CSCI 544: Applied Natural Language Processing
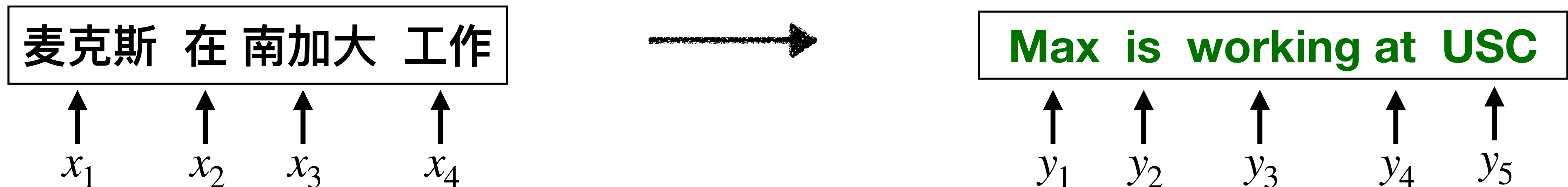
# Self-attention & Transformer

Xuezhe Ma (Max)

# Logistic Points

- **Midterm:**
  - In-person unless approved to participate remotely
  - Remote students: sign up in this Google Sheet https://docs.google.com/spreadsheets/d/1BjmjBs6tceW5RHLj1-nskg9m79Iimq2BzJFVxDYRXcE/edit#gid=0
  - We will email questions to students only in this list
  - Camera and microphone should be open with no virtual background (please find a quite location)
  - Please follow the discussions on the Slack channel, which are sometimes useful

# Recap: Seq2seq Generation

- **Sequence-to-Sequence (Seq2seq) Generation**
  - Input: $X = \{x_1, x_2, \ldots, x_L\}, x_i \in \mathcal{X}$
  - Output: $Y = \{y_1, y_2, \ldots, y_T\}, y_i \in \mathcal{Y}$
  - Model: $p_\theta(Y|X)$
- **Difference from Sequence Labeling**
  - The length of $Y$ can be different from the length of $X$
  - The size of $\mathcal{Y}$ is often much larger

| 麦克斯 在 南加大 工作 |
|---|

$\longrightarrow$

| Max is working at USC |
|---|

$x_1 \quad x_2 \quad x_3 \quad x_4$

$y_1 \quad y_2 \quad y_3 \quad y_4 \quad y_5$

# Recap: Autoregressive Seq2seq Generation

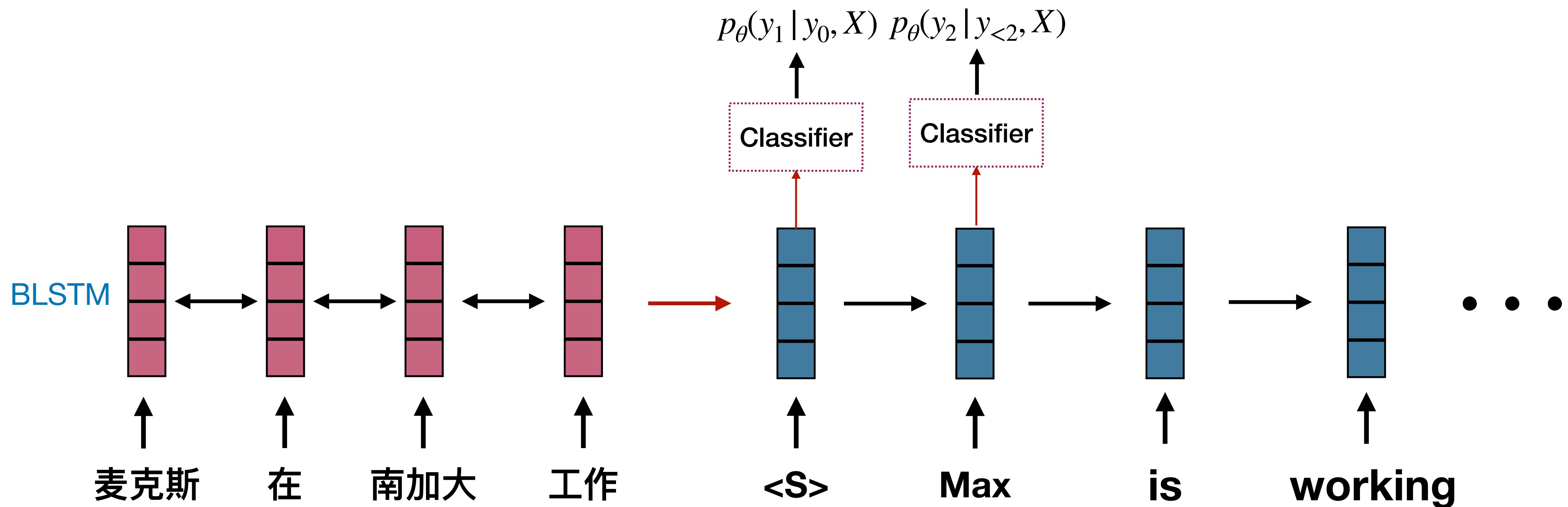- **Autoregressive Factorization:**

$$p_\theta(Y|X) = \prod_{t=1}^{T} p_\theta(y_t | y_{<t}, X)$$

Next Token

history

- Autoregressive factorization is just chain-rule (HMMs, MEMMs)
- Autoregressive factorization does NOT assume any independence
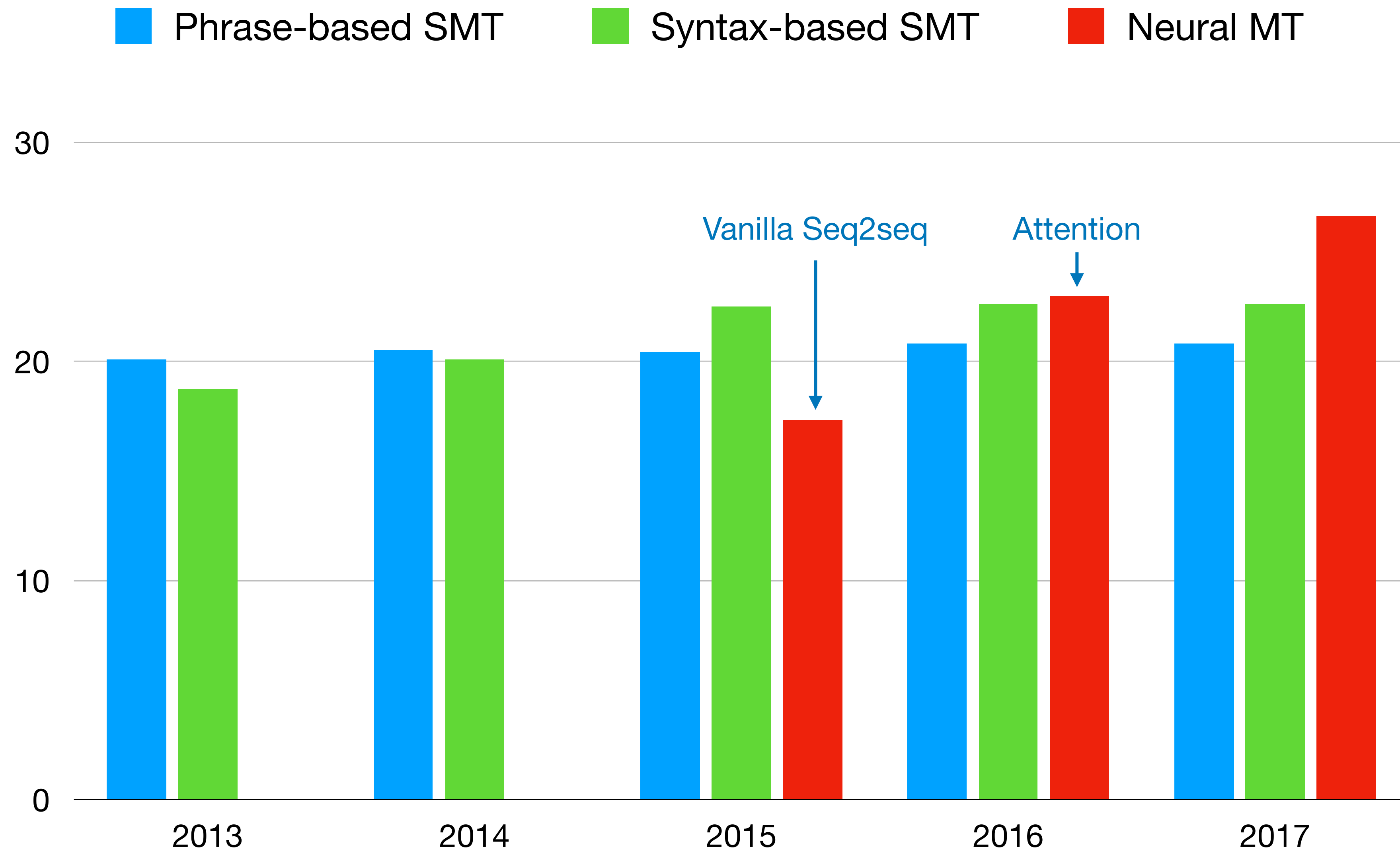- With autoregressive factorization, we need to model each $p_\theta(y_t | y_{<t}, X)$

# Recap: Encoder-Decoder Architecture
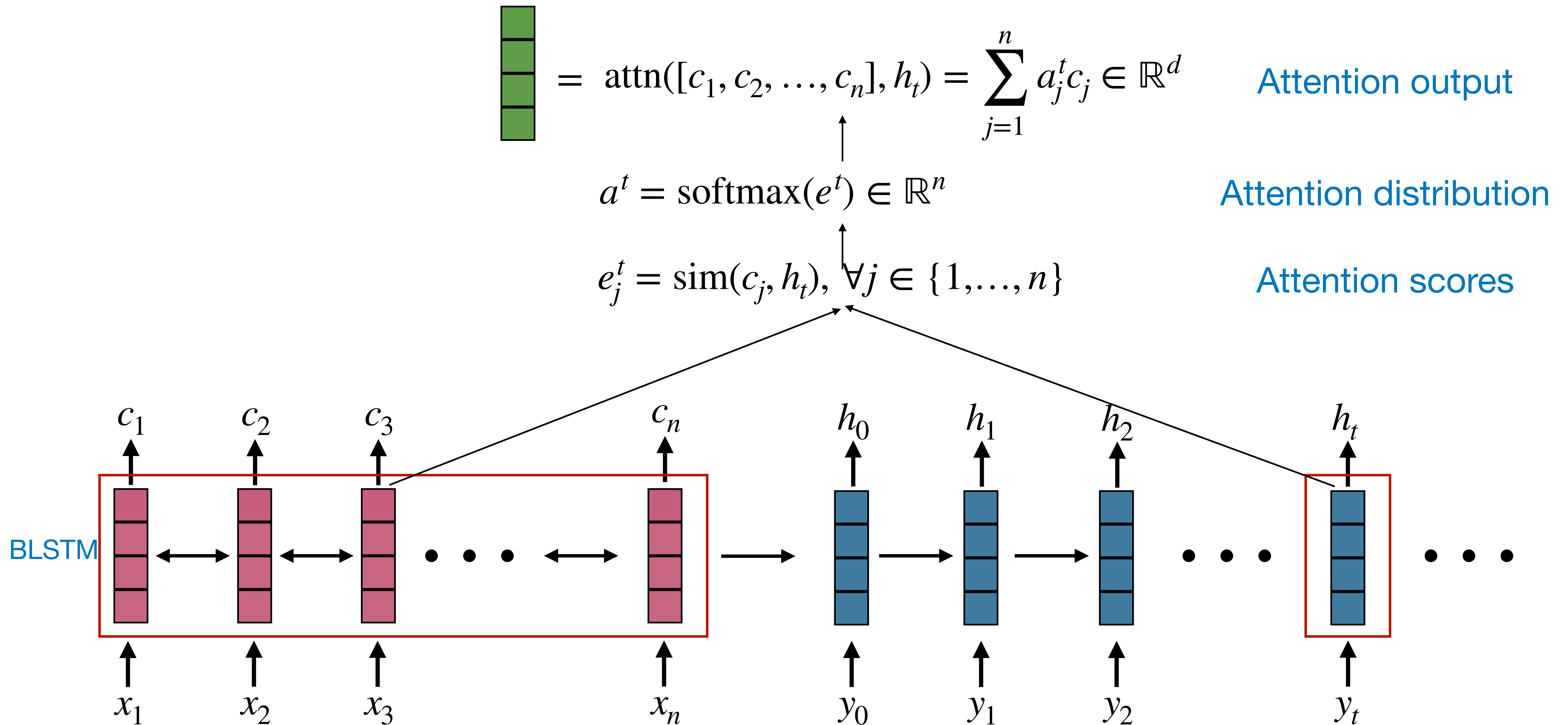
- **Two Components:**
  - Encoder: Convert input sequence into a sequence of vectors
  - Decoder: Convert encoding into a sequence in the output space



$$p_\theta(y_1 | y_0, X) \quad p_\theta(y_2 | y_{<2}, X)$$

# Recap: MT Progress

# Recap: Attention Mechanism

$$\square = \text{attn}([c_1, c_2, \ldots, c_n], h_t) = \sum_{j=1}^{n} a_j^t c_j \in \mathbb{R}^d$$

Attention output

$$a^t = \text{softmax}(e^t) \in \mathbb{R}^n$$

Attention distribution

$$e_j^t = \text{sim}(c_j, h_t), \ \forall j \in \{1, \ldots, n\}$$

Attention scores



$c_1 \quad c_2 \quad c_3 \quad c_n \quad h_0 \quad h_1 \quad h_2 \quad h_t$

BLSTM

$x_1 \quad x_2 \quad x_3 \quad x_n \quad y_0 \quad y_1 \quad y_2 \quad y_t$

# Recap: Types of Attention

- **Dot-product attention** (assumes equal dimensions for $c$ and $h$)

$$\text{sim}(c_j, h_t) = c_j^T h_t$$

- **Multiplicative attention**

$$\text{sim}(c_j, h_t) = c_j^T W h_t, \text{ where } W \text{ is learnable weight matrix}$$

- **Additive attention**

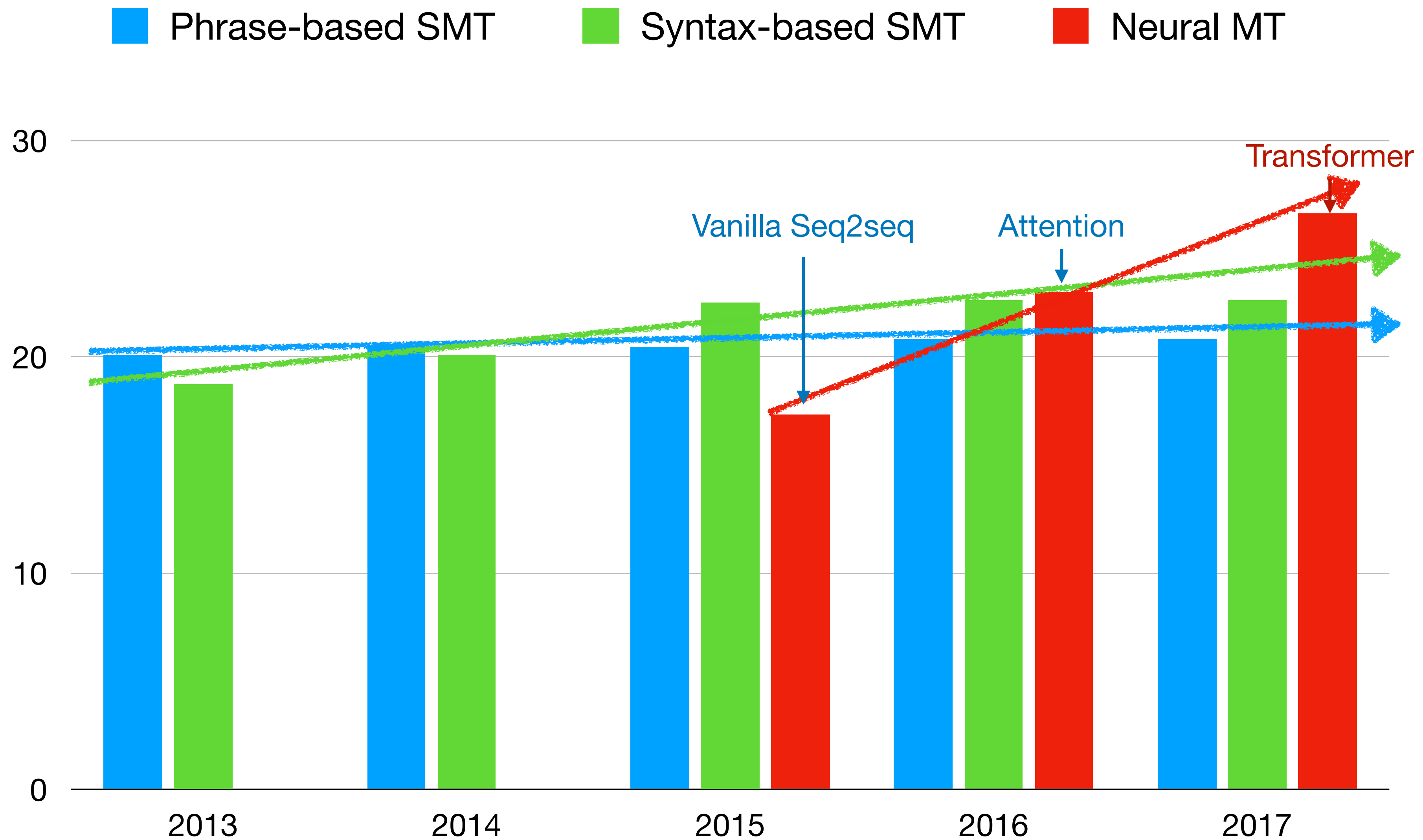$$\text{sim}(c_j, h_t) = v^T \tanh(W_c c_j + W_h h_t)$$

where $W_c$ and $W_h$ are learnable weight matrices and $v$ is a learnable weight vector

# Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation

Table 10: Mean of side-by-side scores on production data

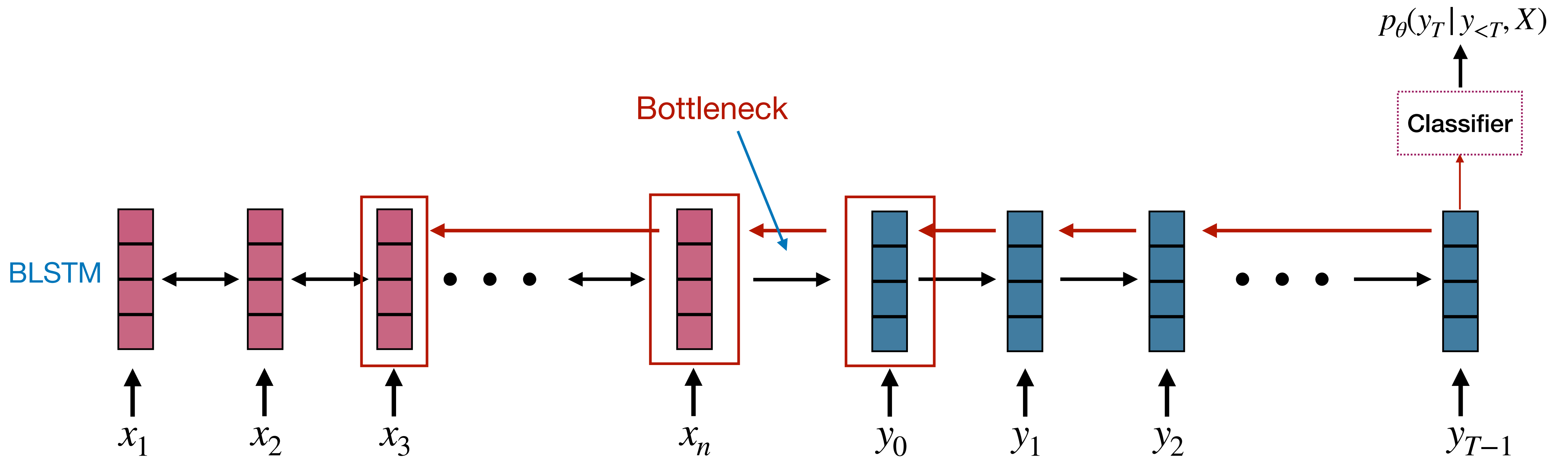|  | PBMT | GNMT | Human | Relative Improvement |
|---|---|---|---|---|
| English → Spanish | 4.885 | 5.428 | 5.504 | 87% |
| English → French | 4.932 | 5.295 | 5.496 | 64% |
| English → Chinese | 4.035 | 4.594 | 4.987 | 58% |
| Spanish → English | 4.872 | 5.187 | 5.372 | 63% |
| French → English | 5.046 | 5.343 | 5.404 | 83% |
| Chinese → English | 3.694 | 4.263 | 4.636 | 60% |

*(Wu et al., 2016)*

9

# MT Progress



Source: Rico Sennrich

# Self-Attention & Transformer
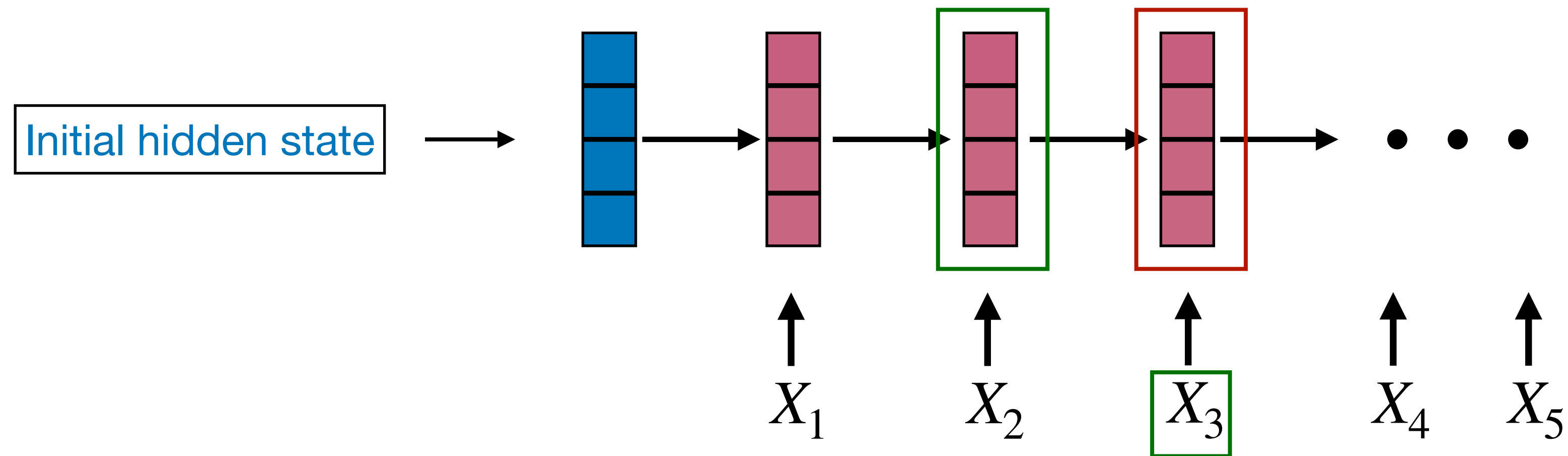
# Revisit: Motivation of Attention Mechanism

- A single encoding vector needs to capture **all the information** about source sentence
- Longer sequences can lead to **vanishing gradients**

$$p_\theta(y_T | y_{<T}, X)$$

Bottleneck

Classifier

BLSTM

$x_1$  $x_2$  $x_3$  $x_n$  $y_0$  $y_1$  $y_2$  $y_{T-1}$

Issues with RNN!

# Issues with RNNs

- **One vector** to memorize all historical information
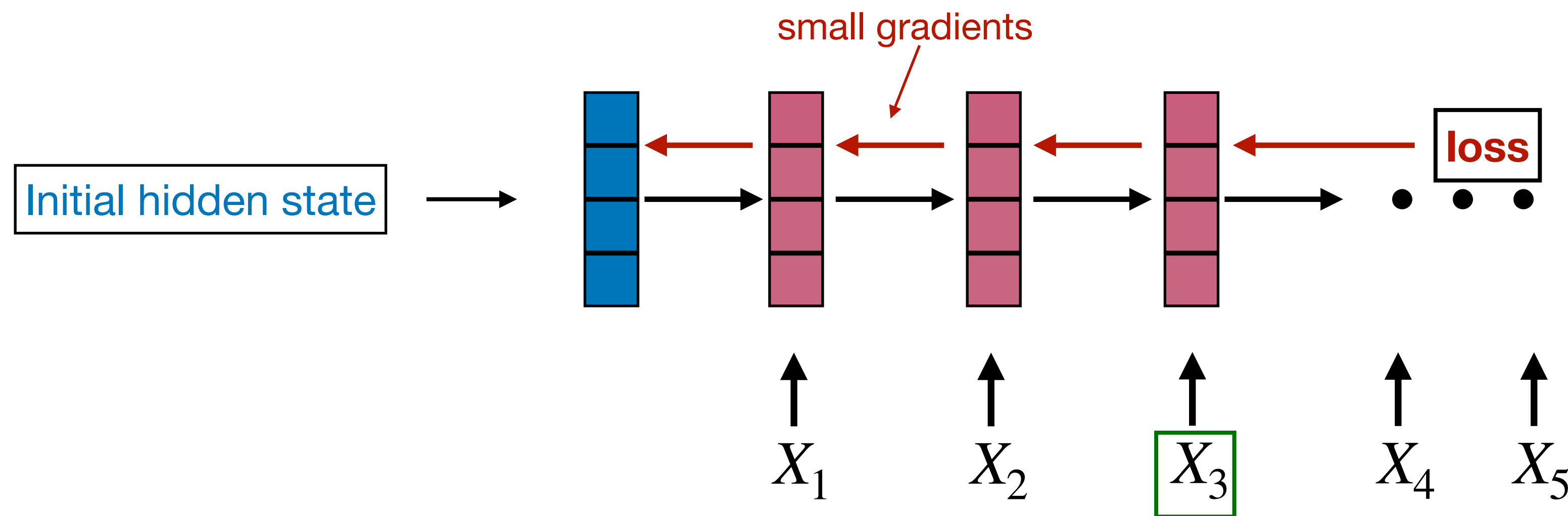
**Inherent Markov Property**

# Issues with RNNs

- **One vector** to memorize all historical information
- Hard to capture long-distance information: **vanishing gradients**



**Inherent Markov Property**

Attention is the key to solve the problem!

small gradients

Initial hidden state

loss

$X_1$  $X_2$  $X_3$  $X_4$  $X_5$

# This Lecture

- **Do we really need RNNs to model the arbitrary context?**
- **Maybe attention is all you need!**

## Attention Is All You Need

**Ashish Vaswani**[*]
Google Brain
avaswani@google.com

**Noam Shazeer**[*]
Google Brain
noam@google.com

**Niki Parmar**[*]
Google Research
nikip@google.com

**Jakob Uszkoreit**[*]
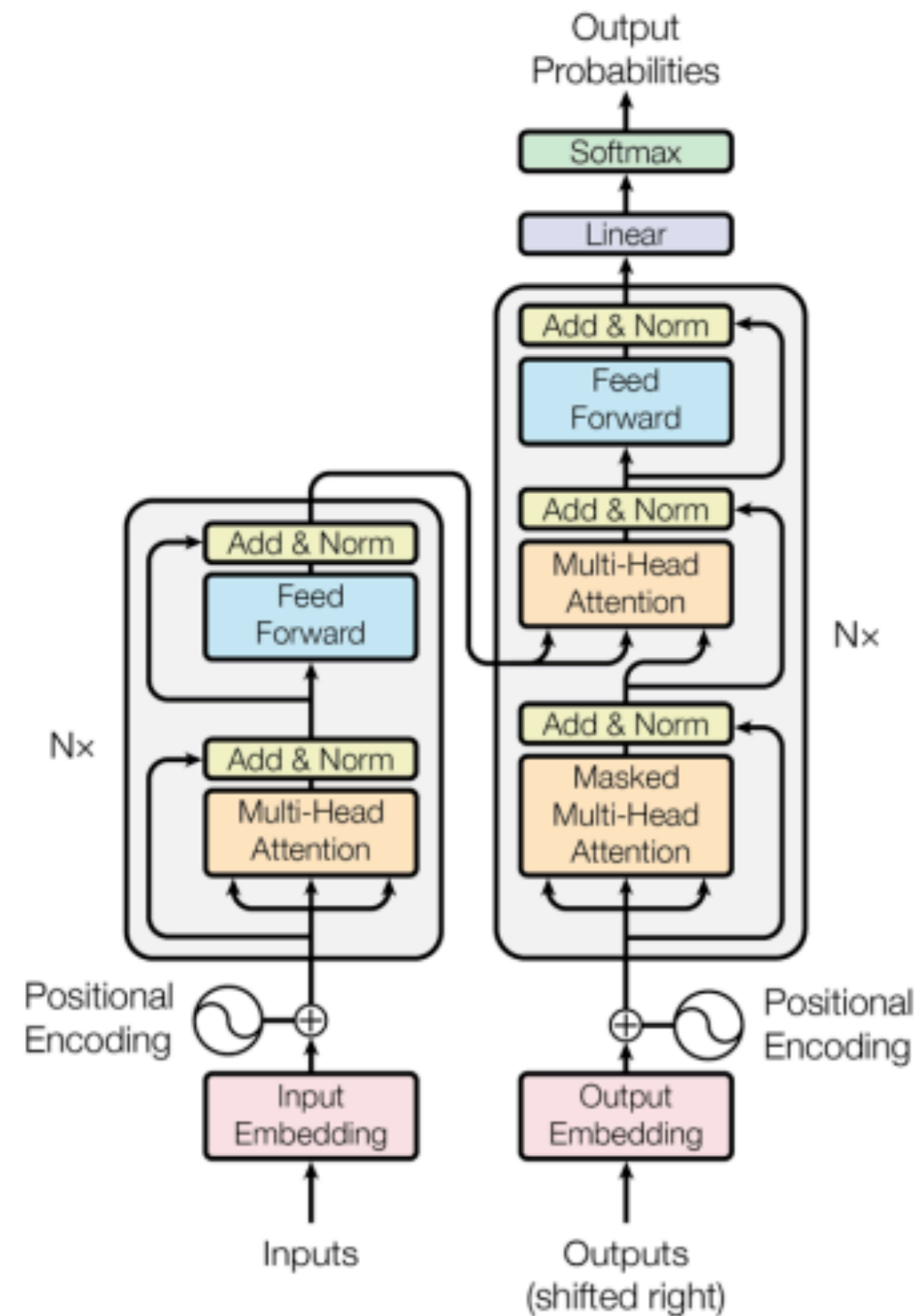Google Research
usz@google.com

**Llion Jones**[*]
Google Research
llion@google.com

**Aidan N. Gomez**[*][†]
University of Toronto
aidan@cs.toronto.edu

**Łukasz Kaiser**[*]
Google Brain
lukaszkaiser@google.com

**Illia Polosukhin**[*][‡]
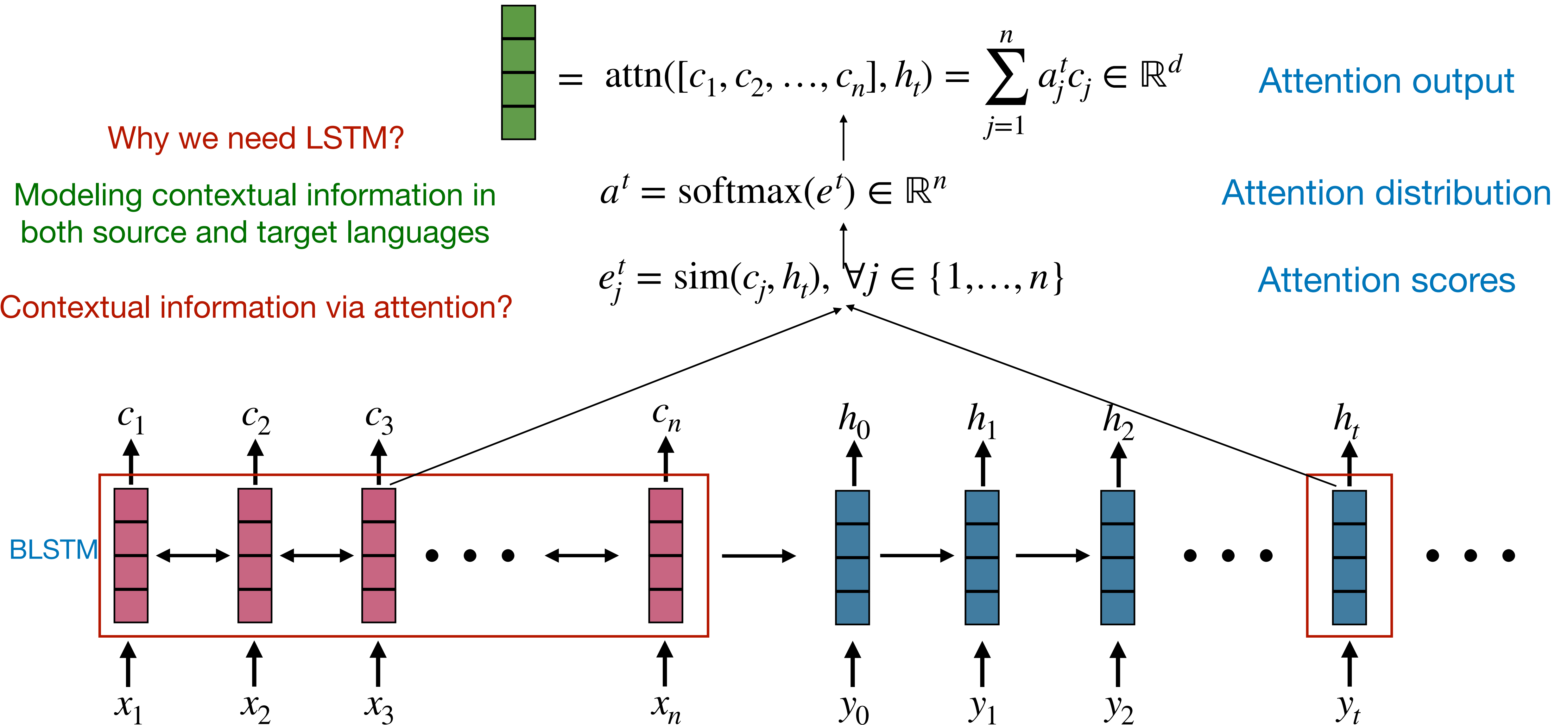illia.polosukhin@gmail.com

# Transformers



- Consists of an encoder and a decoder

- Originally proposed for neural machine translation and later adapted for almost all the NLP tasks
  - For example, BERT only uses the **encoder** of the Transformer architecture (next lecture)

- Both encoder and decoder consist of $N$ layers
  - Each encoder layer has two sub-layers
  - Each decoder layer has three sublayers
  - Key innovation: **multi-head self-attention**

Vaswani et al., 2017: Attention Is All You Need

# Recap: Attention Mechanism

$$= \text{attn}([c_1, c_2, \ldots, c_n], h_t) = \sum_{j=1}^{n} a_j^t c_j \in \mathbb{R}^d$$

Attention output

Why we need LSTM?

Modeling contextual information in both source and target languages

$$a^t = \text{softmax}(e^t) \in \mathbb{R}^n$$

Attention distribution

Contextual information via attention?

$$e_j^t = \text{sim}(c_j, h_t), \ \forall j \in \{1, \ldots, n\}$$

Attention scores

BLSTM

$c_1$ $c_2$ $c_3$ $c_n$ $h_0$ $h_1$ $h_2$ $h_t$

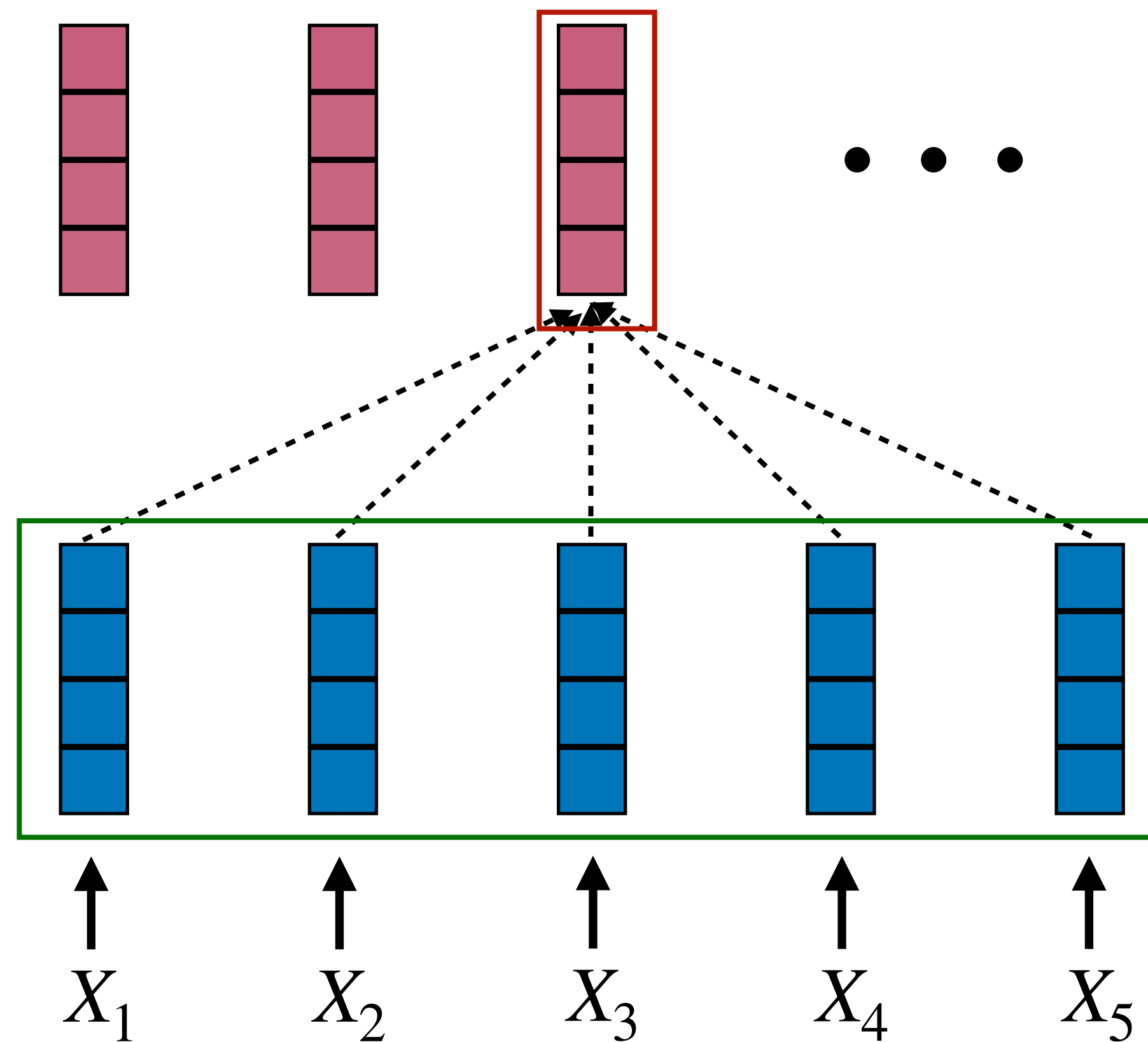$x_1$ $x_2$ $x_3$ $x_n$ $y_0$ $y_1$ $y_2$ $y_t$

# Self-Attention

- **Self-attention: attention within on single sequence**
  - Contexts and queries are drawn from the same source
- **Contextual information via self-attention**



$$X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5$$

# Self-Attention

- **Self-attention: attention within on single sequence**
  - Contexts and queries are drawn from the same source
- **Contextual information via self-attention**



$$X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5$$

# Self-Attention

- **Self-attention: attention within on single sequence**
  - Contexts and queries are drawn from the same source
- **Contextual information via self-attention**

- Capturing long-distance dependencies
- No gradient vanishing



$X_1 \quad X_2 \quad X_3 \quad X_4 \quad X_5$

# Self-attention in equations

- A sequence of input vectors $x_1, \ldots, x_n \in \mathbb{R}^d$

- First, construct a set of queries, keys and values:
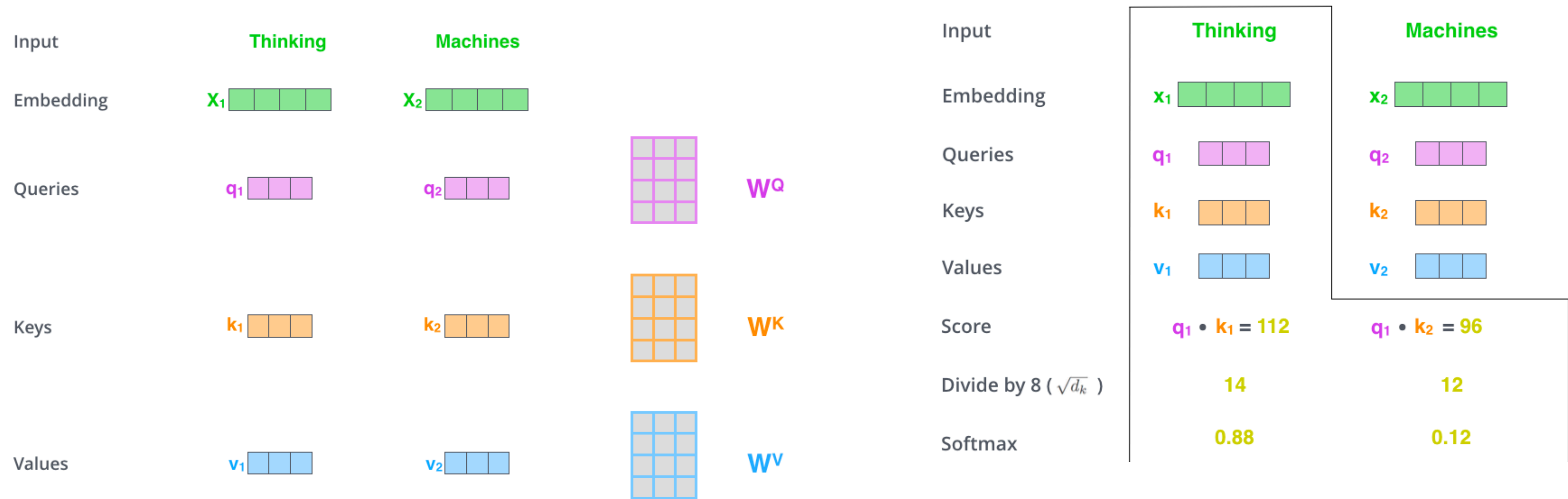
$$q_i = W_Q x_i, \ k_i = W_K x_i, \ v_i = W_V x_i$$

- Second, for each $q_i$, compute attention scores and attention distributions:

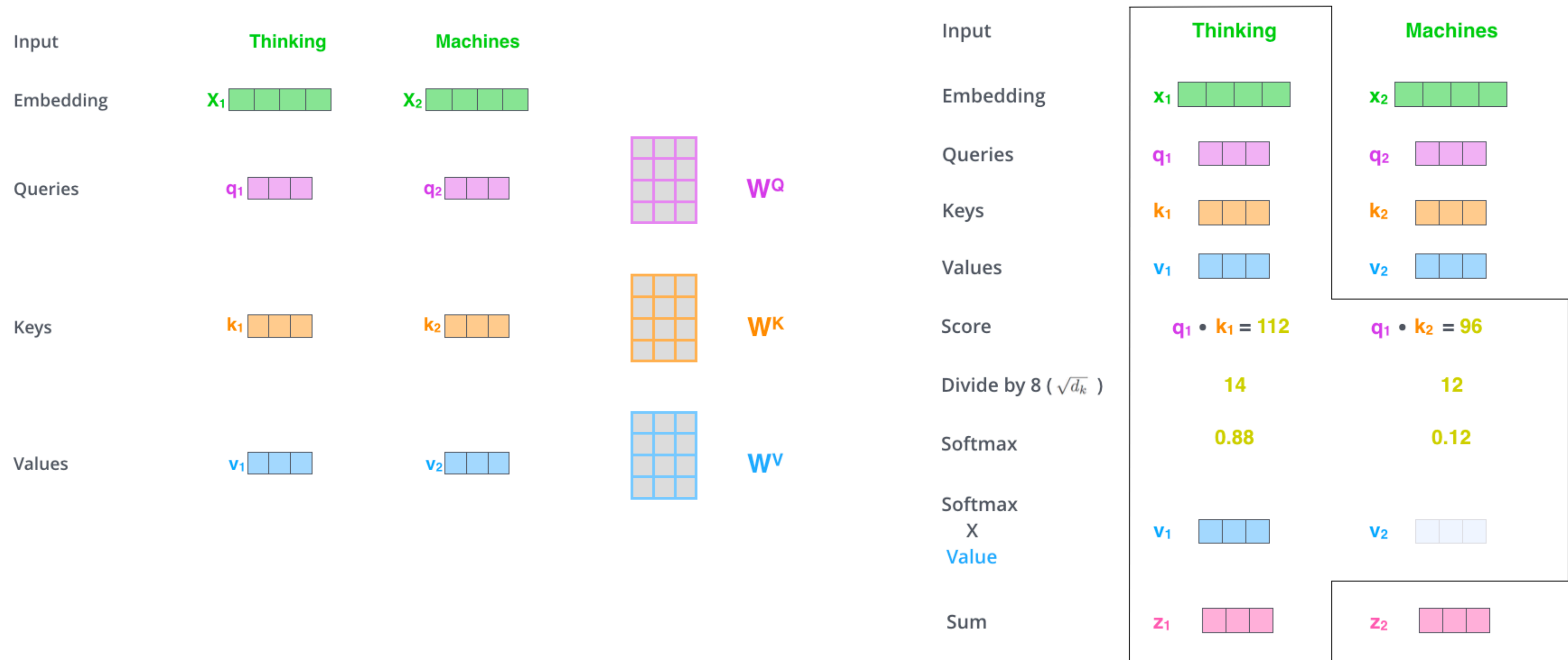$$a_{i,j} = \text{softmax}(\frac{q_i^T k_j}{\sqrt{d}})$$

aka. "scaled dot product"

- Finally, compute the weighted sum:

$$y_i = \sum_{j=1}^{n} a_{i,j} v_j$$

# Why *Scaled* Dot Product?

- **Softmax is sensitive to scale**

If $[x_1, x_2] = [0.1, 0.5], \alpha = 10$

$$\text{softmax}([x_1, x_2]) = \left[ \frac{e^{x_1}}{e^{x_1} + e^{x_2}}, \frac{e^{x_2}}{e^{x_1} + e^{x_2}} \right]$$

$[0.4013, 0.5987]$

$$\text{softmax}([\alpha x_1, \alpha x_2]) = \left[ \frac{e^{\alpha x_1}}{e^{\alpha x_1} + e^{\alpha x_2}}, \frac{e^{\alpha x_2}}{e^{\alpha x_1} + e^{\alpha x_2}} \right]$$

$[0.0180, 0.9820]$

# Self-attention in equations

- A sequence of input vectors $x_1, \ldots, x_n \in \mathbb{R}^d$
- First, construct a set of <span style="color:green">queries</span>, <span style="color:red">keys</span> and <span style="color:blue">values</span>:

$$q_i = W_Q x_i, \; k_i = W_K x_i, \; v_i = W_V x_i$$

- Second, for each $q_i$, compute attention scores and attention distributions:

$$a_{i,j} = \text{softmax}\left(\frac{q_i^T k_j}{\sqrt{d}}\right) \qquad \text{aka. "scaled dot product"}$$

- Finally, compute the weighted sum:

$$y_i = \sum_{j=1}^{n} a_{i,j} v_j$$

# Self-attention: Illustration

# Self-attention: Illustration

# Self-attention: matrix notations

# Self-attention: matrix notations

$$\text{attn}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d}})V$$

# Self-attention: matrix notations

# Attention is *General*

- **Given a set of key and value vectors, and a query vector, attention is a technique to compute a weighted sum of the value vectors, dependent on the query and keys**

  – We sometimes say that the query attends to the values via keys

  – In the NMT vase, each decoder hidden state (query) attends to all the encoder hidden states (keys and values)

- **Intuition**

  – The weighted sum is a *selective summary* of the information contained in the values, where the query and keys determines which values to focus on

  – Attention is a way to obtain a *fixed-size representation* of an arbitrary set of representations (the values), dependent on some other representation (the query)

$$\text{attn}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d}})V$$

# Multi-head Attention

- **Problem with self-attention?**

$$y_i = \sum_{j=1}^{n} a_{i,j} v_j \qquad \text{one set of attention weights } a_i$$

- **It is better to use multiple attention weights instead of one!**
  - Each attention can focus on different positions
- **How to do this? Splits queries, keys, values to multiple heads!**

# Multi-head Attention: Head Split
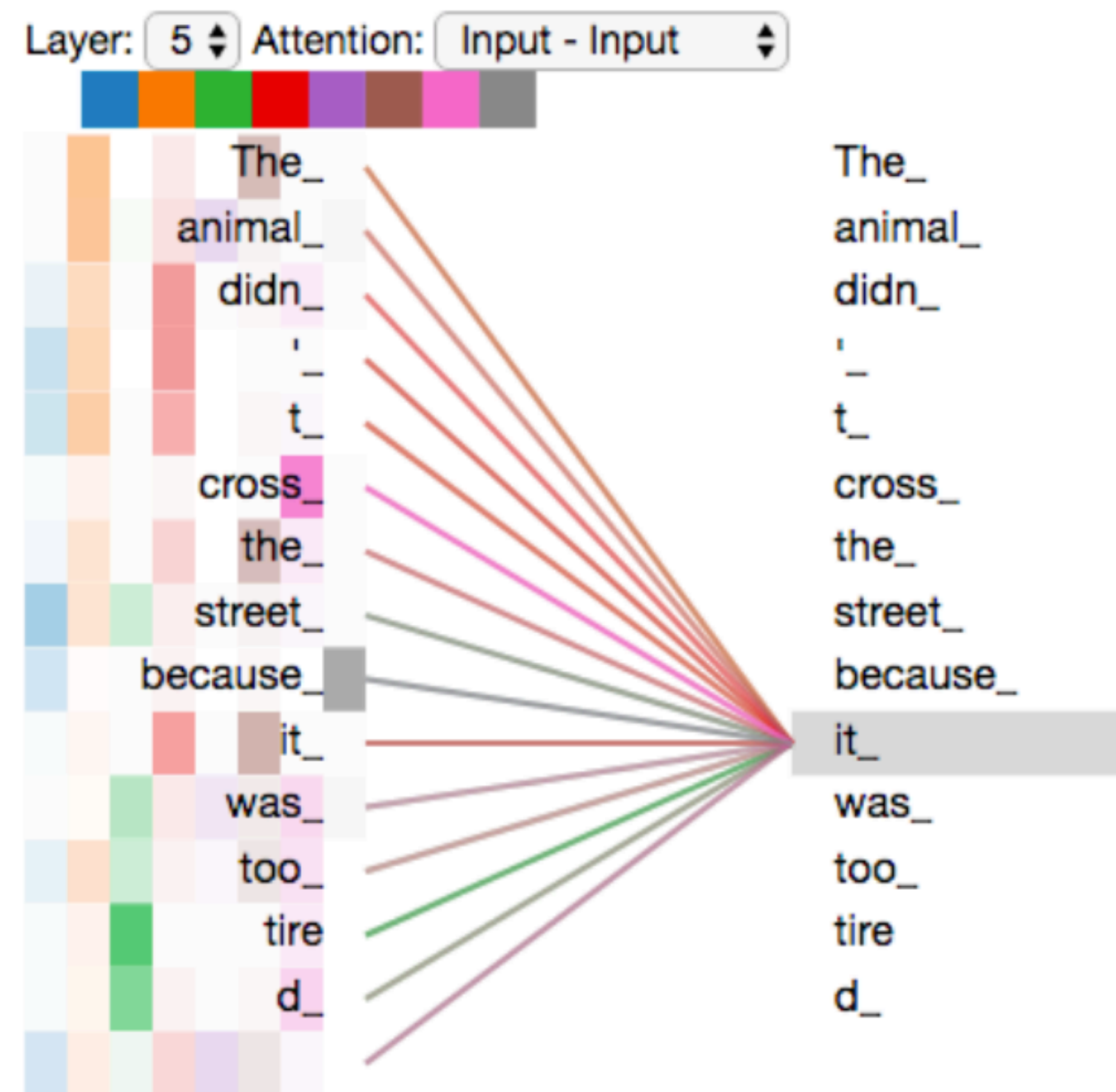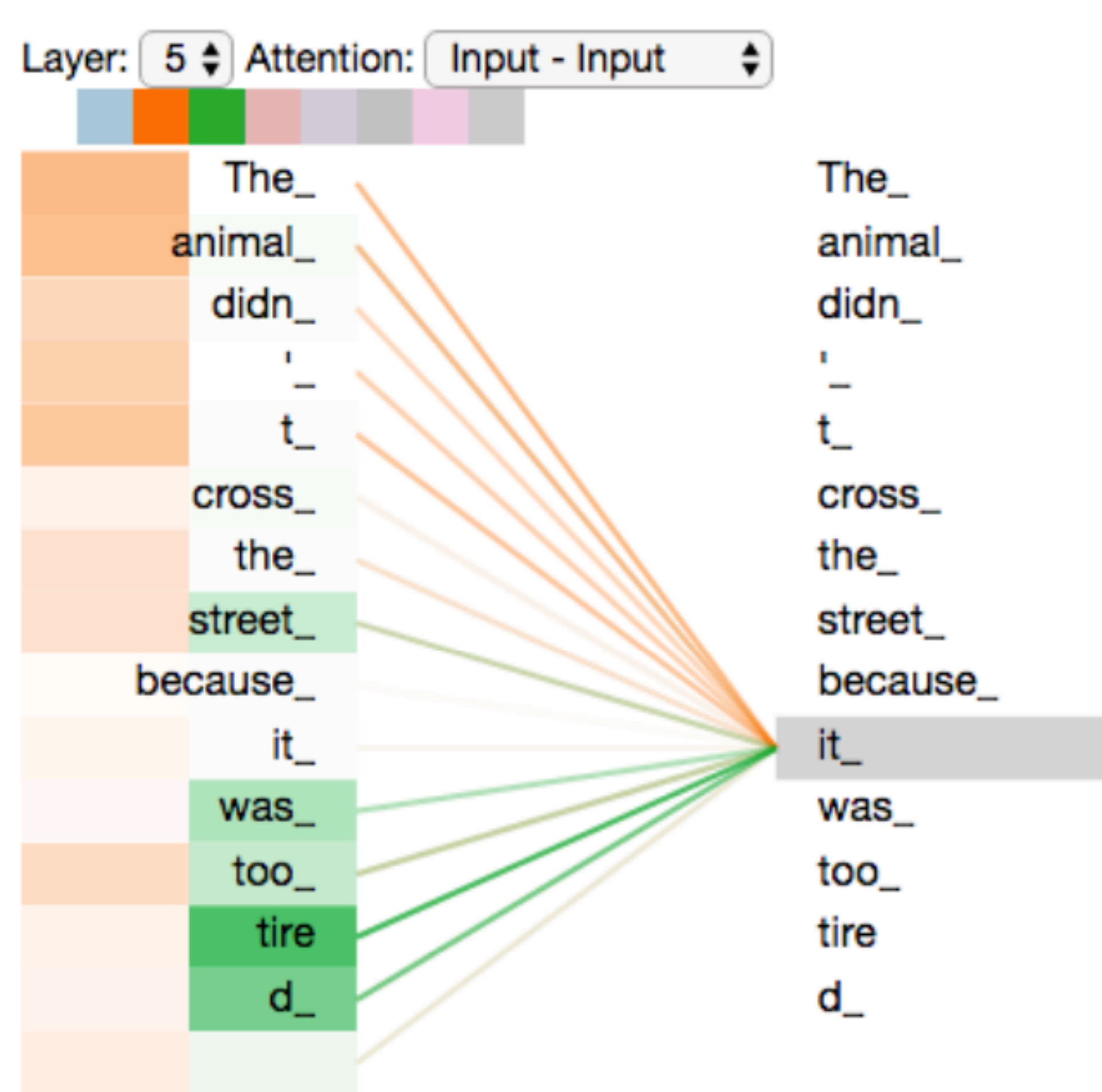


$$h_1 = \text{attn}(Q_1, K_1, V_1) = \text{softmax}(\frac{Q_1 K_1^T}{\sqrt{d/2}})V_1$$

$$h_2 = \text{attn}(Q_2, K_2, V_2) = \text{softmax}(\frac{Q_2 K_2^T}{\sqrt{d/2}})V_2$$

$$Y = \text{concat}(h_1, h_2)W_O$$

# What does multi-head attention learn?

# Transformer Encoder

- **Replacing RNNs with multi-head self-attention**



$$\text{MultiHead}(X) = \text{concat}(h_1, \ldots, h_k)W_O$$
$$h_i = \text{attn}(Q_i, K_i, V_i)$$
$$Q_i = (XW_Q)^i, K_i = (XW_K)^i, V_i = (XW_V)^i$$

$$\text{attn}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d}})V$$

Self-attention does not know the order of the inputs!

# Missing Piece: Positional Information

RNN

Self-attention



$X_1$ $X_2$ $X_3$

$X_1$ $X_2$ $X_3$ $X_4$ $X_5$

# Missing Piece: Positional Information

- **Unlike RNNs, self-attention does <span style="color:red">not</span> build in order information**
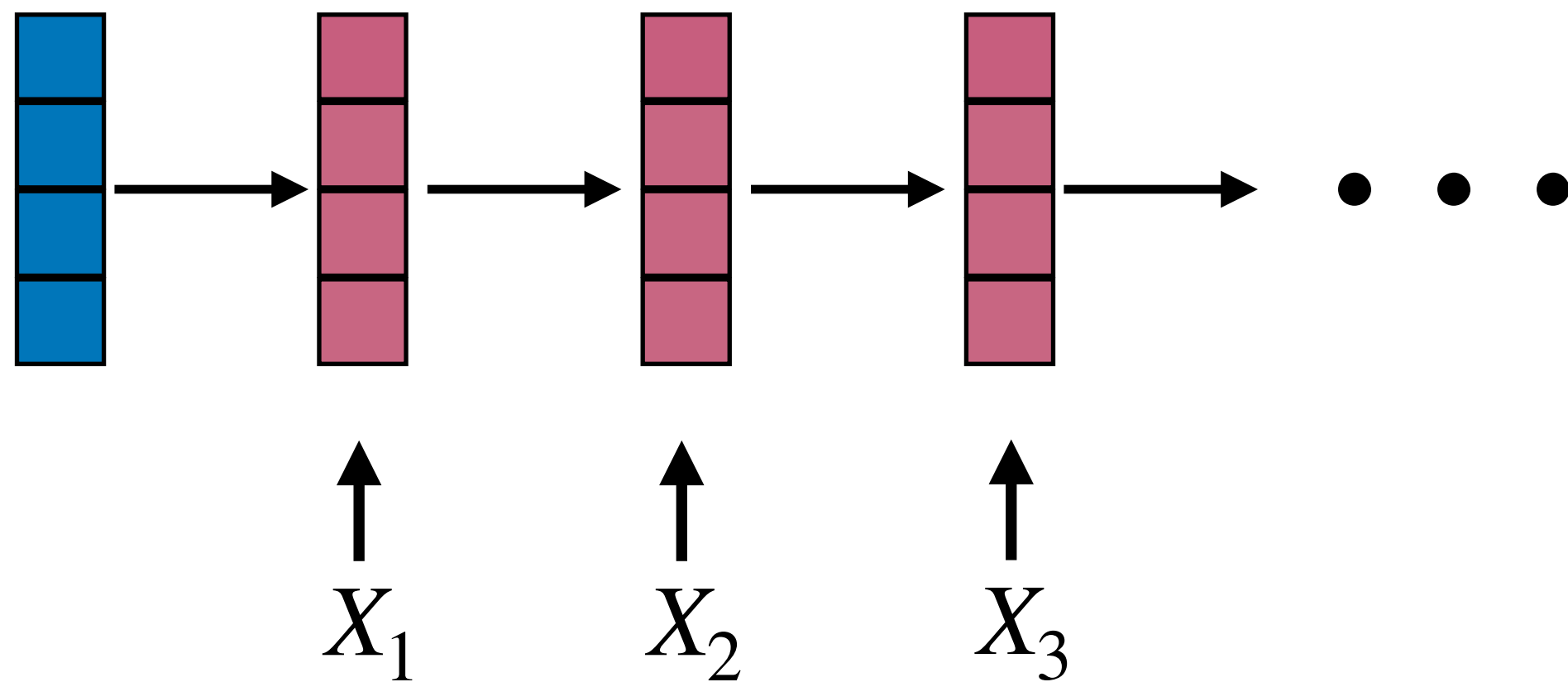  - Encode the order of the sentence into the input $x_1, \ldots, x_n$
- **Solution: add <span style="color:green">positional encoding</span> to the input embeddings**

$$x_i \leftarrow x_i + p_i$$

- **Use sine and cosine functions of different frequencies (not learnable)**

$$p_i = \begin{pmatrix} \sin(i/10000^{2*1/d}) \\ \cos(i/10000^{2*1/d}) \\ \vdots \\ \sin(i/10000^{2*\frac{d}{2}/d}) \\ \cos(i/10000^{2*\frac{d}{2}/d}) \end{pmatrix}$$



Dimension

Index in the sequence

# Adding Nonlinearities

- There is no elementwise nonlinearities in self-attention; stacking more self-attention layers just re-averages value vectors

- Simple fix: add a feed-forward network to post-process each output vector

$$\mathrm{FFN}(\mathbf{x}_i) = W_2 \mathrm{ReLU}(W_1 \mathbf{x}_i + \mathbf{b}_1) + \mathbf{b}_2$$
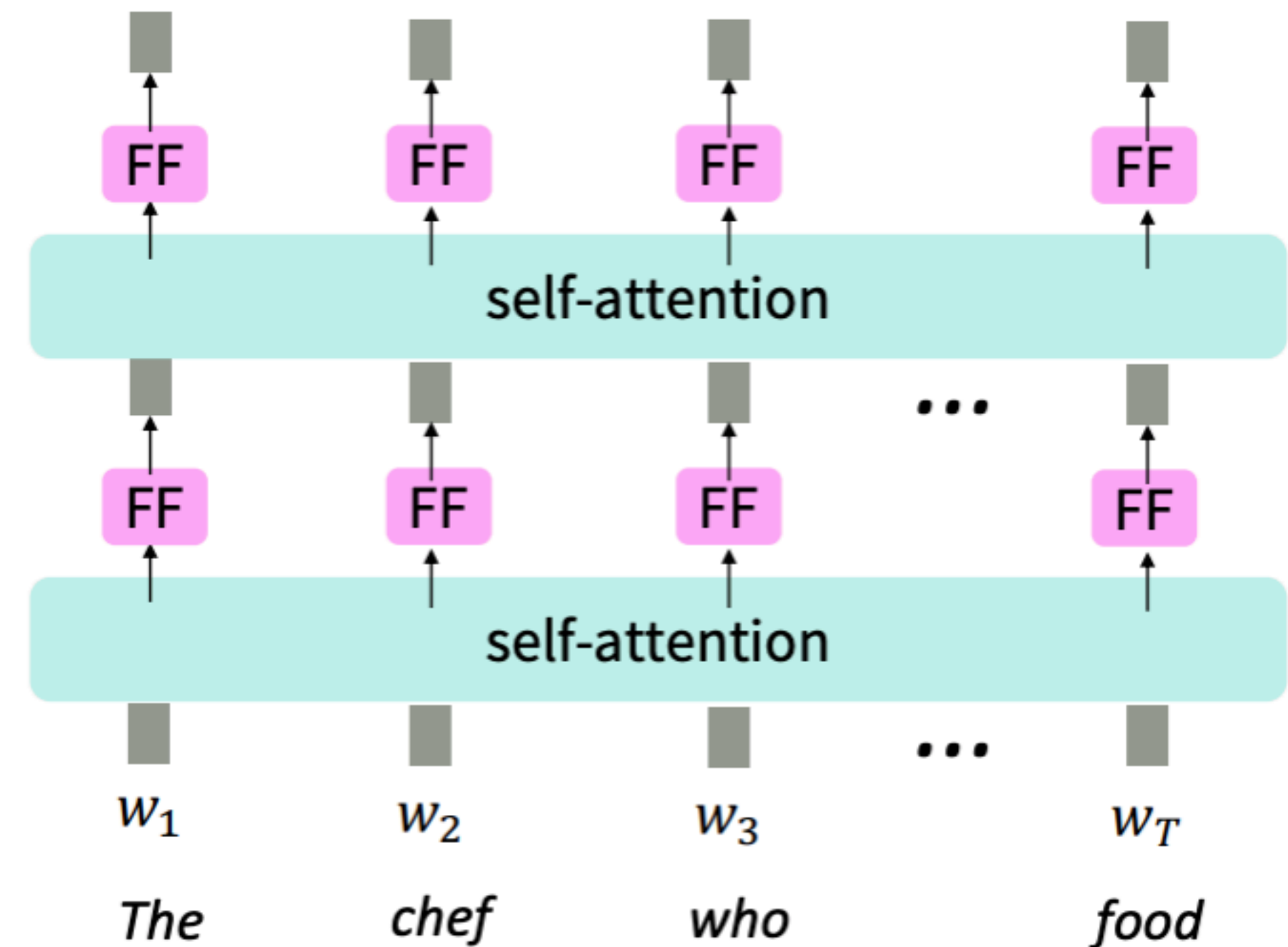
<span style="color:red">A large number of parameters</span>

$$W_1 \in \mathbb{R}^{d_{ff} \times d}, \mathbf{b}_1 \in \mathbb{R}^{d_{ff}}$$

$$W_2 \in \mathbb{R}^{d \times d_{ff}}, \mathbf{b}_2 \in \mathbb{R}^{d}$$

In practice, they use $d_{ff} = 4d$

36

# Feed-Forward Layers

- **Feed-forward layers constitute two-thirds of parameters**
- **Operates as memories of textual patterns** (Gova et al., 2021)

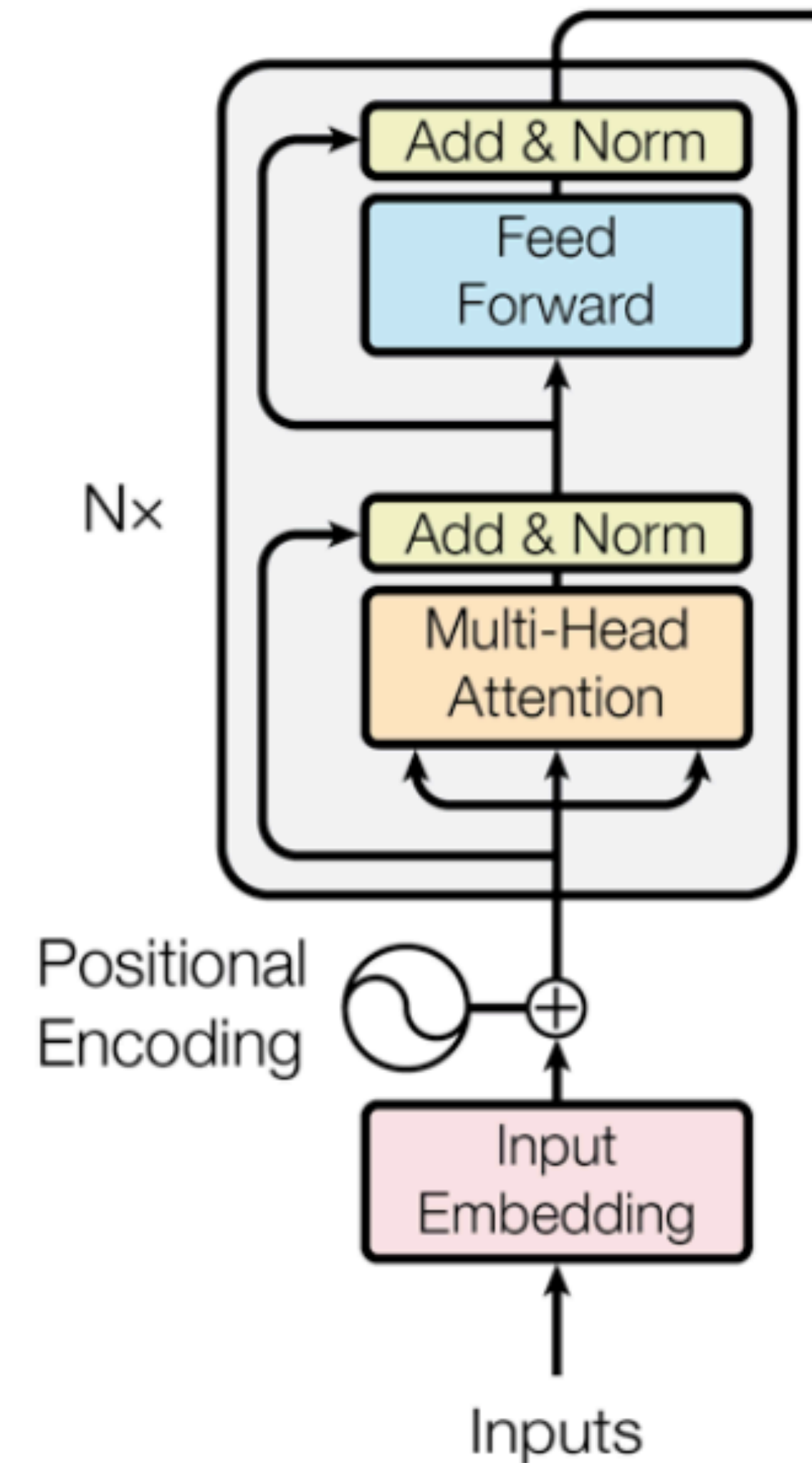| Key | Pattern | Example trigger prefixes |
|---|---|---|
| $\mathbf{k}_{449}^{1}$ | Ends with *"substitutes"* (shallow) | *At the meeting, Elton said that "for artistic reasons there could be no substitutes*<br>*In German service, they were used as substitutes*<br>*Two weeks later, he came off the substitutes* |
| $\mathbf{k}_{2546}^{6}$ | Military, ends with *"base"/"bases"* (shallow + semantic) | *On 1 April the SRSG authorised the SADF to leave their bases*<br>*Aircraft from all four carriers attacked the Australian base*<br>*Bombers flying missions to Rabaul and other Japanese bases* |
| $\mathbf{k}_{2997}^{10}$ | a "part of" relation (semantic) | *In June 2012 she was named as one of the team that competed*<br>*He was also a part of the Indian delegation*<br>*Toy Story is also among the top ten in the BFI list of the 50 films you should* |
| $\mathbf{k}_{2989}^{13}$ | Ends with a time range (semantic) | *Worldwide, most tornadoes occur in the late afternoon, between 3 pm and 7*<br>*Weekend tolls are in effect from 7:00 pm Friday until*<br>*The building is open to the public seven days a week, from 11:00 am to* |
| $\mathbf{k}_{1935}^{16}$ | TV shows (semantic) | *Time shifting viewing added 57 percent to the episode's*<br>*The first season set that the episode was included in was as part of the*<br>*From the original NBC daytime version , archived* |

# Transformer Encoder

- **Each encoder layer has two sub-layers:**
  - A multi-head self-attention layer
  - A feedforward layer
- **Add & Norm:**
  - Add: Residual connection (He et al., 2016)

$$Y \leftarrow Y + X$$

  - Norm: Layer normalization (Ba et al., 2016)

$$Y = \frac{X - \mathrm{E}[X]}{\sqrt{\mathrm{Var}[X] + \epsilon}} * \gamma + \beta$$

In (Vaswani et al., 2017), N = 6

# Question

Which of the following statements is correct?

    (a)  Transformers run faster than LSTMs

    (b)  Transformers are easier to parallelize compared to LSTMs

    (c)  Transformers have less parameters compared to LSTMs

    (d)  Transformers are better at capturing positional information than LSTMs

# Transformer: Pros and Cons

- **Easier to capture dependencies:** we draw attention between every pair of words!

- **Easier to parallelize:**
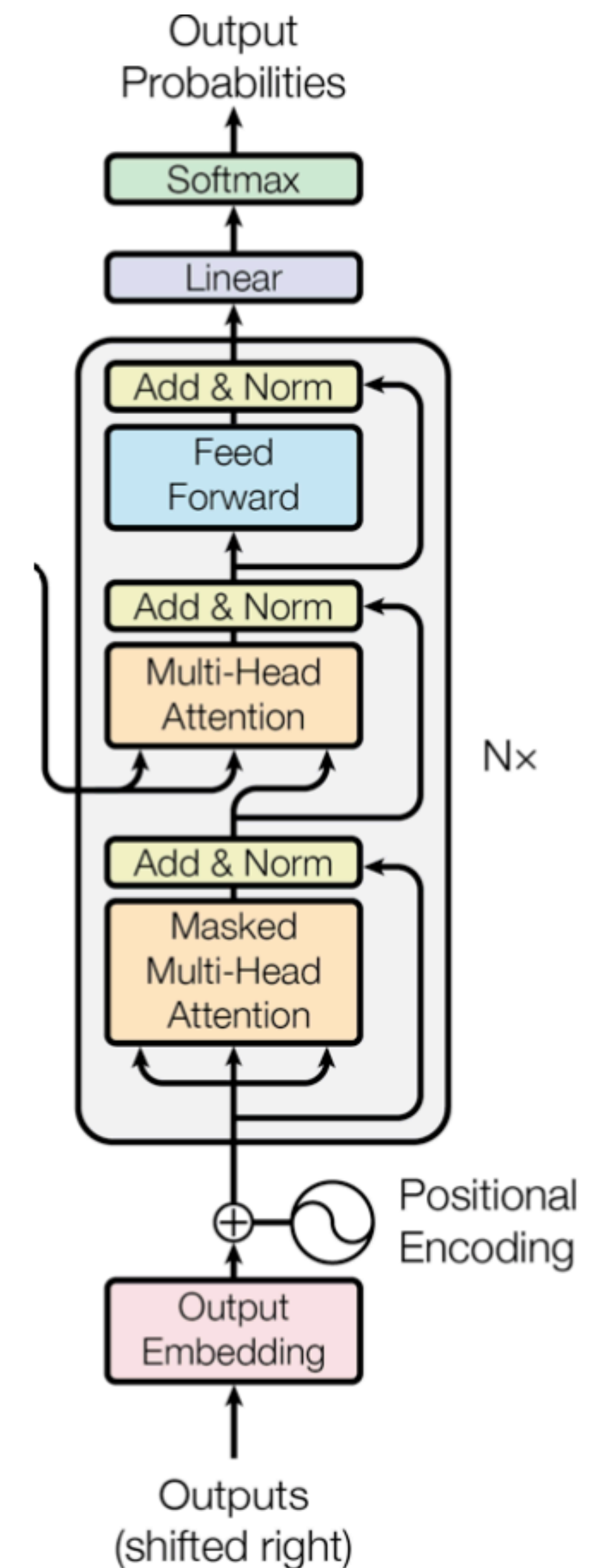
$$\text{MultiHead}(X) = \text{concat}(h_1, \ldots, h_k)W_O$$
$$h_i = \text{attn}(Q_i, K_i, V_i)$$
$$Q_i = (XW_Q)^i, K_i = (XW_K)^i, V_i = (XW_V)^i$$

$$\text{attn}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d}})V$$

- **Quadratic computation in self-attention:**
  - Can be very expensive when the sequence is very long

- **Harder to model positional information**

# Transformer Decoder
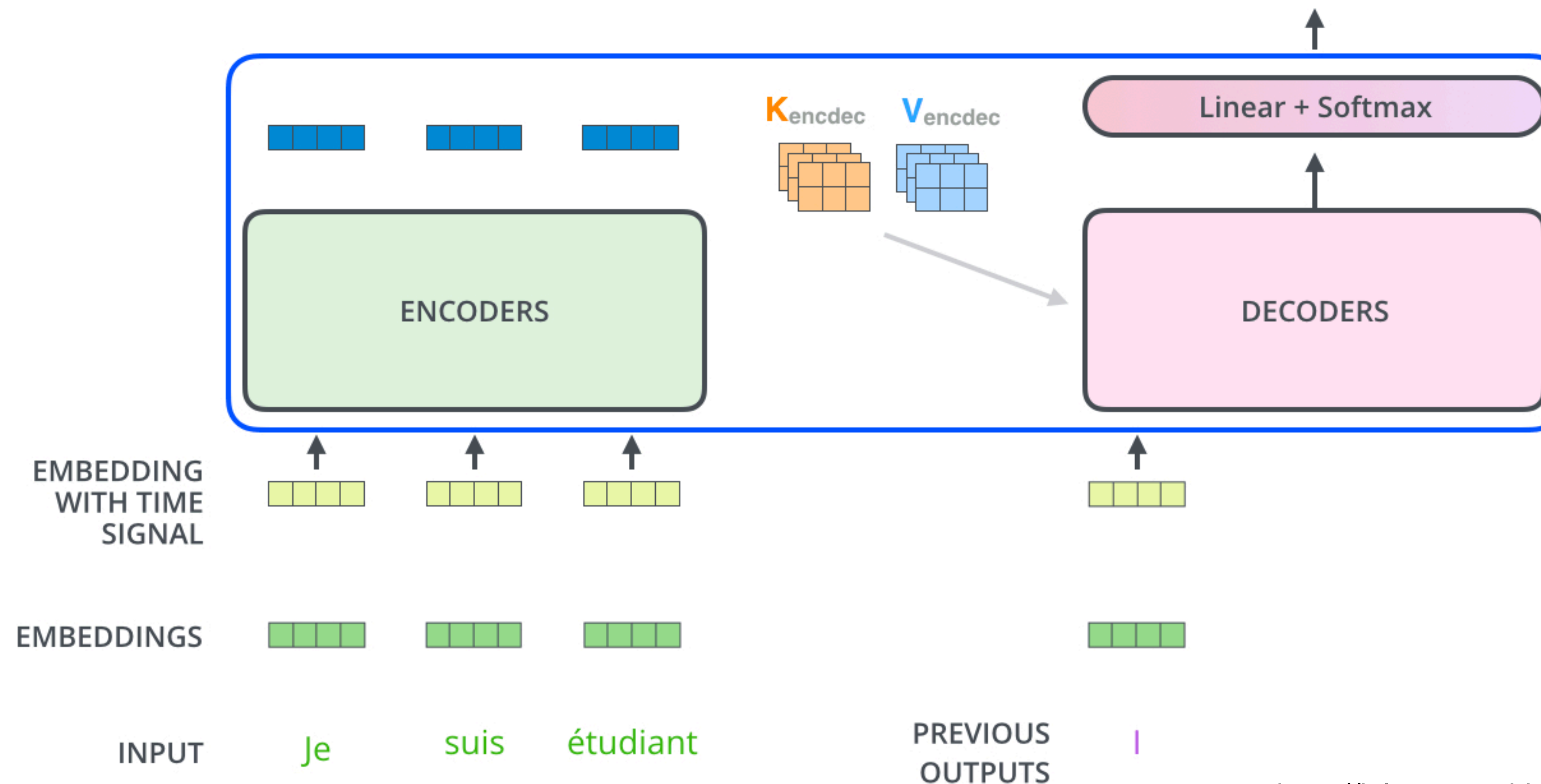
- **Each decoder layer has three sub-layers:**
  - A masked multi-head attention layer
  - A multi-head cross attention layer
  - A feedforward layer
- **Masked multi-head attention**
  - self-attention on the decoder states
- **Multi-head cross attention**
  - Decoder attends to encoder states
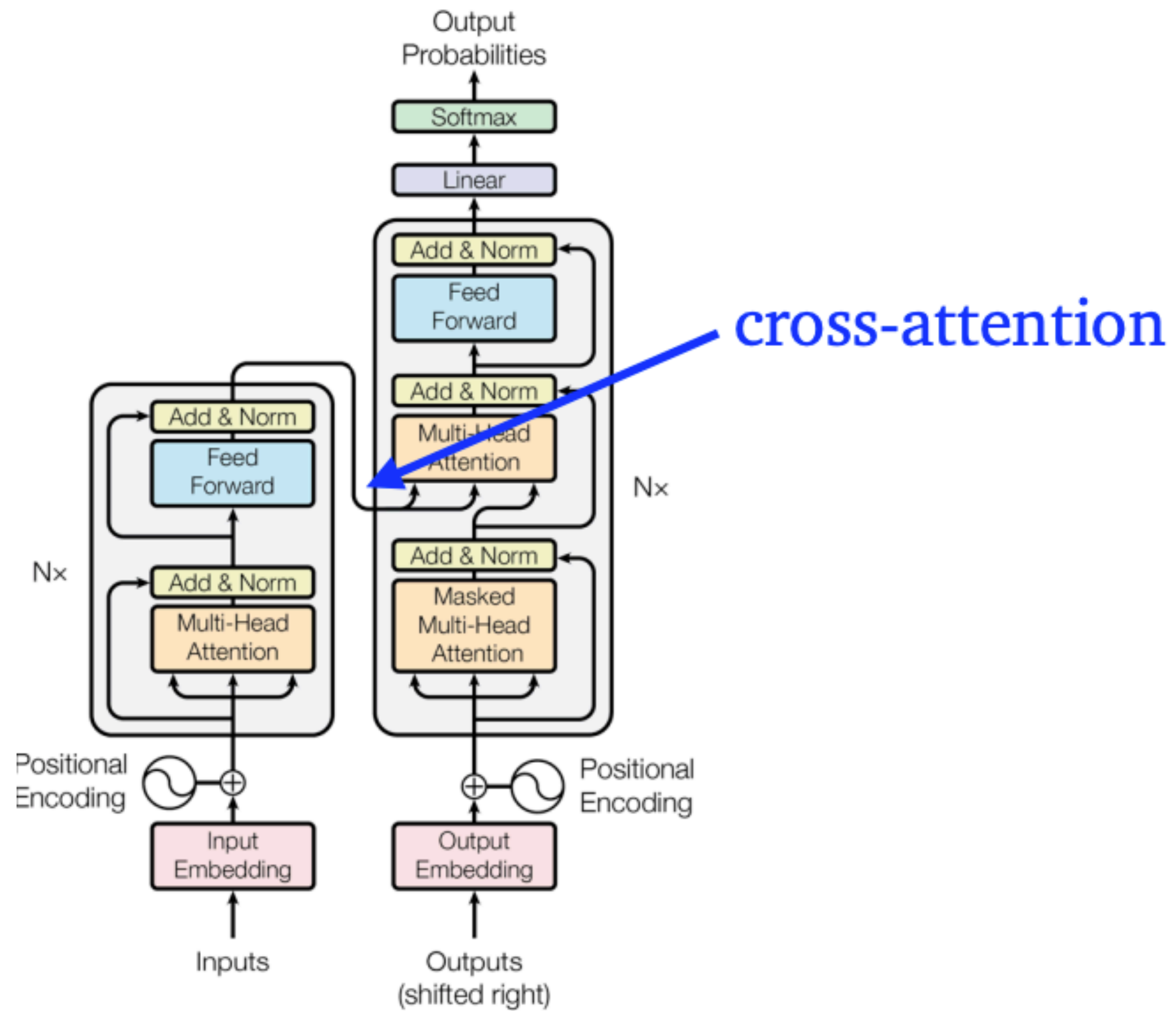  - Encoder: keys/values
  - Decoder: queries



In (Vaswani et al., 2017), N = 6

41

# Multi-head Cross Attention

- **Decoder attends to encoder states**
  - Encoder: keys/values
  - Decoder: queries

# Multi-head Cross Attention



cross-attention

$$q_i = W_Q y_i, \quad k_i = W_K x_i, \ v_i = W_V x_i$$

$$a_{i,j} = \text{softmax}(\frac{q_i^T k_j}{\sqrt{d}})$$

$$y_i' = \sum_{j=1}^{n} a_{i,j} v_j$$

# Masked Multi-Head Attention

- **Key point:** we cannot see the future words in decoder
- **Solution:** for every $q_i$, only attend to $\{(k_j, v_j)\}, j \leq i$



$y_0 \qquad y_1 \qquad y_2 \qquad y_3 \qquad y_4$

# Masked Multi-Head Attention

- **Key point:** we cannot see the future words in decoder
- **Solution:** for every $q_i$, only attend to $\{(k_j, v_j)\}, j \leq i$

# Masked Multi-Head Attention

- **Key point:** we cannot see the future words in decoder
- **Solution:** for every $q_i$, only attend to $\{(k_j, v_j)\}, j \leq i$



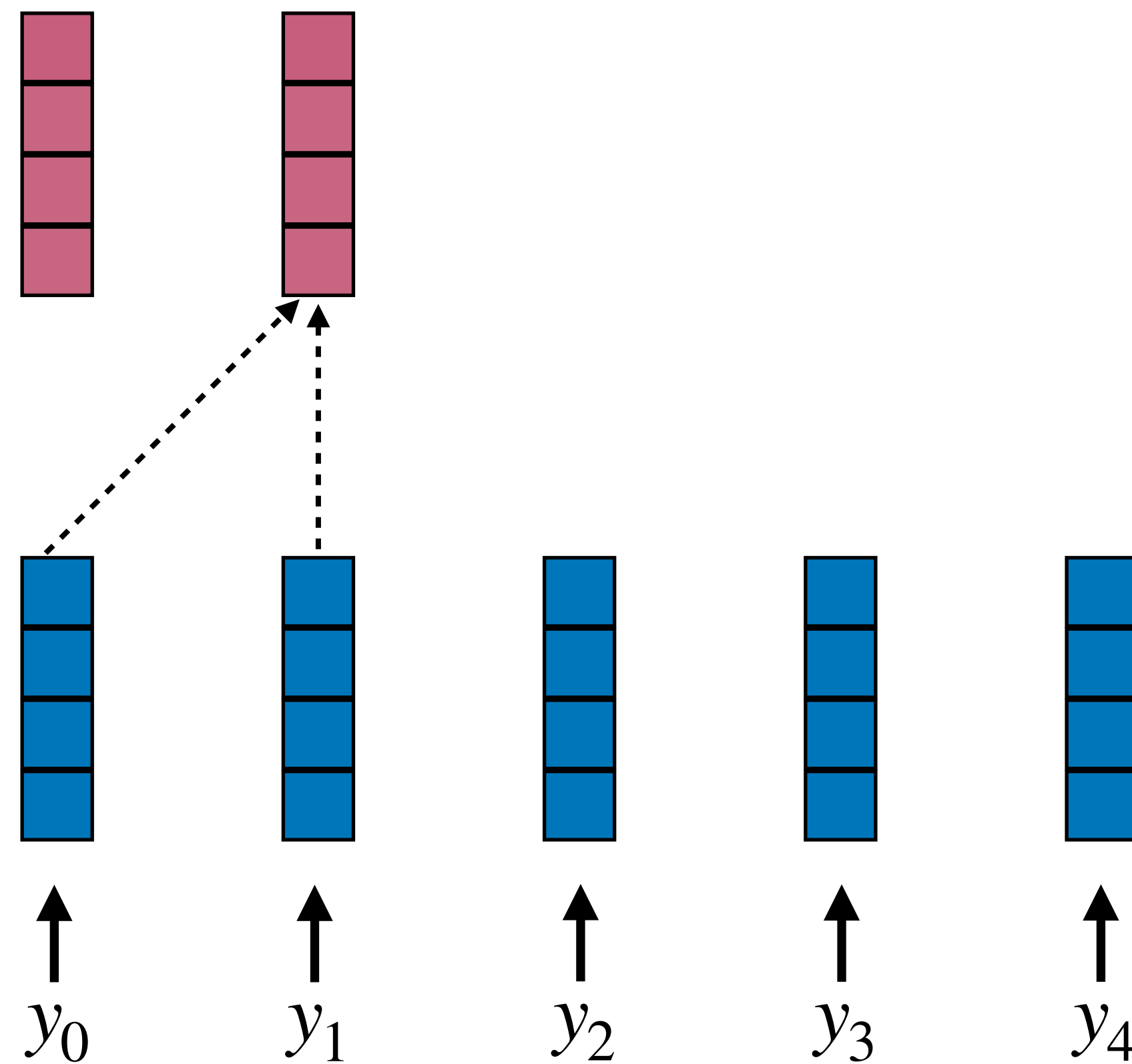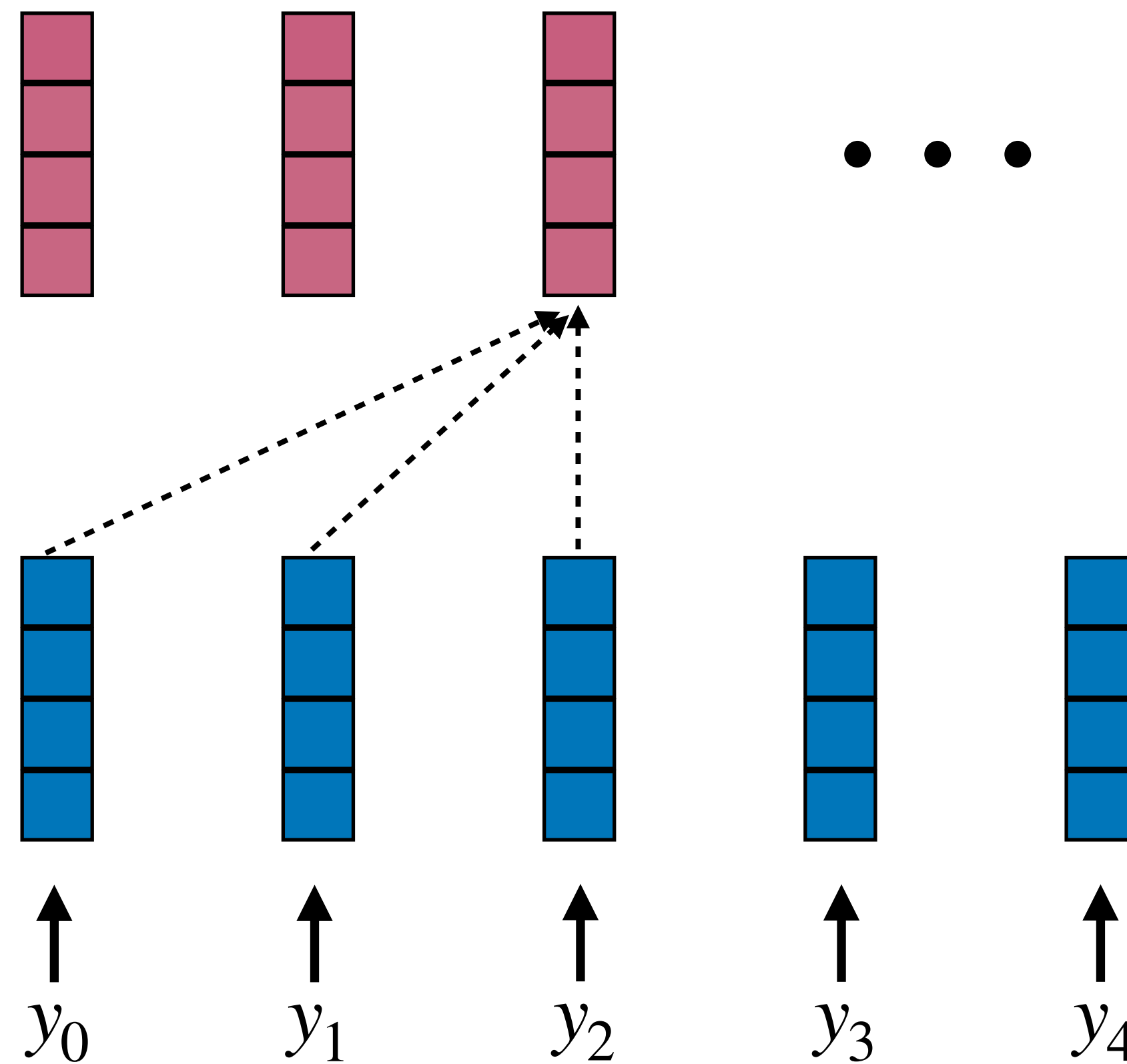cannot be parallel!

$y_0$     $y_1$     $y_2$     $y_3$     $y_4$

# Masked Multi-Head Attention

$$q_i = W_Q x_i, \ k_i = W_K x_i, \ v_i = W_V x_i$$

$$a_{i,j} = \text{softmax}(\frac{q_i^T k_j}{\sqrt{d}})$$



raw attention weights     mask

**Efficient implementation**: compute attention as we normally do, mask out attention to future words by setting attention scores to $-\infty$

```
dot = torch.bmm(queries, keys.transpose(1, 2))

indices = torch.triu_indices(t, t, offset=1)

dot[:, indices[0], indices[1]] = float('-inf')

dot = F.softmax(dot, dim=2)
```

http://peterbloem.nl/blog/transformers

# Putting the pieces together



[predictions!]

Transformer Encoder

Transformer Encoder

[decoder attends to encoder states]

Transformer Decoder

Transformer Decoder

Word Embeddings + Position Representations

Word Embeddings + Position Representations

[input sequence]

[output sequence]

# Transformer: Machine Translation

| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [18] | 23.75 | | | |
| Deep-Att + PosUnk [39] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [38] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [9] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [32] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [39] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [38] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [9] | 26.36 | **41.29** | $7.7 \cdot 10^{19}$ | $1.2 \cdot 10^{21}$ |
| Transformer (base model) | 27.3 | 38.1 | $\mathbf{3.3 \cdot 10^{18}}$ | |
| Transformer (big) | **28.4** | **41.8** | $2.3 \cdot 10^{19}$ | |

Vaswani et al., 2017: Attention Is All You Need

# Transformer: Document Generation

| Model | Test perplexity | ROUGE-L |
|---|---|---|
| $seq2seq\text{-}attention,\ L = 500$ | 5.04952 | 12.7 |
| $Transformer\text{-}ED,\ L = 500$ | 2.46645 | 34.2 |
| $Transformer\text{-}D,\ L = 4000$ | 2.22216 | 33.6 |
| $Transformer\text{-}DMCA,\ no\ MoE\text{-}layer,\ L = 11000$ | 2.05159 | 36.2 |
| $Transformer\text{-}DMCA,\ MoE\text{-}128,\ L = 11000$ | 1.92871 | 37.9 |
| $Transformer\text{-}DMCA,\ MoE\text{-}256,\ L = 7500$ | 1.90325 | 38.8 |

Significant gains compared to seq2seq-attention with LSTMs

Liu et al., 2018: Generating Wikipedia by Summarizing Long Sequences