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

# Recurrent Neural Networks, Transformers

Junxian He

Sep 12, 2025

# Recap: Autoregressive Language Models

$p(\text{the, mouse, ate, the, cheese}) = p(\text{the})$

$p(\text{mouse} \mid \text{the})$

$p(\text{ate} \mid \text{the, mouse})$

$p(\text{the} \mid \text{the, mouse, ate})$

$p(\text{cheese} \mid \text{the, mouse, ate, the}).$

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

$p(x_i \mid x_{i:i-1})$

from

NN

Next Word     Context

# Recap: Neural Language Models

# Recap: Neural Language Models

Neural language models are typically autoregressive
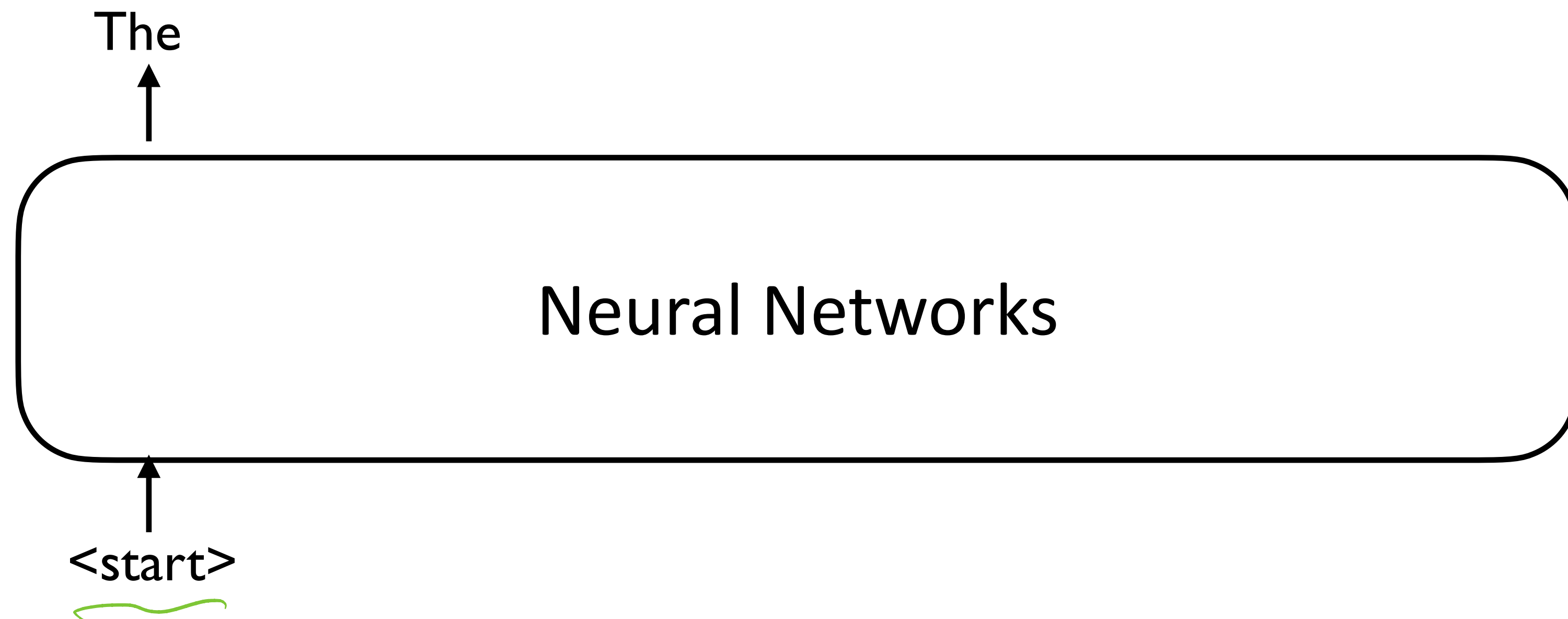
# Recap: Neural Language Models

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

# Recap: Neural Language Models
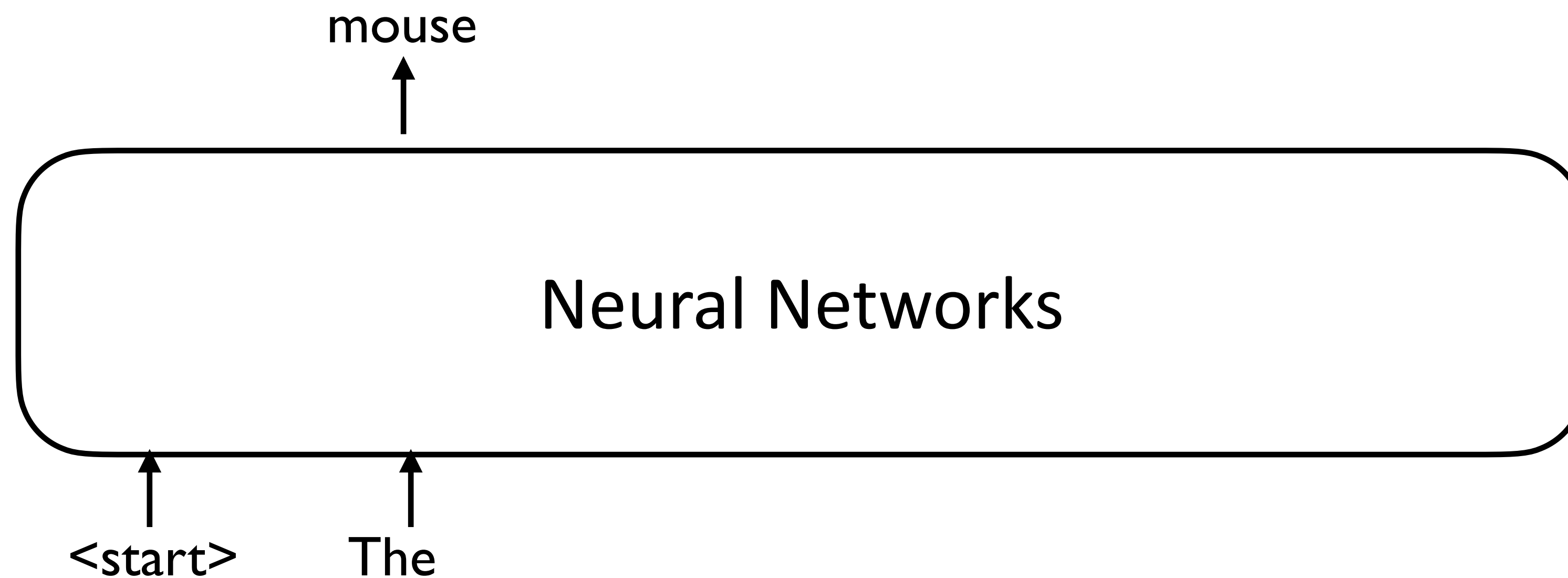
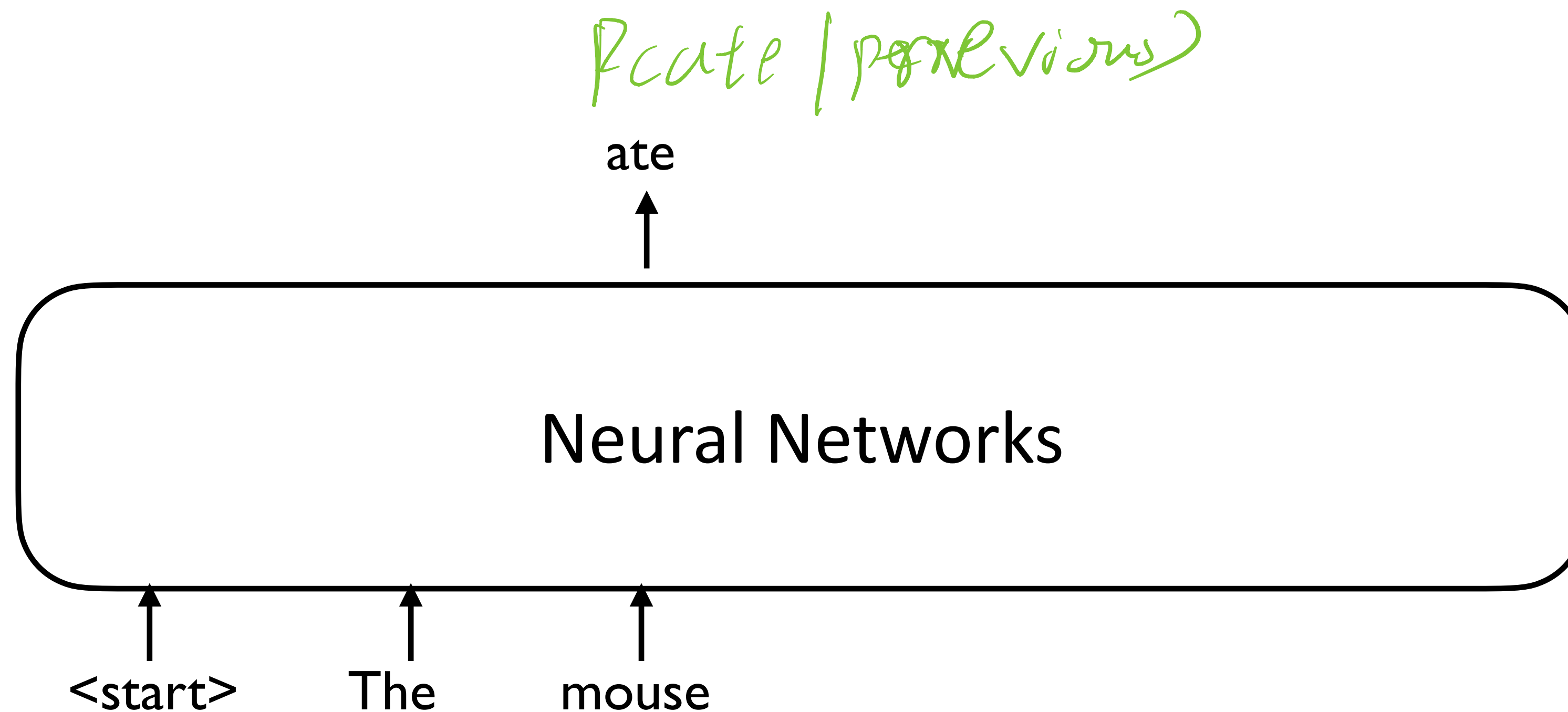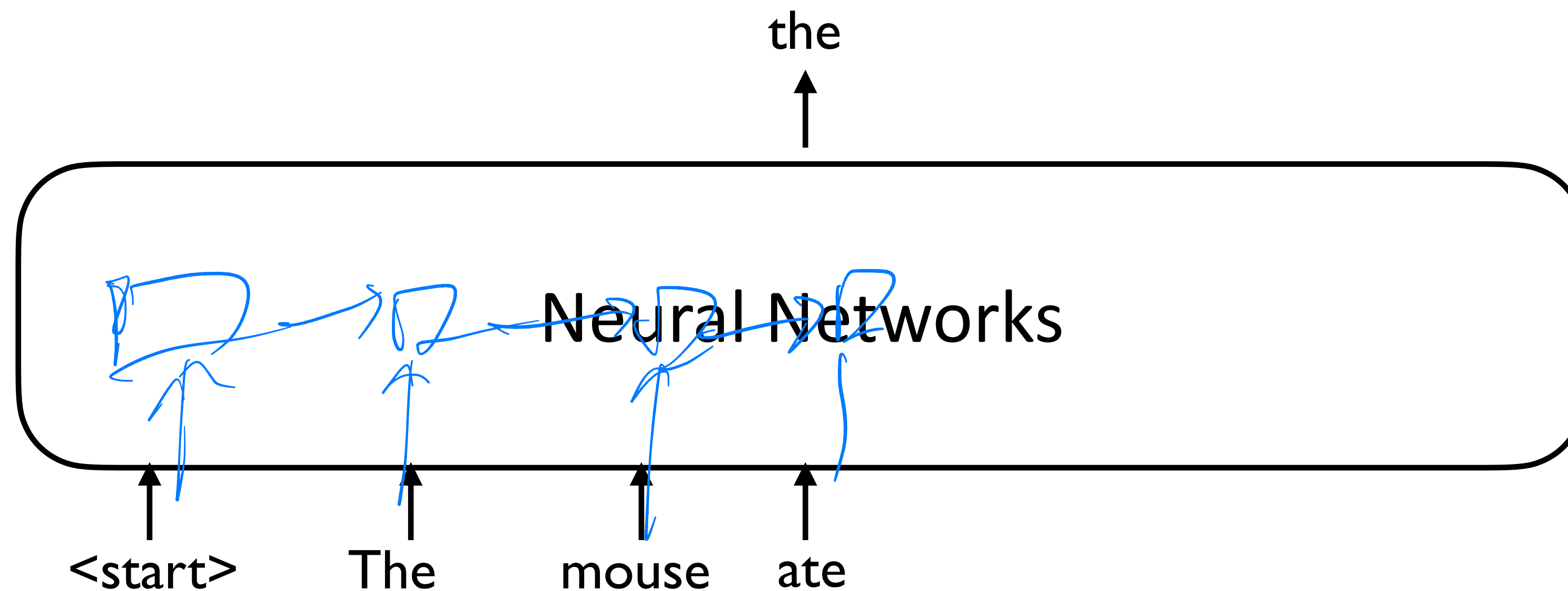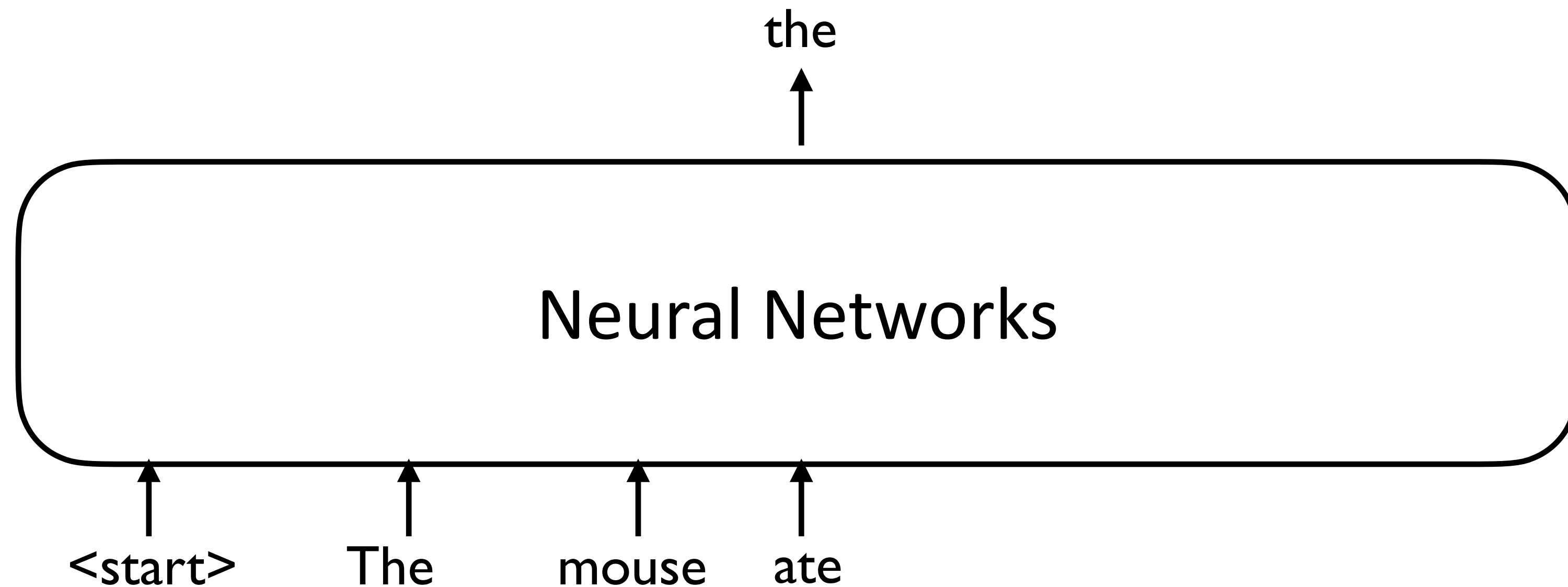Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."
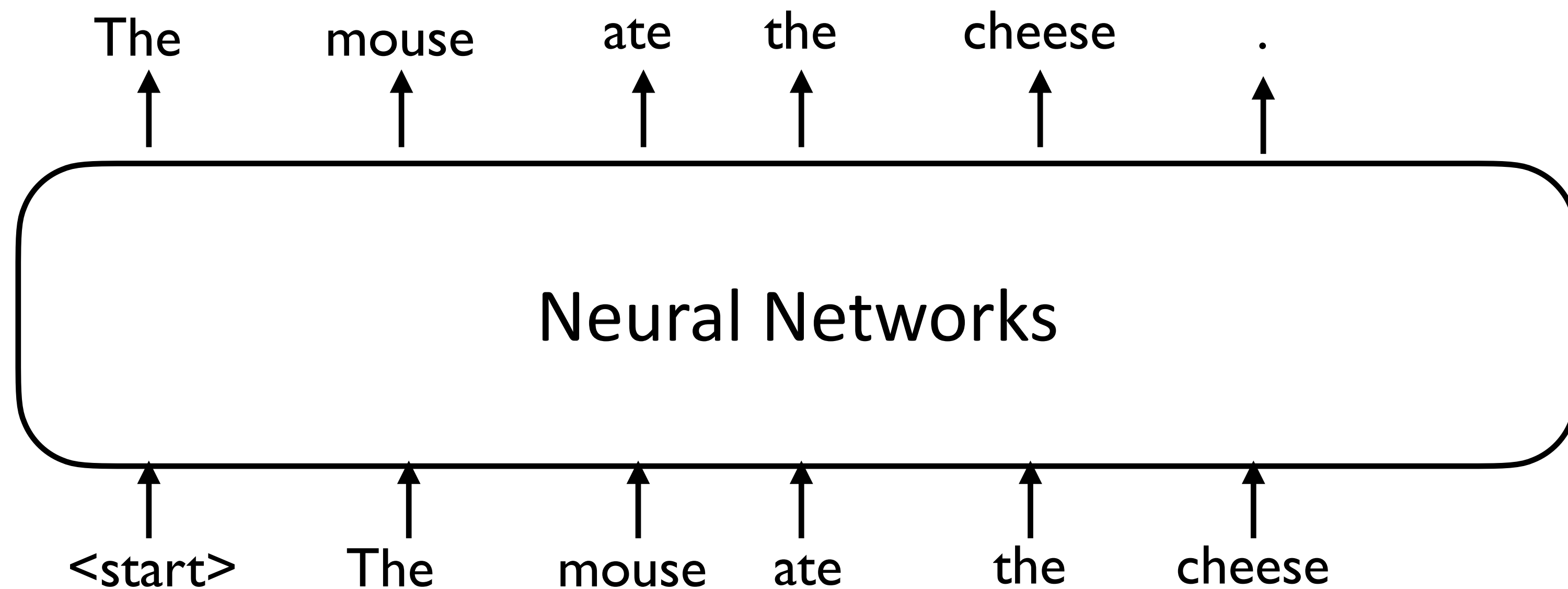
The



Neural Networks

# Recap: Neural Language Models

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

mouse

$\uparrow$

Neural Networks

$\uparrow$      $\uparrow$

&lt;start&gt;      The

# Recap: Neural Language Models

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

$P(ate \mid prev\ previous)$

ate

↑

Neural Networks

↑          ↑          ↑

&lt;start&gt;      The      mouse

# Recap: **Neural Language Models**

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

the

Neural Networks

<start>     The     mouse   ate

# Recap: Neural Language Models

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

the

↑

Neural Networks

↑ ↑ ↑ ↑

&lt;start&gt;     The     mouse    ate

We can compute the loss on every token in parallel

# Recap: Neural Language Models
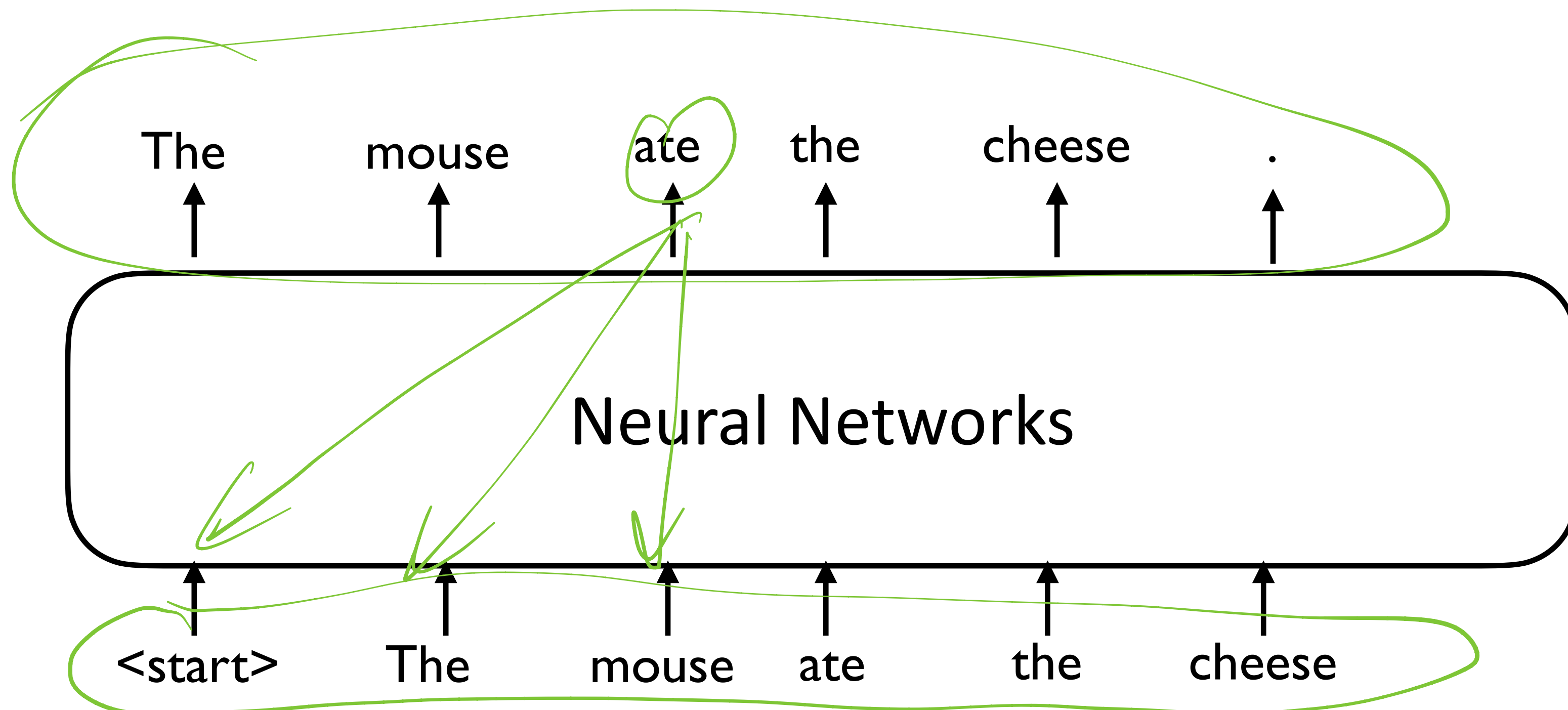
Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."

The    mouse    ate    the    cheese    .

**Neural Networks**

<start>    The    mouse    ate    the    cheese

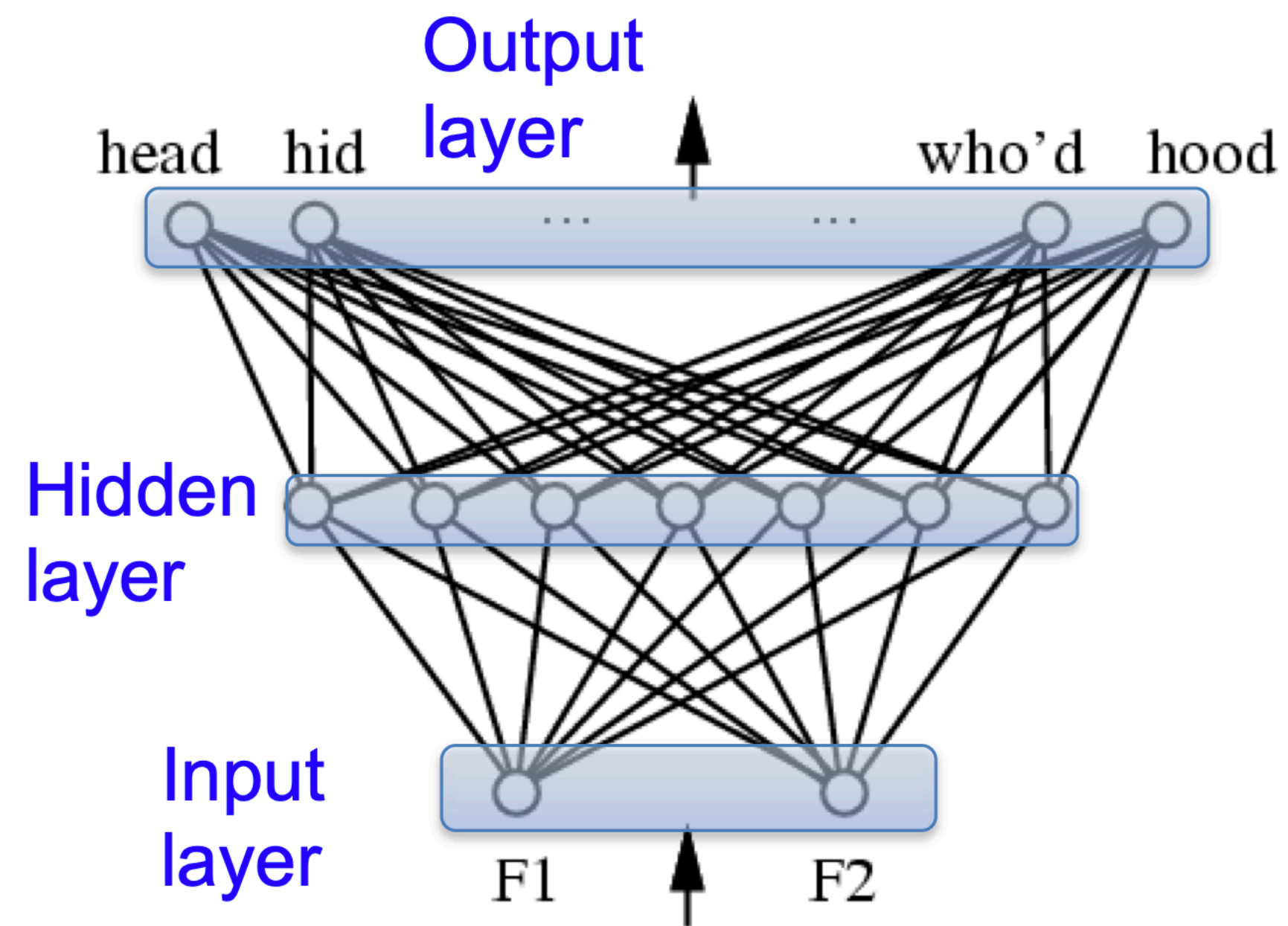# Recap: **Neural Language Models**

Neural language models are typically autoregressive

Data: "The mouse ate the cheese ."
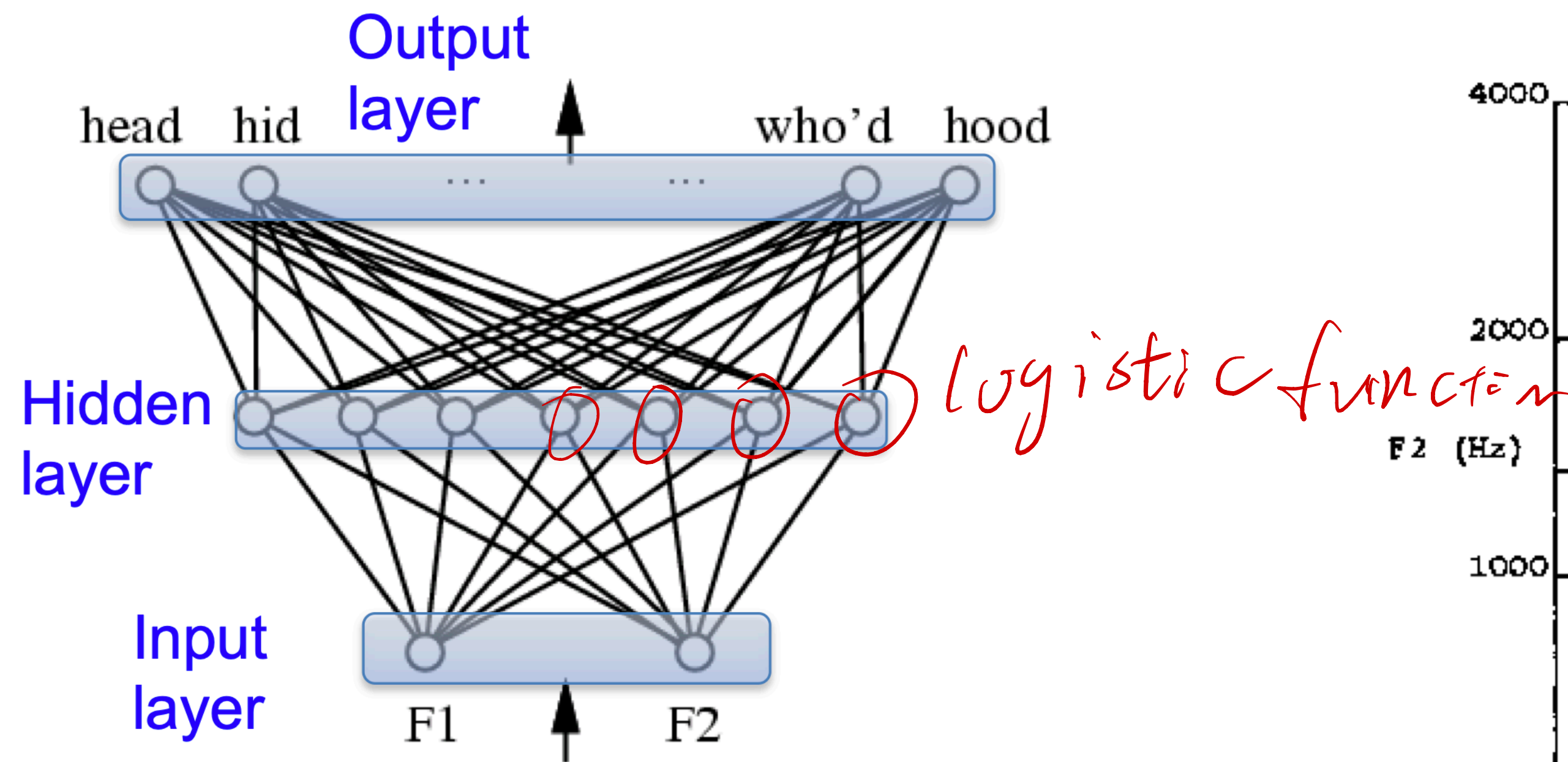


Each prediction only sees the inputs on its left

# Recap: Multilayer Networks of Sigmoid Units

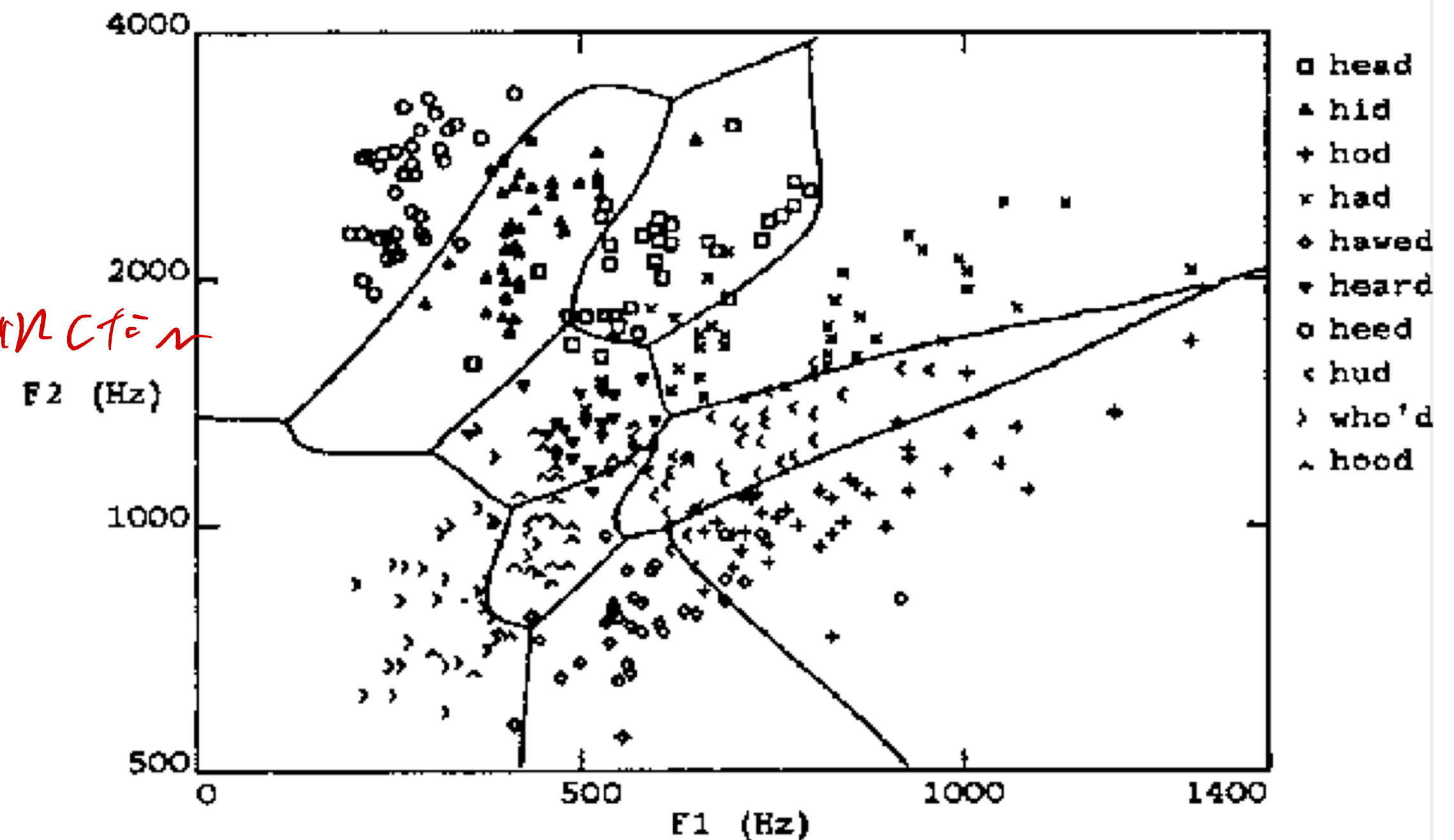# Recap: Multilayer Networks of Sigmoid Units



Two layers of logistic units

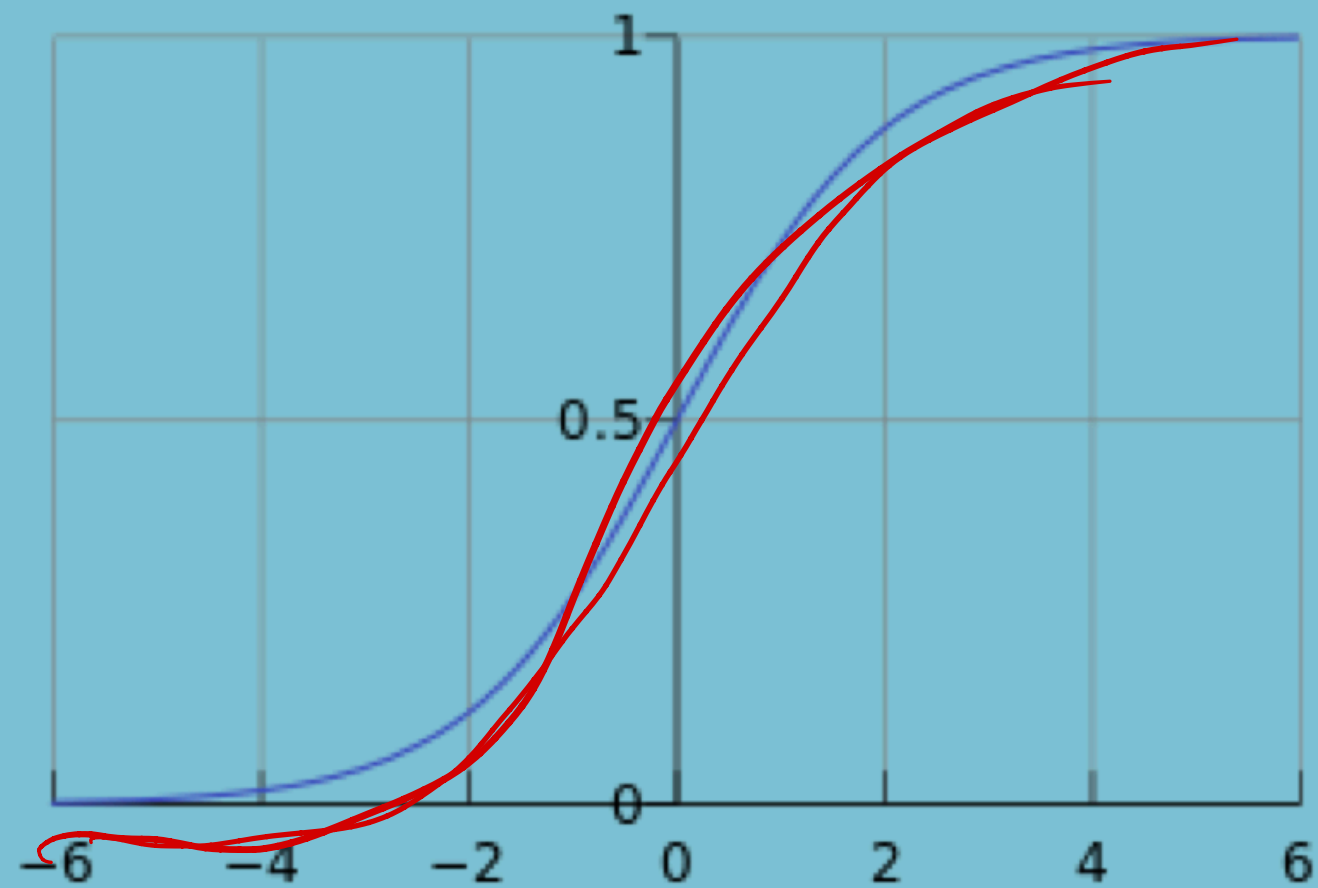# Recap: Multilayer Networks of Sigmoid Units



Output layer

head    hid    ...    ...    who'd    hood

Hidden layer

logistic function

Input layer

F1    F2

Two layers of logistic units



Highly non-linear decision surface

8

# Activation Functions

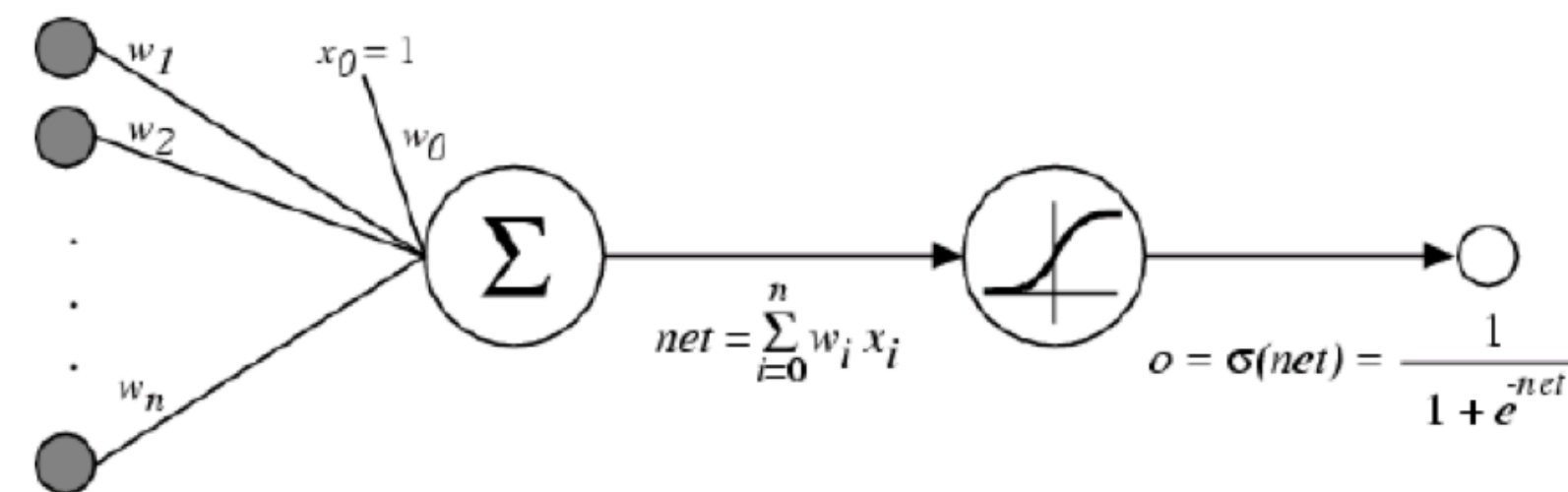$$x \longrightarrow Ax \longrightarrow BAx \longrightarrow CBAx$$

## Sigmoid / Logistic Function

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



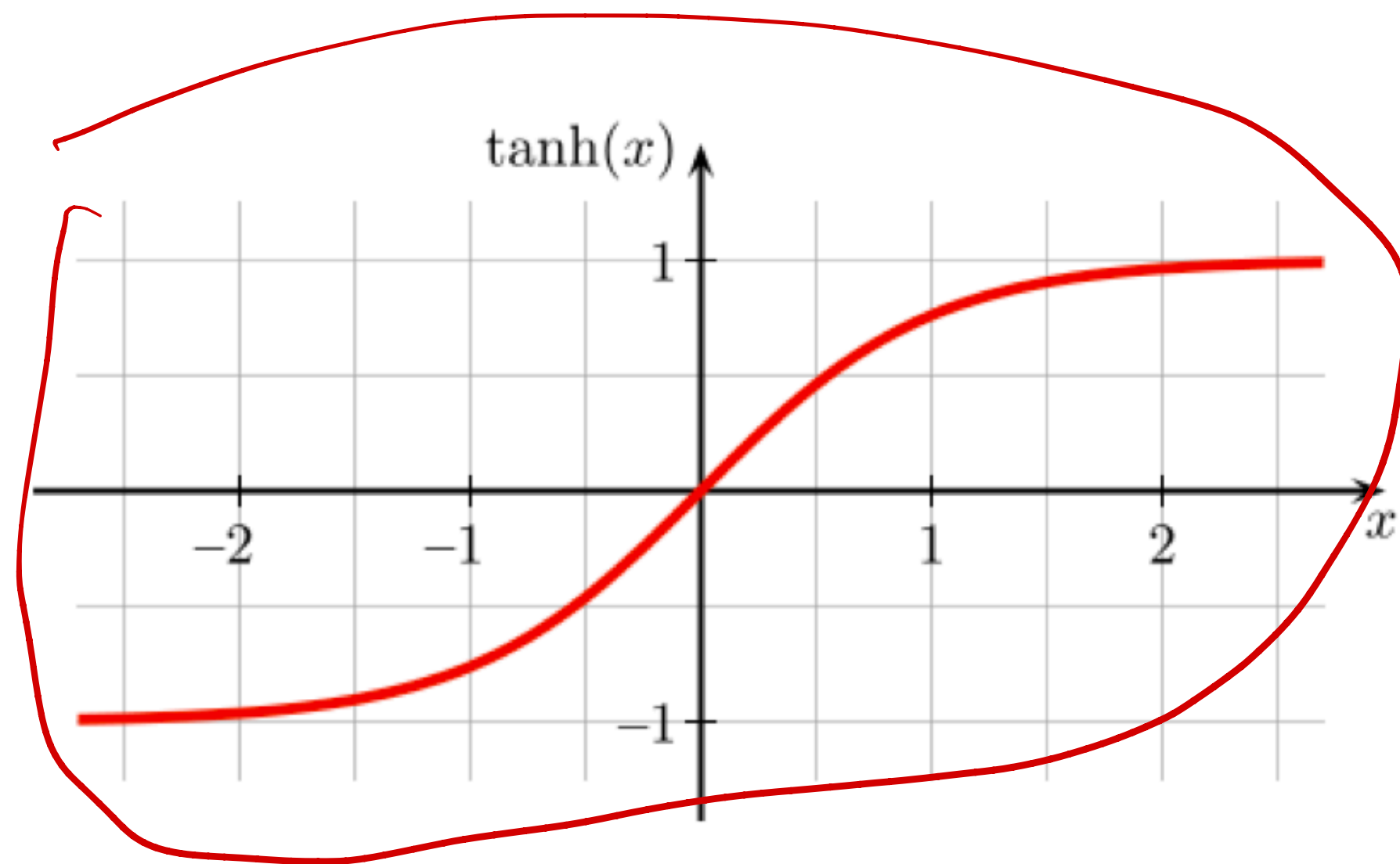So far, we've assumed that the activation function (nonlinearity) is always the sigmoid function…

$$CBA$$
$$||$$
$$Mx$$



$$net = \sum_{i=0}^{n} w_i \, x_i$$

$$o = \sigma(net) = \frac{1}{1+e^{-net}}$$

9

# Tanh

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

*negative*



Alternate 1:
tanh

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

# Activation Function



Understanding the difficulty of training deep feedforward neural networks
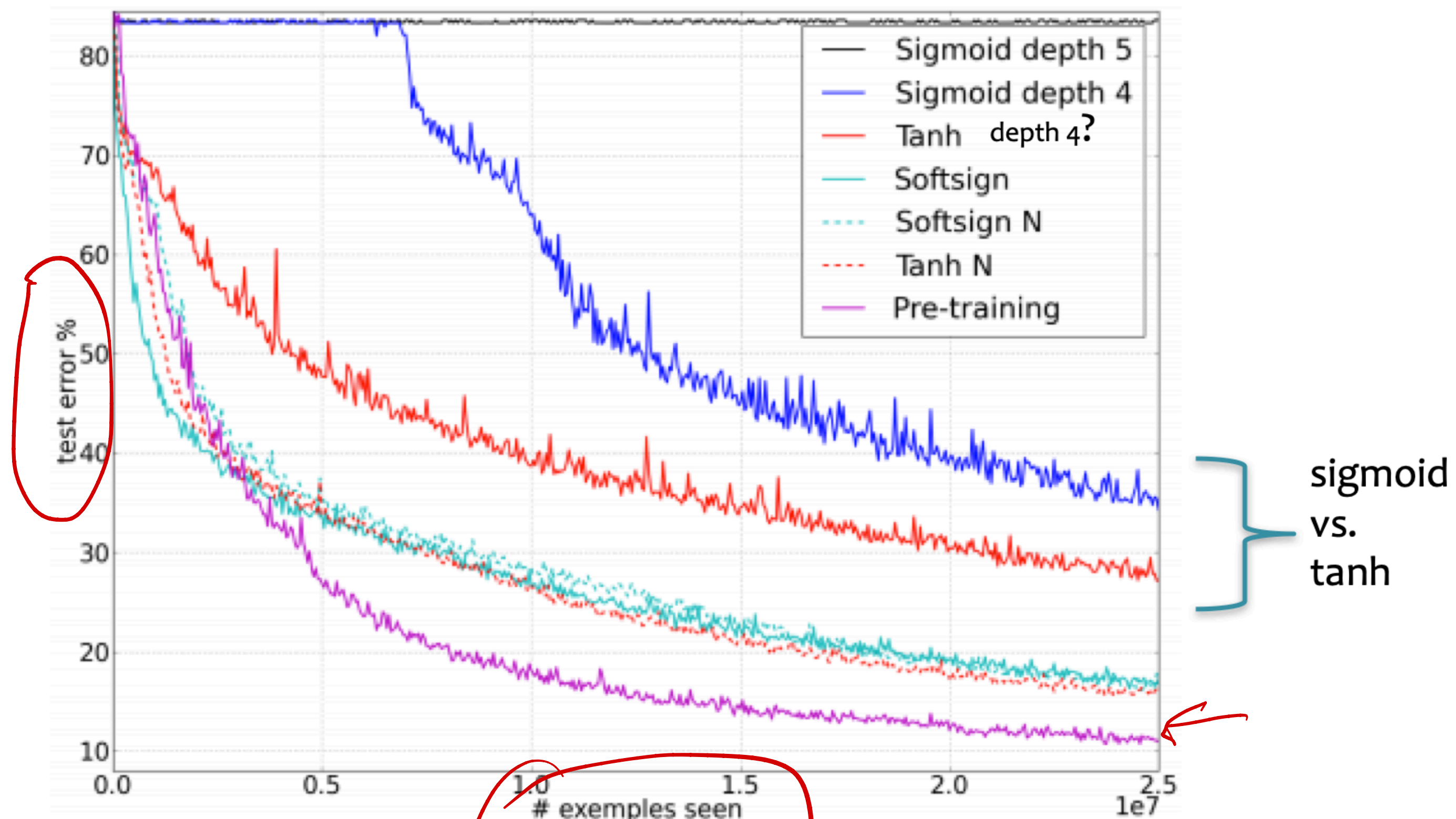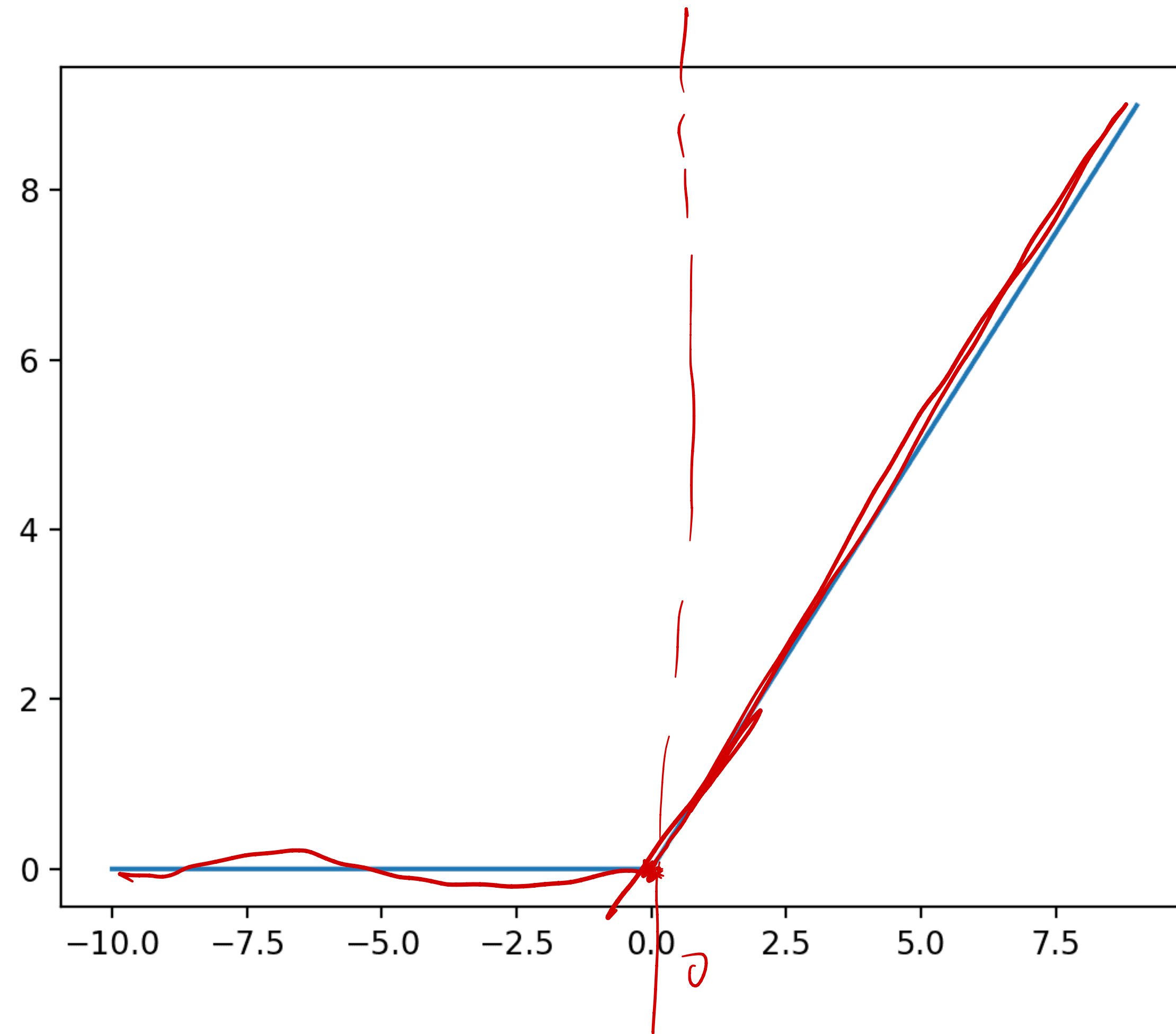
AI Stats 2010

Figure from Glorot & Bentio (2010)

11

# ReLU

$y = max(0, x)$

$Relu(x) = max(0, x)$
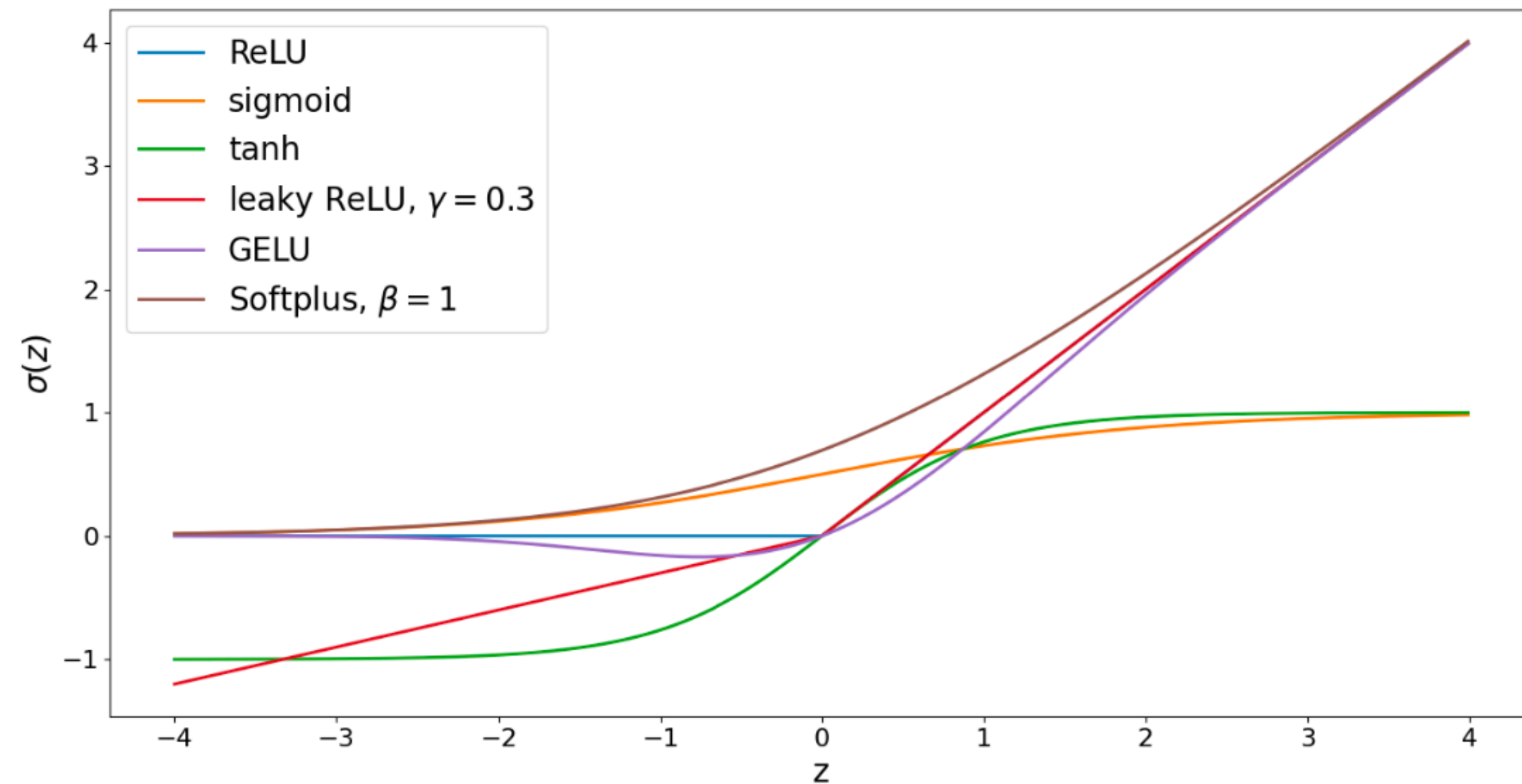
# Other Activation Functions

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

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

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

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

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

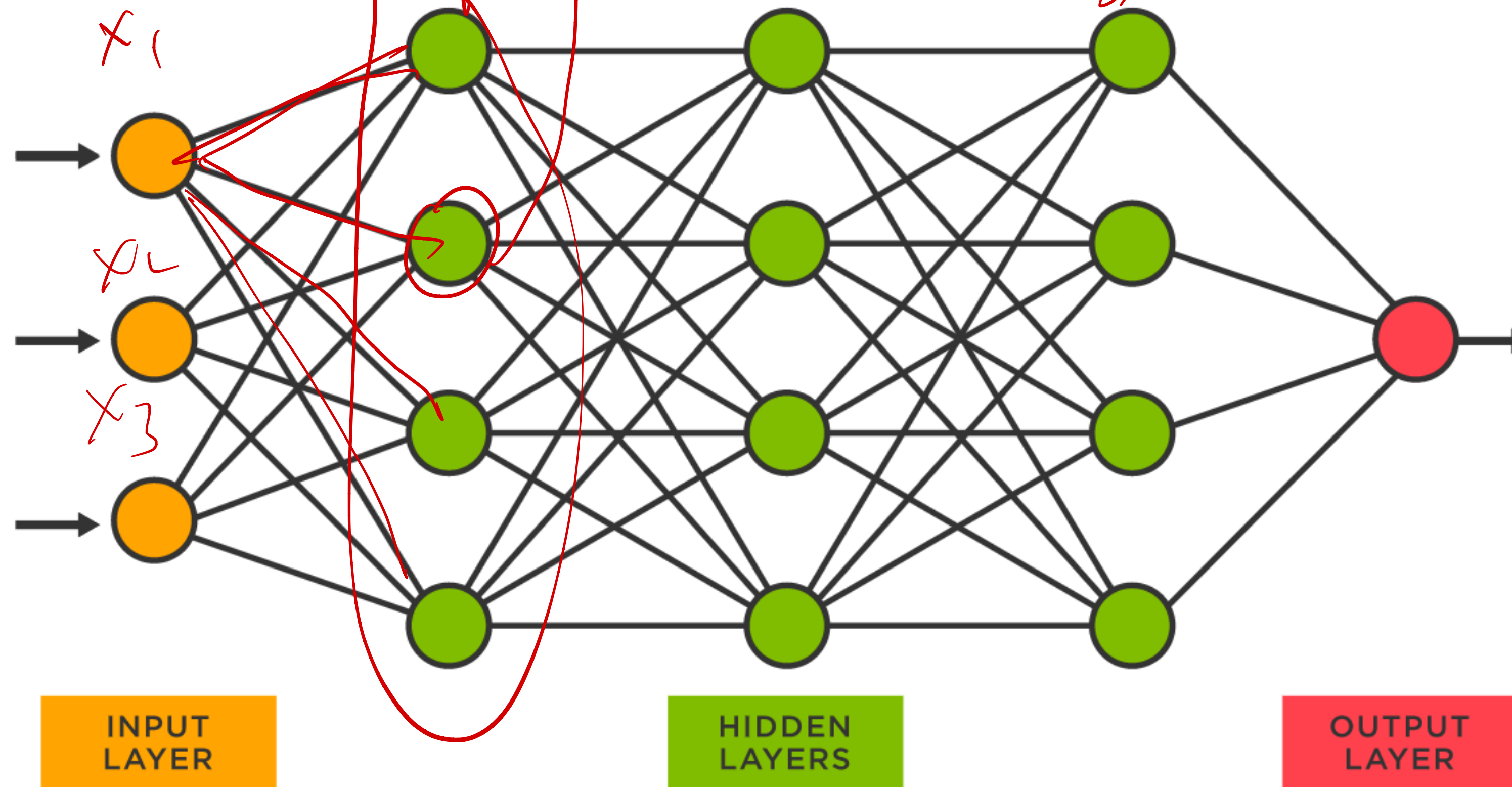# Multilayer Perceptron Neural Networks (MLP)

$\vec{x} = (x_1, x_2, x_3)$

$W\vec{x}$

$\sigma(W_1 x_1 + W_2 x_2 + W_3 x_3)$

$\sigma\sigma(W_1' x_1 + W_2' x_2 + W_3' x_3)$ fully connected

$x_1$

$x_2$

$x_3$

INPUT LAYER

HIDDEN LAYERS

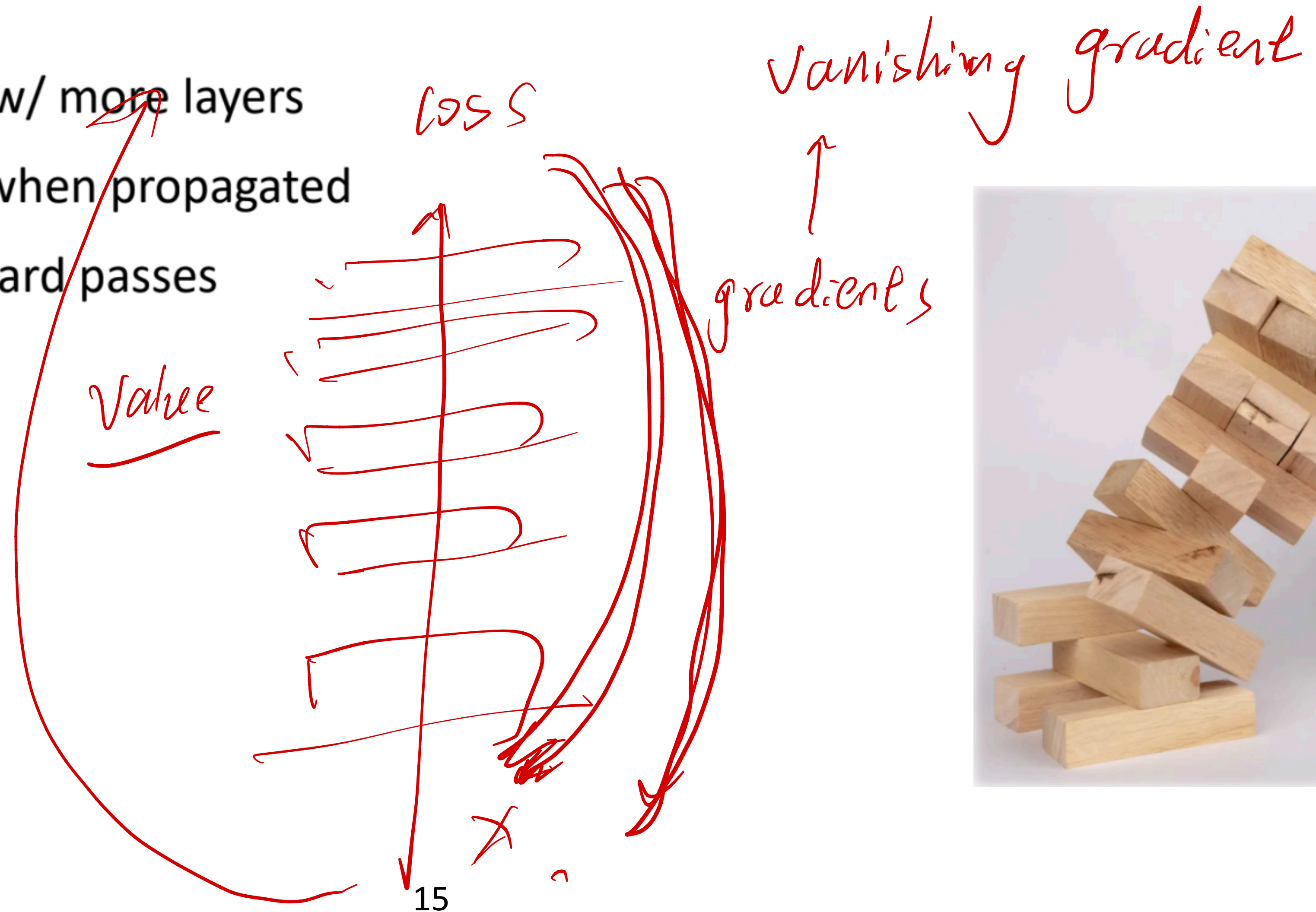OUTPUT LAYER

14

# Residual Connection

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

# 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



Loss

Vanishing gradient

gradients

Value

15

# 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
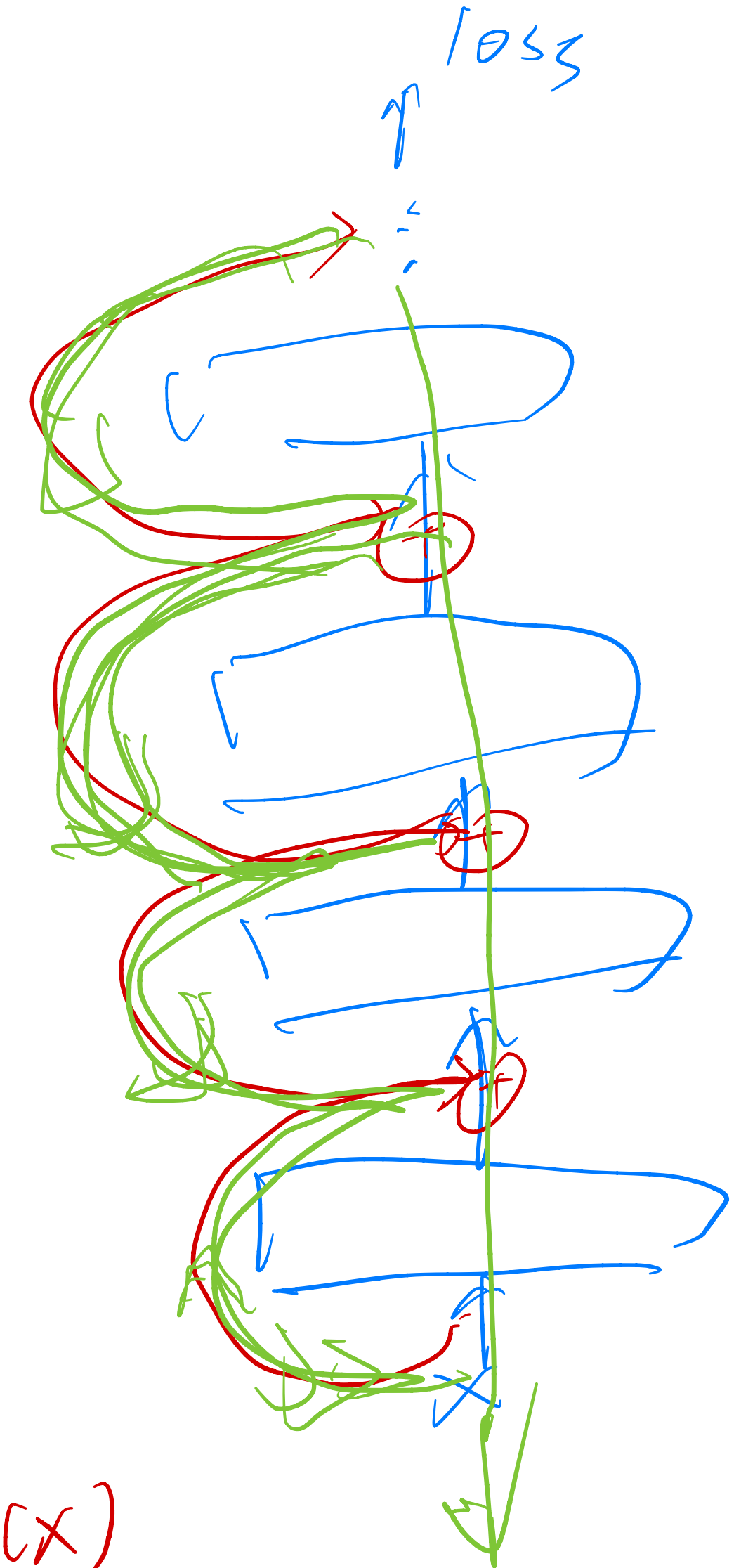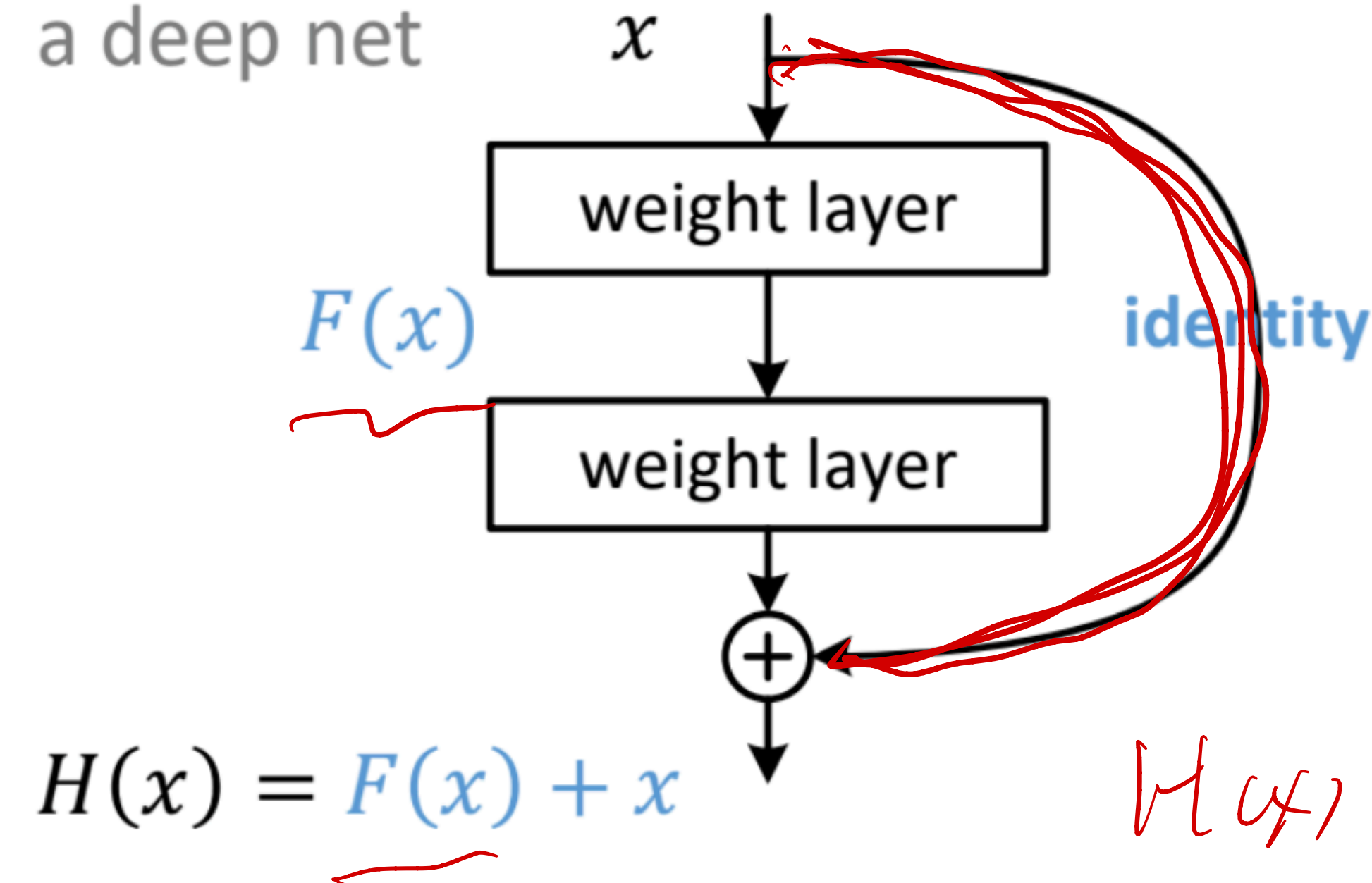
Normalization modules

Deep residual learning

# The Degradation Problem

- Good init + norm enable training deeper models
- Simply stacking more layers?

- Degrade after ~20 layers
- Not overfitting
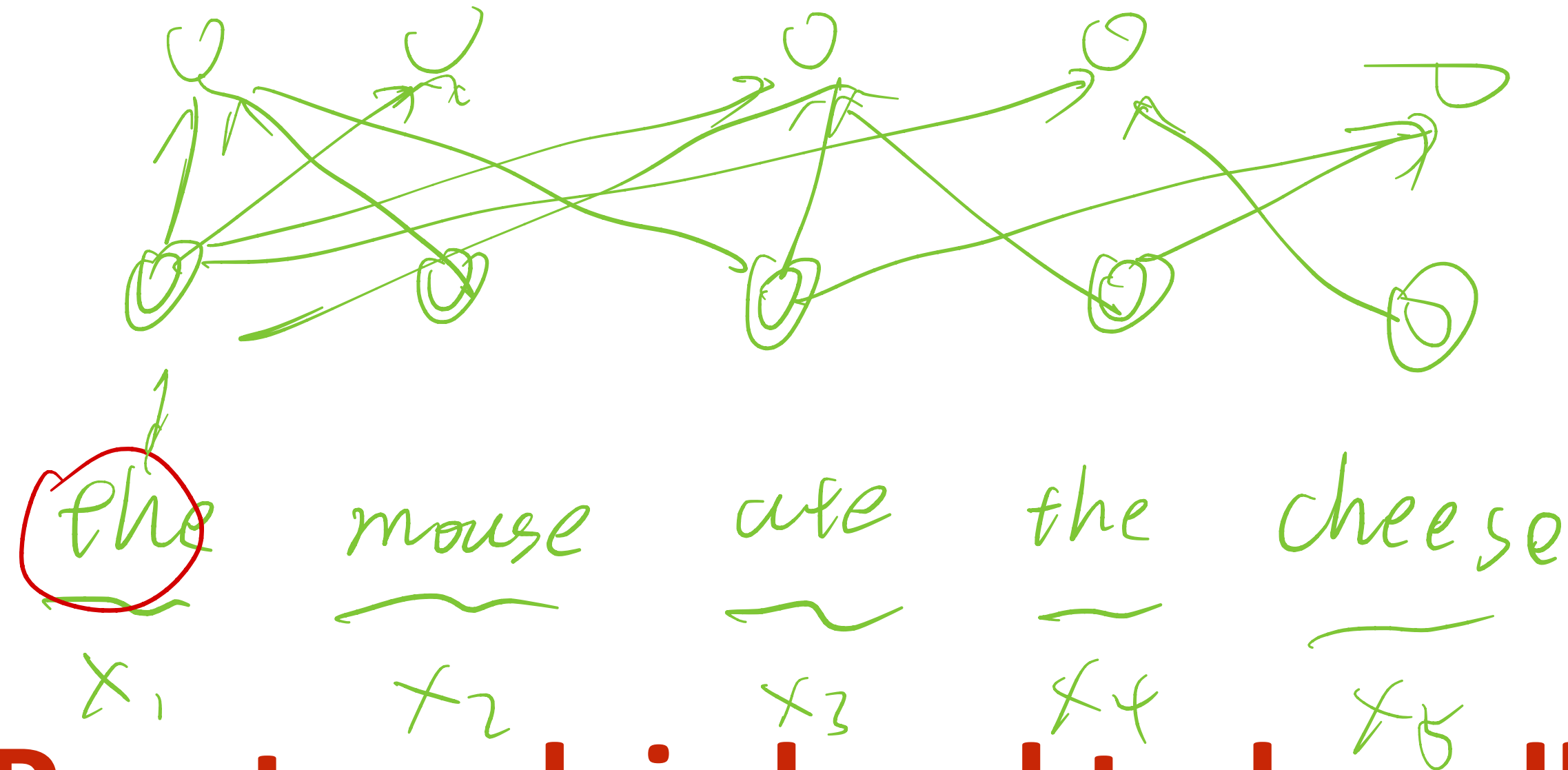- Difficult to train

# Deep Residual Learning

a subnet in
a deep net

$x$

| weight layer |

$F(x)$

identity

| weight layer |

$\bigoplus$

$H(x) = F(x) + x$

ResNet

$H(x) = F(x)$

&

$H(x) = F(x) + x$

loss

the $x_1$  mouse $x_2$  ate $x_3$  the $x_4$  cheese $x_5$

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

word emé

parameter size grows when sequence ⟶ longer

18

word embedding:

I
am
you
study
mouse

dictionary

the

vector
vector

lookup table

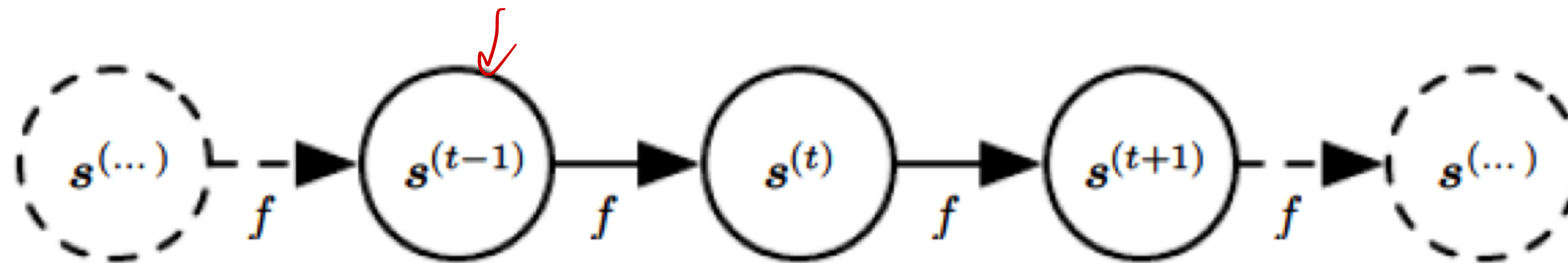$V \times E$  $\theta$

$V \times E$

embedding

dictionary

parameter

# 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

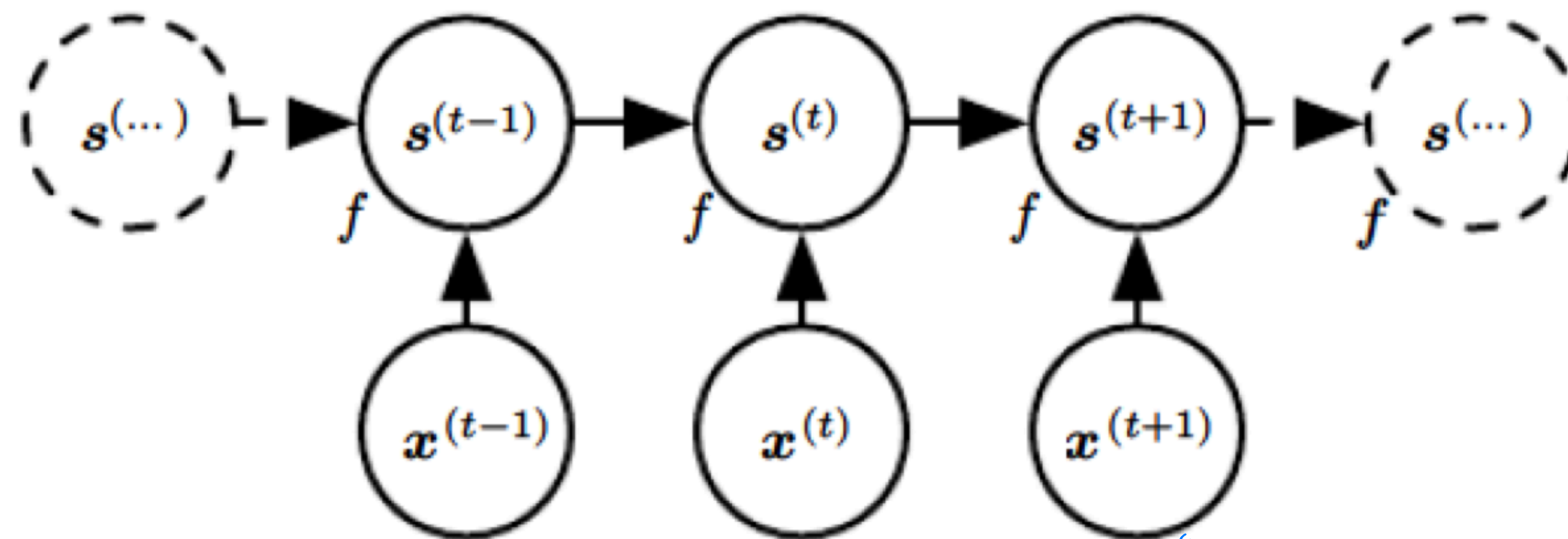- Especially, for natural language processing (NLP)

# Computation Graph



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

next timestep

next word

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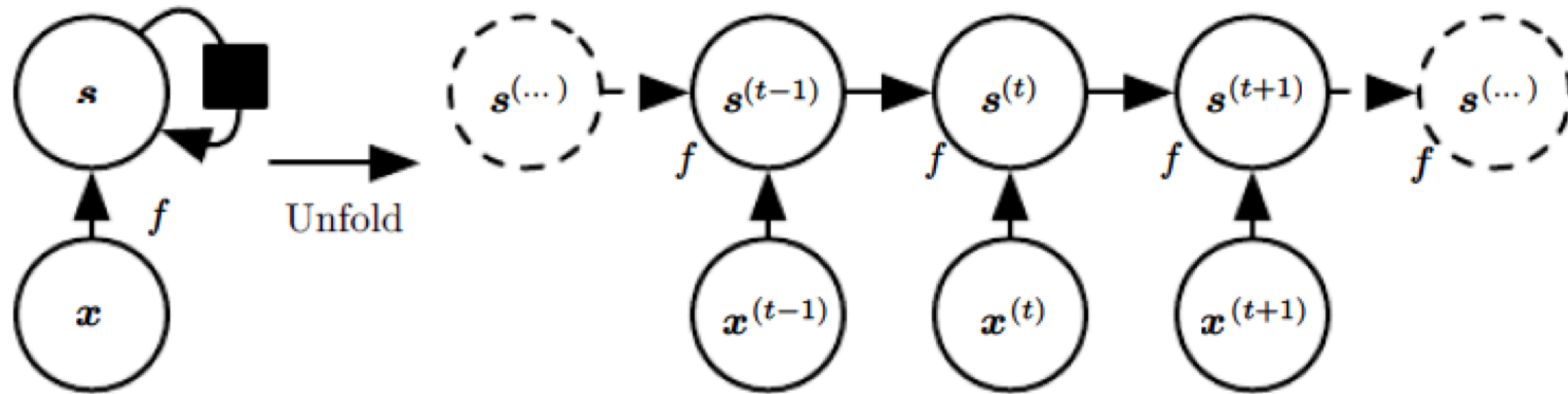Goodfellow, Bengio and Courville

$\theta$  parameter

# Computation Graph



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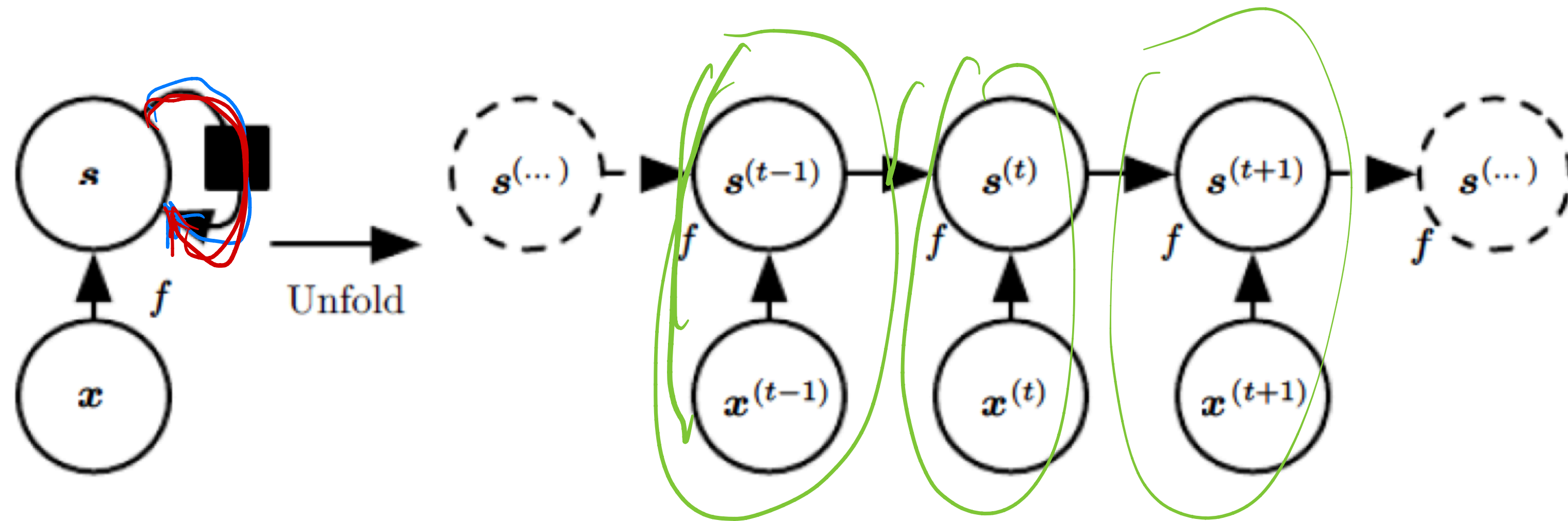
$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

MLP

# Compact view



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

# Compact view



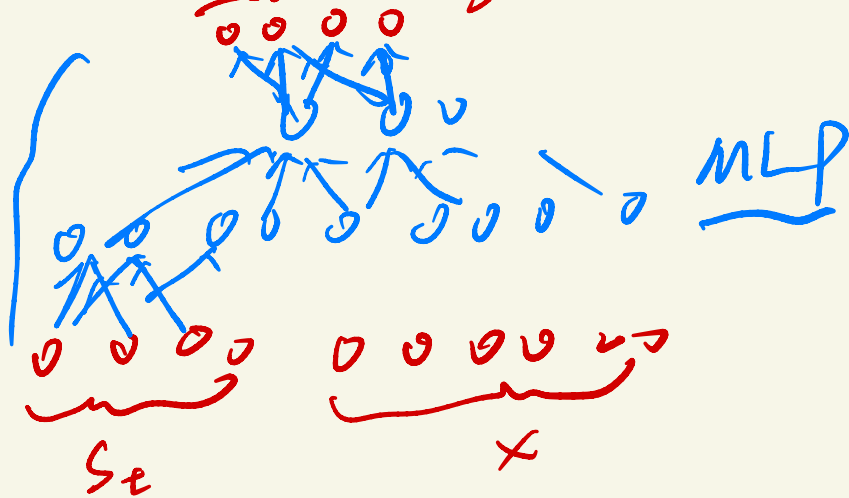$$s^{(t+1)} = f(s^{(t)}, x^{(t+1)}; \theta)$$

Key: the same $f$ and $\theta$ for all time steps

MLP

parameter size doesn't grow

23

$$S^{t+1} = f(S^t, X^t; \Theta)$$

weights

$\Theta$

$s^{t+1}$

MLP

$S_t$

$X$

# Recurrent Neural Networks

# Recurrent Neural Networks

- Use the same computational function and parameters across different time steps of the sequence

# 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

# 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

$P(\text{next word} \mid \text{context})$

# Recurrent Neural Networks



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

$\vec{S}$ dim $[v \ v \ v \ v \ v]$ $\overline{1024}$ ?

mouse

$P(\text{any word} | \text{the})$

$\text{softmax}([v] \cdot S)$ $[\circ$

$([\sqrt{} ] \cdot \vec{S})$ $[0,1] \longrightarrow$

$\text{dice site} \times 1024$

dim= dict size

?

vector of dice

$O$

size

$\boxed{NN} \longrightarrow h \xrightarrow{\quad V \quad} \vec{0} \xrightarrow{\quad softmax() \quad} \boxed{P(\text{next word})}$

goal

dict size

$\downarrow$

RNN

transformers

# Recurrent Neural Networks



Math formula:

$$a^{(t)} = b + Ws^{(t-1)} + Ux^{(t)}$$
$$s^{(t)} = \tanh(a^{(t)})$$
$$o^{(t)} = c + Vs^{(t)}$$
$$\hat{y}^{(t)} = \mathrm{softmax}(o^{(t)})$$

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

26

# Recurrent Neural Networks



Math formula:
$$a^{(t)} = b + Ws^{(t-1)} + Ux^{(t)}$$
$$s^{(t)} = \tanh(a^{(t)})$$
$$o^{(t)} = c + Vs^{(t)}$$
$$\hat{y}^{(t)} = \text{softmax}(o^{(t)})$$

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

mouse ate        the

shared

not shared

the   mouse ate the cheese

# Sequence-to-Sequence Learning

Example of Neural Machine Translation

# Sequence-to-Sequence Learning

Example of Neural Machine Translation

# Sequence-to-Sequence Learning

Example of Neural Machine Translation

# Sequence-to-Sequence Learning

Example of Neural Machine Translation

LM

Encoding of the source sentence.
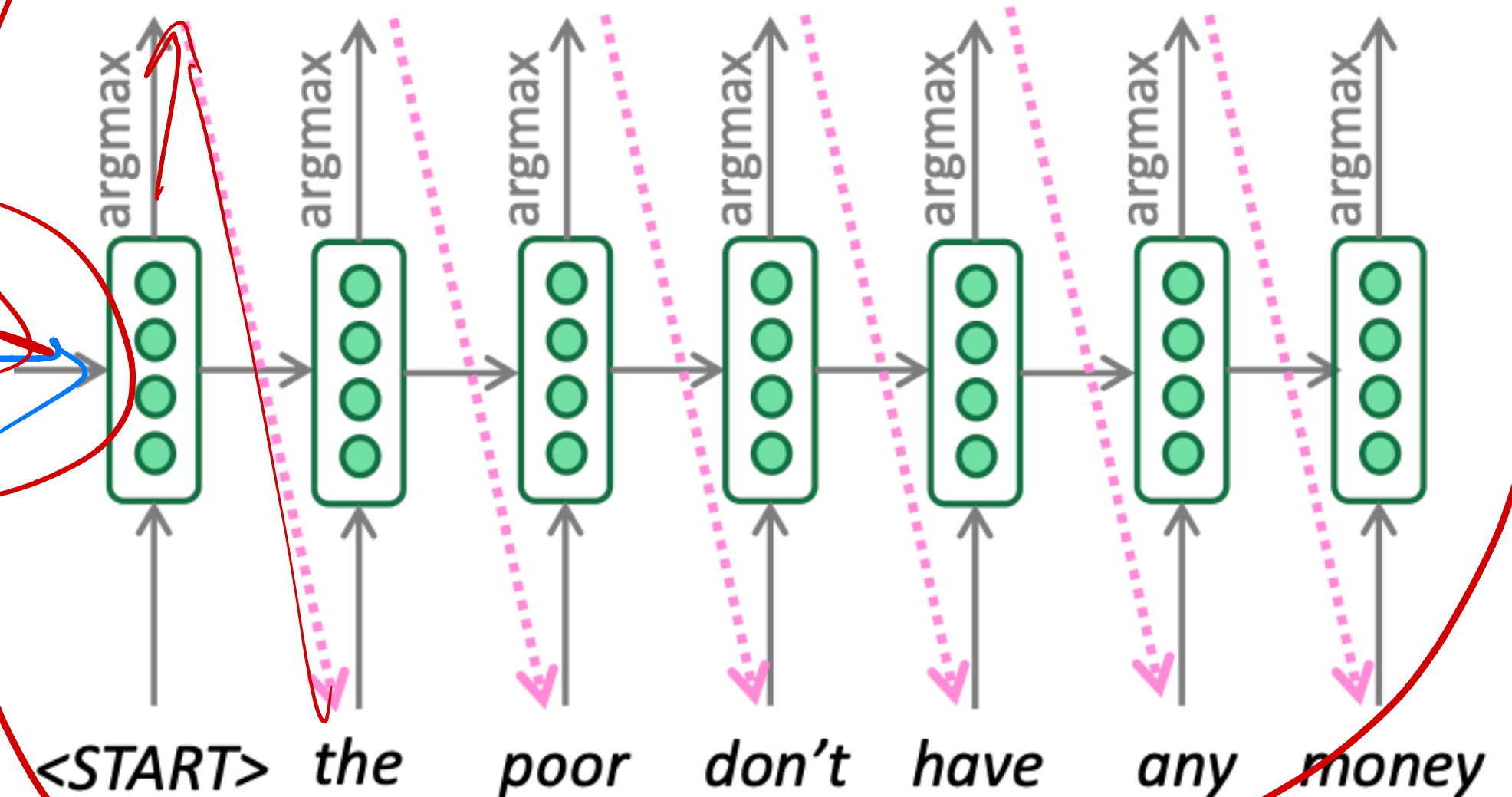Provides initial hidden state
for Decoder RNN.

one RNN



Encoder RNN

les   pauvres   sont   démunis

Source sentence (input)

Encoder RNN produces
an encoding of the
source sentence.

Target sentence (output)

the   poor   don't   have   any   money   <END>

argmax

<START>   the   poor   don't   have   any   money

Decoder RNN

# Sequence-to-Sequence Learning

Example of Neural Machine Translation



Encoding of the source sentence. Provides initial hidden state for Decoder RNN.

Encoder RNN

Target sentence (output)

the    poor    don't    have    any    money    <END>

argmax

Decoder RNN

les    pauvres    sont    démunis

<START>    the    poor    don't    have    any    money

Source sentence (input)

Encoder RNN produces an encoding of the source sentence.

Decoder RNN is a Language Model that generates target sentence conditioned on encoding.

27

Same model

les pauvres --- < stores the poor ---

# RNN Language Model

# RNN Language Model

*parallel computation*

*sequential*

Target sentence (output)



the    poor    don't    have    any    money    <END>

<START>    the    poor    don't    have    any    money

Decoder RNN

$\int_A^B$

sequential

28

# Transformer



Output
Probabilities

Softmax

Linear

Add & Norm

Feed
Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Masked
Multi-Head
Attention

Add & Norm

Feed
Forward

N×

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

Inputs

Positional
Encoding

Output
Embedding

Outputs
(shifted right)

*encoder*

*decoder*

*machine translation*

Vaswani et al. Attention is All You Need. NeurIPS 2017.

29

# Encoder

# Decoder



Output Probabilities

*encoder*

*decoder*

LM

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

N×

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

target

31

# Transformer Encoder



add    Norm

fully connected

residual

Add & Norm

Feed
Forward

Nx

Add & Norm

Multi-Head
Attention

Positional
Encoding

Input
Embedding

32

# Transformer Encoder



Self-attention

# Transformer Encoder



MLP

Self-attention

32

# Transformer Encoder



Residual connection

MLP

Self-attention

# What is Attention

Scaled Dot-Product Attention



Q: Query
K: key
V: value

attention

attend

We are from CS deparment, and we are taking class

attention weight

RNN

attention weight dot product

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \cdot \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum x_i y_i$$

query (class) $\cdot$ key (CS) $\longrightarrow$ Weight

How much attend

query
key
value

query
key
value

query
key
value

We are from CS department and we are taking class

key (department)

[weight] $\in$ [0, 1]

weight ( class, cs )

weigh ( class, depures )

weight ( class, we )

)

;

;

weight $[0, 1)$

probabily "class"

attend "cs)

Softmax ( ) = attention weight

influence

atln weight for each word

$w_i$

~~ctos~~

$\sum_i w_i v_i$ $\xrightarrow{\text{influence}}$

$\alpha 9 \times \overrightarrow{value}$

CS

class

# What is Attention

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

Scaled Dot-Product Attention



Q: Query
K: key
V: value

# What is Attention

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

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

Scaled Dot-Product Attention



Q: Query
K: key
V: value

# What is Attention

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

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

## Scaled Dot-Product Attention



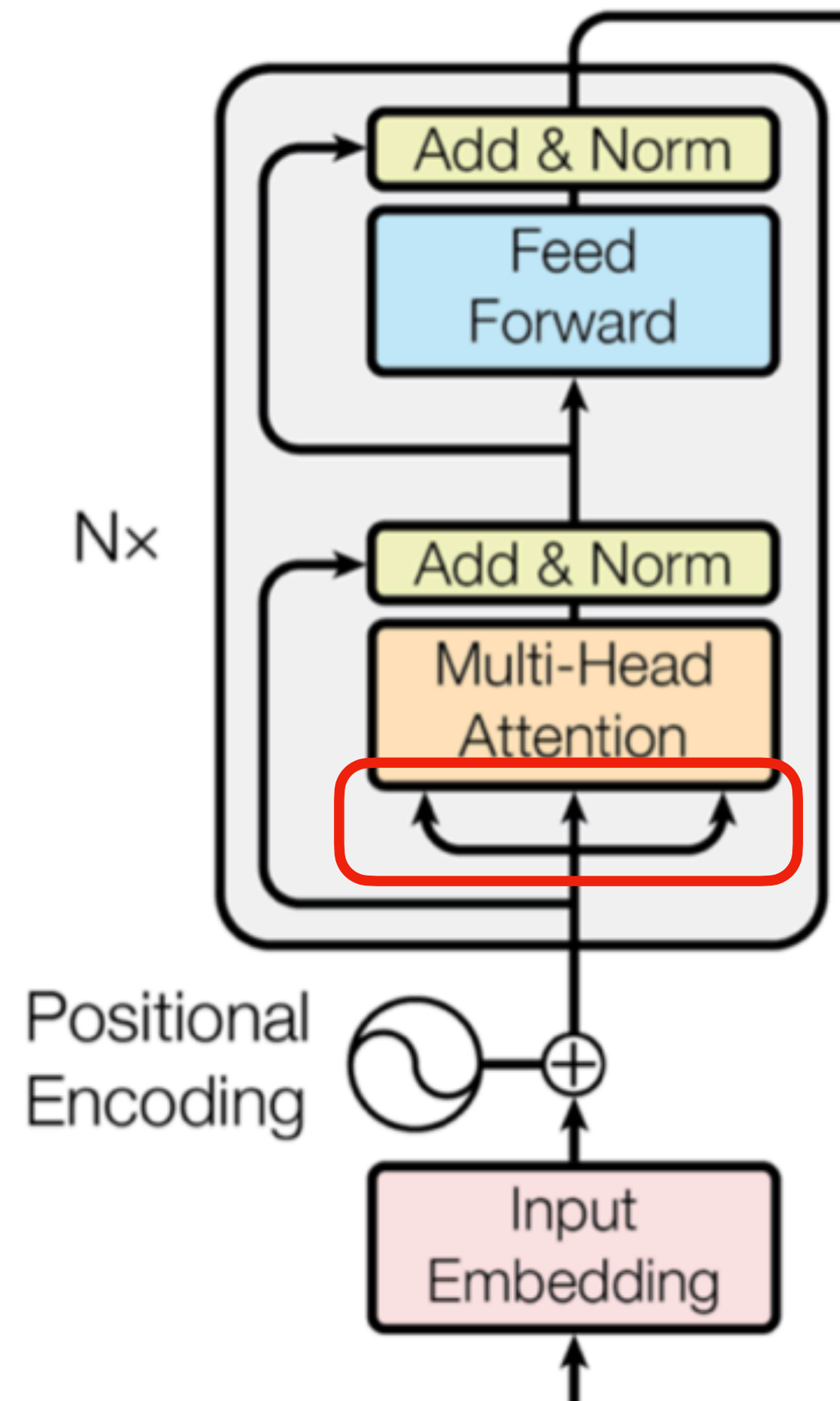Q: Query
K: key
V: value

Attention weight = softmax($QK^T$)

# What is Attention

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

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

Scaled Dot-Product Attention



Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Q: Query
K: key
V: value

# What is Attention

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

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

## Scaled Dot-Product Attention



Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Scaled Attention weight = softmax($\dfrac{QK^T}{\sqrt{d_k}}$)

Q: Query
K: key
V: value

33

# What is Attention

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

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

Scaled Dot-Product Attention



Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Scaled Attention weight = softmax($\dfrac{QK^T}{\sqrt{d_k}}$)     Shape is mxn

Q: Query
K: key
V: value

33

# What is Attention

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

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

Scaled Dot-Product Attention



Attention weight = softmax($QK^T$)

Dot-products grow large in magnitude

Scaled Attention weight = softmax($\frac{QK^T}{\sqrt{d_k}}$)  Shape is mxn

Attention weight represents the strength to "attend" values V

Q: Query
K: key
V: value

33

# What is Attention

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

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

Scaled Dot-Product Attention



Attention weight = softmax$(QK^T)$

Dot-products grow large in magnitude

Scaled Attention weight = softmax$(\dfrac{QK^T}{\sqrt{d_k}})$    Shape is mxn

Attention weight represents the strength to "attend" values V

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

Q: Query
K: key
V: value

33

# Q, K, V



What are Q, K, V in the transformer

# Self-Attention



Input is X

Jay Alammar. The Illustrated Transformer.

# Self-Attention



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

Input is X

35

Jay Alammar. The Illustrated Transformer.

# Self-Attention



Input is X

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

Jay Alammar. The Illustrated Transformer.

# Self-Attention

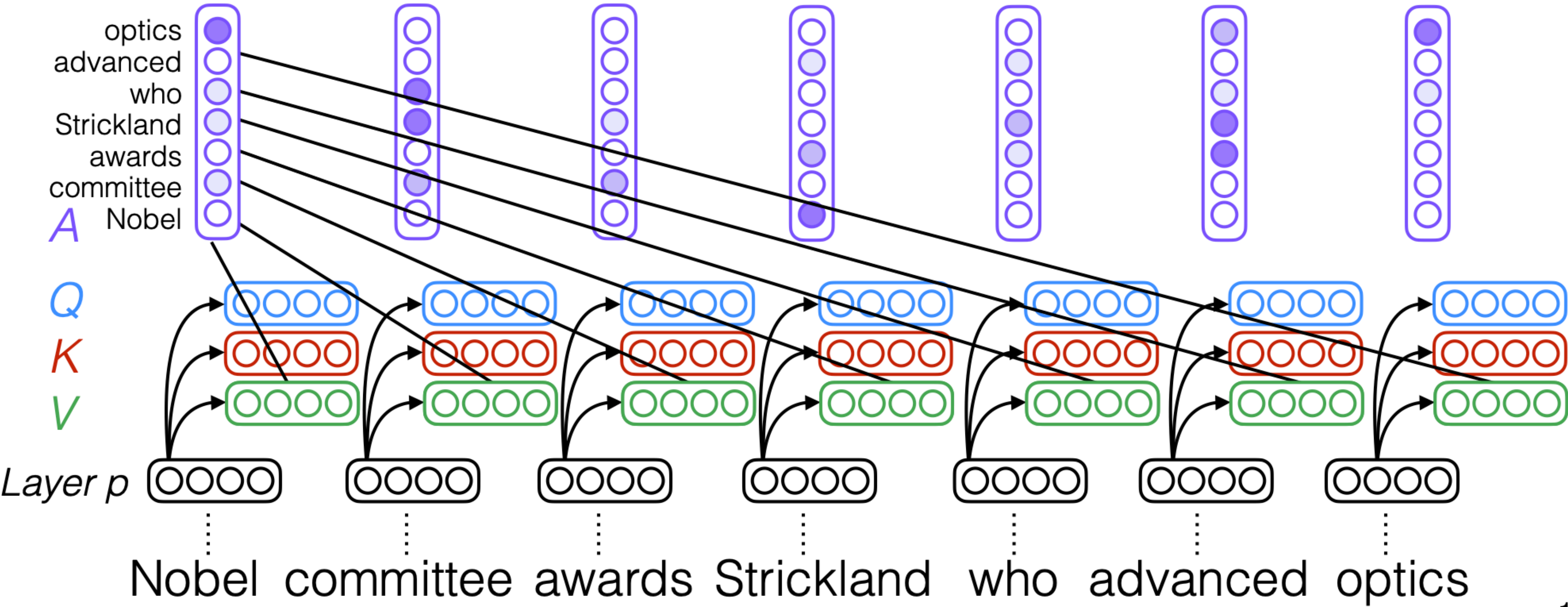# Self-Attention

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



36

# Self-Attention

# Self-Attention



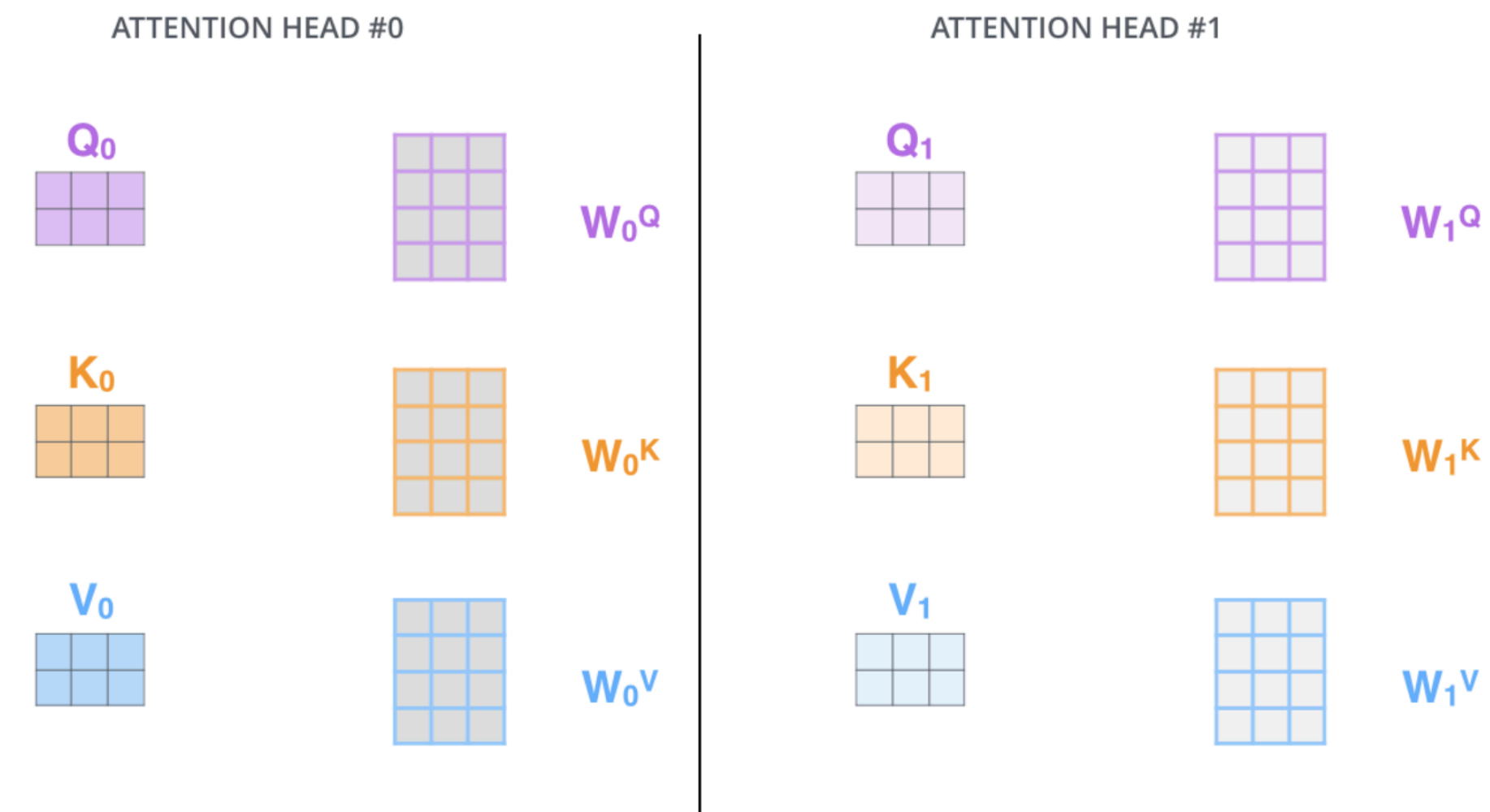Attention weight on every word in the sequence

optics
advanced
who
Strickland
awards
committee
Nobel

$Q$
$K$
$V$

*Layer p*

Nobel  committee  awards  Strickland  who  advanced  optics

# Self-Attention

# Self-Attention

# Multi-Head Attention
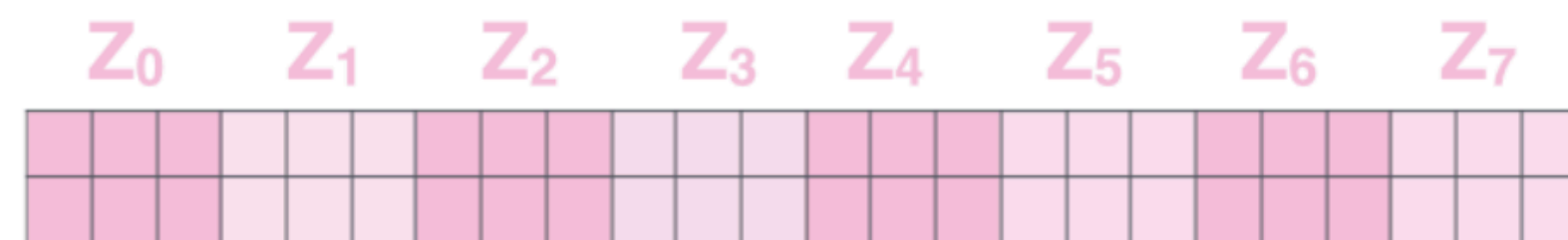


Multi-Head Attention

# Multi-Head Self-Attention

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

$Q_0$      $W_0^Q$      $Q_1$      $W_1^Q$

$K_0$      $W_0^K$      $K_1$      $W_1^K$

$V_0$      $W_0^V$      $V_1$      $W_1^V$

42

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

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

X

$W^O$

Jay Alammar. The Illustrated Transformer.

# Multi-Head Self-Attention

1) Concatenate all the attention heads

$Z_0$  $Z_1$  $Z_2$  $Z_3$  $Z_4$  $Z_5$  $Z_6$  $Z_7$

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

X

$W^O$

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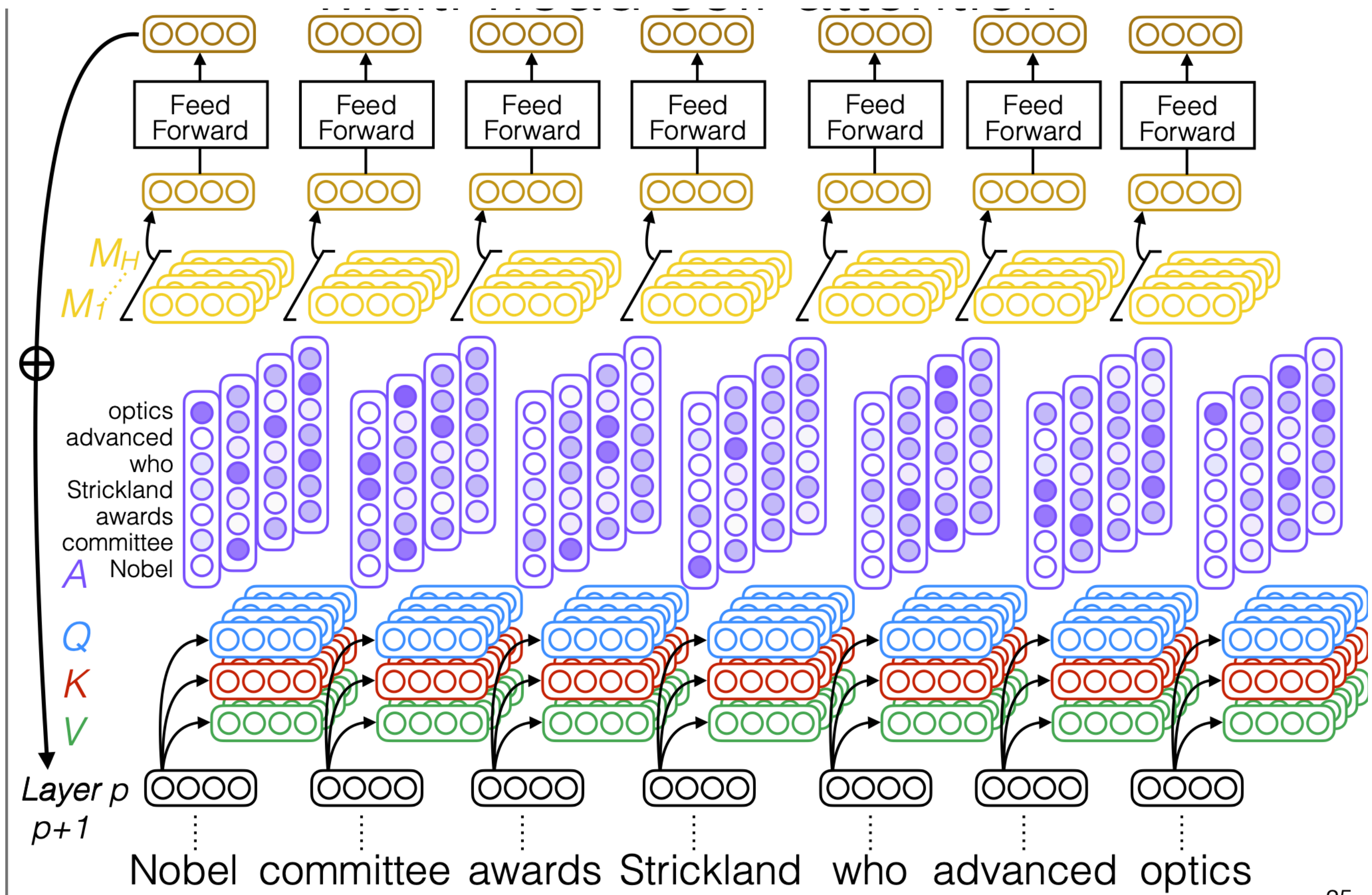
Z

=

Jay Alammar. The Illustrated Transformer.

# Multi-head Self-Attention

# Multi-head Self-Attention
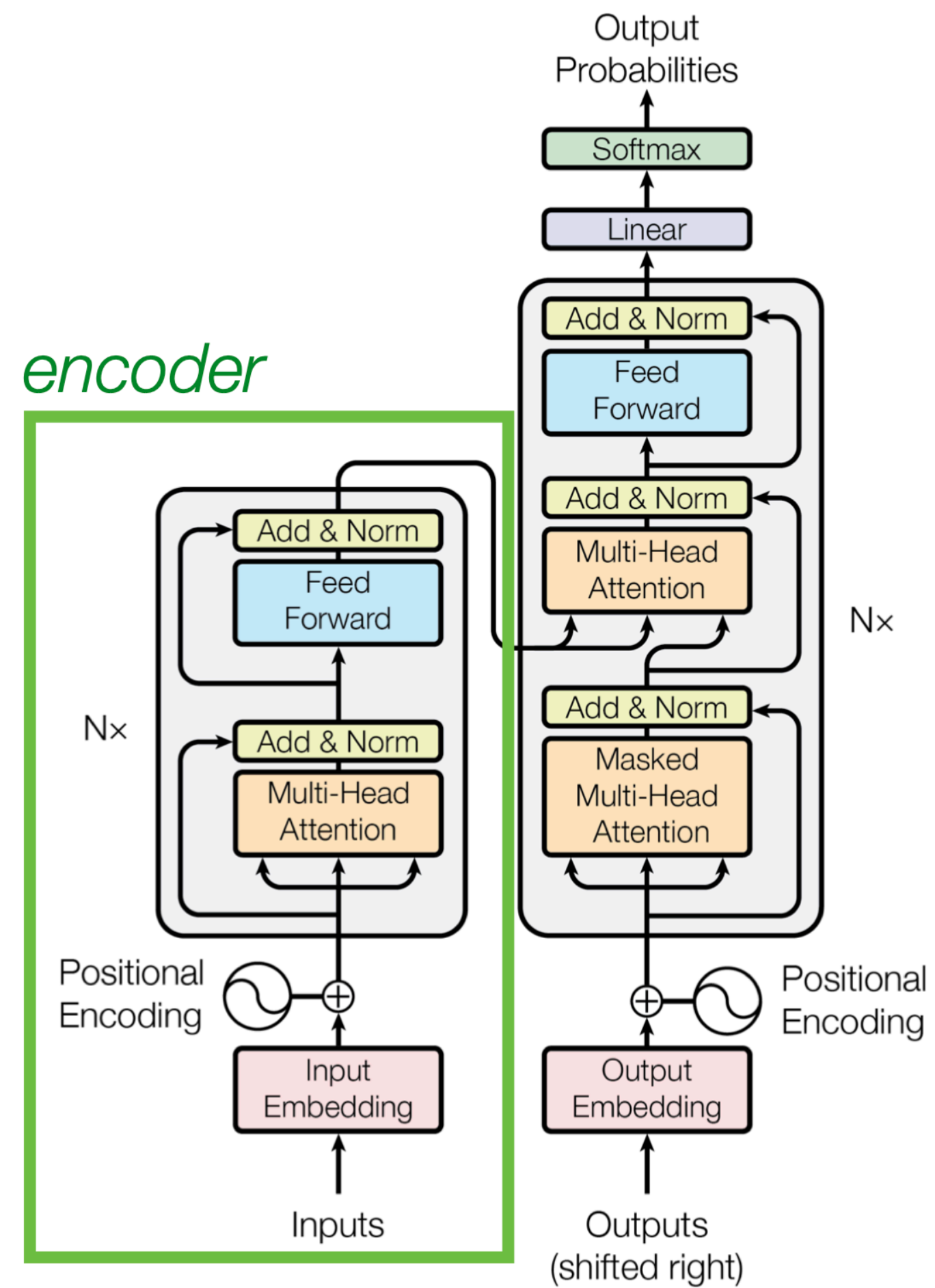


Concat and output projection

45

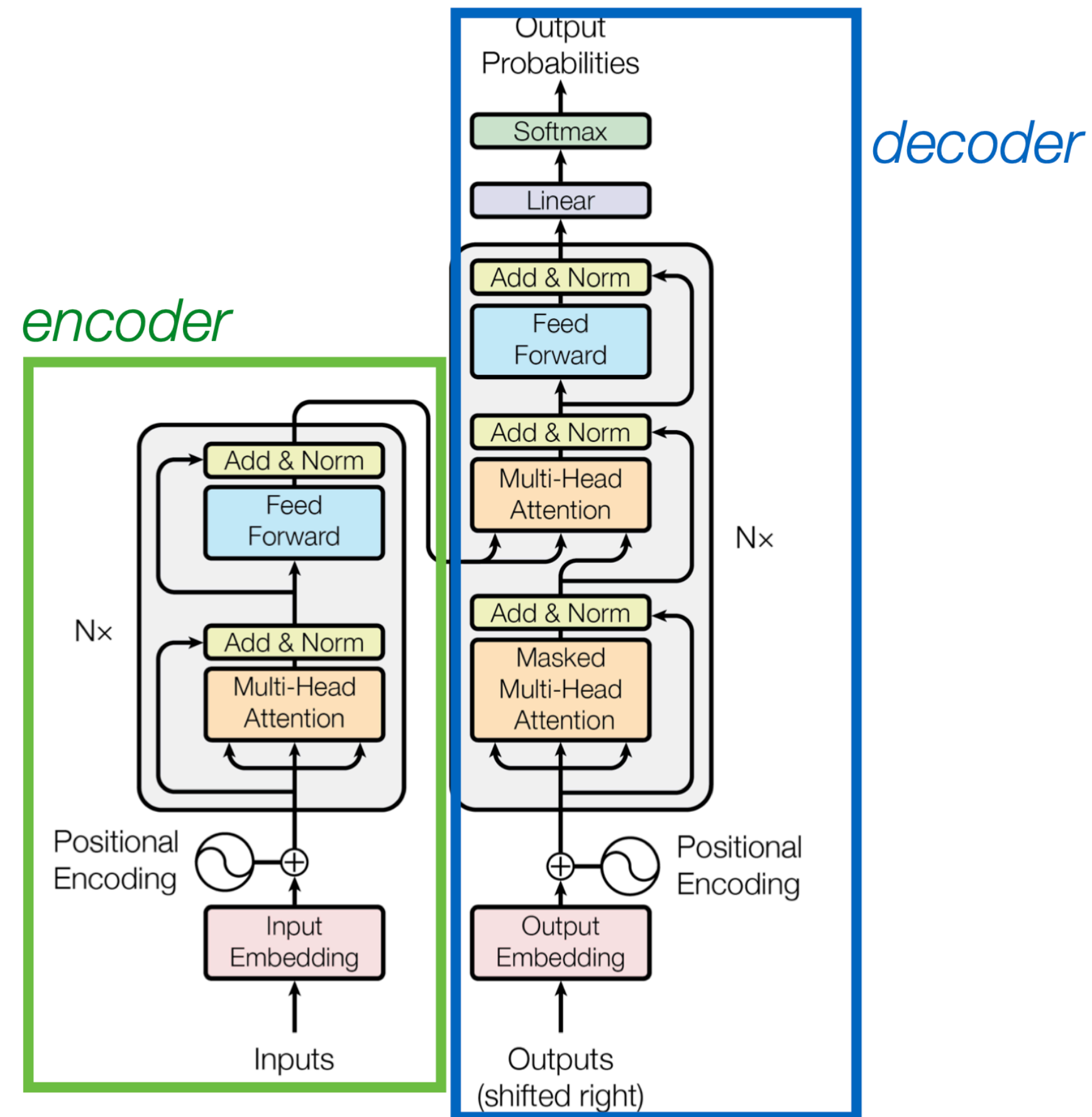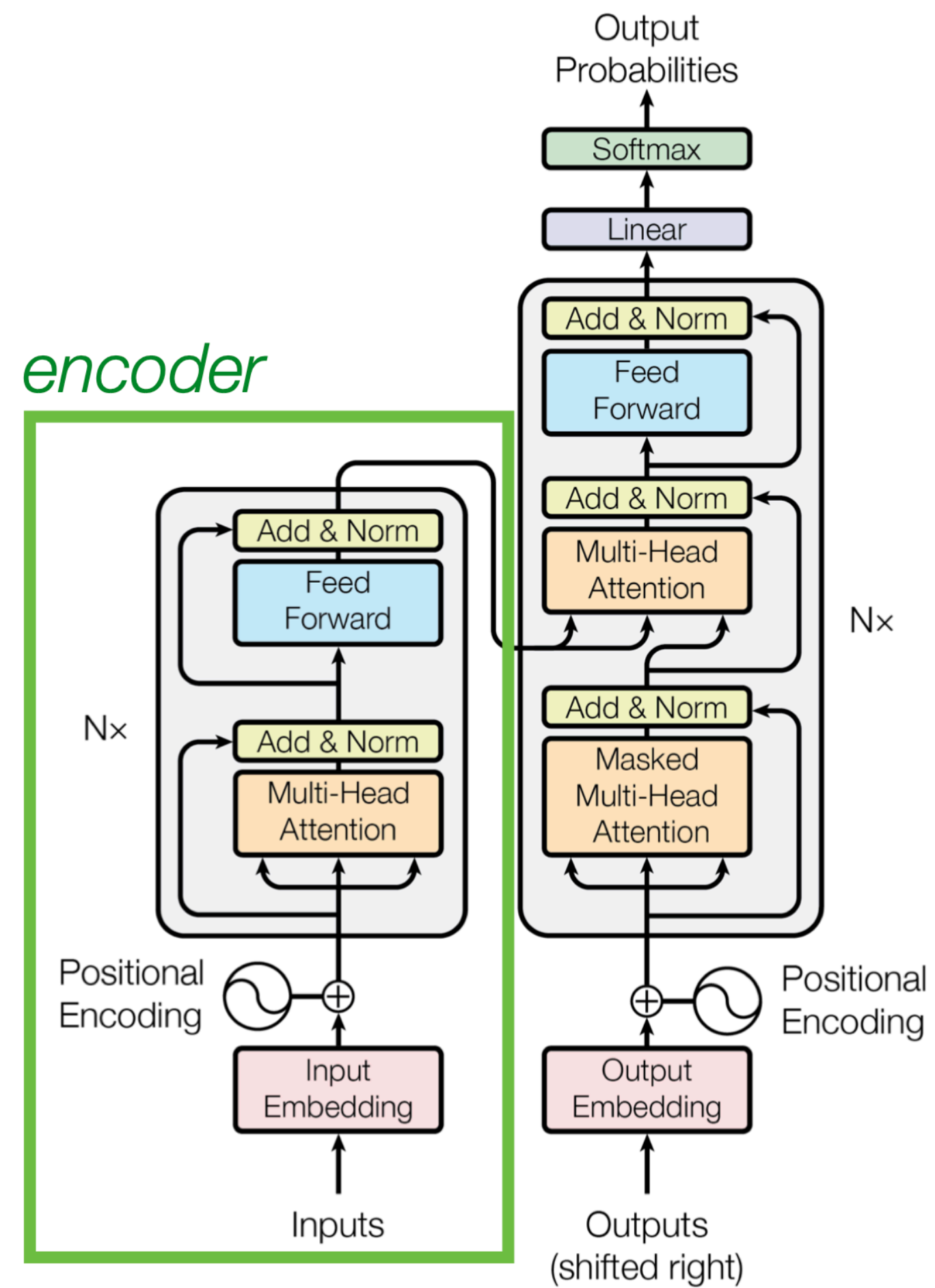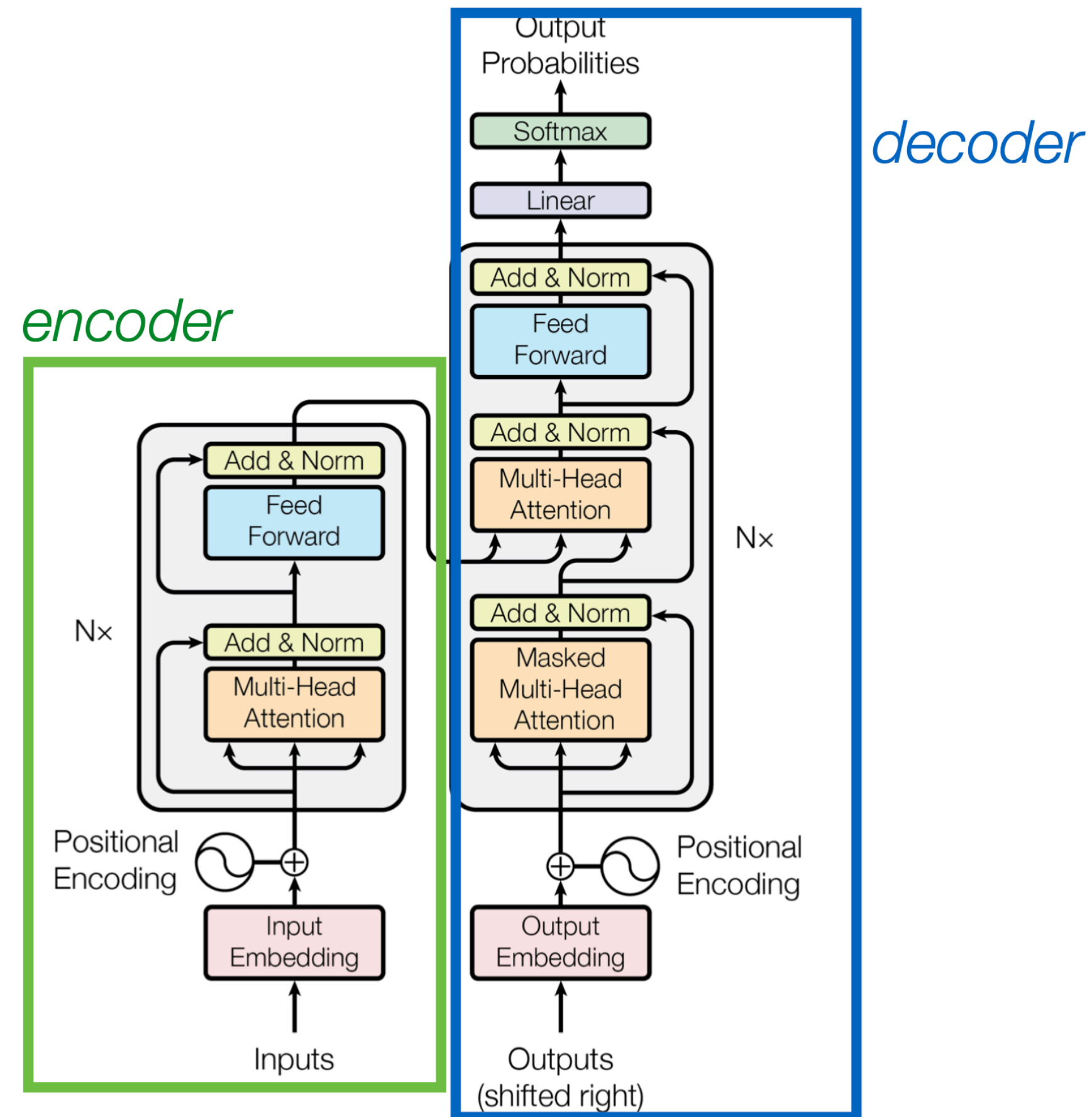# Multi-head Self-Attention + FFN



46

# Transformer Encoder

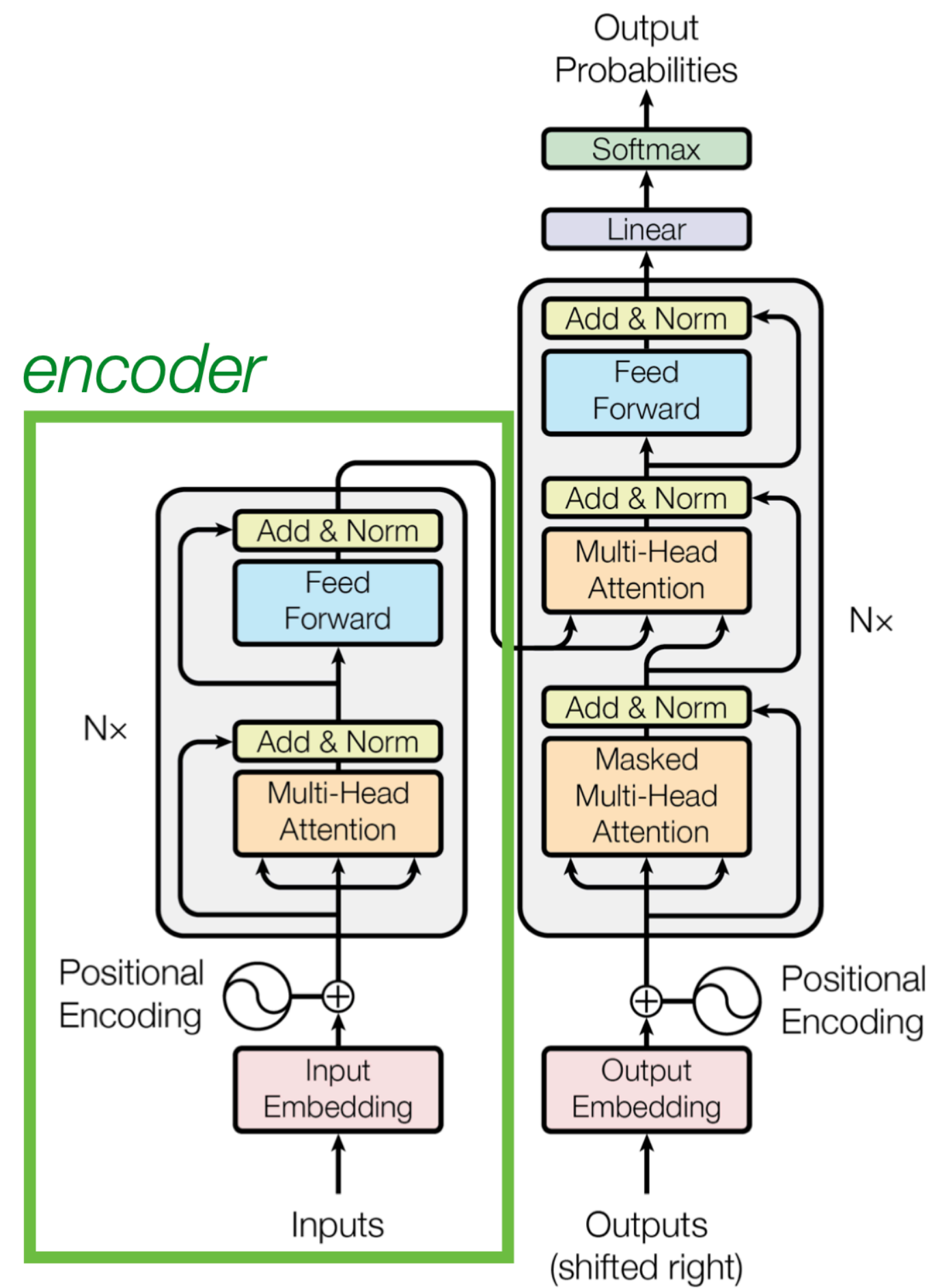Currently we only cover the encoder side

# Transformer Encoder

Currently we only cover the encoder side
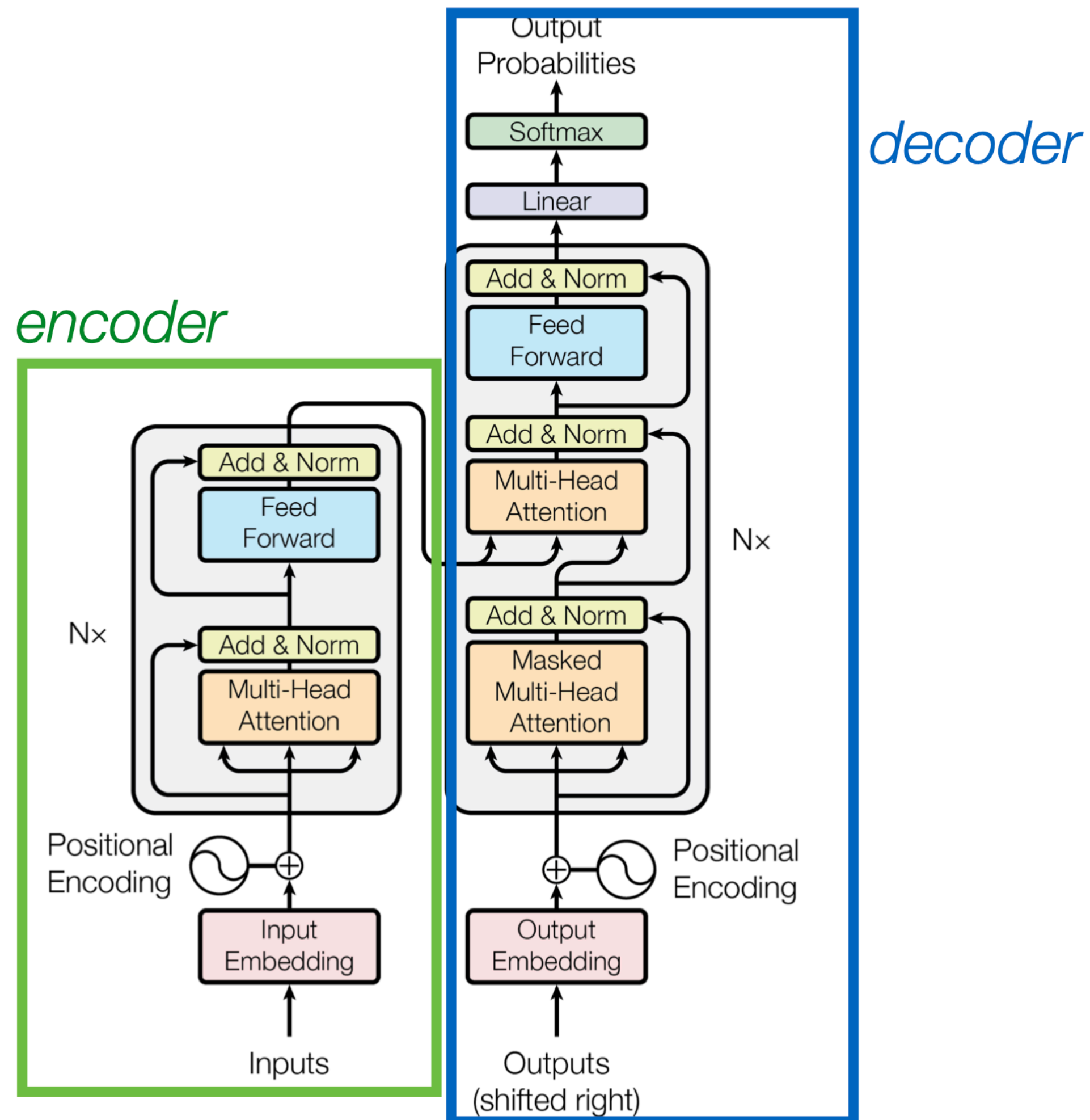


47

# Transformer Encoder
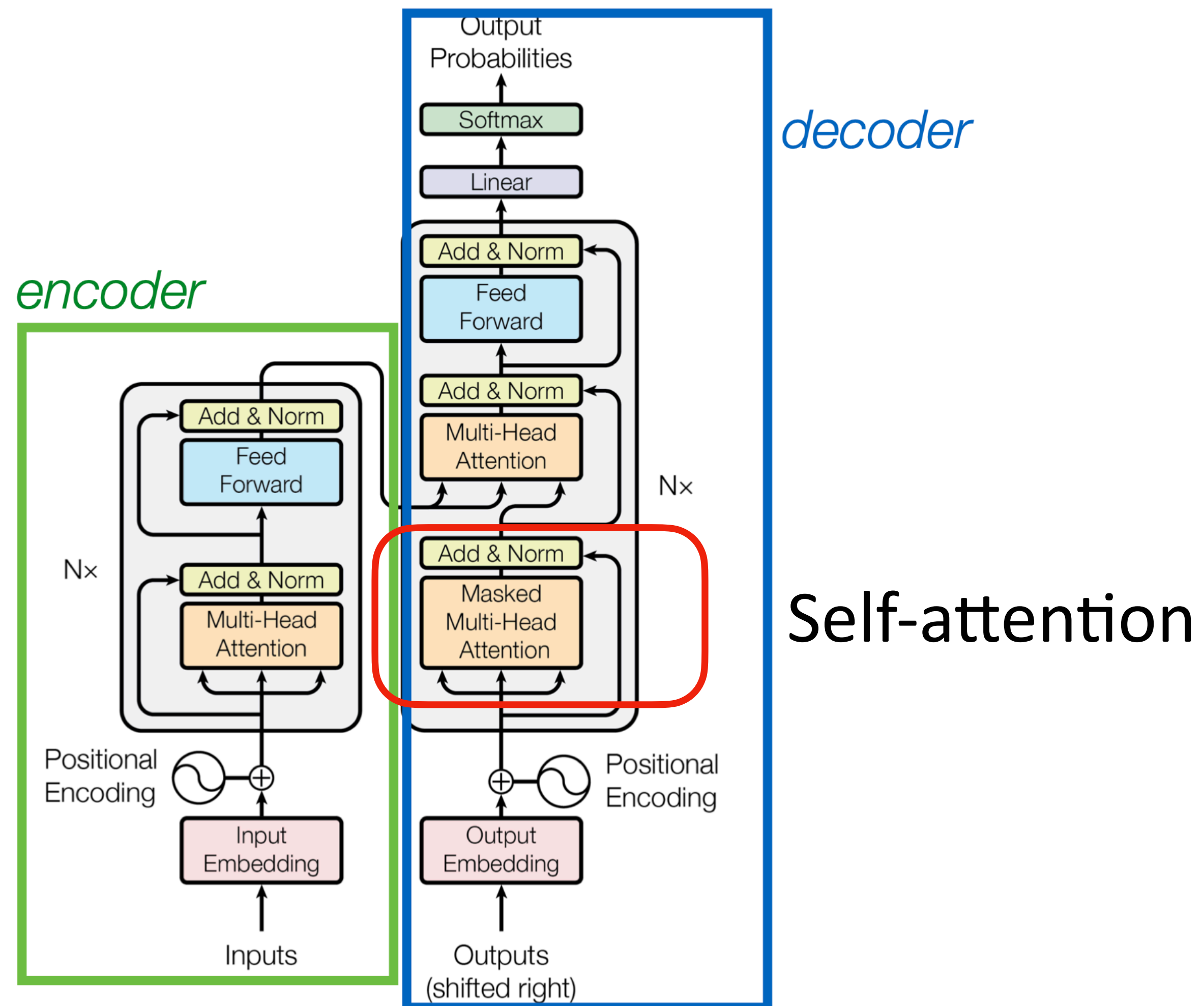
Currently we only cover the encoder side



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
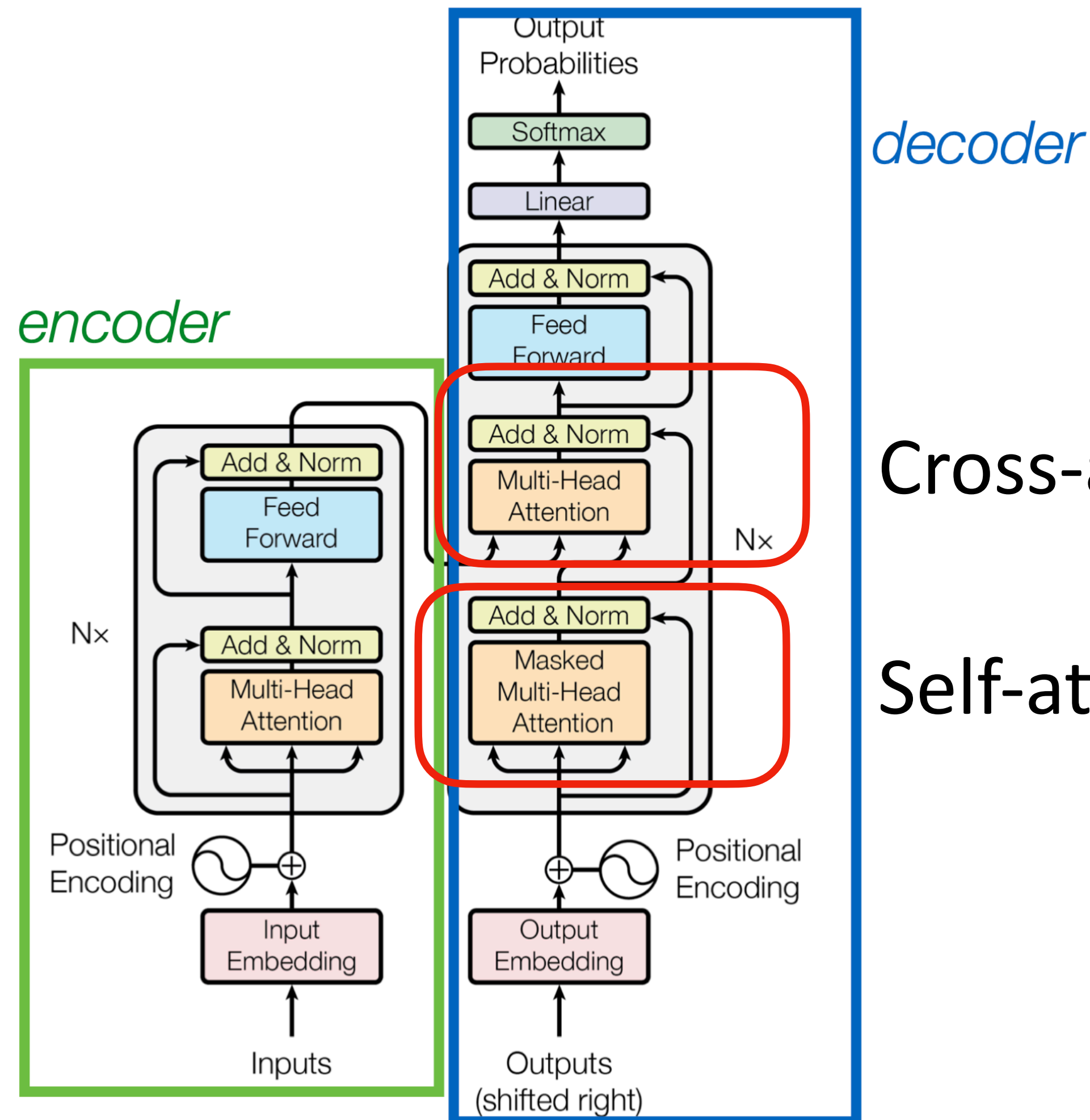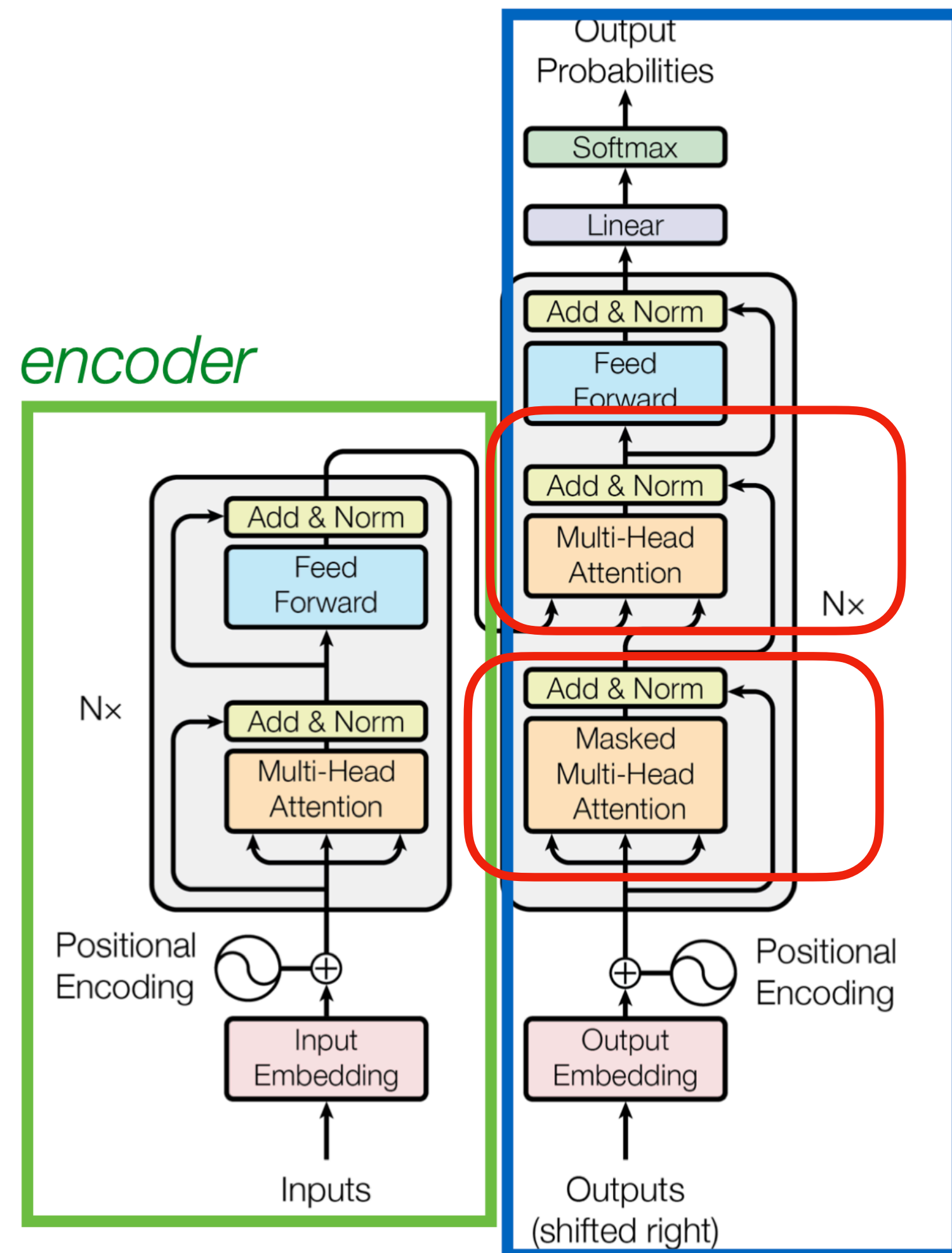
# Transformer Decoder in Seq2Seq

# Transformer Decoder in Seq2Seq



*encoder*

*decoder*

Self-attention

48

# Transformer Decoder in Seq2Seq

# Transformer Decoder in Seq2Seq



*decoder*

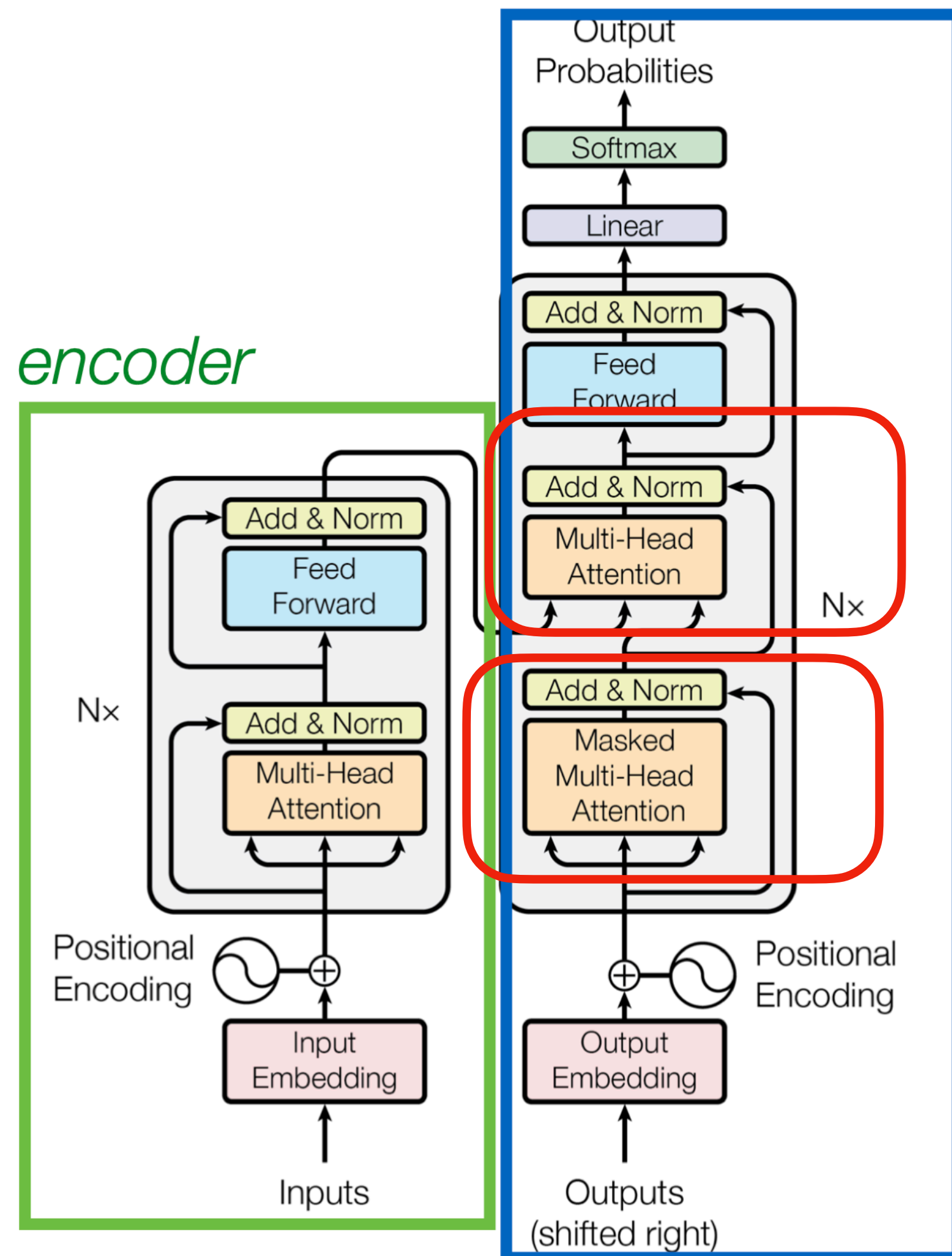*encoder*

Cross-attention

Self-attention
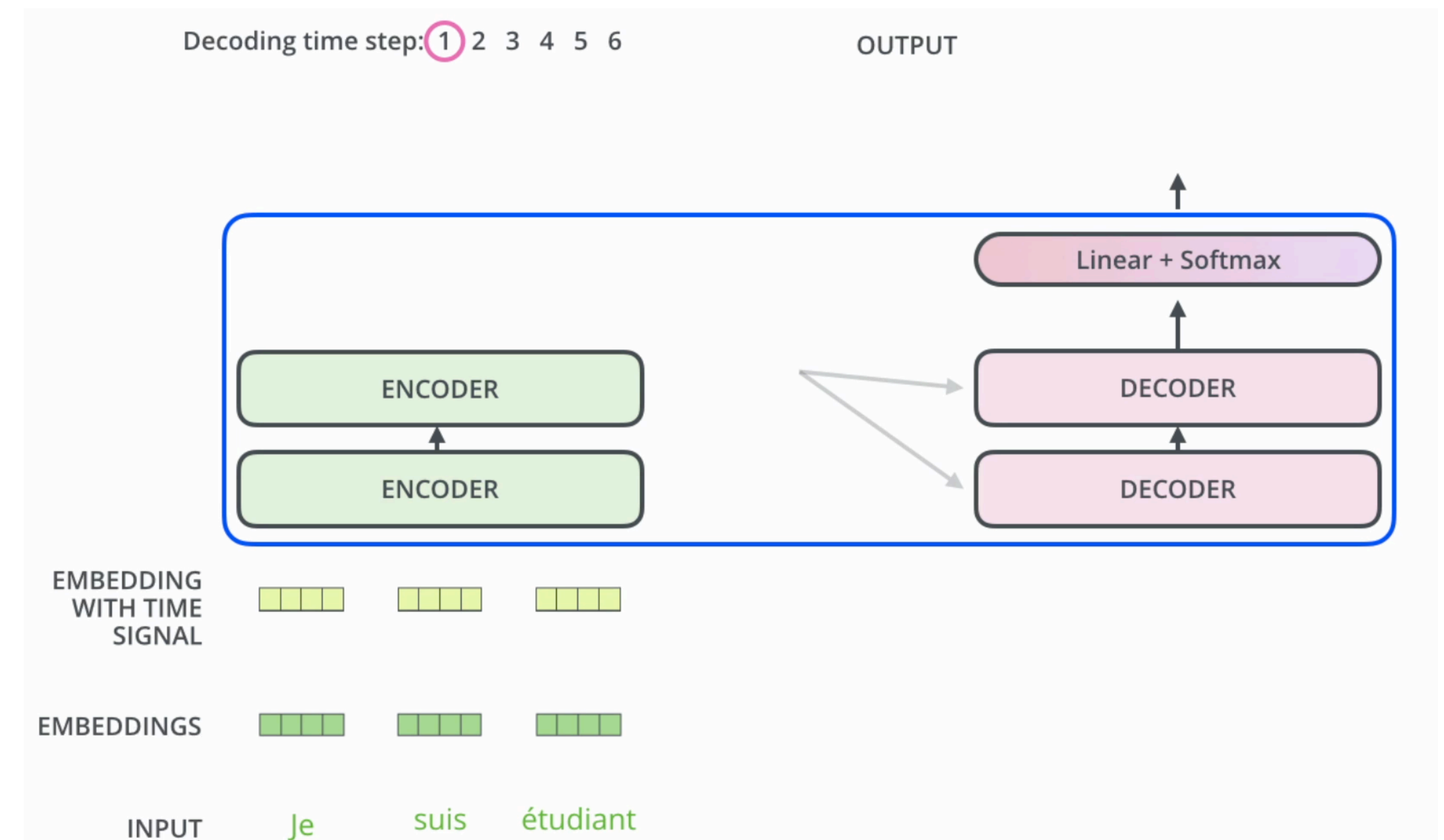
Cross-attention uses the output of encoder as input

48

# Transformer Decoder in Seq2Seq



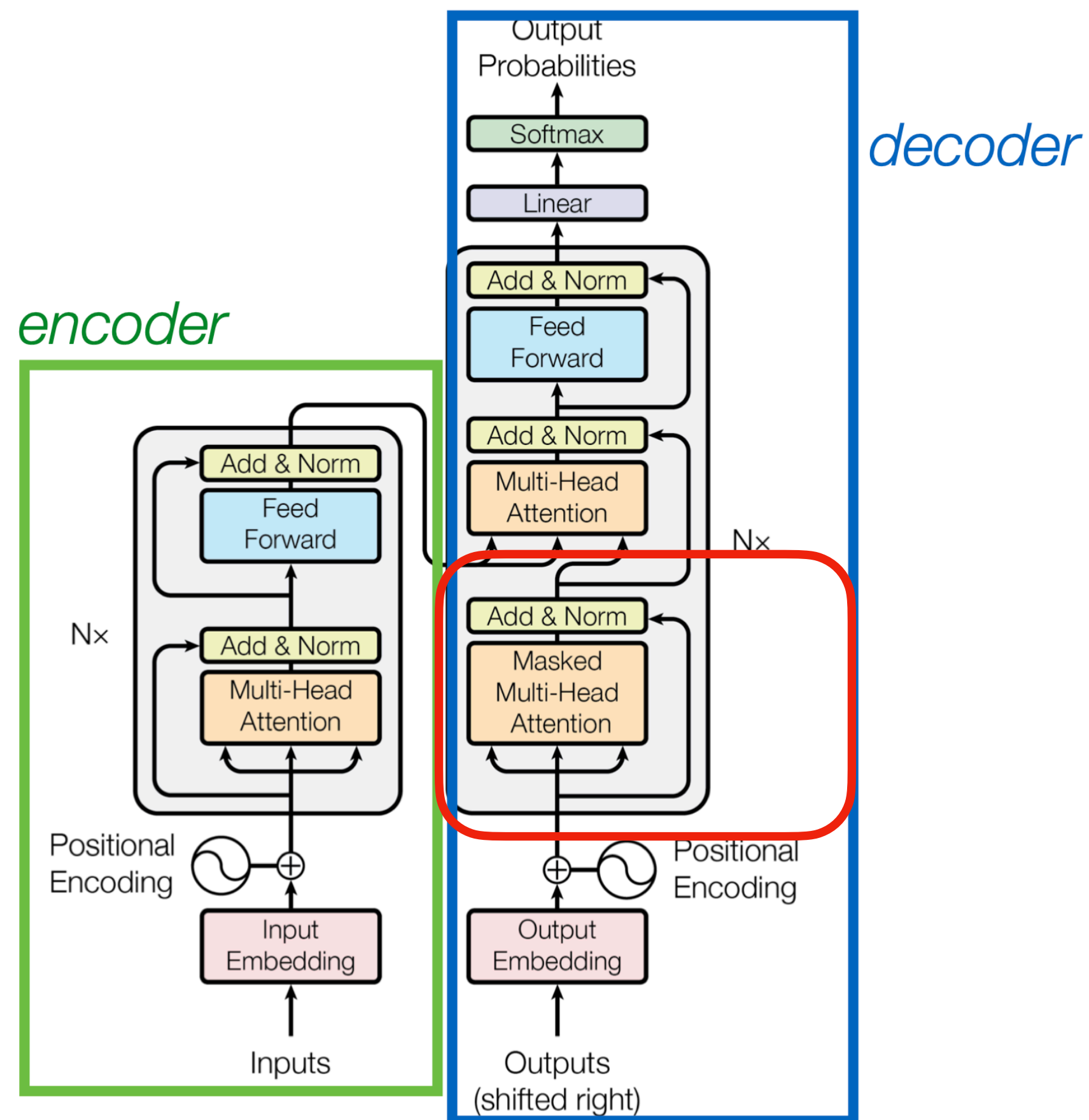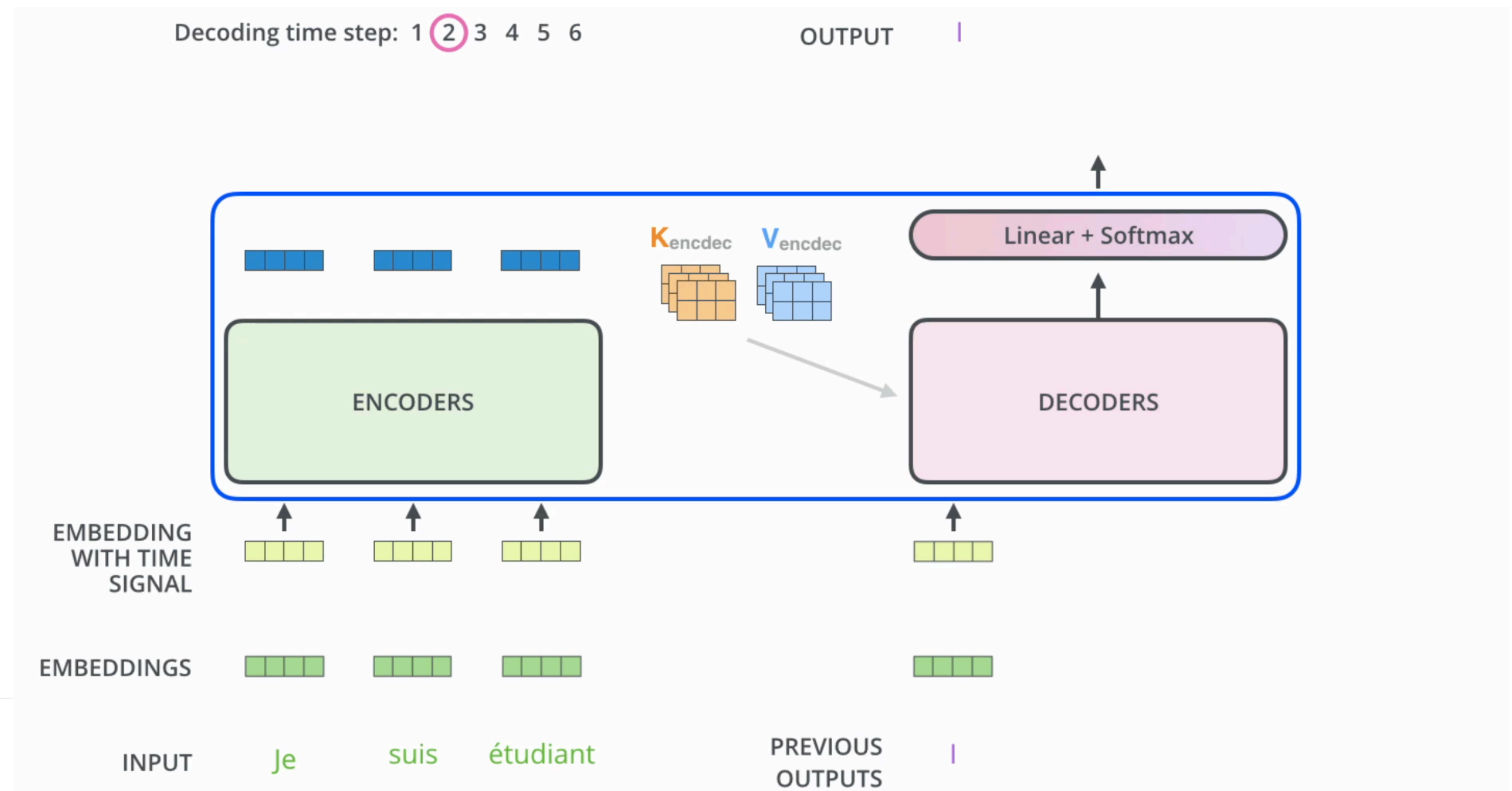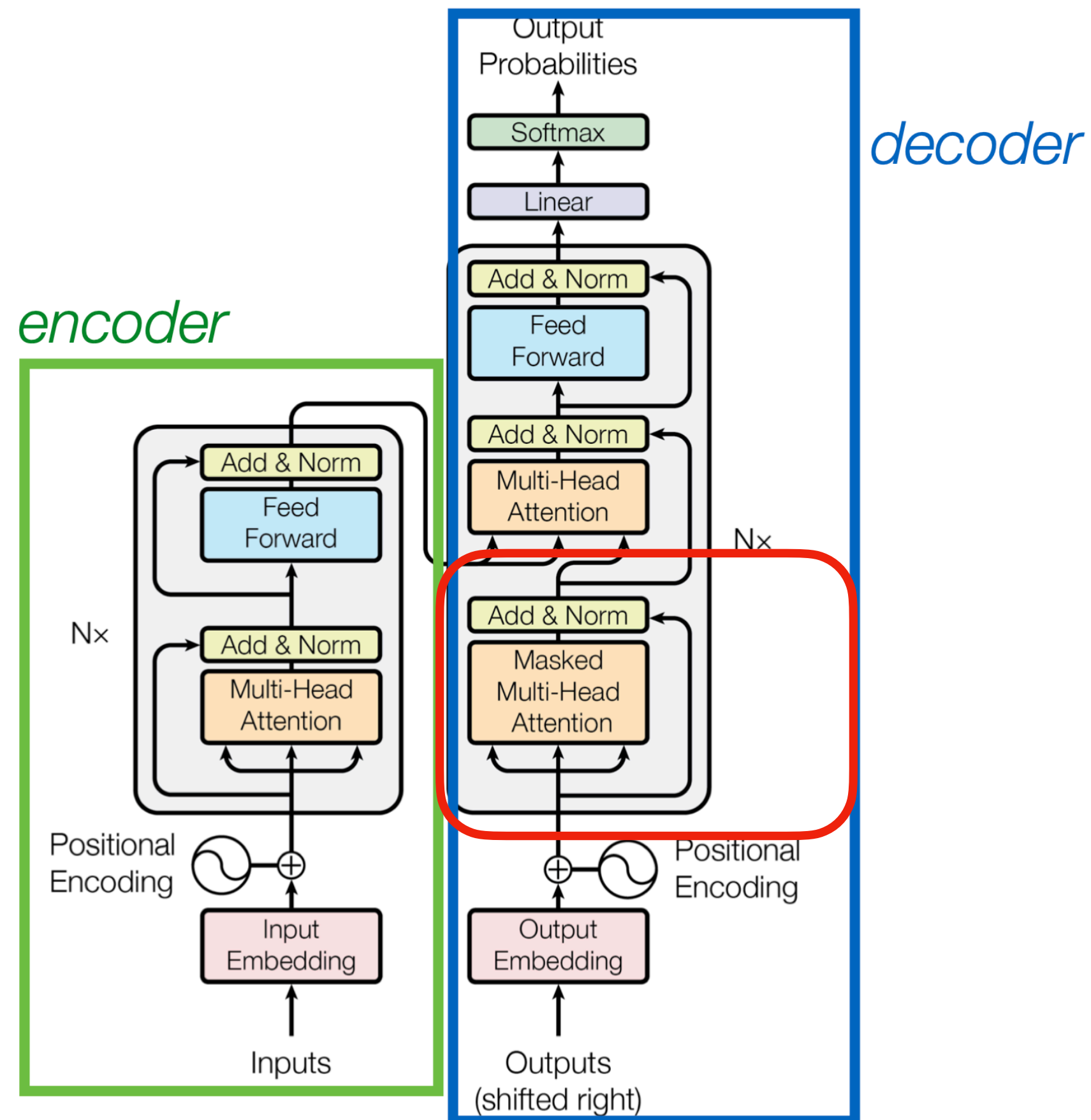Cross-attention

Self-attention

Cross-attention uses the output of encoder as input

48

# Masked Attention

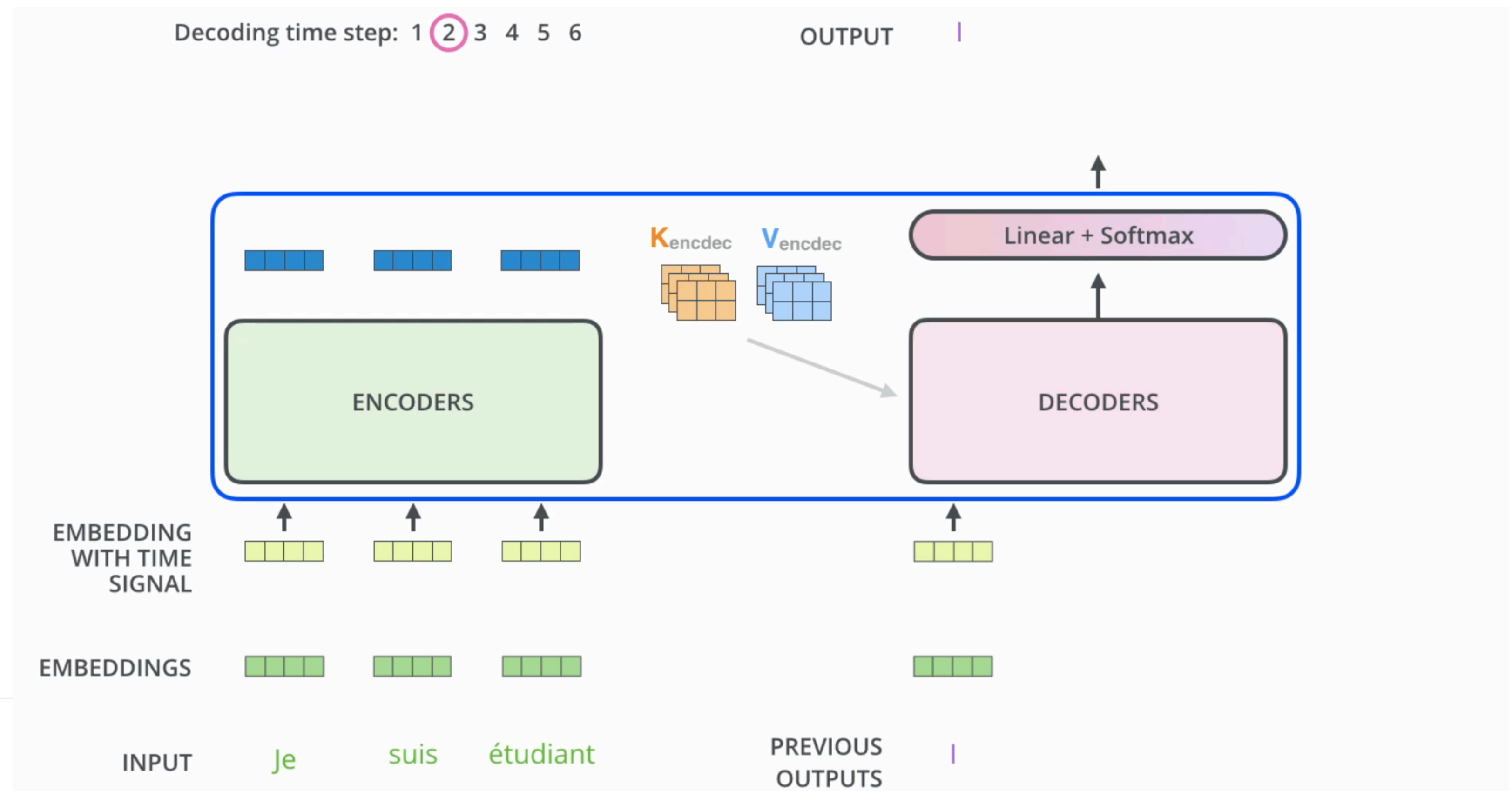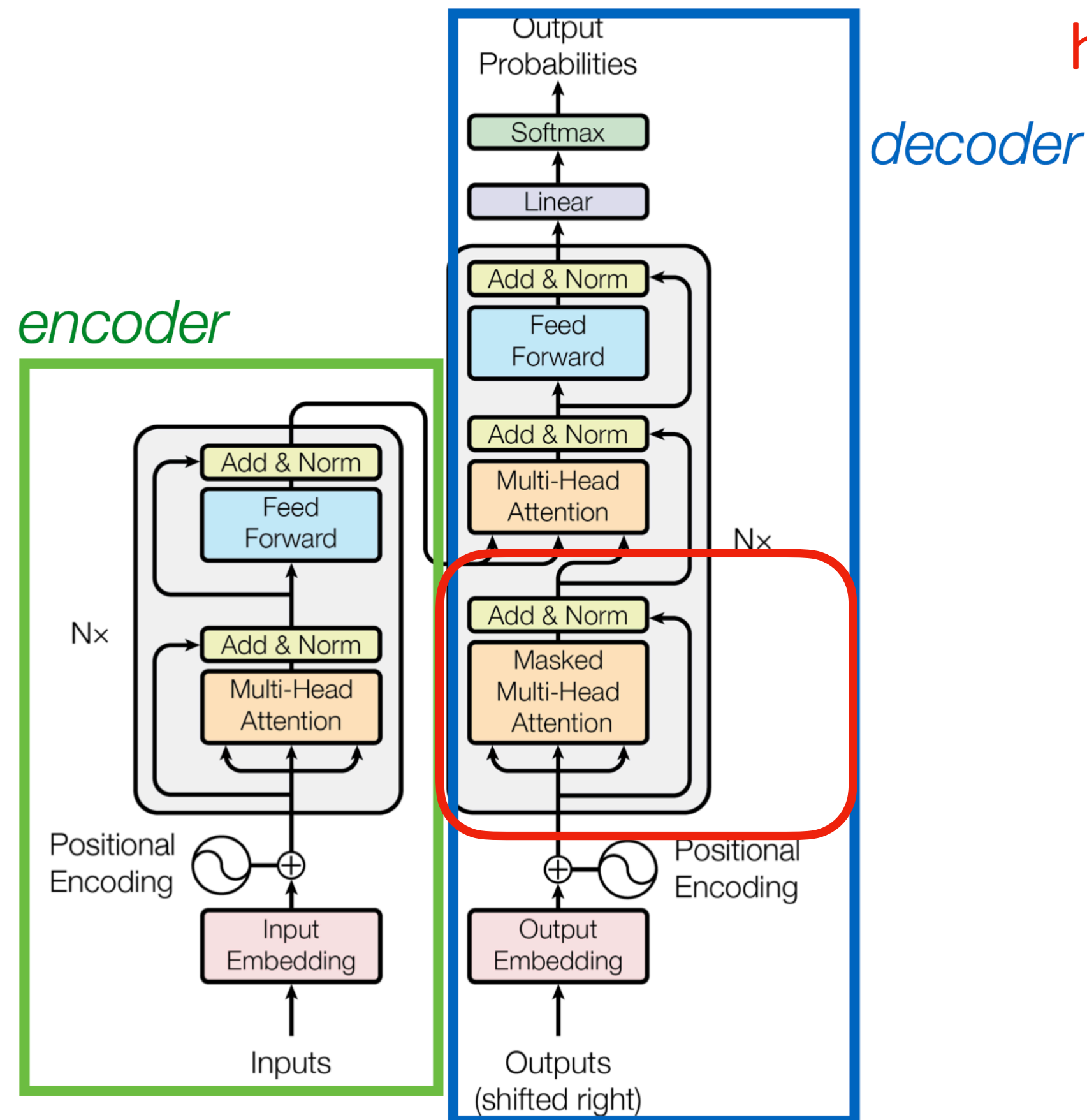# Masked Attention

# 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



50

# 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



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
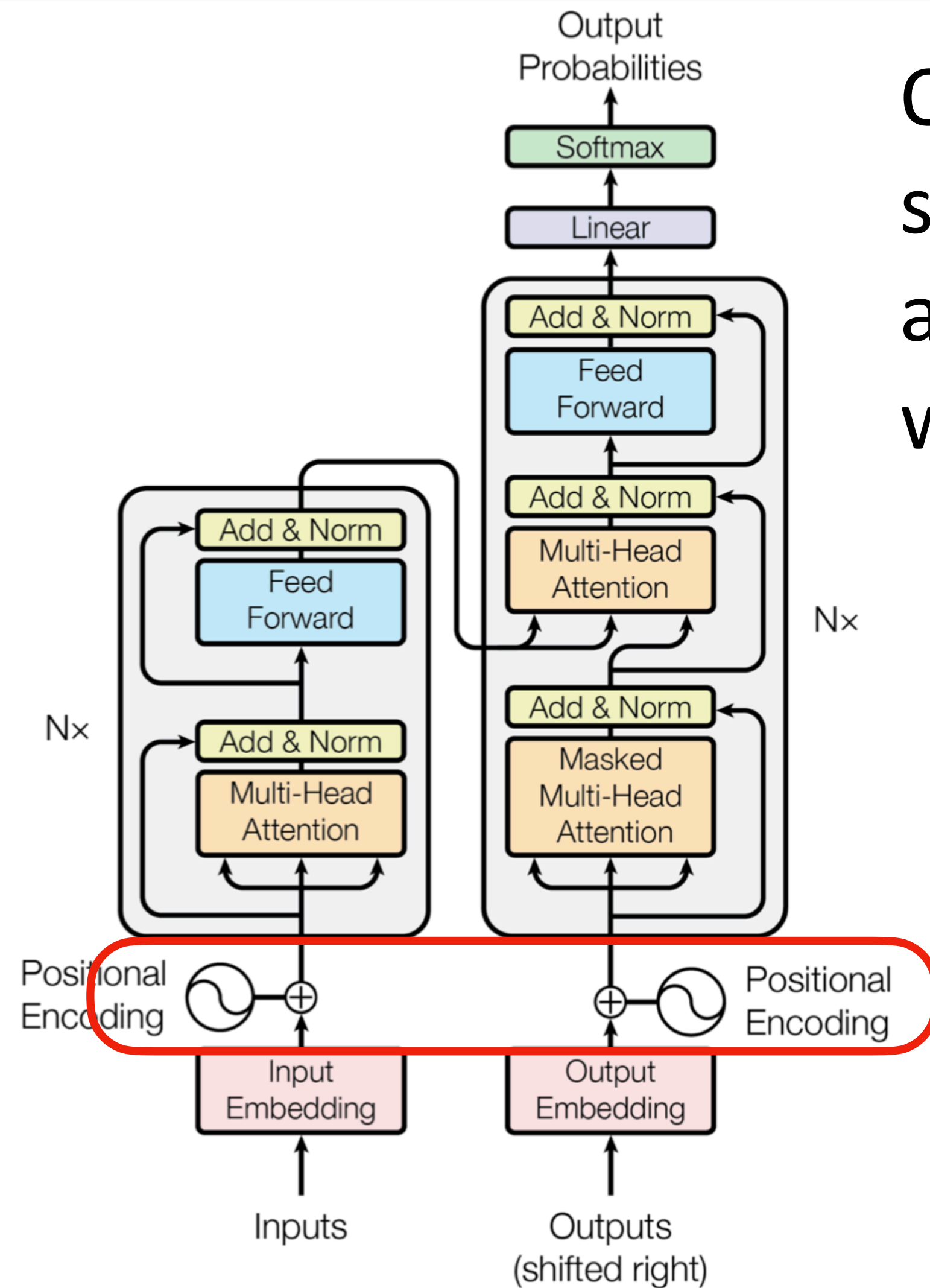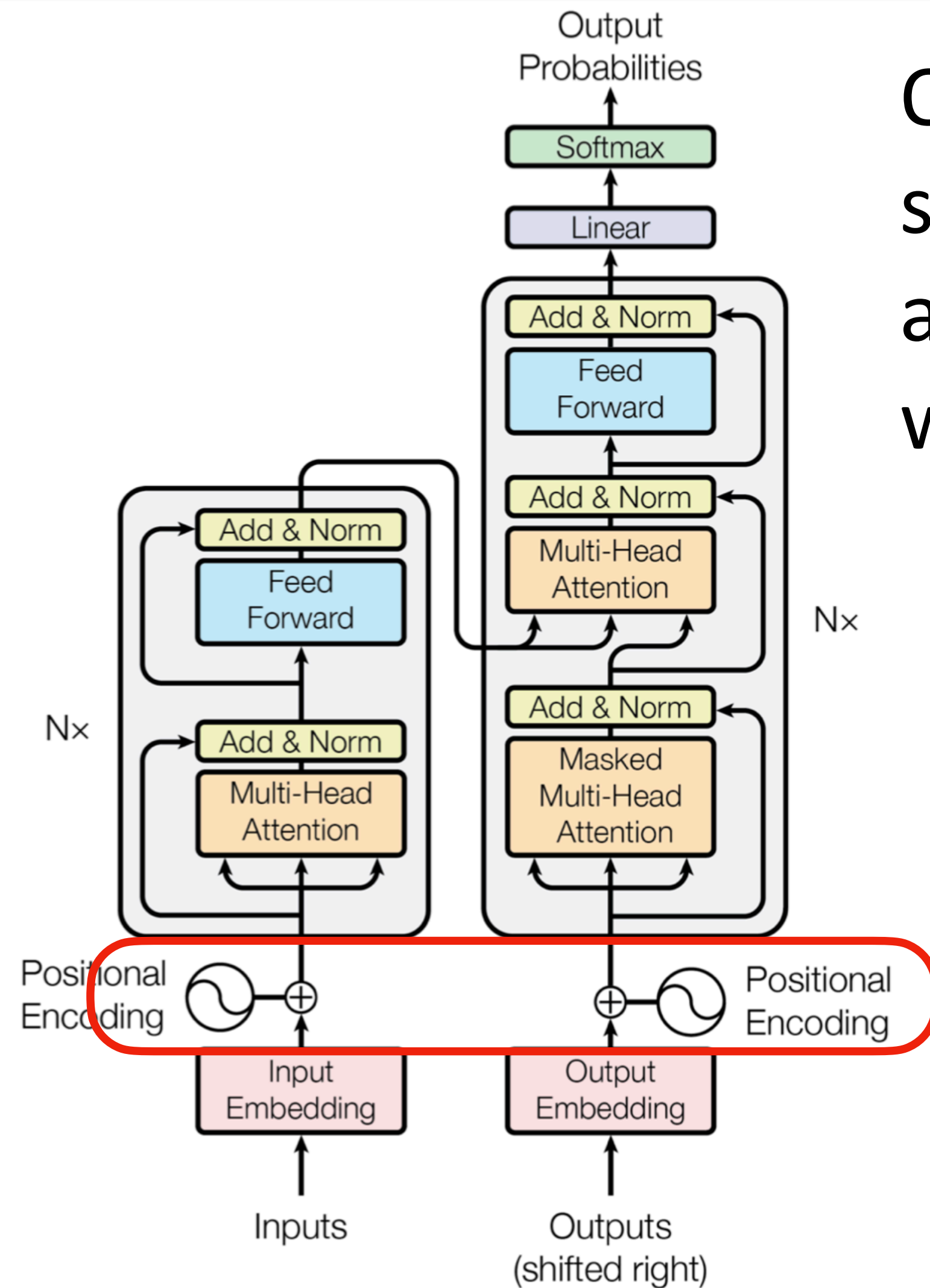
# Transformer Positional Encoding

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

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

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.

# 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

# 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!