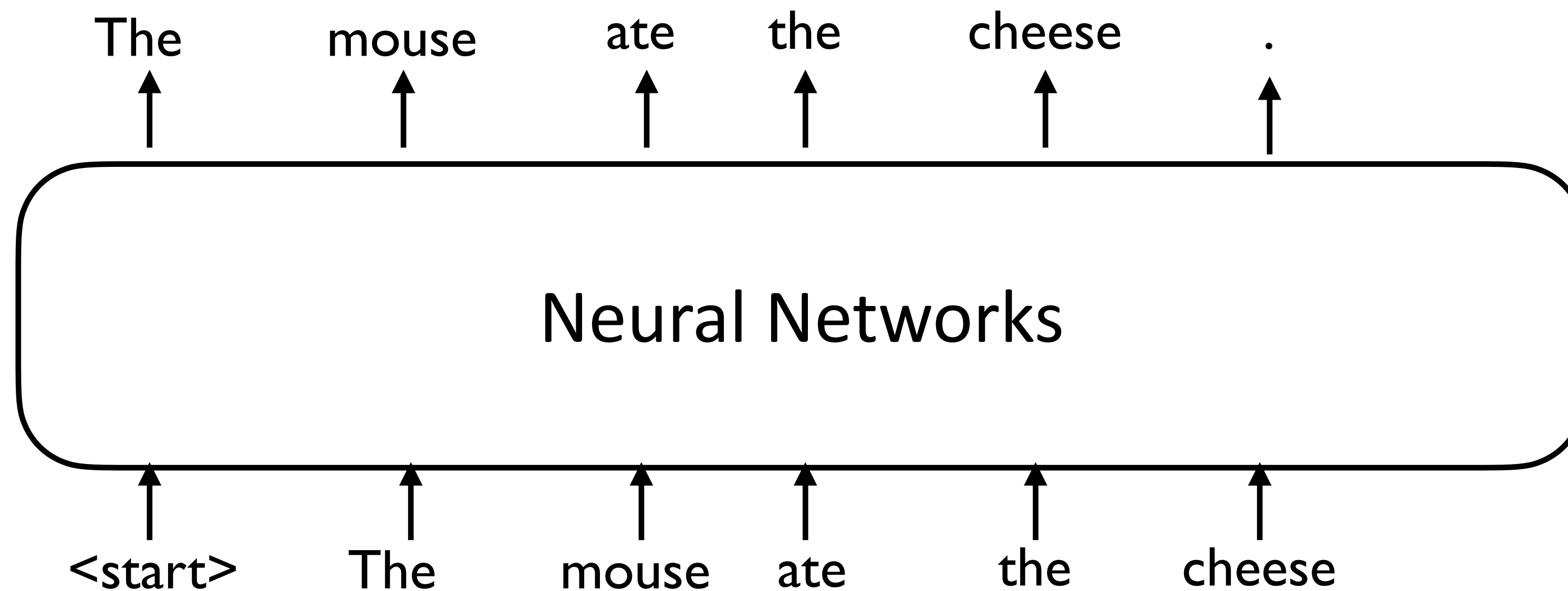


We can compute the loss on every token in parallel

Recap: Neural Language Models

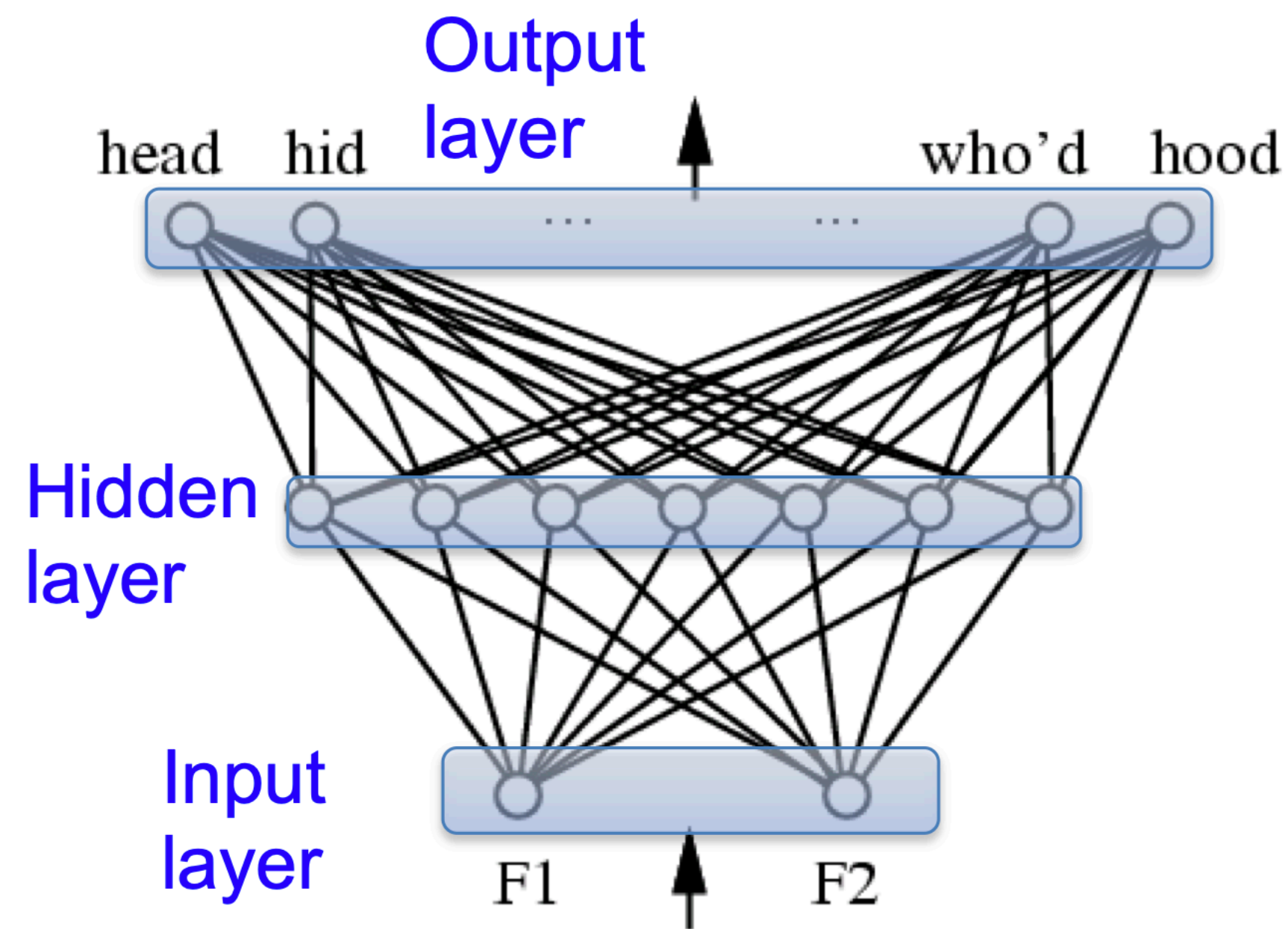
Neural language models are typically autoregressive

Data: “The mouse ate the cheese .”

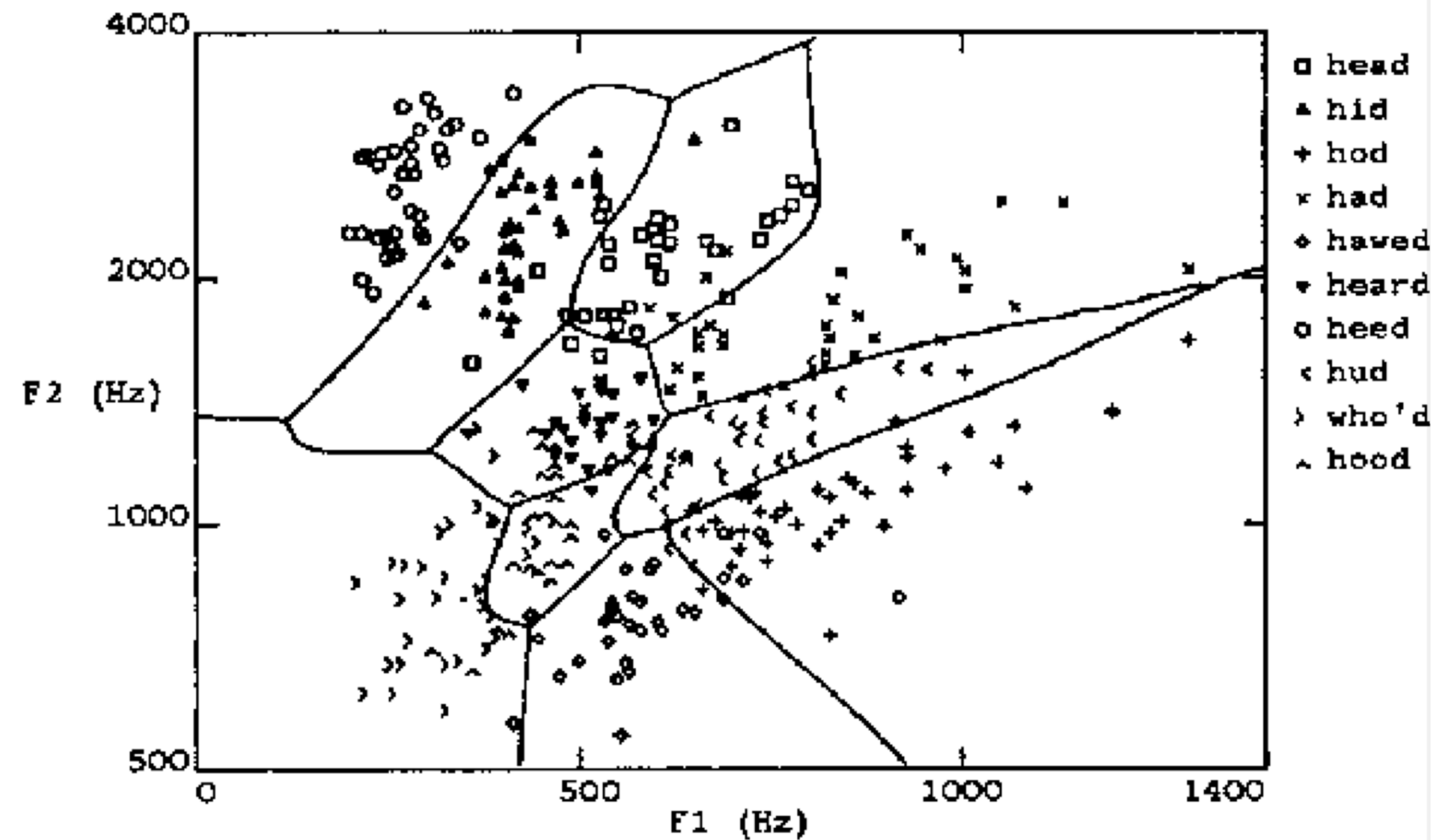


Each prediction only sees the inputs on its left

Recap: Multilayer Networks of Sigmoid Units



Two layers of logistic units

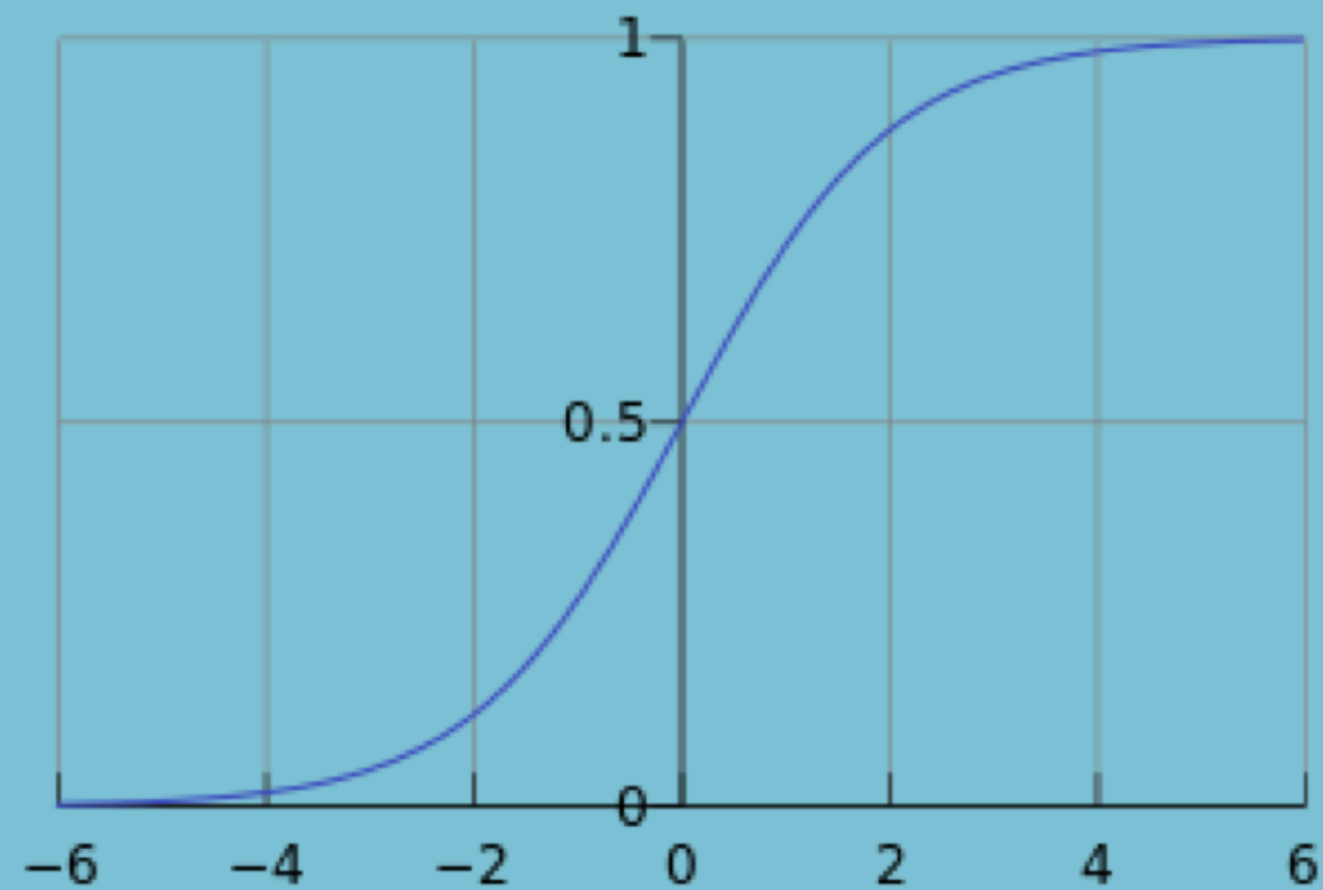


Highly non-linear decision surface

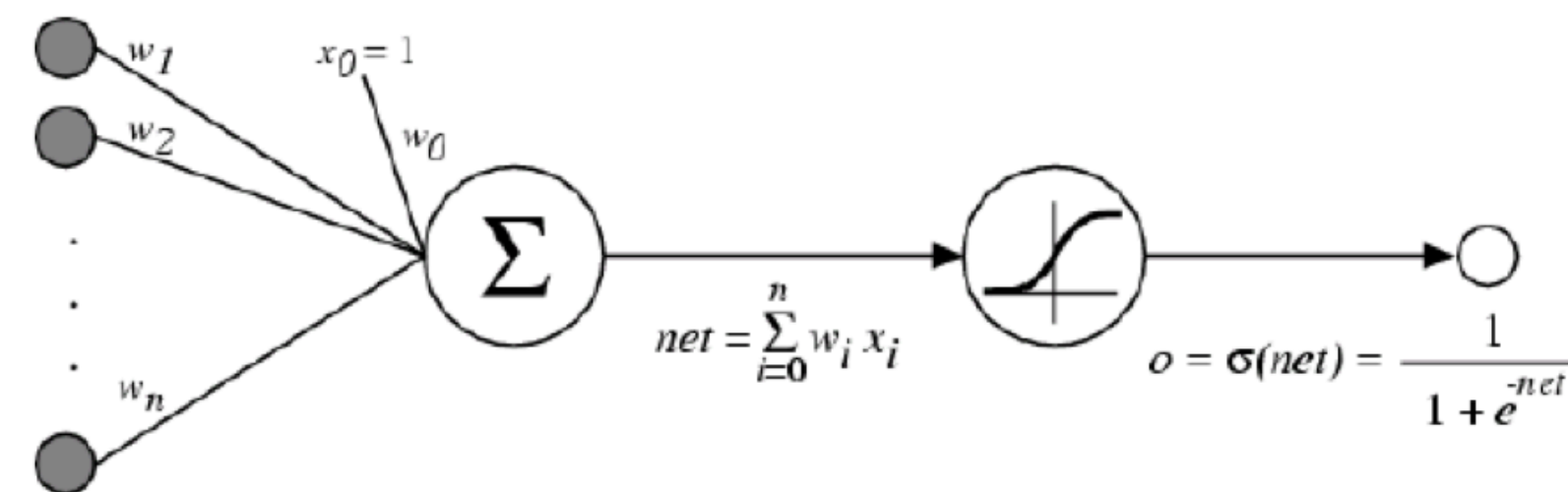
Activation Functions

Sigmoid / Logistic Function

$$\text{logistic}(u) \equiv \frac{1}{1 + e^{-u}}$$

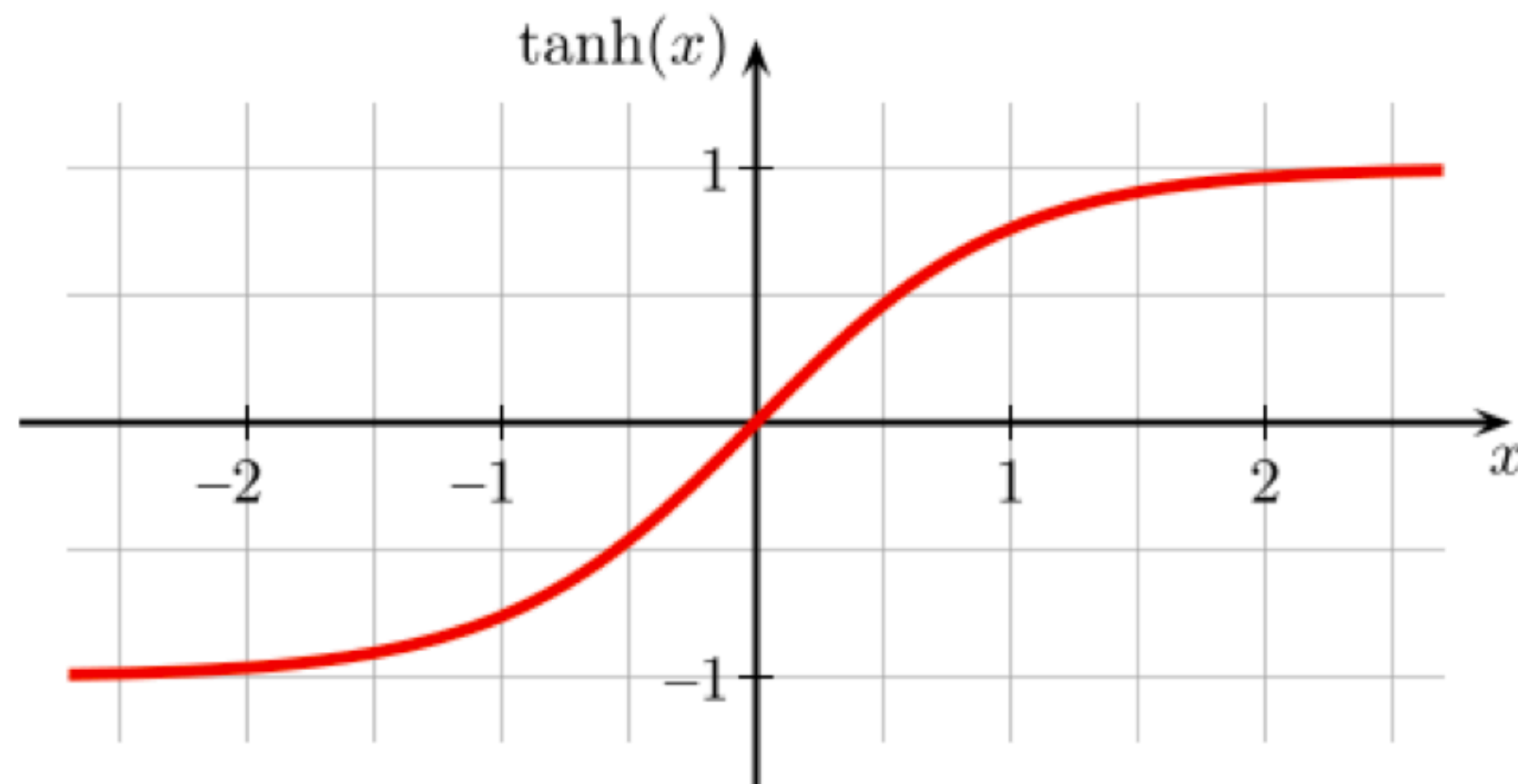


So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function...



Tanh

- A new change: modifying the nonlinearity
 - The logistic is not widely used in modern ANNs



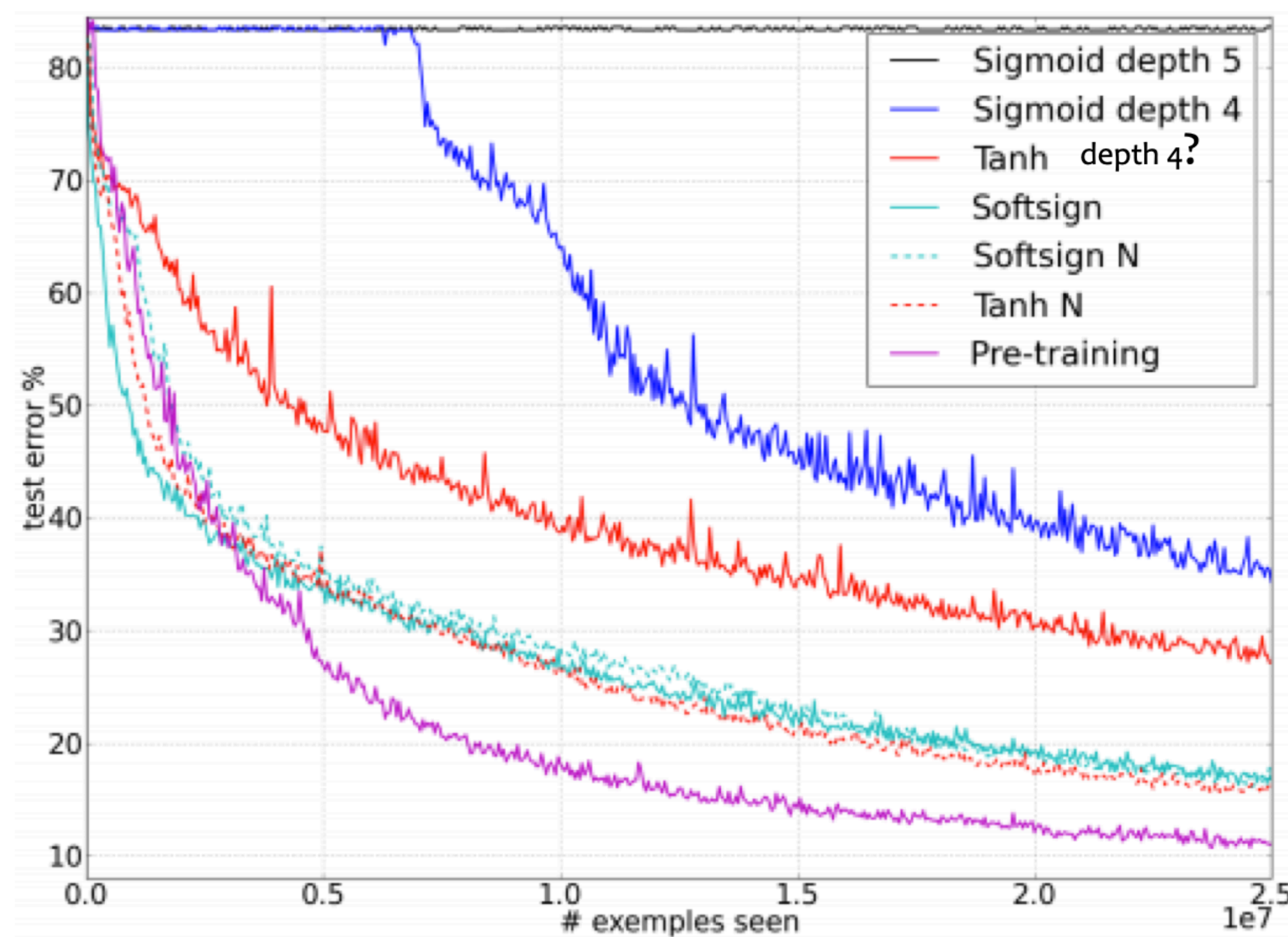
Alternate 1:
tanh

Like logistic function but
shifted to range $[-1, +1]$

Activation Function

Understanding the difficulty of training deep feedforward neural networks

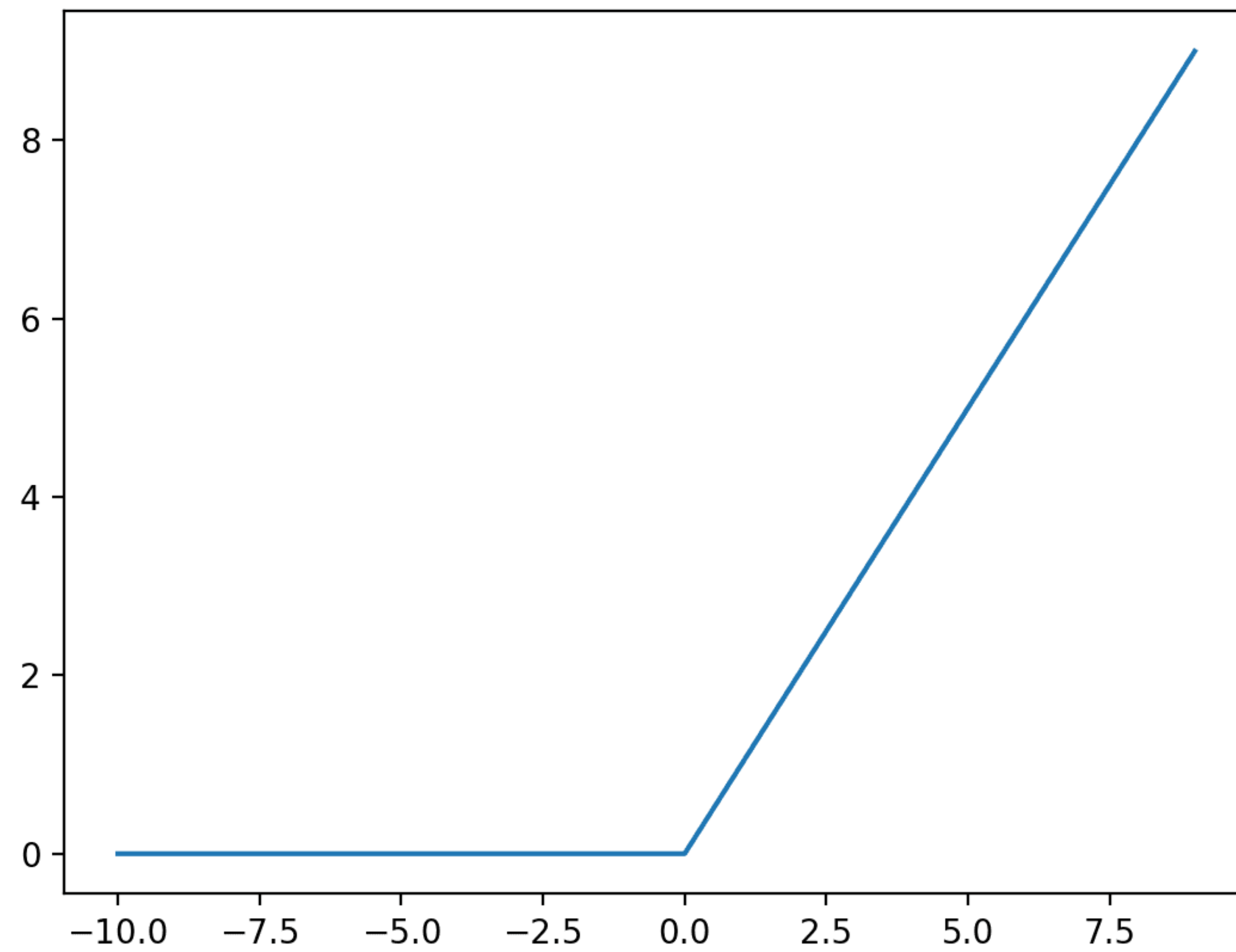
AI Stats 2010



} sigmoid
vs.
tanh

Figure from Glorot & Bentio (2010)

ReLU



Other Activation Functions

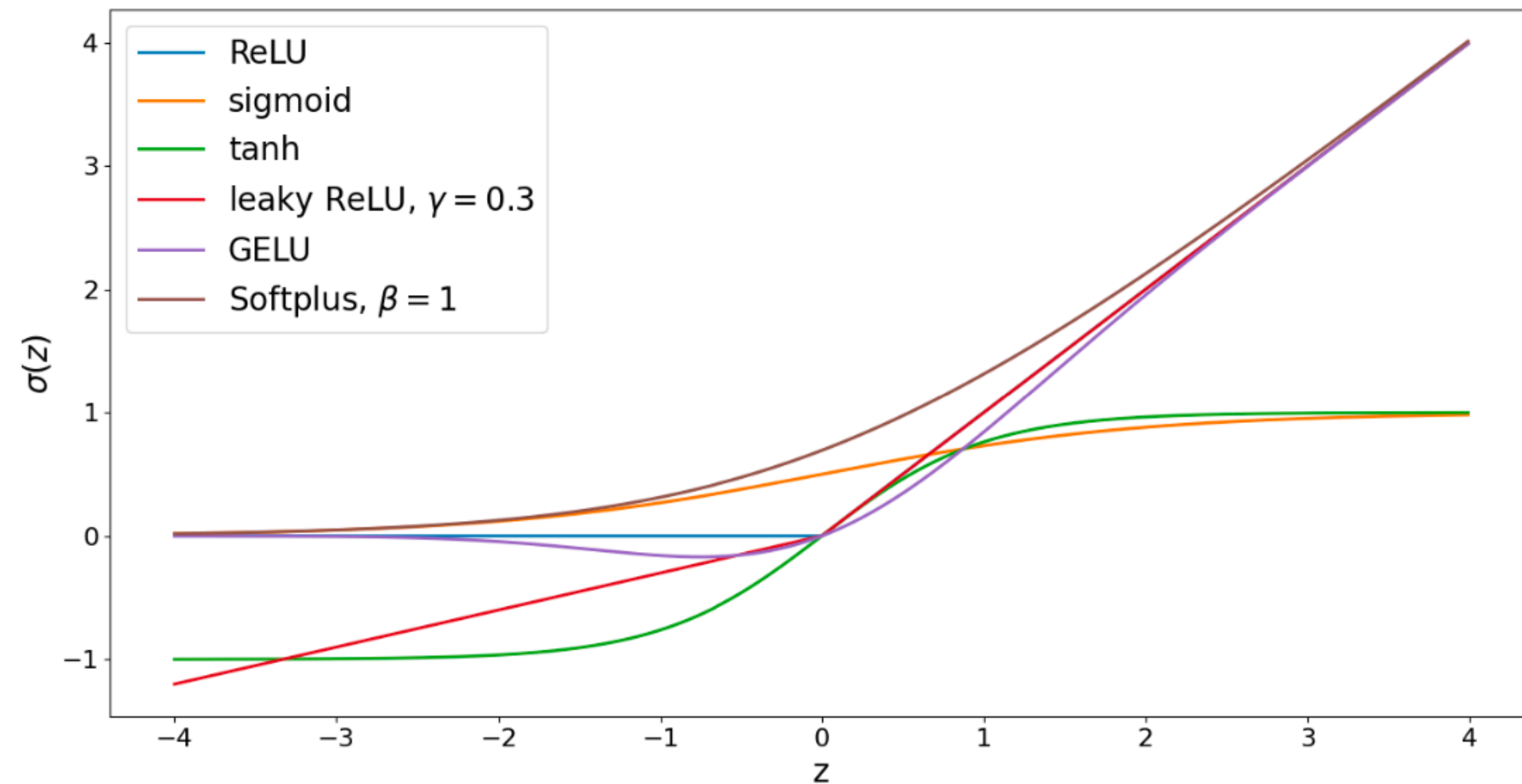
$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (\text{sigmoid})$$

$$\sigma(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad (\text{tanh})$$

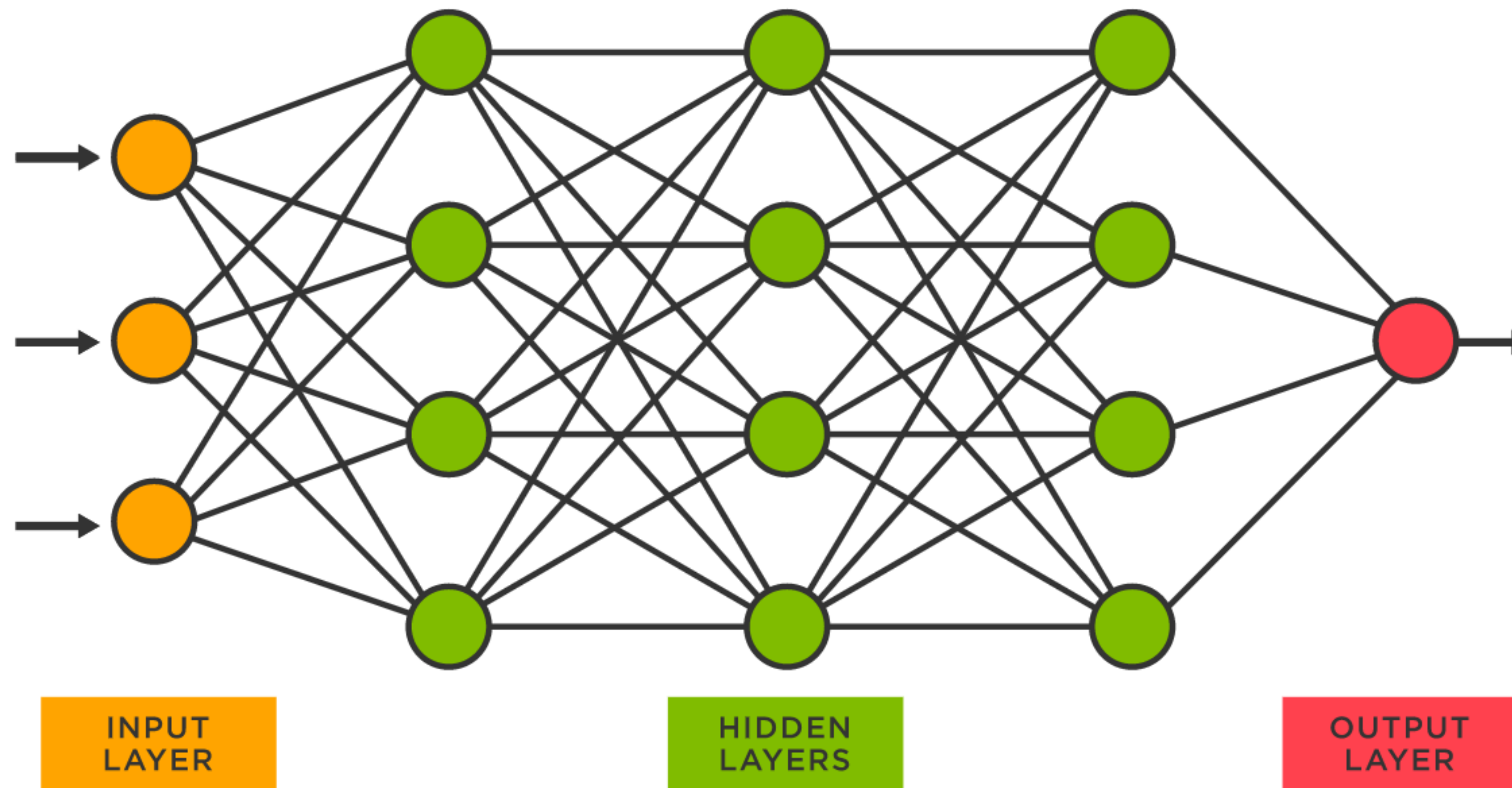
$$\sigma(z) = \max\{z, \gamma z\}, \gamma \in (0, 1) \quad (\text{leaky ReLU})$$

$$\sigma(z) = \frac{z}{2} \left[1 + \operatorname{erf}\left(\frac{z}{\sqrt{2}}\right) \right] \quad (\text{GELU})$$

$$\sigma(z) = \frac{1}{\beta} \log(1 + \exp(\beta z)), \beta > 0 \quad (\text{Softplus})$$



Multilayer Perceptron Neural Networks (MLP)



Residual Connection

We want deeper and deeper NNs, but going deep is difficult

- Troubles accumulate w/ more layers
- Signals get distorted when propagated
- in forward and backward passes

Commonly used techniques to train “Deep” NNs:

Weight initialization

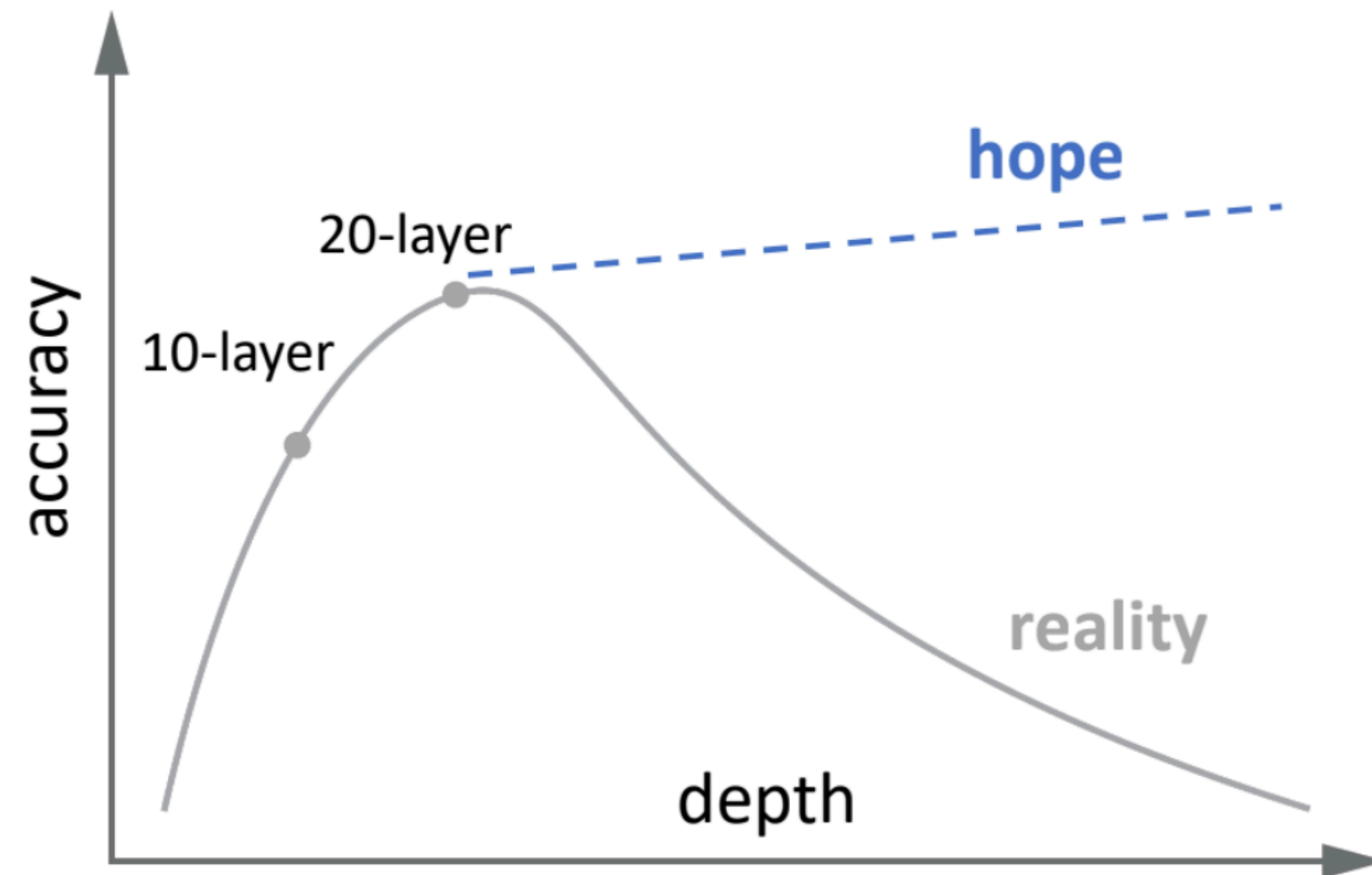
Normalization modules

Deep residual learning



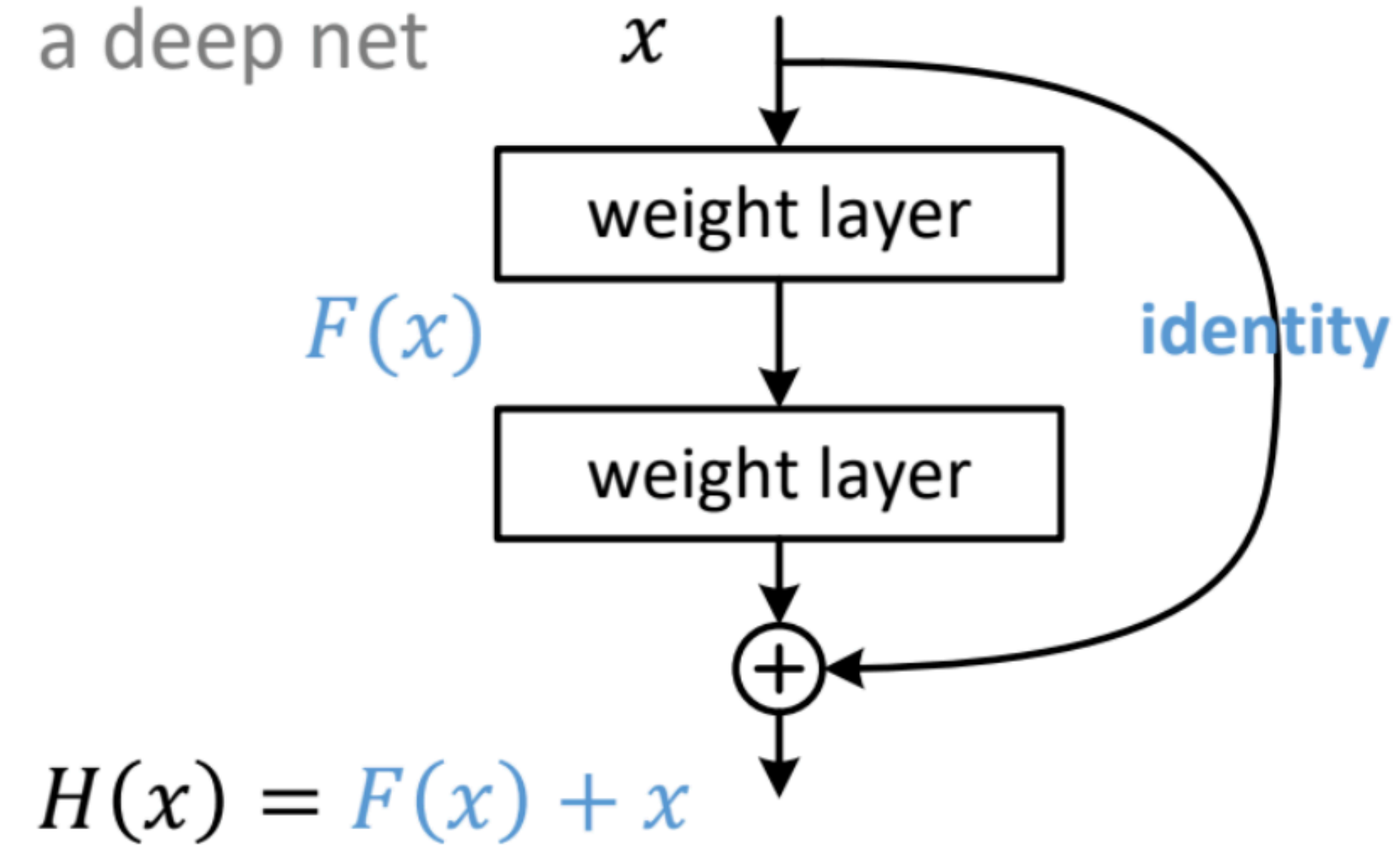
The Degradation Problem

- Good init + norm enable training deeper models
- Simply stacking more layers?
- Degrade after ~20 layers
- Not overfitting
- Difficult to train



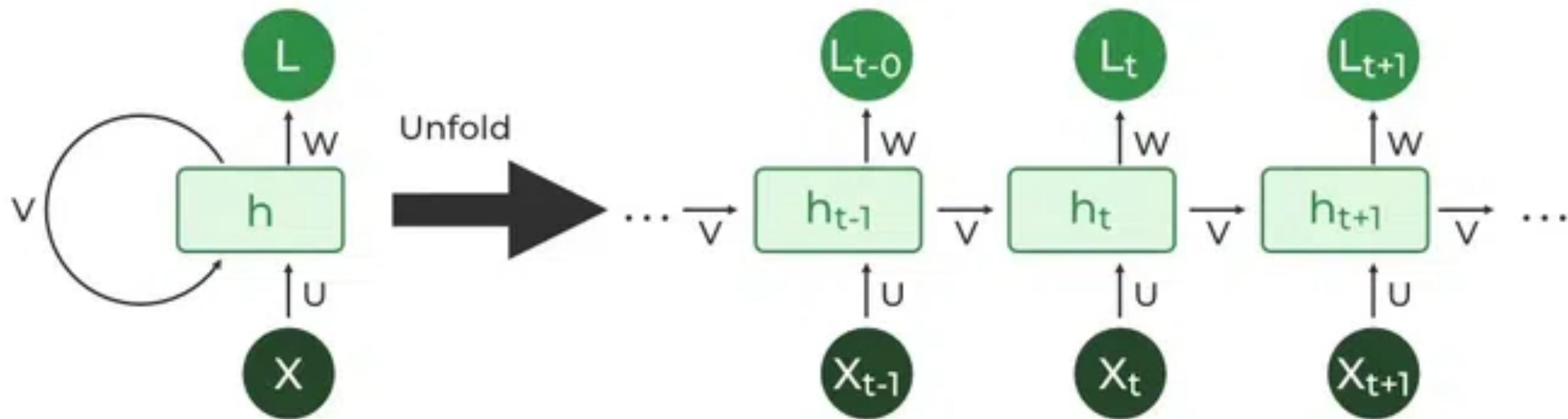
Deep Residual Learning

a subnet in
a deep net



**MLP network is hard to handle
sequence data with varying length**

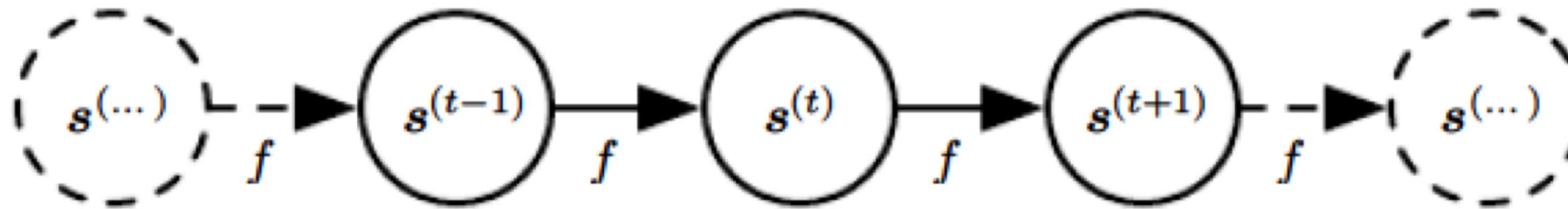
Recurrent Neural Networks (RNNs)



Recurrent Neural Networks

- Dates back to (Rumelhart *et al.*, 1986)
- A family of neural networks for handling sequential data, which involves variable length inputs or outputs
- Especially, for natural language processing (NLP)

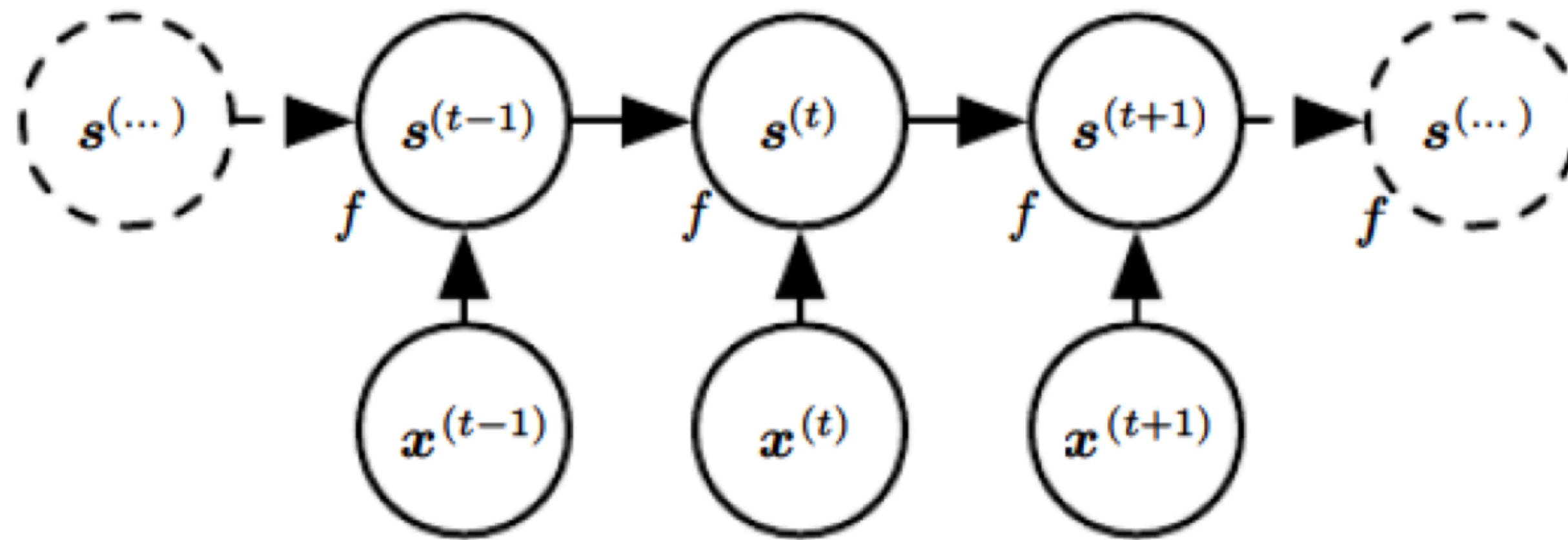
Computation Graph



$$s^{(t+1)} = f(s^{(t)}; \theta)$$

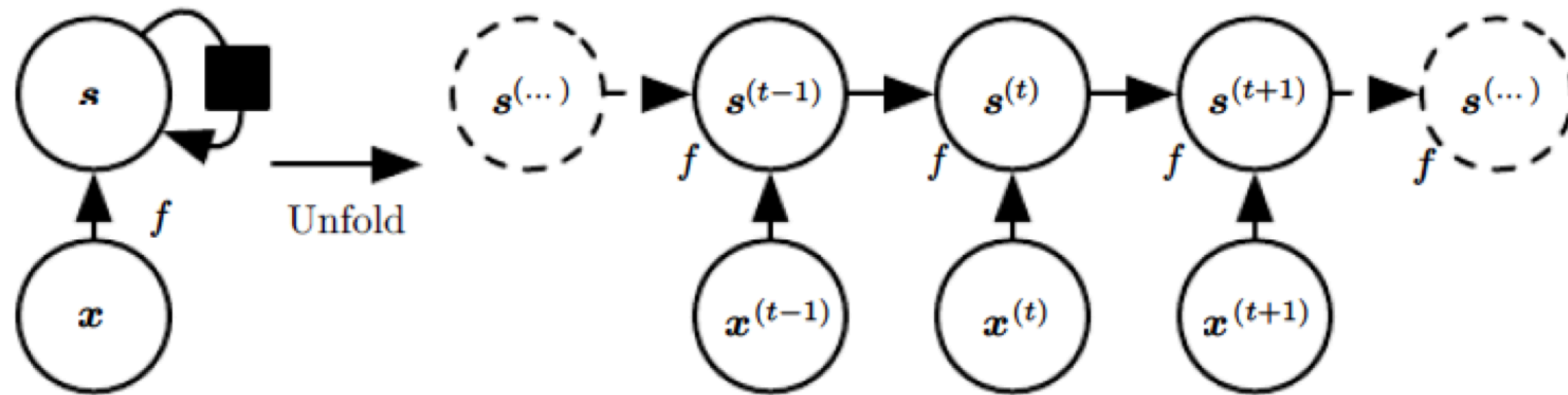
Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

Computation Graph



$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

Compact view



$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

Key: the same f and θ
for all time steps

Recurrent Neural Networks

- Use **the same** computational function and parameters across different time steps of the sequence
- Each time step: takes the input entry **and the previous hidden state** to compute the output entry
- Loss: typically computed every time step

Recurrent Neural Networks

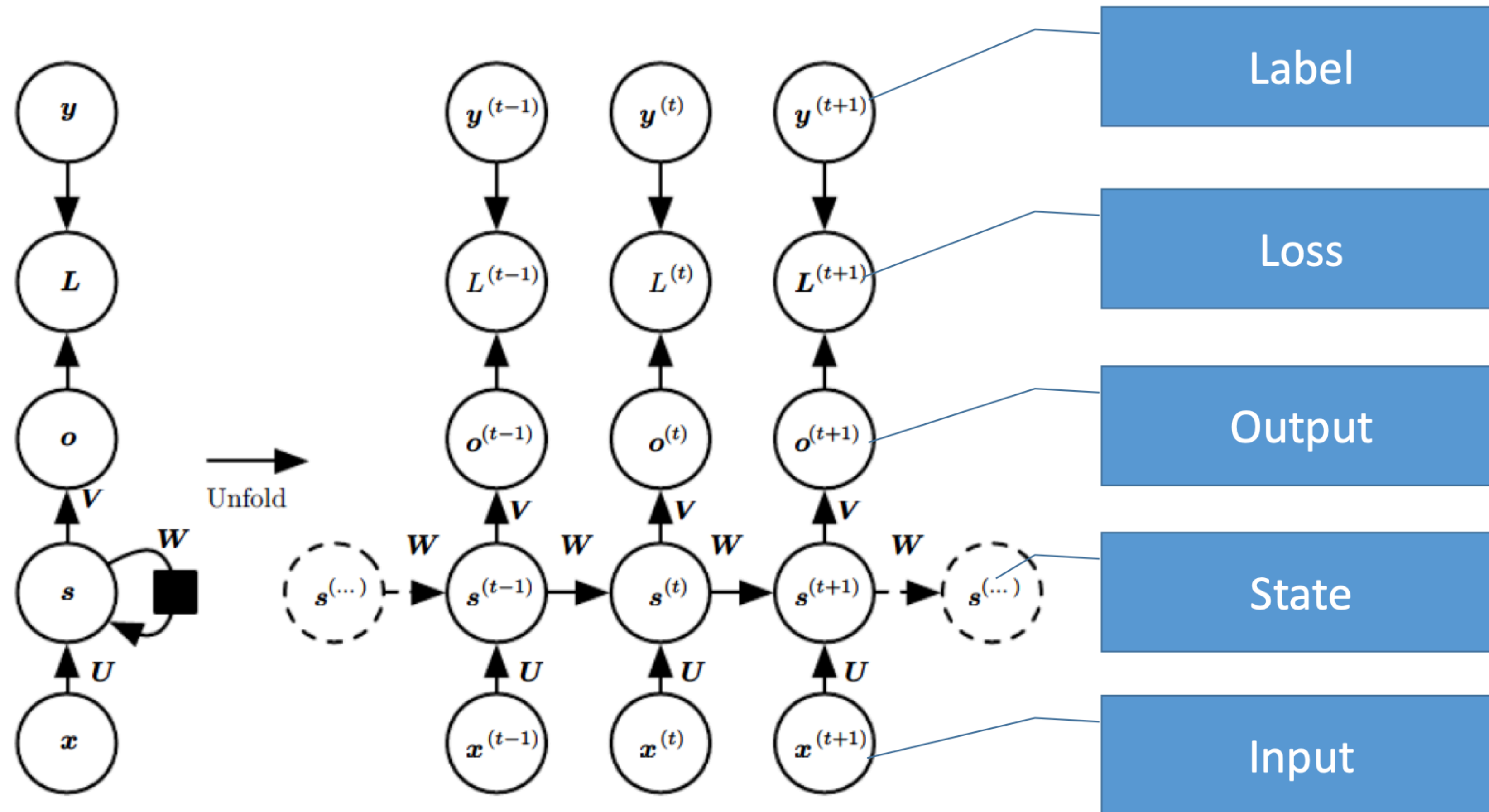


Figure from *Deep Learning*, by Goodfellow, Bengio and Courville

Recurrent Neural Networks

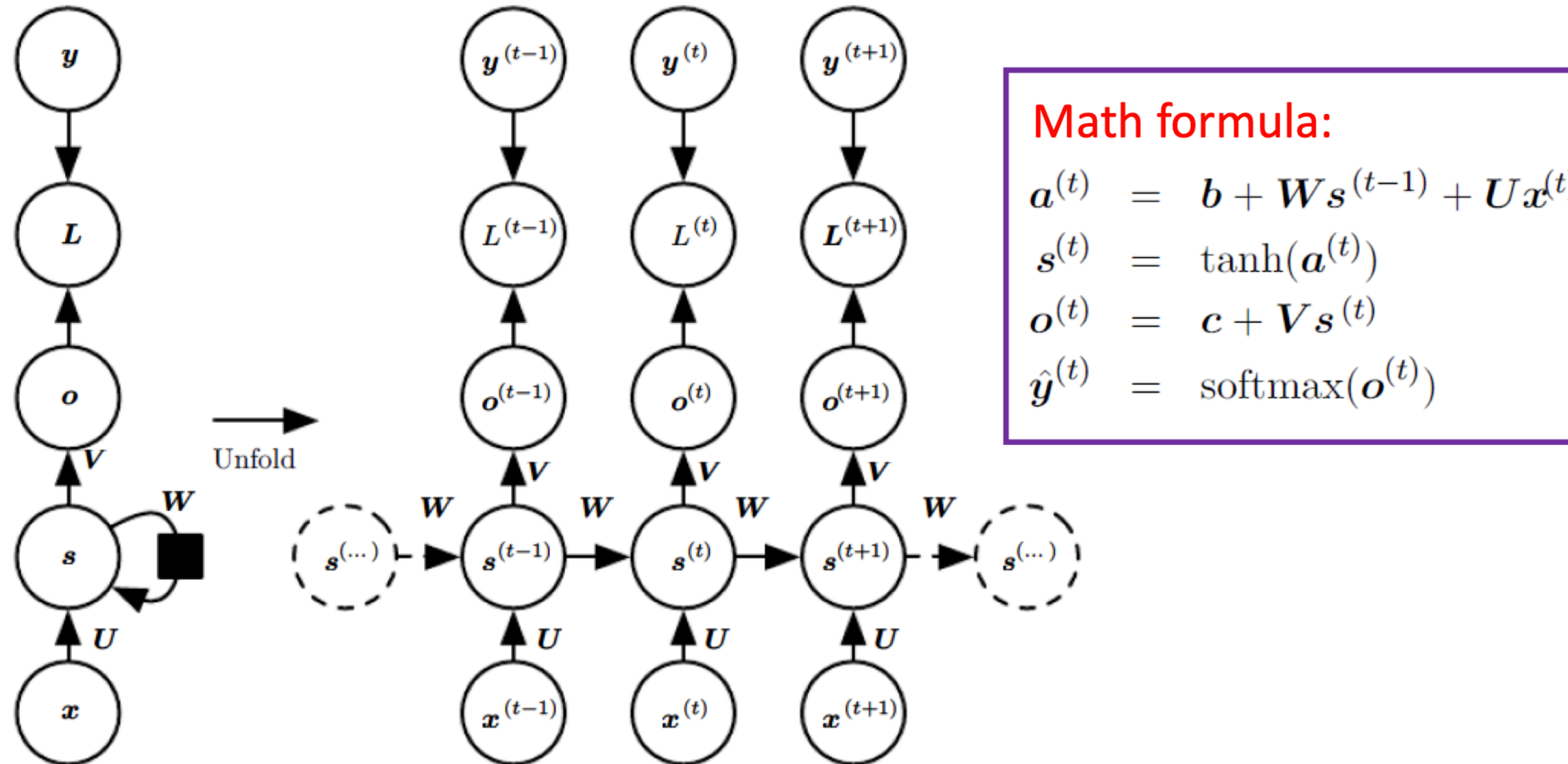
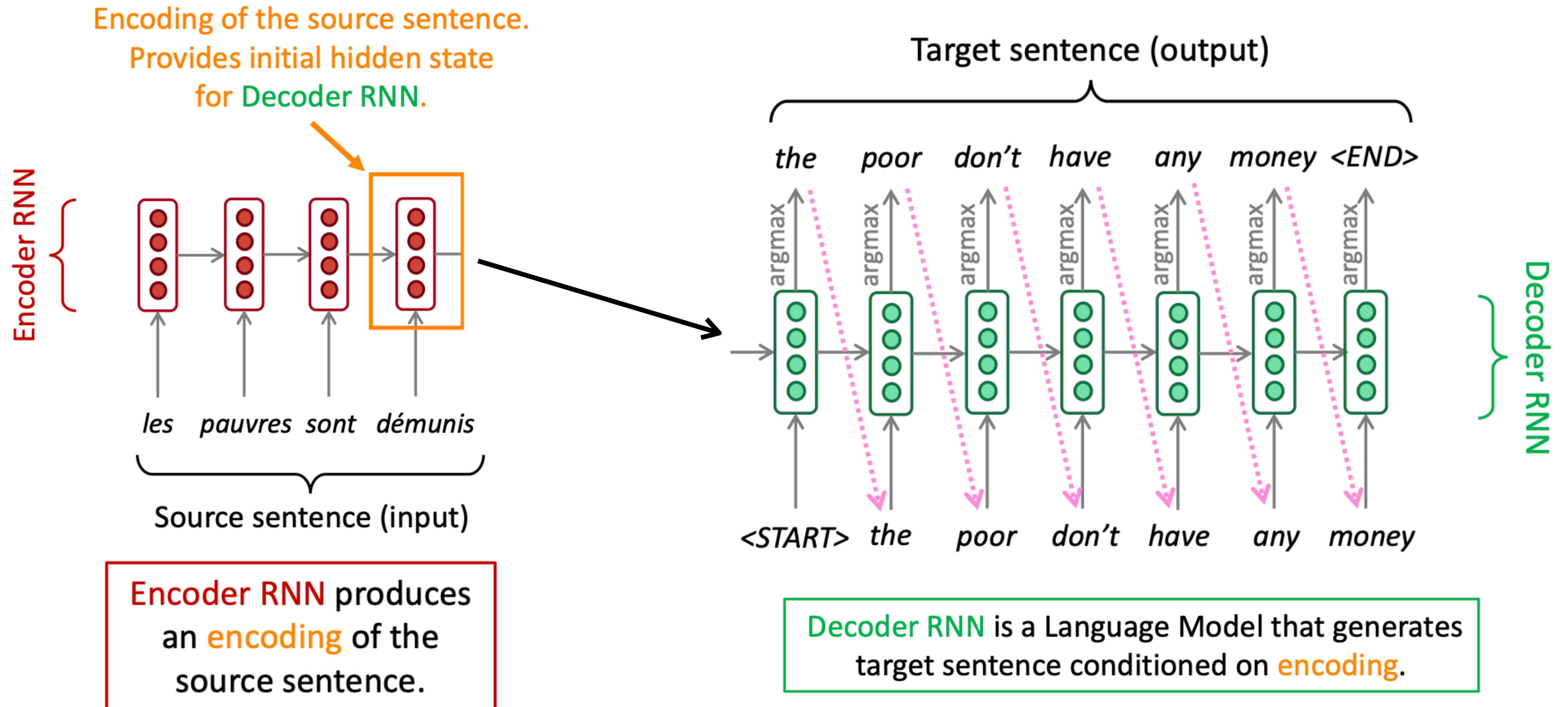


Figure from *Deep Learning*,
Goodfellow, Bengio and Courville

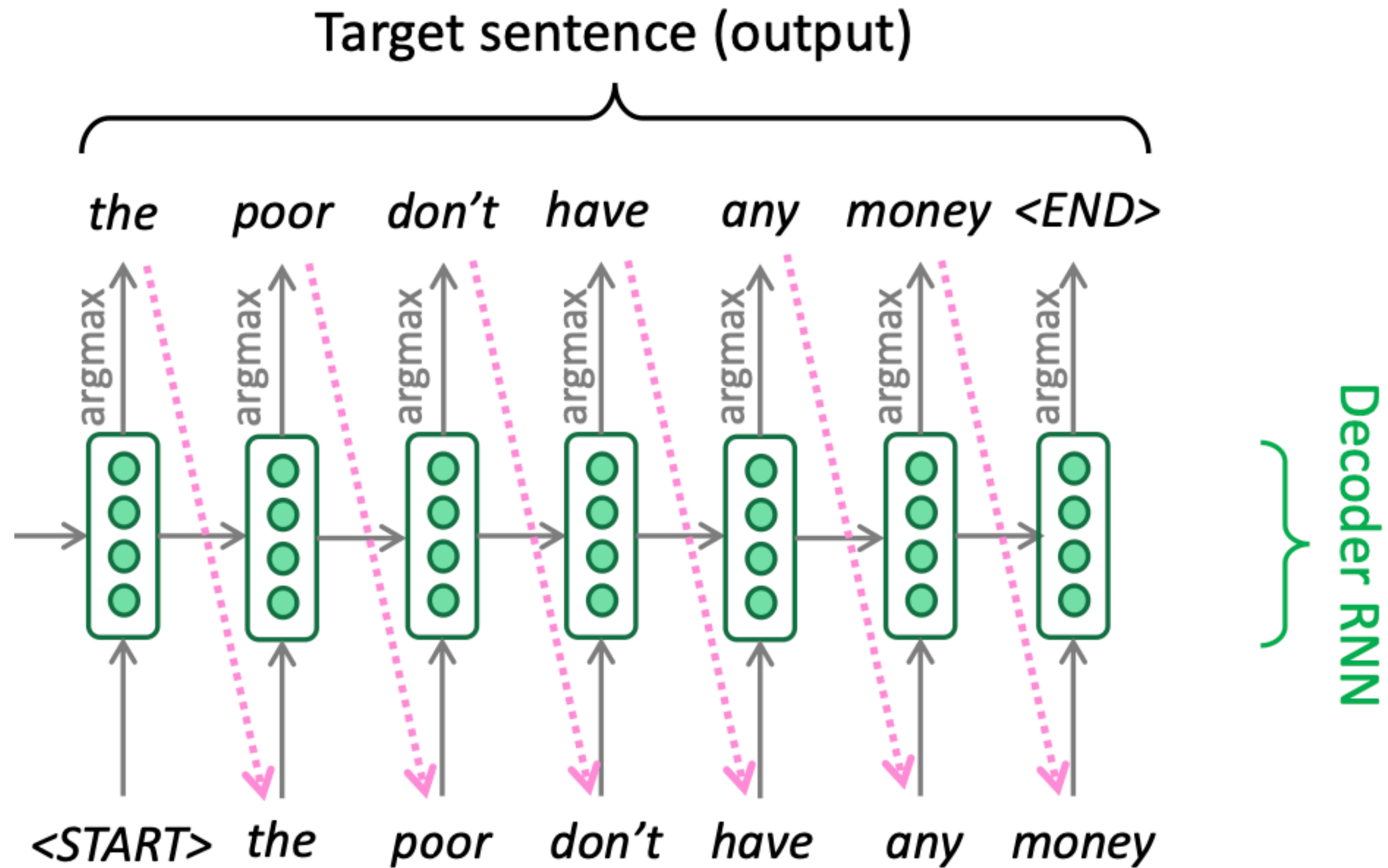
There are many variants of RNNs since the functional form to compute $s^{(t)}$ can vary, e.g., LSTM

Sequence-to-Sequence Learning

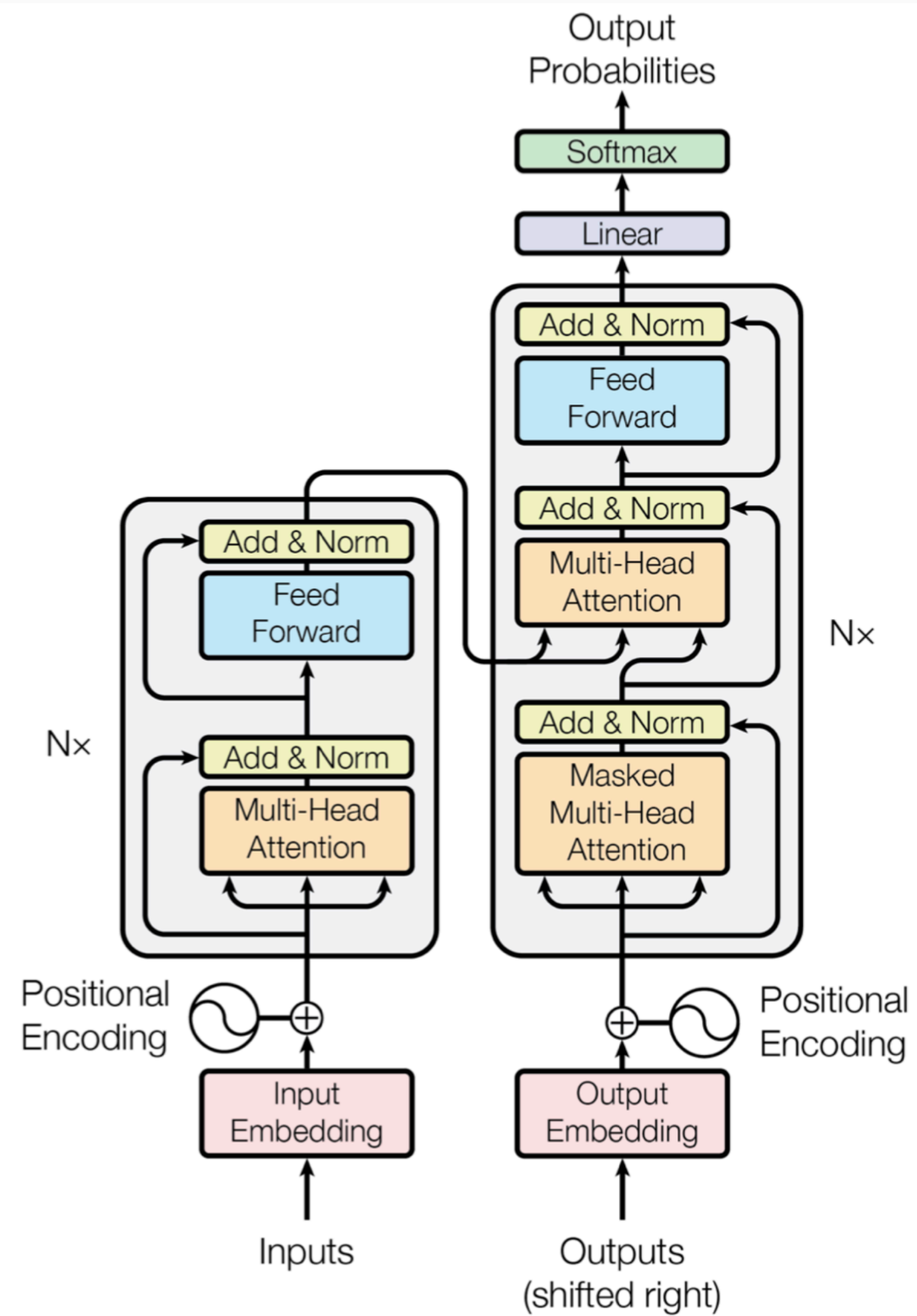
Example of Neural Machine Translation



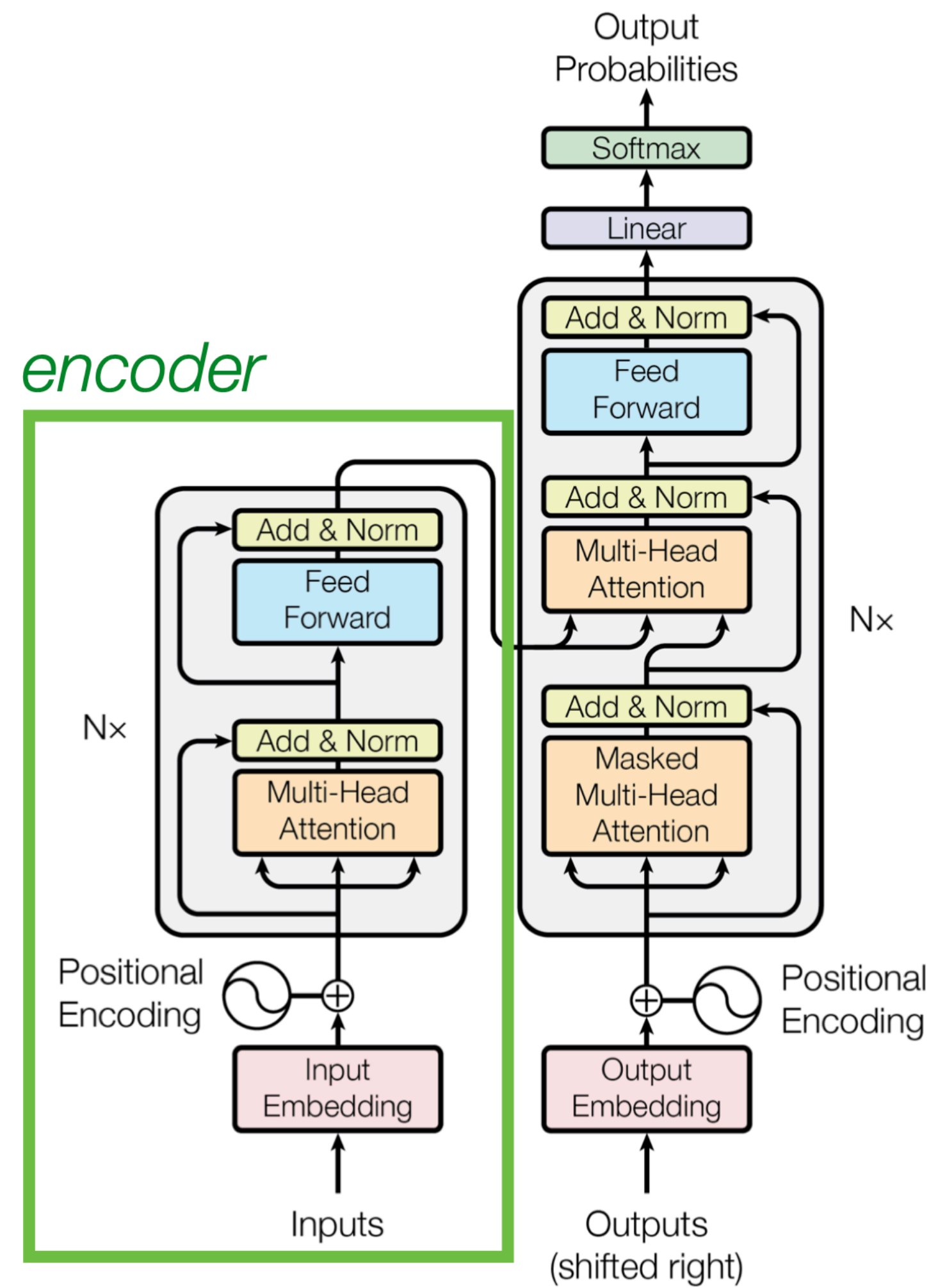
RNN Language Model



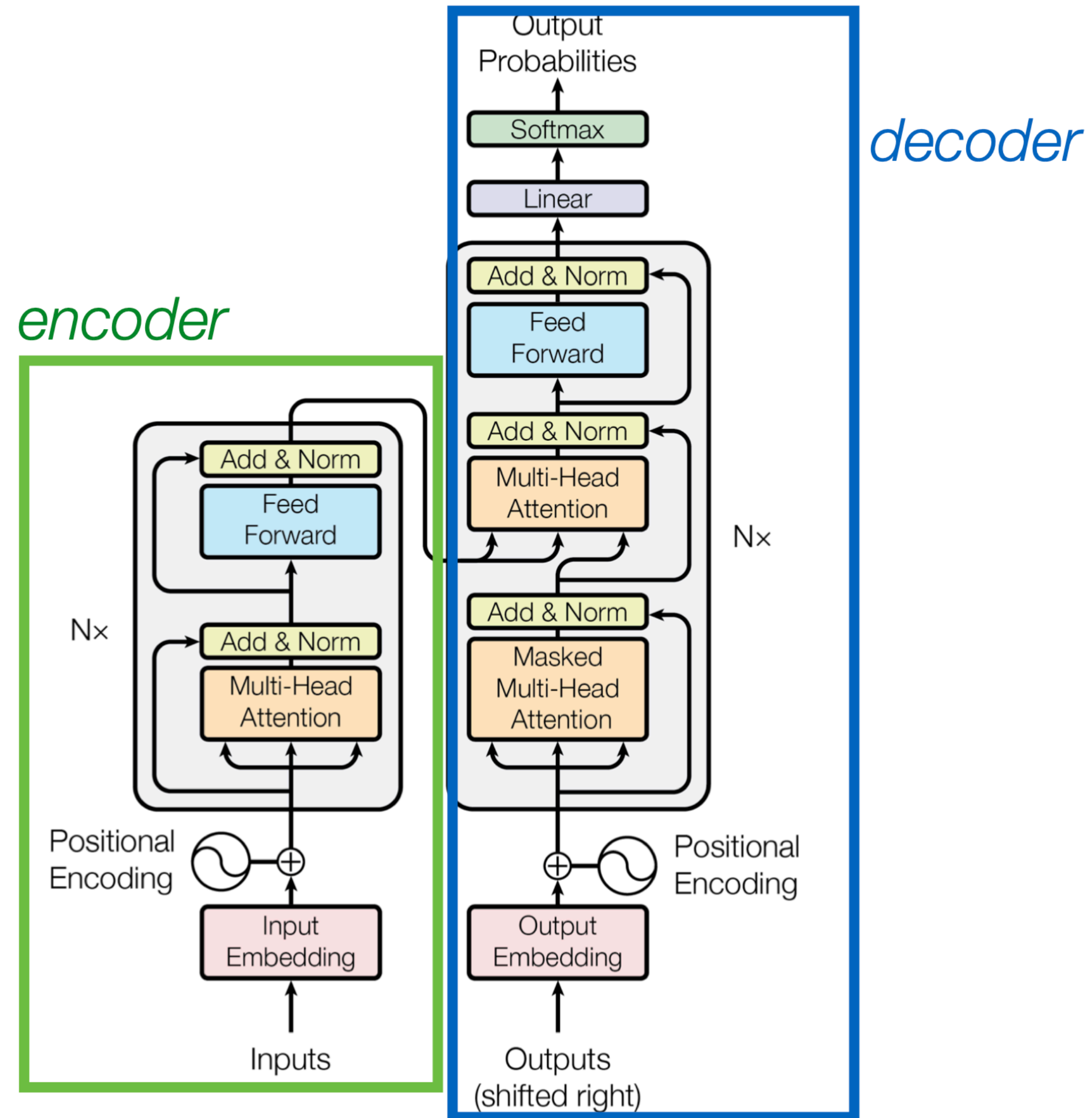
Transformer



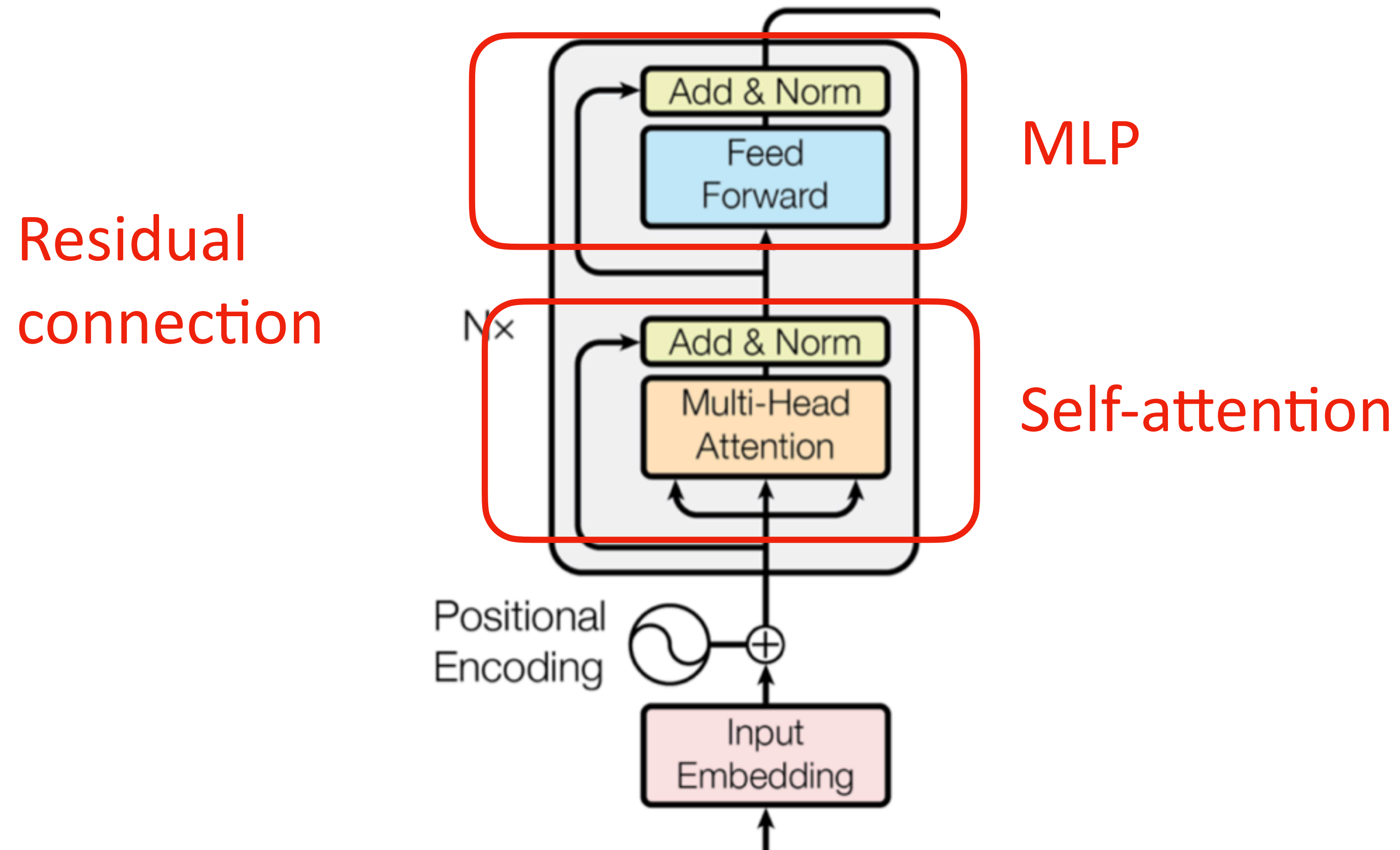
Encoder



Decoder



Transformer Encoder

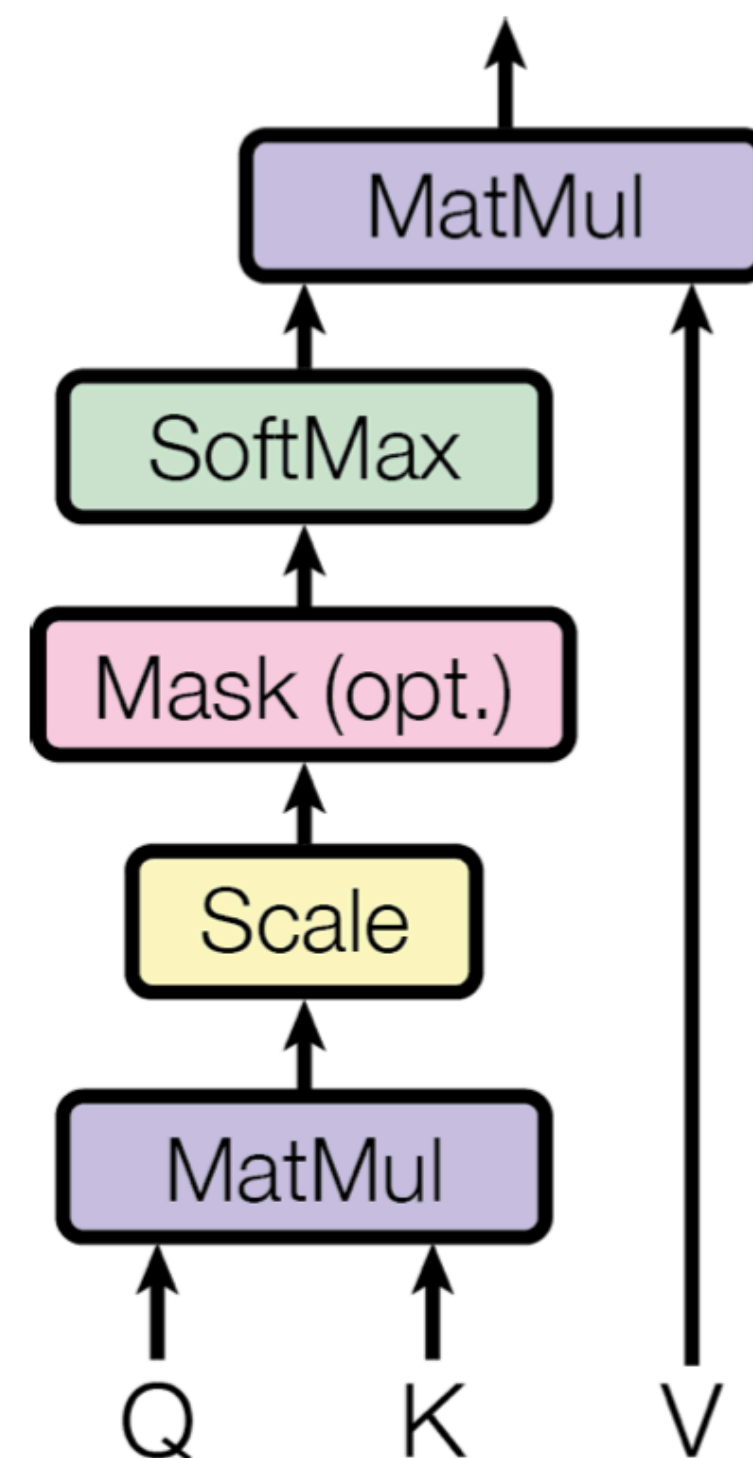


What is Attention

$$Q \in R^{n \times d} \quad K \in R^{m \times d} \quad V \in R^{m \times d}$$

We have n queries, m (key, value) pairs

Scaled Dot-Product Attention



Q: Query
K: key
V: value

$$\text{Attention weight} = \text{softmax}(QK^T)$$

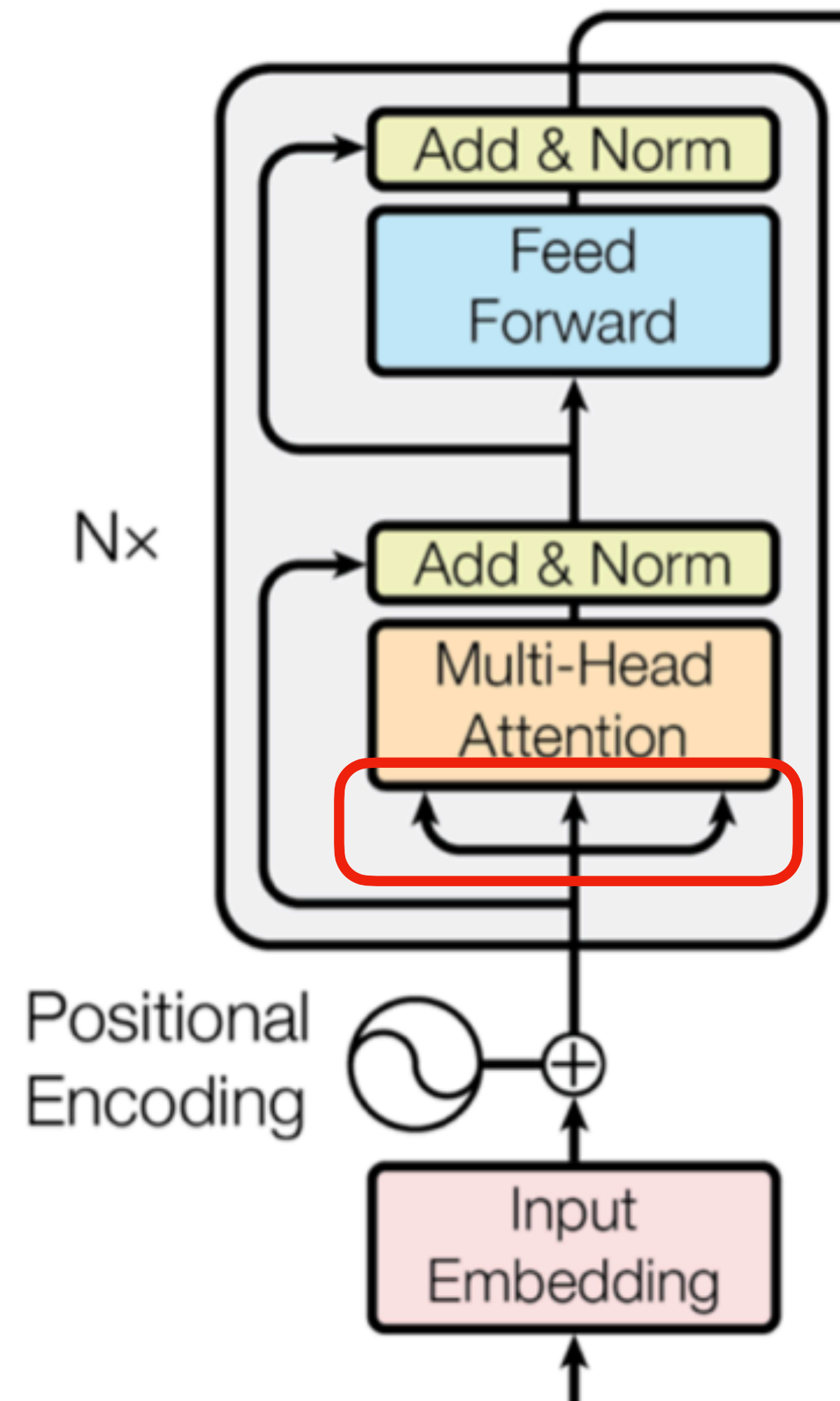
Dot-products grow large in magnitude

$$\text{Scaled Attention weight} = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \quad \text{Shape is } m \times n$$

Attention weight represents the strength to “attend” values V

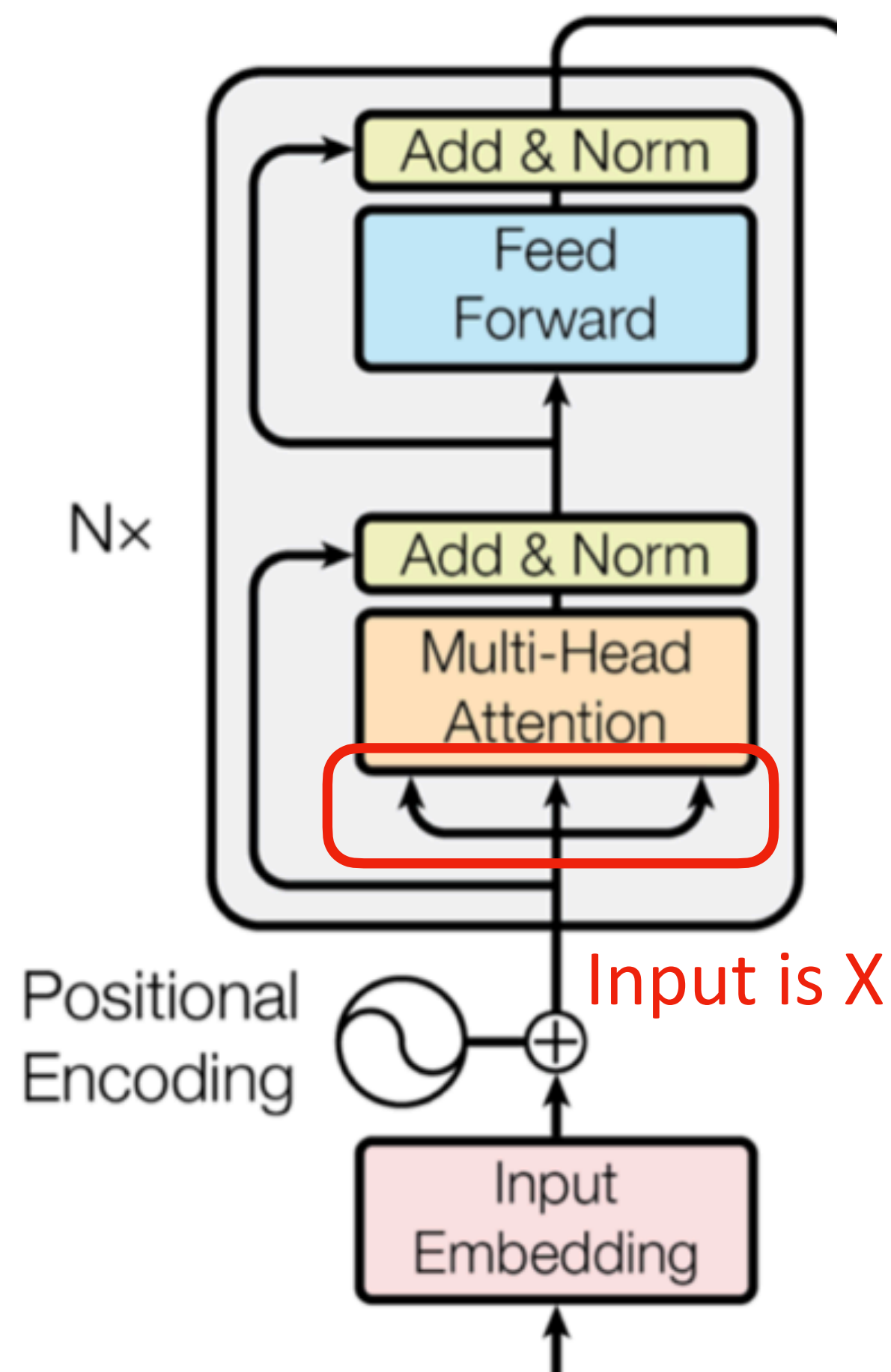
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$

Q, K, V

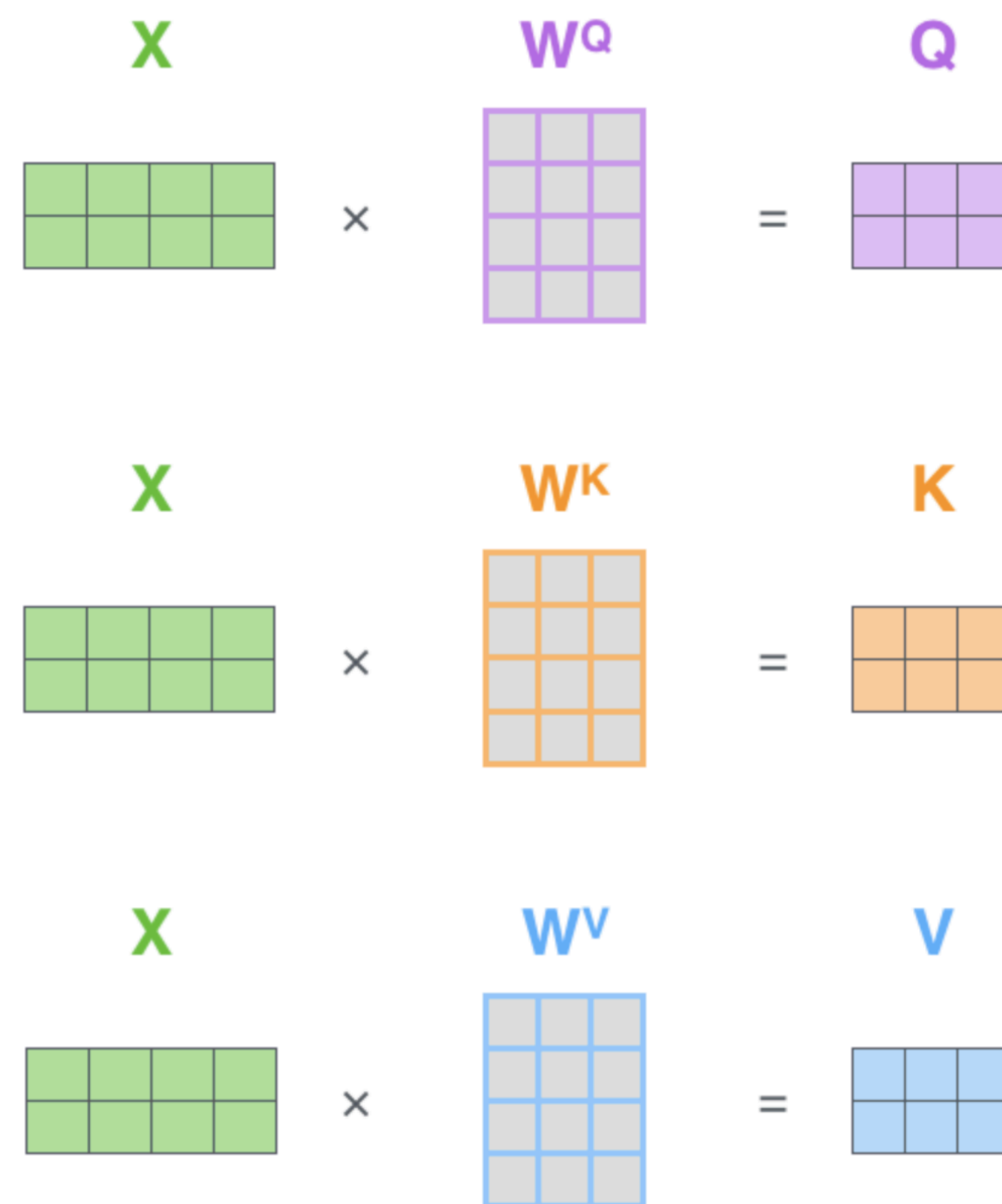


What are Q, K, V in the transformer

Self-Attention



35



35

Query, key, and value are from the same input, thus it is called “self”-attention

$$\text{softmax}\left(\frac{Q \times K^T}{\sqrt{d_k}}\right) V$$

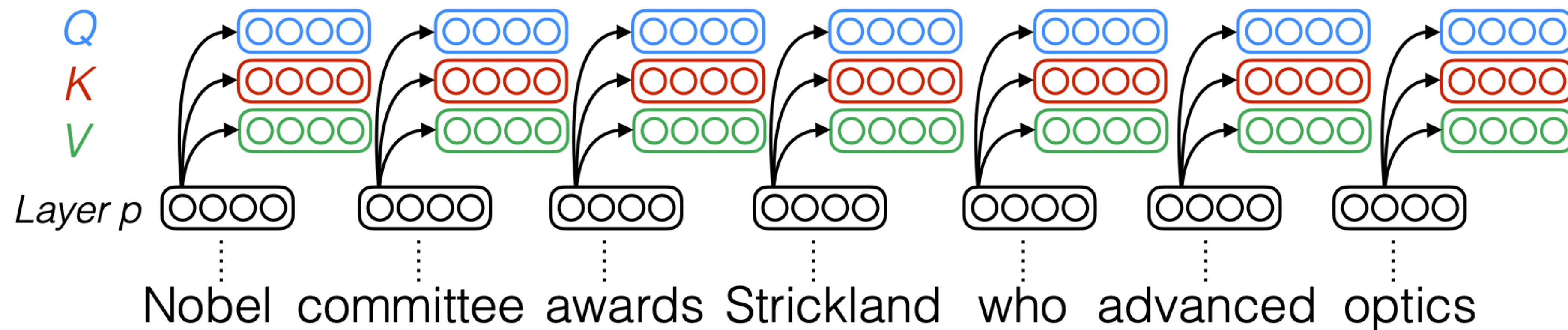
$$= Z$$

The diagram shows the calculation of the self-attention output Z (pink 2x3 grid). It involves the dot product of the Query matrix Q (purple 2x3 grid) and the transpose of the Key matrix K^T (orange 3x2 grid), scaled by $\sqrt{d_k}$, followed by a softmax operation and multiplication with the Value matrix V (blue 2x3 grid).

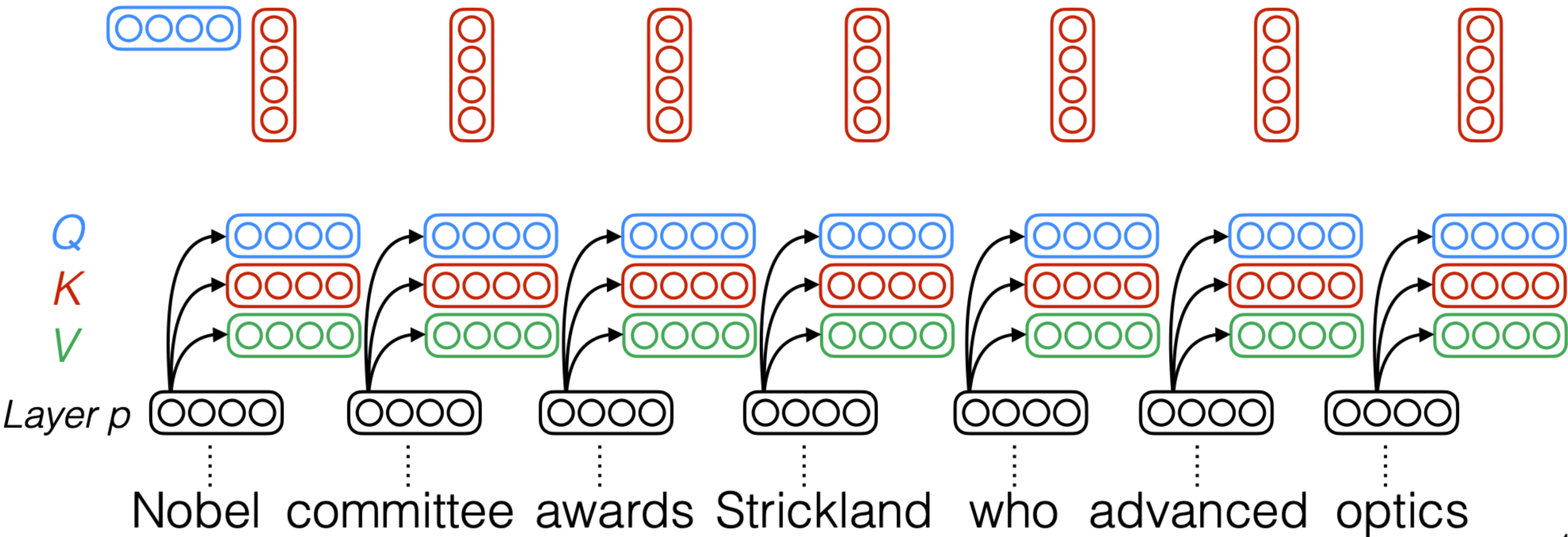
Jay Alammar. The Illustrated Transformer.

Self-Attention

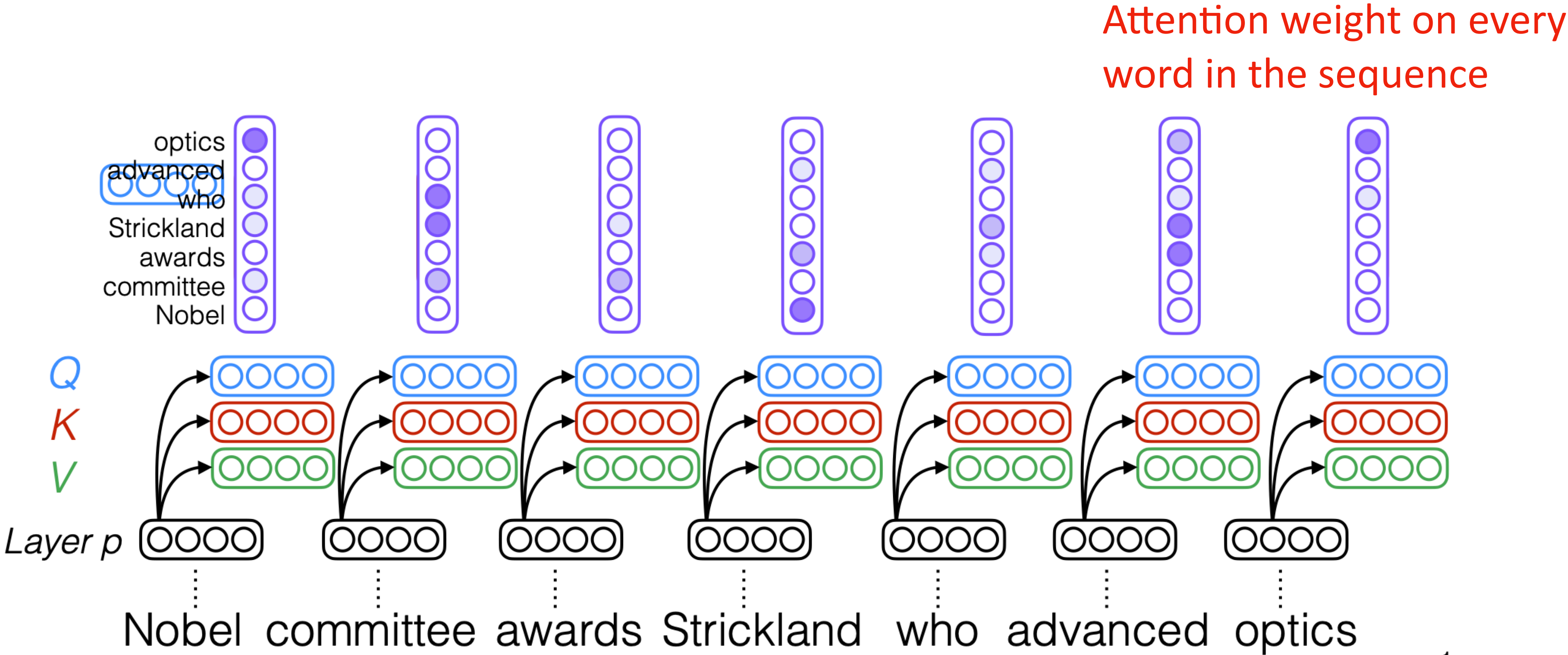
At each step, the attention computation attends to all steps in the input example



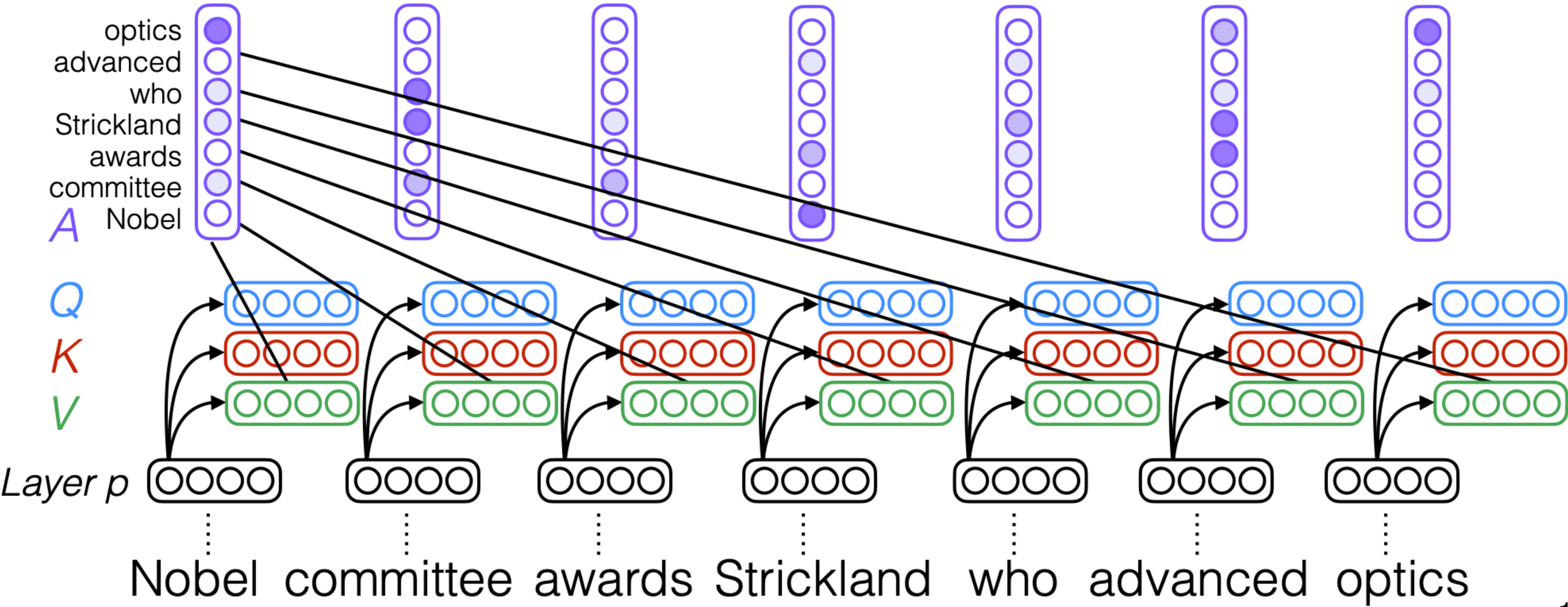
Self-Attention



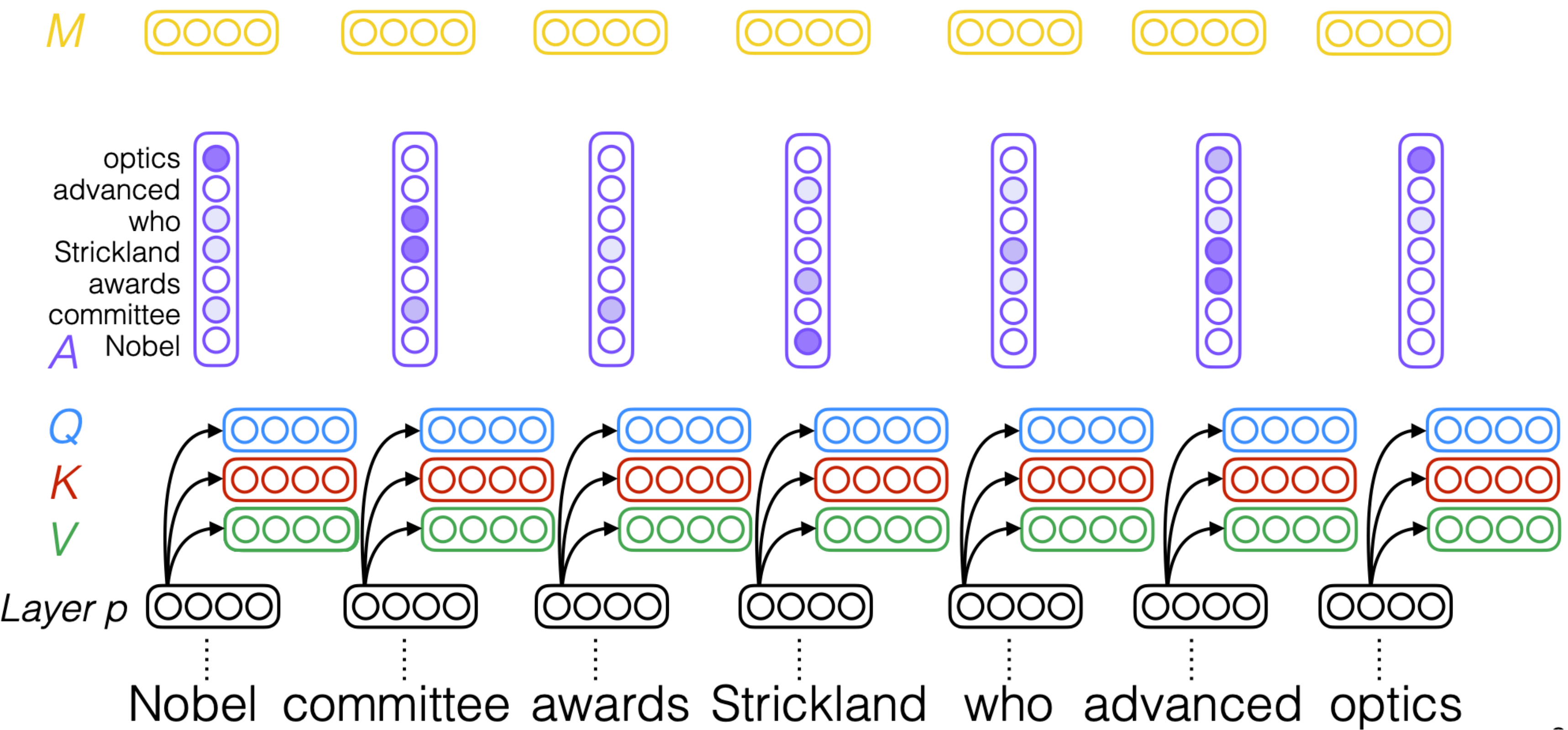
Self-Attention



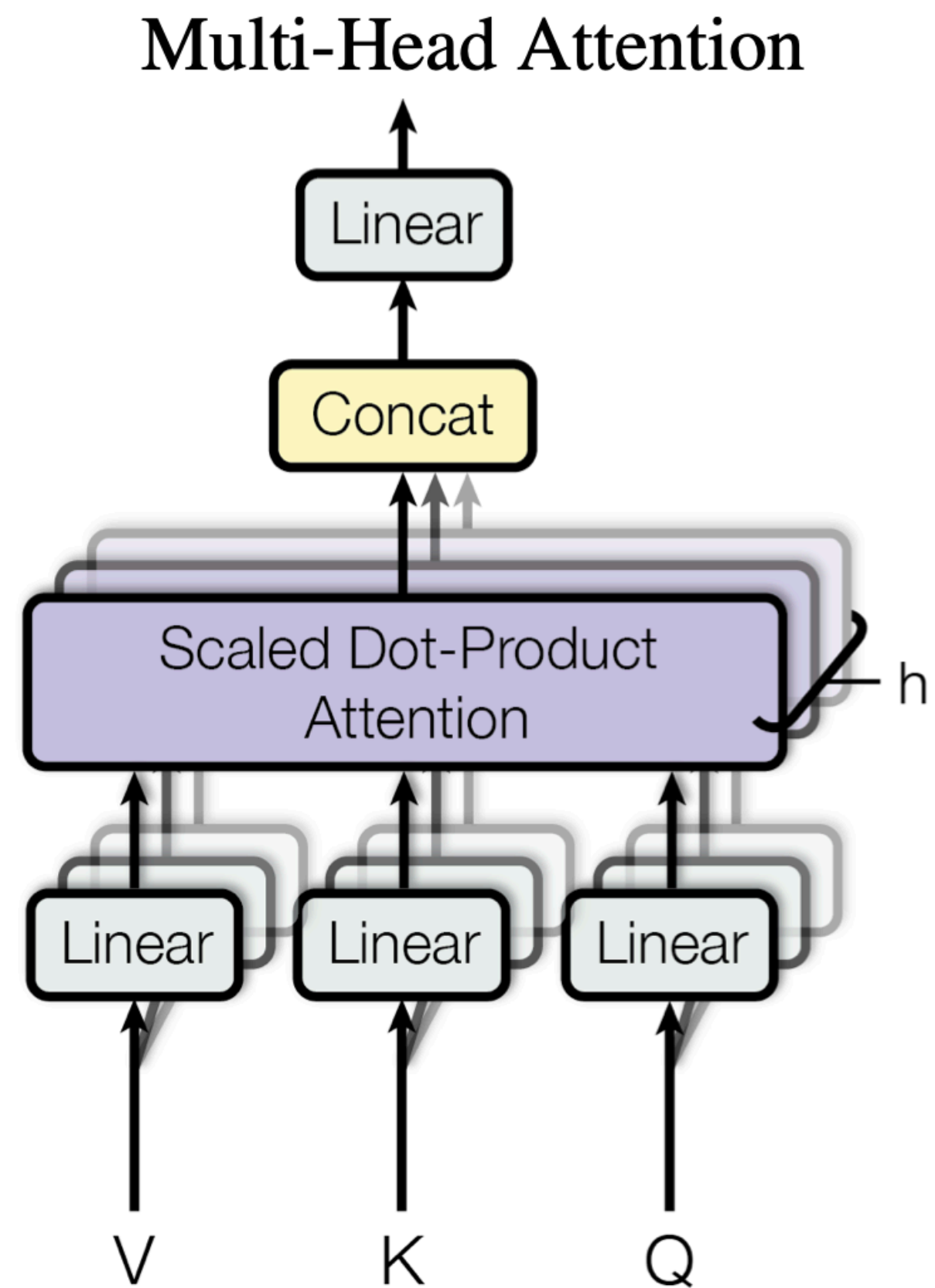
Self-Attention



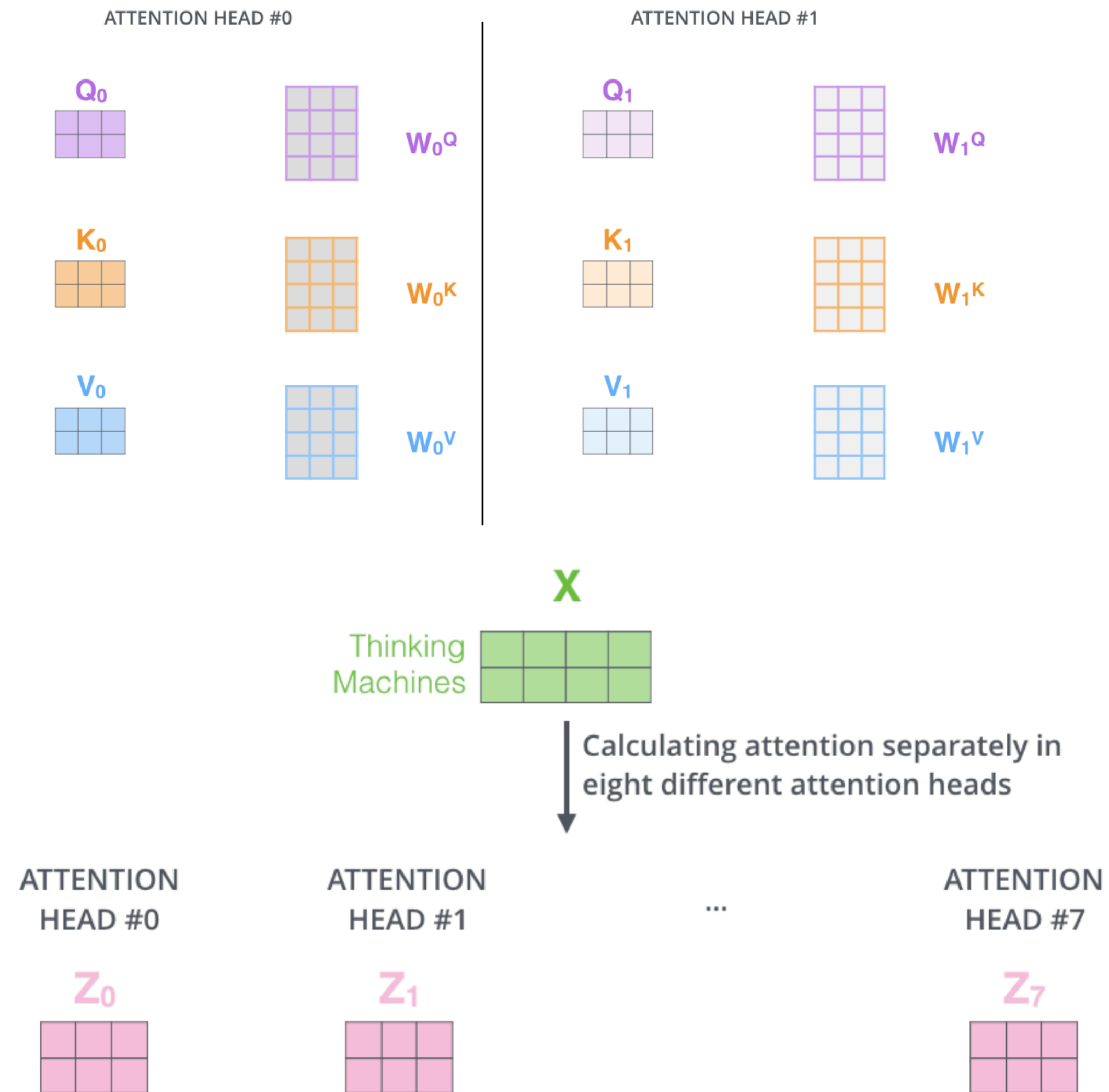
Self-Attention



Multi-Head Attention



Multi-Head Self-Attention



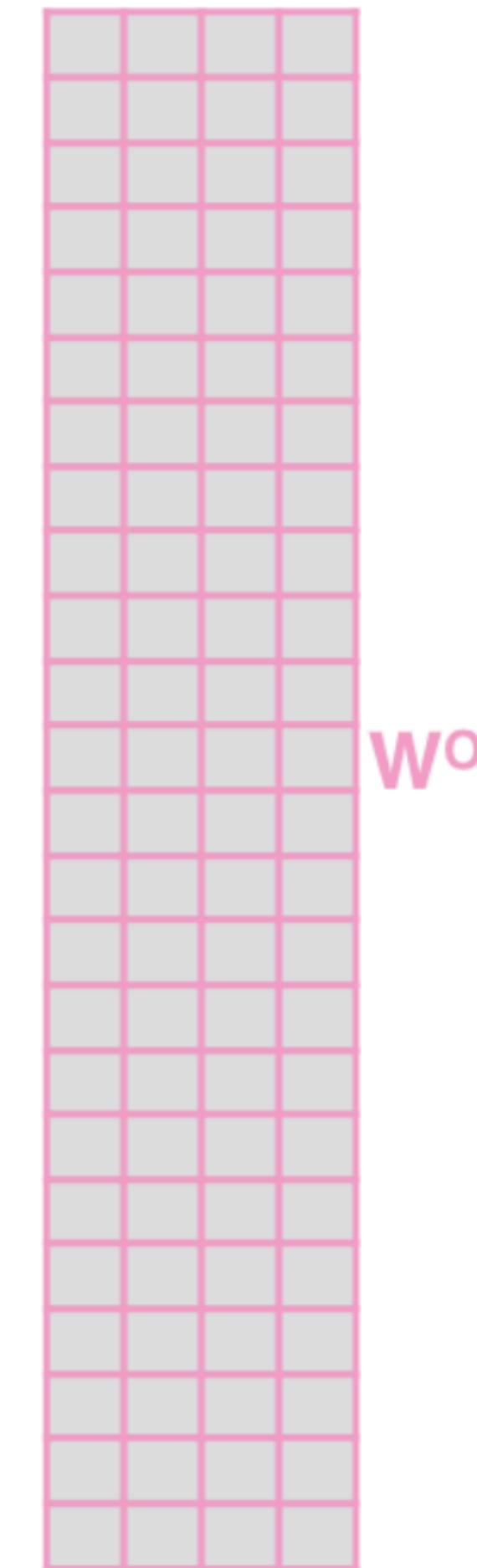
Multi-Head Self-Attention

1) Concatenate all the attention heads

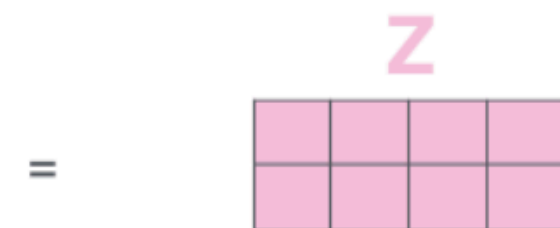


2) Multiply with a weight matrix W^O that was trained jointly with the model

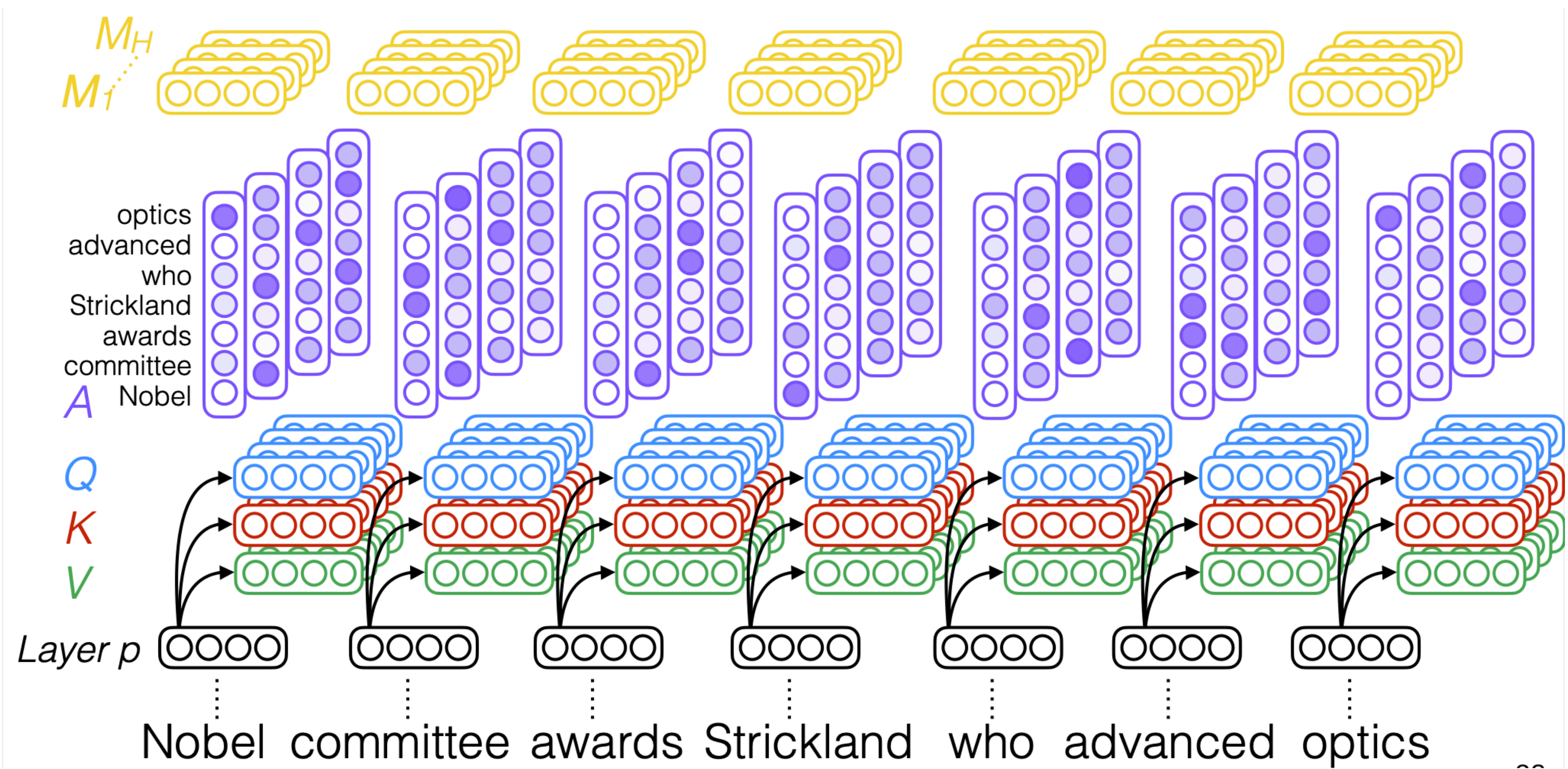
\times



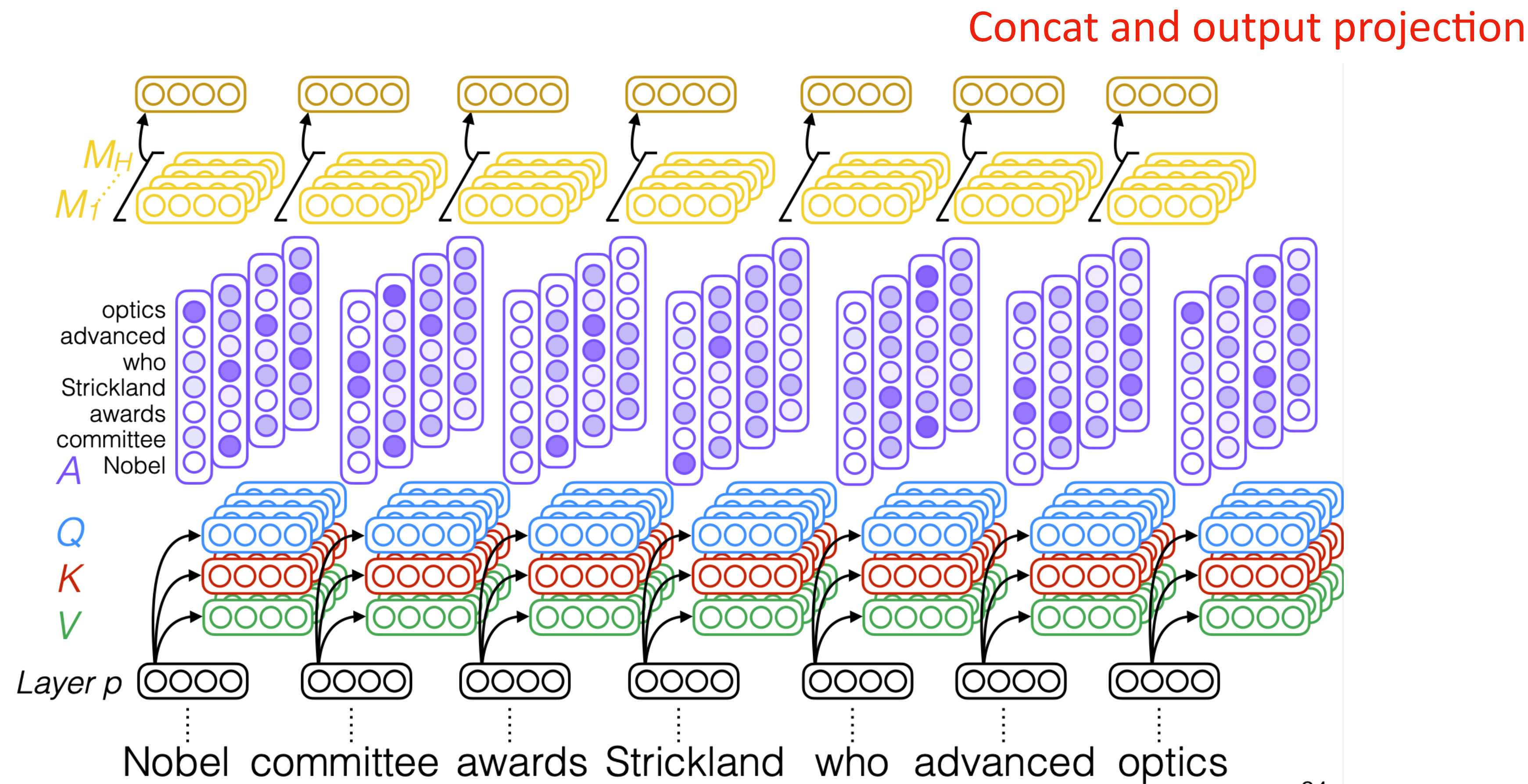
3) The result would be the Z matrix that captures information from all the attention heads. We can send this forward to the FFNN



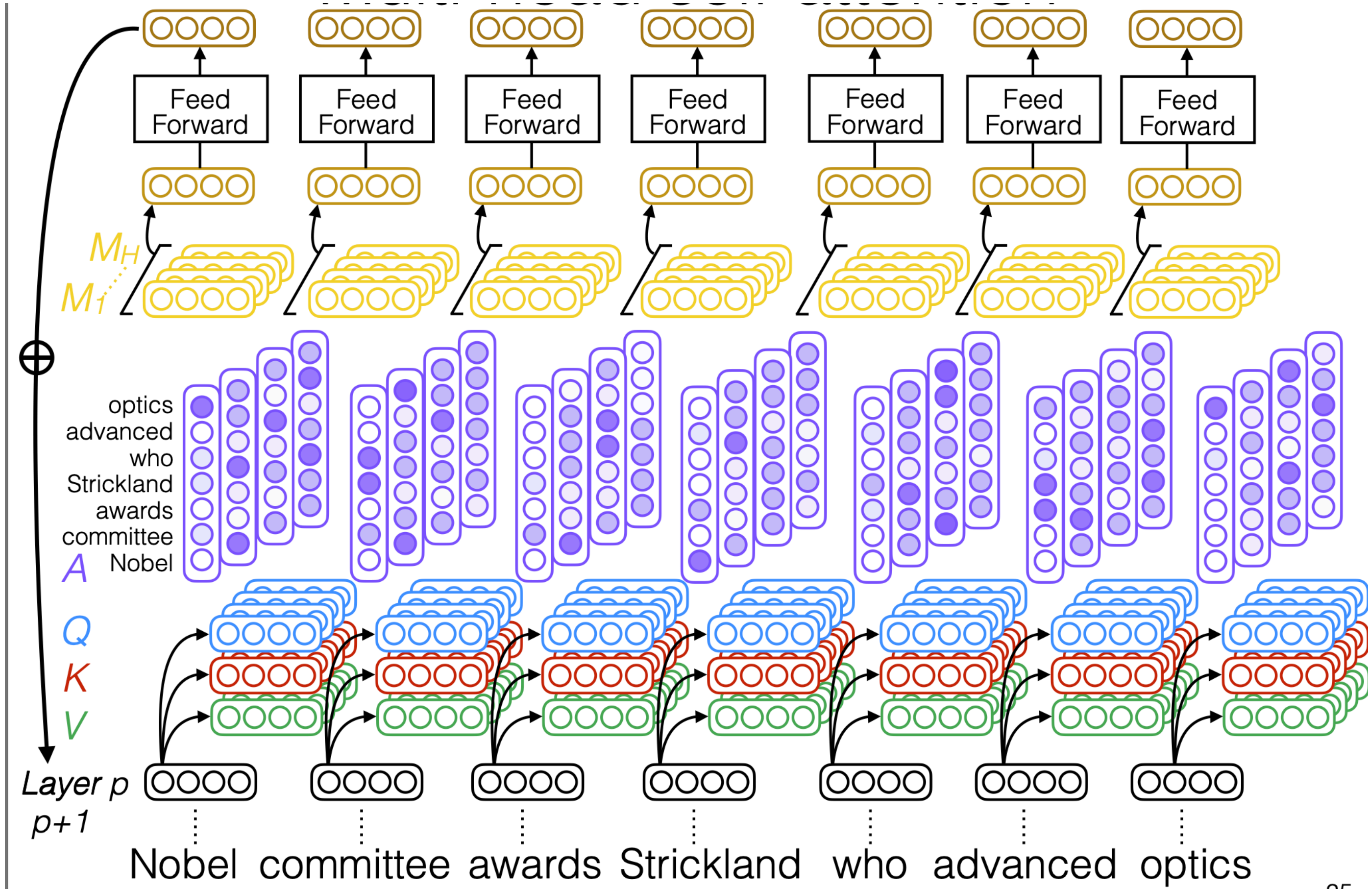
Multi-head Self-Attention



Multi-head Self-Attention

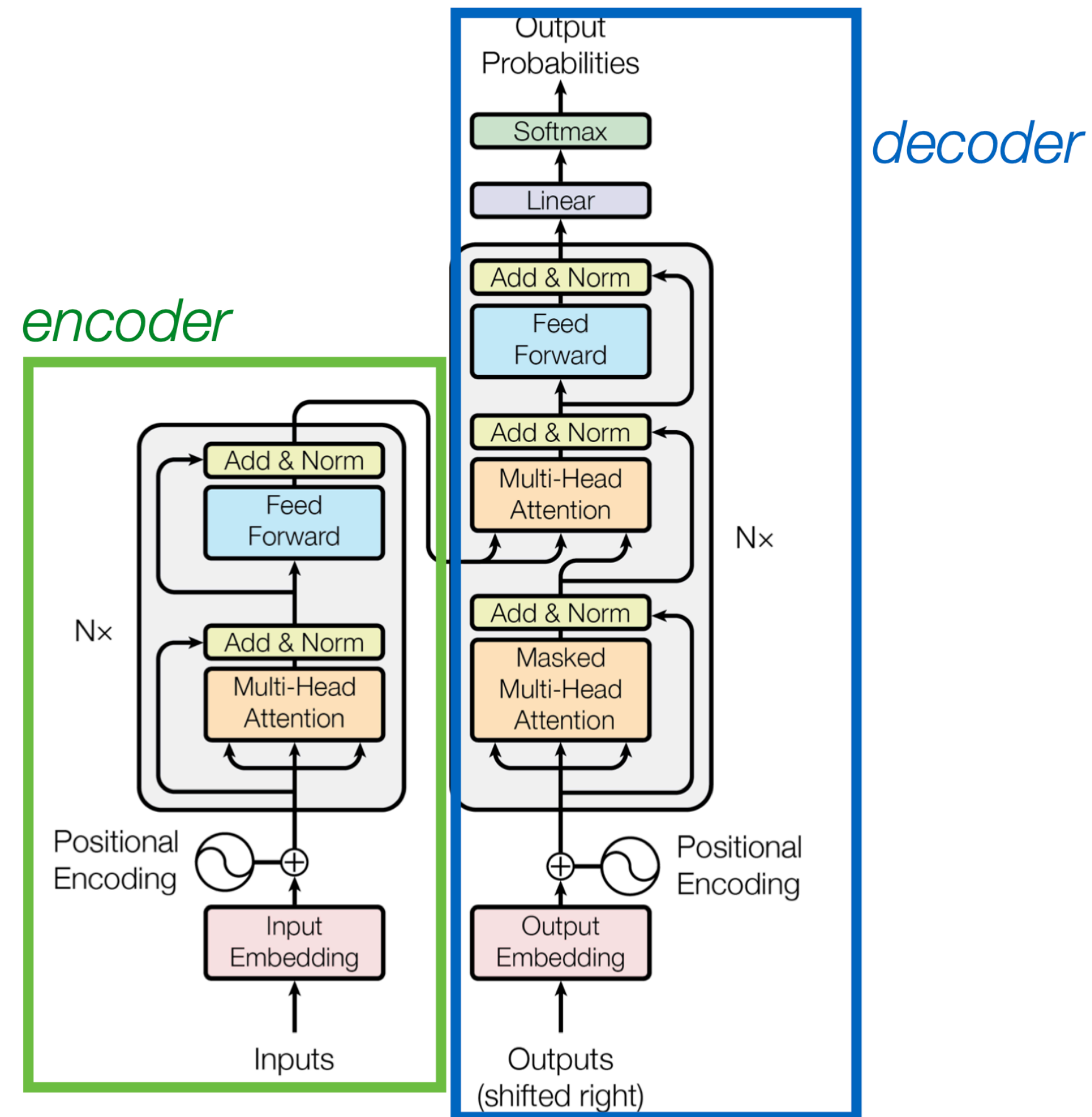
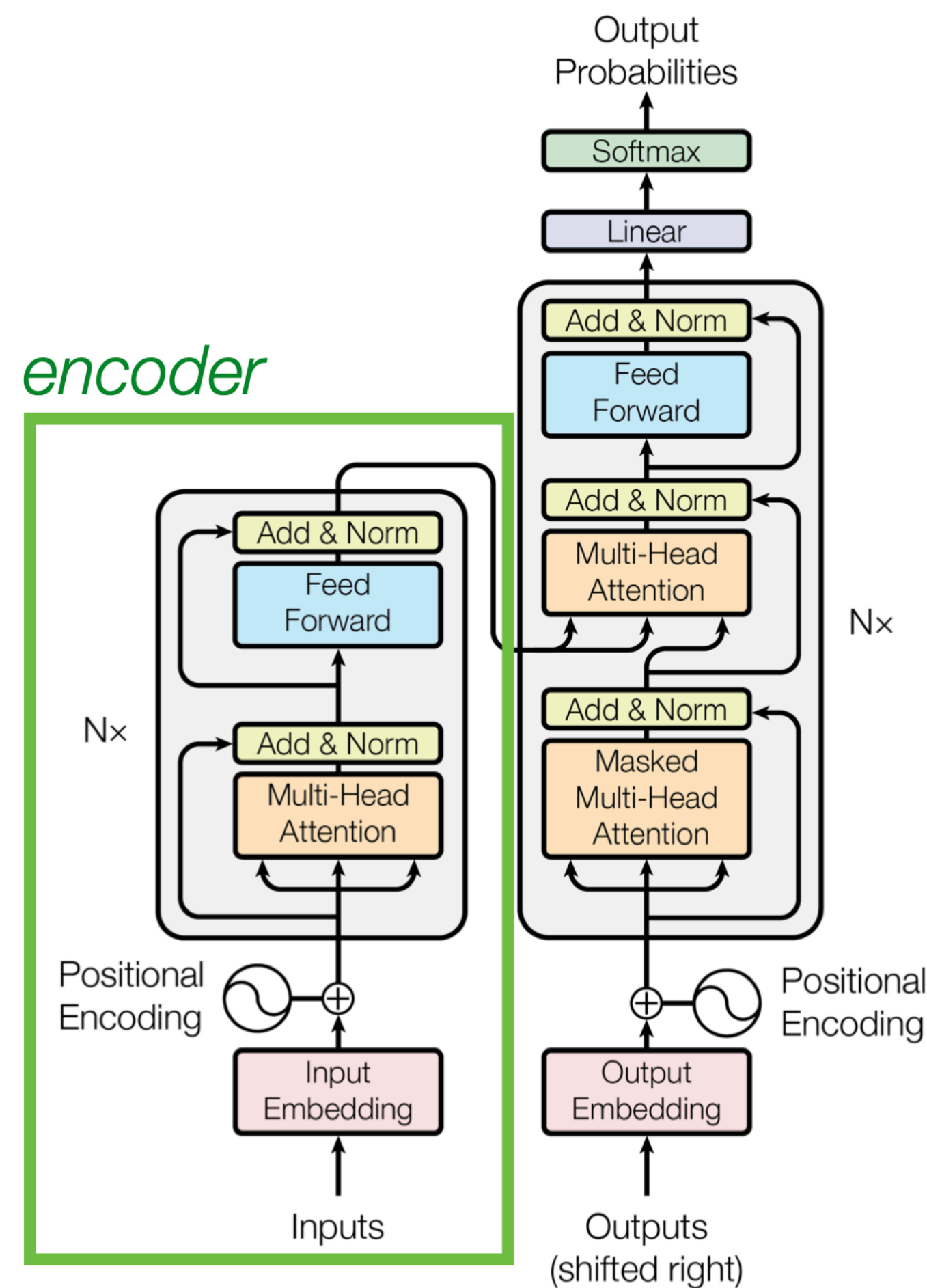


Multi-head Self-Attention + FFN



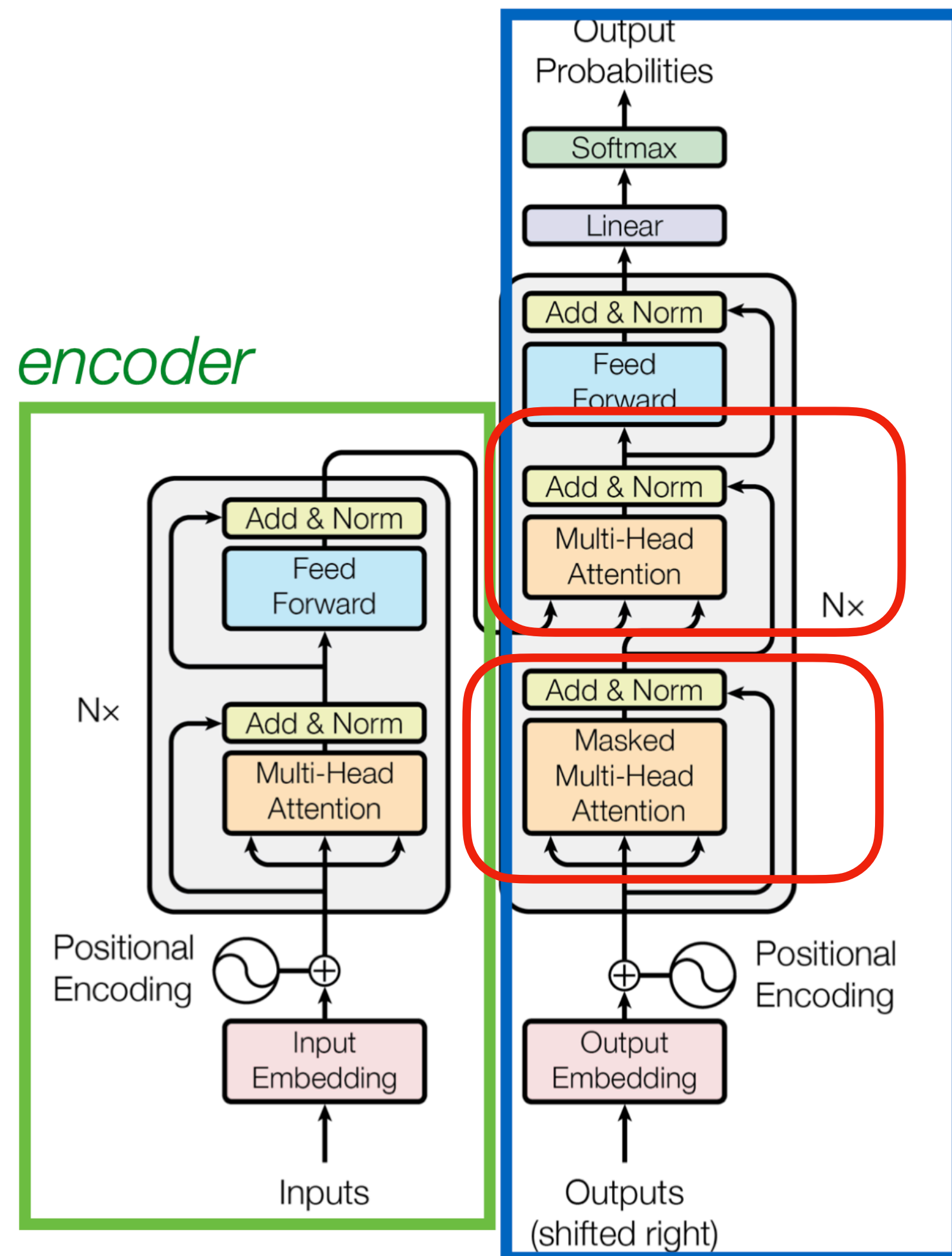
Transformer Encoder

Currently we only cover the encoder side



This encoder-decoder arch is originally proposed as a seq2seq arch, for classification tasks, often only encoder is used. And language models often only have a decoder

Transformer Decoder in Seq2Seq

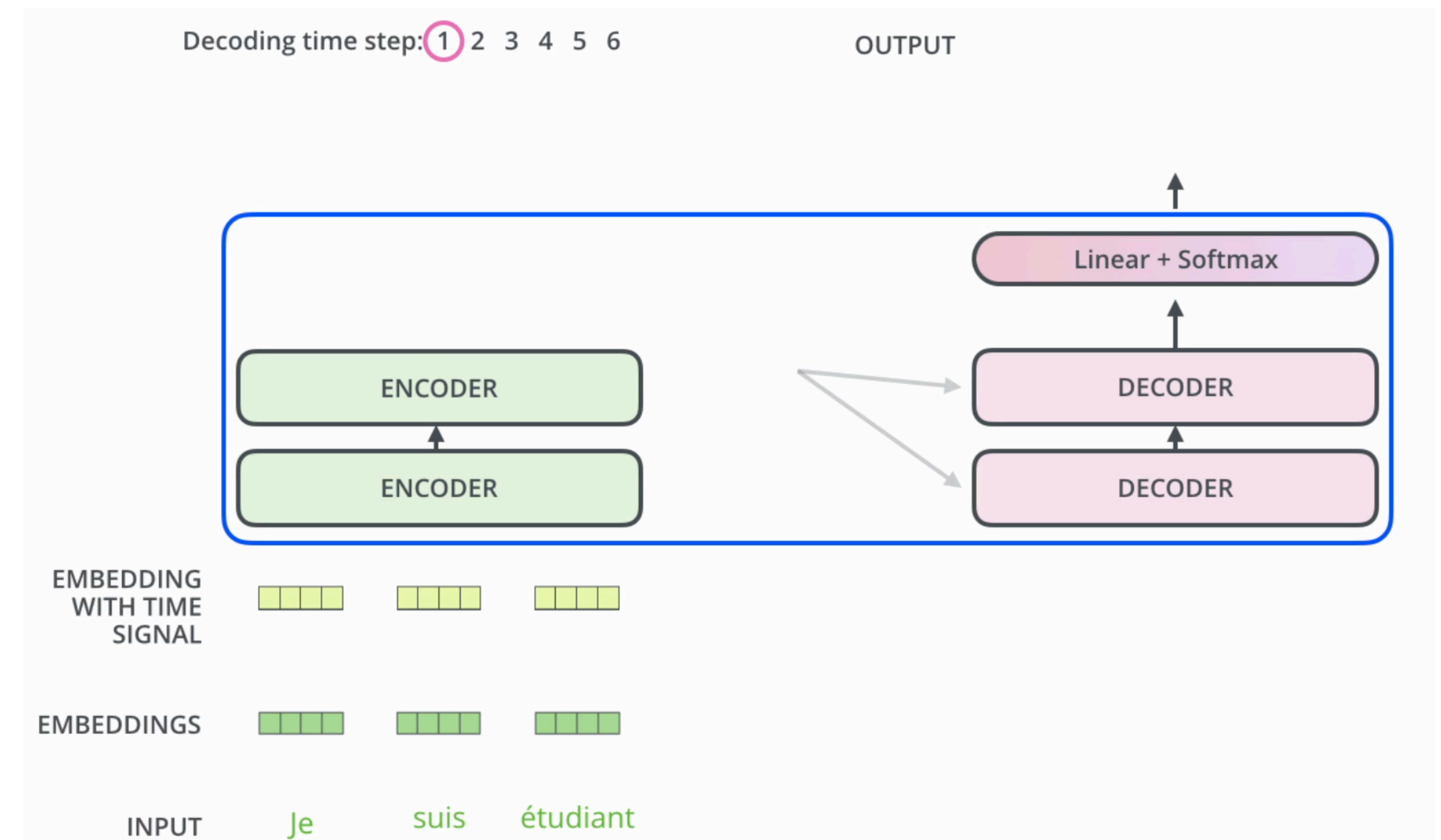


decoder

Cross-attention

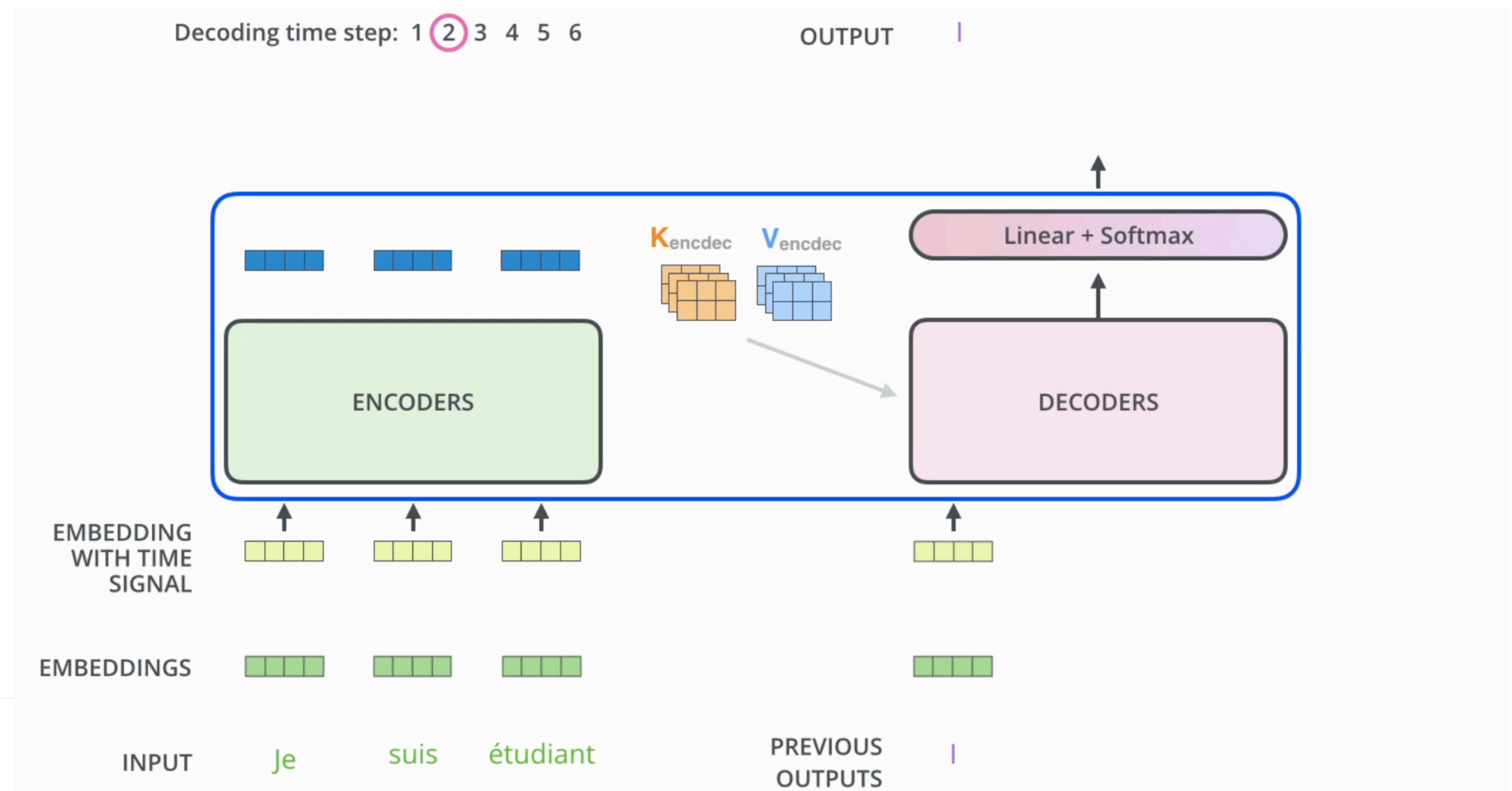
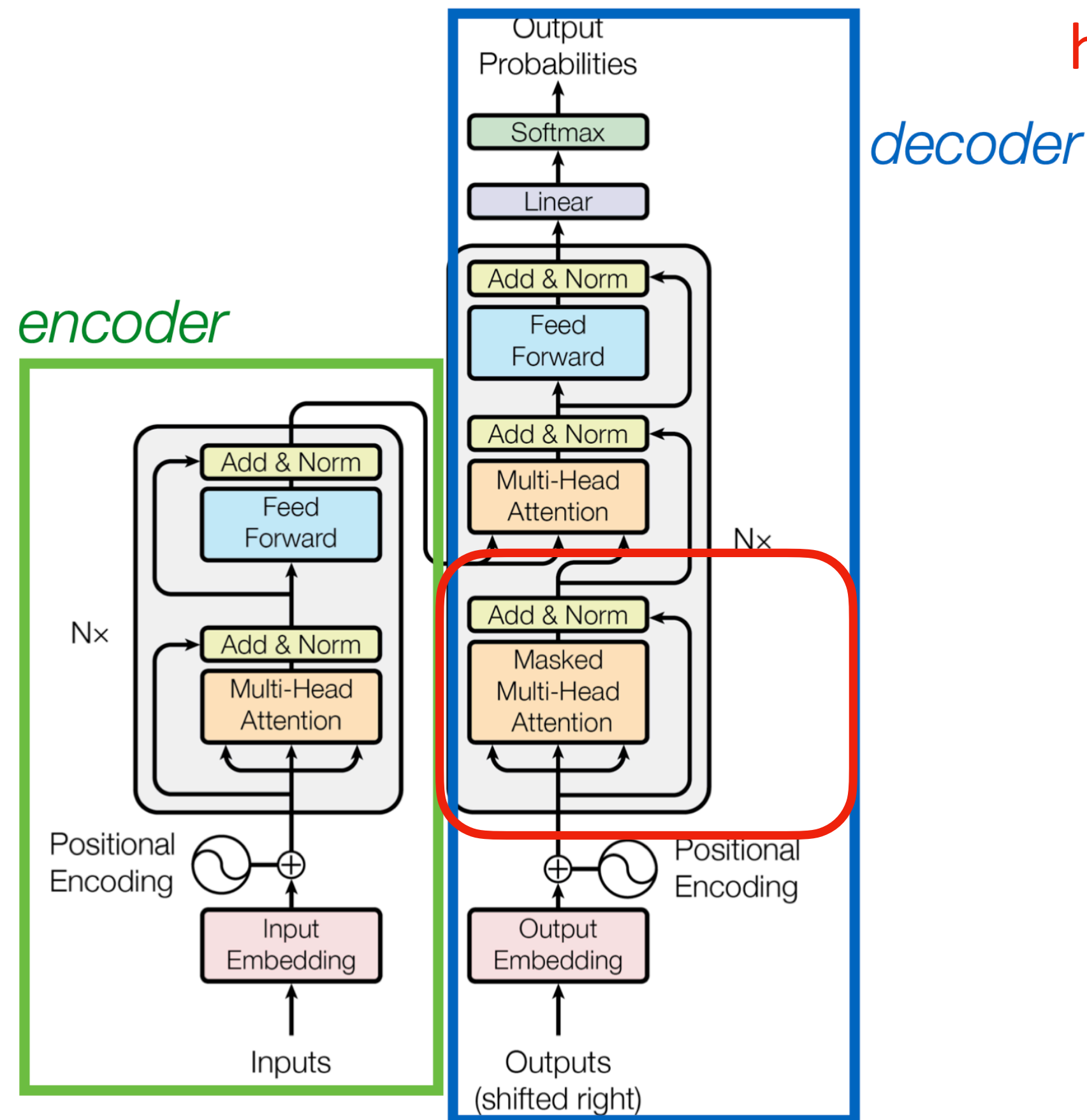
Self-attention

Cross-attention uses the output of encoder as input

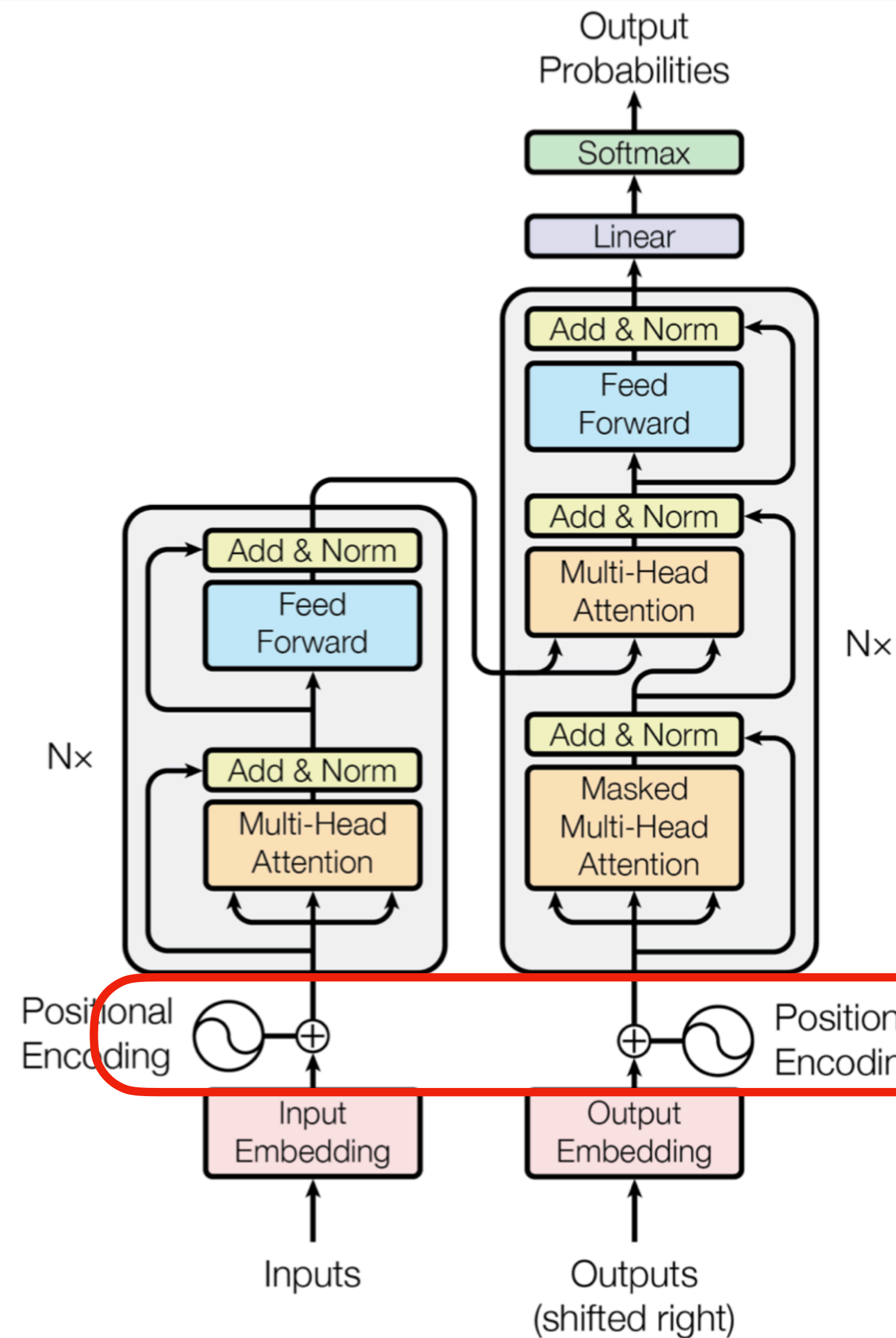


Masked Attention

Typical attention attends to the entire sequence, while masked attention only attends to the ones on the left because future words have not been generated



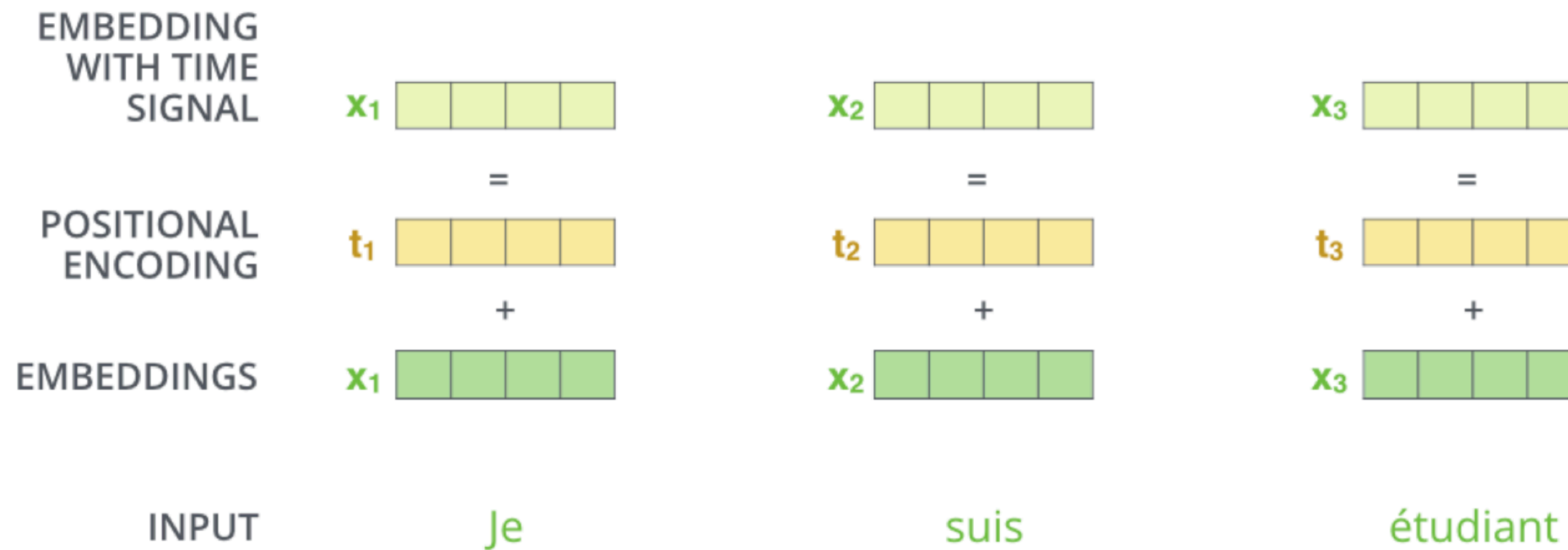
Position Embeddings



Question: If we shuffle the order of words in the sequence, will that change the attention output and feed forward output of the corresponding word?

Position embeddings are added to each word embedding, otherwise our model is unaware of the position of a word

Positional Encoding



Transformer Positional Encoding

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

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

Positional encoding is a 512d vector
 i = a particular dimension of this vector
 pos = dimension of the word
 $d_{model} = 512$

Complexity

Layer Type	Complexity per Layer	Sequential Operations
Self-Attention	$O(n^2 \cdot d)$	$O(1)$
Recurrent	$O(n \cdot d^2)$	$O(n)$
Convolutional	$O(k \cdot n \cdot d^2)$	$O(1)$
Self-Attention (restricted)	$O(r \cdot n \cdot d)$	$O(1)$

n is sequence length, d is embedding dimension.

Restricted self-attention means not attending all words in the sequence, but only a restricted field

Square complexity of sequence length is a major issue for transformers to deal with long sequence

Thank You!