

# 10-701: Introduction to Machine Learning

## Lecture 13 – Attention & Transformers

Henry Chai

2/28/24

# Front Matter

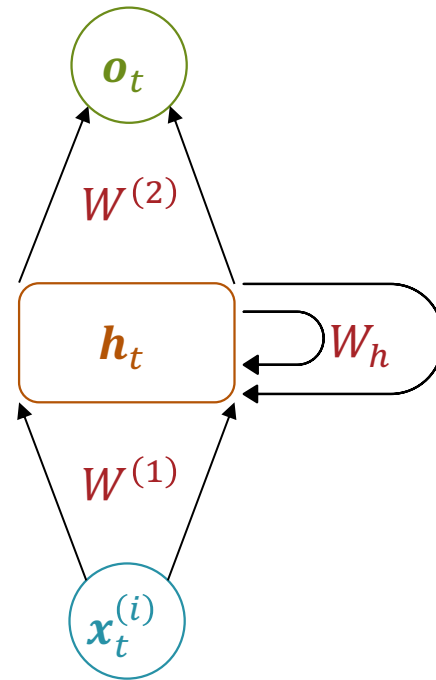
- Announcements
  - HW3 released 2/19, due 2/28 (today!) at 11:59 PM
  - HW4 released 2/28 (today!), due 3/15 (after break) at 11:59 PM
  - Project details will be released 3/1 (Friday)
    - **You must work in groups of 2 or 3 on the project**
- Recommended Readings
  - Zhang, Lipton, Li & Smola, [Chapters 9 & 10](#)

# Recurrent Neural Networks

- Neural networks are frequently applied to inputs with some inherent temporal or sequential structure (e.g., text or video) of variable length
- Idea: use the information from previous parts of the input to inform subsequent predictions
- Insight: the hidden layers learn a useful representation (relative to the task)
- Approach: incorporate the output from earlier hidden layers into later ones.

# Recurrent Neural Networks

$$\mathbf{h}_t = \left[ 1, \theta \left( W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left( W^{(2)} \mathbf{h}_t \right)$$



- Training dataset consists of (input **sequence**, label **sequence**) pairs, potentially of varying lengths

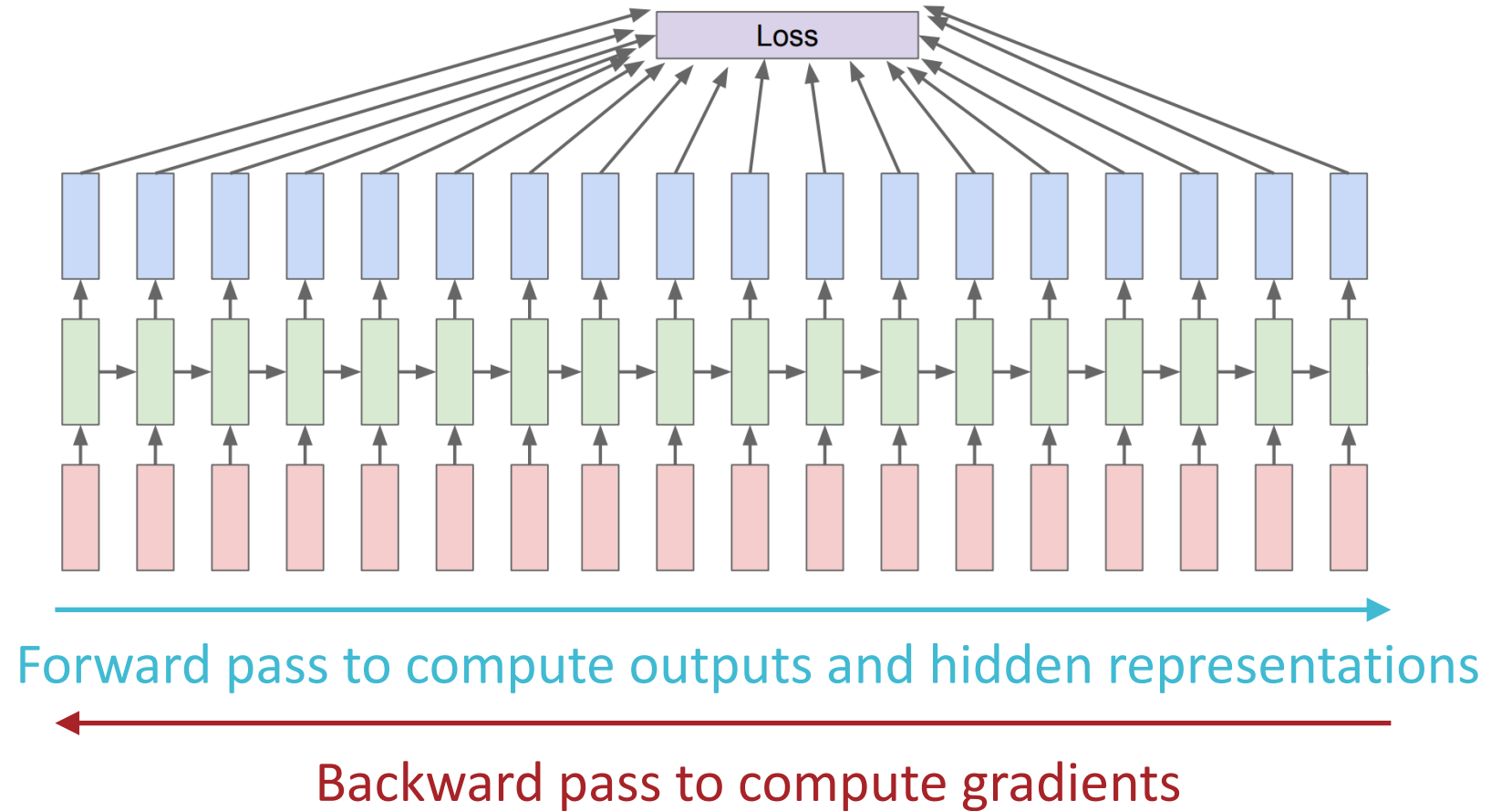
$$\mathcal{D} = \left\{ \left( \mathbf{x}^{(n)}, \mathbf{y}^{(n)} \right) \right\}_{n=1}^N$$

$$\mathbf{x}^{(n)} = \left[ \mathbf{x}_1^{(n)}, \dots, \mathbf{x}_{T_n}^{(n)} \right]$$

$$\mathbf{y}^{(n)} = \left[ \mathbf{y}_1^{(n)}, \dots, \mathbf{y}_{T_n}^{(n)} \right]$$

- This model requires an initial value for the hidden representation,  $\mathbf{h}_0$ , typically a vector of all zeros

# Training RNNs: Challenges



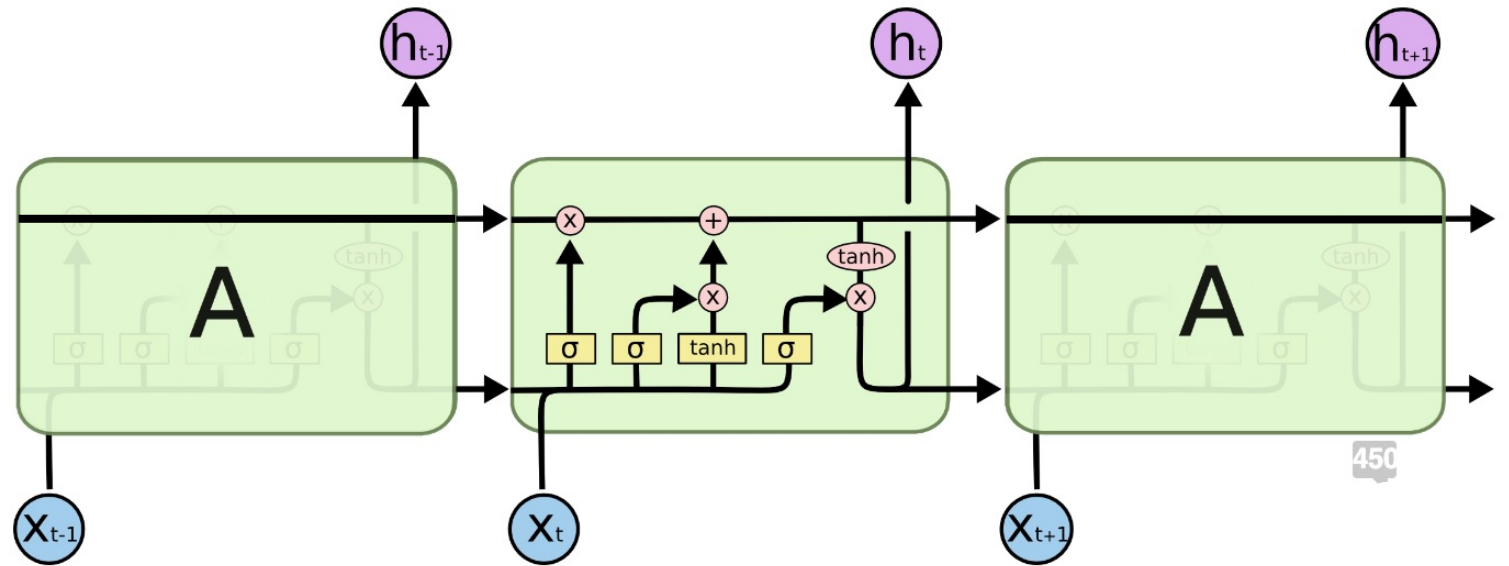
- Issue: as the sequence length grows, the gradient is more likely to explode or vanish

# Long Short-Term Memory (Hochreiter & Schmidhuber, 1997)

- LSTM networks address the vanishing gradient problem by replacing hidden layers with *memory cells*
- Each cell still computes a hidden representation but also maintains a separate internal *state*,  $C_t$
- The flow of information through a cell is manipulated by three *gates*:
  - An input gate,  $I_t$ , that controls how much the state looks like the normal RNN hidden layer
  - An output gate,  $O_t$ , that “releases” the hidden representation to later timesteps
  - A forget gate,  $F_t$ , that determines if the previous memory cell’s state affects the current internal state

# Long Short-Term Memory (Hochreiter & Schmidhuber, 1997)

- LSTM networks address the vanishing gradient problem by replacing hidden layers with *memory cells*
- Each cell still computes a hidden representation but also maintains a separate internal state,  $C_t$



- The internal state allows information to move through time without needing to affect the hidden representations!

# Applications of LSTMs



**2018:** [OpenAI](#) used LSTM trained by policy gradients to beat humans in the complex video game of Dota 2,<sup>[11]</sup> and to control a human-like robot hand that manipulates physical objects with unprecedented dexterity.<sup>[10][54]</sup>

**2019:** [DeepMind](#) used LSTM trained by policy gradients to excel at the complex video game of [Starcraft II](#).<sup>[12][54]</sup>



# Key Takeaways

- Recurrent neural networks use contextual information to reason about sequential data
  - Can still be learned using backpropagation → backpropagation through time
  - Susceptible to exploding/vanishing gradients for long training sequences
  - LSTMs allow contextual information to reach later timesteps without directly affecting intermediate hidden representations

# Language Models

1. Convert raw text into *embeddings*

$$\mathbf{x}^{(i)} = \left[ \underline{\mathbf{x}_1^{(i)}}, \dots, \underline{\mathbf{x}_{T_i}^{(i)}} \right]$$

2. Learn or approximate a joint probability distribution over sequences

$$P(\mathbf{x}^{(i)}) = P\left(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}\right)$$

3. Sample from the implied conditional distribution to generate new sequences

$$P\left(\mathbf{x}_{T_i+1} \mid \mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}\right) = \frac{P\left(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}, \mathbf{x}_{T_i+1}\right)}{P\left(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}\right)}$$

# Language Models

1. Convert raw text into *embeddings*

$$\mathbf{x}^{(i)} = [\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}]$$

2. Learn or approximate a joint probability distribution over sequences

$$P(\mathbf{x}^{(i)}) = P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)})$$

- Use the chain rule of probability: predict the next word based on the previous words in the sequence

$$\begin{aligned} P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}) &= P(\mathbf{x}_1^{(i)}) \\ &\cdot P(\mathbf{x}_2^{(i)} | \mathbf{x}_1^{(i)}) \\ &\cdot P(\mathbf{x}_3^{(i)} | \mathbf{x}_1^{(i)}, \mathbf{x}_2^{(i)}) \\ &\cdot P(\mathbf{x}_{T_i}^{(i)} | \mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i-1}^{(i)}) \end{aligned}$$

# Language Models

1. Convert raw text into *embeddings*

$$\mathbf{x}^{(i)} = [\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}]$$

2. Learn or approximate a joint probability distribution over sequences

$$P(\mathbf{x}^{(i)}) = P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)})$$

- ~~Use the chain rule of probability~~ Just throw an RNN at it!

$$\begin{aligned} P(\mathbf{x}^{(i)}) &= P(\mathbf{x}_1^{(i)}) \\ &\quad * P(\mathbf{x}_2^{(i)} \mid \mathbf{x}_1^{(i)}) \\ &\quad \vdots \\ &\quad * P(\mathbf{x}_{T_i}^{(i)} \mid \mathbf{x}_{T_i-1}^{(i)}, \dots, \mathbf{x}_1^{(i)}) \end{aligned}$$

# RNN Language Models

1. Convert raw text into *embeddings*

$$\mathbf{x}^{(i)} = [\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)}]$$

2. Learn or approximate a joint probability distribution over sequences

$$P(\mathbf{x}^{(i)}) = P(\mathbf{x}_1^{(i)}, \dots, \mathbf{x}_{T_i}^{(i)})$$

- ~~Use the chain rule of probability~~ Just throw an RNN at it!

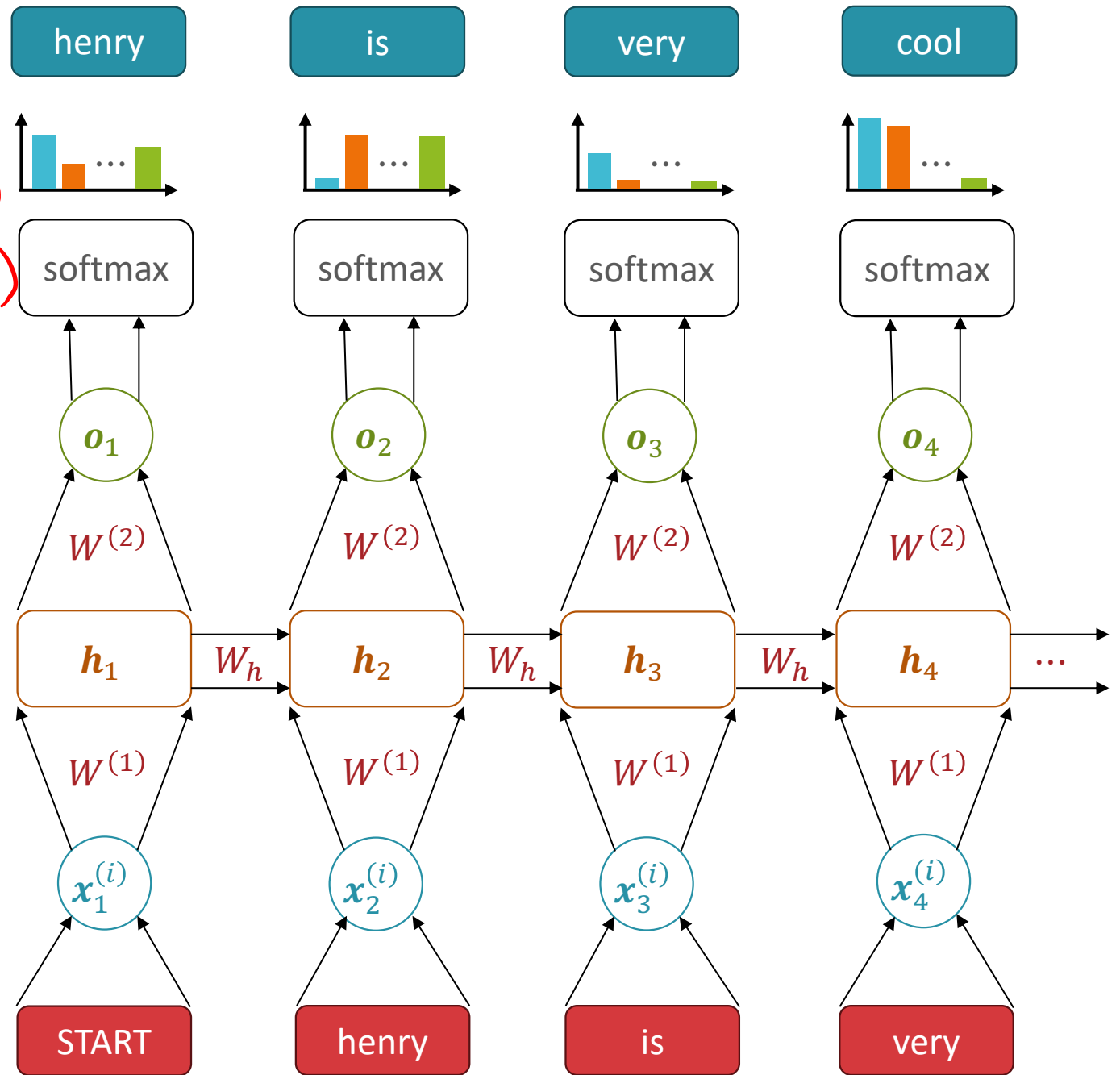
$$\begin{aligned} P(\mathbf{x}^{(i)}) &\approx \mathbf{o}_1(\mathbf{x}_1^{(i)}) \\ &* \mathbf{o}_2(\mathbf{x}_2^{(i)}, \mathbf{h}_1(\mathbf{x}_1^{(i)})) \\ &\vdots \\ &* \mathbf{o}_{T_i}(\mathbf{x}_{T_i}^{(i)}, \mathbf{h}_{T_i-1}(\mathbf{x}_{T_i-1}^{(i)}, \dots, \mathbf{x}_1^{(i)})) \end{aligned}$$

# RNN Language Models: Training

Target sequence (try to predict the next word)

$$\frac{\exp(o_i[a])}{\sum_b \exp(o_i[b])}$$

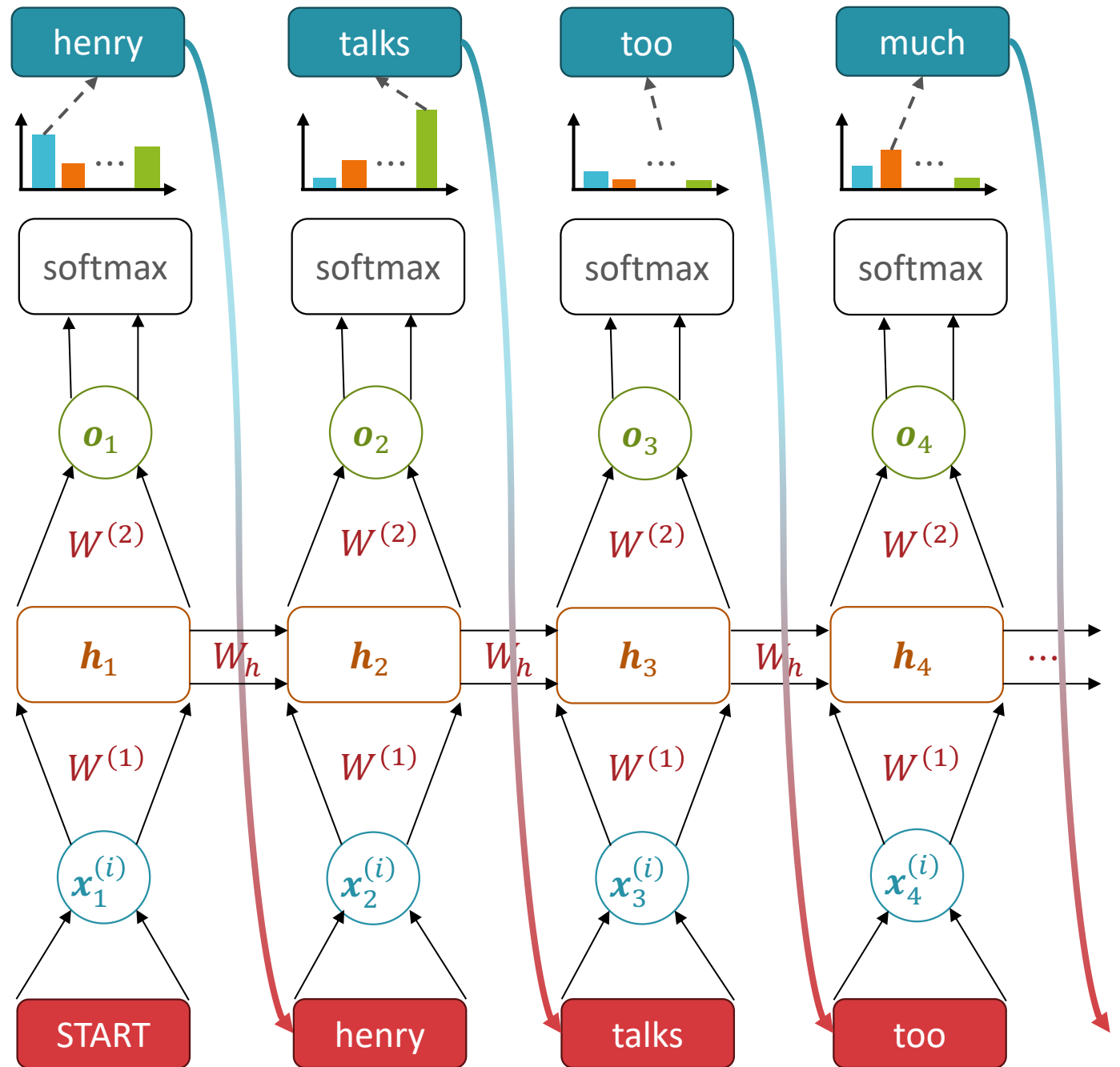
Input sequence



Generated sequence (use each token as the input to the next timestep)

# RNN Language Models: Sampling

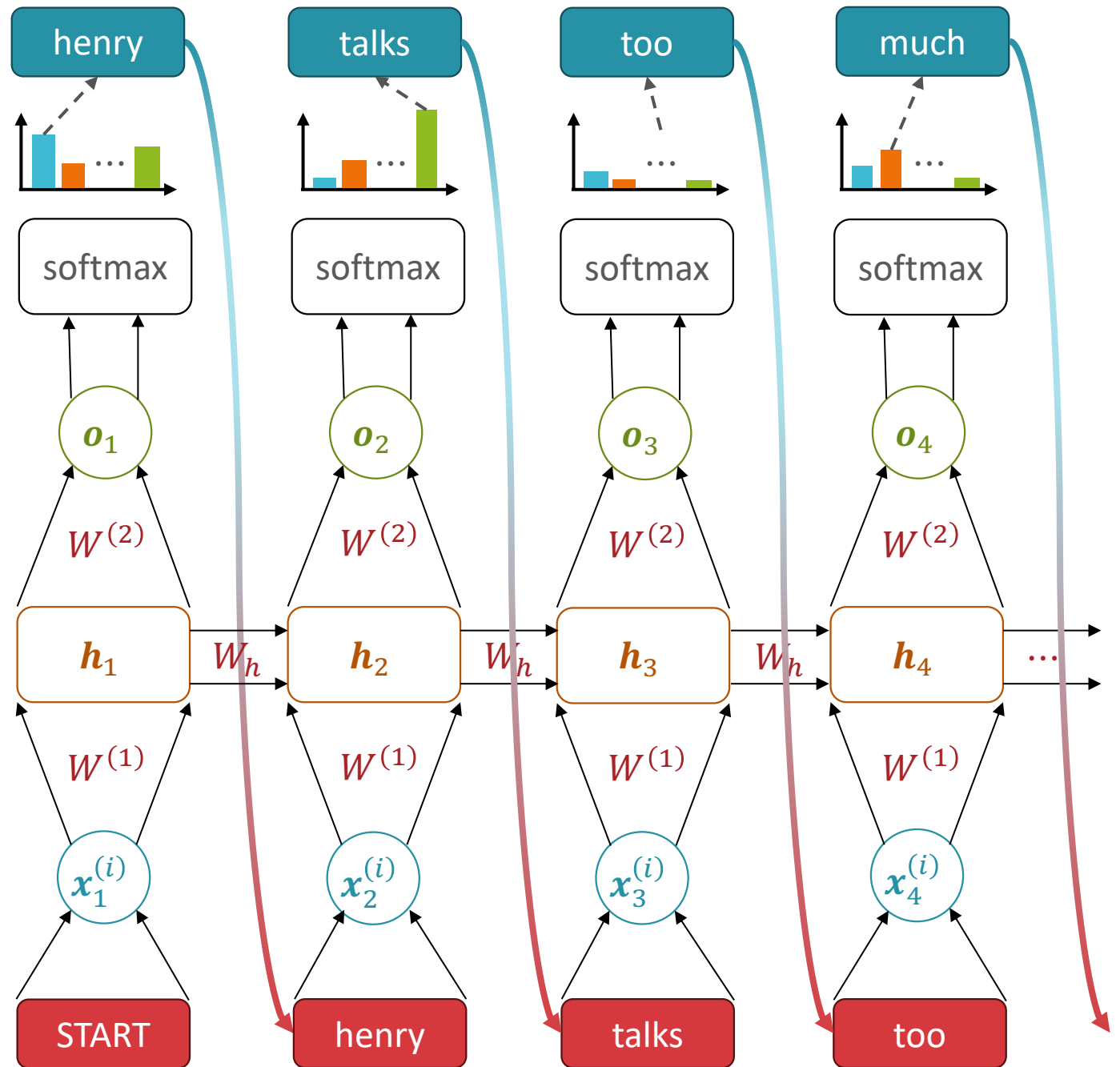
Input sequence



Generated sequence (use each token as the input to the next timestep)

Aside:  
Sampling from these distributions to get the next word is not always the best thing to do

Input sequence



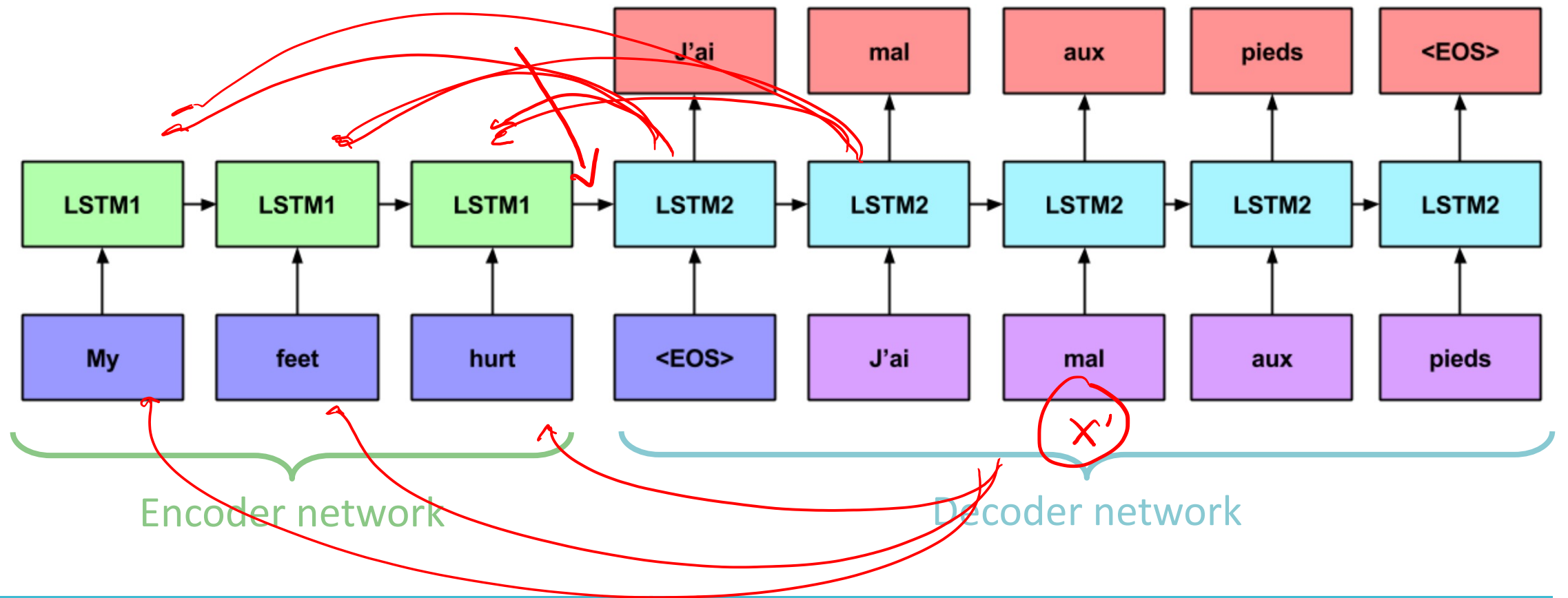


# RNN Language Models: Pros & Cons

- Pros:
  - Can handle arbitrary sequence lengths without having to increase model size (i.e., # of learnable parameters)
  - Trainable via backpropagation
- Cons
  - Vanishing/exploding gradients
  - Does not consider information from later timesteps
    - Can be addressed by bidirectional RNNs
  - Computation is inherently sequential
  - "You can't cram the meaning of a whole %&!\$# sentence into a single \$&!#\* vector!" – Ray Mooney, UT Austin

# RNN Language Models: Pros & Cons

- Pros:
  - Can handle arbitrary sequence lengths without having to increase model size (i.e., # of learnable parameters)
  - Trainable via backpropagation
- Cons
  - Vanishing/exploding gradients
  - Does not consider information from later timesteps
    - Can be addressed by bidirectional RNNs
  - Computation is inherently sequential
  - The entire sequence up to some timestep is represented using just one vector (or two vectors in an LSTM)



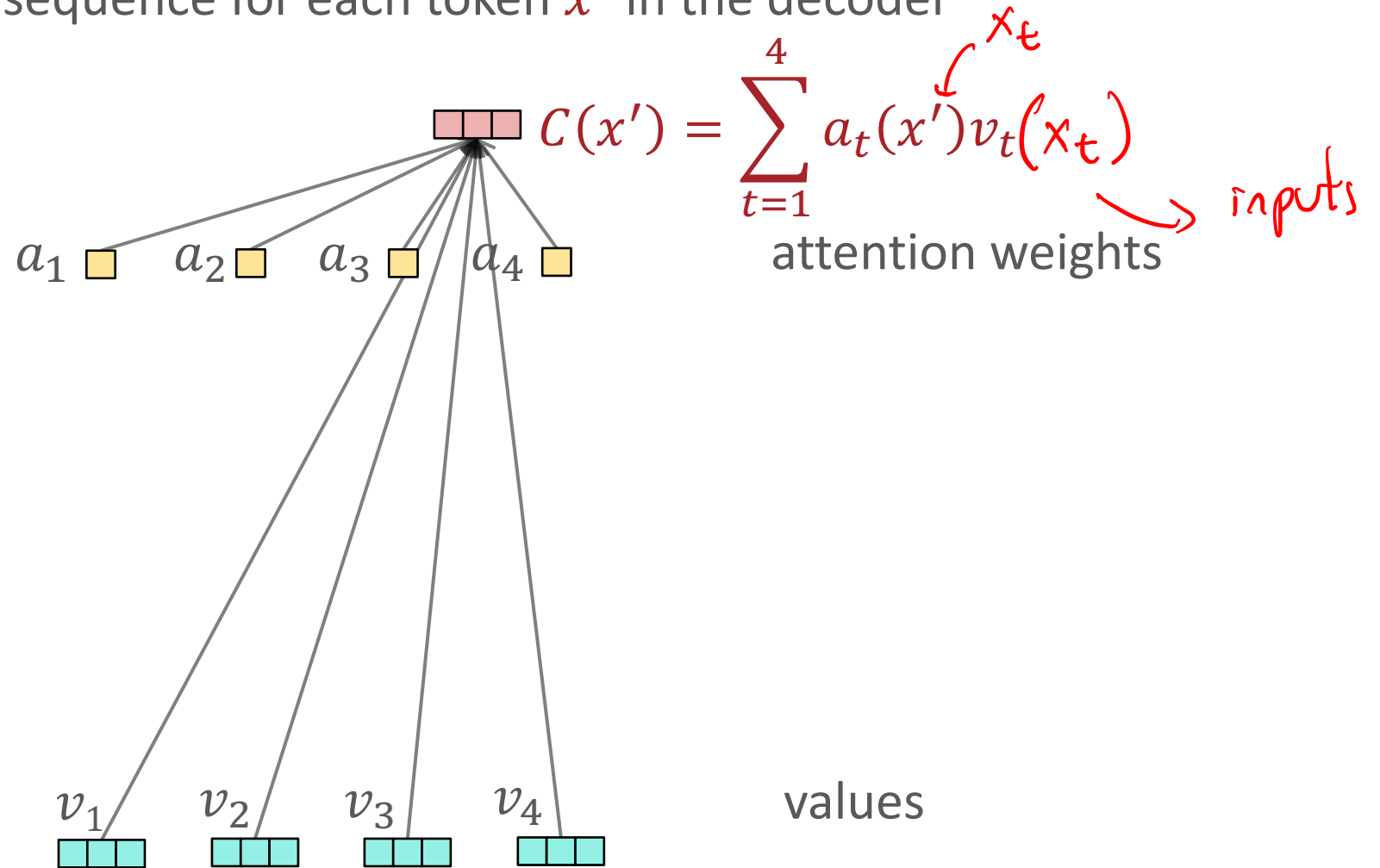
# Encoder-Decoder Architectures (Sutskever et al., 2014)

# Attention

- Approach: compute a representation of the input sequence for each token  $x'$  in the decoder
- Idea: allow the decoder to learn which tokens in the input to “pay attention to” i.e., put more weight on

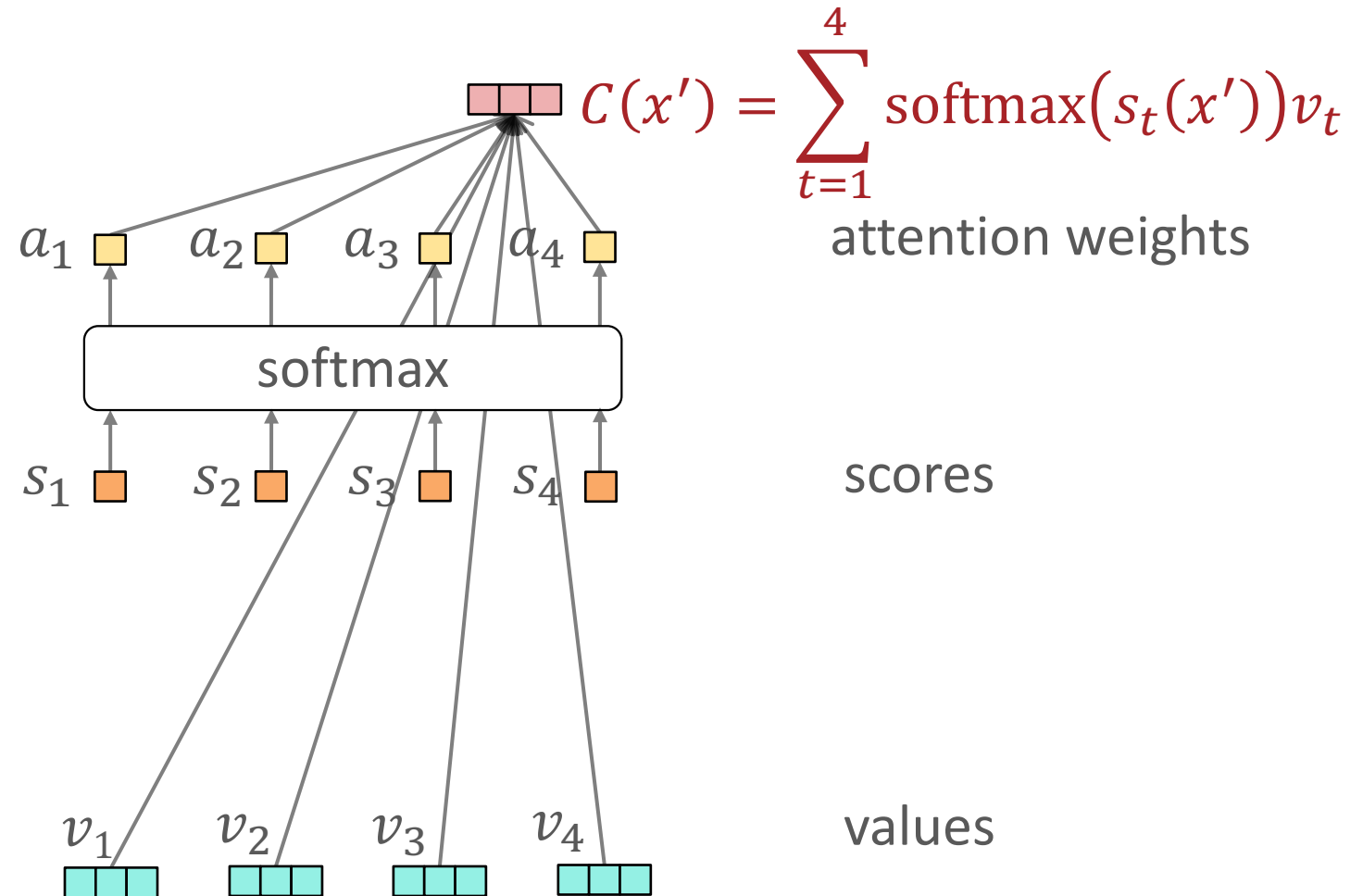
# Attention

- Approach: compute a representation of the input sequence for each token  $x'$  in the decoder



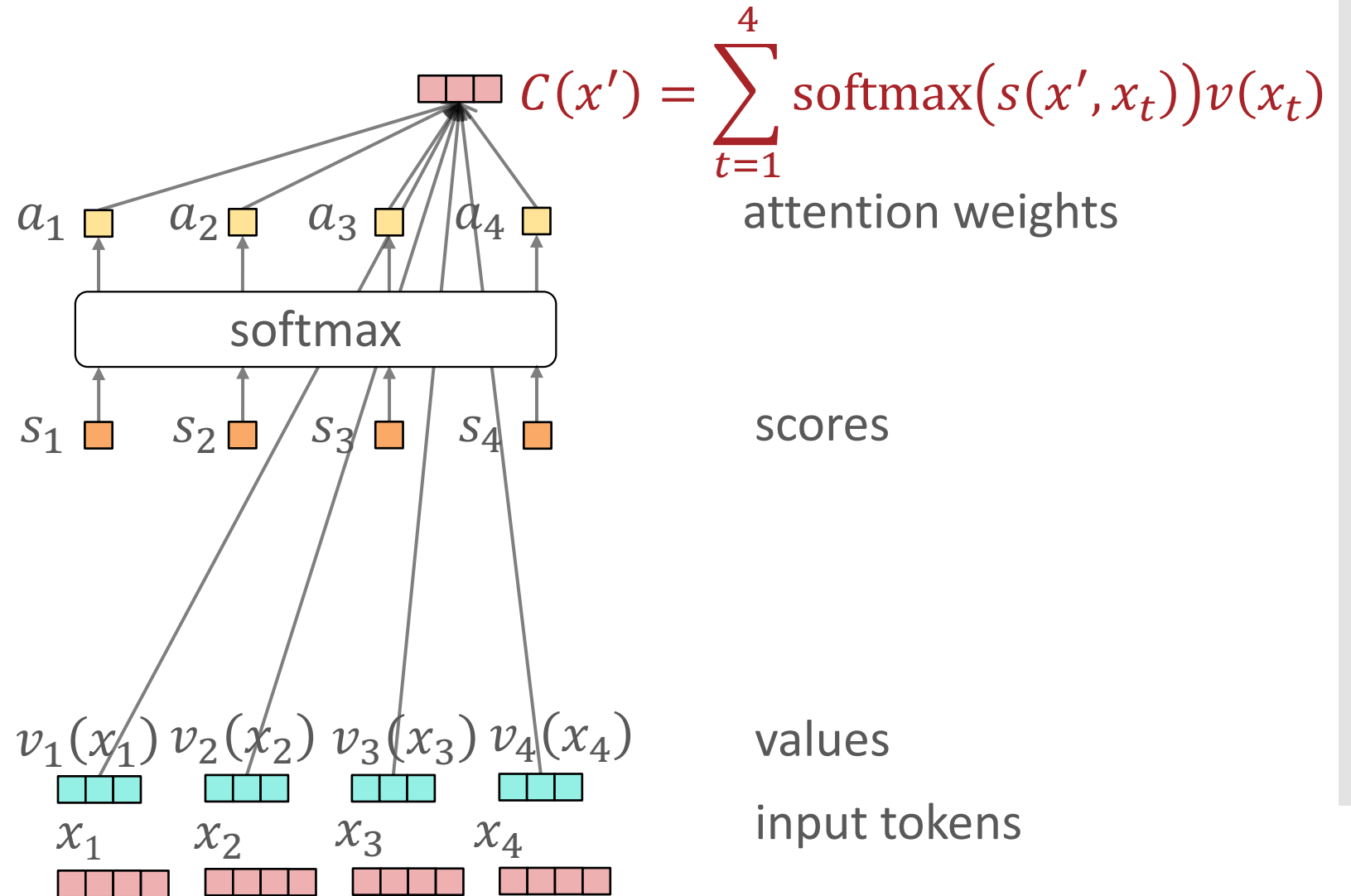
# Attention

- Approach: compute a representation of the input sequence for each token  $x'$  in the decoder



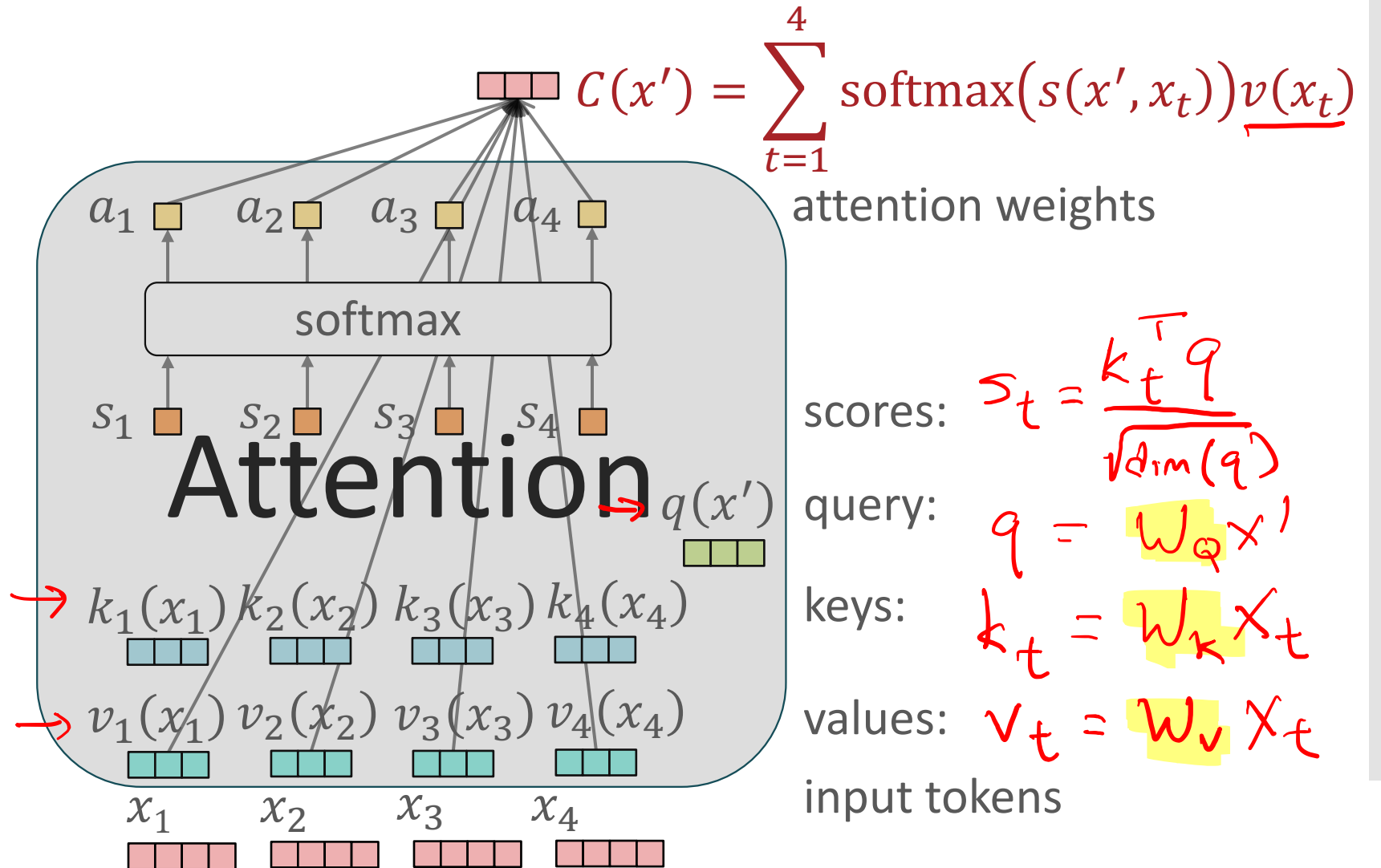
# Attention

- Approach: compute a representation of the input sequence for each token  $x'$  in the decoder

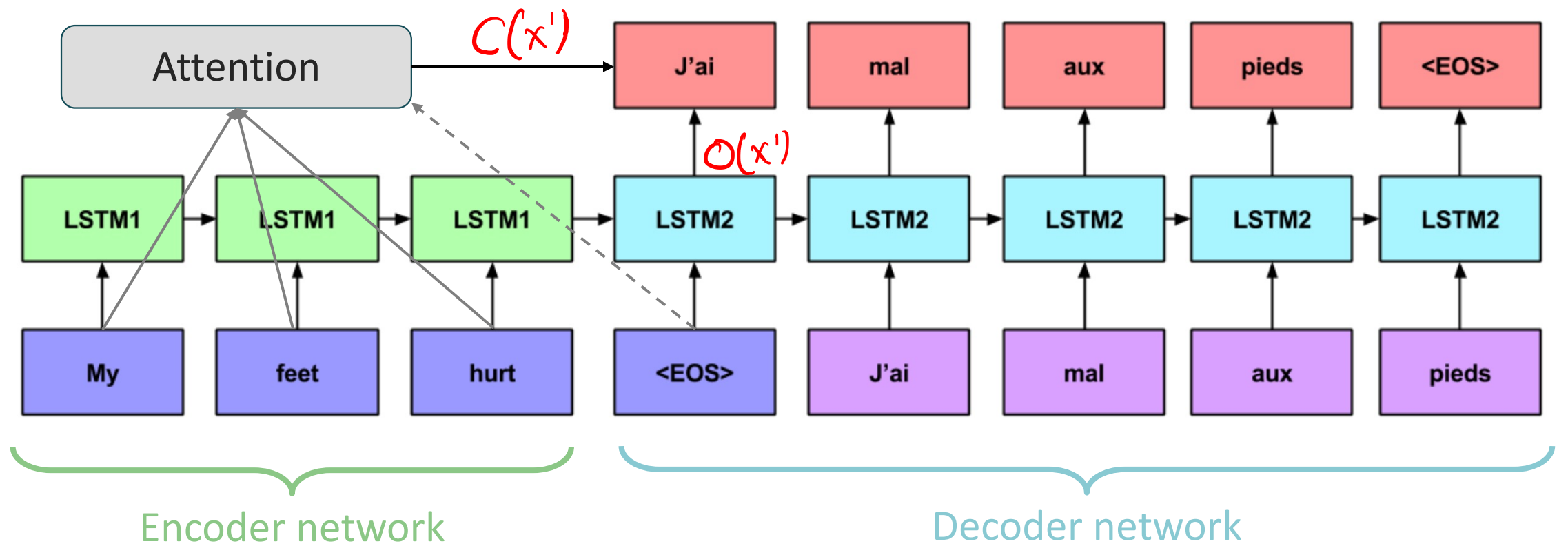


→ Scaled  
Dot-product  
Attention

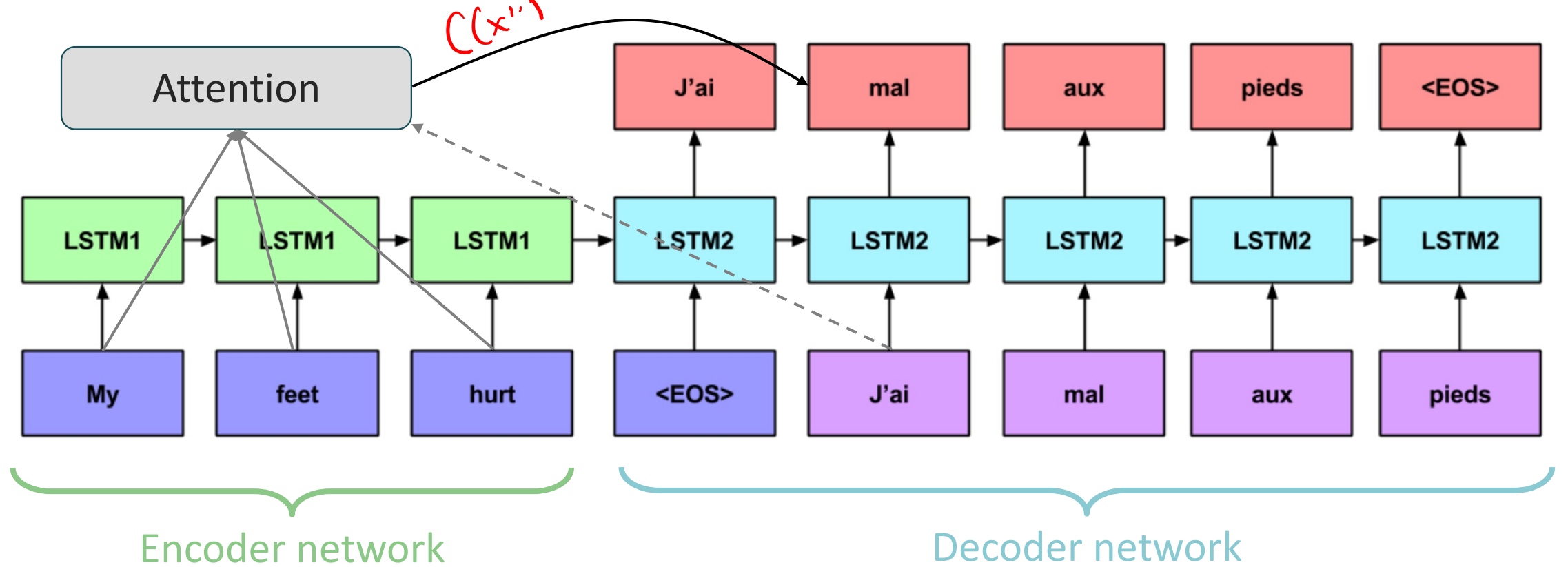
- Approach: compute a representation of the input sequence for each token  $x'$  in the decoder



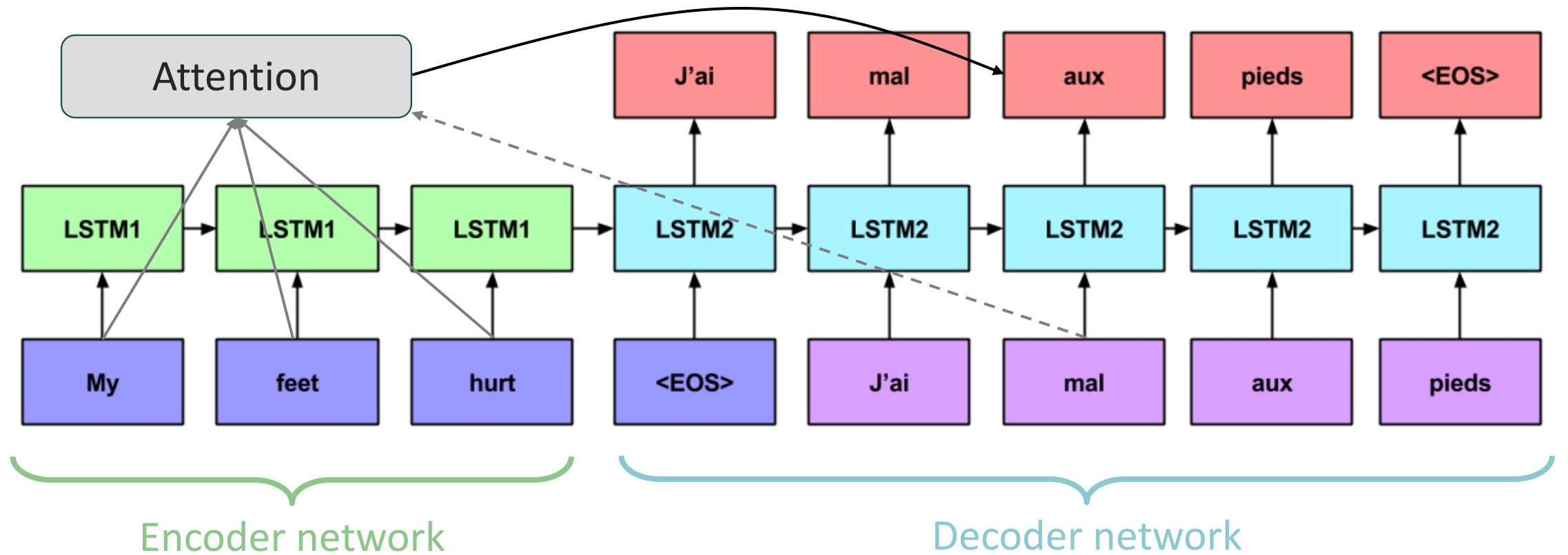




# Encoder-Decoder Architectures with Attention

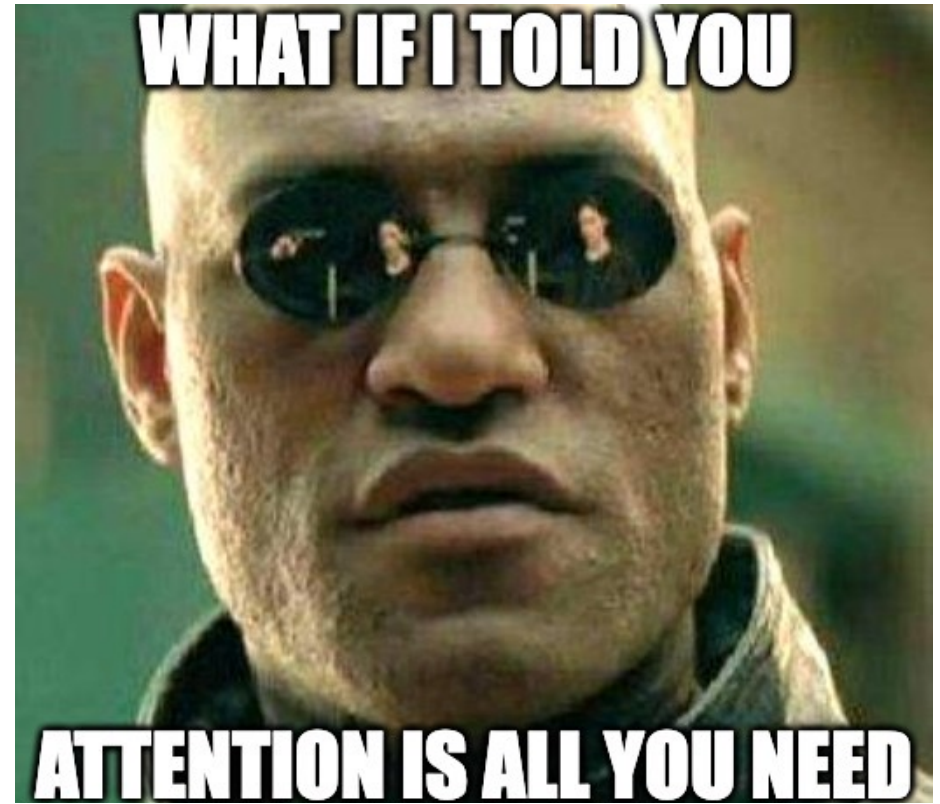


# Encoder-Decoder Architectures with Attention



# Encoder-Decoder Architectures with Attention

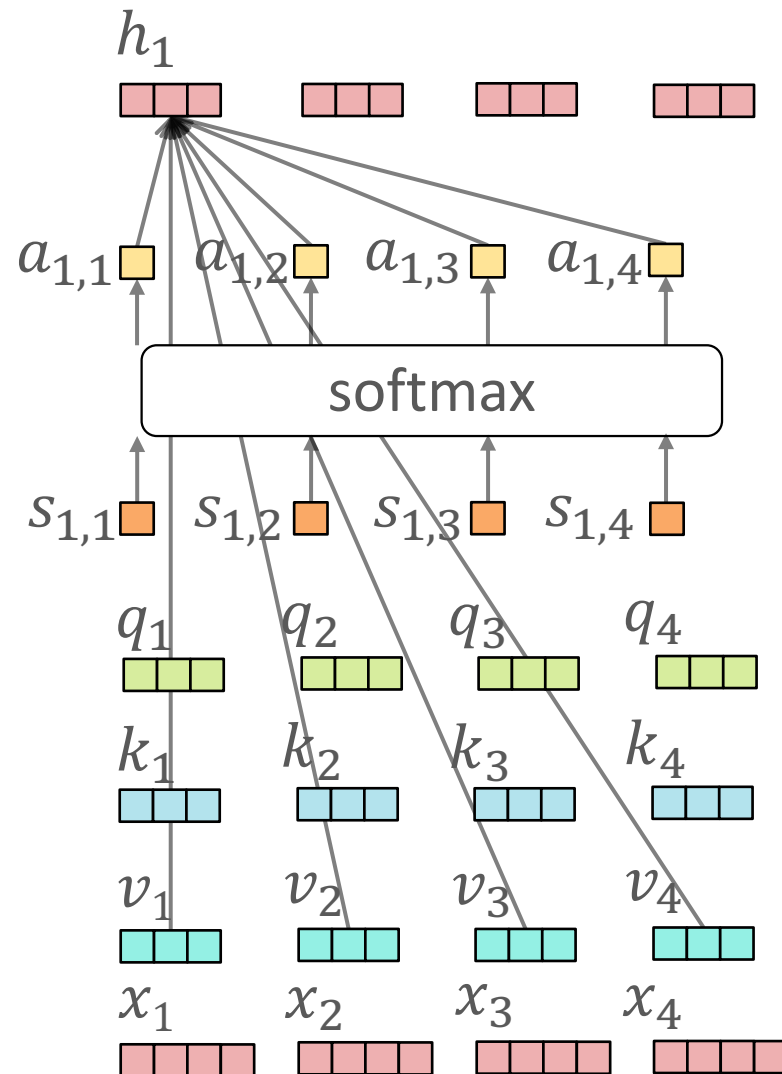
Attention



# ~~Encoder-Decoder Architectures~~ with Attention (Vaswani et al., 2017)

# Scaled Dot-product Self-attention

- Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$h_1 = \sum_{j=1}^4 \text{softmax}(s_{1,j}) v_j$$

attention weights

scores:

queries:  $q_t = W_Q x_t$

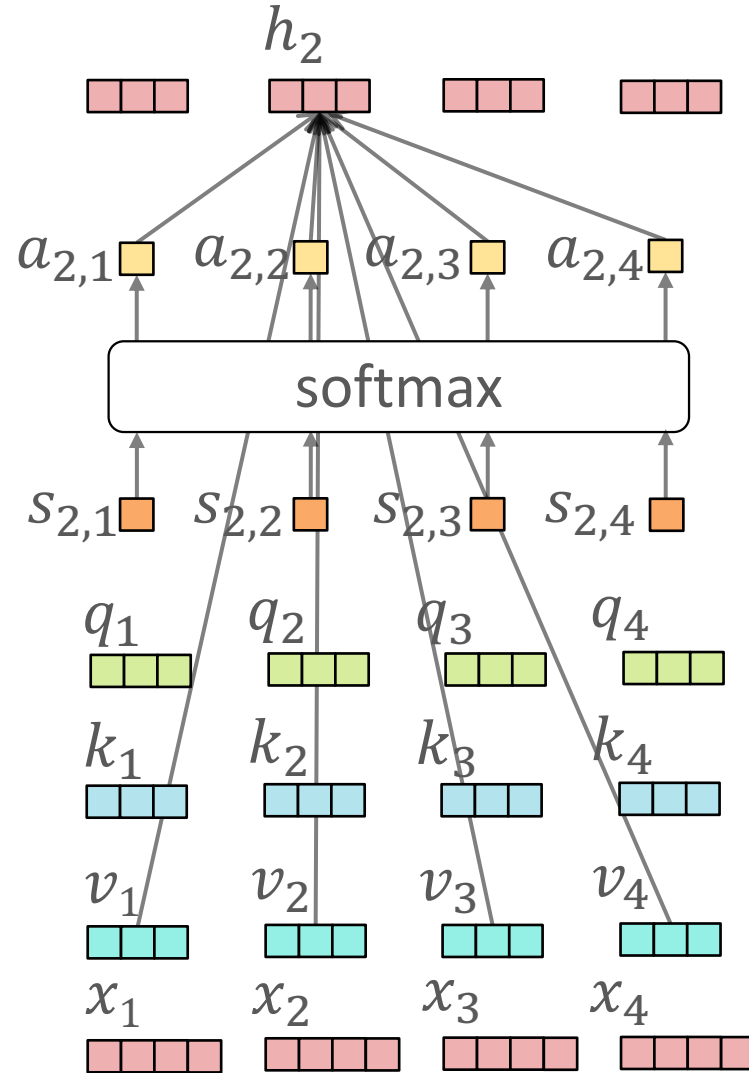
keys:  $k_t = W_K x_t$

values:  $v_t = W_V x_t$

input tokens

# Scaled Dot-product Self-attention

- Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$h_2 = \sum_{j=1}^4 \text{softmax}(s_{2,j}) v_j$$

attention weights

scores:

queries:  $q_t = W_Q x_t$

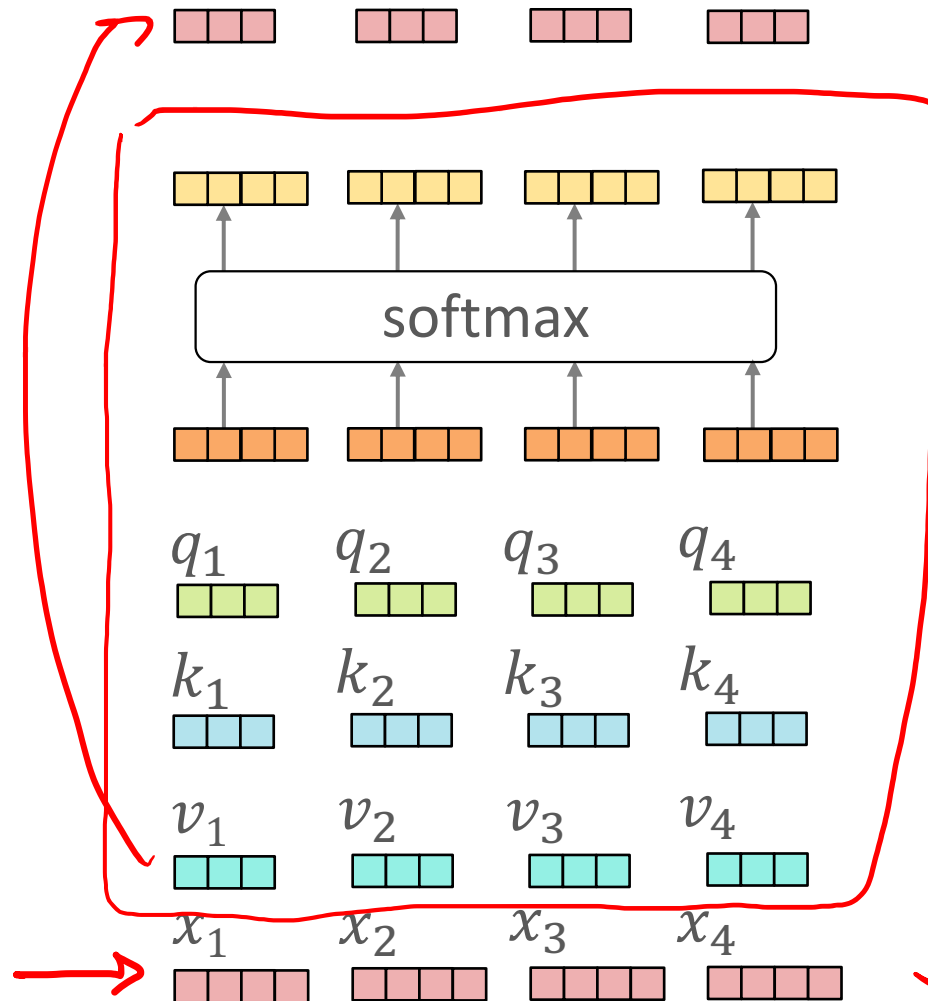
keys:  $k_t = W_K x_t$

values:  $v_t = W_V x_t$

input tokens

# Scaled Dot-product Self-attention: Matrix Form

- Approach: compute a representation for each token in the *input sequence* by attending to all the input tokens



$$H = \text{softmax}(S)V$$

attention weights  
 $\uparrow$  applied row-wise to  $S$   
 $\mathbb{R}^{N \times d_k}$   $\mathbb{R}^{d_k \times N}$   $\in \mathbb{R}^{N \times N}$

scores:  $S = \frac{QK^T}{\sqrt{d_k}}$   
 queries:  $Q = XW_Q \in \mathbb{R}^{N \times d_k}$

keys:  $K = XW_K \in \mathbb{R}^{N \times d_k}$

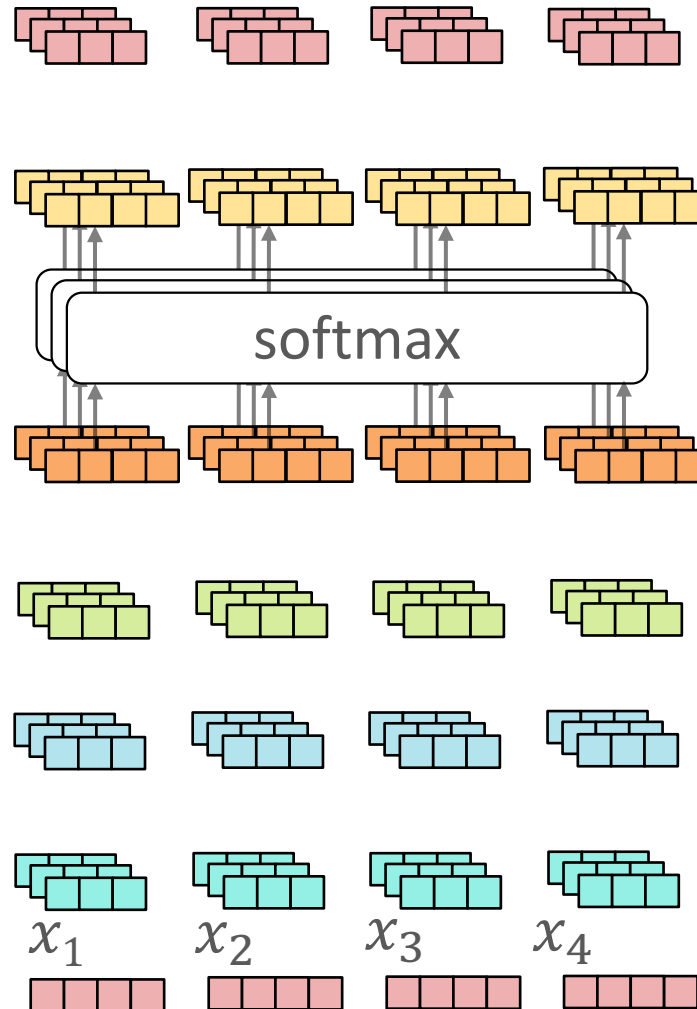
values:  $V = XW_V \in \mathbb{R}^{N \times d_v}$

design matrix:  $X \in \mathbb{R}^{N \times D}$

$W_k \in \mathbb{R}^{D \times d_k}$

# Multi-head Scaled Dot-product Self-attention

- Idea: just like we might want multiple convolutional filters in a convolutional layer, we might want multiple attention weights to learn different relationships between tokens!



$$H^{(h)} = \text{softmax}(S^{(h)})V^{(h)}$$

attention weights

$$\text{scores: } S^{(h)} = \frac{Q^{(h)}K^{(h)T}}{\sqrt{d_k^{(h)}}}$$

$$\text{queries: } Q^{(h)} = XW_Q^{(h)}$$

$$\text{keys: } K^{(h)} = XW_K^{(h)}$$

$$\text{values: } V^{(h)} = XW_V^{(h)}$$

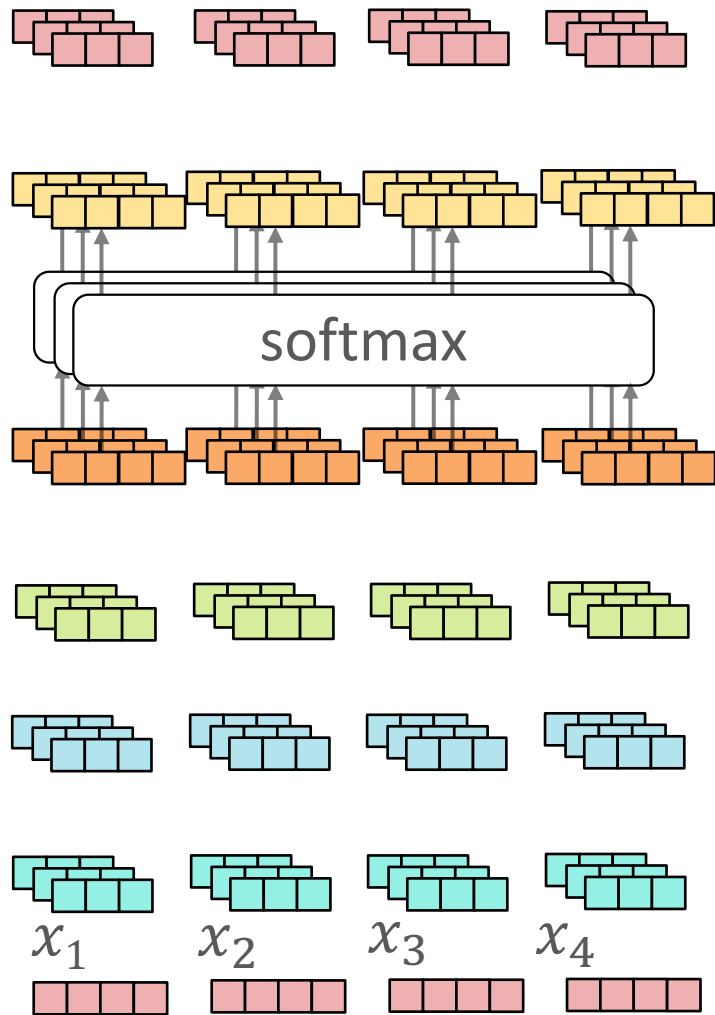
design matrix:  $X$



Key Takeaway:  
 All of this computation is

1. differentiable
2. highly parallelizable!

- Idea: just like we might want multiple convolutional filters in a convolutional layer, we might want multiple attention weights to learn different relationships between tokens!



$$H^{(h)} = \text{softmax}(S^{(h)})V^{(h)}$$

attention weights

scores: 
$$S^{(h)} = \frac{Q^{(h)}K^{(h)T}}{\sqrt{d_k^{(h)}}}$$

queries: 
$$Q^{(h)} = XW_Q^{(h)}$$

keys: 
$$K^{(h)} = XW_K^{(h)}$$

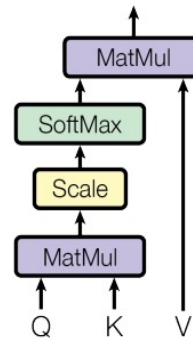
values: 
$$V^{(h)} = XW_V^{(h)}$$

design matrix:  $X$

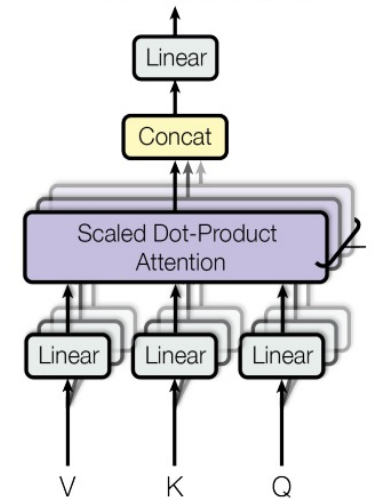
# Multi-head Scaled Dot-product Self-attention

- Idea: just like we might want multiple convolutional filters in a convolutional layer, we might want multiple attention weights to learn different relationships between tokens!

Scaled Dot-Product Attention



Multi-Head Attention

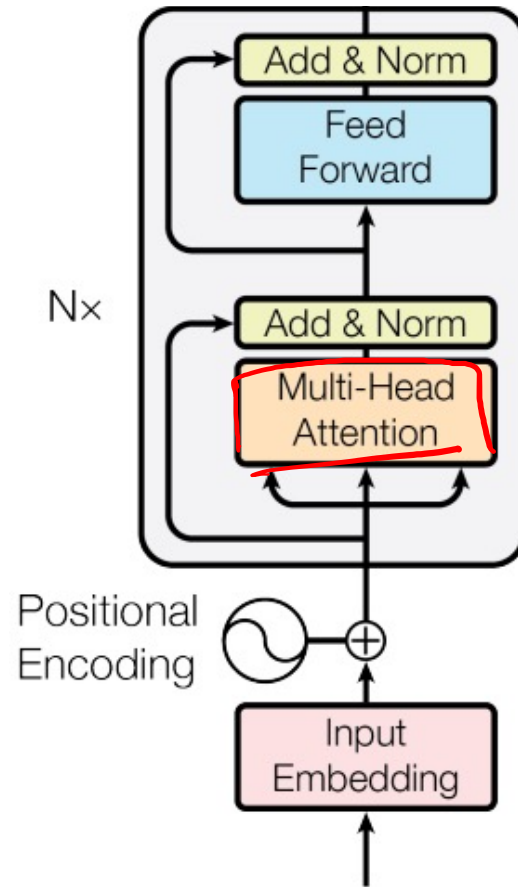


- The outputs from all the attention heads are concatenated together to get the final representation

$$H = [H^{(1)}, H^{(2)}, \dots, H^{(h)}]$$

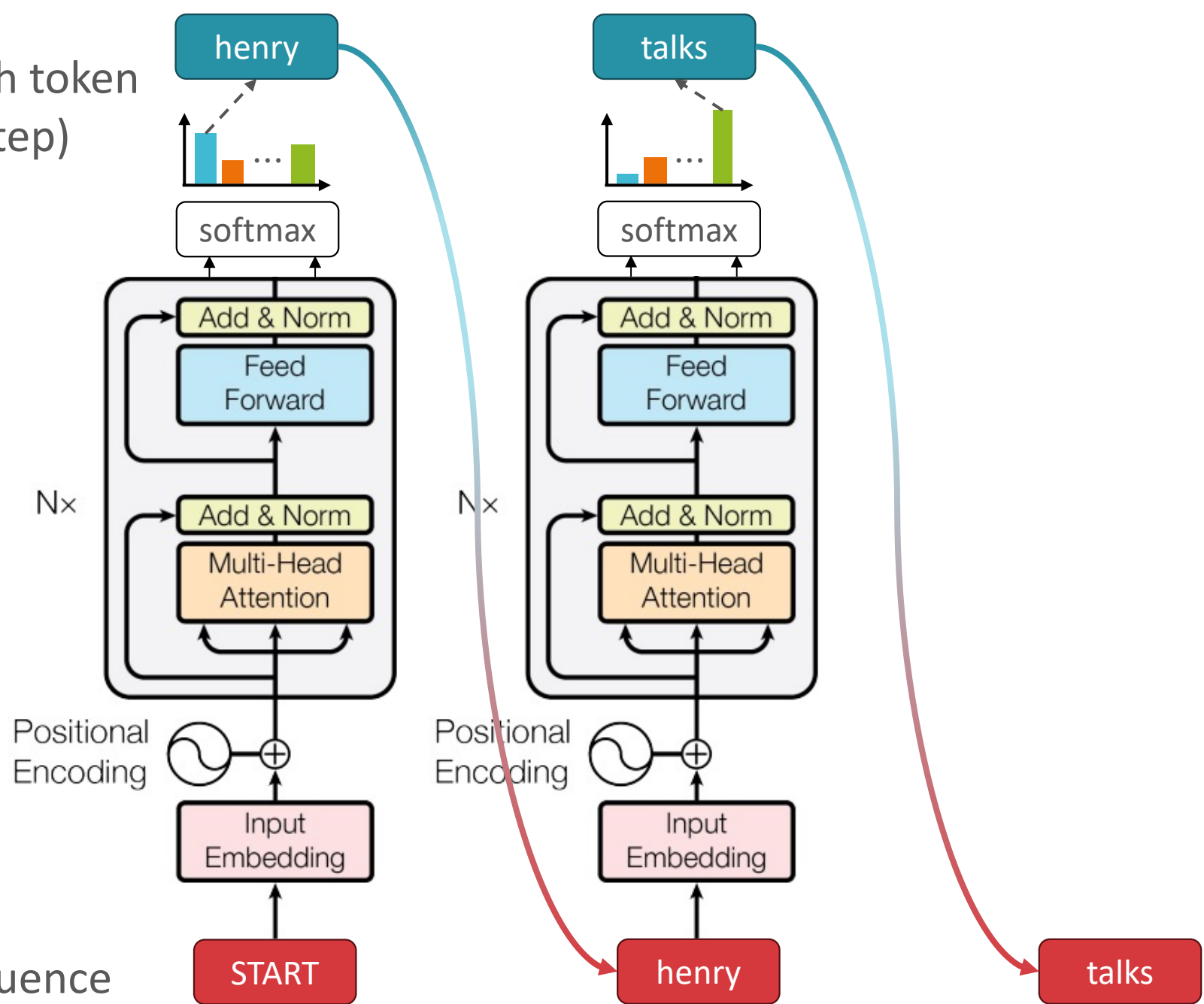
- Common architectural choice:  $d_v = D/h \rightarrow |H| = D$

# Transformers



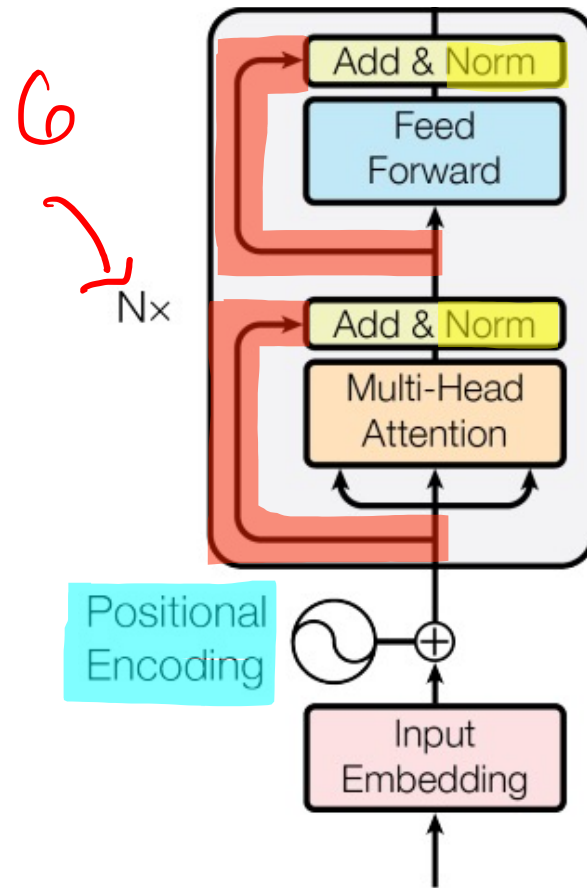
Generated sequence (use each token as the input to the next timestep)

# Transformer Language Models



Input sequence

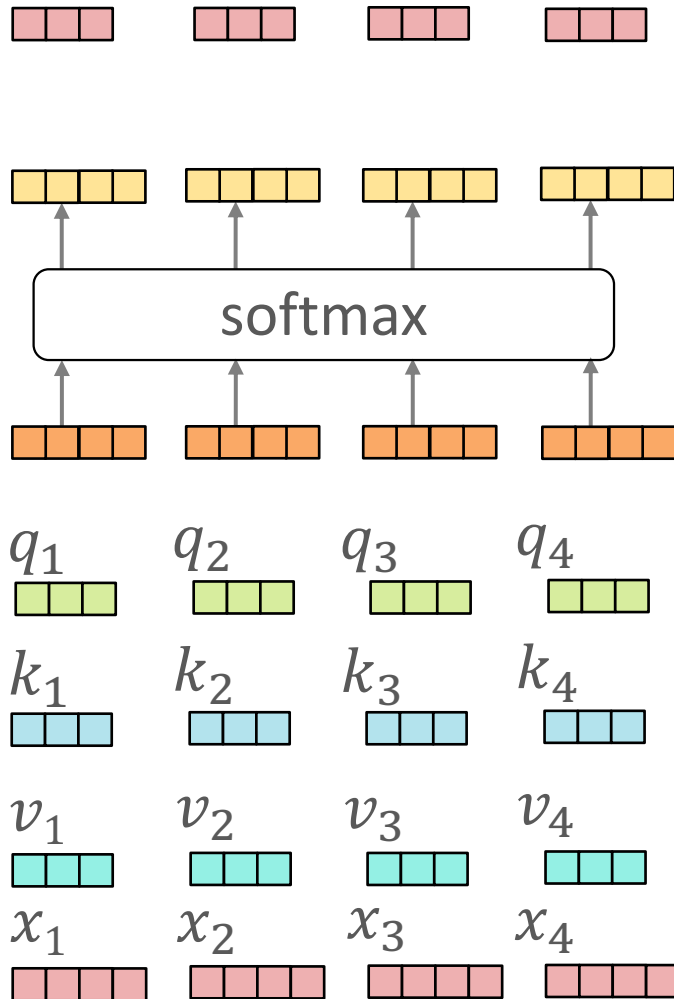
# Transformers



- In addition to multi-head attention, transformer architectures use
  1. Positional encodings
  2. Layer normalization
  3. Residual connections
  4. A fully-connected feed-forward network

# Scaled Dot-product Self-attention: Matrix Form

- Issue: if all tokens attend to every token in the sequence, then how does the model infer the order of tokens?



$$H = \text{softmax}(S)V \in \mathbb{R}^{N \times d_v}$$

attention weights

$$\text{scores: } S = \frac{QK^T}{\sqrt{d_k}} \in \mathbb{R}^{N \times N}$$

$$\text{queries: } Q = XW_Q \in \mathbb{R}^{N \times d_k}$$

$$\text{keys: } K = XW_K \in \mathbb{R}^{N \times d_k}$$

$$\text{values: } V = XW_V \in \mathbb{R}^{N \times d_v}$$

$$\text{design matrix: } X \in \mathbb{R}^{N \times D}$$

# Positional Encodings

- Issue: if all tokens attend to every token in the sequence, then how does the model infer the order of tokens?
- Idea: add a position-specific embedding  $p_t$  to the token embedding  $x_t$

$$x'_t = x_t + p_t$$

- Positional encodings can be
  - *fixed* i.e., some predetermined function of  $t$  or *learned* alongside the token embeddings
  - *absolute* i.e., only dependent on the token's location in the sequence or *relative* to the query token's location

# Layer Normalization

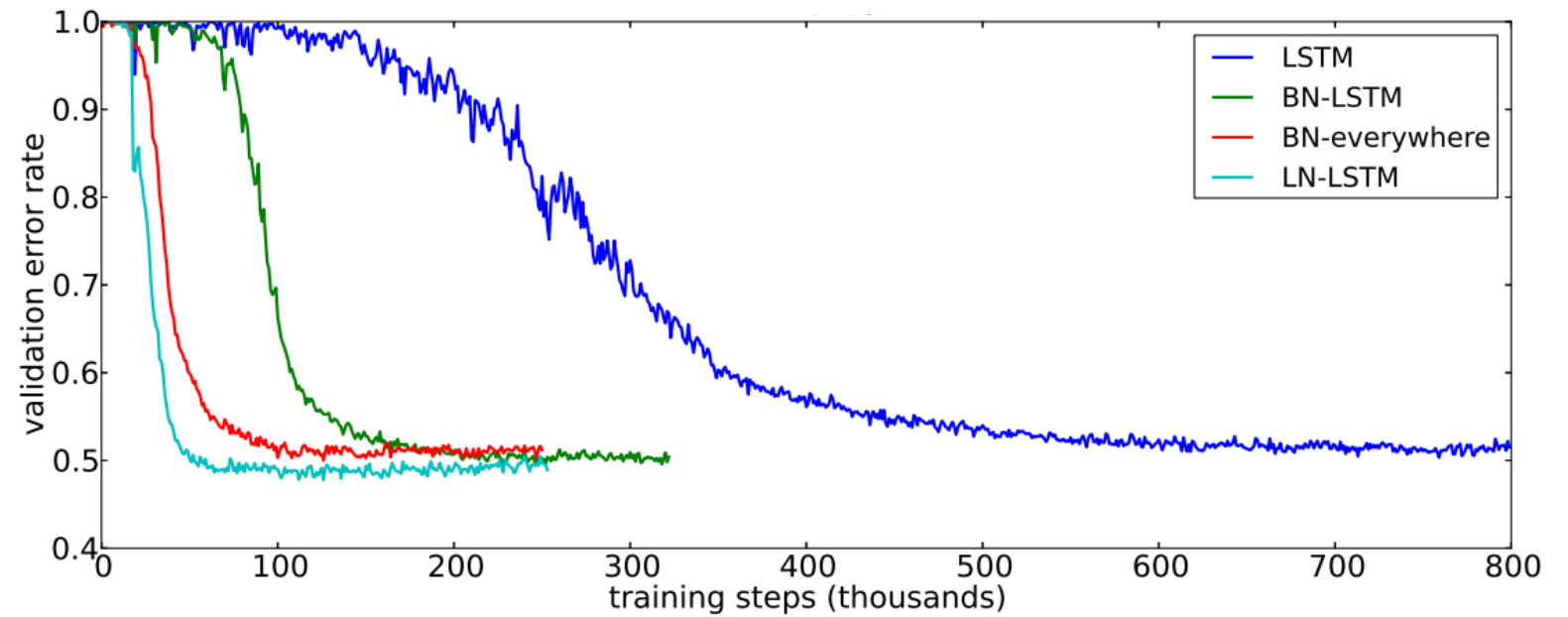
- Issue: for certain activation functions, the weights in later layers are **highly sensitive** to changes in the earlier layers
  - Small changes to weights in early layers are amplified so weights in deeper layers have to deal with massive dynamic ranges → slow optimization convergence
- Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$

$$H' = \gamma \left( \frac{H - \mu}{\sigma} \right) + \beta$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the values in the vector  $H$



# Layer Normalization



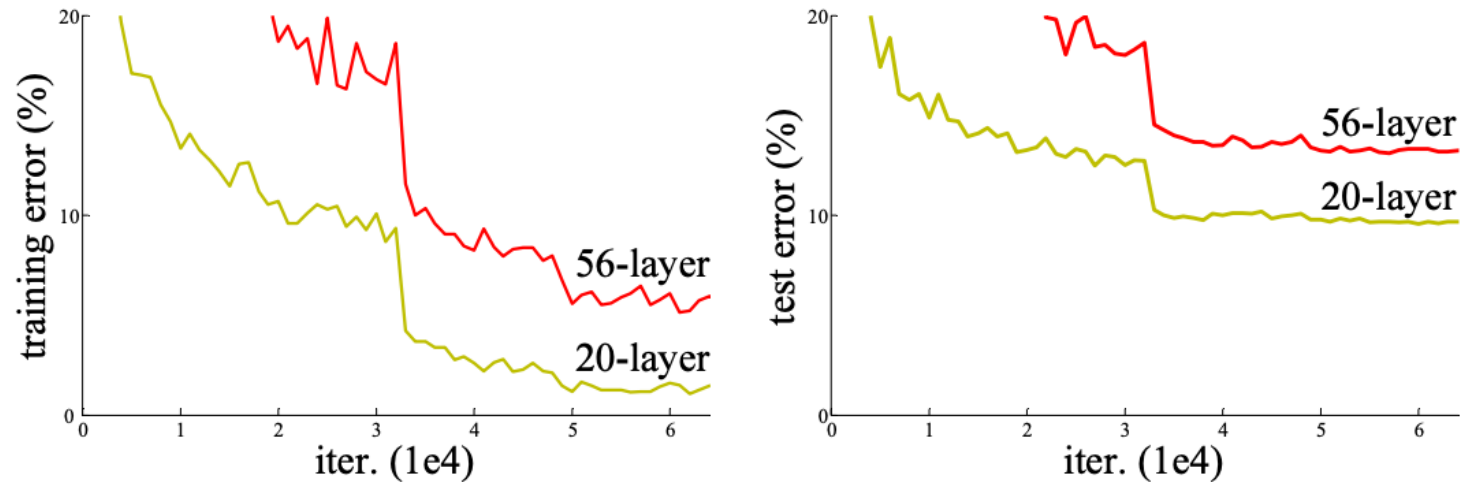
- Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$

$$H' = \gamma \left( \frac{H - \mu}{\sigma} \right) + \beta$$

where  $\mu$  is the mean and  $\sigma$  is the standard deviation of the values in the vector  $H$

# Residual Connections

- Observation: early deep neural networks suffered from the “degradation” problem where adding more layers actually made performance worse!



- Wait but this is ridiculous: if the later layers aren't helping, couldn't they just learn the identity transformation???
- Insight: neural network layers actually have a hard time learning the identity function

# Residual Connections

- Observation: early deep neural networks suffered from the “degradation” problem where adding more layers actually made performance worse!
- Idea: add the input embedding back to the output of a layer

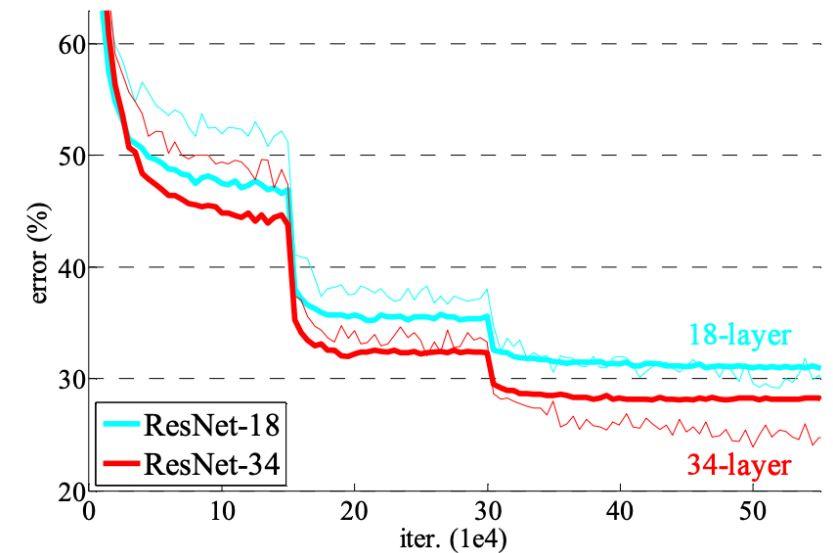
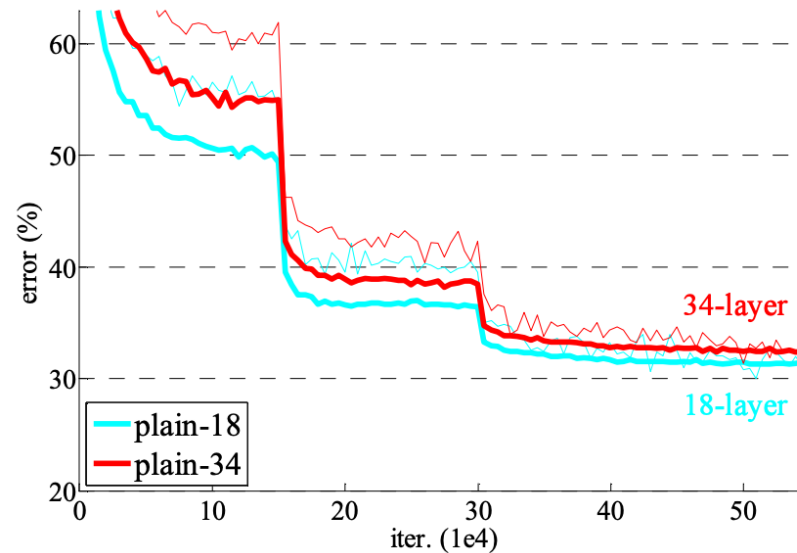
$$H' = H(x^{(i)}) + x^{(i)}$$

- Suppose the target function is  $f$ 
  - Now instead of having to learn  $f(x^{(i)})$ , the hidden layer just needs to learn the residual  $r = f(x^{(i)}) - x^{(i)}$
  - If  $f$  is the identity function, then the hidden layer just needs to learn  $r = 0$ , which is easy for a neural network!

# Residual Connections

- Observation: early deep neural networks suffered from the “degradation” problem where adding more layers actually made performance worse!
- Idea: add the input embedding back to the output of a layer

$$H' = H(x^{(i)}) + x^{(i)}$$



# Key Takeaways

- Language models fit joint probability distributions to sequences of inputs
  - Can be sampled from to generate text
- Attention allows information to directly pass between every pair of tokens
  - Attention can be used in conjunction with RNNs/LSTMs
  - However, (self-)attention can also be used in isolation
- Transformers consist of multi-head attention layers with residual connections, layer normalization and fully-connected layers