# Seq2seq and Attention

## Lena Voita

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer ⟶ o Idea

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

o Idea

o Self-Attention

o Masked Self-Attention

o Multi-Head Attention

o Model Architecture

# Idea: Attention is All You Need

| | Seq2seq without attention | Seq2seq with attention |
|---|---|---|
| processing within encoder | RNN/CNN | RNN/CNN |
| processing within decoder | RNN/CNN | RNN/CNN |
| decoder-encoder interaction | static fixed-sized vector | attention |

# Idea: Attention is All You Need

| | Seq2seq without attention | Seq2seq with attention | Transformer |
|---|---|---|---|
| processing within encoder | RNN/CNN | RNN/CNN | attention |
| processing within decoder | RNN/CNN | RNN/CNN | attention |
| decoder-encoder interaction | static fixed-sized vector | attention | attention |

# Idea: Attention is All You Need

# Idea: Attention is All You Need

## Encoder

Who is doing:

What they are doing:

The animation is from the Google AI blog post.

# Idea: Attention is All You Need

<span style="color:green">Encoder</span>

<u>Who</u> is doing:

- all source tokens

<u>What</u> they are doing:

The animation is from the <u>Google AI blog post</u>.

# Idea: Attention is All You Need

Who is doing:

- all source tokens

What they are doing:

- look at each other

- update representations

repeat
N times

The animation is from the Google AI blog post.

# Idea: Attention is All You Need

<span style="color:#a01040">Decoder</span>

<u>Who</u> is doing:

<u>What</u> they are doing:

The animation is from the <u>Google AI blog post</u>.

# Idea: Attention is All You Need

## Decoder

Who is doing:

- target token at the current step

What they are doing:

The animation is from the Google AI blog post.

# Idea: Attention is All You Need

Who is doing:

- target token at the current step

What they are doing:

- looks at previous target tokens

- looks at source representations

- update representation

repeat
N times

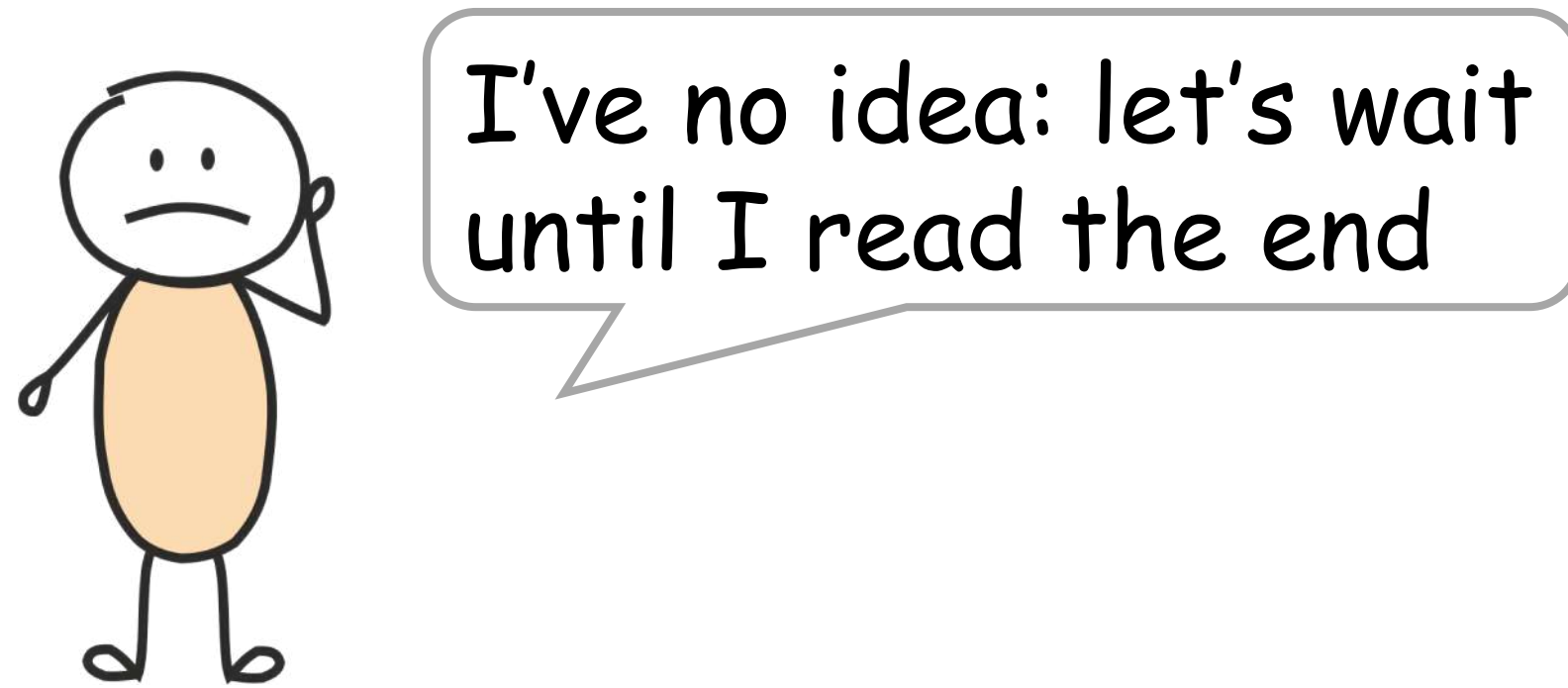The animation is from the Google AI blog post.

# Why can this be better than RNNs?

I arrived at the bank after crossing the ...        ...street?   ...river?

What does **bank** mean in this sentence?

# Why can this be better than RNNs?

I arrived at the bank after crossing the ...     ...street?   ...river?
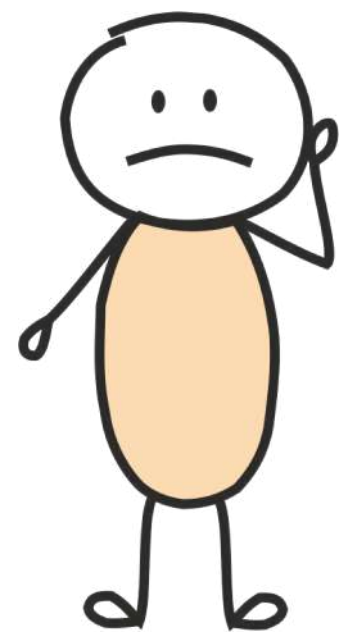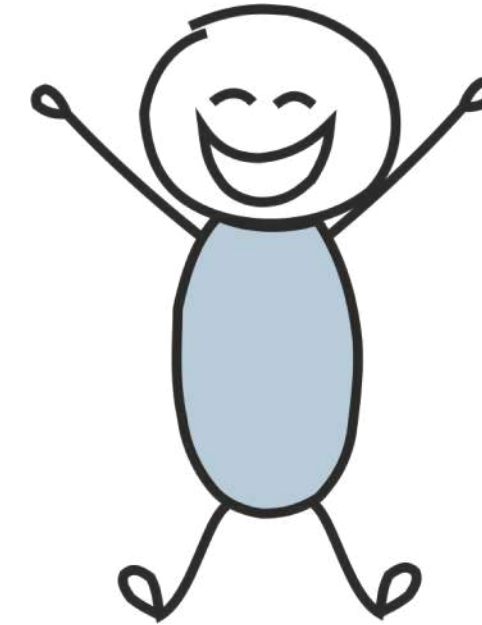
What does **bank** mean in this sentence?

I've no idea: let's wait until I read the end

RNNs

O(N) steps to process a sentence with length N

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer   ⟶  
  - Idea
  - Self-Attention
  - Masked Self-Attention
  - Multi-Head Attention
  - Model Architecture

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer  →  ○ Idea
                 ○ Self-Attention
                 ○ Masked Self-Attention
                 ○ Multi-Head Attention
                 ○ Model Architecture

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

# Self-Attention: Why "Self"?

Decoder-encoder attention is looking

- **from**: one current decoder state
- **at**: all encoder states

# Self-Attention: Why "Self"?

Decoder-encoder attention is looking
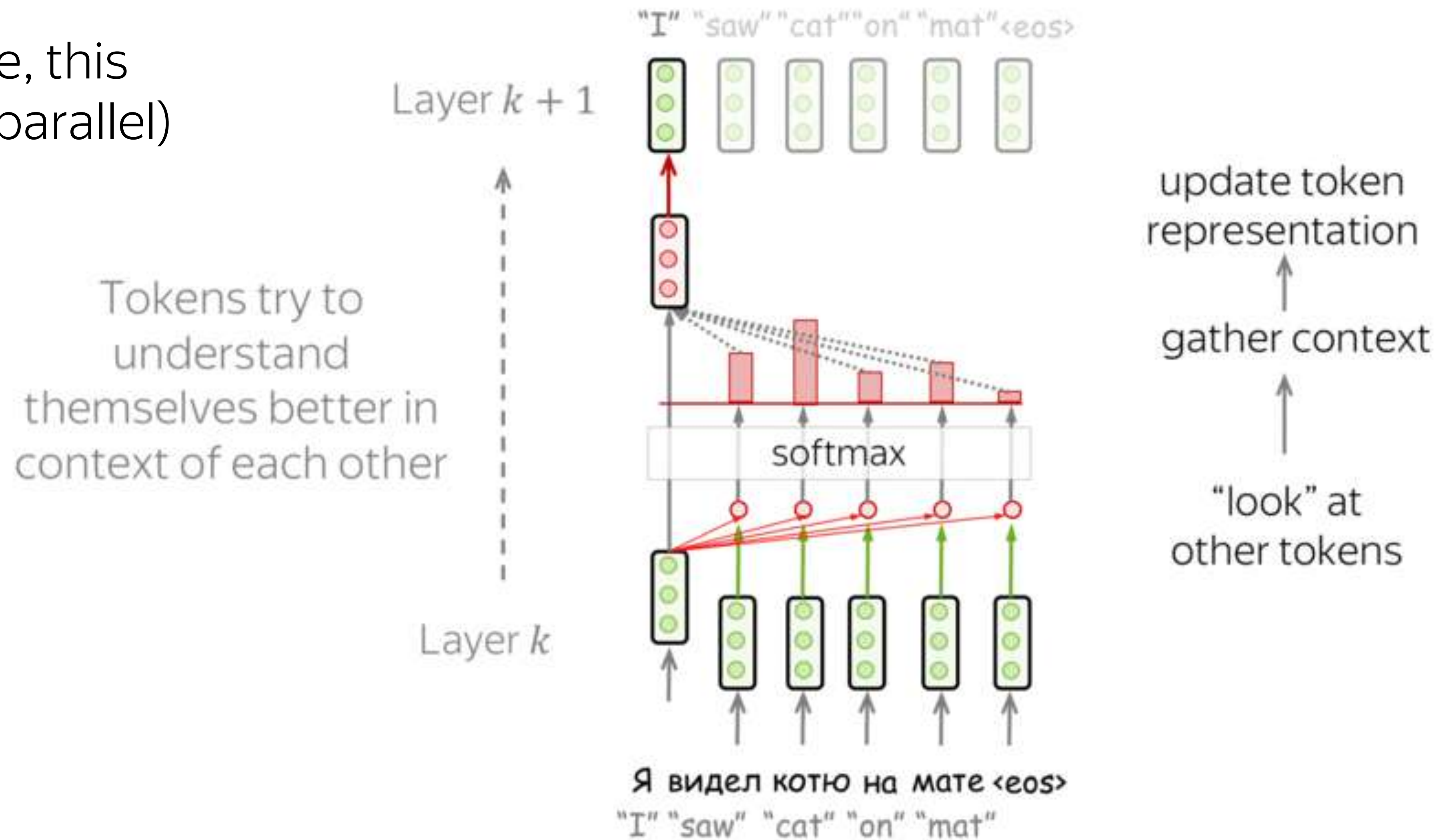
- **from**: one current decoder state

- **at**: all encoder states

Self-attention is looking

- **from**: each state from a set of states

- **at**: all other states in the same set

# Self-Attention: The "Look at Each Other" Part

(in practice, this happens in parallel)



Tokens try to understand themselves better in context of each other

Layer $k + 1$

"I" "saw" "cat""on" "mat"<eos>

softmax

Layer $k$

Я видел котю на мате <eos>

"I" "saw" "cat" "on" "mat"

update token representation

gather context

"look" at other tokens

# Query, Key, Value

Each vector receives three representations ("roles")

$$\left[W_Q\right] \times \bigg| \bigg| = \bigg| \bigg|$$

Query: vector **from** which the attention is looking

**"Hey there, do you have this information?"**
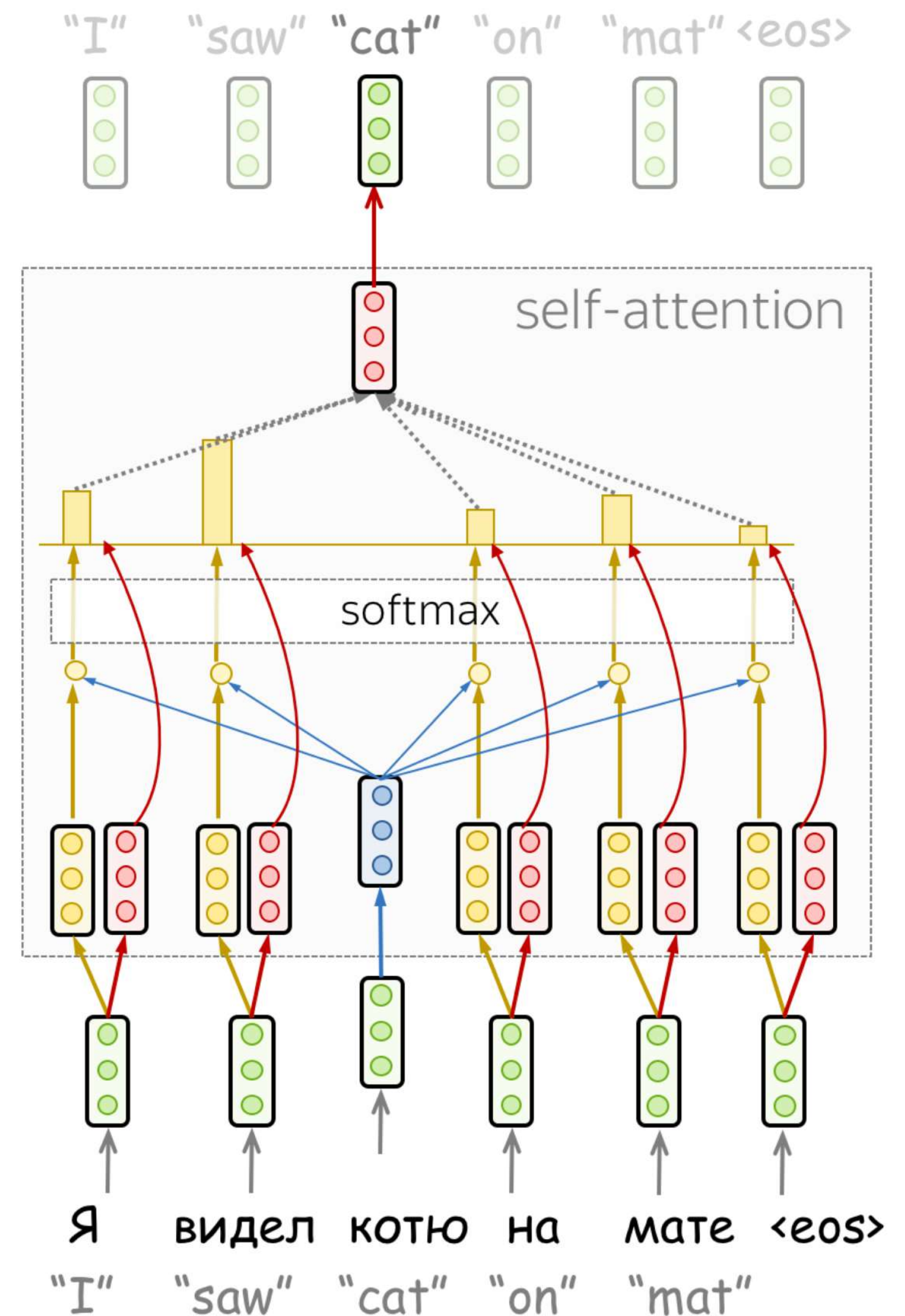
$$\left[W_K\right] \times \bigg| \bigg| = \bigg| \bigg|$$

Key: vector **at** which the query looks to compute weights

**"Hi, I have this information – give me a large weight!"**

$$\left[W_V\right] \times \bigg| \bigg| = \bigg| \bigg|$$

Value: their weighted sum is attention output

**"Here's the information I have!"**



"I"  "saw" "cat" "on"  "mat" <eos>

self-attention

softmax

Я  видел  котю  на  мате  <eos>
"I"  "saw"  "cat"  "on"  "mat"

# Query, Key, Value

Scaled dot-product

$$\frac{1}{\sqrt{d_k}} \cdot \boxed{\underset{q}{\circ \circ \circ}} \times \boxed{\underset{k^T}{\circ \circ}}$$

$$\text{score}(q, k) = \frac{q\,k^T}{\sqrt{d_k}}$$

"I"    "saw"   "cat"   "on"   "mat"   <eos>

self-attention

softmax

Я    видел   котю   на    мате   <eos>

"I"   "saw"   "cat"   "on"   "mat"

# Query, Key, Value

$$Attention(q, k, v) = softmax\left(\frac{qk^T}{\sqrt{d_k}}\right)v$$

Attention weights

from    to

vector dimensionality of K, V

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer         ⟶

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

 

- o Idea
- o Self-Attention
- o Masked Self-Attention
- o Multi-Head Attention
- o Model Architecture

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer  ⟶  o  Idea

  o  Self-Attention

- Subword Segmentation: BPE  o  Masked Self-Attention

  o  Multi-Head Attention

- 🔍 Analysis and Interpretability  o  Model Architecture

# Masked Self-Attention: "Don't Look Ahead"

In the decoder, we forbid looking at future tokens – we don't know them

Note: in training, decoder processes all target tokens at once – without masks, it would see future



update token representation

↑

gather context

↑

"look" at the previous tokens
(future tokens are masked out)

softmax

<bos> I saw a cat .

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer ——————————→
  - o Idea
  - o Self-Attention
  - o Masked Self-Attention
  - o Multi-Head Attention
  - o Model Architecture

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer ──────────────→    o   Idea

                                            o   Self-Attention

- Subword Segmentation: BPE    o   Masked Self-Attention

                                            o   Multi-Head Attention

- Analysis and Interpretability    o   Model Architecture

# Multi-Head Attention
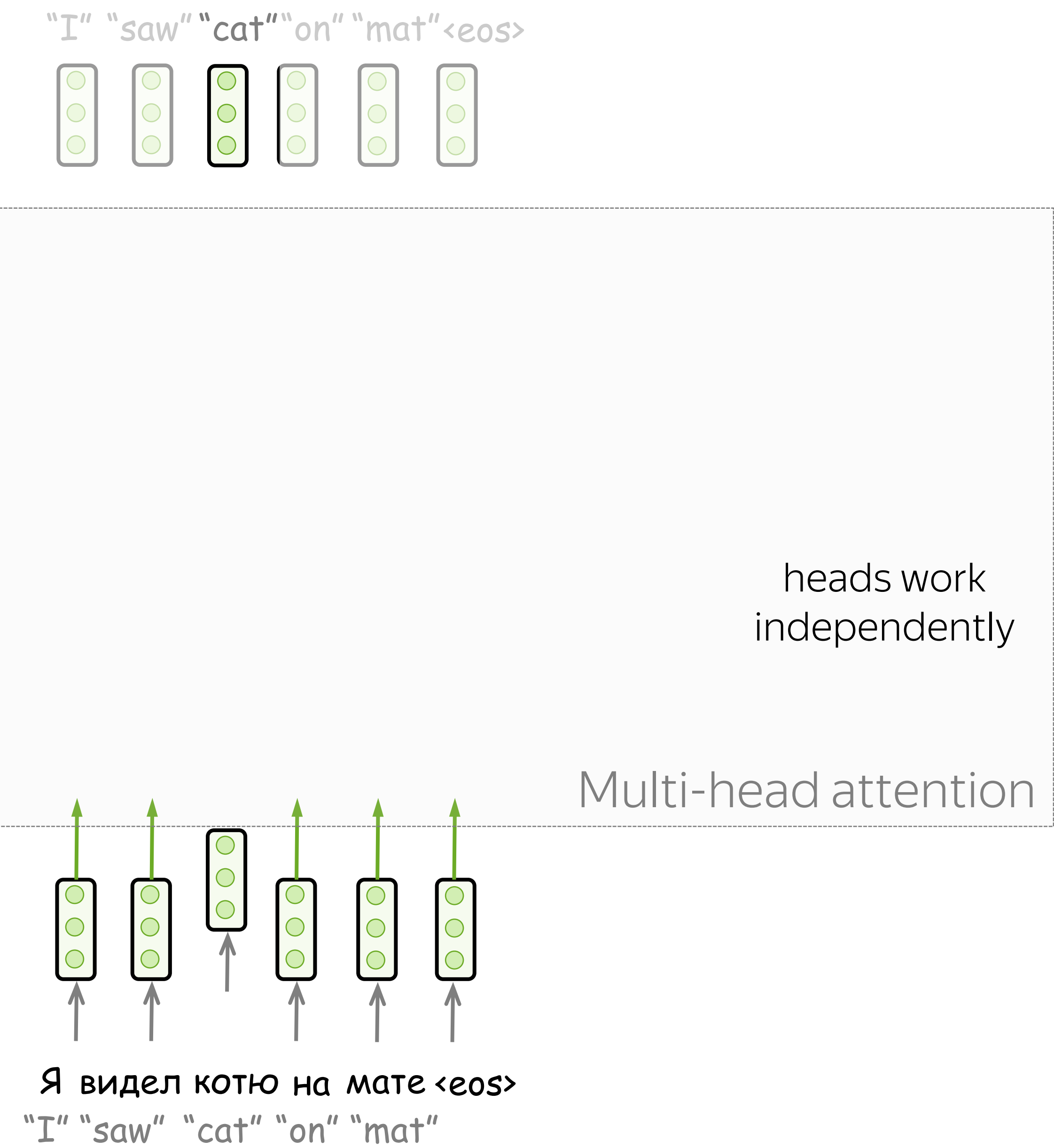
We need to track many
different things at once!

Она руководит **новым** проектом
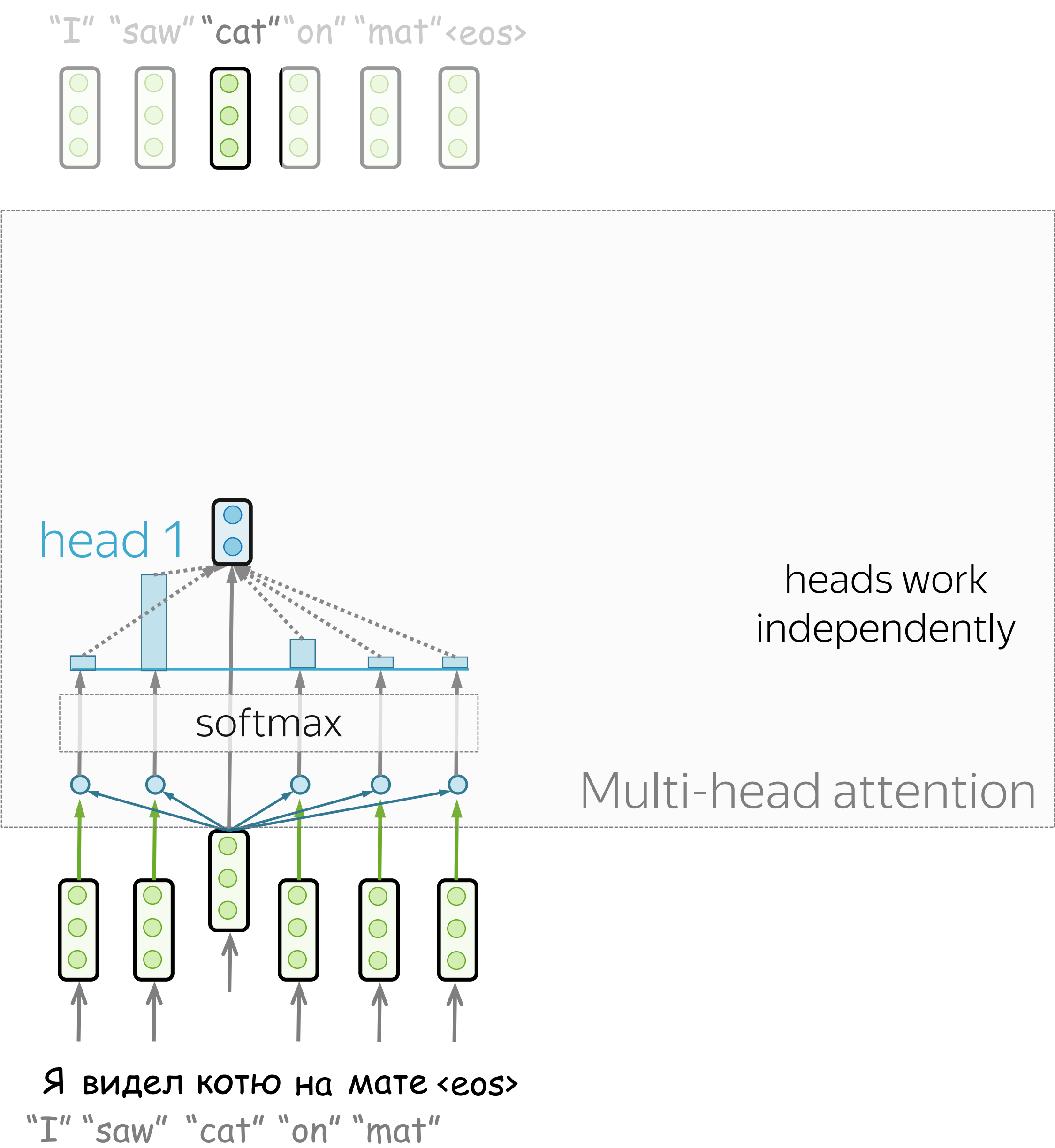
- Gender agreement
- Case government
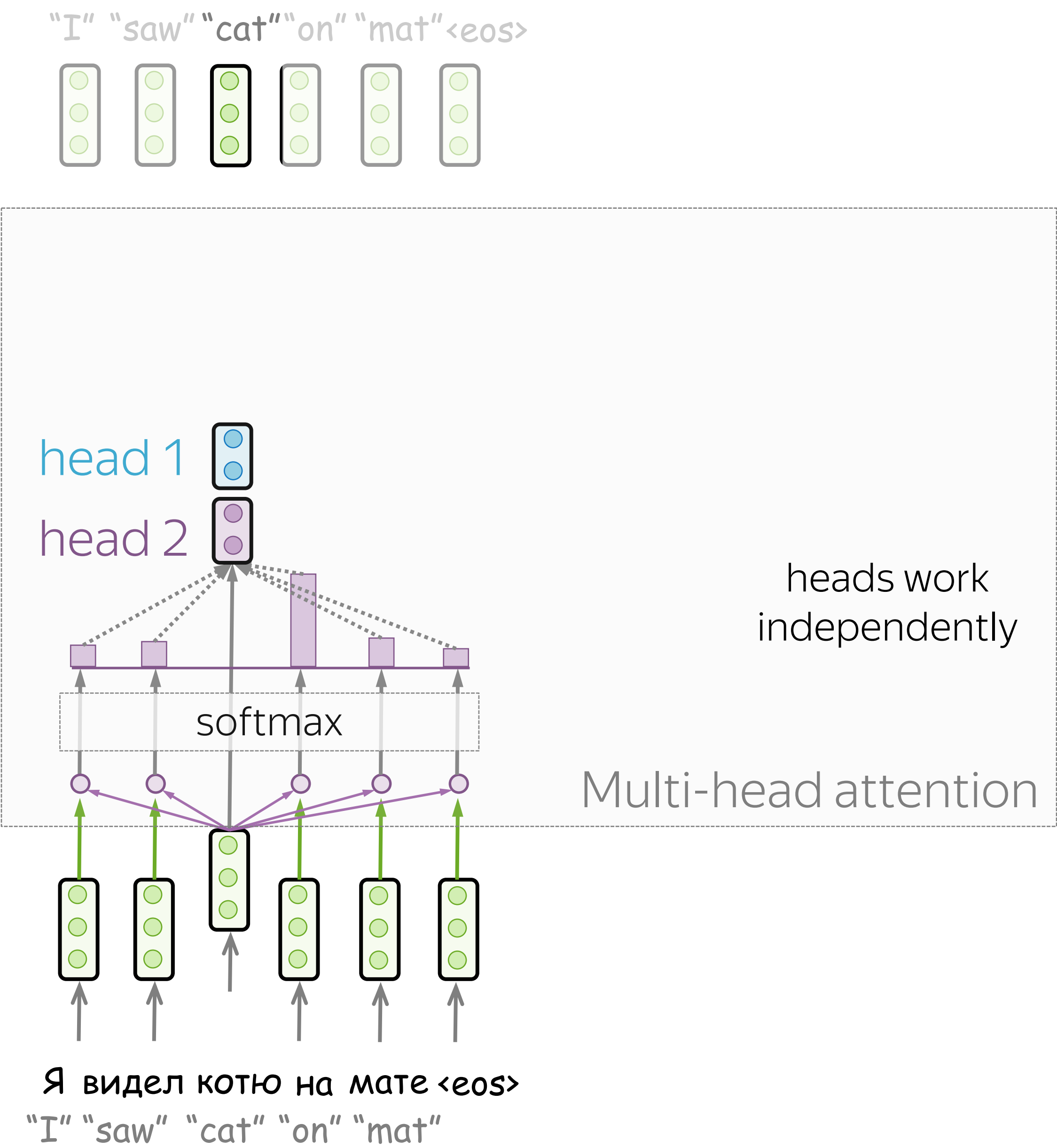- Lexical preferences
- ...

The example is from David Talbot

# Multi-Head Attention



"I" "saw" "cat" "on" "mat" <eos>

heads work
independently

Multi-head attention

Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention

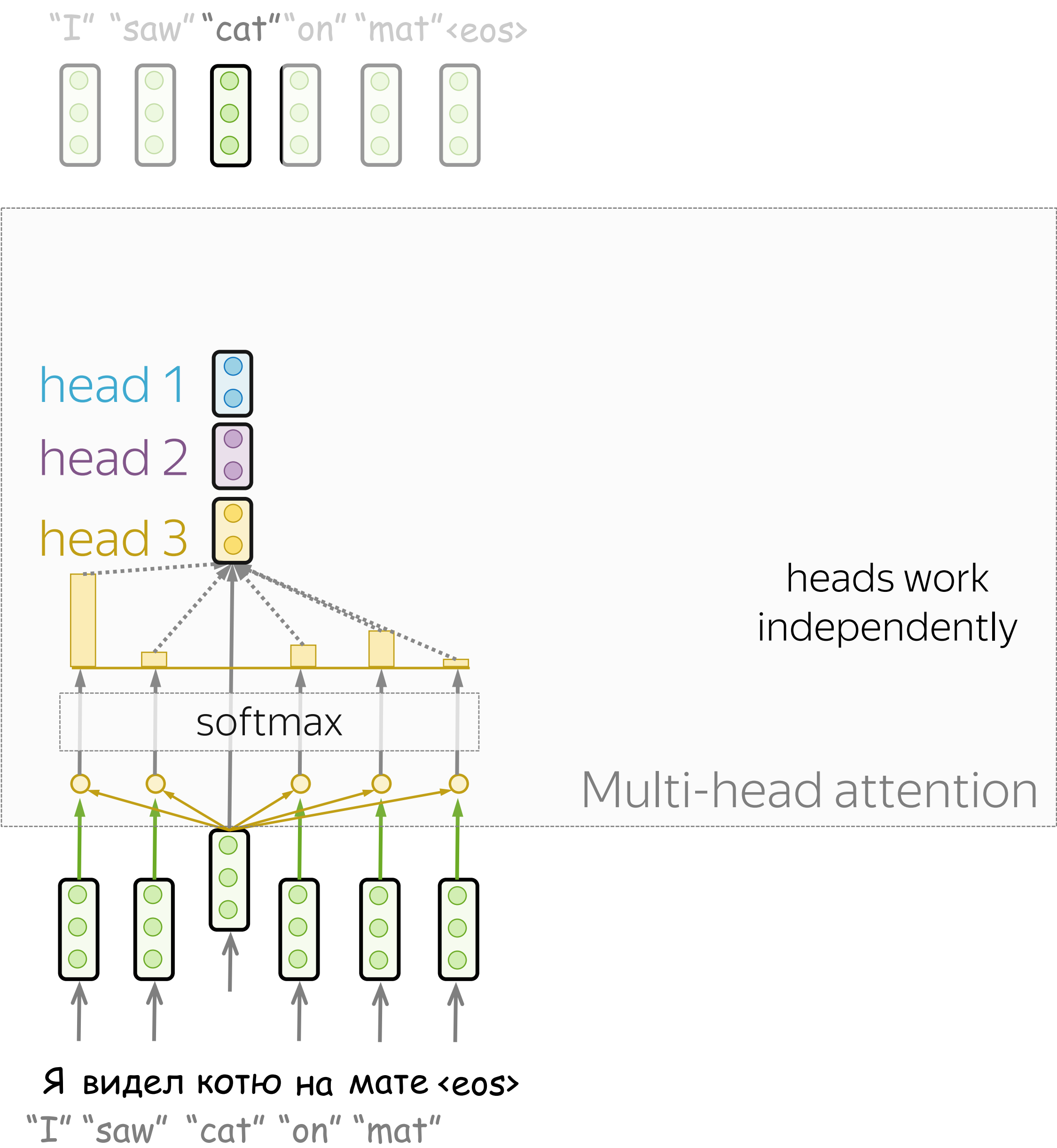# Multi-Head Attention

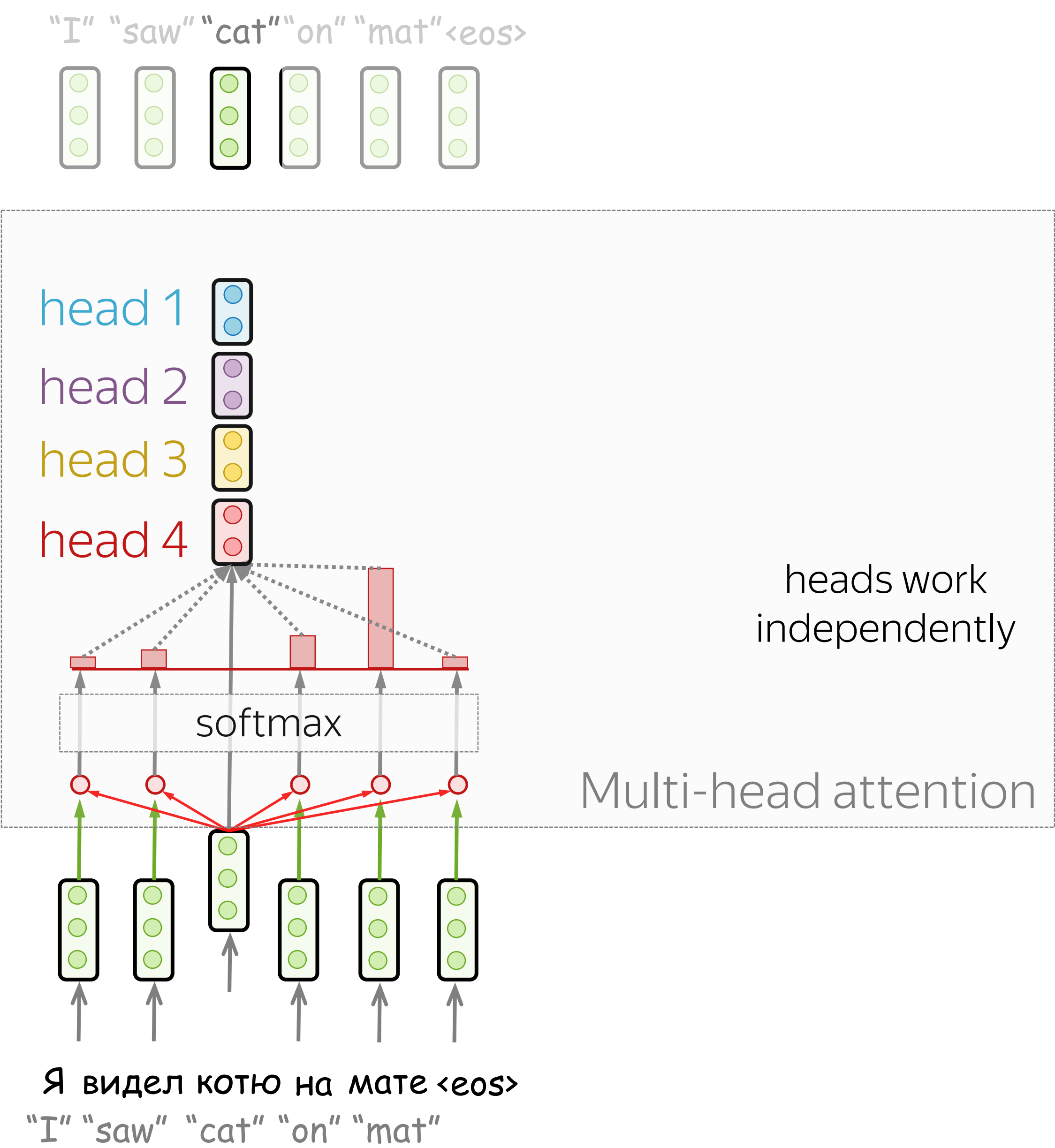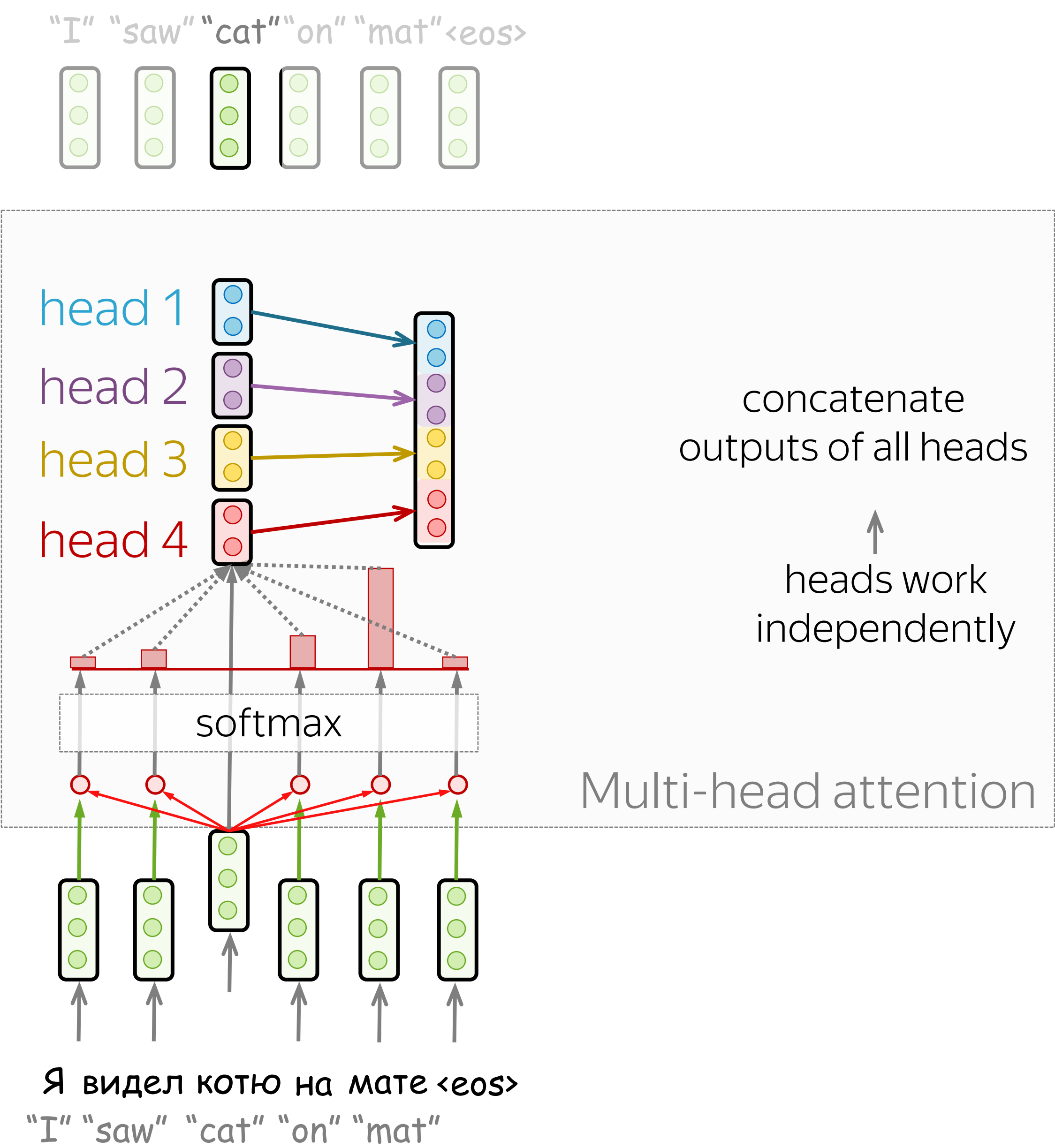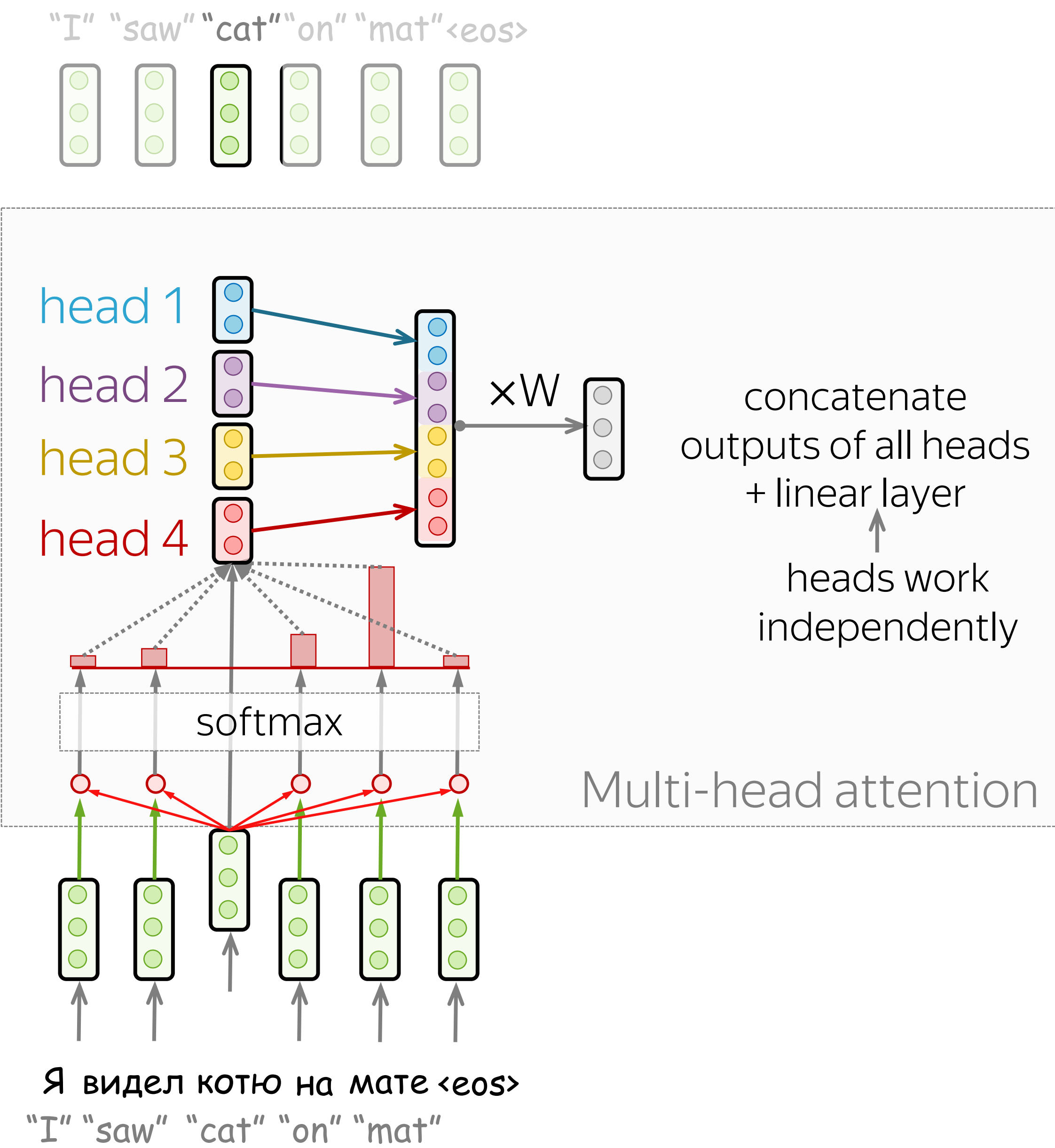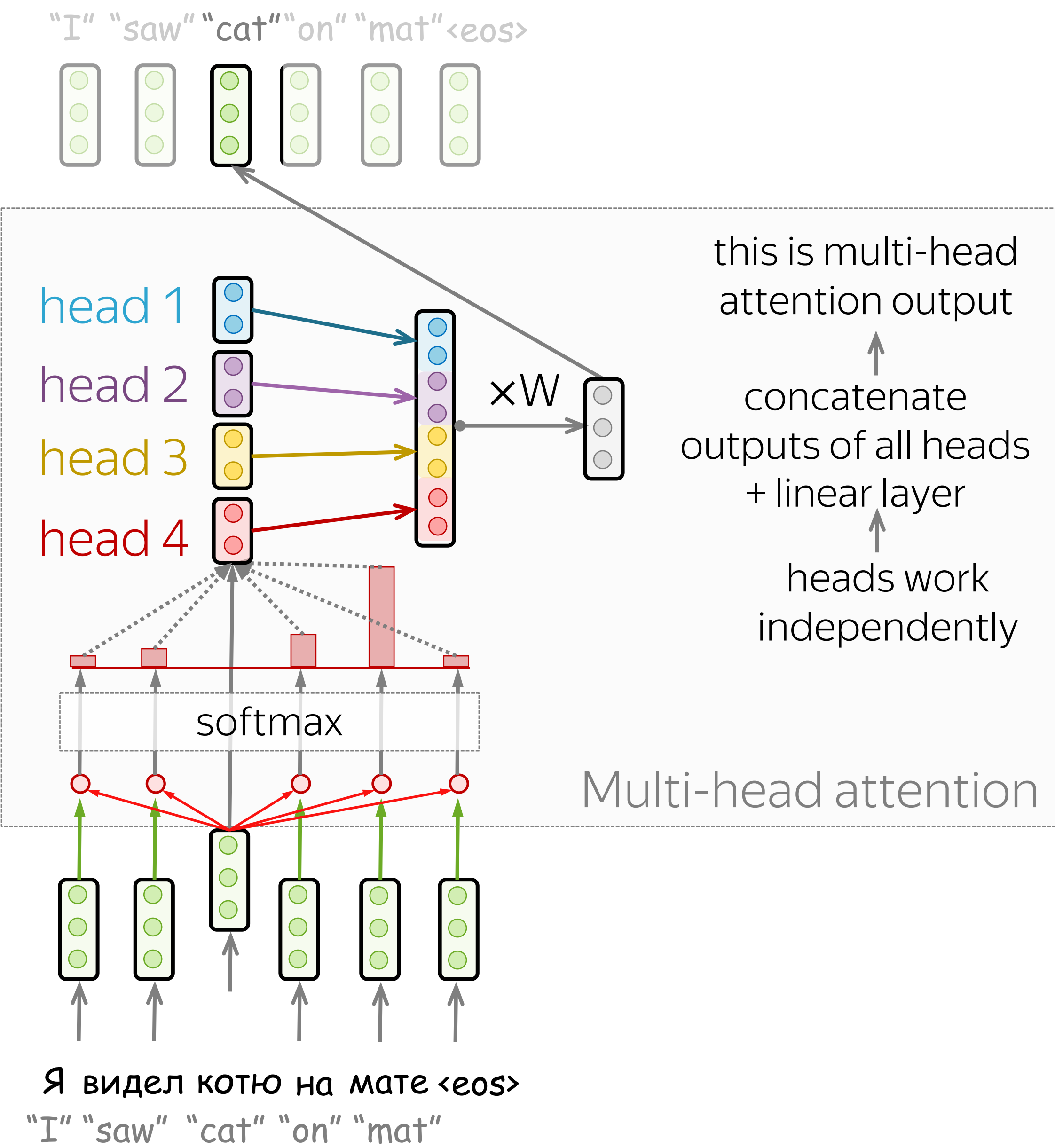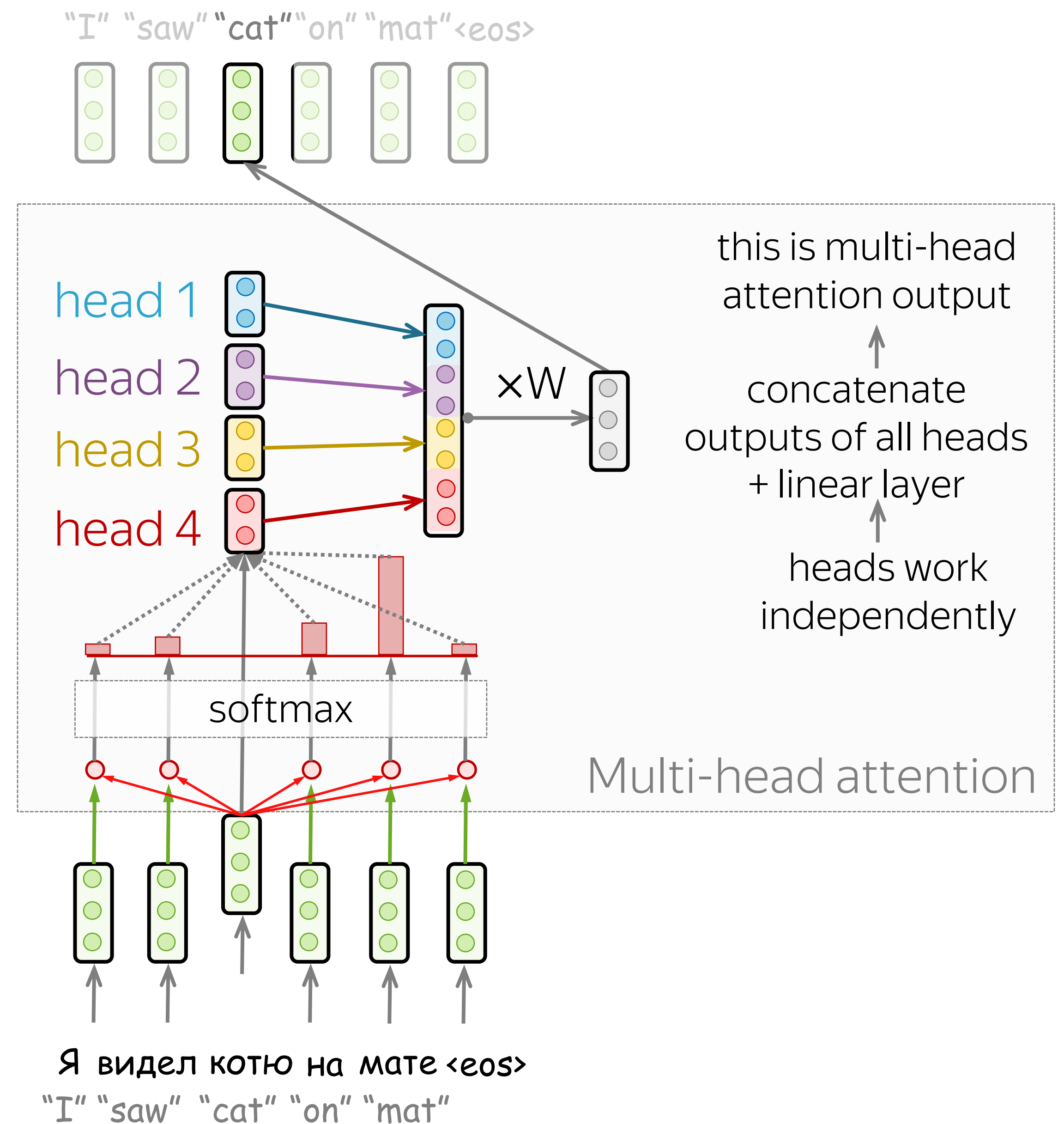# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention

# Multi-Head Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \ldots, \text{head}_n)W_o,$$

$$\text{head}_i = \text{Attention}(QW_Q^i, KW_K^i, VW_V^i)$$

# Multi-Head Attention

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_n)W_o,$$

$$\text{head}_i = \text{Attention}(QW_Q^i, KW_K^i, VW_V^i)$$



query, key, value for:

← head$_1$
← head$_2$
← head$_3$
← head$_4$

Split equally into number of heads parts

$W_Q$  $W_K$  $W_V$

"I" "saw" "cat" "on" "mat" <eos>

head 1
head 2
head 3
head 4

×W

this is multi-head attention output
↑
concatenate outputs of all heads + linear layer
↑
heads work independently

softmax

Multi-head attention

Я видел котю на мате <eos>
"I" "saw" "cat" "on" "mat"

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer     ⟶

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

o Idea

o Self-Attention

o Masked Self-Attention

o Multi-Head Attention

o Model Architecture

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer     →     
  - Idea

  - Self-Attention

  - Masked Self-Attention

  - Multi-Head Attention

  - Model Architecture

- Subword Segmentation: BPE

- Analysis and Interpretability

# Transformer

# Transformer

# 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

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

Encoder self-attention:
tokens look at each other

queries, keys, values
are computed from
encoder states

# Transformer



Feed-forward network:
after taking information from other tokens, take a moment to think and process this information

Encoder self-attention:
tokens look at each other

queries, keys, values are computed from encoder states

# Transformer

Output
Probabilities

Softmax

Linear

Add & Norm

Feed
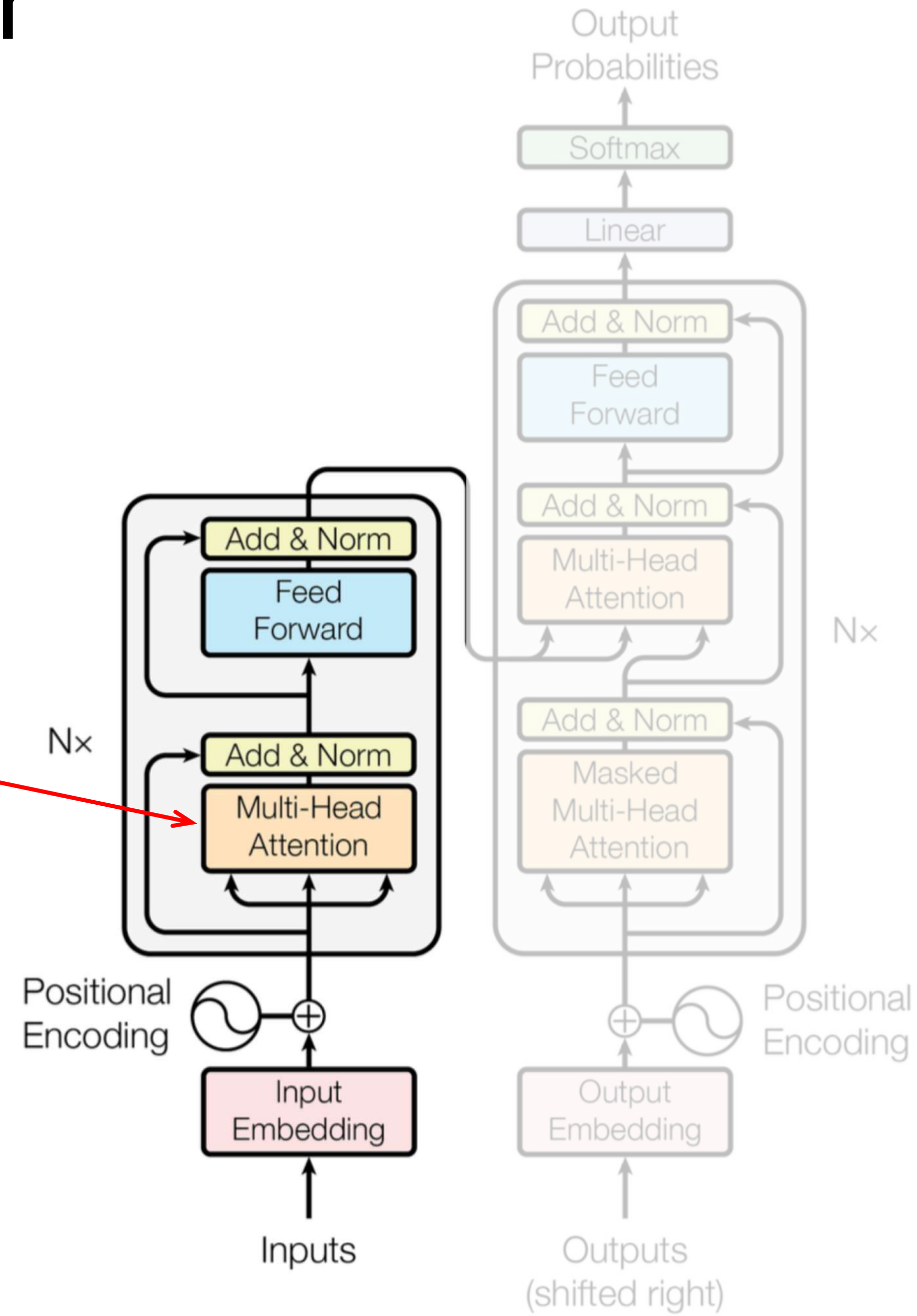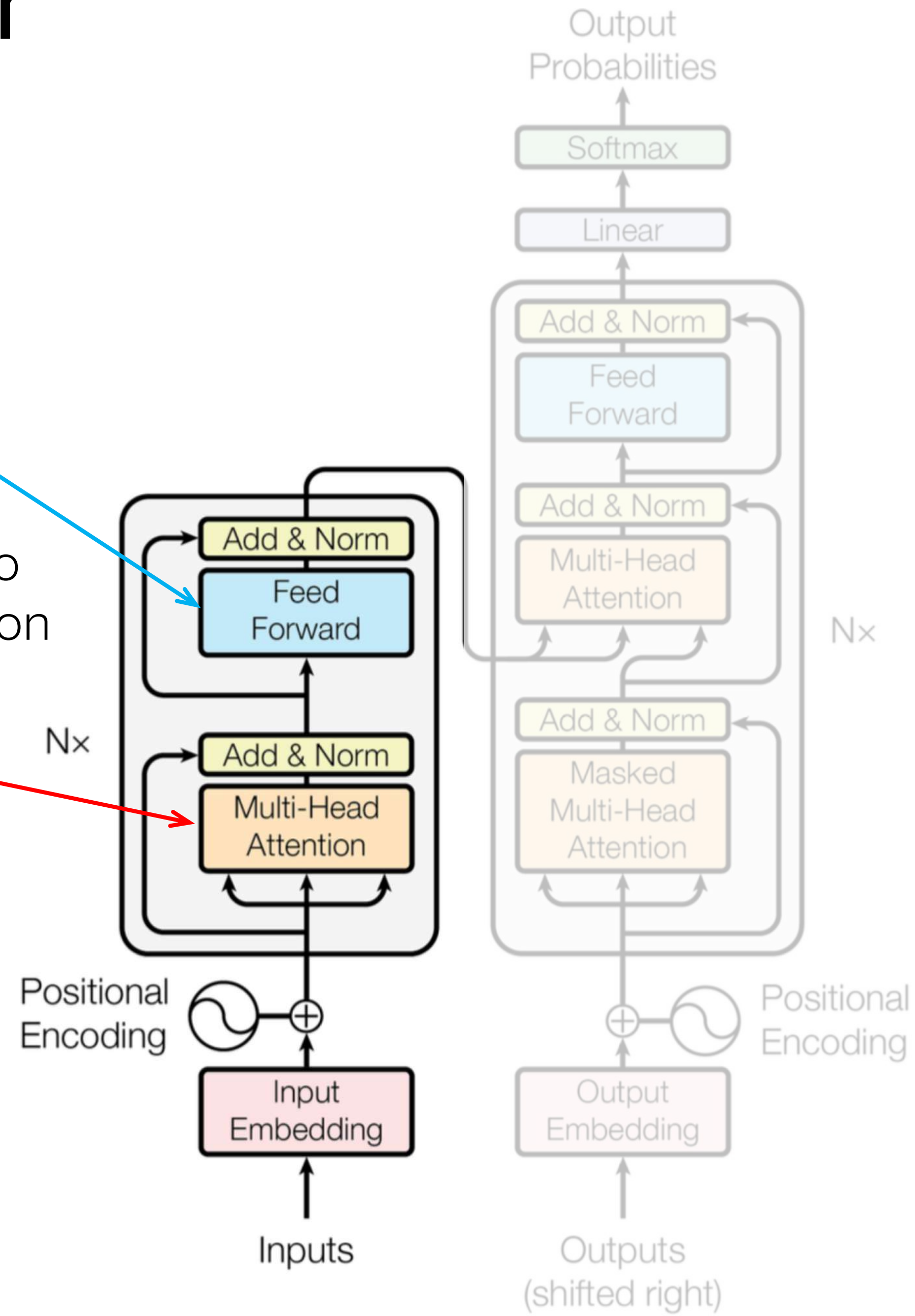Forward

Add & Norm

Multi-Head
Attention

N×

Add & Norm

Masked
Multi-Head
Attention

Feed-forward network:
after taking information from
other tokens, take a moment to
think and process this information

Add & Norm

Feed
Forward

N×

Encoder self-attention:
tokens look at each other

Add & Norm

Multi-Head
Attention

queries, keys, values
are computed from
encoder states

Positional
Encoding

Positional
Encoding

Input
Embedding

Output
Embedding

Inputs

Outputs
(shifted right)

# Transformer



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Output Embedding

Outputs (shifted right)

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Positional Encoding

Input Embedding

Inputs

Feed-forward network:
after taking information from other tokens, take a moment to think and process this information

Encoder self-attention:
tokens look at each other
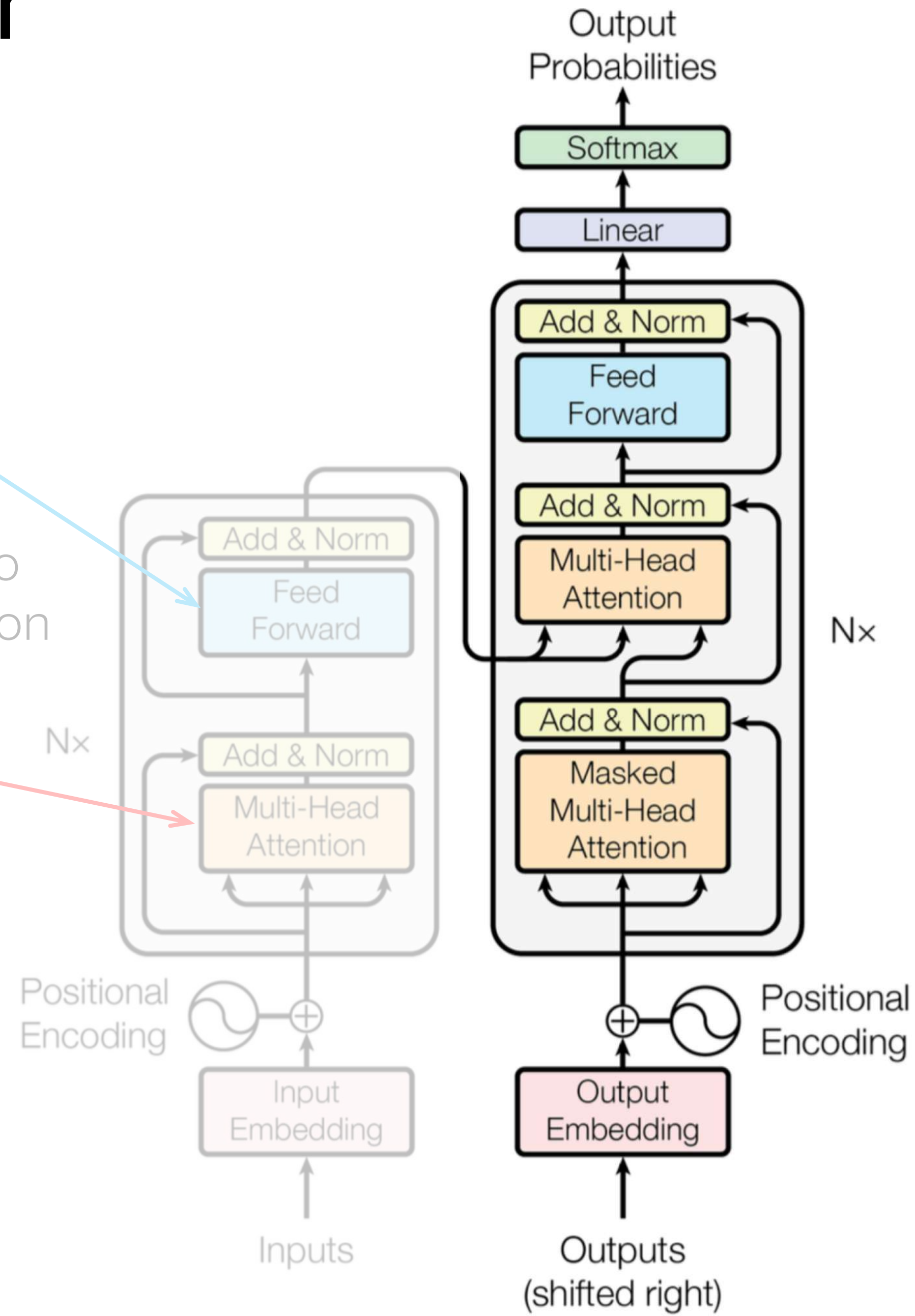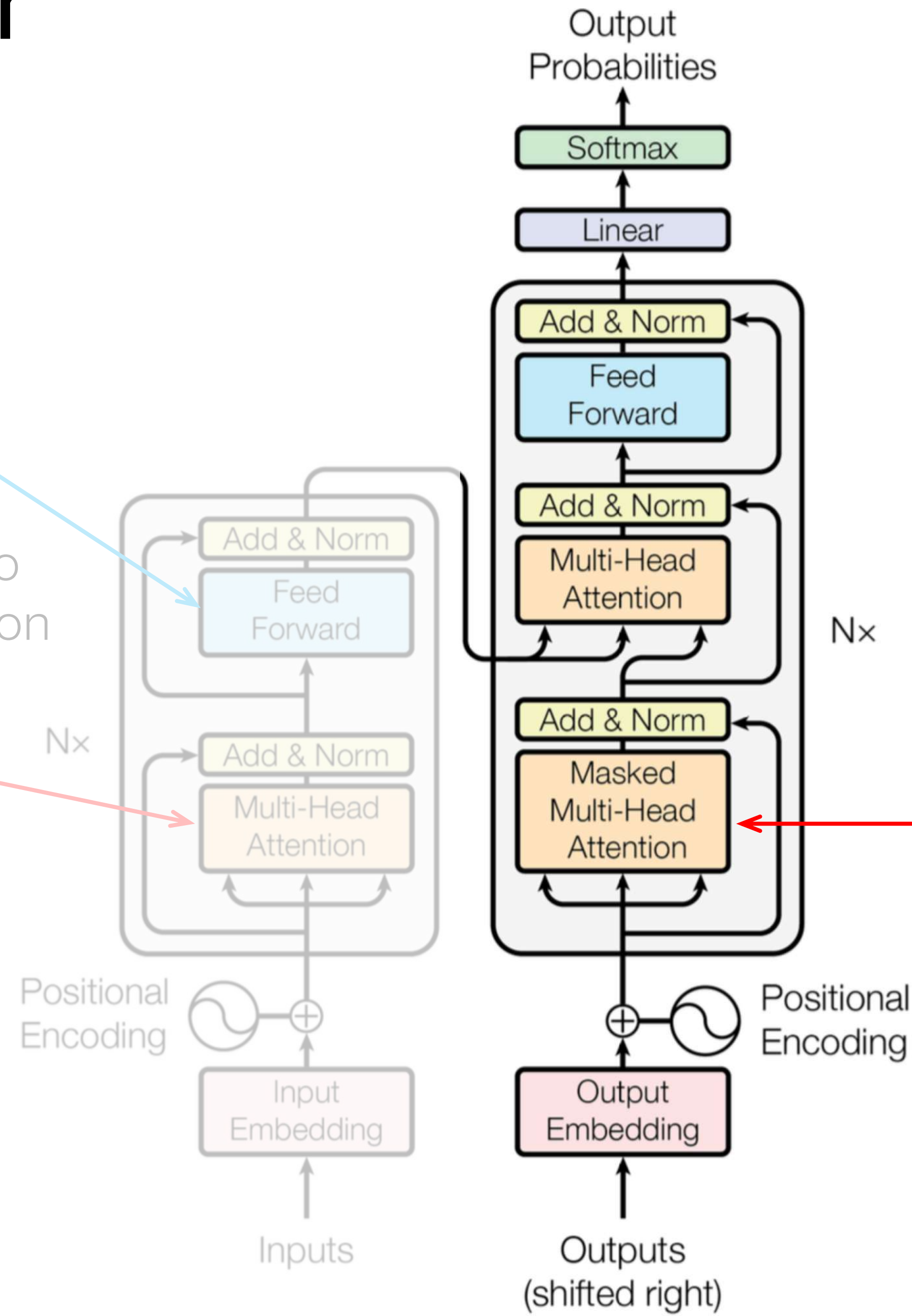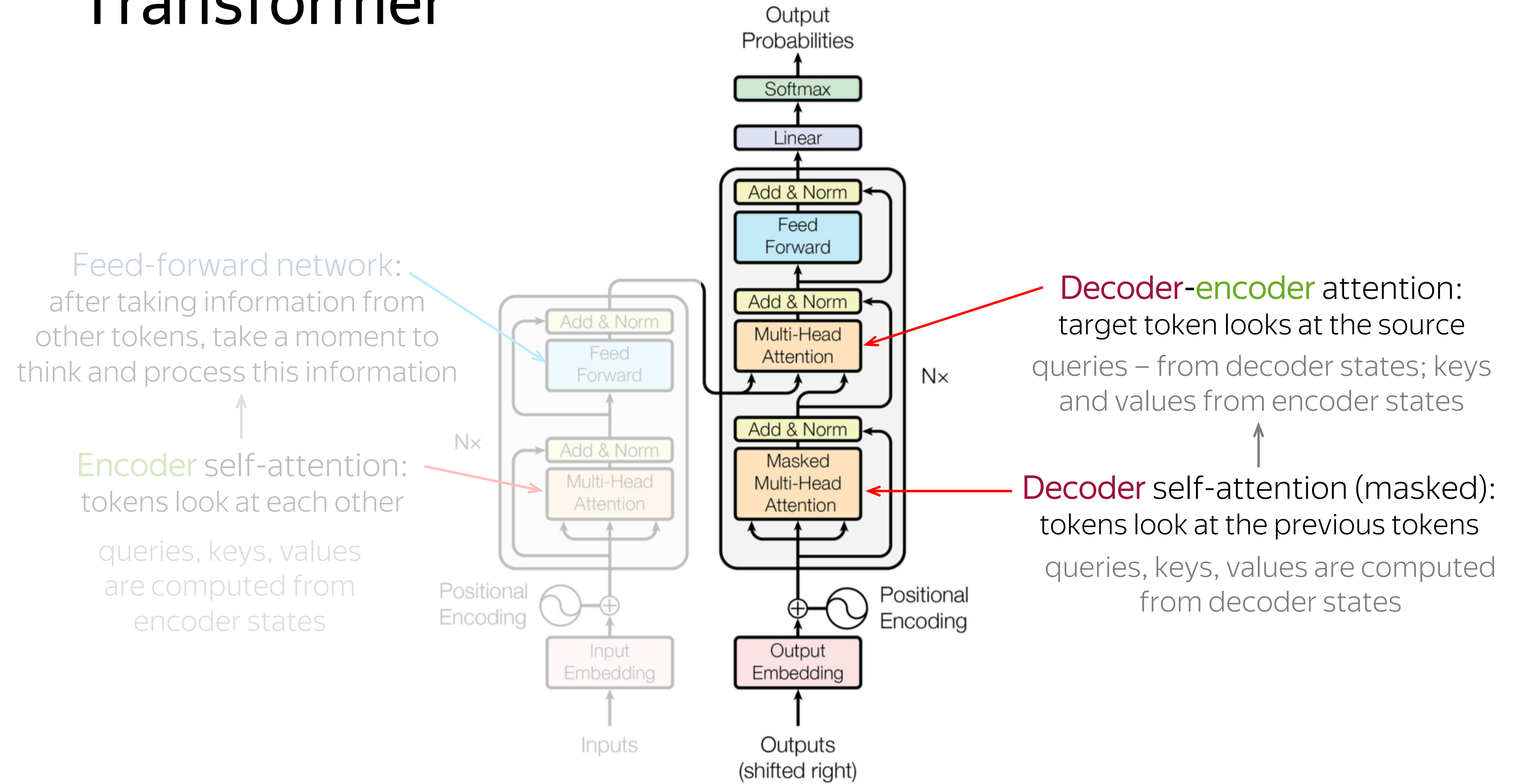
queries, keys, values are computed from encoder states

Decoder self-attention (masked):
tokens look at the previous tokens

queries, keys, values are computed from decoder states

# Transformer

Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Positional Encoding

Output Embedding

Outputs (shifted right)

Positional Encoding

Input Embedding

Inputs

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

**Feed-forward network:**
after taking information from other tokens, take a moment to think and process this information

**Encoder self-attention:**
tokens look at each other

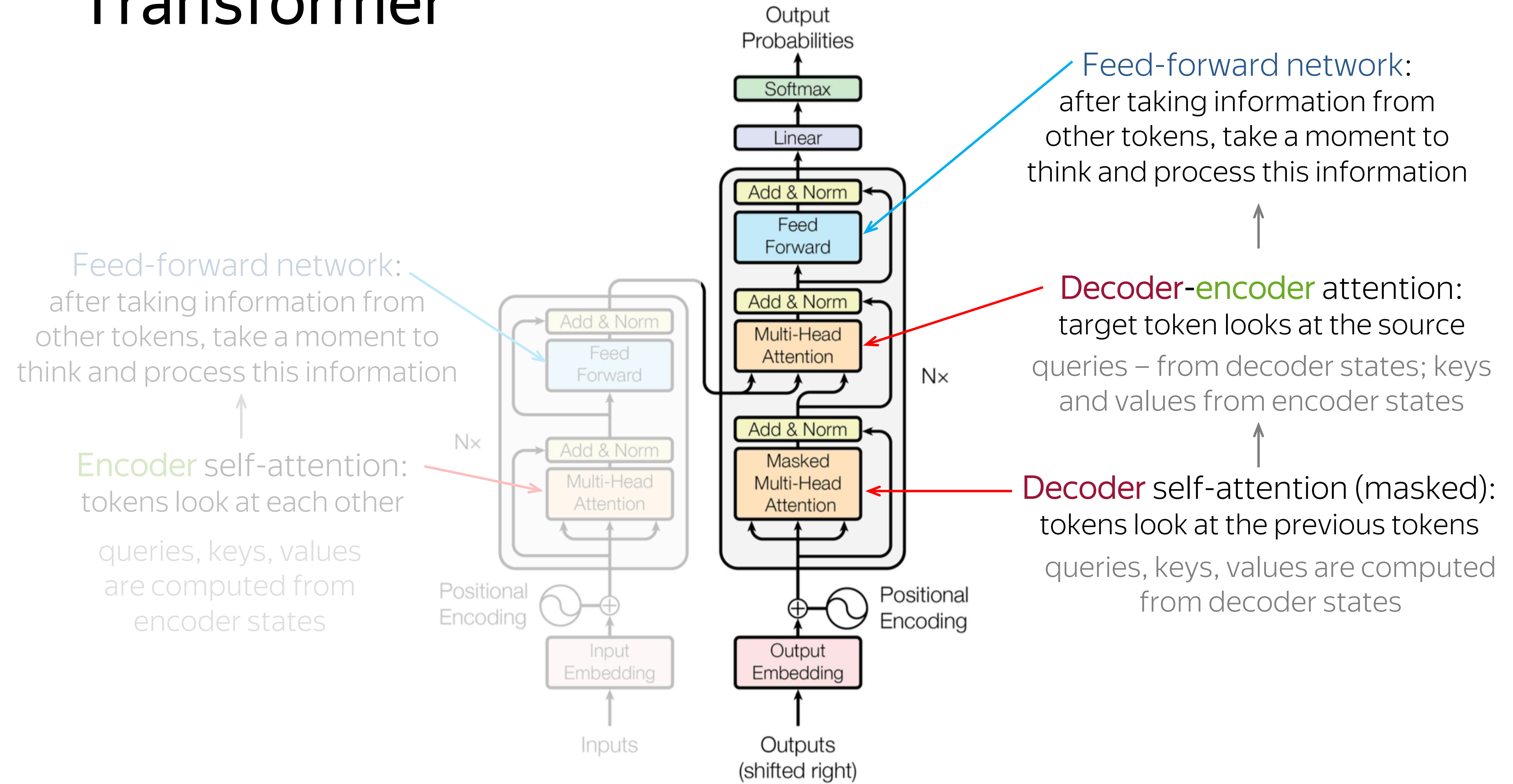queries, keys, values are computed from encoder states

**Decoder-encoder attention:**
target token looks at the source

queries – from decoder states; keys and values from encoder states

**Decoder self-attention (masked):**
tokens look at the previous tokens
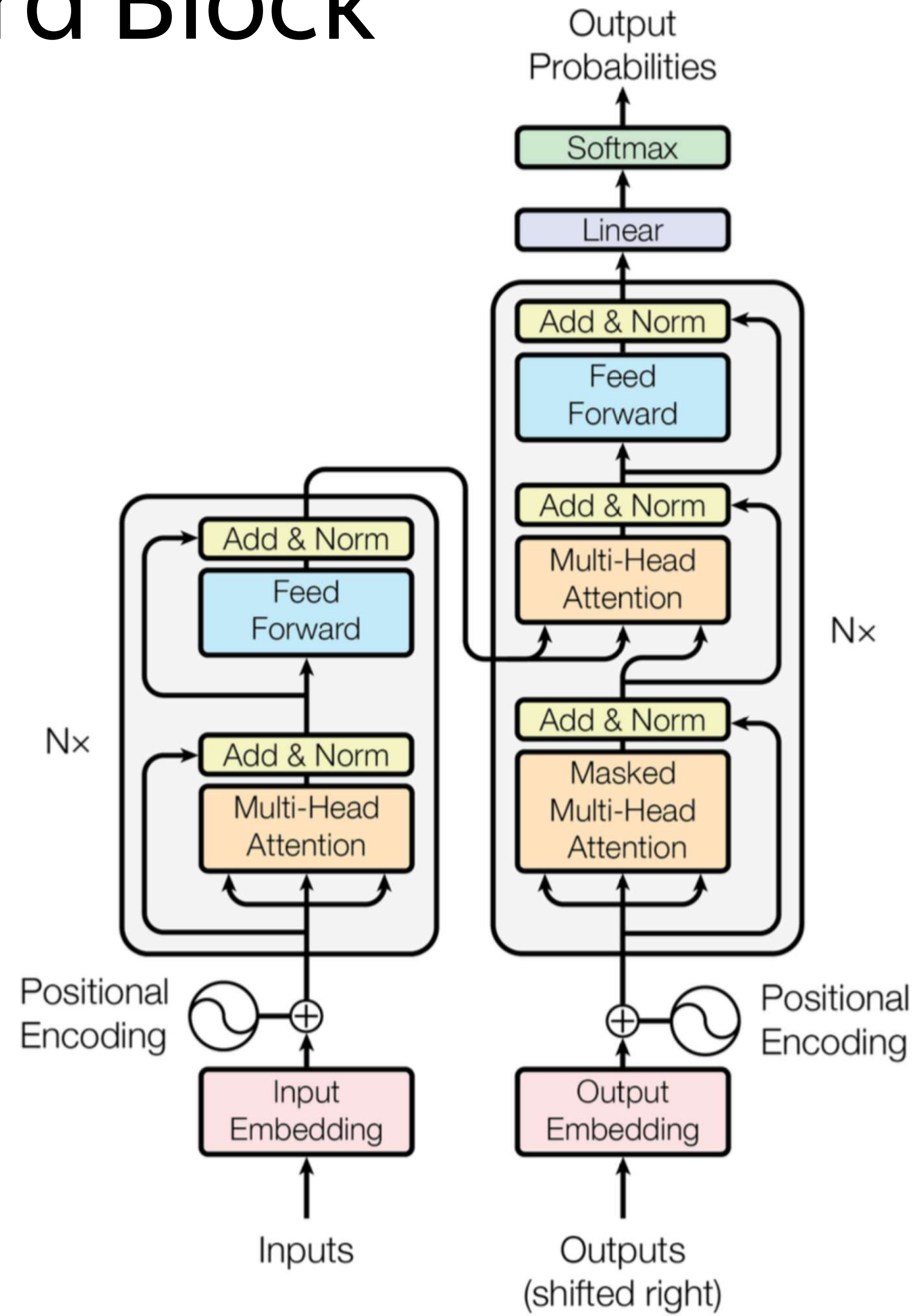
queries, keys, values are computed from decoder states
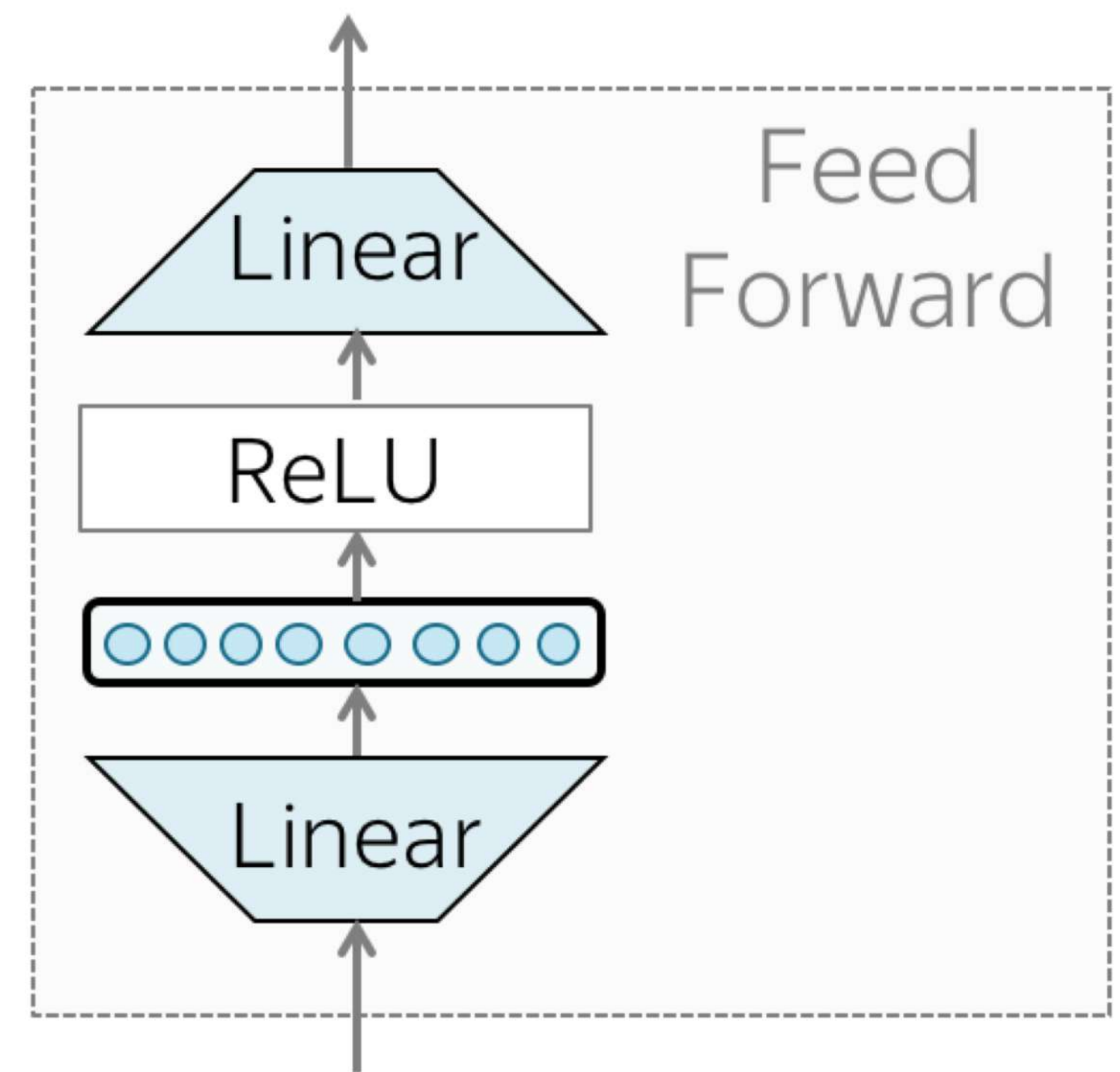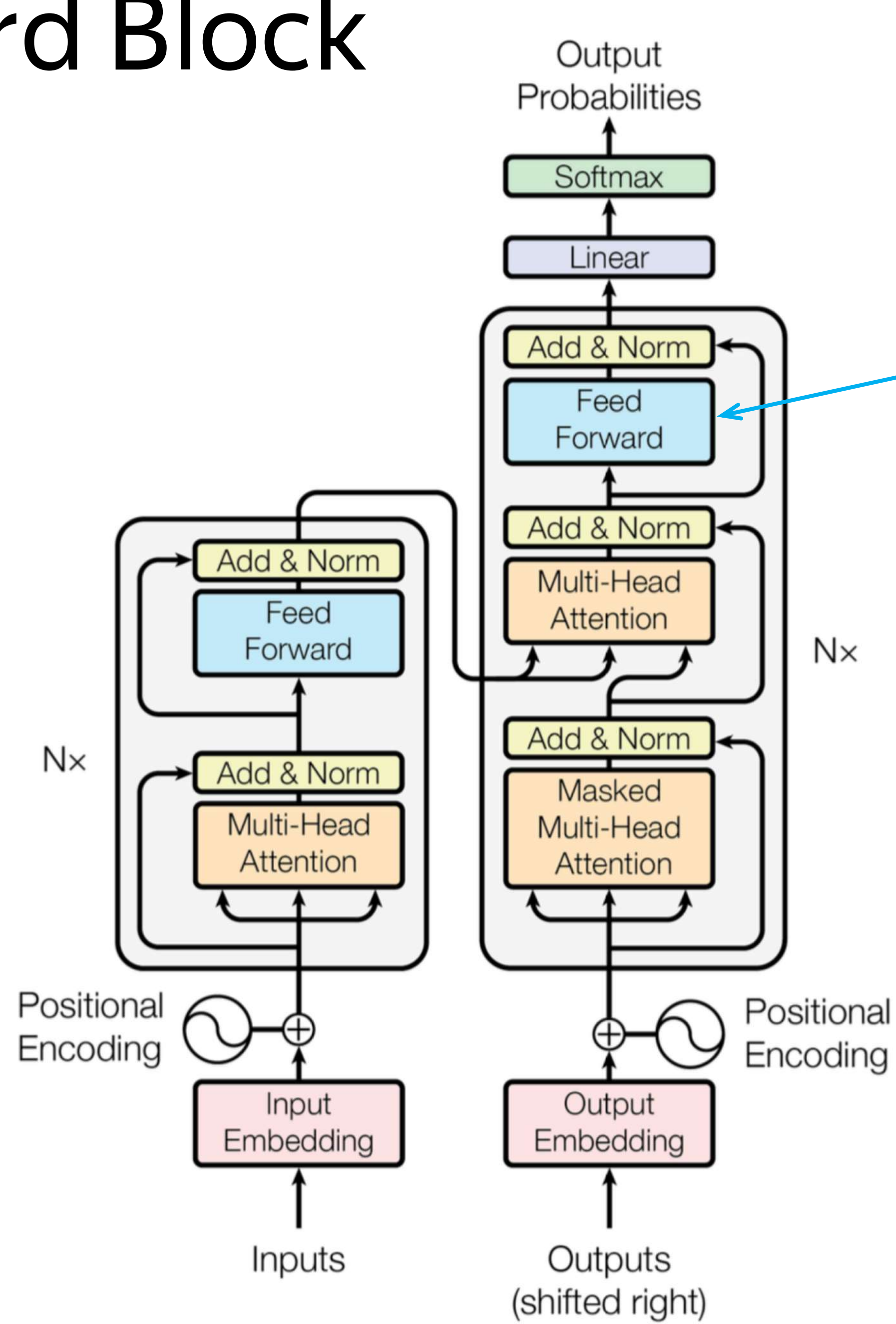
# Transformer



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Nx

Positional Encoding

Output Embedding

Positional Encoding

Outputs (shifted right)

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Positional Encoding

Input Embedding

Inputs

**Feed-forward network:** after taking information from other tokens, take a moment to think and process this information

**Decoder-encoder attention:** target token looks at the source

queries – from decoder states; keys and values from encoder states

**Decoder self-attention (masked):** tokens look at the previous tokens

queries, keys, values are computed from decoder states

**Feed-forward network:** after taking information from other tokens, take a moment to think and process this information

**Encoder self-attention:** tokens look at each other

queries, keys, values are computed from encoder states

# Feed Forward Block

# Feed Forward Block
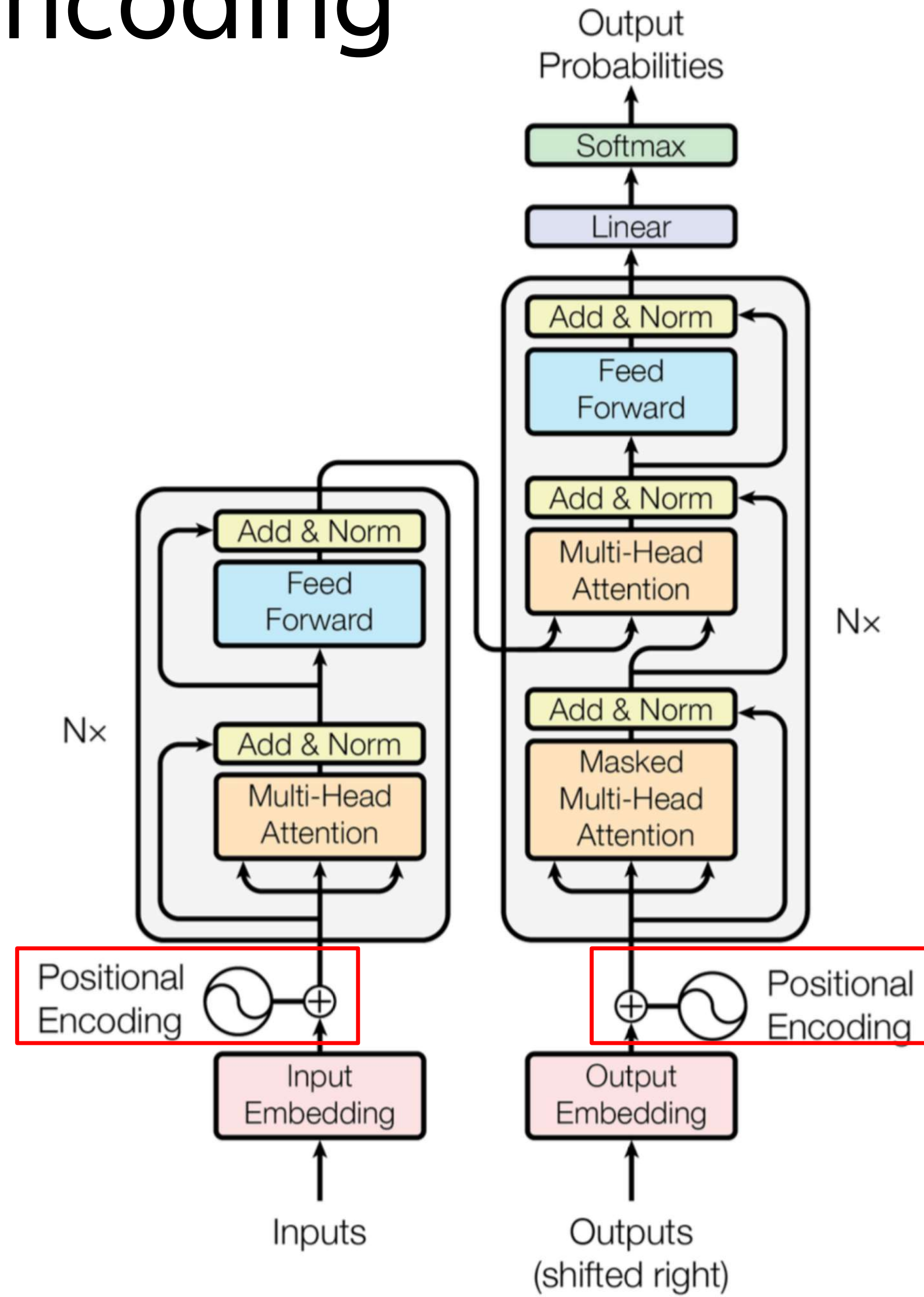


$$FFN(x) = \max(0, xW_1 + b_1)W_2 + b_2$$

# Residual Connections



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Nx

Add & Norm

Feed Forward

Nx

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

output

**block** with some layers

input

Residual connection: add a block's input to its output

# Layer Normalization



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Masked Multi-Head Attention

N×

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

N×

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Layer Normalization

$$\vdots \quad \cdot scale + bias \quad \vdots$$

$$\vdots \quad \frac{1}{\sqrt{\sigma_k}} \cdot \left[ \quad - \mu_k \right] = \quad \vdots$$

$$\ldots \quad h_{k-1} \qquad h_k \qquad h_{k+1} \quad \ldots$$

# Positional Encoding

# Positional Encoding

Problem:

- without recurrence or convolution, the model knows nothing about position

Solution:

- encode position explicitly and add

...

Encoder

"token **x** on position **k**" →

Input is <u>sum of two embeddings</u>: for token and position →

tokens →

Я
"I"    видел
       "saw"    котю
                "cat"    <eos>

positions → 0  1  2  3

# Positional Encoding

Fixed encodings:

$$\text{PE}_{pos,2i} = \sin(pos/10000^{2i/d_{model}}),$$

$$\text{PE}_{pos,2i+1} = \cos(pos/10000^{2i/d_{model}})$$

pos – position, i - dimension

"token **x** on position **k**" →

Input is <u>sum of two embeddings</u>: for token and position →

tokens →   Я        видел    котю    <eos>
           "I"      "saw"    "cat"

positions →   0        1        2        3

...

Encoder

# Transformer

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer

- Subword Segmentation: BPE

- Analysis and Interpretability

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer

- Subword Segmentation: BPE   →   NLP Course For You: read here

- Analysis and Interpretability

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer

- Subword Segmentation: BPE

- 🔍 Analysis and Interpretability

# Analysis Methods

Model-specific:

- Looking at model components
- …

Model-agnostic:

# Analysis Methods

Model-specific:

- Looking at model components
- …

In the previous lectures:

- Convolutional filters of classifiers

- Neurons in RNN/CNN LMs

Model-agnostic:

# Analysis Methods

Model-specific:

- Looking at model components
- ...

In the previous lectures:

- Convolutional filters of classifiers
- Neurons in RNN/CNN LMs

Today:

- Heads in Multi-Head Attention

Model-agnostic:

# Analysis Methods

Model-specific:

- Looking at model components
- ...

In the previous lectures:

- Convolutional filters of classifiers
- Neurons in RNN/CNN LMs

Today:

- Heads in Multi-Head Attention

Model-agnostic:

In the previous lecture:

- Look at the predictions: contrastive evaluation of specific phenomena

# Analysis Methods

Model-specific:

- Looking at model components
- ...

In the previous lectures:

- Convolutional filters of classifiers
- Neurons in RNN/CNN LMs

Today:

- Heads in Multi-Head Attention

Model-agnostic:

In the previous lecture:

- Look at the predictions: contrastive evaluation of specific phenomena

Today:

- Probing: What do representations capture?

# Analysis Methods

**Model-specific:**

- Looking at model components
- …

In the previous lectures:

- Convolutional filters of classifiers

- Neurons in RNN/CNN LMs

Today:

- Heads in Multi-Head Attention

**Model-agnostic:**

In the previous lecture:

- Look at the predictions: contrastive evaluation of specific phenomena

Today:

- Probing: What do representations capture?

# Multi-Head Attention: What are the Heads Doing?

Heads which on average contribute more to generated translations ("important heads") play interpretable roles:

- Positional

- Syntactic

- Attention to rare tokens

Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

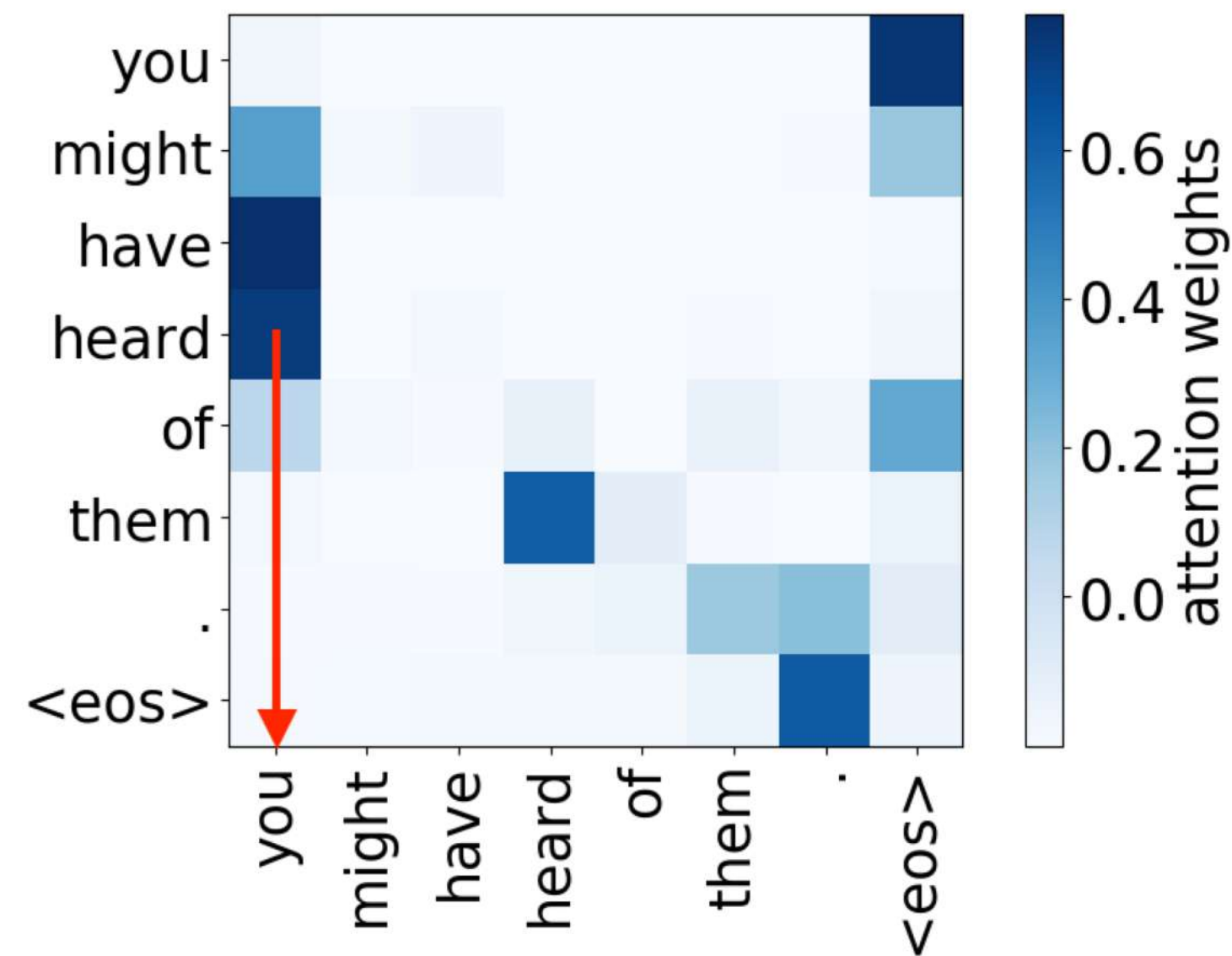# Positional Heads: Attention to Neighbors



Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

# Positional Heads: Attention to Neighbors



Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

# Syntactic Heads: Track Dependencies

- verb->subject
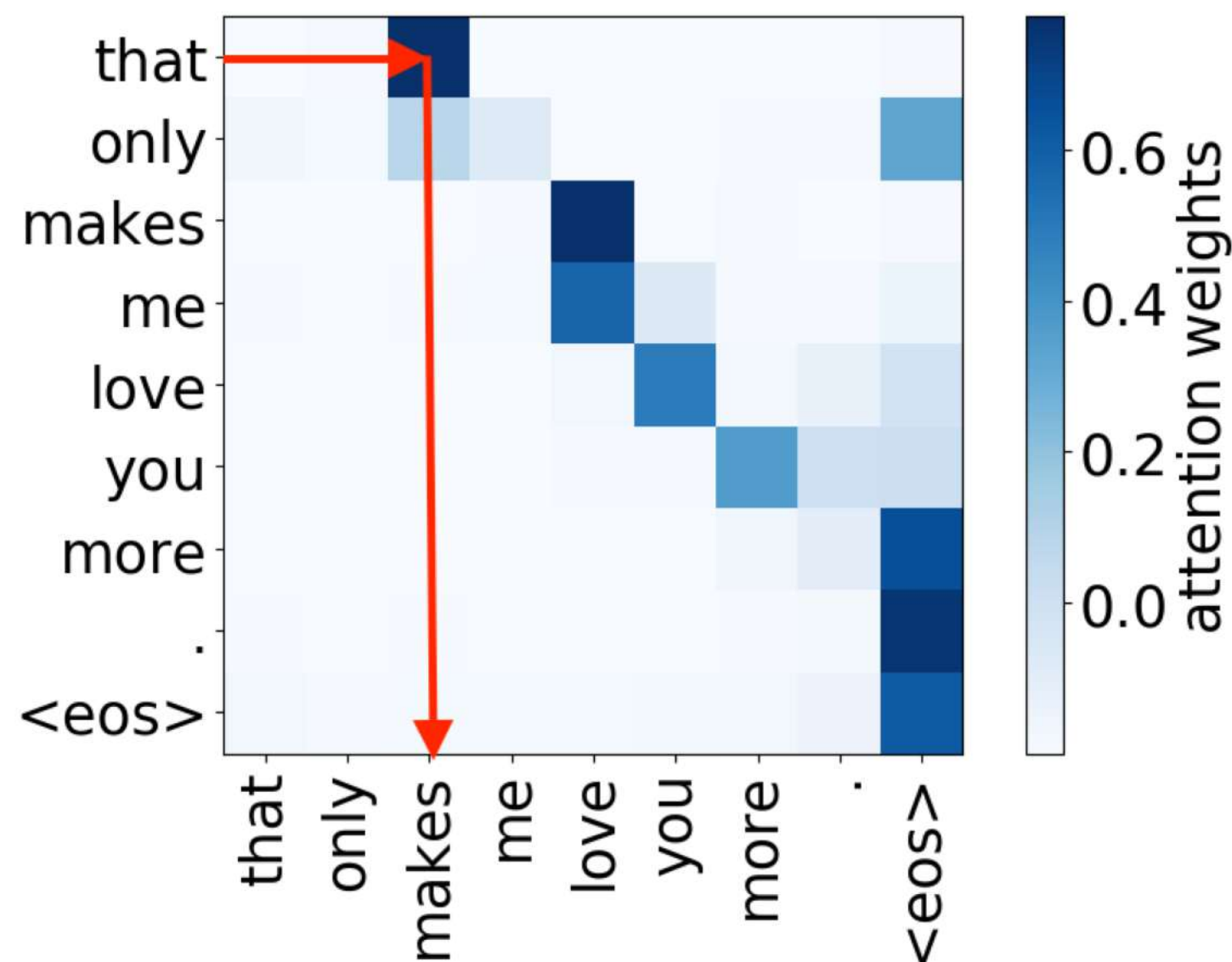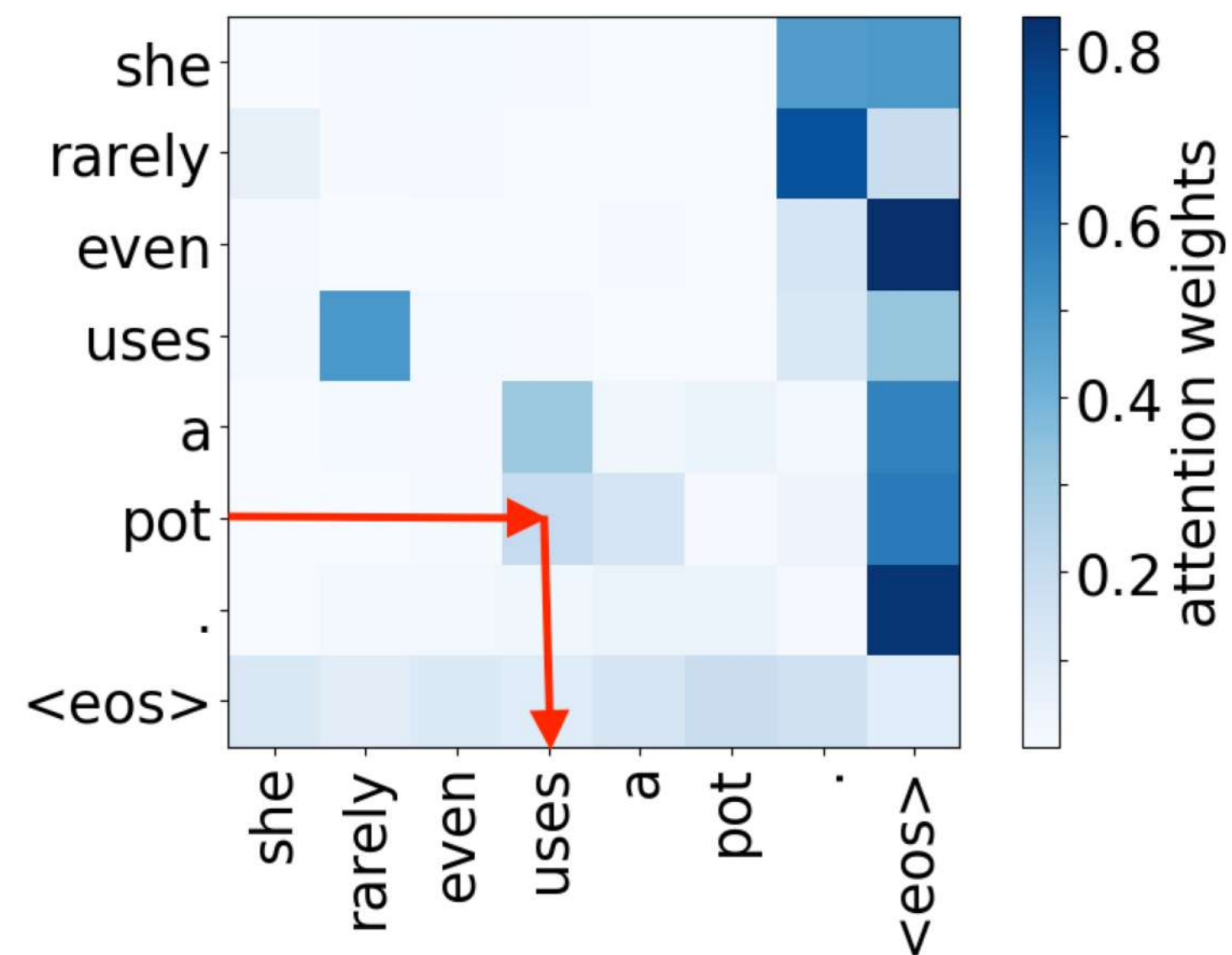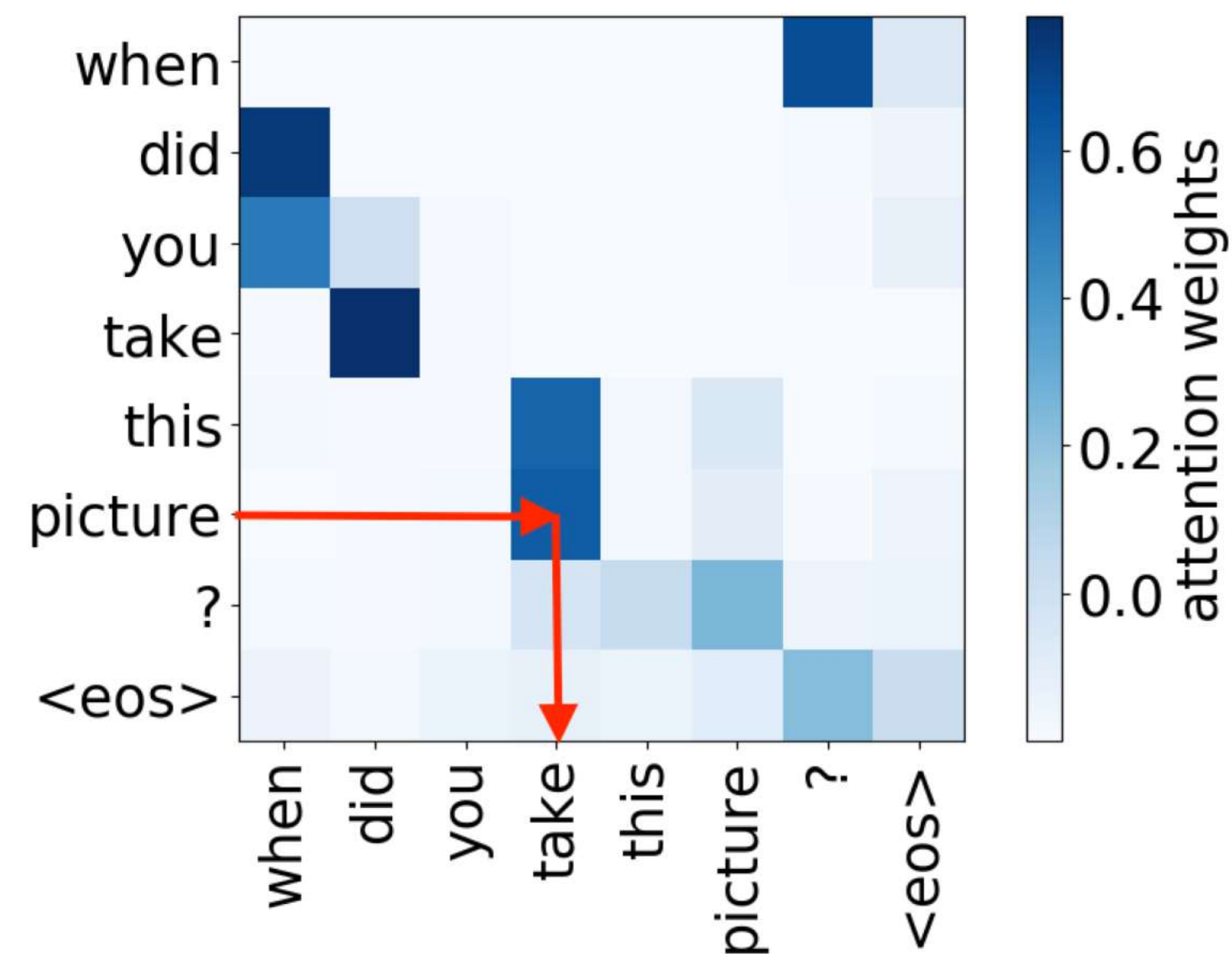
Она руководит **новым** проектом

- Gender agreement
- Case government
- Lexical preferences
- …



Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

# Syntactic Heads: Track Dependencies
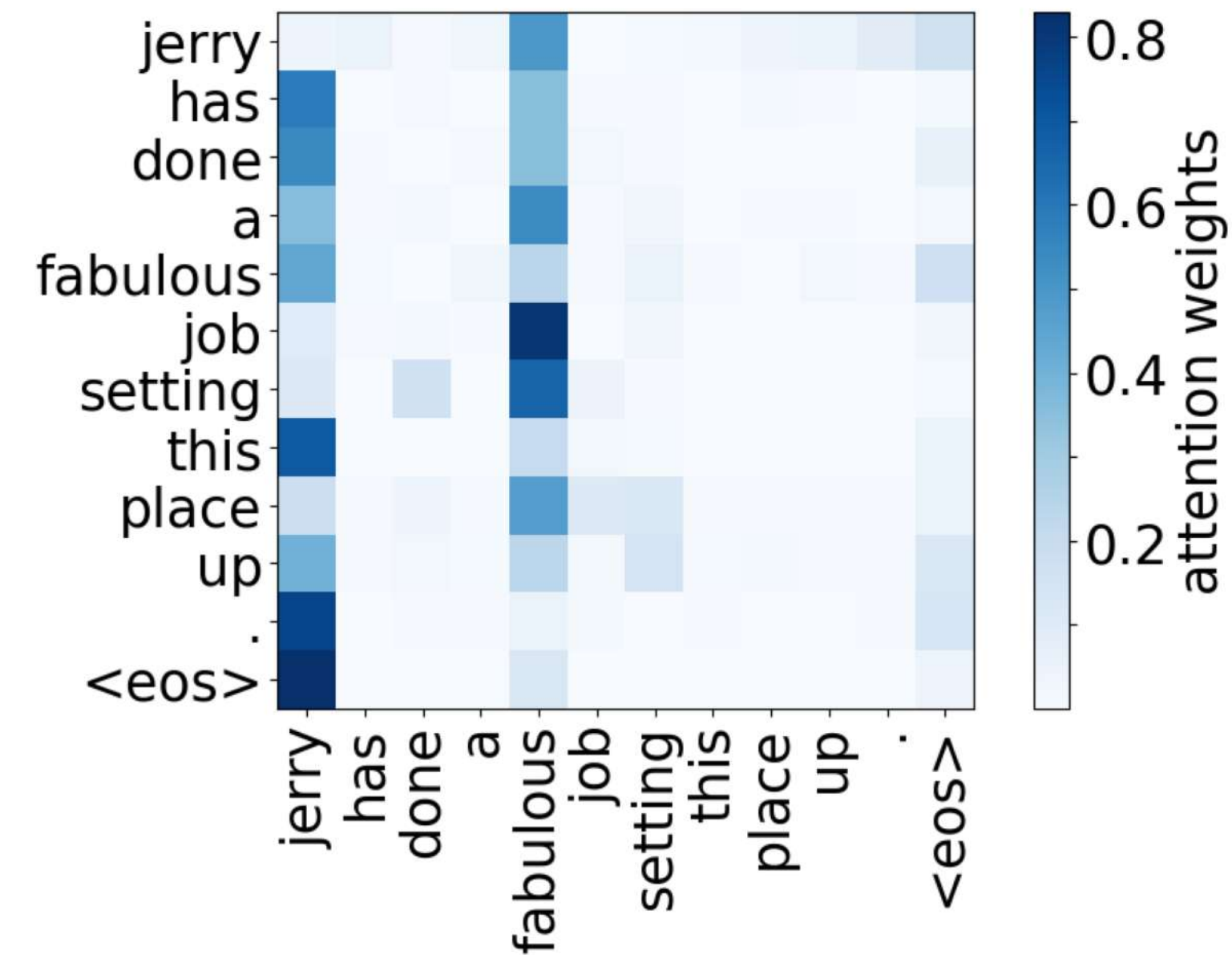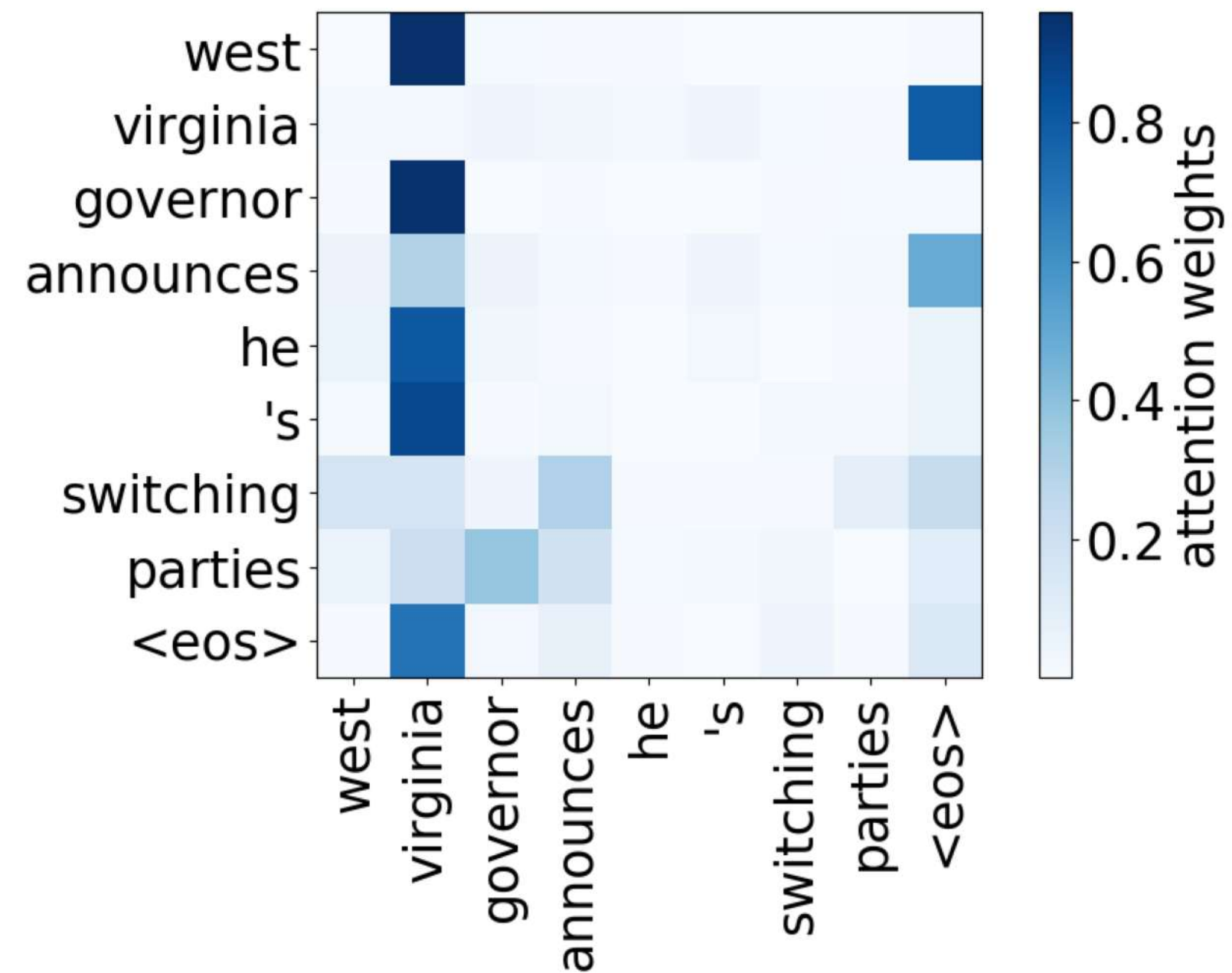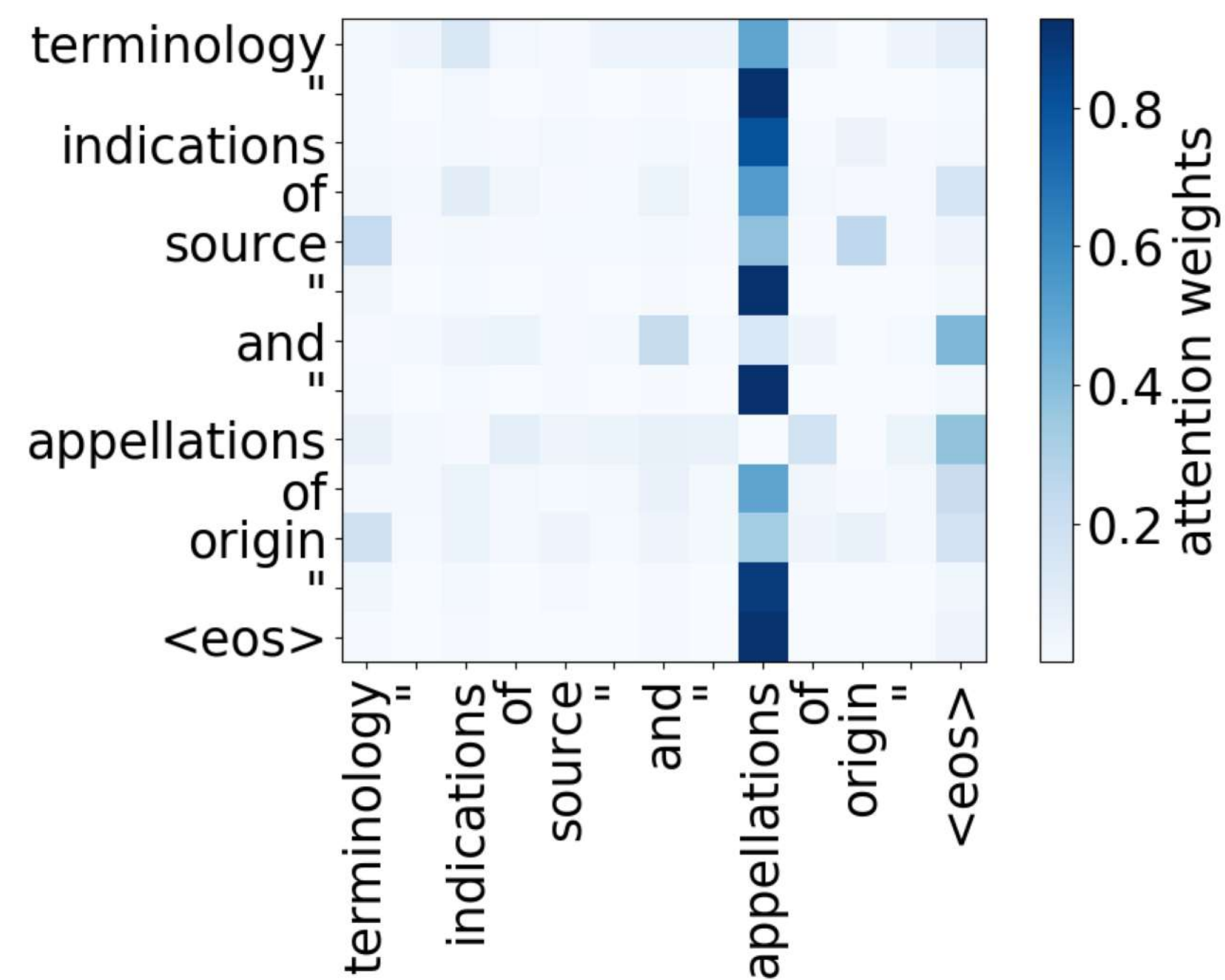
- subject->verb



Она руководит **новым** проектом

- Gender agreement
- Case government
- Lexical preferences
- ...

Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

# Syntactic Heads: Track Dependencies

- object->verb



Она руководит **новым** проектом

- Gender agreement
- Case government
- Lexical preferences
- ...

Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

# Attention to Rare Tokens



Paper: Analyzing Multi-Head Self-Attention:
Specialized Heads Do the Heavy Lifting, the Rest Can Be Pruned

# Analysis Methods

Model-specific:

- Looking at model components
- ...

In the previous lectures:

- Convolutional filters of classifiers

- Neurons in RNN/CNN LMs

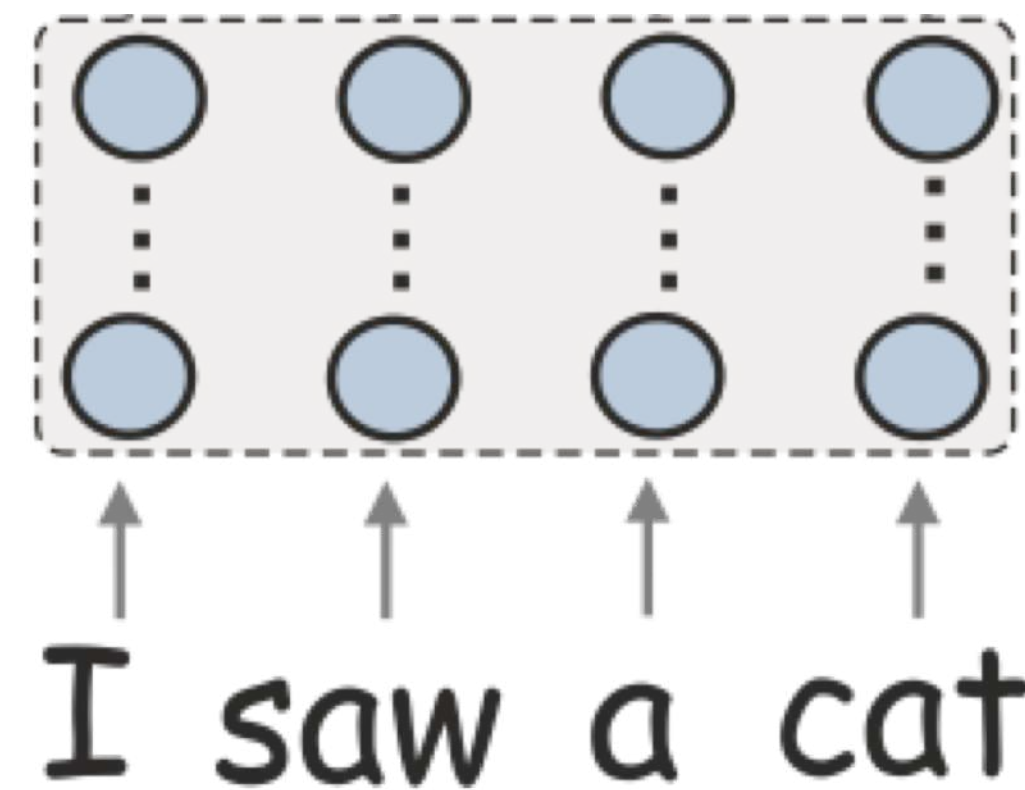Today:

- Heads in Multi-Head Attention

Model-agnostic:

In the previous lecture:

- Look at the predictions: contrastive evaluation of specific phenomena

Today:

- Probing: What do representations capture?

# Analysis Methods

## Model-specific:

- Looking at model components
- ...

In the previous lectures:

- Convolutional filters of classifiers
- Neurons in RNN/CNN LMs

Today:

- Heads in Multi-Head Attention

## Model-agnostic:

In the previous lecture:

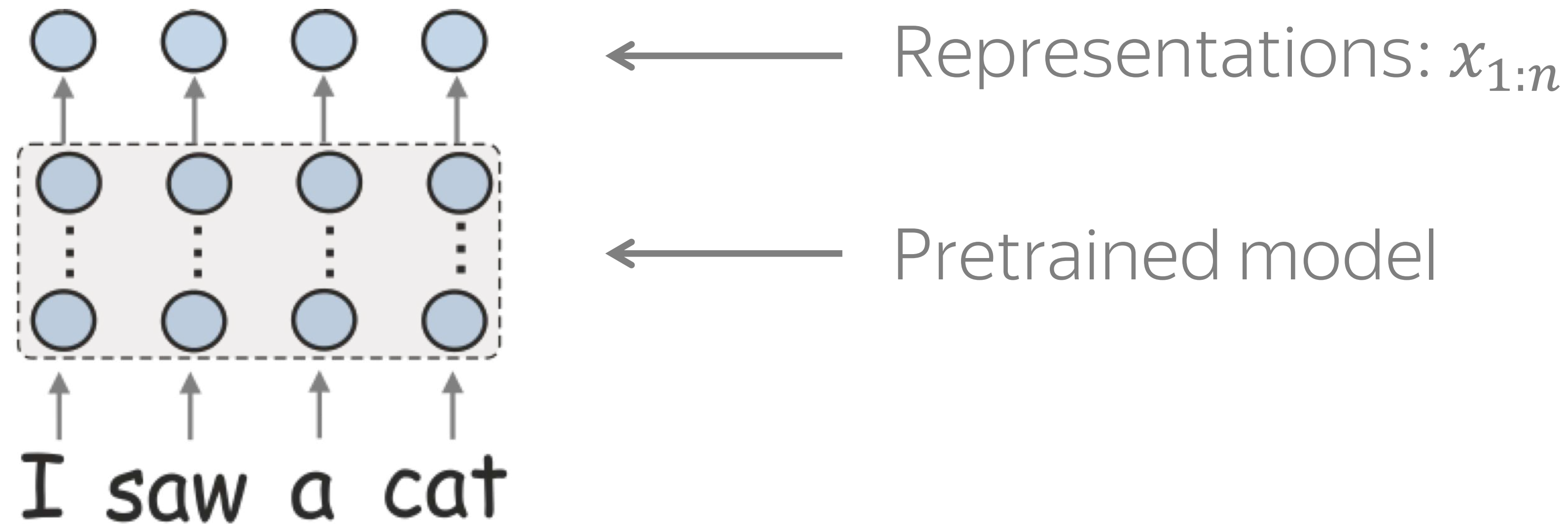- Look at the predictions: contrastive evaluation of specific phenomena

Today:

- Probing: What do representations capture?

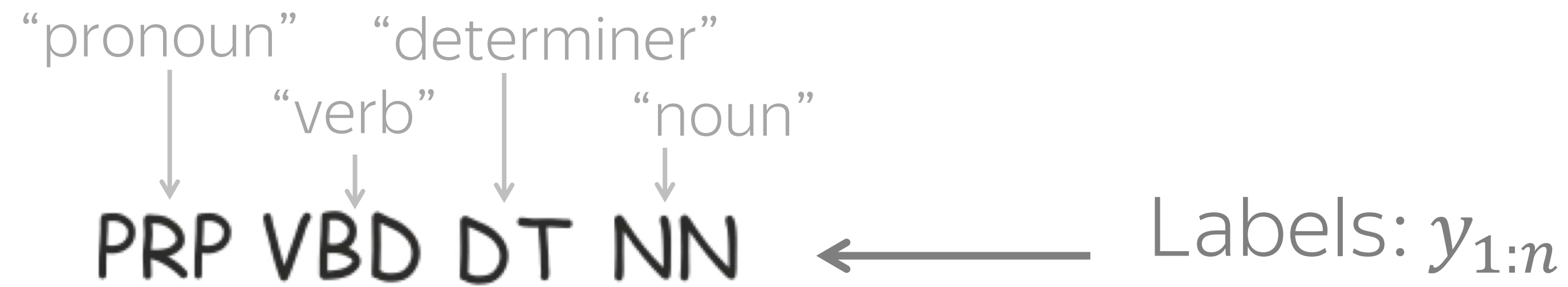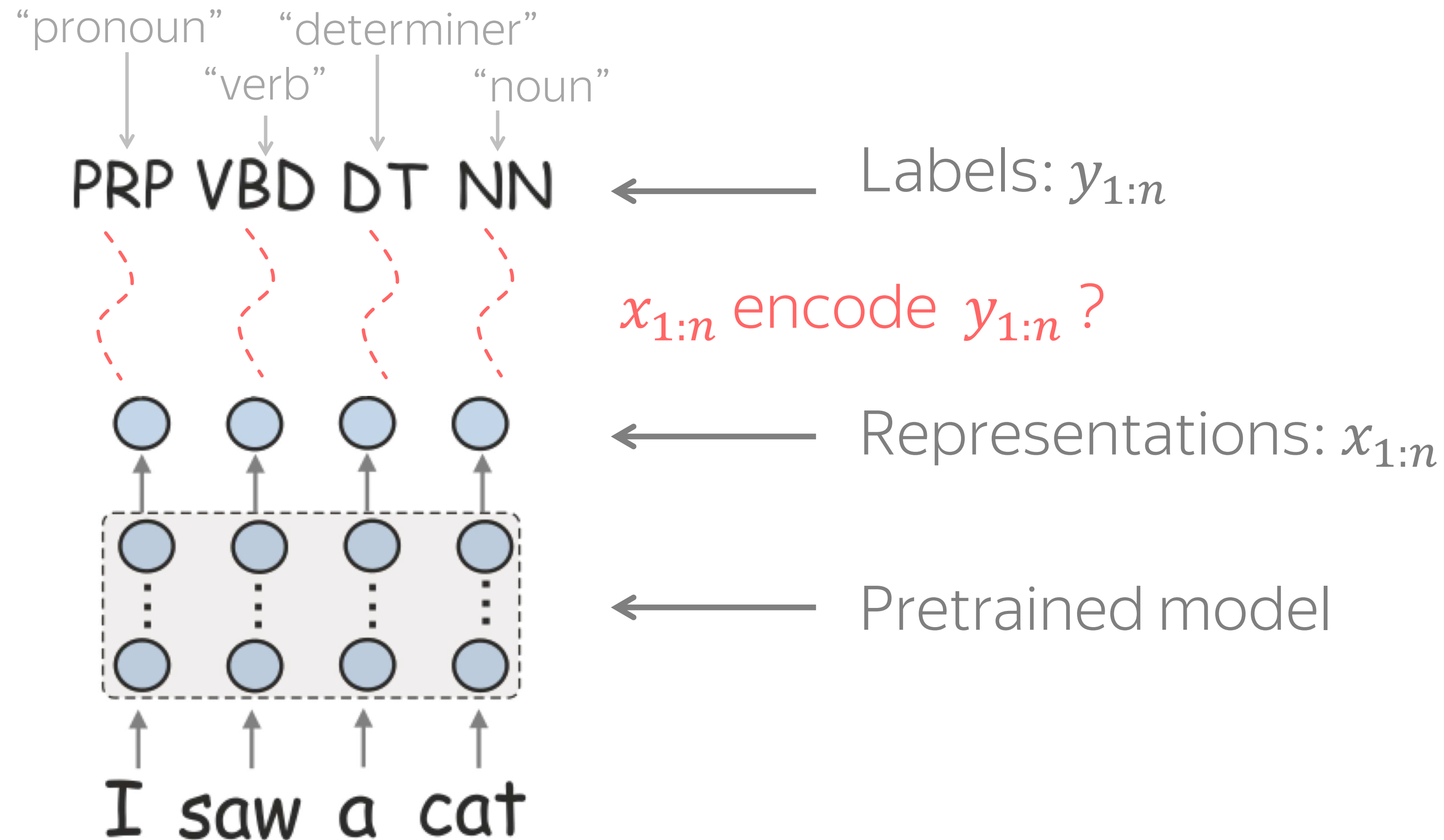# How to understand if a model captures a linguistic property?



Pretrained model

I saw a cat

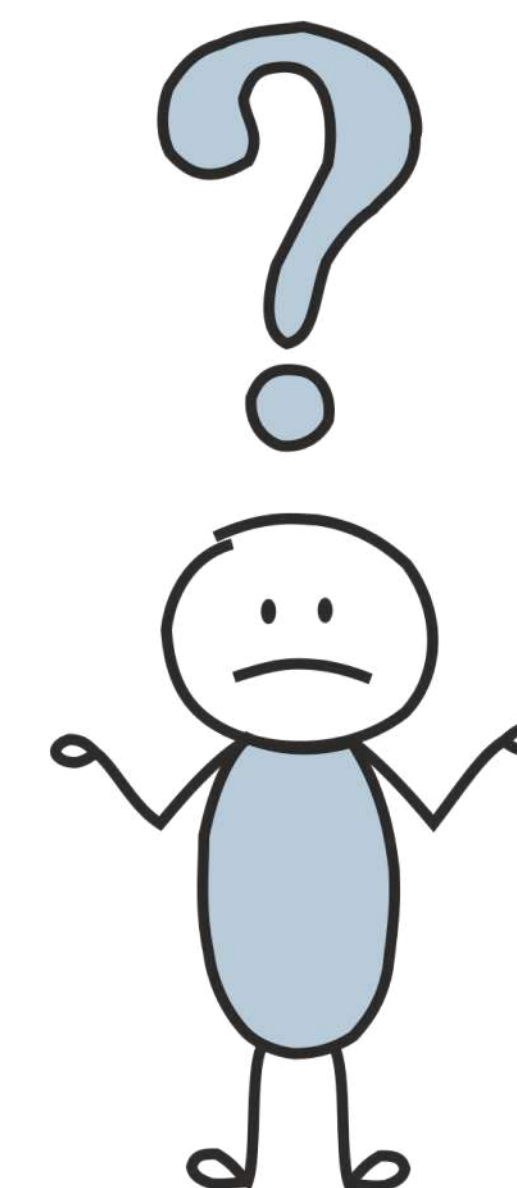# How to understand if a model captures a linguistic property?



Representations: $x_{1:n}$

Pretrained model

I saw a cat

# How to understand if a model captures a linguistic property?



"pronoun"   "determiner"
       "verb"       "noun"

PRP VBD DT NN  ⟵  Labels: $y_{1:n}$

Representations: $x_{1:n}$

Pretrained model

I saw a cat

# How to understand if a model captures a linguistic property?



"pronoun" "verb" "determiner" "noun"

PRP VBD DT NN ⟵ Labels: $y_{1:n}$

$x_{1:n}$ encode $y_{1:n}$ ?

Representations: $x_{1:n}$

Pretrained model

I saw a cat

# How to understand if a model captures a linguistic property?

# Standard Probing: train a classifier, use its accuracy



PRP VBD DT NN ← Labels: $y_{1:n}$

← Probing classifier

← Representations: $x_{1:n}$

← Pretrained model

I saw a cat

# What Do NMT Models Learn About Morphology?

- Take NMT models for different language pairs

- Look at the encoder layers



POS accuracy by Representation Layer

Legend: layer 0, layer 1, layer 2
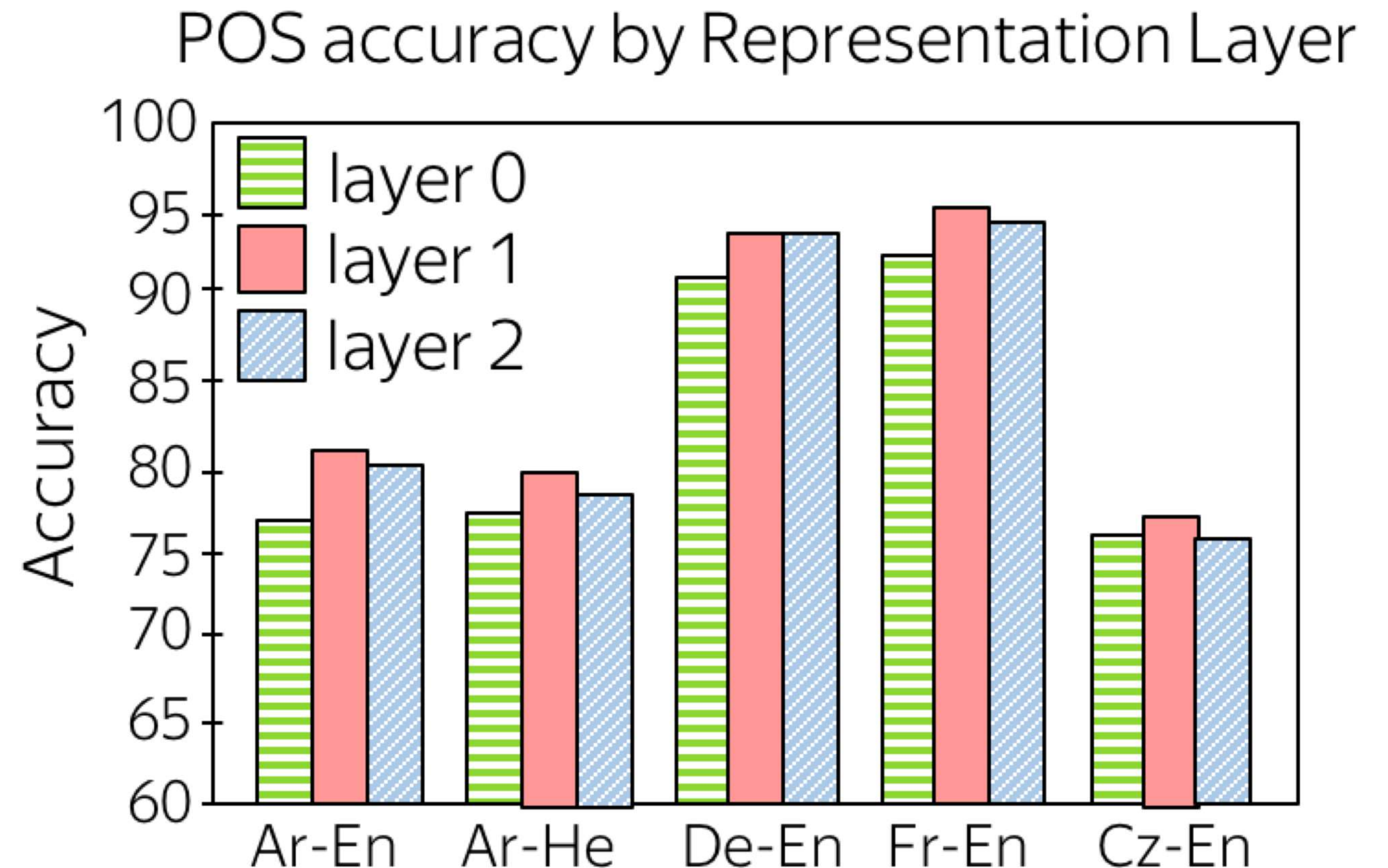
Paper: What do Neural Machine Translation Models Learn about Morphology?

# What Do NMT Models Learn About Morphology?

- Take NMT models for different language pairs

- Look at the encoder layers

Results:

- Encoding helps: layer 0 (word embeddings) is the worst



POS accuracy by Representation Layer

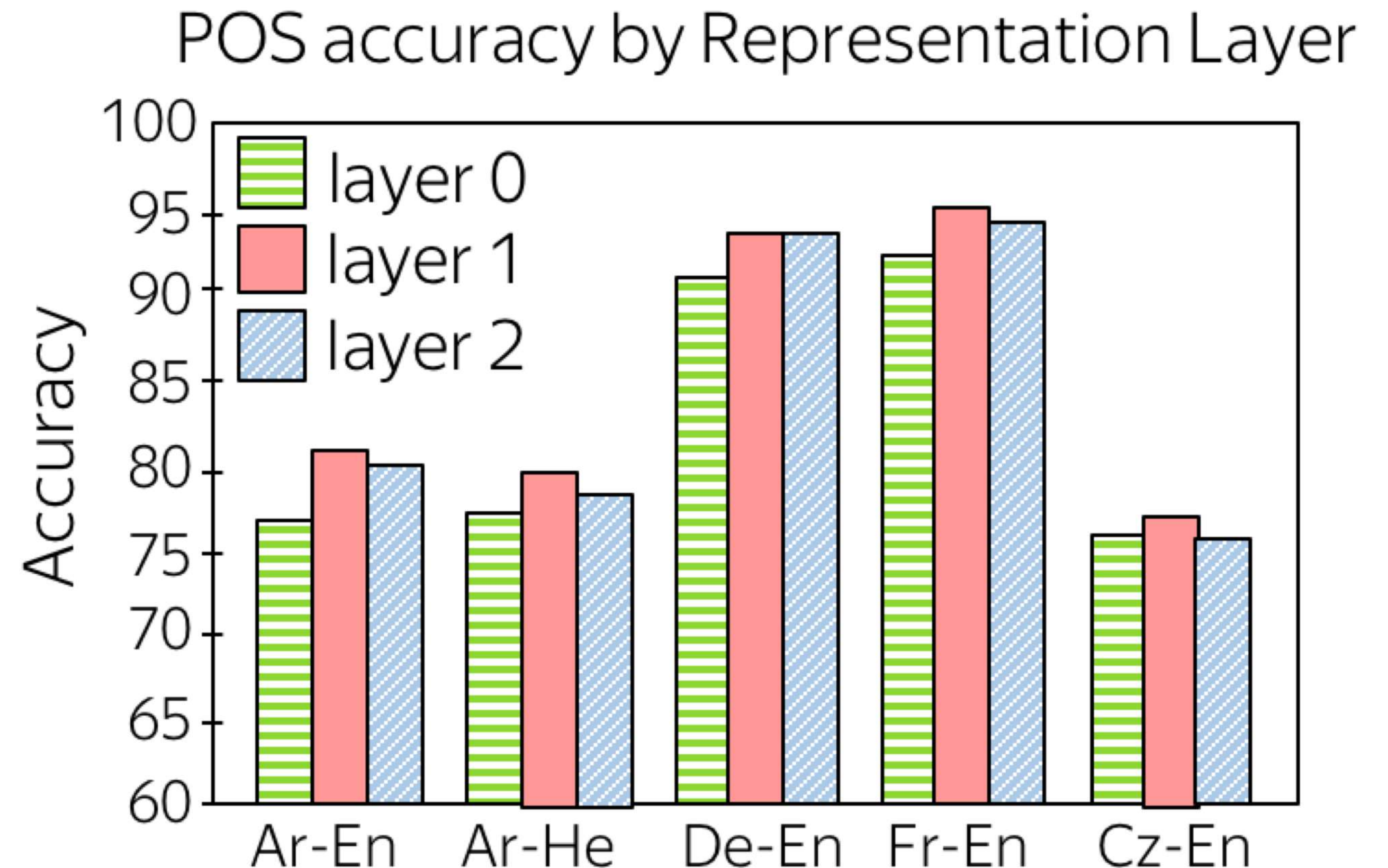Paper: <u>What do Neural Machine Translation Models Learn about Morphology?</u>

# What Do NMT Models Learn About Morphology?

- Take NMT models for different language pairs

- Look at the encoder layers

Results:

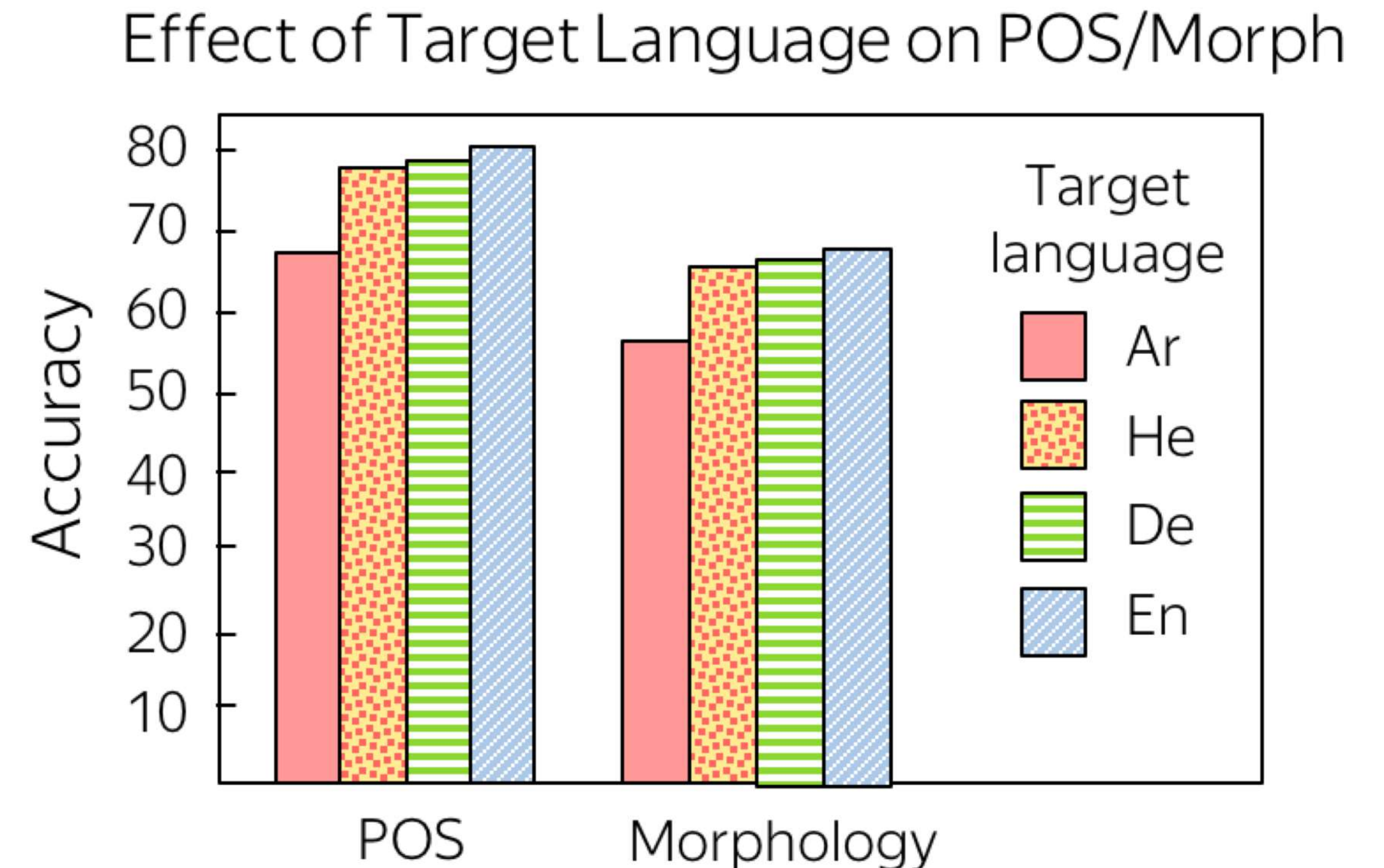- Encoding helps: layer 0 (word embeddings) is the worst

- Layer 1 is better than layer 2



POS accuracy by Representation Layer

Paper: <u>What do Neural Machine Translation Models Learn about Morphology?</u>

# What Do NMT Models Learn About Morphology?

- Take NMT models with the same source language and different target languages

- Look at the encoder
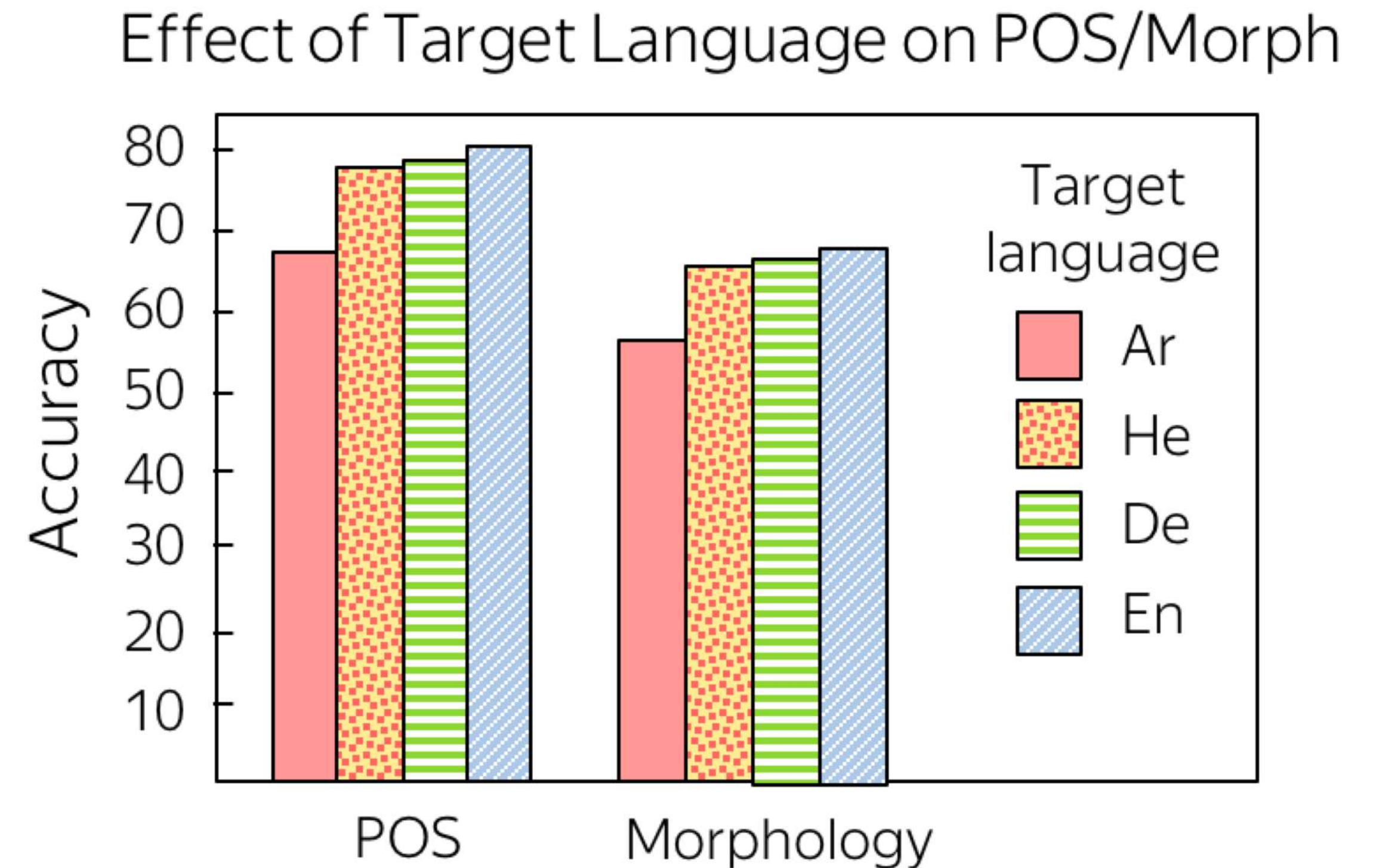


Effect of Target Language on POS/Morph

Paper: <u>What do Neural Machine Translation Models Learn about Morphology?</u>

# What Do NMT Models Learn About Morphology?

- Take NMT models with the same source language and different target languages

- Look at the encoder

Results:

- Surprising: weaker target morphology leads to stronger encoder



Effect of Target Language on POS/Morph

Paper: What do Neural Machine Translation Models Learn about Morphology?

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer

- Subword Segmentation: BPE

- Analysis and Interpretability

# What is going to happen:

- Seq2seq Basics

- Attention

- Transformer

- Subword Segmentation: BPE

-  Analysis and Interpretability

 Learn more in the NLP Course For You → This is up to You!

# Thank you!

Lena Voita

PhD student, Uni Edinburgh & Uni Amsterdam

Facebook PhD Fellow in NLP

✉  lena-voita@hotmail.com

🗔  https://lena-voita.github.io

🐦  @lena_voita

🐱  lena-voita