# From Attention to Transformers
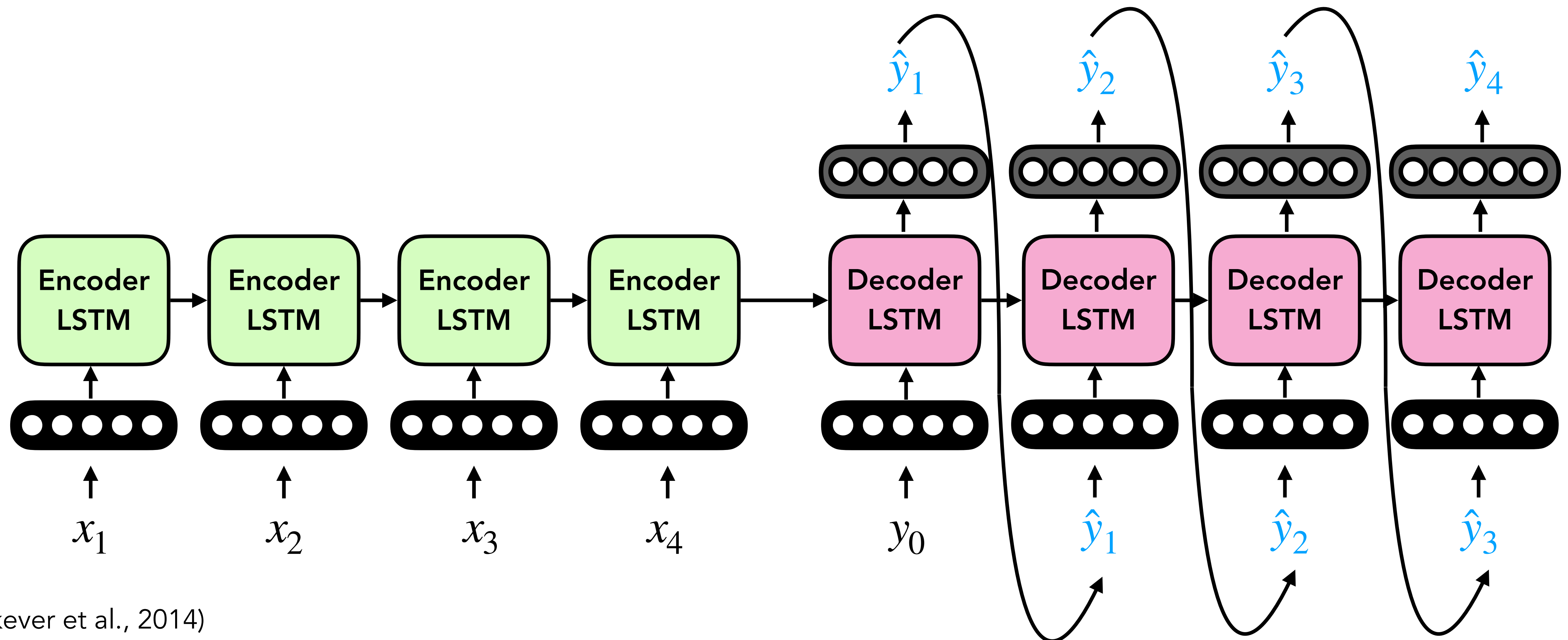
Antoine Bosselut

# Announcements

- **Assignment #1** released last Friday, March 10th

  - **Due Friday, March 24, 2023 @ 11:59 PM**

  - **Please post questions to Ed Discussion Board, so others can benefit from your queries!**

# Today's Outline

- **Quick Recap:** Attention

- **Supercharging Attention:** Transformers

- **Generating sequences:** Decoding from sequence to sequence models
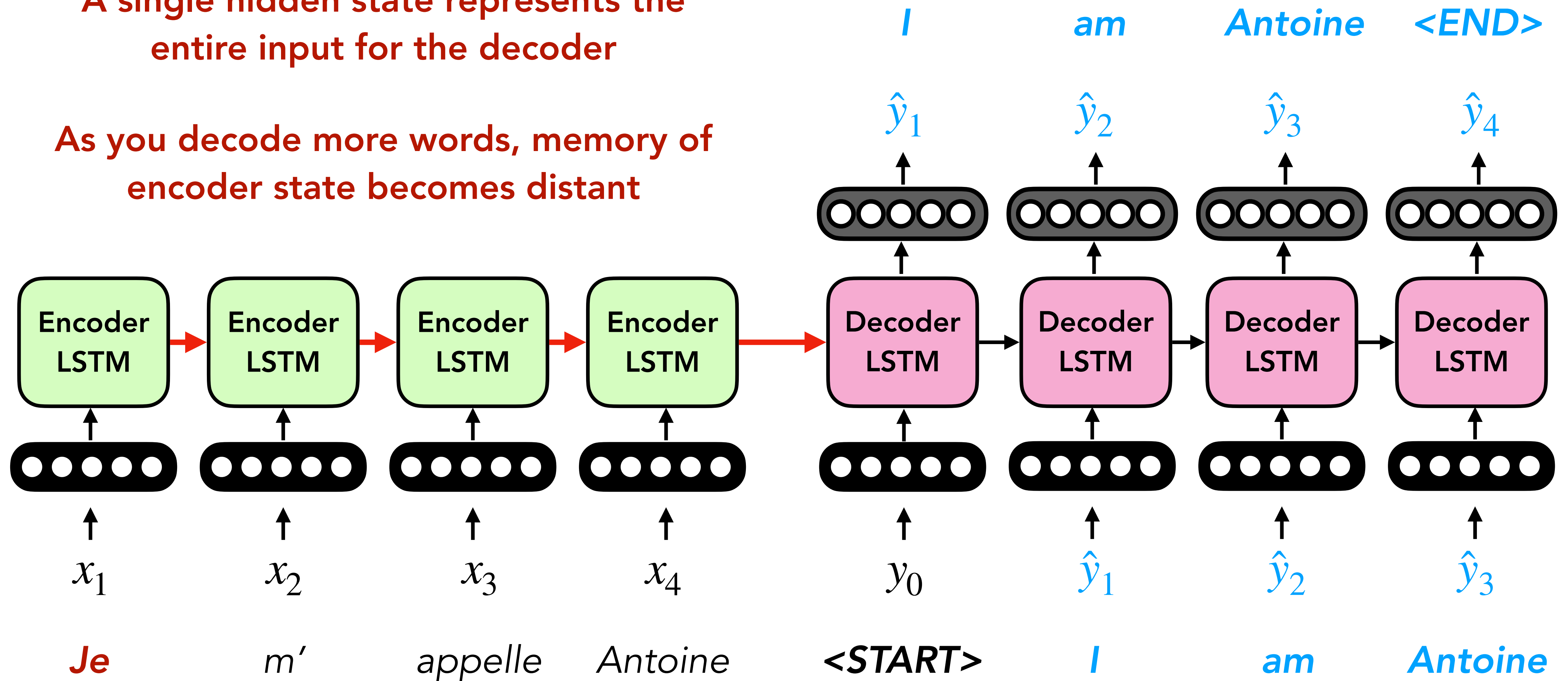
# Encoder-Decoder Models

- Encode a sequence fully with one model (**encoder**) and use its representation to seed a second model that decodes another sequence (**decoder**)
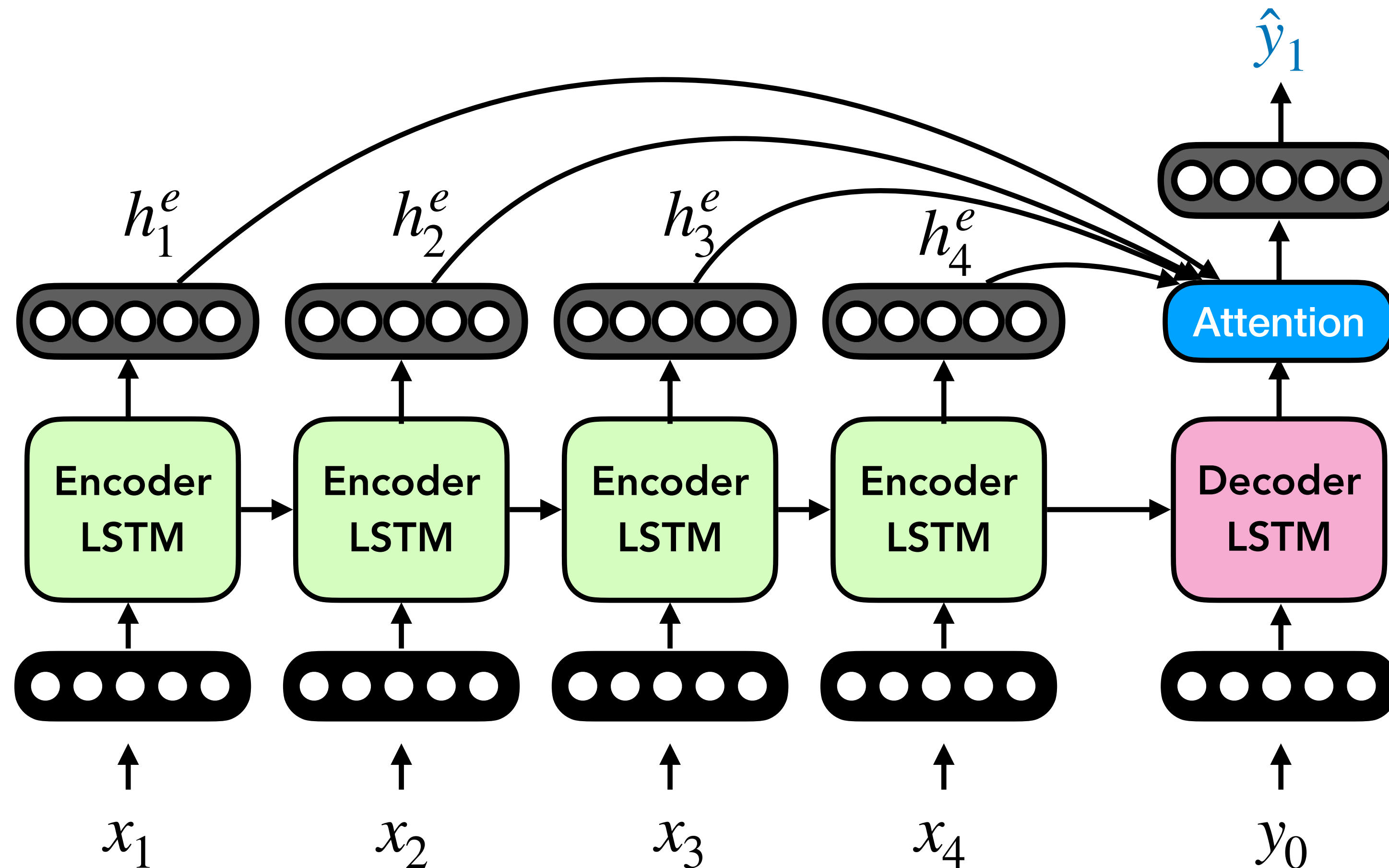


(Sutskever et al., 2014)

# Toy Example



A single hidden state represents the entire input for the decoder

As you decode more words, memory of encoder state becomes distant

# Attentive Encoder-Decoder Models



- **Recall:** Attention reduces this temporal bottleneck!

- **Intuition**: focus on different parts of the input at each time step

- **Idea:** Use the output of the Decoder LSTM to compute an **attention** (i.e., a mixture) over all the $h_t^e$ outputs of the encoder LSTM

(Bahdanau et al., 2015)

# Attention Function

- **Compute** pairwise similarity between each encoder hidden state and decoder hidden state ("idea of what to decode")

$h_t^e$ = encoder output hidden states

$h_t^d$ = decoder output hidden state

**Also known as a "keys"**

**Also known as a "query"**

$$a_1 = f\left( \begin{array}{c} h_1^e \end{array}, \begin{array}{c} h_1^d \end{array} \right) \quad a_2 = f\left( \begin{array}{c} h_2^e \end{array}, \begin{array}{c} h_1^d \end{array} \right) \quad a_3 = f\left( \begin{array}{c} h_3^e \end{array}, \begin{array}{c} h_1^d \end{array} \right) \quad a_4 = f\left( \begin{array}{c} h_4^e \end{array}, \begin{array}{c} h_1^d \end{array} \right)$$

- **We have a single query vector for multiple key vectors**

# Attention Function

| Attention Function | Formula |
|---|---|
| Multiplicative | $a = h^e \mathbf{W} h^d$ |
| Linear | $a = v^T \phi(\mathbf{W}[h^e; h^d])$ |
| Scaled Dot Product | $a = \dfrac{(\mathbf{W}h^e)^T (\mathbf{U}h^d)}{\sqrt{d}}$ |

# Attention Function

- **Compute** pairwise similarity between each encoder hidden state and decoder hidden state ("idea of what to decode")

$$a_1 = f\left( h_1^e, h_1^d \right) \qquad a_2 = f\left( h_2^e, h_1^d \right) \qquad a_3 = f\left( h_3^e, h_1^d \right)$$

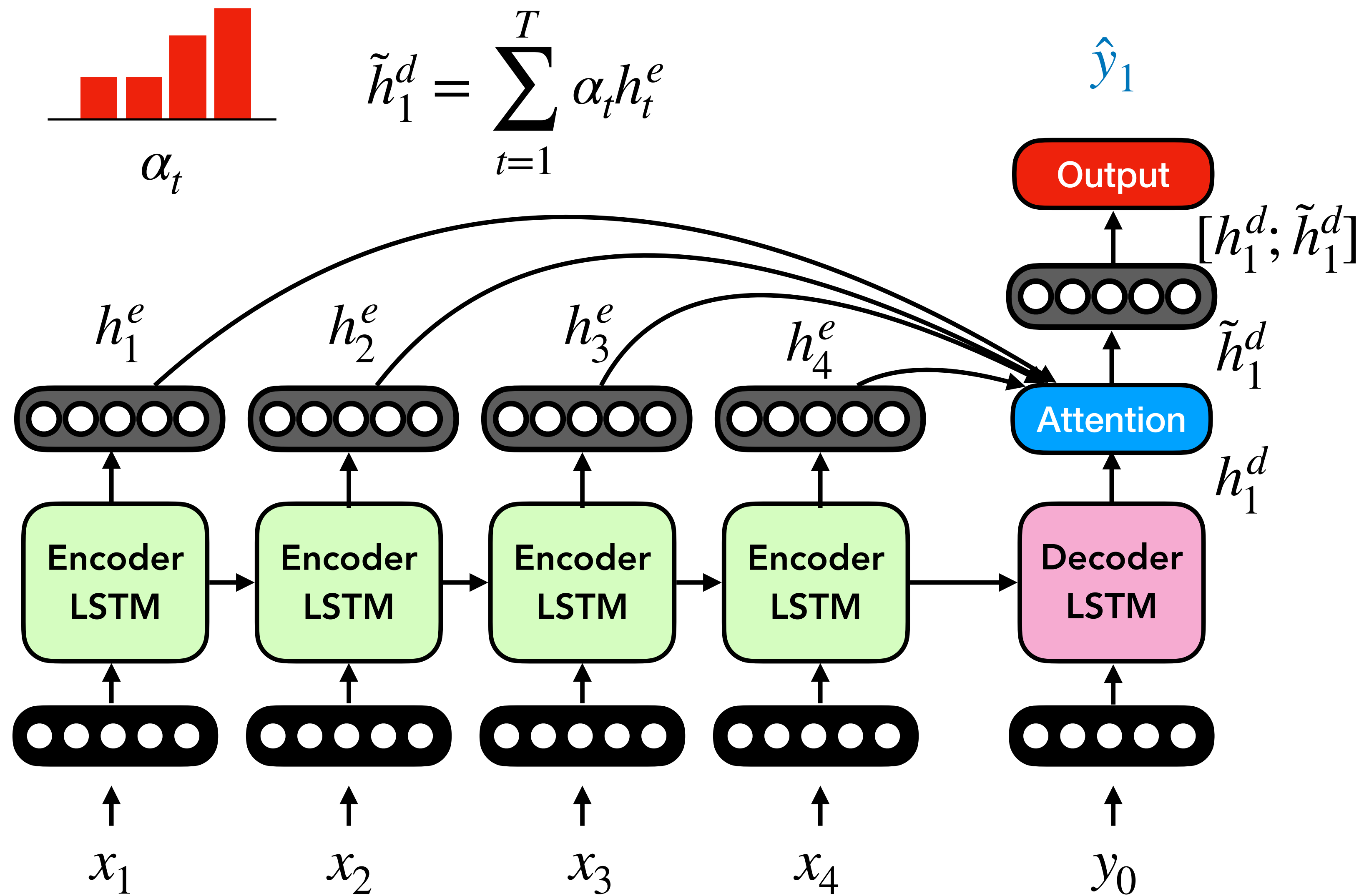"key"  "query"                "key"  "query"                "key"  "query"

- **Convert** pairwise similarity scores to probability **distribution** (using softmax!) over encoder hidden states and compute weighted average:

**Softmax!**  $\boxed{\alpha_t = \dfrac{e^{a_t}}{\sum_j e^{a_j}}}$ $\rightarrow$ $\alpha_t$ $\rightarrow$ $\tilde{h}_1^d = \sum_{t=1}^{T} \alpha_t h_t^e$ **Here $h_t^e$ is known as the "value"**

# Attentive Encoder-Decoder Models



$$\tilde{h}_1^d = \sum_{t=1}^{T} \alpha_t h_t^e$$

$\alpha_t$

$\hat{y}_1$

Output

$[h_1^d; \tilde{h}_1^d]$

$\tilde{h}_1^d$

$h_1^e \quad h_2^e \quad h_3^e \quad h_4^e$

Attention

$h_1^d$

Encoder LSTM → Encoder LSTM → Encoder LSTM → Encoder LSTM → Decoder LSTM

$x_1 \quad x_2 \quad x_3 \quad x_4 \quad y_0$

**Intuition:** $\tilde{h}_1^d$ contains information about encoder hidden states that got **high** attention from the decoder

(Bahdanau et al., 2015)

# Attentive Encoder-Decoder Models



(Bahdanau et al., 2015)

# Attentive Encoder-Decoder Models



(Bahdanau et al., 2015)

# Attentive Encoder-Decoder Models

$\alpha_t$

**and the next one…**

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$

Output   Output   **Output**

$; \tilde{h}_3^d]$

**At each time step, a different mixture of the encoder states is computed by the decoder**

"The decoder attends to different outputs of the encoder"

**Enc**
**LS**

And so forth…

$x_1$   $x_2$   $x_3$   $x_4$   $y_0$   $\hat{y}_1$   $\hat{y}_2$

(Bahdanau et al., 2015)

# Attention Recap

- **Main Idea:** Decoder computes a weighted sum of encoder outputs

  - Compute pairwise score between each encoder hidden state and initial decoder hidden state

- Many possible functions for computing scores (dot product, bilinear, etc.)

- **Temporal Bottleneck Fixed! Direct link** between decoder and encoder states

  - Helps with vanishing gradients and modelling long-term dependencies!

- Attention is **agnostic** to the type of RNN used in the encoder and decoder!

# Question

Do any other inefficiencies remain in our sequence to sequence pipelines?

# Encoder is still Recurrent

- **Encoder:** Recurrent functions can't be parallelized because previous state needs to be computed to encode next one



- **Problem: Encoder hidden states must still be computed in series**
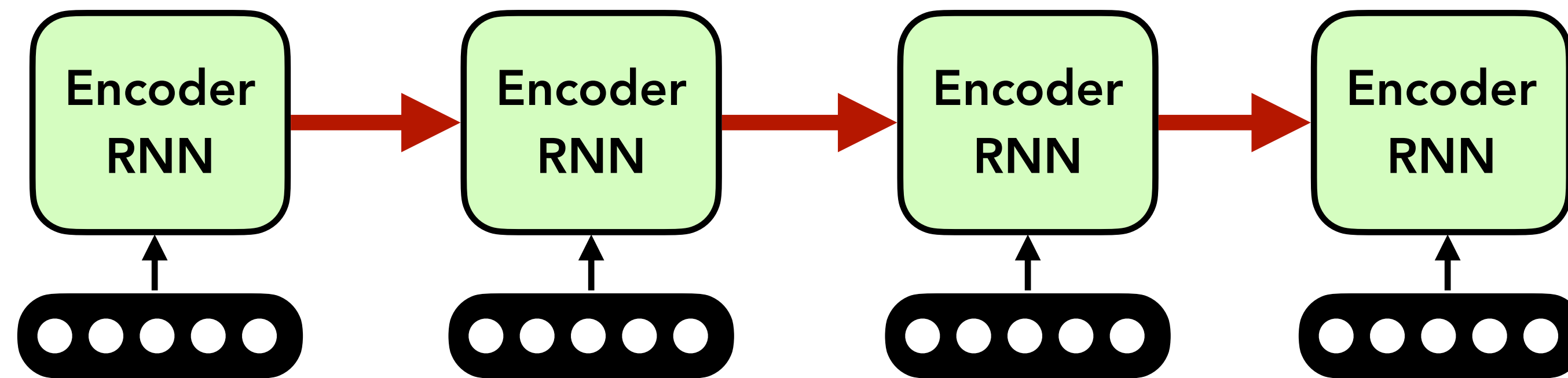
# Encoder is still Recurrent

- **Encoder:** Recurrent functions can't be parallelized because previous state needs to be computed to encode next one



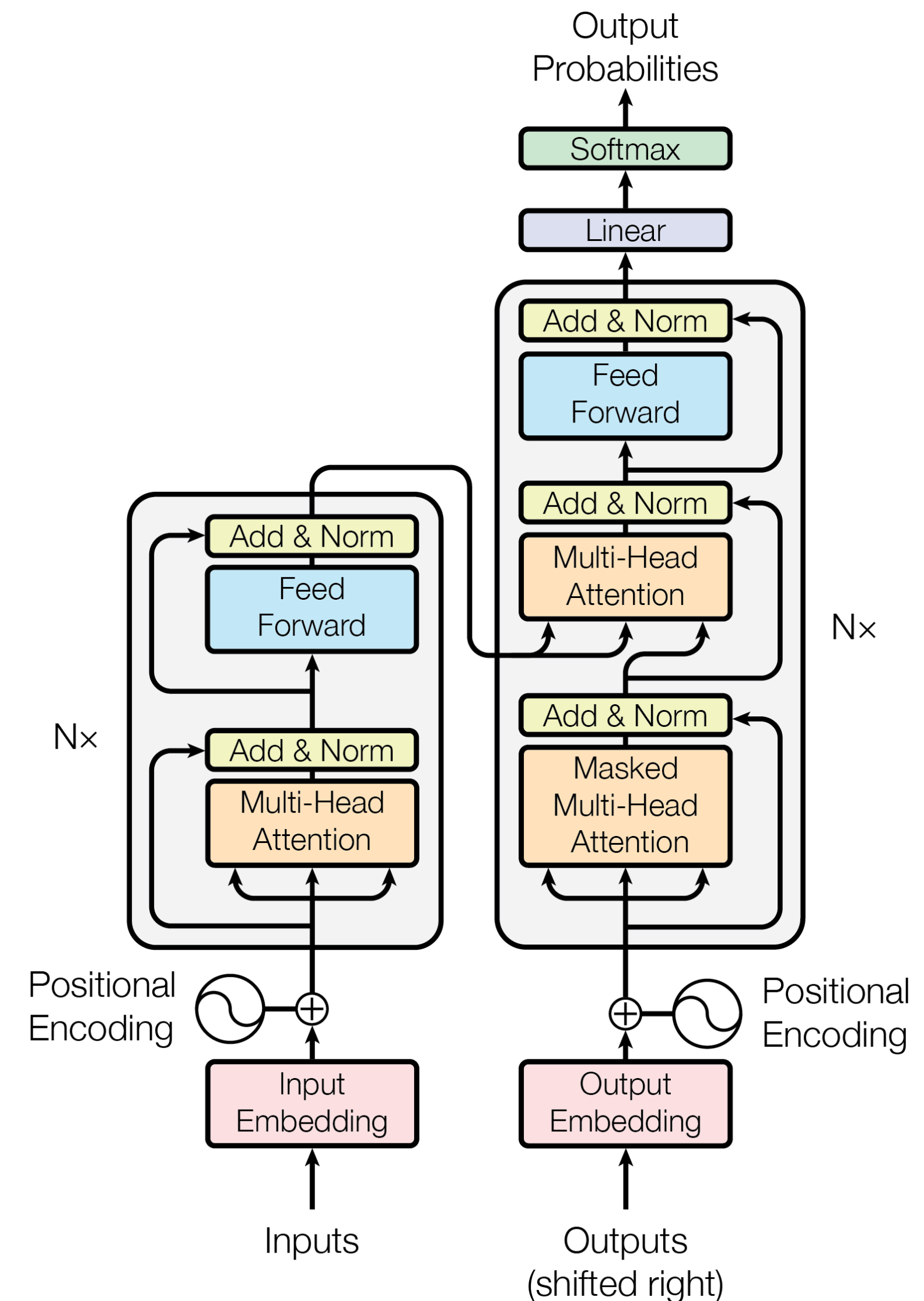- **Problem: Encoder hidden states must still be computed in series**

**Who can think of a task where this might be a problem?**

# Solution:
# **Transformers!**

# Full Transformer

- Made up of encoder and decoder

- Both encoder and decoder made up of multiple cascaded transformer blocks

  - slightly different architecture in encoder and decoder transformer blocks

- Blocks generally made up **multi-headed attention** layers (self-attention) and **feedforward** layers

- No recurrent computations!

  **Encode sequences with self-attention**



Output Probabilities

Softmax

Linear

Add & Norm

Feed Forward

Add & Norm

Multi-Head Attention

Add & Norm

Feed Forward

N×

Add & Norm

Multi-Head Attention

N×

Add & Norm

Masked Multi-Head Attention

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

(Vaswani et al., 2017)

# Self-Attention Toy Example

$h_t^\ell$ = encoder hidden state at time step $t$ at layer $\ell$



$h_1^0$      $h_2^0$      $h_3^0$      $h_4^0$

**"key"**   **"key"**   **"query"**   **"key"**

# Recap: Attention with RNNs

- **Compute** pairwise similarity between each encoder hidden state and decoder hidden state ("idea of what to decode")

$$a_1 = f\left( \begin{array}{c} h_1^e \end{array}, \begin{array}{c} h_1^d \end{array} \right) \qquad a_2 = f\left( \begin{array}{c} h_2^e \end{array}, \begin{array}{c} h_1^d \end{array} \right) \qquad a_3 = f\left( \begin{array}{c} h_3^e \end{array}, \begin{array}{c} h_1^d \end{array} \right)$$

$h_1^e$ **"key"**  $h_1^d$ **"query"** $\qquad$ $h_2^e$ **"key"** $h_1^d$ **"query"** $\qquad$ $h_3^e$ **"key"** $h_1^d$ **"query"**
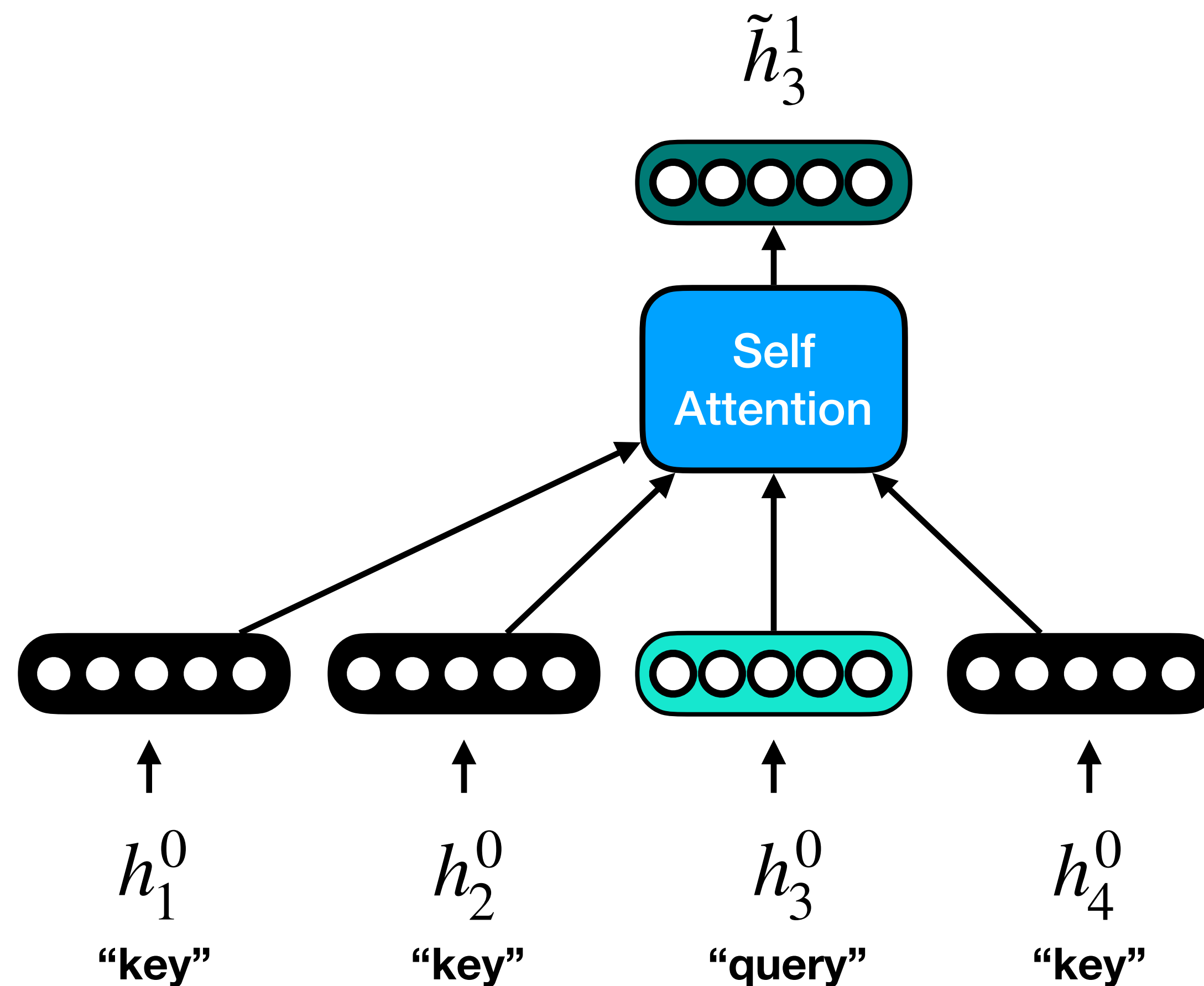
- **Convert** pairwise similarity scores to probability **distribution** (using softmax!) over encoder hidden states and compute weighted average:

**Softmax!** $\boxed{\alpha_t = \dfrac{e^{a_t}}{\sum_j e^{a_j}}} \rightarrow \underset{\alpha_t}{\blacksquare} \rightarrow \tilde{h}_1^d = \sum_{t=1}^{T} \alpha_t h_t^e$ **Here $h_t^e$ is known as the "value"**
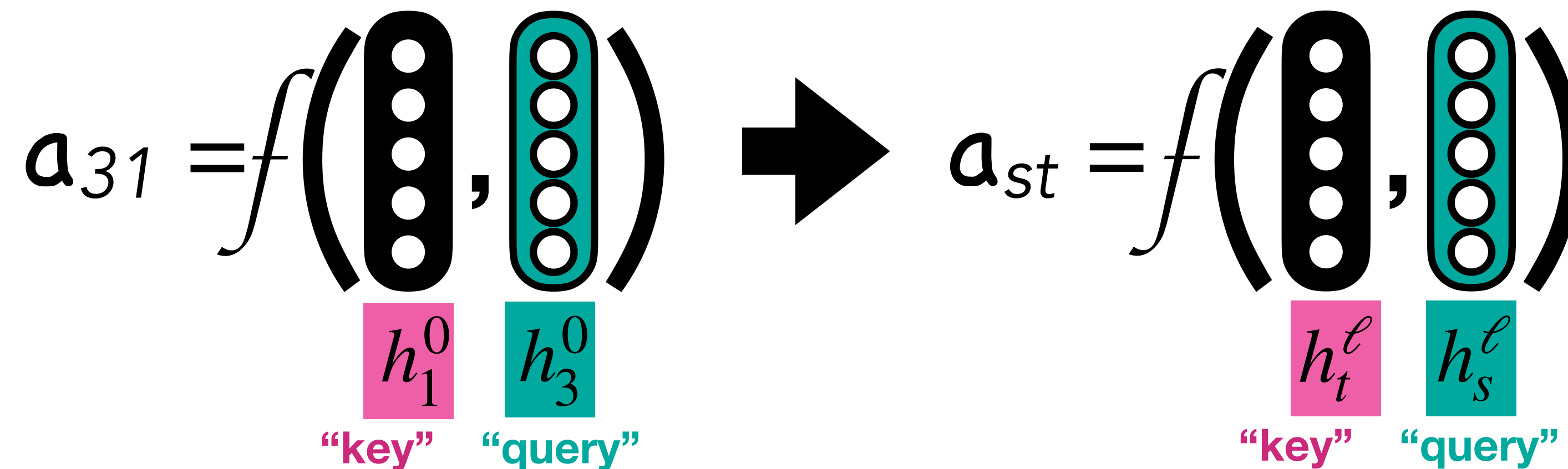
# Self-Attention Toy Example

# Self-Attention Toy Example

$h_t^\ell$ = encoder hidden state at time step $t$ at layer $\ell$

$$a_{31} = f\left( \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \end{array}, \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \right) \quad \Rightarrow \quad a_{st} = f\left( \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \end{array}, \begin{array}{c} \bullet \\ \bullet \\ \bullet \\ \bullet \end{array} \right)$$

$h_1^0$  $h_3^0$         $h_t^\ell$  $h_s^\ell$

**"key"**  **"query"**        **"key"**  **"query"**

$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}} \qquad \alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}} \qquad \tilde{h}_s^\ell = \sum_{t=1}^{T} \alpha_{st}(\mathbf{W}^V h_t^\ell)$$

**Compute pairwise scores**      **Get attention distribution**      **Attend to values to get weighted sum**

# Self-Attention Toy Example

$h_t^\ell$ = encoder hidden state at time step $t$ at layer $\ell$

$$a_{31} = f\left(\begin{matrix}\bullet\\\bullet\\\bullet\\\bullet\\\bullet\end{matrix}, \begin{matrix}\circ\\\circ\\\circ\\\circ\\\circ\end{matrix}\right) \quad\blacktriangleright\quad a_{st} = f\left(\begin{matrix}\bullet\\\bullet\\\bullet\\\bullet\\\bullet\end{matrix}, \begin{matrix}\circ\\\circ\\\circ\\\circ\\\circ\end{matrix}\right)$$

$h_1^0$   $h_3^0$

"key"   "query"

$h_t^\ell$   $h_s^\ell$

"key"   "query"

{1, …, t, …, T}
includes s!

$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}}$$
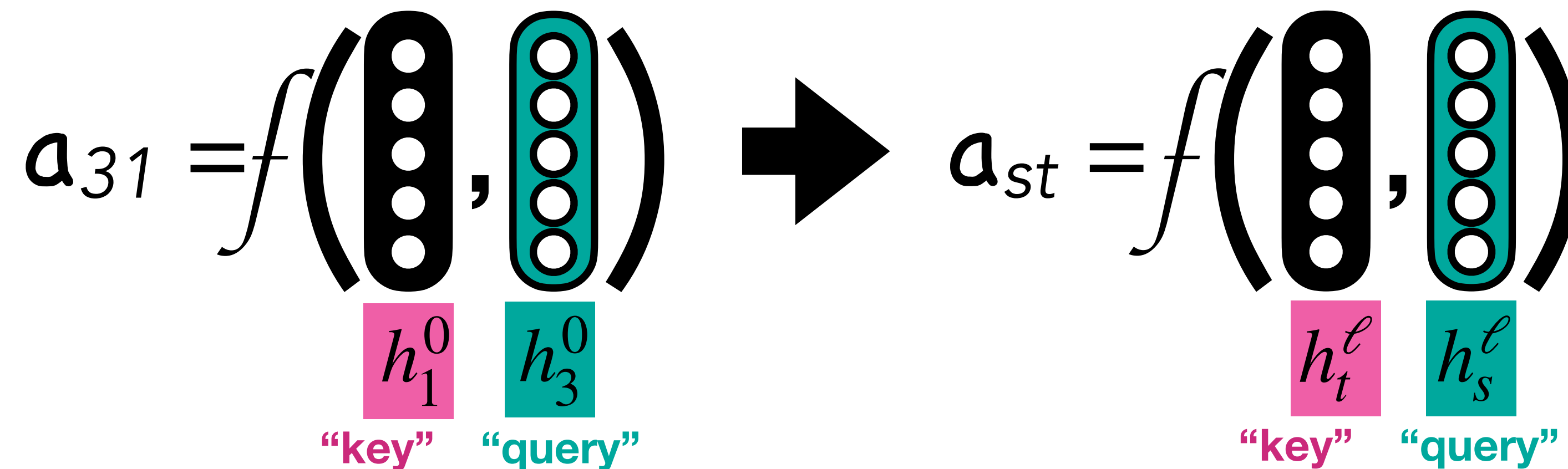
$$\alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}}$$

$$\tilde{h}_s^\ell = \sum_{t=1}^{T} \alpha_{st} (\mathbf{W}^V h_t^\ell)$$

Self-attention!

**Compute pairwise scores**

**Get attention distribution**

**Attend to values to get weighted sum**
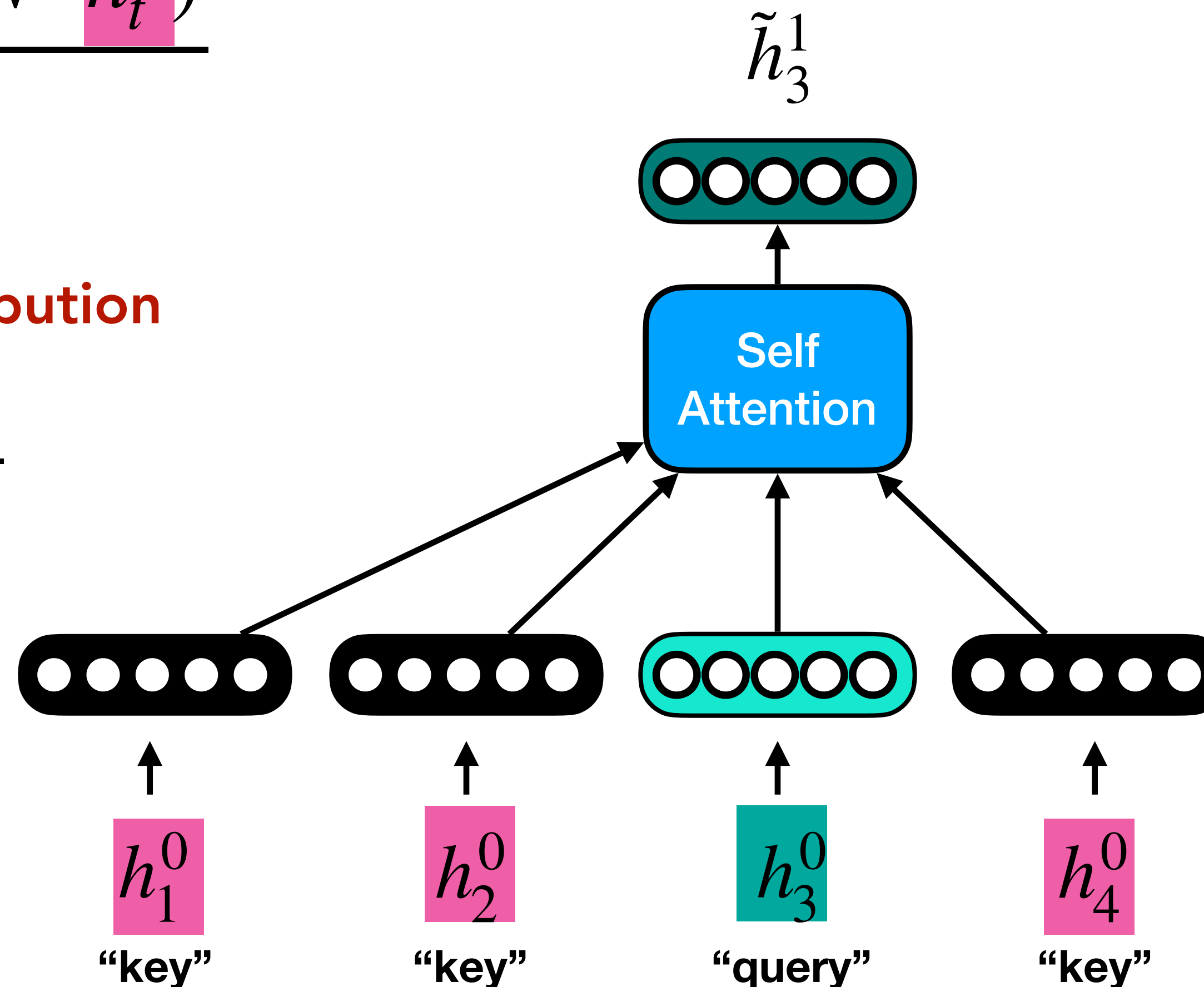
# Self-Attention Toy Example

**Compute pairwise scores**

$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T(\mathbf{W}^K h_t^\ell)}{\sqrt{d}}$$

**Attend to values to get weighted sum**

$$\tilde{h}_s^\ell = \sum_{t=1}^{T} \alpha_{st}(\mathbf{W}^V h_t^\ell)$$

**Get attention distribution**

$$\alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}}$$

$\tilde{h}_3^1$

Self Attention

$h_1^0$    $h_2^0$    $h_3^0$    $h_4^0$

"key"    "key"    "query"    "key"

# Self-Attention Toy Example

Compute pairwise scores
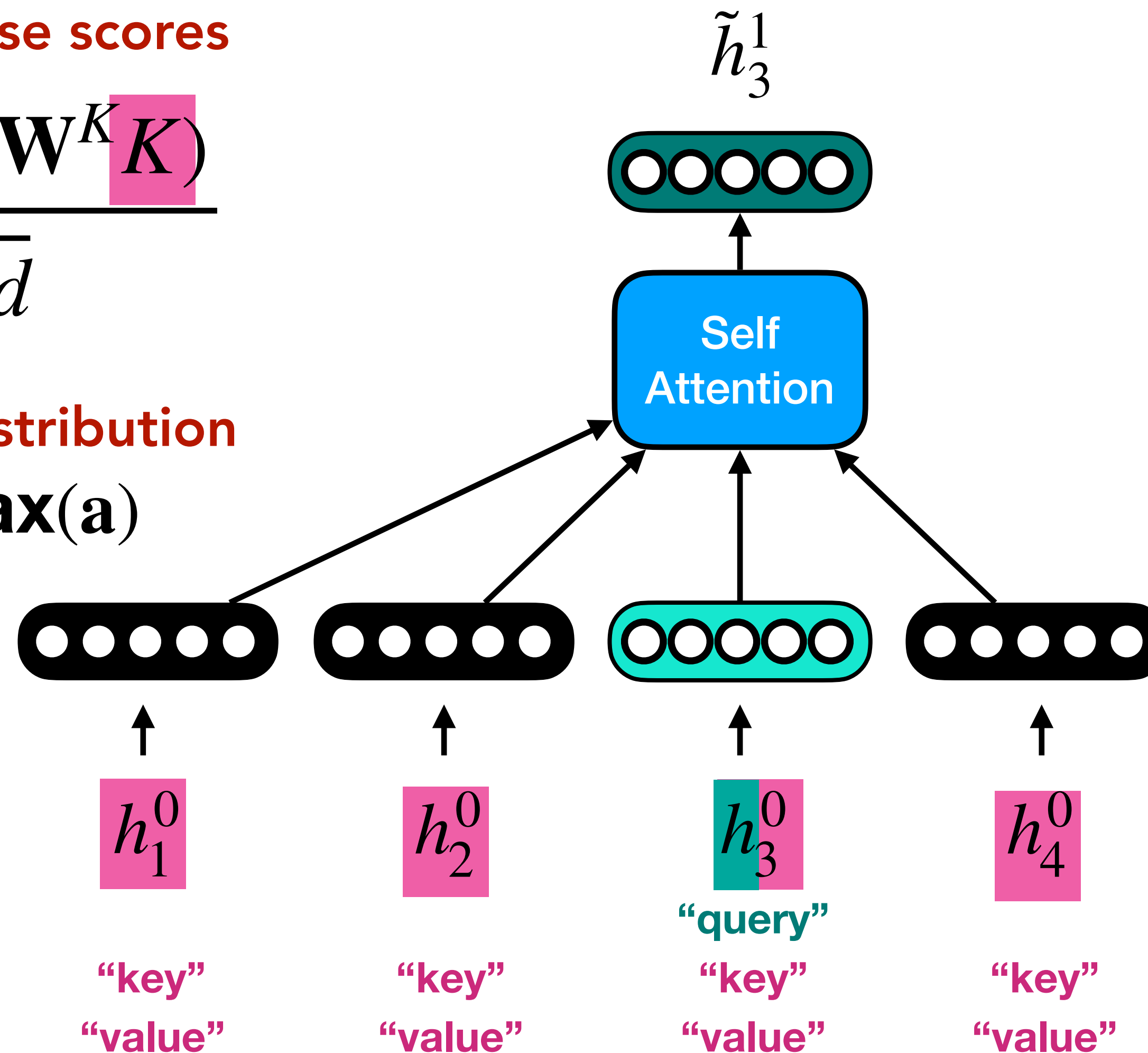
$$\mathbf{a} = \frac{(\mathbf{W}^Q q)(\mathbf{W}^K K)}{\sqrt{d}}$$

Get attention distribution

$$\alpha = \mathbf{softmax}(\mathbf{a})$$

$\tilde{h}_3^1$

Self Attention

Attend to values to get weighted sum

$$\tilde{h}^\ell = W^O \alpha\big(V \mathbf{W}^V\big)$$

"query" $\quad q = h_s^\ell$

"values"

$K = V = \{h_t^\ell\}_{t=0}^T$

"keys"

$h_1^0 \qquad h_2^0 \qquad h_3^0 \qquad h_4^0$

"query"

"key"     "key"     "key"     "key"

"value"    "value"    "value"    "value"

For each attention computation, every element is a key and value, and one element is a query

# Self-Attention Toy Example

Compute pairwise scores

$$\mathbf{a} = \frac{(\mathbf{W}^Q q)(\mathbf{W}^K K)}{\sqrt{d}}$$

Get attention distribution

$$\alpha = \mathbf{softmax}(\mathbf{a})$$

$$\tilde{h}_2^1$$

Self Attention

Attend to values to get weighted sum

$$\tilde{h}^\ell = W^O \alpha\left(V \mathbf{W}^V\right)$$

"query" $\quad q = h_s^\ell$

"values"

$$K = V = \{h_t^\ell\}_{t=0}^T$$

"keys"

$$h_1^0 \qquad h_2^0 \qquad h_3^0 \qquad h_4^0$$

"query"

"key"  "key"  "key"  "key"
"value" "value" "value" "value"

For each attention computation, every element is a key and value, and one element is a query
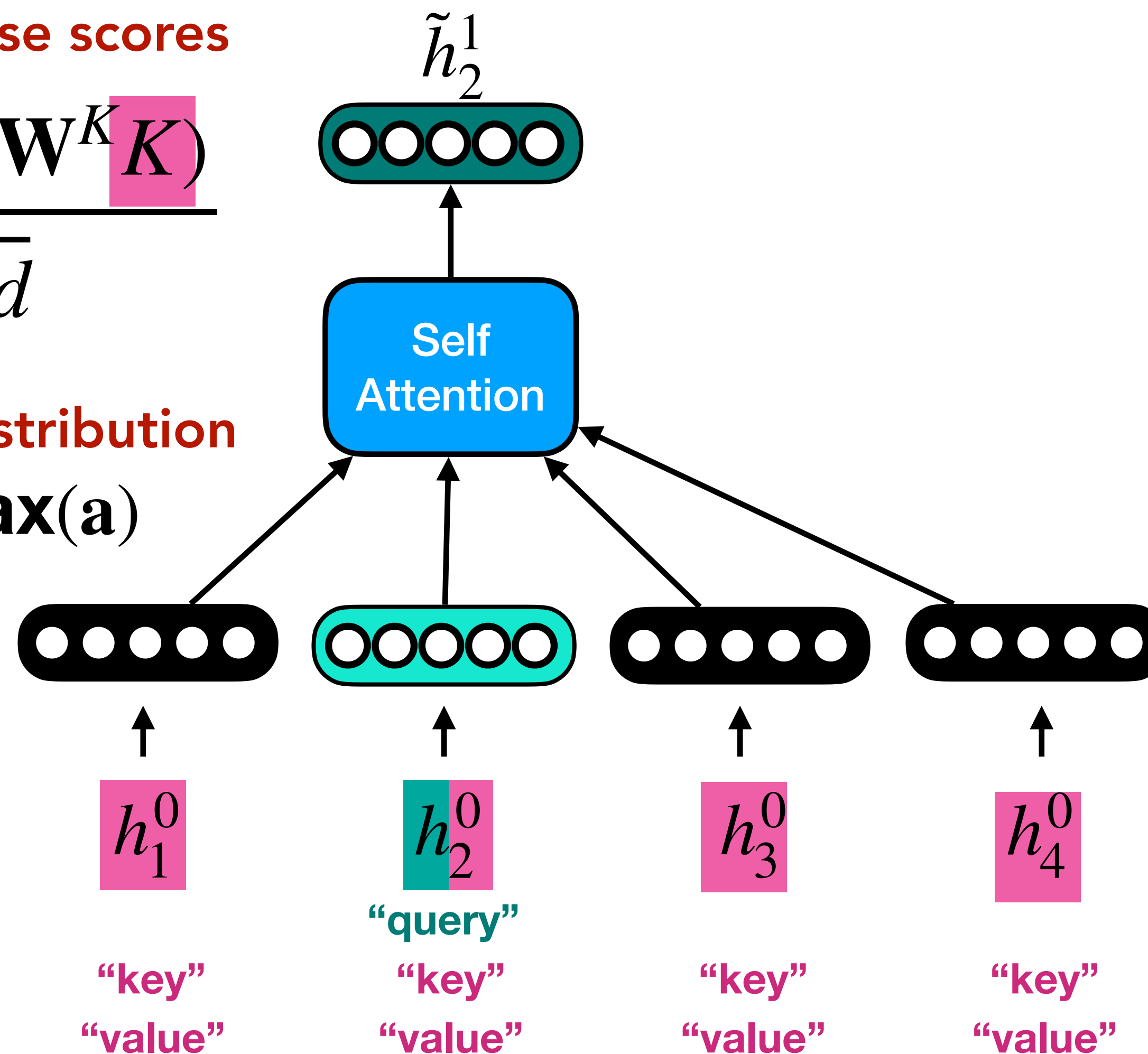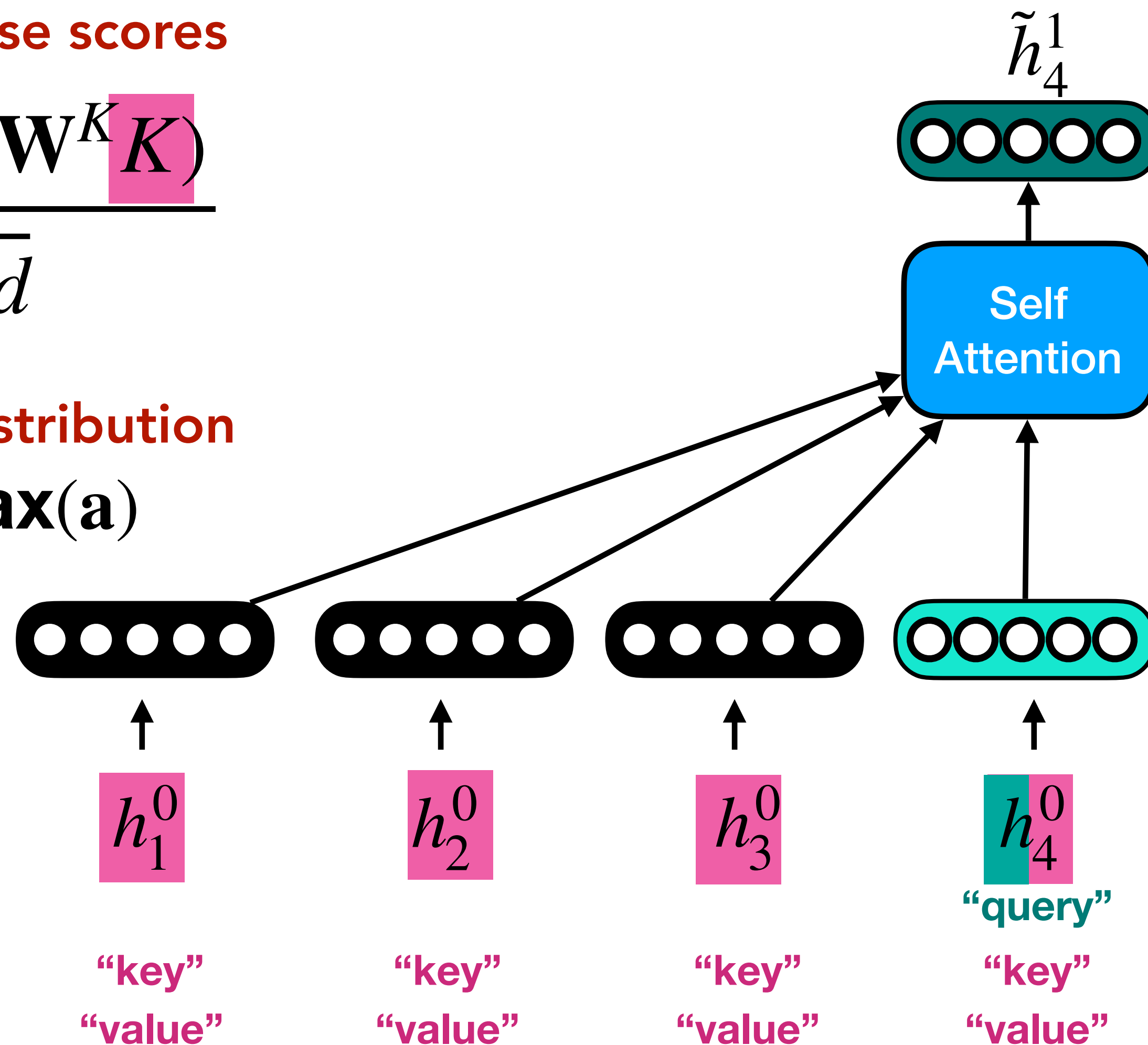
# Self-Attention Toy Example

**Compute pairwise scores**

$$\mathbf{a} = \frac{(\mathbf{W}^Q q)(\mathbf{W}^K K)}{\sqrt{d}}$$

**Get attention distribution**

$$\alpha = \mathbf{softmax}(\mathbf{a})$$

$\tilde{h}_4^1$

Self Attention

**Attend to values to get weighted sum**

$$\tilde{h}^\ell = W^O \alpha(V \mathbf{W}^V)$$

"query" $\quad q = h_s^\ell$

"values"

$K = V = \{h_t^\ell\}_{t=0}^T$

"keys"

$h_1^0 \qquad h_2^0 \qquad h_3^0 \qquad h_4^0$

"query"

"key" "value"  "key" "value"  "key" "value"  "key" "value"

**For each attention computation, every element is a key and value, and one element is a query**

# Self-Attention Toy Example



$$\tilde{h}_1^1 = \mathbf{Attention}\left(h_1^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

$$\tilde{h}_1^2 = \mathbf{Attention}\left(h_2^0, \{h_t^0\}_{t=0}^{t=3}\right)$$
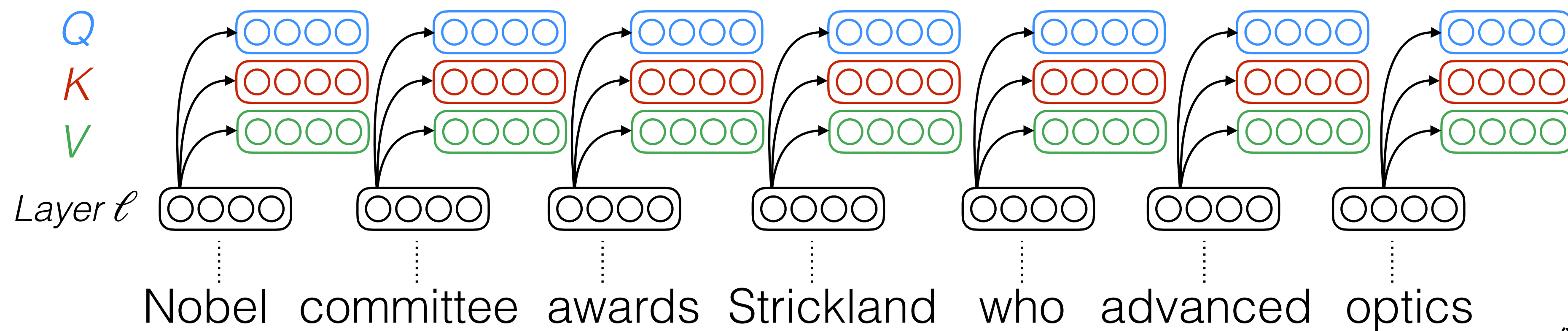
$$\tilde{h}_1^3 = \mathbf{Attention}\left(h_3^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

$$\tilde{h}_1^4 = \mathbf{Attention}\left(h_4^0, \{h_t^0\}_{t=0}^{t=3}\right)$$

# Self-attention (in encoder)
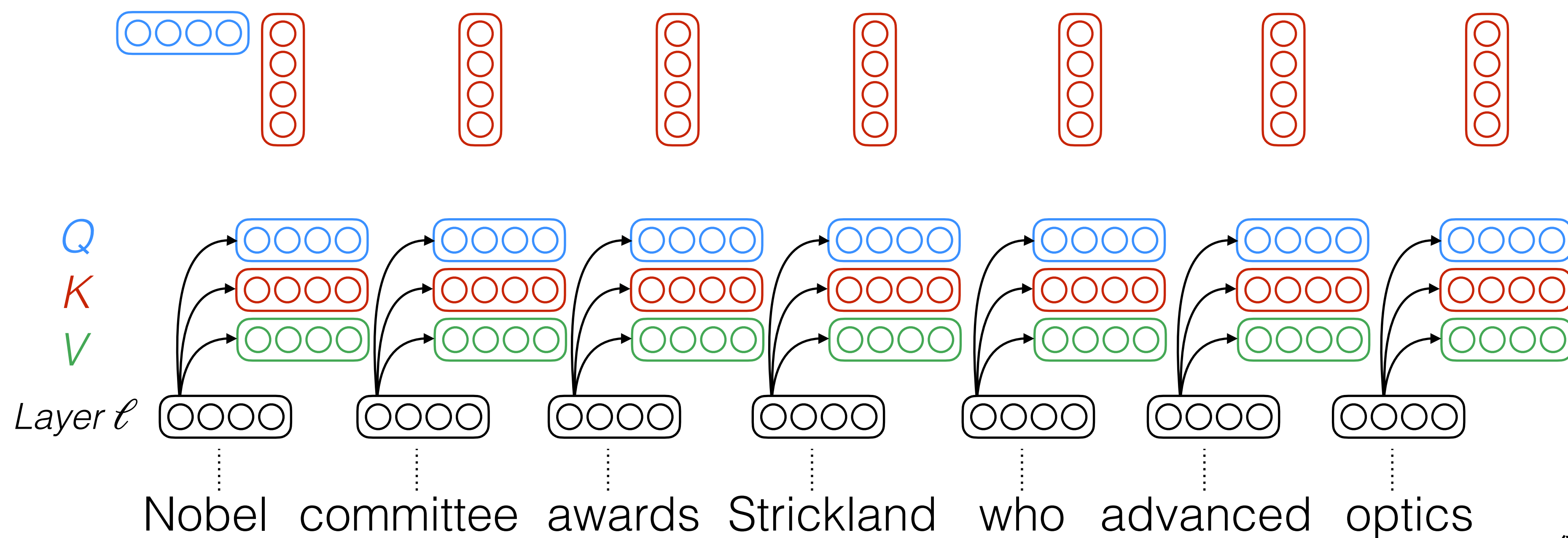


Q
K
V
Layer ℓ

Nobel  committee  awards  Strickland  who  advanced  optics

(Vaswani et al., 2017)

# Self-attention (in encoder)

Q
K
V

*Layer ℓ*

Nobel  committee  awards  Strickland  who  advanced  optics

(Vaswani et al., 2017)

# Self-attention (in encoder)



optics
advanced
who
Strickland
awards
committee
Nobel

*Q*

*K*

*V*

*Layer ℓ*

Nobel  committee  awards  Strickland  who  advanced  optics

(Vaswani et al., 2017)

# Self-attention (in encoder)



optics
advanced
who
Strickland
awards
committee
Nobel

$A$

$Q$

$K$

$V$

*Layer* $\ell$

Nobel    committee    awards    Strickland    who    advanced    optics

(Vaswani et al., 2017)

# Self-attention (in encoder)



Layer $\ell$

Nobel  committee  awards  Strickland  who  advanced  optics

(Vaswani et al., 2017)

# Self-attention (in encoder)



optics
advanced
who
Strickland
awards
committee
Nobel

*A*

*Q*

*K*

*V*

*Layer ℓ*

Nobel   committee   awards   Strickland   who   advanced   optics

(Vaswani et al., 2017)

# Self-attention (in encoder)



$M$

optics
advanced
who
Strickland
awards
committee
Nobel
$A$

$Q$
$K$
$V$

*Layer* $\ell$

Nobel   committee   awards   Strickland   who   advanced   optics

(Vaswani et al., 2017)

# Self-attention (in encoder)



*M*

optics
advanced
who
Strickland
awards
committee
Nobel

*A*

*Q*

*K*

*V*

*Layer ℓ*

Nobel   committee   awards   Strickland   who   advanced   optics

(Vaswani et al., 2017)

# Multi-head self-attention



$M_H$
$M_1$

$A$

optics
advanced
who
Strickland
awards
committee
Nobel

$Q$
$K$
$V$

*Layer* $\ell$

Nobel    committee    awards    Strickland    who    advanced    optics

(Vaswani et al., 2017)

# Multi-head self-attention



$M_H$
$M_1$

optics
advanced
who
Strickland
awards
committee
Nobel

$A$

$Q$
$K$
$V$

Layer $\ell$

Nobel    committee    awards    Strickland    who    advanced    optics

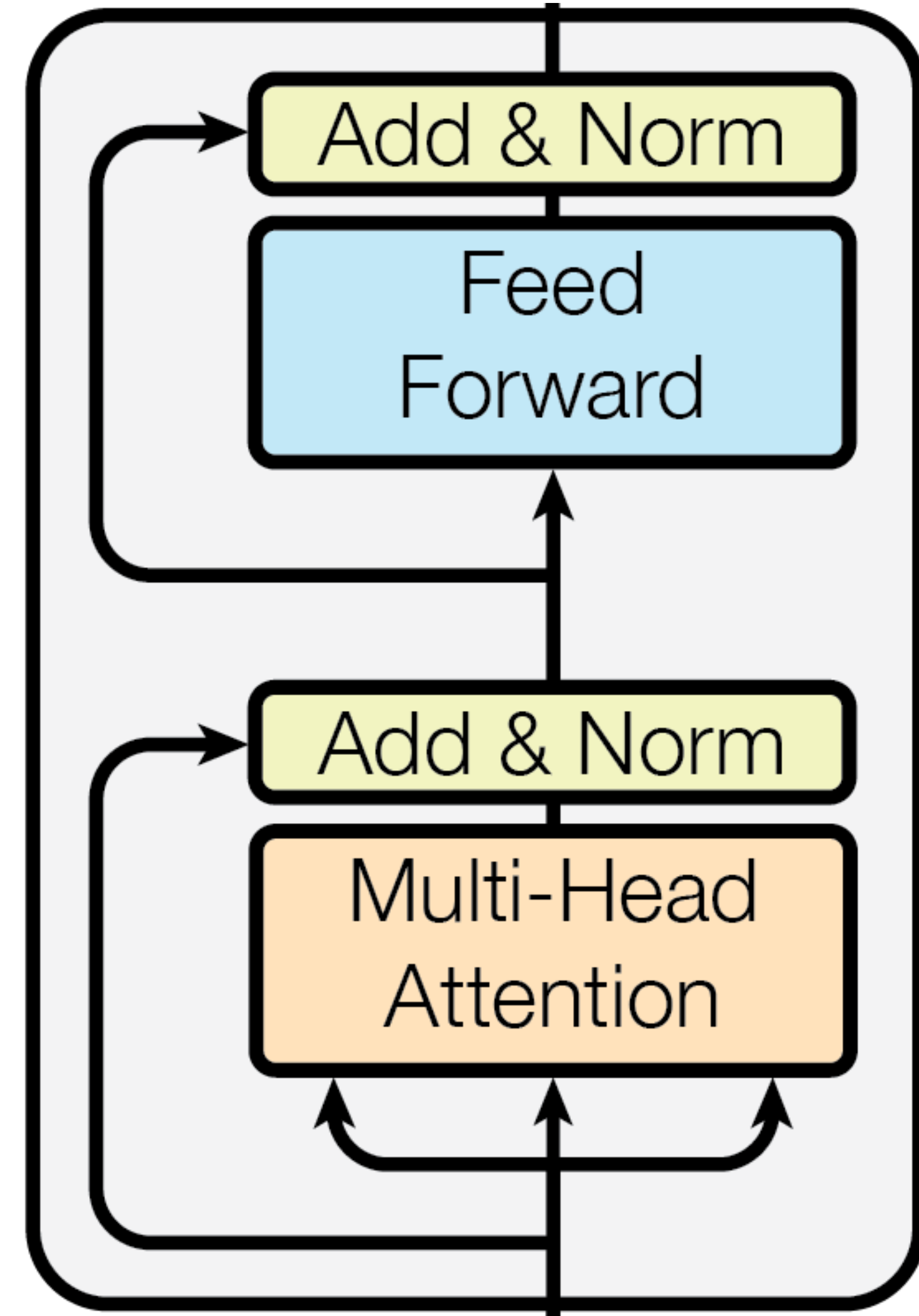(Vaswani et al., 2017)

# Question

What are two advantages of self-attention over recurrent models?

# Transformer Block

- Multi-headed attention is the main innovation of the transformer model!



Vaswani et al., 2017

# Transformer Block

- Multi-headed attention is the main innovation of the transformer model!

- Each block also composed of:

  - a layer normalisations

  - a feedforward network

  - residual connections



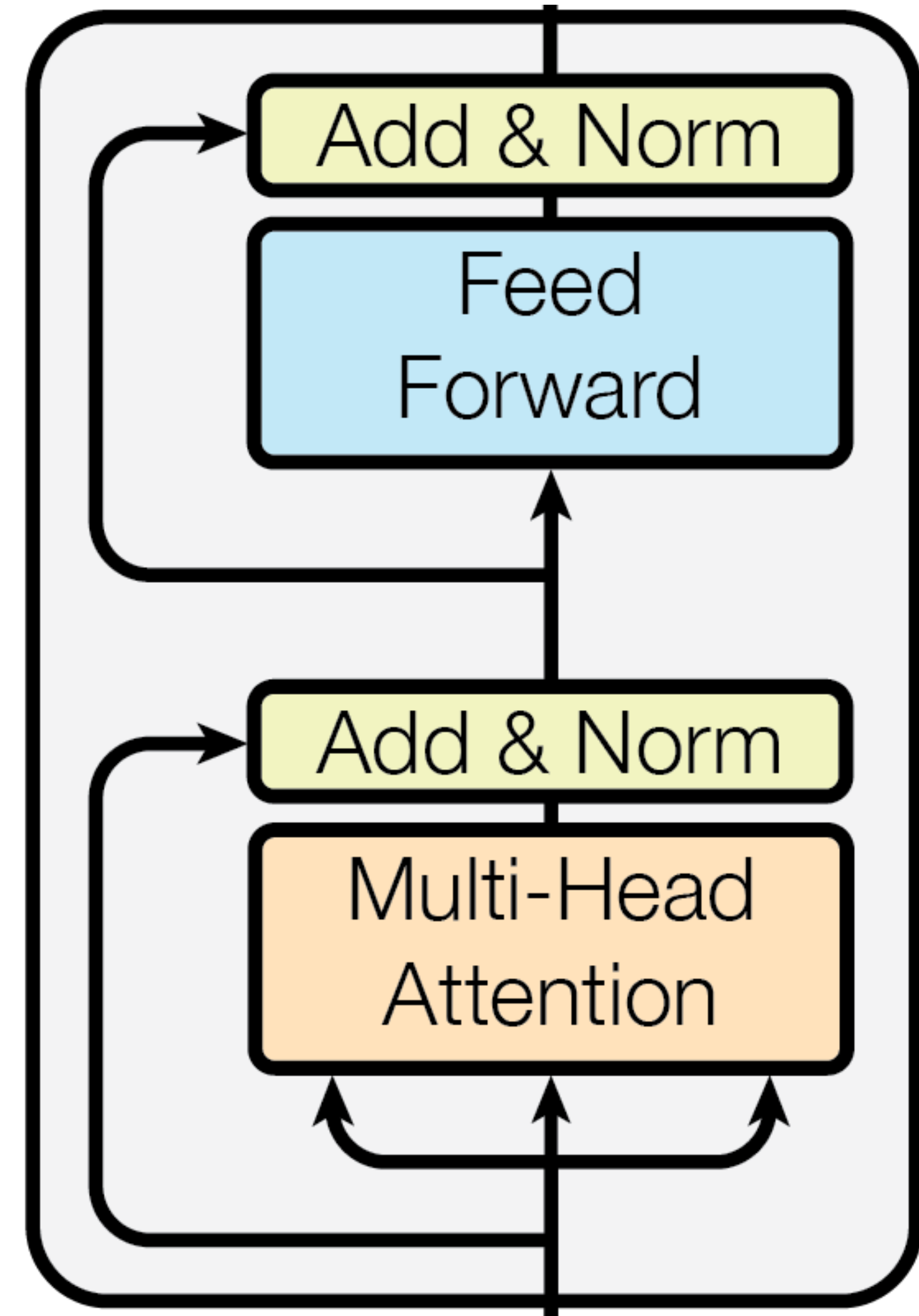Vaswani et al., 2017

# LayerNorm & Residual Connections

- **Layer Normalisation**

  - Normalize the outputs of different modules

  $$y = \frac{x - \mathbf{E}[x]}{\sqrt{\mathrm{Var}[x] + \epsilon}} * \gamma + \beta$$
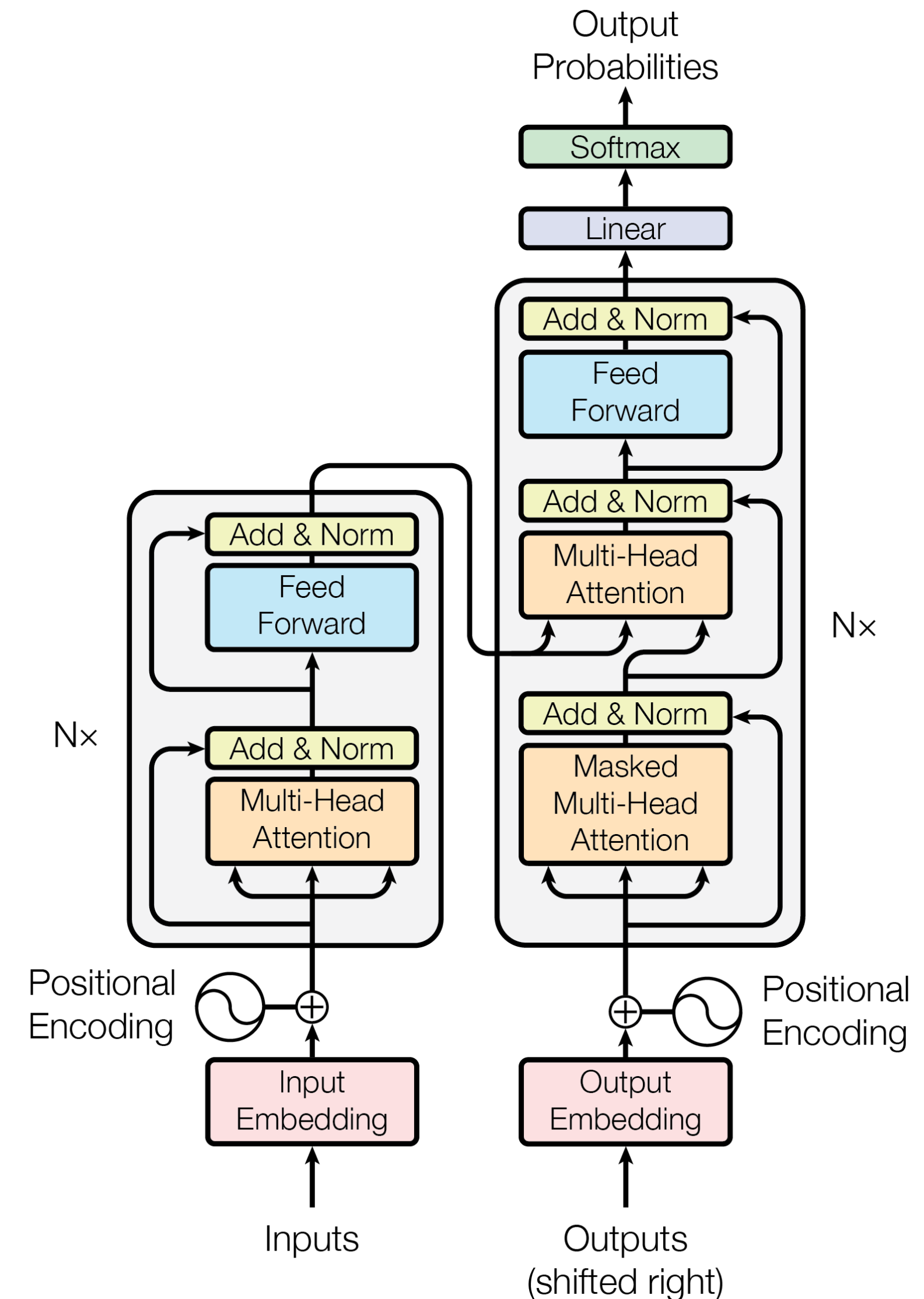
- **Residual Connections**

  - Add the input of a module to its output

  - $\mathrm{LayerNorm}(x + \mathrm{Sublayer}(x))$

# Full Transformer

- Full transformer encoder is multiple cascaded transformer blocks

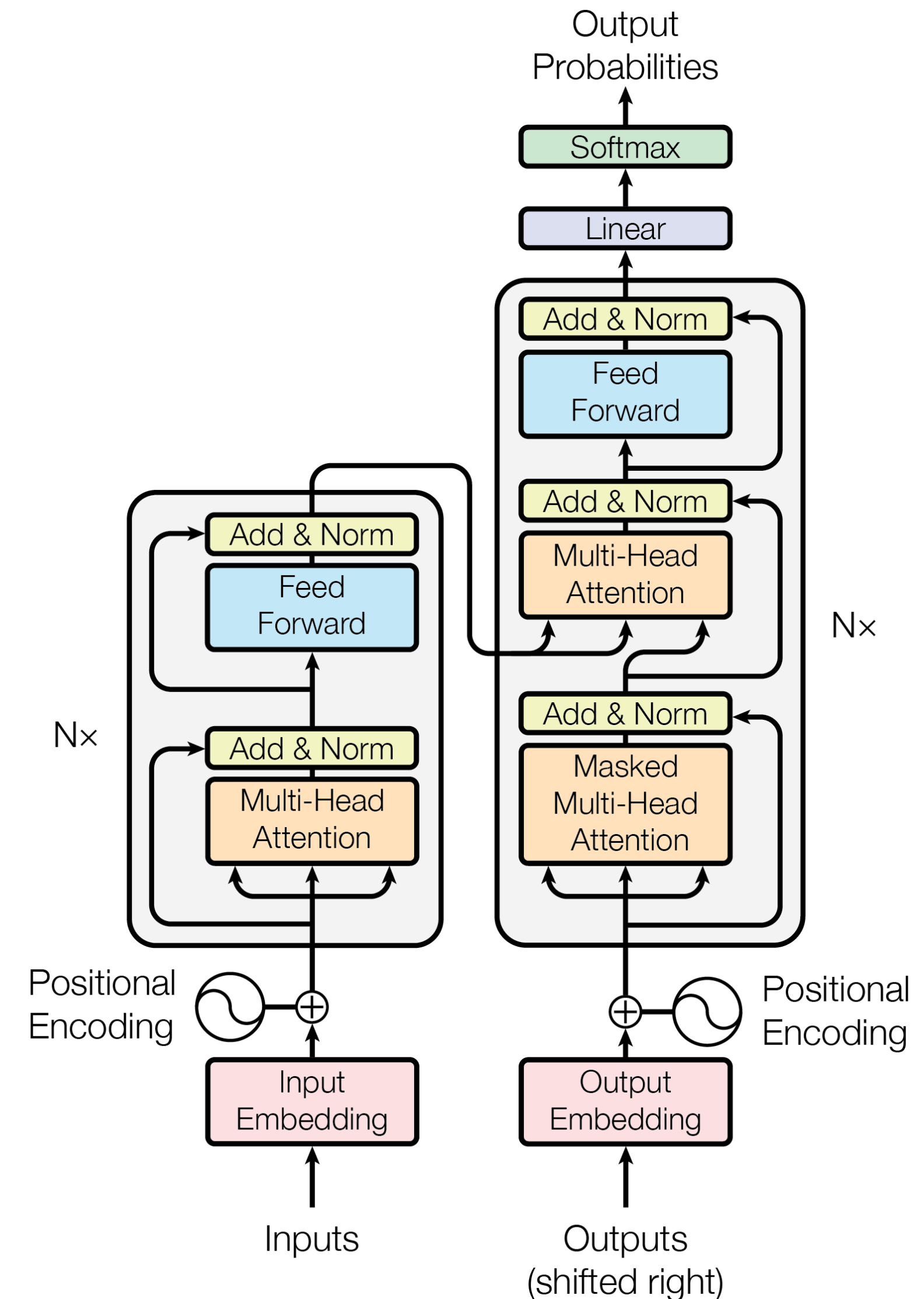  - **build up compositional representations of inputs**



Vaswani et al., 2017

# Full Transformer

- Full transformer encoder is multiple cascaded transformer blocks

  - **build up compositional representations of inputs**

- Transformer decoder (right) similar to encoder

  - First layer of block is **masked** multi-headed attention

  - Second layer is multi-headed attention over *final-layer* encoder outputs **(cross-attention)**

  - Third layer is feed-forward network

Output Probabilities

Softmax

Linear

Add & Norm
Feed Forward

Add & Norm
Multi-Head Attention

Add & Norm
Feed Forward

Nx

Add & Norm
Multi-Head Attention

Add & Norm
Masked Multi-Head Attention

Nx

Positional Encoding

Positional Encoding

Input Embedding

Output Embedding

Inputs

Outputs (shifted right)

Vaswani et al., 2017

# Question

What is an issue with self-attention for the decoder?

# Masked Multi-headed Attention

- Self-attention can attend to any token in the sequence

- For the decoder, **you don't want tokens to attend to future tokens**

  - Decoder used to generate text (i.e., machine translation)

# Masked Multi-headed Attention

- Self-attention can attend to any token in the sequence

- For the decoder, **you don't want tokens to attend to future tokens**

  - Decoder used to generate text (i.e., machine translation)



raw attention weights          mask

$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}}$$ ➡️ $$a_{st} := a_{st} - \infty \; ; \; s < t$$
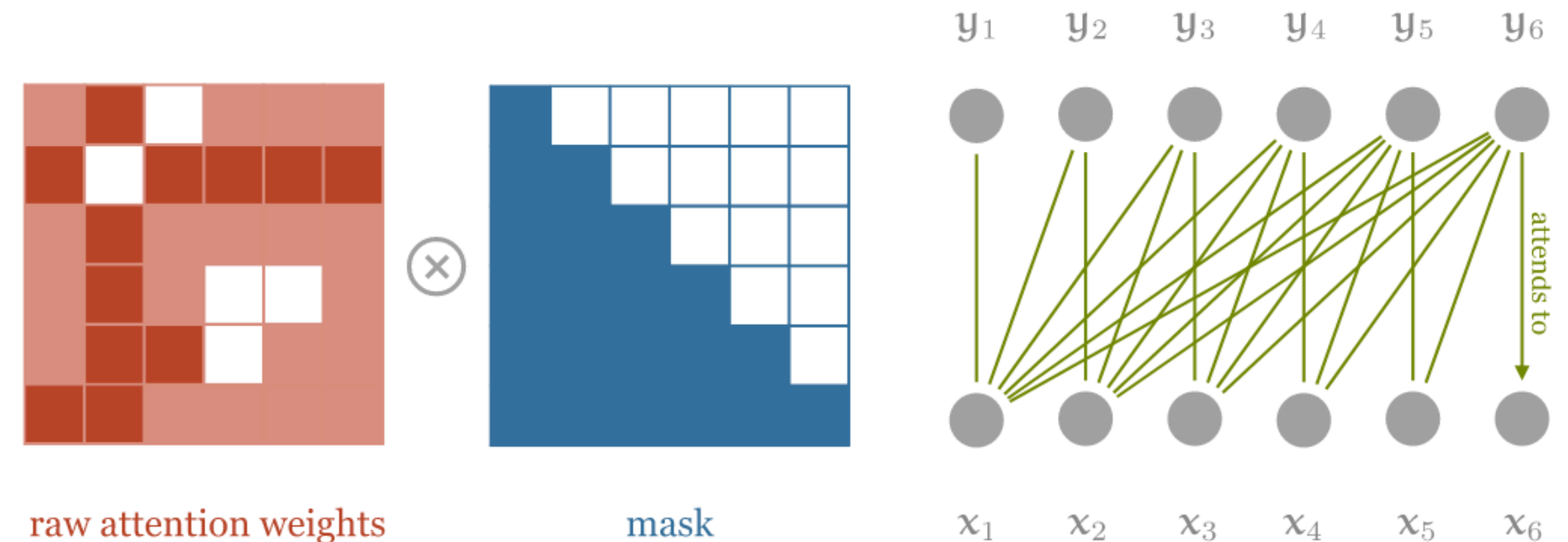
# Masked Multi-headed Attention

- Self-attention can attend to any token in the sequence

- For the decoder, **you don't want tokens to attend to future tokens**

  - Decoder used to generate text (i.e., machine translation)

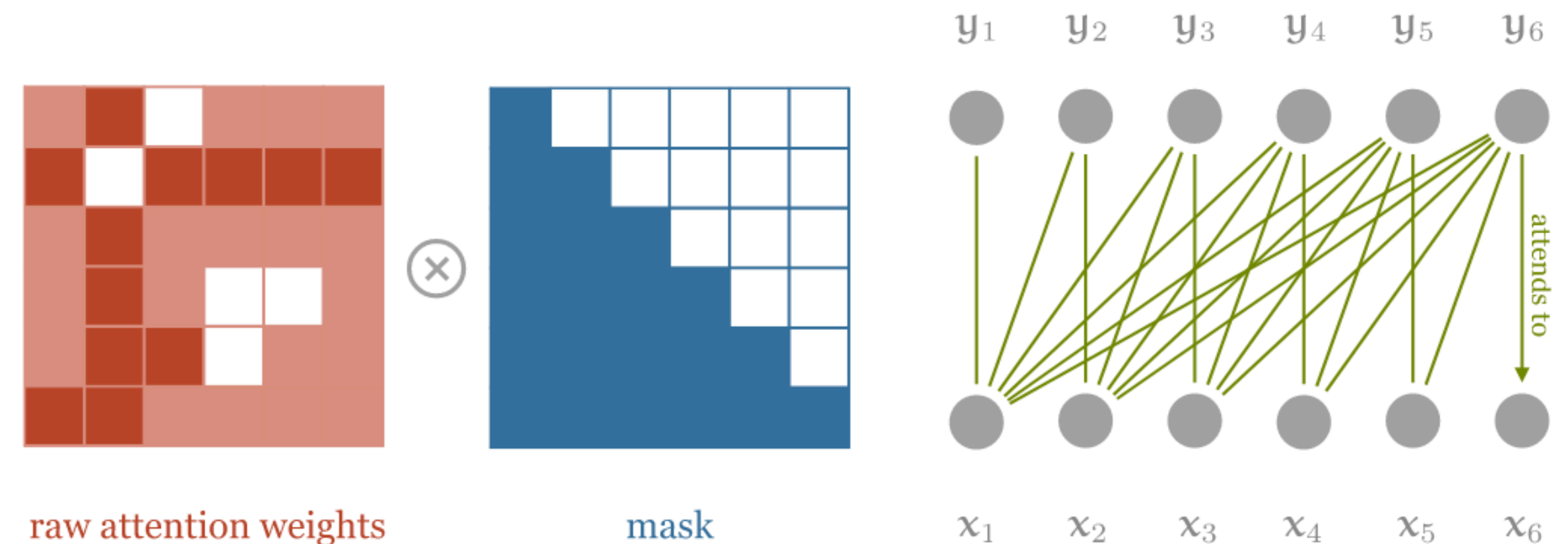**Mask the attention scores of future tokens so their attention = 0**



raw attention weights          mask          $x_1$  $x_2$  $x_3$  $x_4$  $x_5$  $x_6$

$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}} \quad \Longrightarrow \quad a_{st} := a_{st} - \infty \; ; \; s < t \quad \Longrightarrow \quad \alpha_{st} = \frac{e^{a_{st}}}{\sum_j e^{a_{sj}}} = 0$$

# Cross-attention

- **Cross attention** is the same classical attention as in the RNN encoder-decoder model

- The query to the attention function is the output of the masked multi-headed attention in the decoder (i.e., a decoder state)

- The keys and values are the output of the **final** encoder transformer

- Once again, a representation from the decoder is used to **attend** to the encoder outputs



Vaswani et al., 2017

# Full Transformer

- Full transformer encoder is multiple cascaded transformer blocks

  -

- Tr

  -

  Second layer is multi-headed attention over encoder outputs **(cross-attention)**

  - Third layer is feed-forward network

**Recurrent models provided word order information**

**Does self-attention provide word order information?**

Output Probabilities

Softmax

Linear

Multi-Head Attention

Attention

Positional Encoding

Input Embedding

Inputs

Positional Encoding

Output Embedding

Outputs (shifted right)

Vaswani et al., 2017

# Position Embeddings

Positional
Encoding

⊕

Input
Embedding

↑

Inputs

⊕

Positional
Encoding

Output
Embedding

↑

Outputs
(shifted right)

Vaswani et al., 2017

# Position Embeddings



Positional Encoding → Input Embedding ← Inputs

Positional Encoding → Output Embedding ← Outputs (shifted right)

- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position

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



Dimension

Index in the sequence

Vaswani et al., 2017

# Position Embeddings



Positional Encoding → Input Embedding ← Inputs

Positional Encoding → Output Embedding ← Outputs (shifted right)
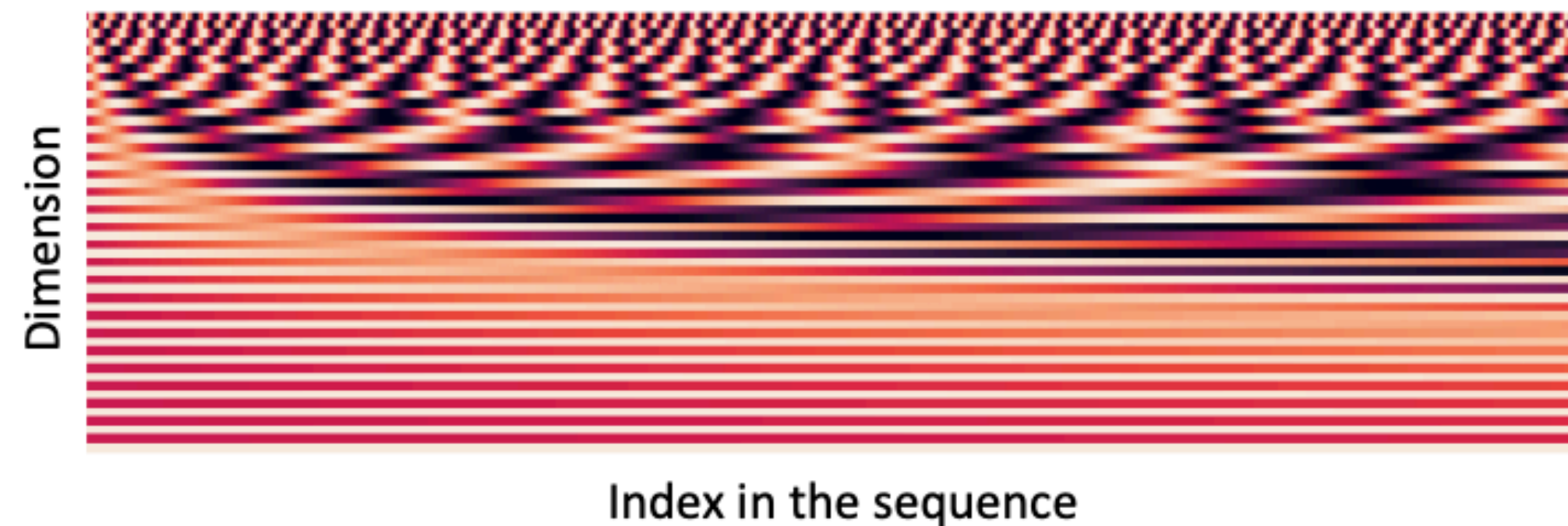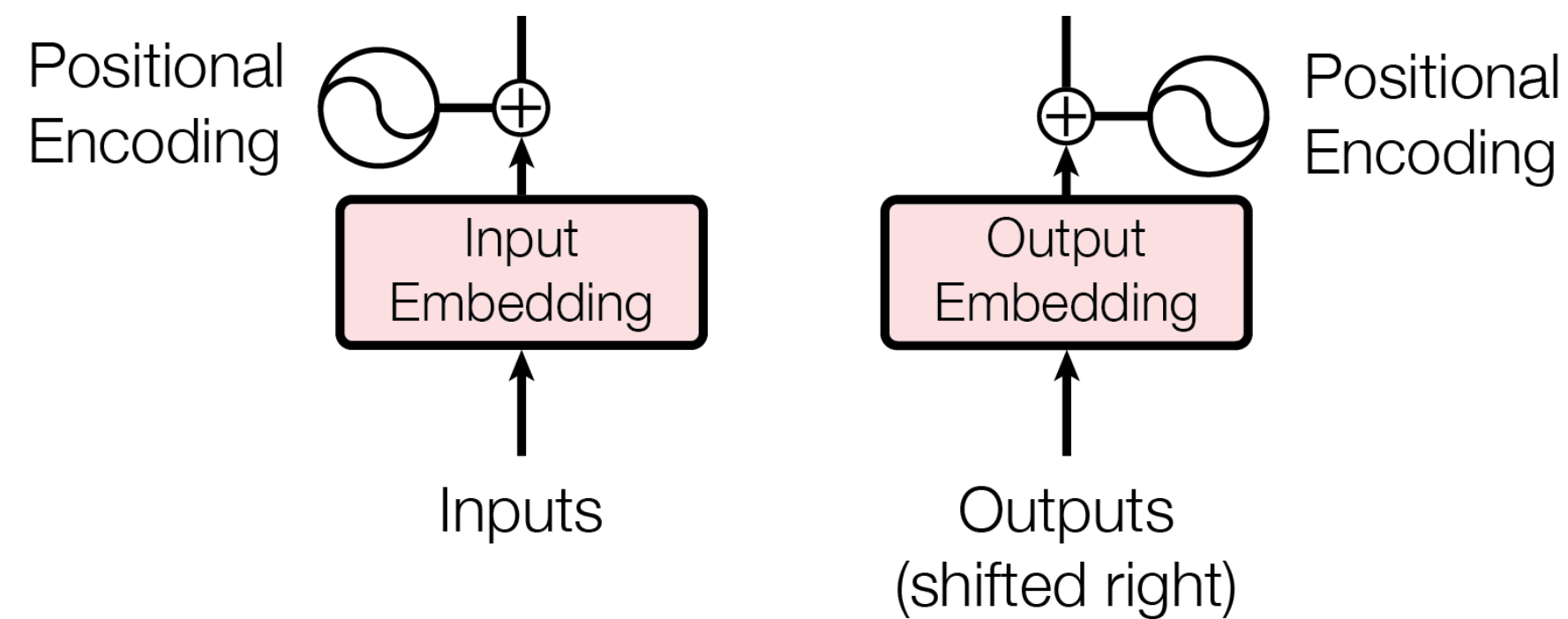
- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position

- **In practice, easiest is to learn position embeddings from scratch**

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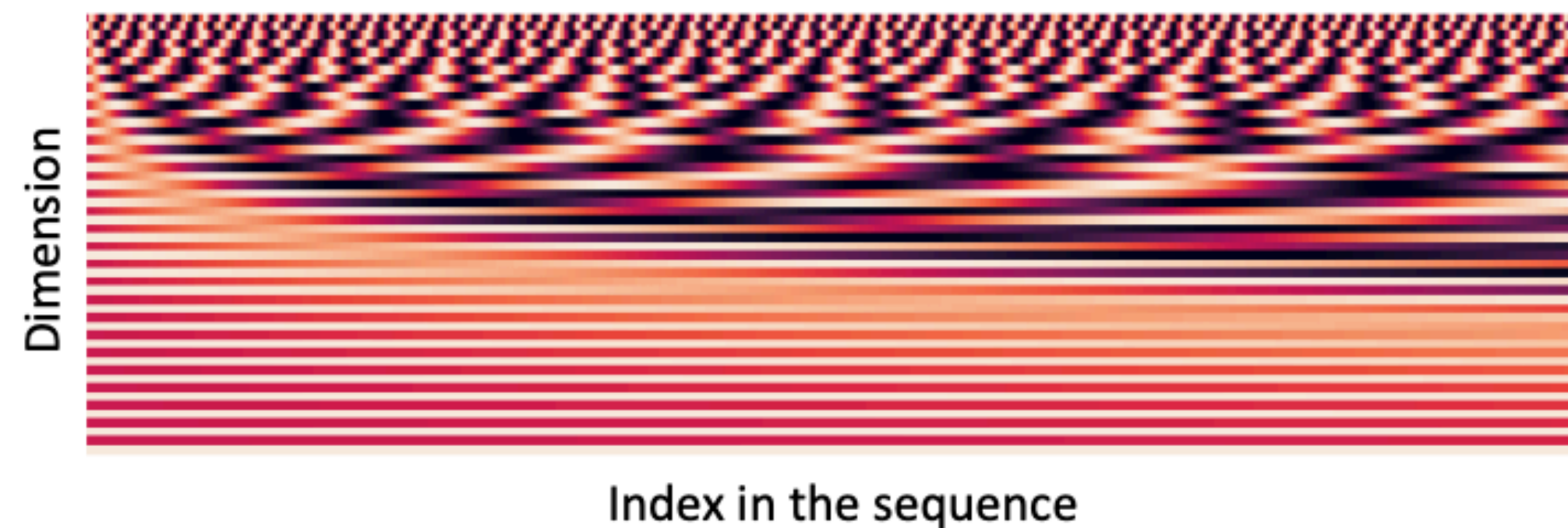


Dimension

Index in the sequence

Vaswani et al., 2017

# Question

What might be a disadvantage of using learned position embeddings?

Poor generalisation to sequences longer than the maximum position embedding you have learned

# Position Embeddings

Lots of potential for new methods that generalise to longer sequences

Position embeddings remain an active area of research

- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position

- **In practice, easiest is to learn position embeddings from scratch**

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

Index in the sequence

Vaswani et al., 2017

# Performance: Machine Translation
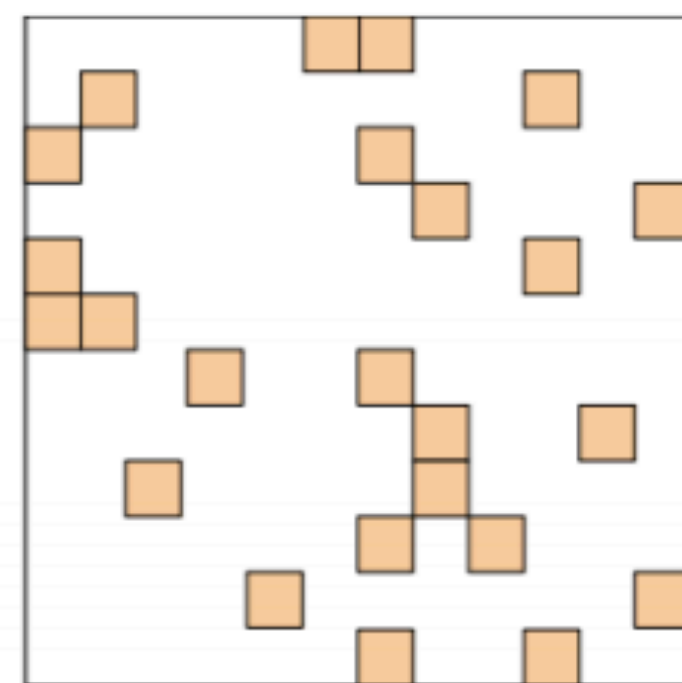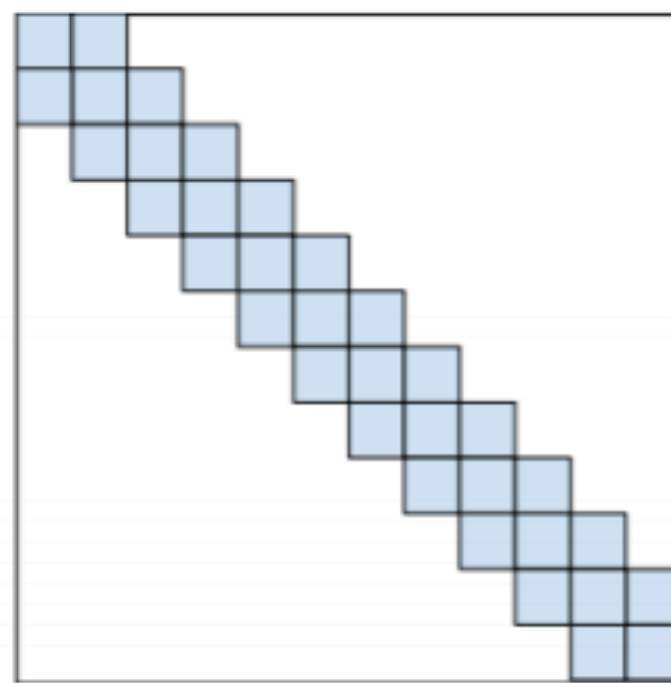
| Model | BLEU | | Training Cost (FLOPs) | |
|---|---|---|---|---|
| | EN-DE | EN-FR | EN-DE | EN-FR |
| ByteNet [15] | 23.75 | | | |
| Deep-Att + PosUnk [32] | | 39.2 | | $1.0 \cdot 10^{20}$ |
| GNMT + RL [31] | 24.6 | 39.92 | $2.3 \cdot 10^{19}$ | $1.4 \cdot 10^{20}$ |
| ConvS2S [8] | 25.16 | 40.46 | $9.6 \cdot 10^{18}$ | $1.5 \cdot 10^{20}$ |
| MoE [26] | 26.03 | 40.56 | $2.0 \cdot 10^{19}$ | $1.2 \cdot 10^{20}$ |
| Deep-Att + PosUnk Ensemble [32] | | 40.4 | | $8.0 \cdot 10^{20}$ |
| GNMT + RL Ensemble [31] | 26.30 | 41.16 | $1.8 \cdot 10^{20}$ | $1.1 \cdot 10^{21}$ |
| ConvS2S Ensemble [8] | 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** | **41.0** | $2.3 \cdot 10^{19}$ | |

(Vaswani et al., 2017)

# Question

What could be a disadvantage of transformers over RNNs?
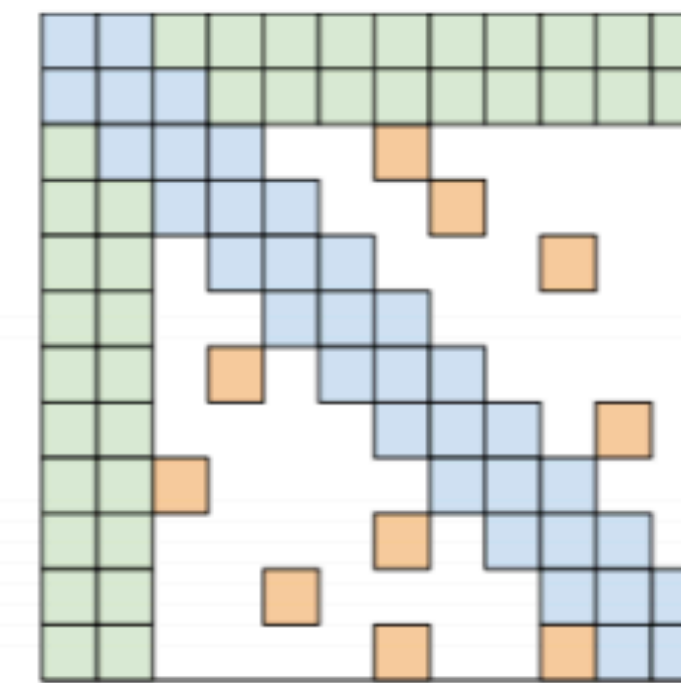


(a) Random attention     (b) Window attention     (c) Global Attention     (d) BIGBIRD

# Other Resources of Interest

- The Annotated Transformer

  - https://nlp.seas.harvard.edu/2018/04/03/attention.html

- The Illustrated Transformer

  - https://jalammar.github.io/illustrated-transformer/

- Only basics presented here today! Many modifications to initial transformers exist

# Recap

- **Temporal Bottleneck**: **Vanishing gradients** stop many RNN architectures from learning **long-range dependencies**

- **Parallelisation Bottleneck:** RNN states depend on previous time step hidden state, so must be **computed in series**

- **Attention**: Direct connections between output states and inputs (solves temporal bottleneck)

- **Self-Attention**: Remove recurrence, allowing parallel computation

- Modern **Transformers** use attention, but require position embeddings to capture sequence order

# Decoding from Neural Models

Antoine Bosselut

# Encoder-Decoder Models

- Encode a sequence fully with one model (**encoder**) and use its representation to seed a second model that decodes another sequence (**decoder**)

- Decoder is *autoregressive*, generates one word at a time (like an LM)

# Decoding: Main Idea

- At each time step $t$, our model computes a vector of scores for each token in our vocabulary, $S \in \mathbb{R}^V$:

$$S = f\left(\{y_{<t}\}\right)$$

> $f(.)$ is your decoder

- Then, we compute a probability distribution $P$ over these scores (with a softmax):

$$P\left(y_t = w \,\middle|\, \{y_{<t}\}\right) = \frac{\exp(S_w)}{\sum_{w' \in V} \exp(S_{w'})}$$

- Decoding algorithm defines a function to select a token from this distribution:

$$\hat{y}_t = g\left(P(y_t | \hat{y}_{<t})\right)$$

> $g(.)$ is your decoding algorithm

# Decoding: Main Idea

- Decoding algorithm defines a function to select a token from this distribution

$$\hat{y}_t = g\left(P(y_t \mid \hat{y}_{<t})\right)$$

# Optional: Encoder Input

- Decoding algorithm defines a function to select a token from this distribution

$$\hat{y}_t = g\big(P(y_t \,|\, X, \hat{y}_{<t})\big)$$

# Greedy methods: Argmax Decoding

$$\hat{y}_t = \underset{w \in V}{\textbf{argmax}}\ P(y_t = w \mid \{y\}_{<t})$$

- $g$ = select the token with the highest probability:

He wanted to go to the → **Decoder** →

restroom
grocery
store
airport
pub
gym
bathroom
game
beach
hospital
doctor

…

# Greedy methods: Argmax Decoding

$$\hat{y}_t = \textbf{argmax } P(y_t = w \,|\, \{y\}_{<t})$$

Select **highest scoring** token

**What's a potential problem with argmax decoding?**

- $g$ = selec

He wanted to go to the → **Decoder** →

store
airport
pub
gym
bathroom
game
beach
hospital
doctor

...

# Issues with argmax decoding

- In argmax decoding, we cannot revise prior decisions

  - *les pauvres sont démunis (the poor don't have any money)*
  - *→ the ____*
  - *→ the poor ____*
  - *→ the poor are ____*

- Potential leads to sequences that are

  - **Ungrammatical**

  - **Unnatural**

  - **Nonsensical**

  - **Incorrect**

  the beam size *k* is usually 5-10

# Beam Search

- *les pauvres sont démunis (the poor don't have any money)*
- *→ the ____*
- *→ the poor ____*
- *→ the poor <span style="color:red">are</span> ____*

- **Beam Search**: Explore several different hypotheses instead of just one

  - Track of the *b* highest scoring sequences at each decoder step instead of just one

  - Score at each step: $\sum_{t=1}^{j} \log P(\hat{y}_t | \hat{y}_1, \ldots, \hat{y}_{t-1}, X)$
    the beam size *k* is usually 5-10

  - *b* is called the **beam size**

# Beam Search

Beam size = 2

$\log \textcolor{magenta}{P}(\hat{y}_1 | y_0)$

the   -1.05

<START>

a   -1.39

# Beam Search

Beam size = 2

$$\sum_{t=1}^{2} \log P(\hat{y}_t | \hat{y}_0, \ldots, \hat{y}_{t-1})$$

```
                        ┌─────────┐
                   ┌───▶│  poor   │   -1.90
        ┌──────┐   │    └─────────┘
        │ the  │───┤    ┌─────────┐
        └──────┘   └───▶│ people  │   -2.3
          ▲             └─────────┘
          │
 ┌─────────┐
 │ <START> │
 └─────────┘
          │             ┌─────────┐
          ▼        ┌───▶│  poor   │   -1.54
        ┌──────┐   │    └─────────┘
        │  a   │───┤    ┌─────────┐
        └──────┘   └───▶│ person  │   -3.2
                        └─────────┘
```

# Beam Search

Beam size = 2

$$\sum_{t=1}^{3} \log P(\hat{y}_t | y_0, \hat{y}_1, \ldots, \hat{y}_{t-1})$$



| | | | are | -2.42 |
| the | poor | | don't | -2.13 |
| | people | | | |

<START>

| | poor | | person | -3.12 |
| a | | | but | -3.53 |
| | person | | | |

# Beam Search

Beam size = 2

```
                                          ┌──────────┐
                                          │  always  │  -3.82
                                          └──────────┘
                                          ┌──────────┐
                               ┌──────────┤   not    │  -2.67
                               │   are    ├──────────┘
                    ┌──────────┤          │
          ┌─────────┤   poor   └──────────┤  ┌──────────┐
          │   the   ├──────────┐          │  │   have   │  -3.32
<START>───┤         │  people  │  don't   ├──┴──────────┘
          │         └──────────┘          │  ┌──────────┐
          │                               └──┤   take   │  -3.61
          │         ┌──────────┐             └──────────┘
          │         │   poor   ├──────────┐  ┌──────────┐
          └─────────┤          │          │  │  person  │
              a     ├──────────┤  poor     ──┤──────────┘
                    │  person  │          │  │   but    │
                    └──────────┘          └──┴──────────┘
                                          person
```

and so on…

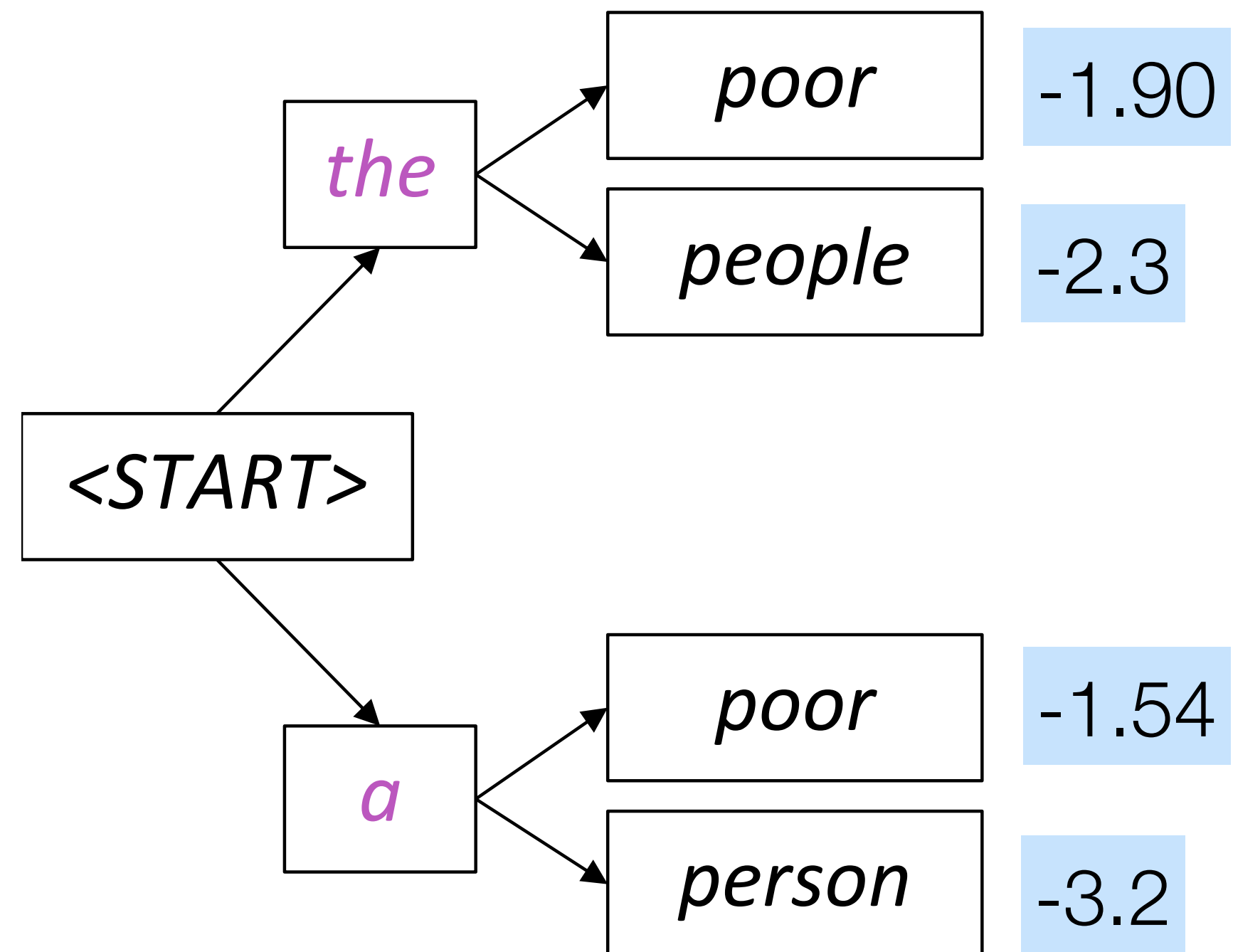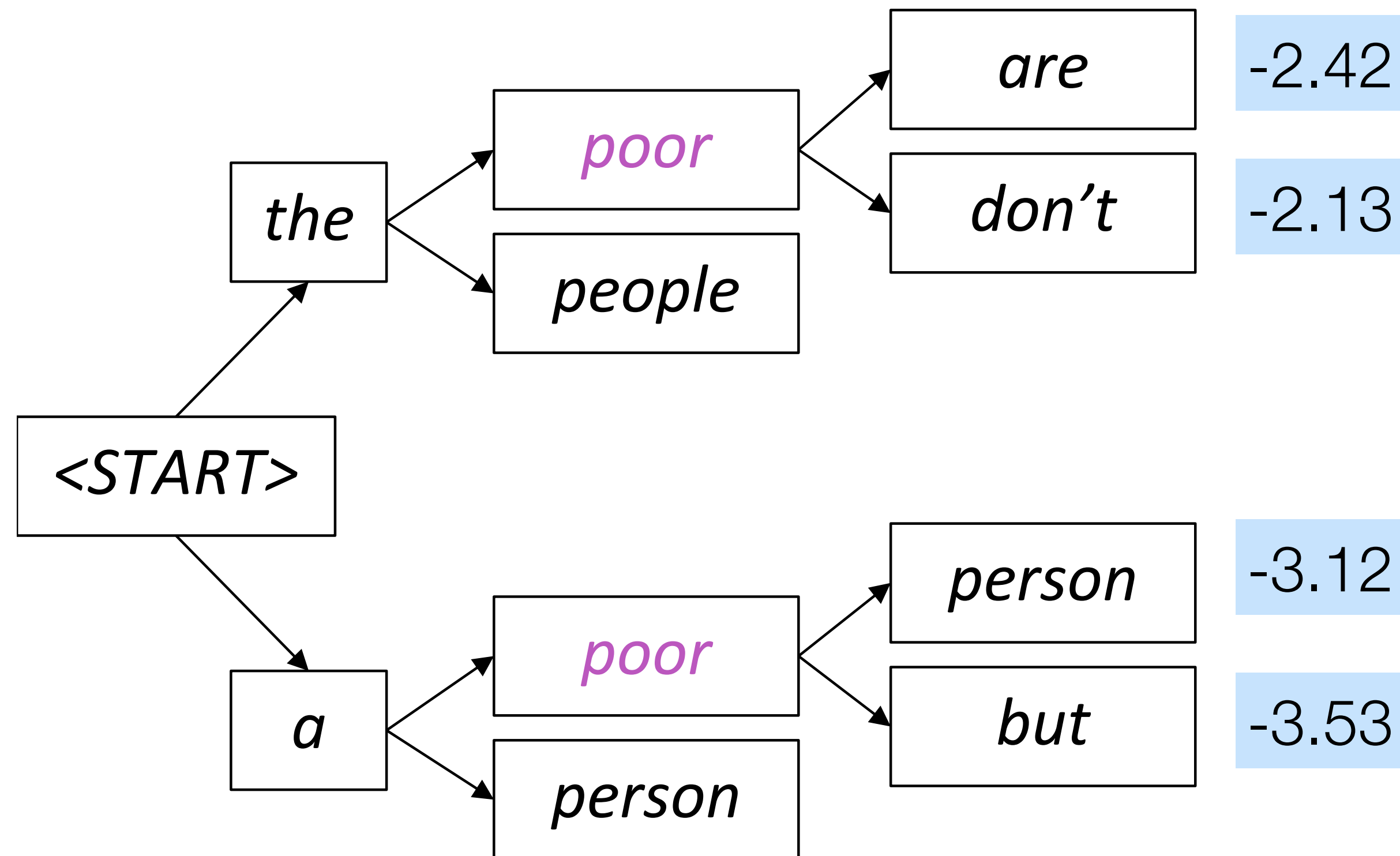$$\sum_{t=1}^{j} \log P(\hat{y}_t \mid \hat{y}_1, \ldots, \hat{y}_{t-1})$$
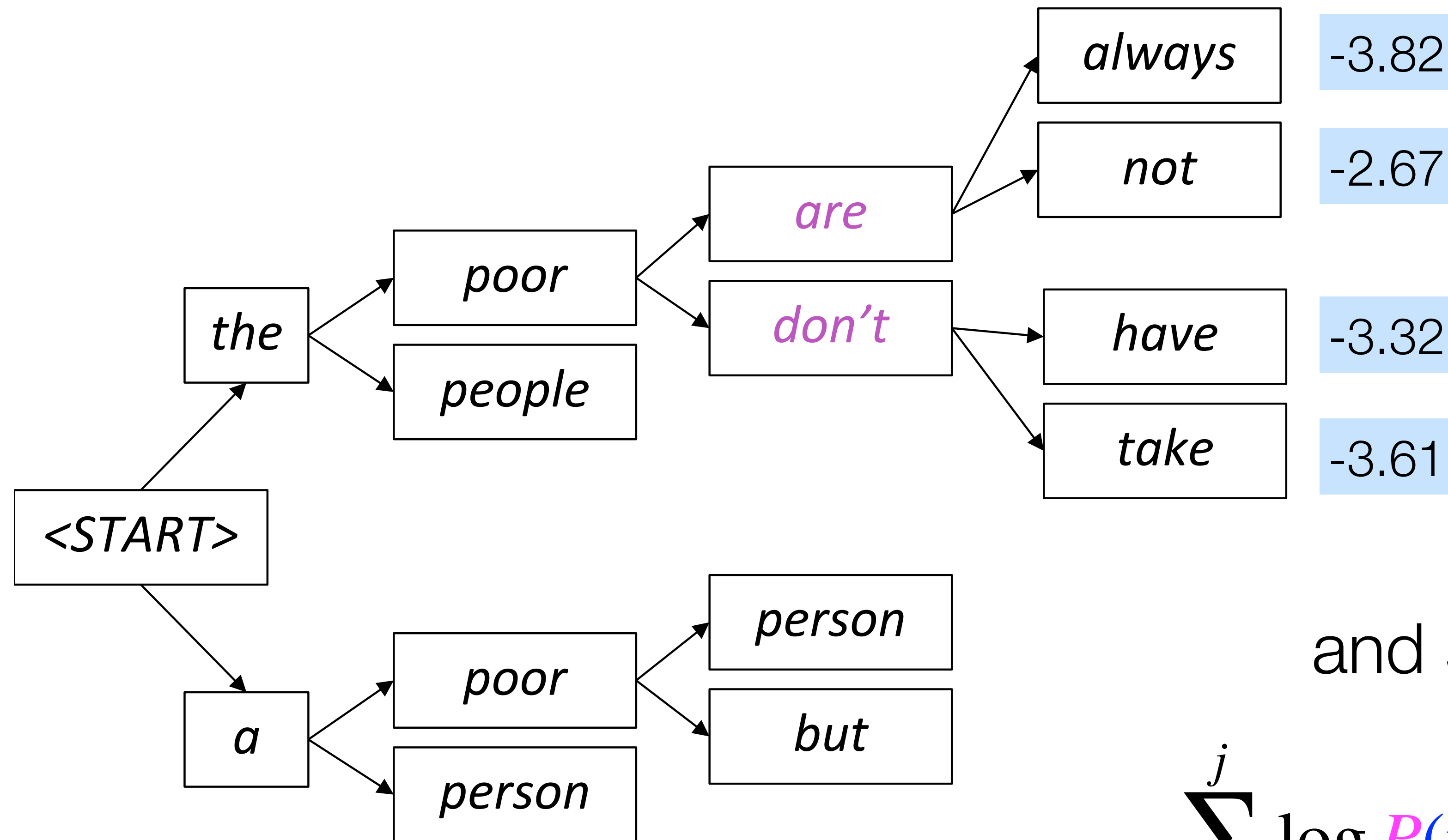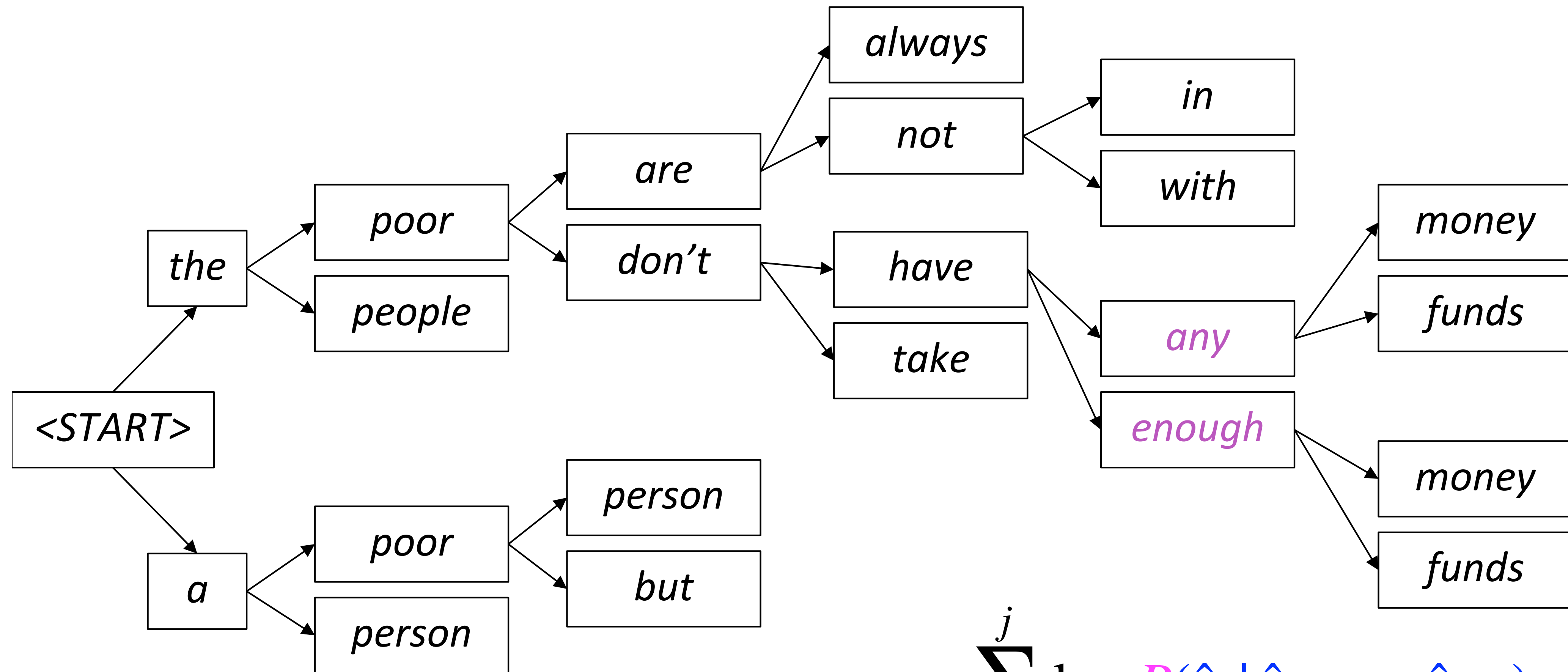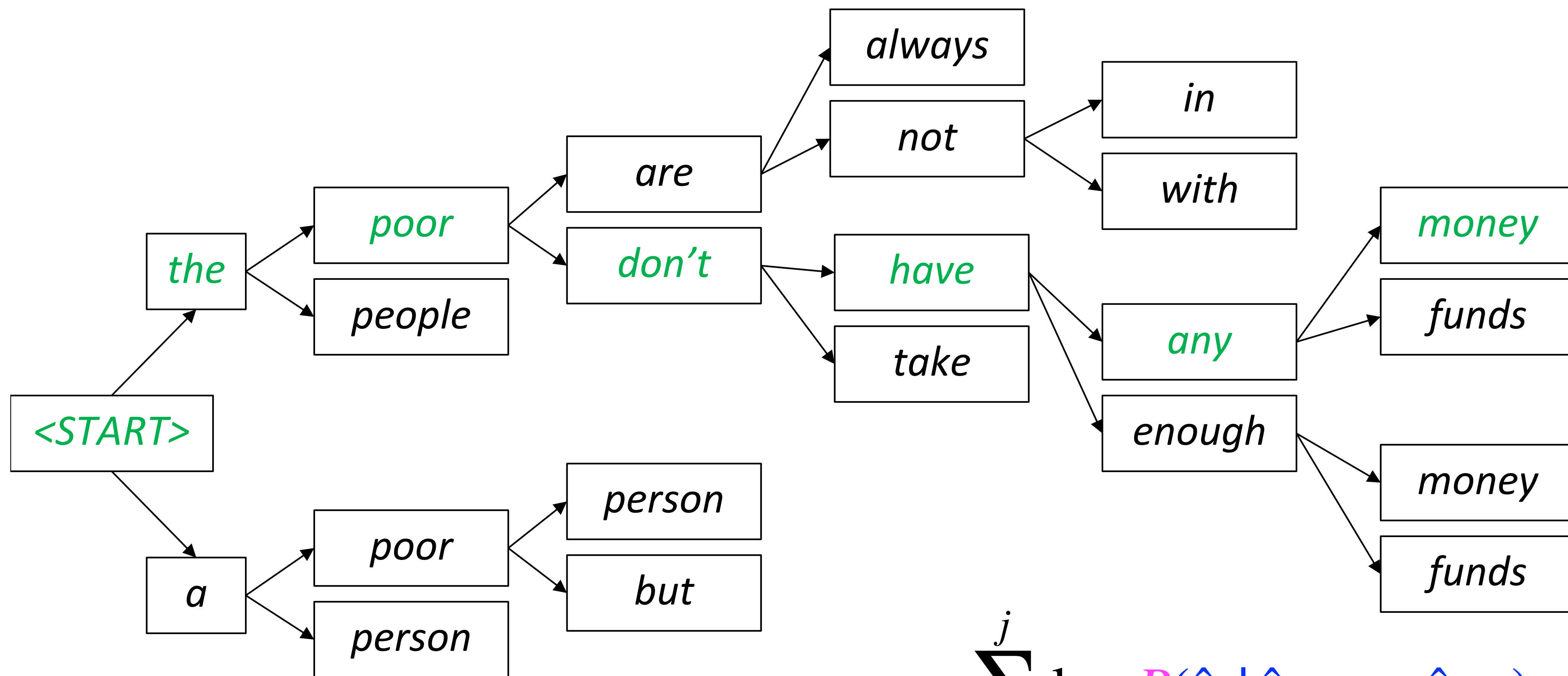
# Beam Search

Beam size = 2



$$\sum_{t=1}^{j} \log P(\hat{y}_t \mid \hat{y}_1, \ldots, \hat{y}_{t-1})$$

# Beam Search

Beam size = 2



$$\sum_{t=1}^{j} \log P(\hat{y}_t \mid \hat{y}_1, \ldots, \hat{y}_{t-1})$$

# Beam Search

- Different hypotheses may produce <END> token at different time steps

  - When a hypothesis produces <END>, stop expanding it and place it aside

- Continue beam search until:

  - All $b$ beams (hypotheses) produce <END> OR

  - Hit max decoding limit T

- Select top hypotheses using the *normalized* likelihood score

$$\frac{1}{T} \sum_{t=1}^{T} \log P(\hat{y}_t \,|\, \hat{y}_1, \ldots, \hat{y}_{t-1}, X)$$

  - Otherwise shorter hypotheses have higher scores

# What do you think might happen if we increase the beam size?

# Effect of beam size

- Small beam size *b* has similar issues as argmax decoding

  - **Outputs that are ungrammatical, unnatural, nonsensical, incorrect**

  - b=1 is the same as argmax decoding

- Larger beam size *b* reduces some of these problems

  - Potentially much more computationally expensive

  - Outputs tend to get shorter and more generic

# Looking Forward

- **Tomorrow**: Pretraining Transformers - GPT

- **Next week:** Pretraining masked language models (BERT), Transfer Learning

- **Exercise Session:** Transformers, Decoding

# References

- Bahdanau, D., Cho, K., & Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *CoRR, abs/1409.0473.*

- Vaswani, A., Shazeer, N.M., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., & Polosukhin, I. (2017). Attention is All you Need. *ArXiv, abs/1706.03762.*

- Wu et al., Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. arxiv 2016