

10-701: Introduction to Machine Learning

# Lecture 13 - Attention & Transformers

Hoda Heidari

\* Slides adopted from F24 offering of 10701 by Henry Chai.

# 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 representation from earlier hidden layers into later ones.

# Recurrent Neural Networks

- Data points consists of (input **sequence**, label **sequence**) pairs, potentially of **varying lengths**

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

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

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

# Recurrent Neural Networks

- RNNs process inputs one time step at a time, using **recurrence**:

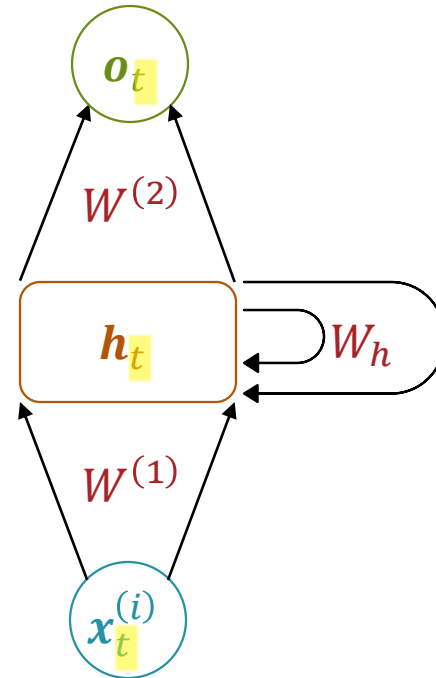
$$\mathbf{h}_t = \left[ 1, \theta \left( W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta(W^{(2)} \mathbf{h}_t)$$

Where  $\mathbf{h}_t$  serves as a summary or latent representation of the sequence up to time  $t$ .

- The same parameters  $W^{(1)}$ ,  $W_h$  and  $W^{(2)}$  are reused at every step.
- We can unroll the RNN for as many time steps as the sequence requires.
- So, at training and inference time, the RNN can run for different numbers of steps depending on the input length.

# Recurrent Neural Networks

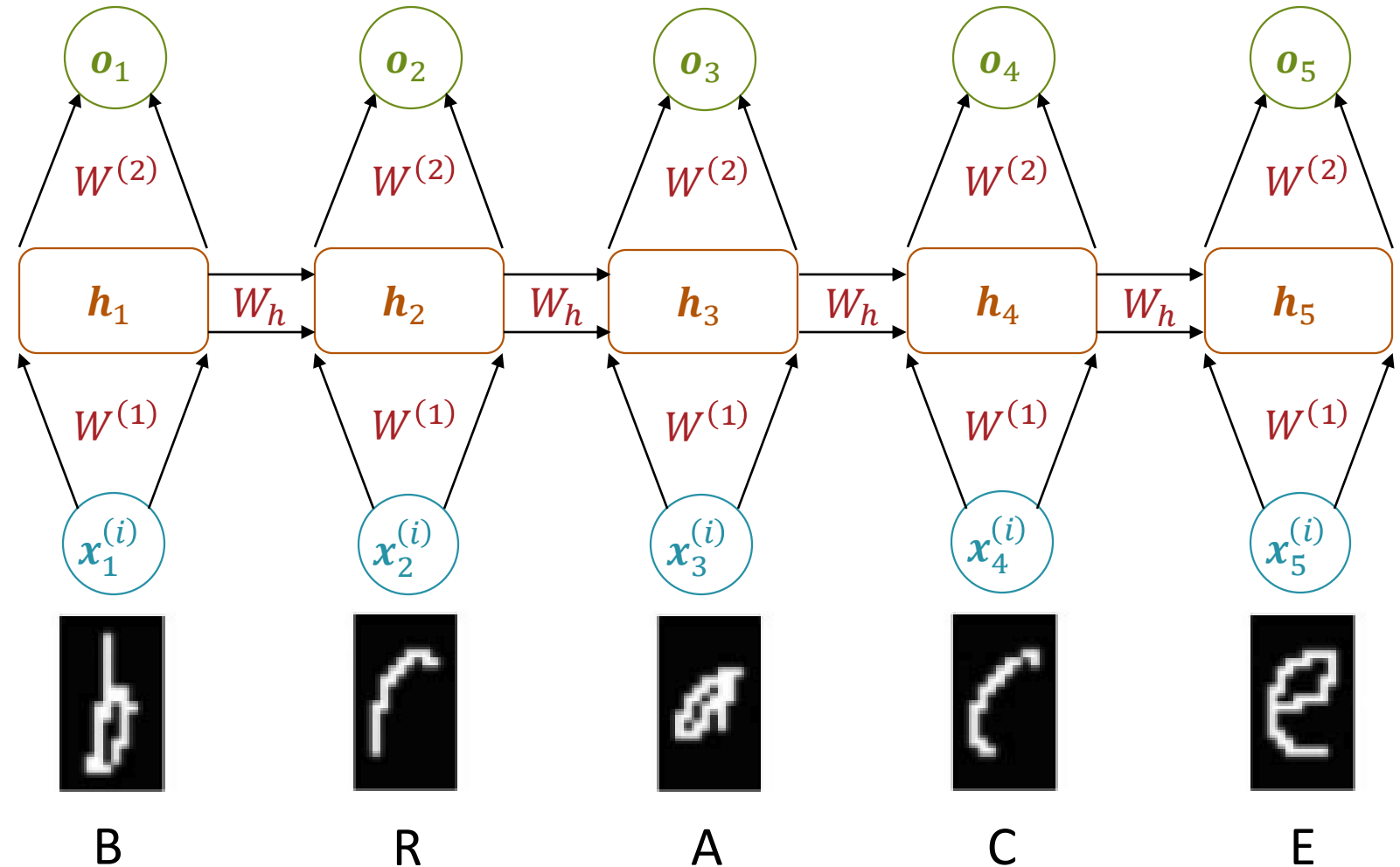
$$\mathbf{h}_t = \left[ 1, \theta \left( W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta(W^{(2)} \mathbf{h}_t)$$



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

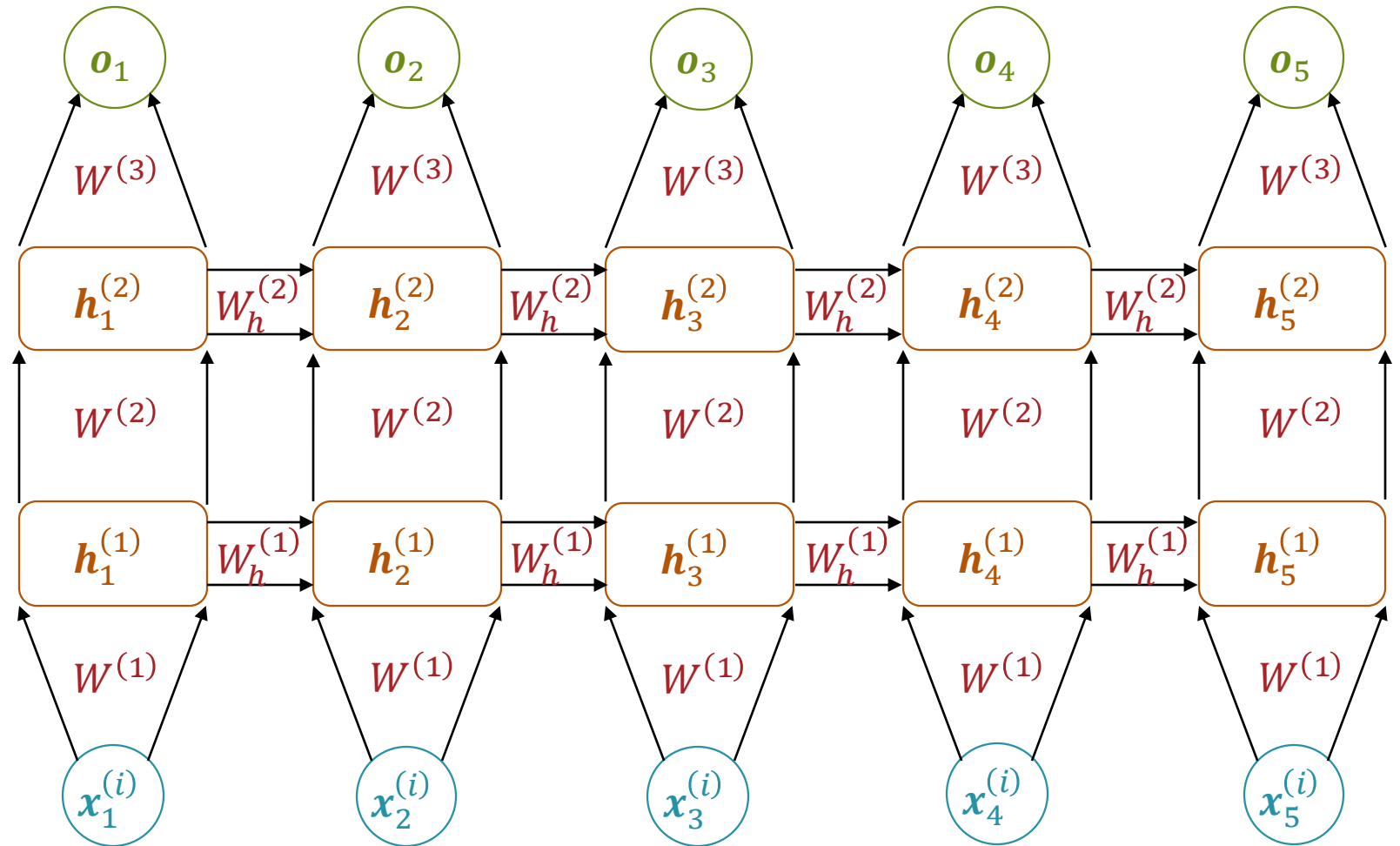
# Unrolling Recurrent Neural Networks

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



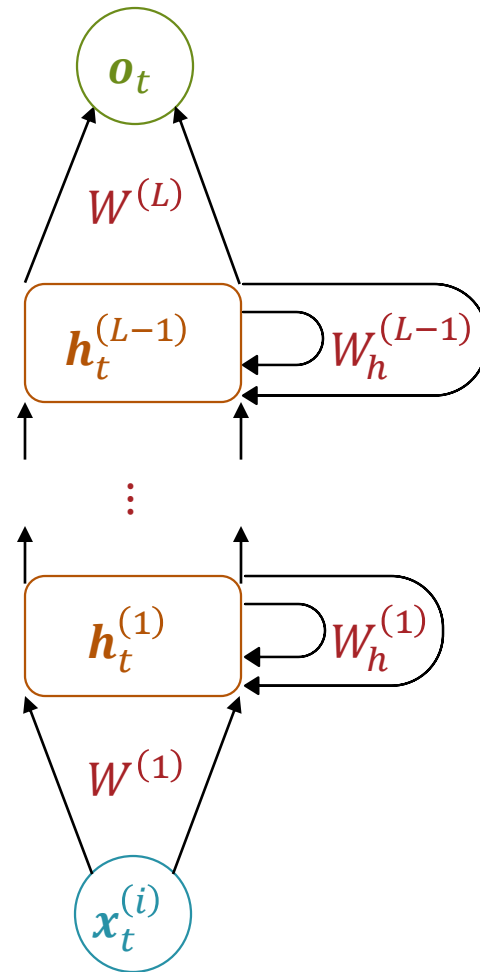
# Deep Recurrent Neural Networks

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



# Deep Recurrent Neural Networks

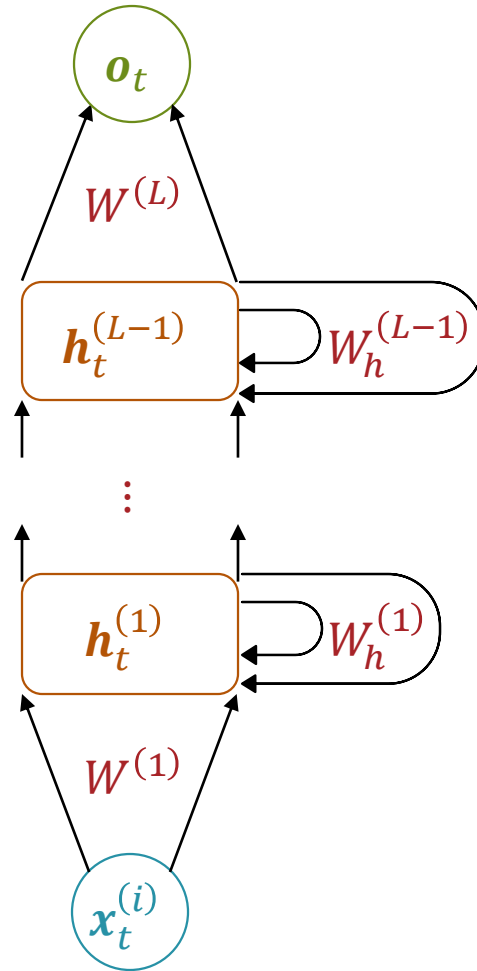
$$\mathbf{h}_t^{(l)} = \left[ 1, \theta \left( W^{(l)} \mathbf{h}_t^{(l-1)} + W_h^{(l)} \mathbf{h}_{t-1}^{(l)} \right) \right]^T \text{ and } \mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left( W^{(L)} \mathbf{h}_t^{(L-1)} \right)$$





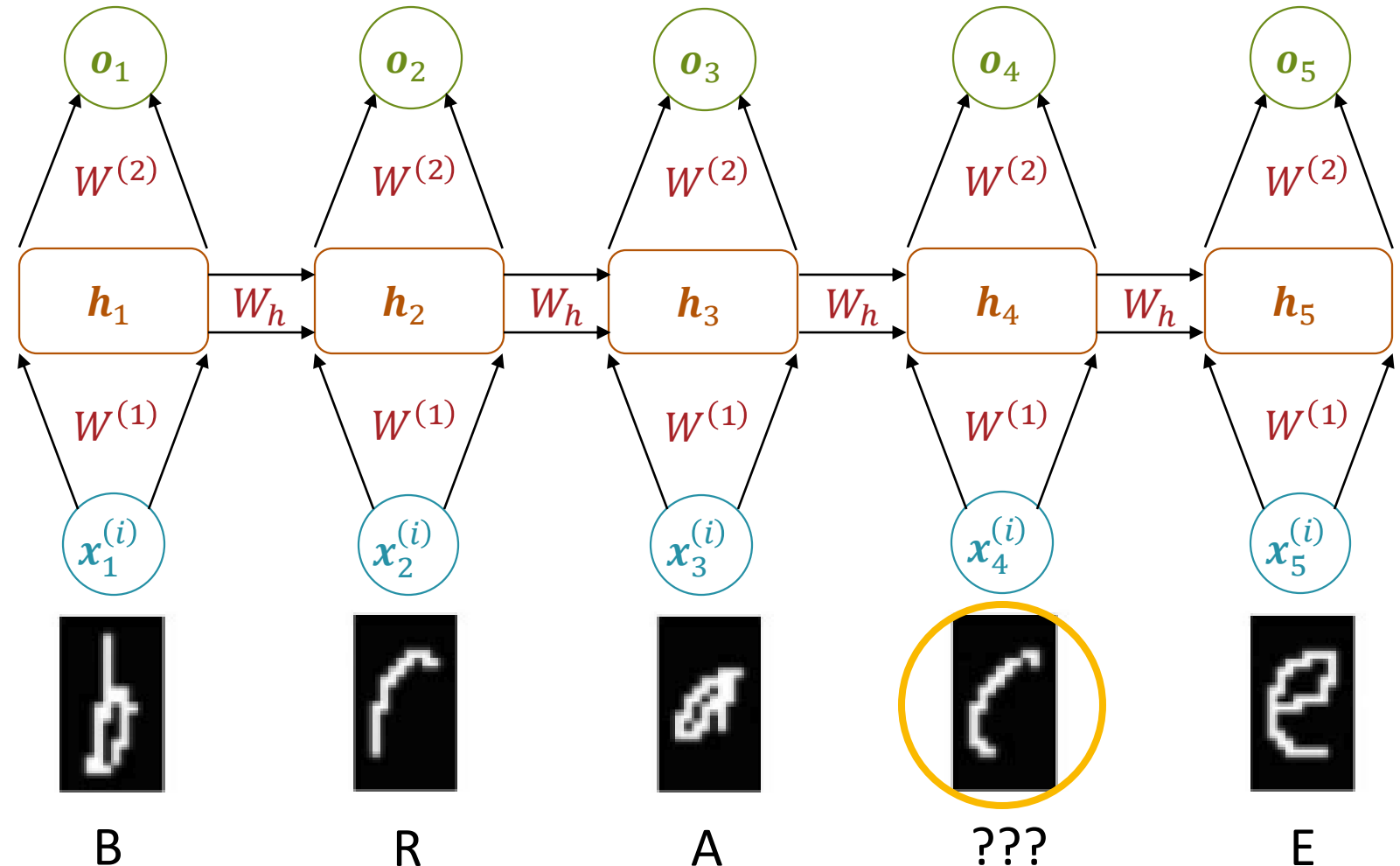
But why do we only pass information forward?  
What if later time steps have useful information as well?

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



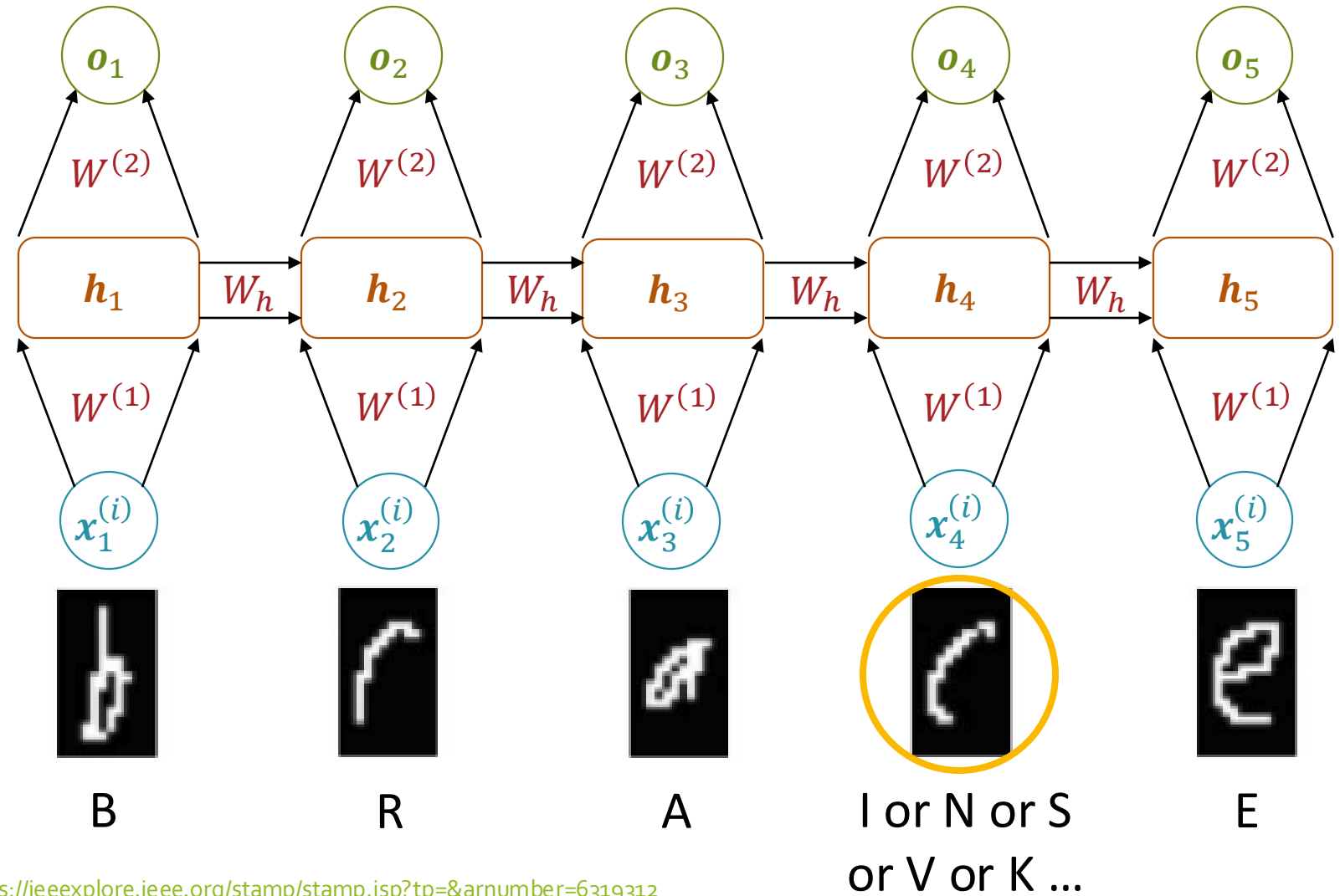
But why do we only pass information forward?  
What if later time steps have useful information as well?

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



But why do we only pass information forward?  
What if later time steps have useful information as well?

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

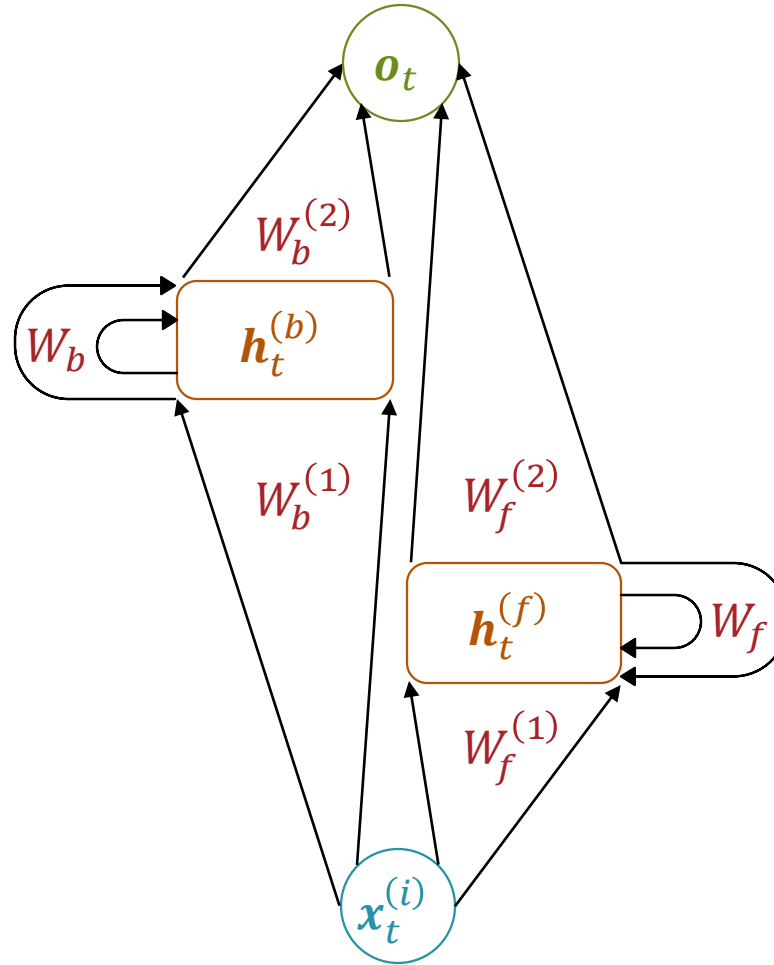


# Bidirectional Recurrent Neural Networks

- Bidirectional Recurrent Neural Networks (BiRNNs) capture **context from both the past and the future** of a sequence.
- A BiRNN has two RNNs:
  - one  $\mathbf{h}_t^{(f)}$  processes the sequence **forward in time**
  - one  $\mathbf{h}_t^{(b)}$  processes it **backward in time**
  - The combination contains information from the entire sequence centered around position  $t$ .

# Bidirectional Recurrent Neural Networks

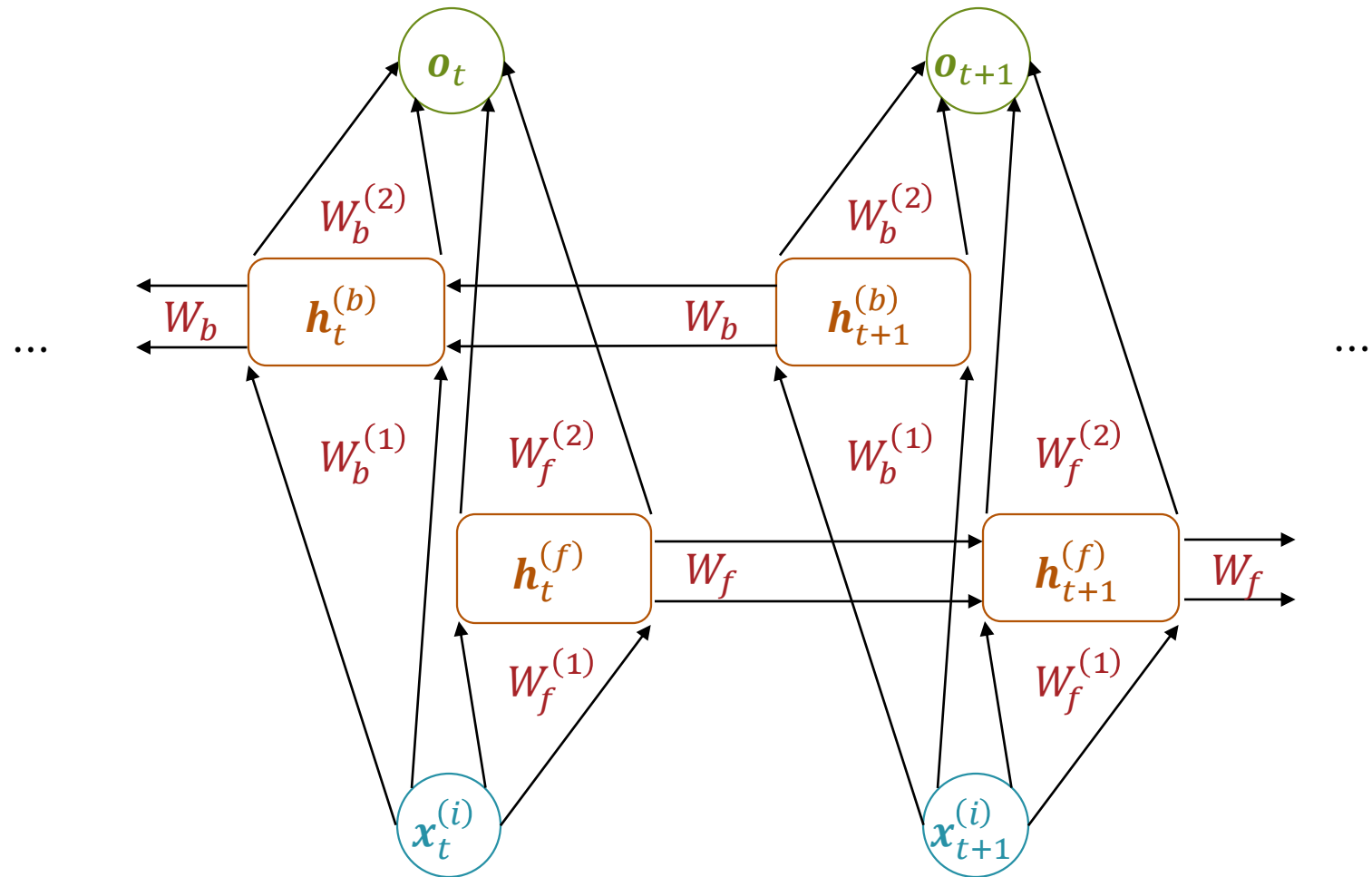
$$\mathbf{h}_t^{(f)} = \left[ 1, \theta \left( W_f^{(1)} \mathbf{x}_t^{(i)} + W_f \mathbf{h}_{t-1}^{(f)} \right) \right]^T \text{ and } \mathbf{h}_t^{(b)} = \left[ 1, \theta \left( W_b^{(1)} \mathbf{x}_t^{(i)} + W_b \mathbf{h}_{t+1}^{(b)} \right) \right]^T$$
$$\mathbf{o}_t = \hat{y}_t^{(i)} = \theta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



# Unrolling Bidirectional Recurrent Neural Networks

$$\mathbf{o}_t = \hat{y}_t^{(i)} = \theta \left( W_f^{(2)} \mathbf{h}_t^{(f)} + W_b^{(2)} \mathbf{h}_t^{(b)} \right)$$

$$\mathbf{h}_t^{(f)} = \left[ 1, \theta \left( W_f^{(1)} \mathbf{x}_t^{(i)} + W_f \mathbf{h}_{t-1} \right) \right]^T \text{ and } \mathbf{h}_t^{(b)} = \left[ 1, \theta \left( W_b^{(1)} \mathbf{x}_t^{(i)} + W_b \mathbf{h}_{t+1} \right) \right]^T$$



# Inference via Bidirectional Recurrent Neural Networks

- Inference when the entire sequence is available (e.g., sequence labeling, sentence-level classification)
  1. run forward RNN on the full sequence;
  2. Run backward RNN on the reversed sequence;
  3. Combine representations at time t:  $[\mathbf{h}_t^{(f)}, \mathbf{h}_t^{(b)}]$
  4. Feed the combination to a classifier

$$\hat{y}_t^{(i)} = \theta \left( W_f^{(2)} \mathbf{h}_t^{(f)} + W_b^{(2)} \mathbf{h}_t^{(b)} \right)$$

# Training RNNs

- A (deep/bidirectional) RNN simply represents a (somewhat complicated) computation graph
  - Weights ( $W^{(1)}$ ,  $W_h$  and  $W^{(2)}$ ) are shared between different timesteps, significantly reducing the number of parameters to be learned!
- Can be trained using (stochastic) gradient descent/backpropagation → “backpropagation through time”



# Loss Functions for RNNs

- **Sequence-to-sequence prediction:** Token-wise cross-entropy averaged across time steps

$$\mathcal{L} = \frac{1}{T} \sum_t CE(y_t, \text{softmax}(Wh_t))$$

- **Regression over sequences** (e.g., forecasting, speech features): MSE between predicted and target sequence

$$\mathcal{L} = \frac{1}{T} \sum_t \|y_t - \hat{y}_t\|^2$$

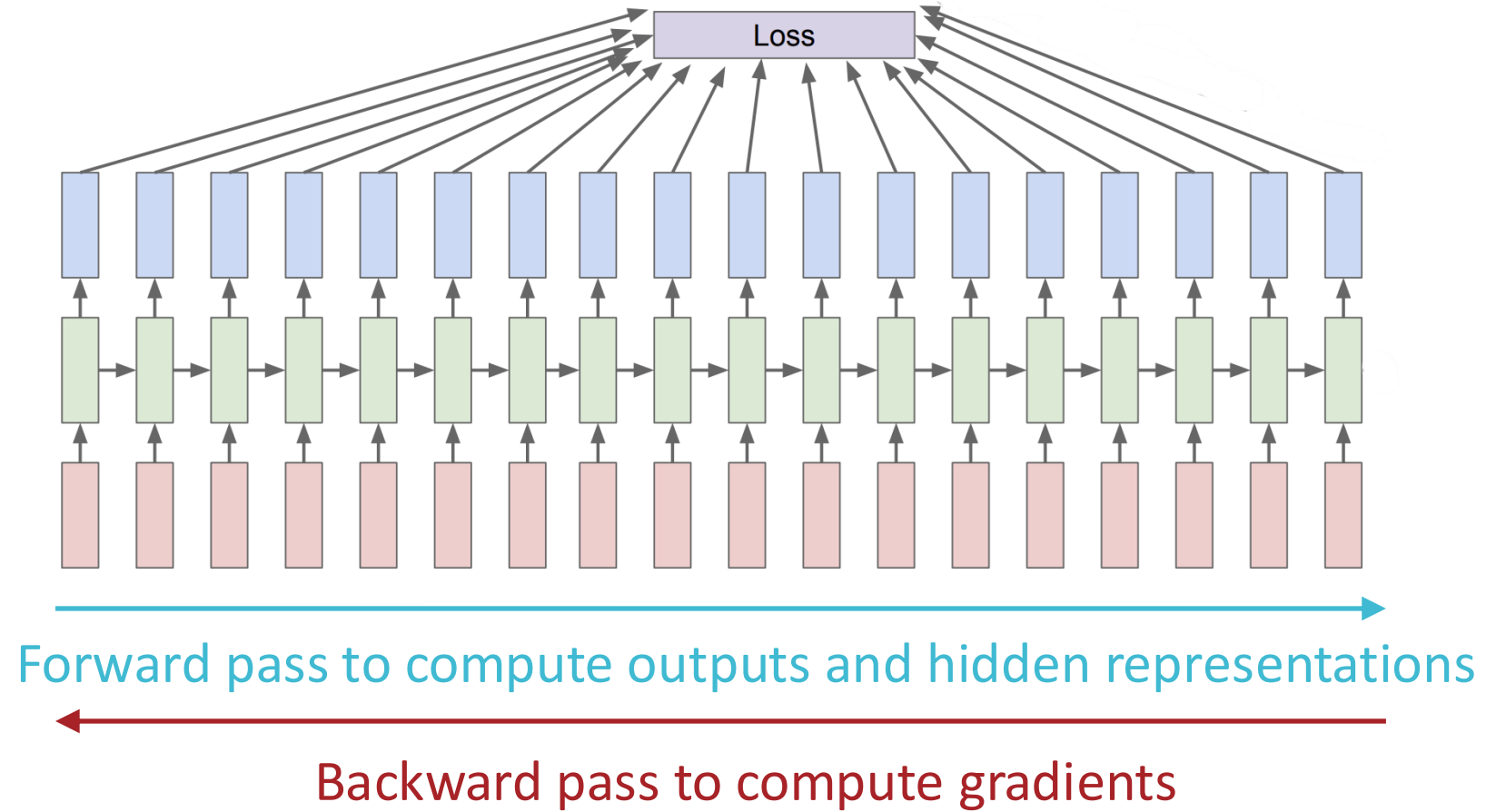
- **Sequence classification** (one label per sequence, e.g., sentiment classification): Cross-entropy on the final (or pooled) hidden state

$$\mathcal{L} = CE(y, \text{softmax}(Wh_T))$$

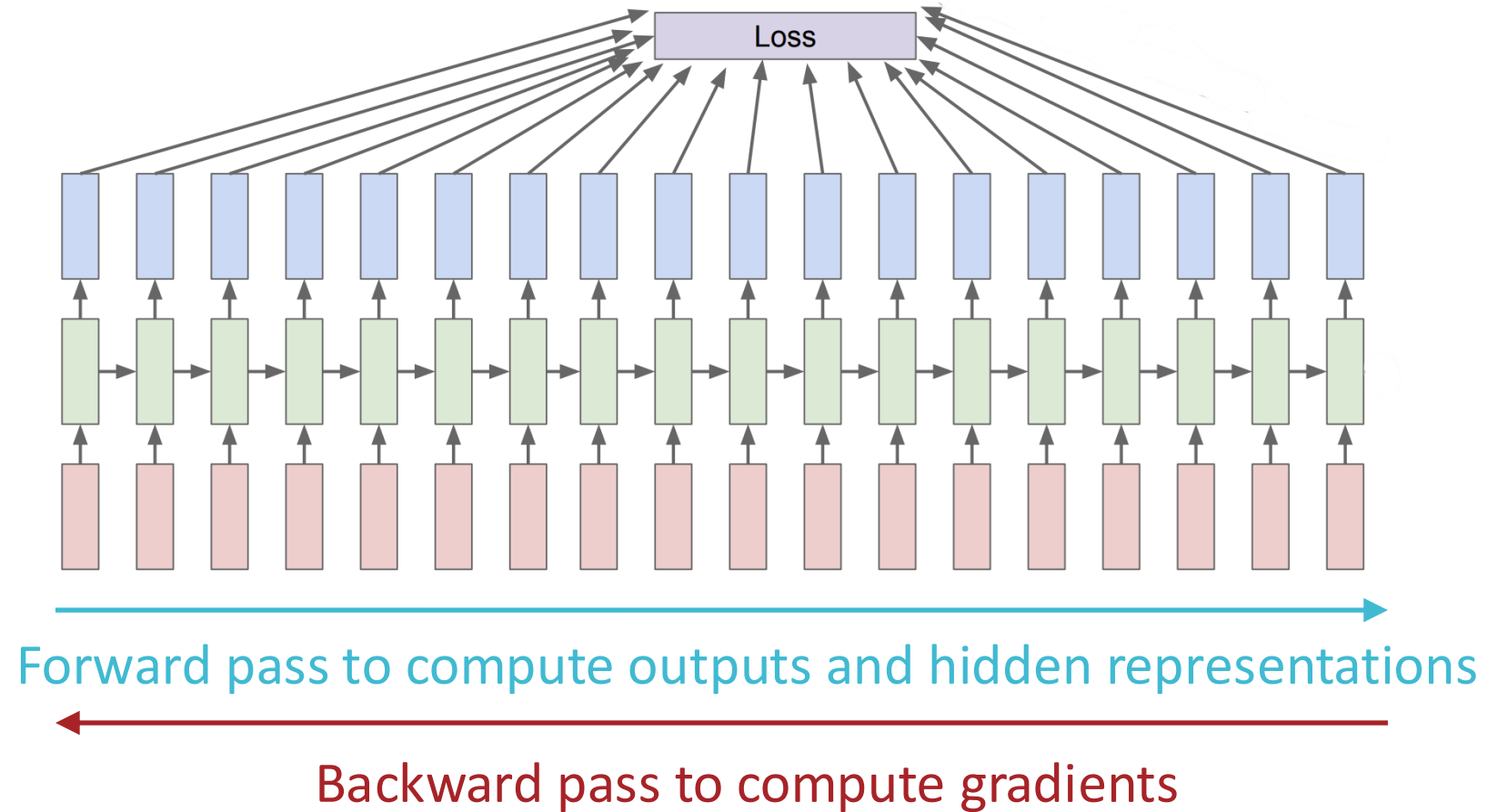
# Backprop Through Time

- Each hidden state  $h_t$  influences not only its immediate output  $y_t$ , but also all future hidden states  $h_{t+1}, h_{t+2}, \dots$
- Thus, each parameter ( $W^{(1)}$ ,  $W_h$  and  $W^{(2)}$ ) affects the loss indirectly **through time**.
- So, during training, need to propagate the gradient back through all those time steps.
- To train the RNN, we unroll it over the sequence, treating it like a deep feed-forward network with  $T$  layers — one per time step — all sharing the same parameters.
- Then, we apply standard backpropagation over this unrolled network.

# Training RNNs



# Training RNNs: Challenges

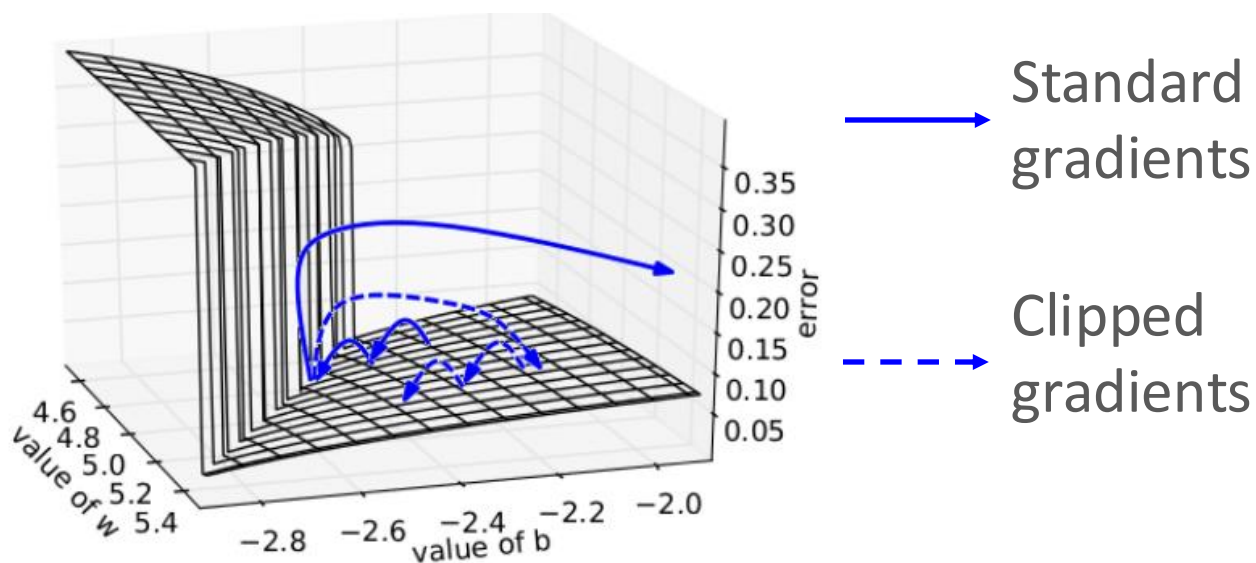


- Issue: as the sequence length grows, the gradient is more likely to explode or vanish

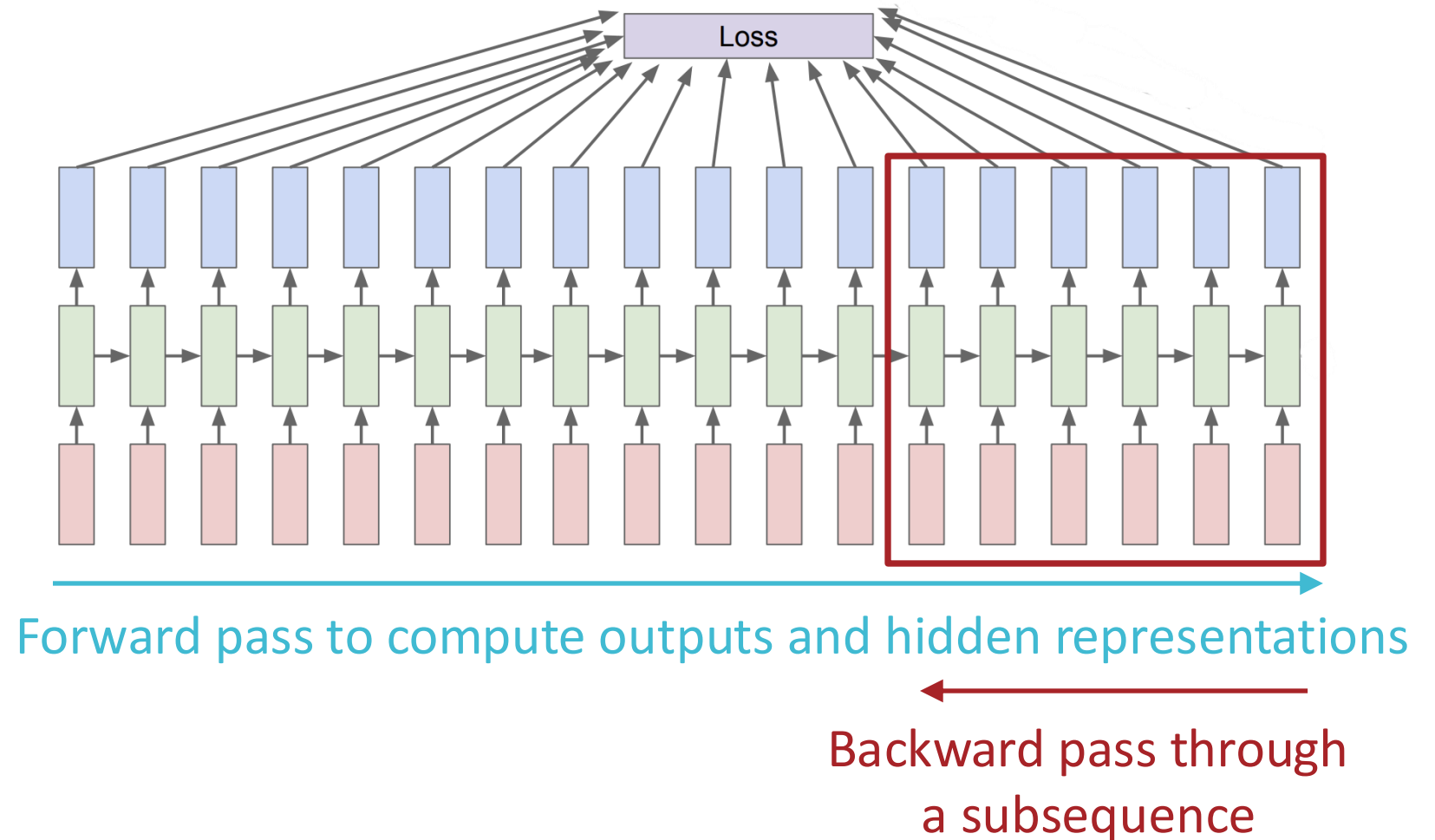
# Gradient Clipping (Pascanu et al., 2013)

- Common strategy to deal with exploding gradients: if the magnitude of the gradient ever exceeds some threshold, simply scale it down to the threshold

$$G = \begin{cases} \nabla_W \ell^{(i)} & \text{if } \|\nabla_W \ell^{(i)}\|_2 \leq \tau \\ \left( \frac{\tau}{\|\nabla_W \ell^{(i)}\|_2} \right) \nabla_W \ell^{(i)} & \text{otherwise} \end{cases}$$



# Truncated Backpropagation Through Time



- Idea: limit the number of time steps to backprop through

# 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  $\mathbf{h}_t$  but also maintains a separate internal *state*,  $\mathbf{C}_t$
- The flow of information through a cell is manipulated by three *gates*:
  - An **input** gate,  $\mathbf{I}_t$ , that controls how much the state looks like the normal RNN hidden layer
  - An **output** gate,  $\mathbf{O}_t$ , that “releases” the hidden representation to later timesteps
  - A **forget** gate,  $\mathbf{F}_t$ , 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  $\mathbf{h}_t$  but also maintains a separate internal *state*,  $\mathbf{C}_t$
- Gates are implemented as **sigmoid**s: a value of 0 would be a fully *closed* gate and 1 would be fully *open*

$$I_t = \sigma \left( W_{ix} \mathbf{x}_t^{(i)} + W_{ih} \mathbf{h}_{t-1} \right)$$

$$O_t = \sigma \left( W_{ox} \mathbf{x}_t^{(i)} + W_{oh} \mathbf{h}_{t-1} \right)$$

$$F_t = \sigma \left( W_{fx} \mathbf{x}_t^{(i)} + W_{fh} \mathbf{h}_{t-1} \right)$$

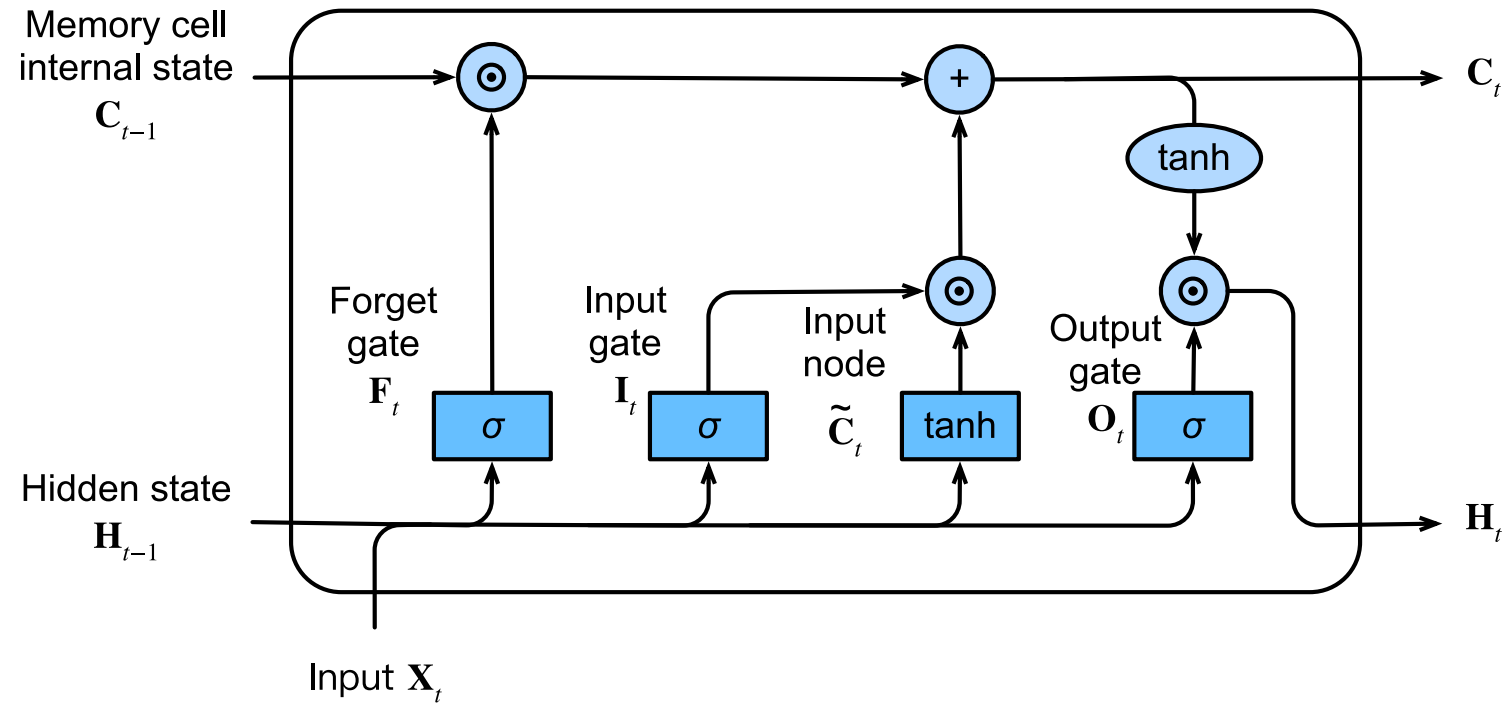
$$\mathbf{C}_t = F_t \odot \mathbf{C}_{t-1} + I_t \odot \theta \left( W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right)$$

$$\mathbf{h}_t = \mathbf{C}_t \odot \mathbf{O}_t$$



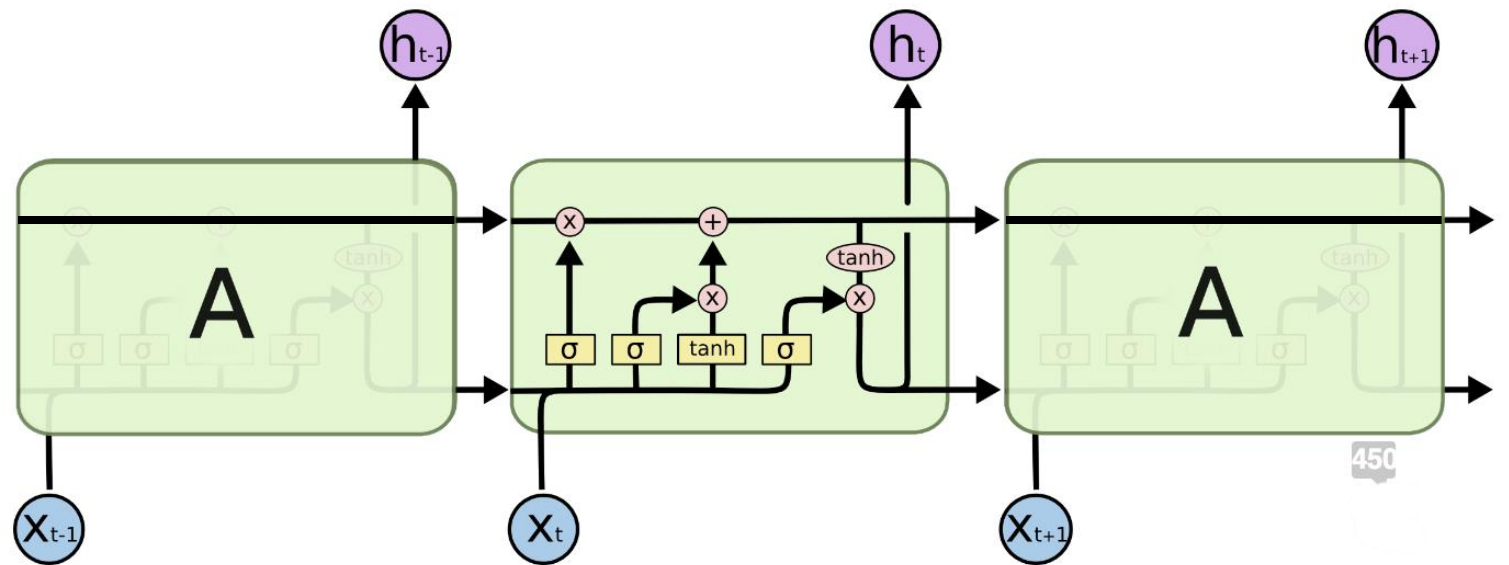
# Long Short-Term Memory (Hochreiter & Schmidhuber, 1997)

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



# 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  $h_t$  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!

# LSTM and Gradients

- In a plain RNN,  $\mathbf{h}_t = \left[ 1, \theta \left( W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right) \right]^T$  so the gradient w.r.t. an old state involves repeated multiplication by  $W_h$  and by  $\theta$ . Over long time, that product tends to 0 (vanish) or  $\infty$  (explode).
- Instead, the LSTM keeps a *side channel* (additive update + gate values near 1)

$$C_t = F_t \odot C_{t-1} + I_t \odot \theta \left( W^{(1)} \mathbf{x}_t^{(i)} + W_h \mathbf{h}_{t-1} \right)$$

- The old memory contributes **linearly**, thus gradients are less fragile.
- The backprop signal for something at time  $T$  reaching back to time  $t \ll T$  can travel mostly along the **cell-state skip path**.
- If gates permit it (learned behavior), that path is close to identity — so the gradient can survive many steps.

# 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.

# Key Takeaways

- Recurrent neural networks use contextual information to reason about sequential data.
  - Can still be used for long-range dependency →  
Transformers replaced LSTM in most large-scale NLP because attention scales better and captures long-range dependency without recurrence.
- LSTMs allow earlier information to reach later timesteps without directly affecting intermediate hidden representations.

# 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

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

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

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

# 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

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



# 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

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

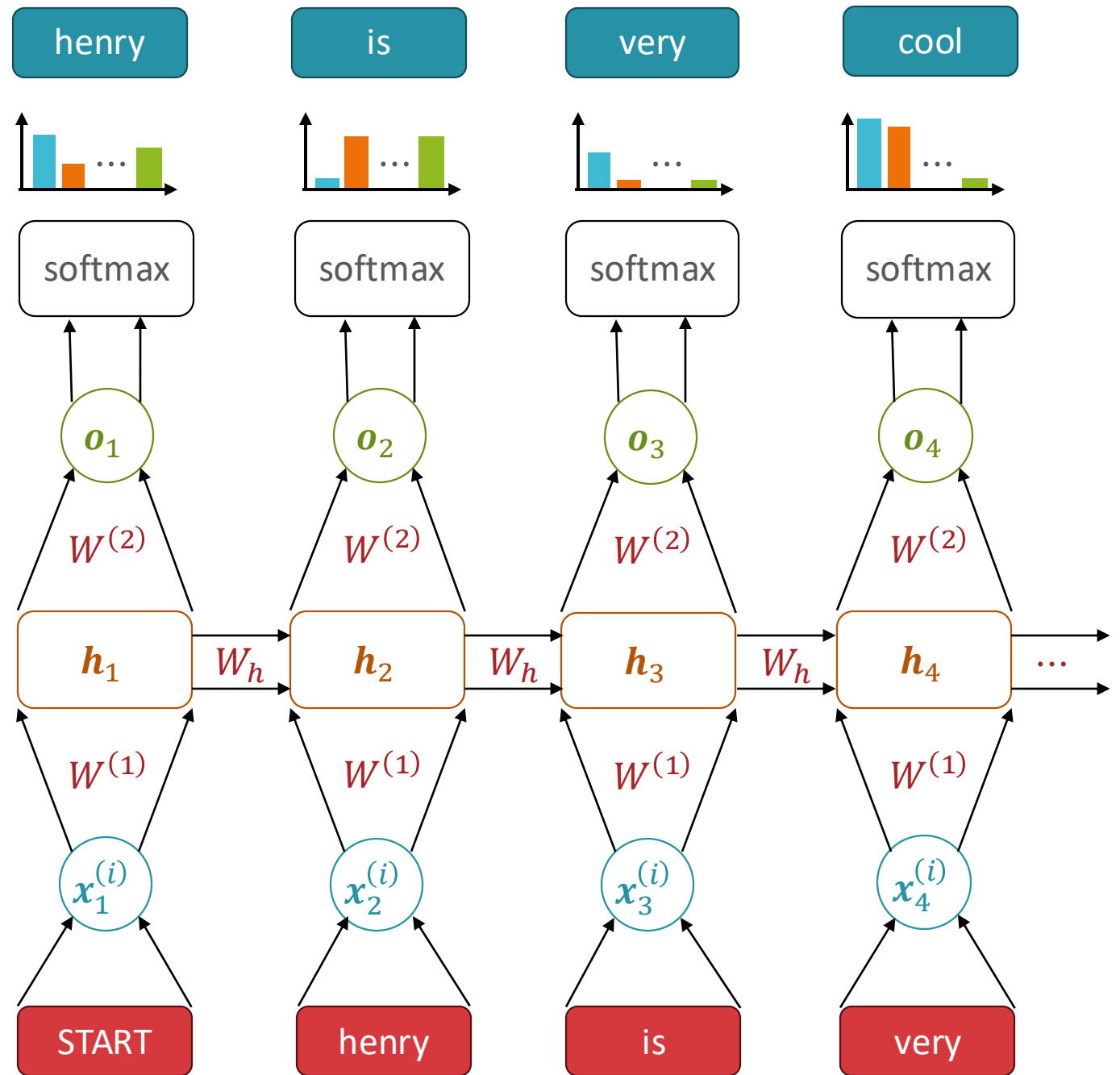
- ~~Use the chain rule of probability~~ Just throw an RNN at it!

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

# RNN Language Models: Training

Target sequence (try to  
predict the next word)

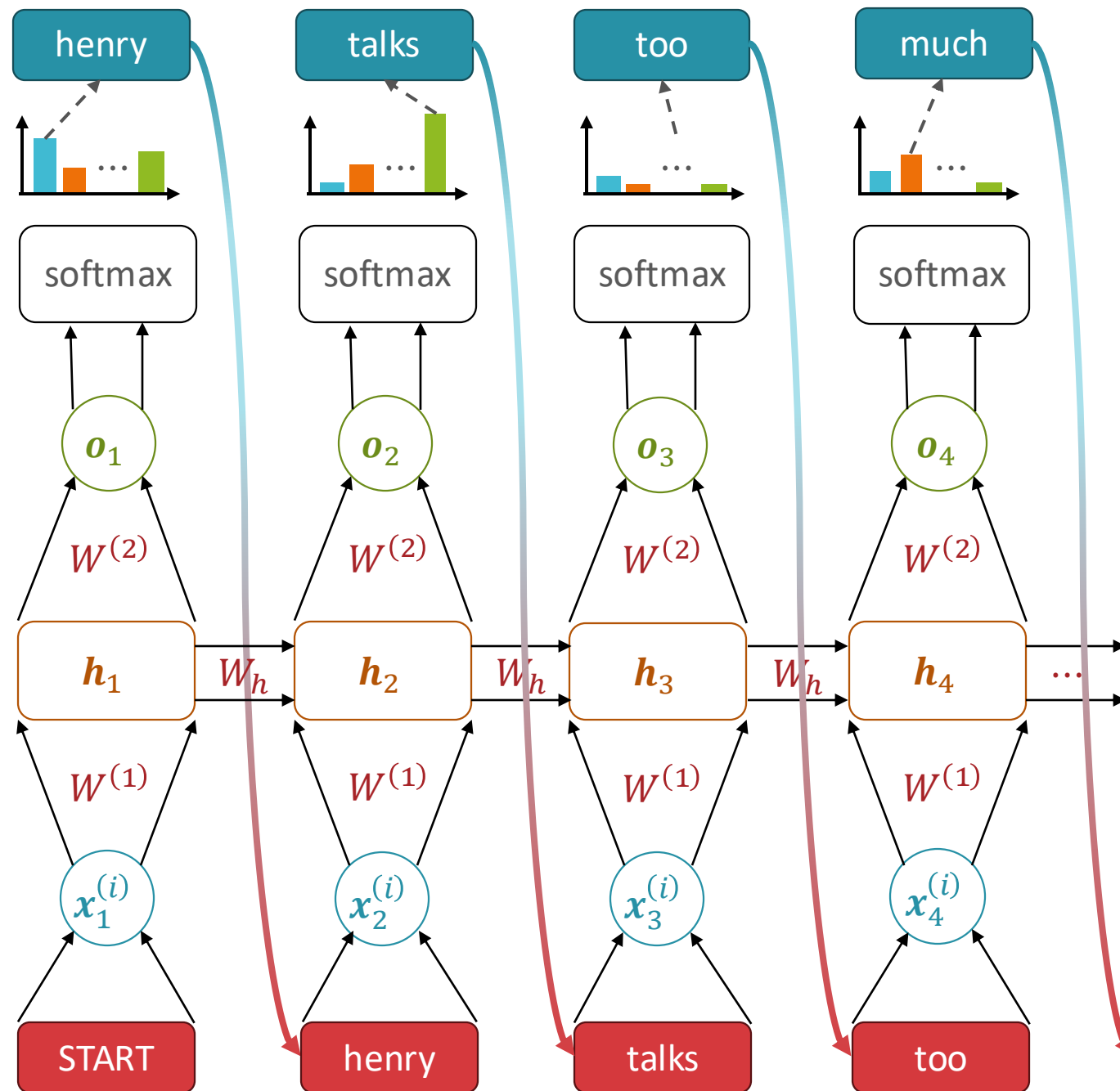
Input sequence



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

# RNN Language Models: Sampling

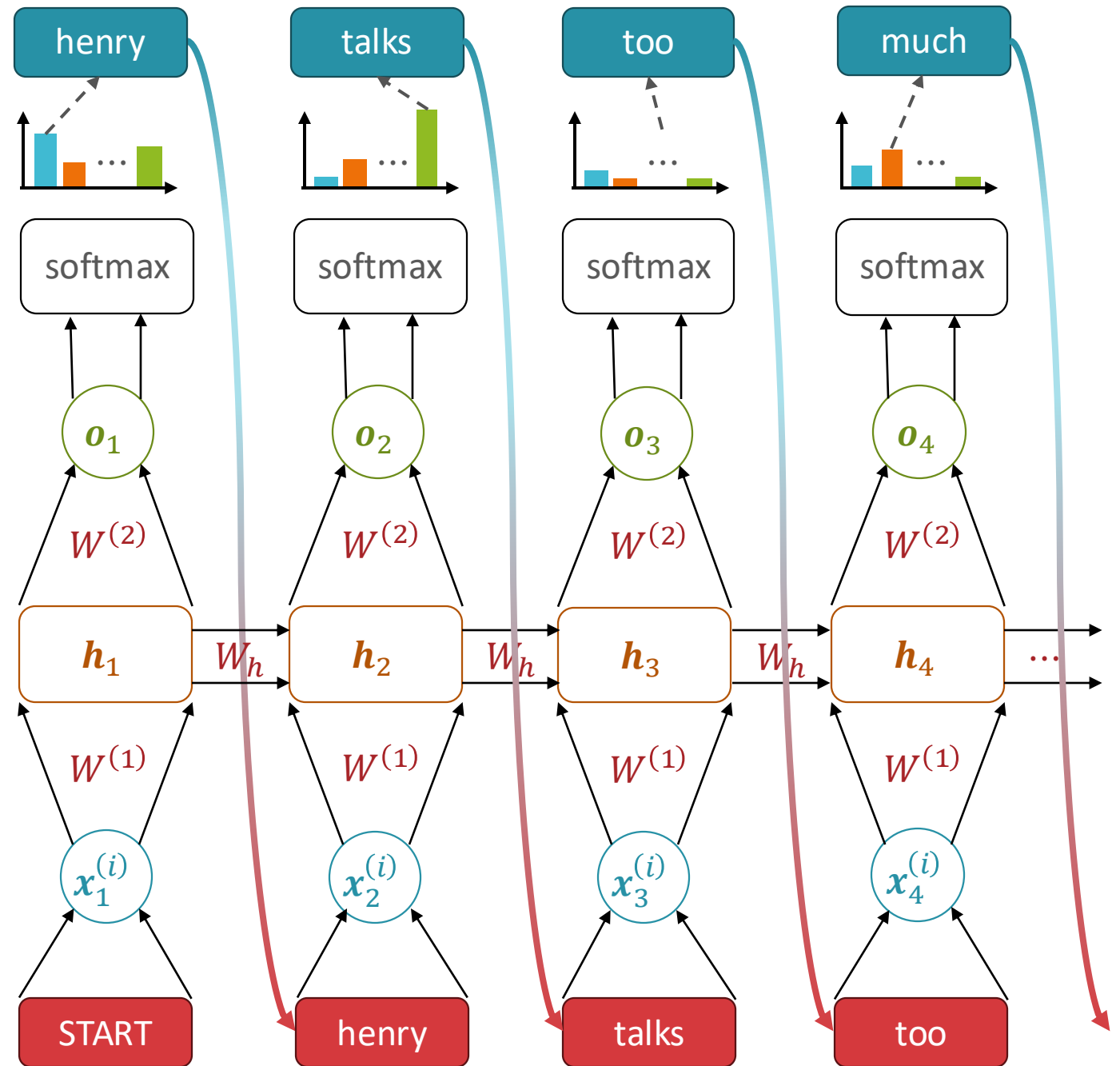
Input sequence



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



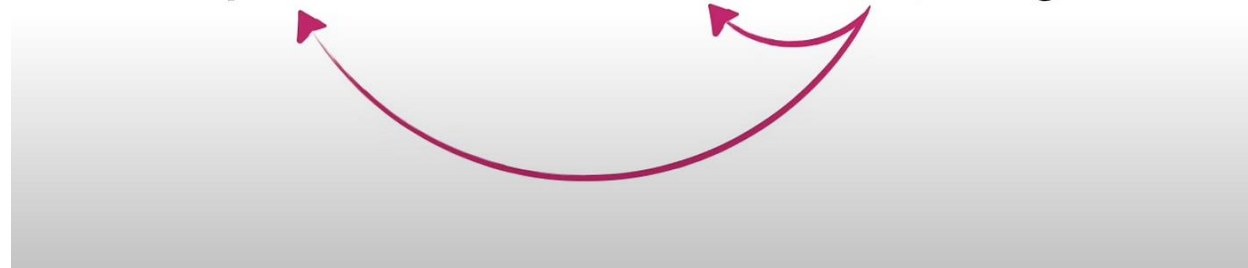
# RNN Language Models: Pros & Cons

- Pros:
  - Can handle arbitrary sequence lengths without having to increase model size (i.e., # of learnable parameters)
  - Trainable via backpropagation
- Cons
  - Vanishing/exploding gradients
    - Can be addressed by LSTMs
  - 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)

# Transformers

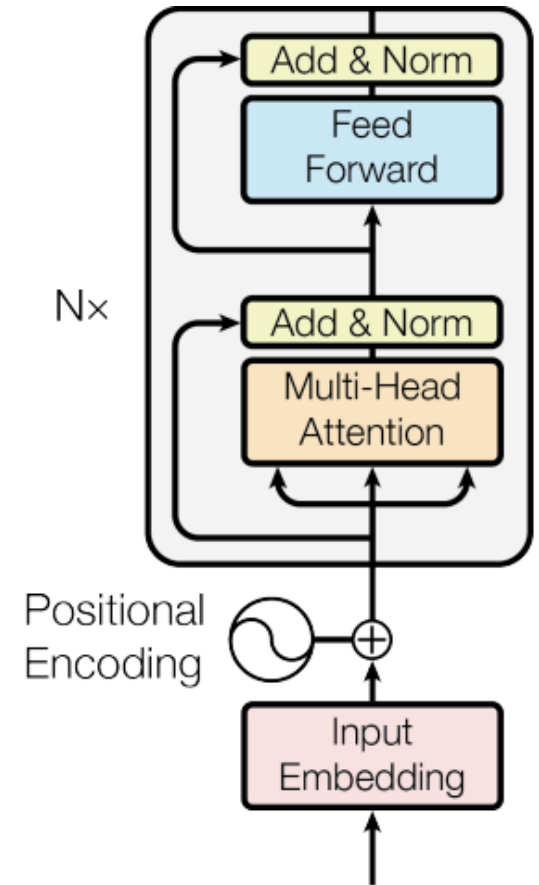
- Instead of reading text word-by-word like an RNN/LSTM, a transformer looks at the entire sentence at once and asks:  
“For this word, which other words in the sentence matter right now, and by how much?”
- This mechanism is called **self-attention**: every word computes links to other words and assigns strengths to them.

The **pizza** came out of the **oven** and **it** tasted good!



# Transformers

1. **Embed words:** Words are turned into vectors (points in a space that roughly encode semantic meaning).
2. **Attention Scoring:** Each word looks at all other words and asks: “How relevant are you to me given the task at hand?”
3. **Blend information:** Each word builds a new representation by taking a weighted mix of the other words, weighted by their scores.
4. **Repeat:** Stack that many times-- layers repeat the same pattern, so deeper layers see richer relational structure/representations.



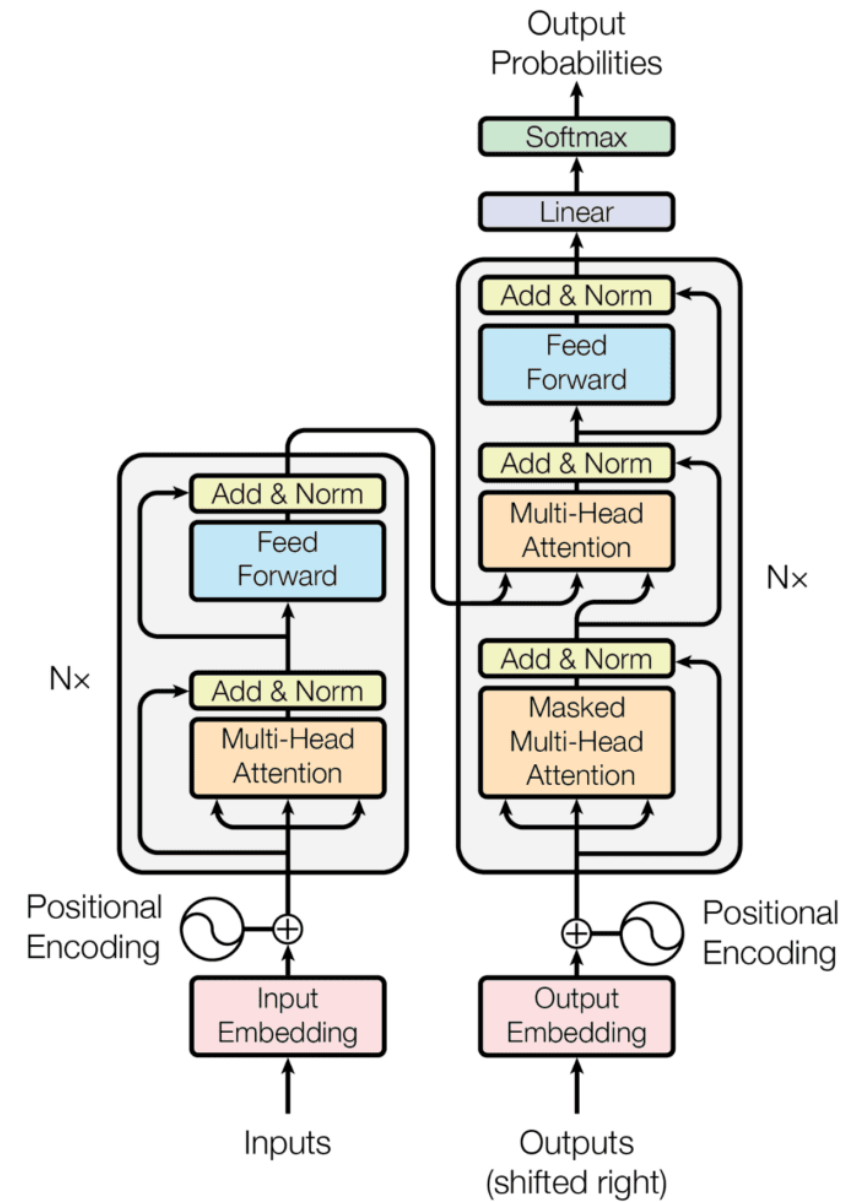


# Encoder vs. Decoder

## 5. Generate/understand text:

A decoder side produces next tokens by repeatedly applying the same attention process.

- **Encoder:** map an input sequence to representations, fed into a decoder.
- **Decoder:** receives the output of the encoder and the decoder output at the previous time step to generate new output.



# Transformers vs. RNNs

RNNs	Transformers
RNNs only get <b>cumulative context</b> from the past; they do not have easy access to distant words unless encoded in hidden state.	In a transformer, at each layer, every token sees every other token ( <b>global context</b> ).
RNNs can forget over long distances even with LSTMs.	Transformers use <b>self-attention</b> to directly connect any word to any other word
RNNs/LSTMs read sequences <i>step-by-step</i> ( $x_1 \rightarrow x_2 \rightarrow x_3 \dots$ ), so they can't parallelize across time.	Transformers consume the whole sequence at once and compute attention in parallel $\rightarrow$ efficiency.

# The Attention Mechanism

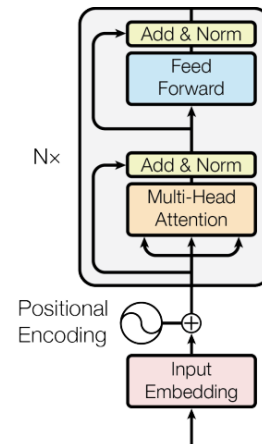
**Analogy:** attention as soft lookup in a dictionary

1. The token  $x'$  broadcasts its **query**  $q = w_Q^T x'$
2. Every other token  $x_t$  offers up its **key**  $k_t = W_K x_t$
3. The model computes a *similarity* between the query and keys--scoring how relevant each token is to the query.

$$s_t(x', x_t) = \frac{k_t^T q}{\sqrt{\text{length}(q)}}$$

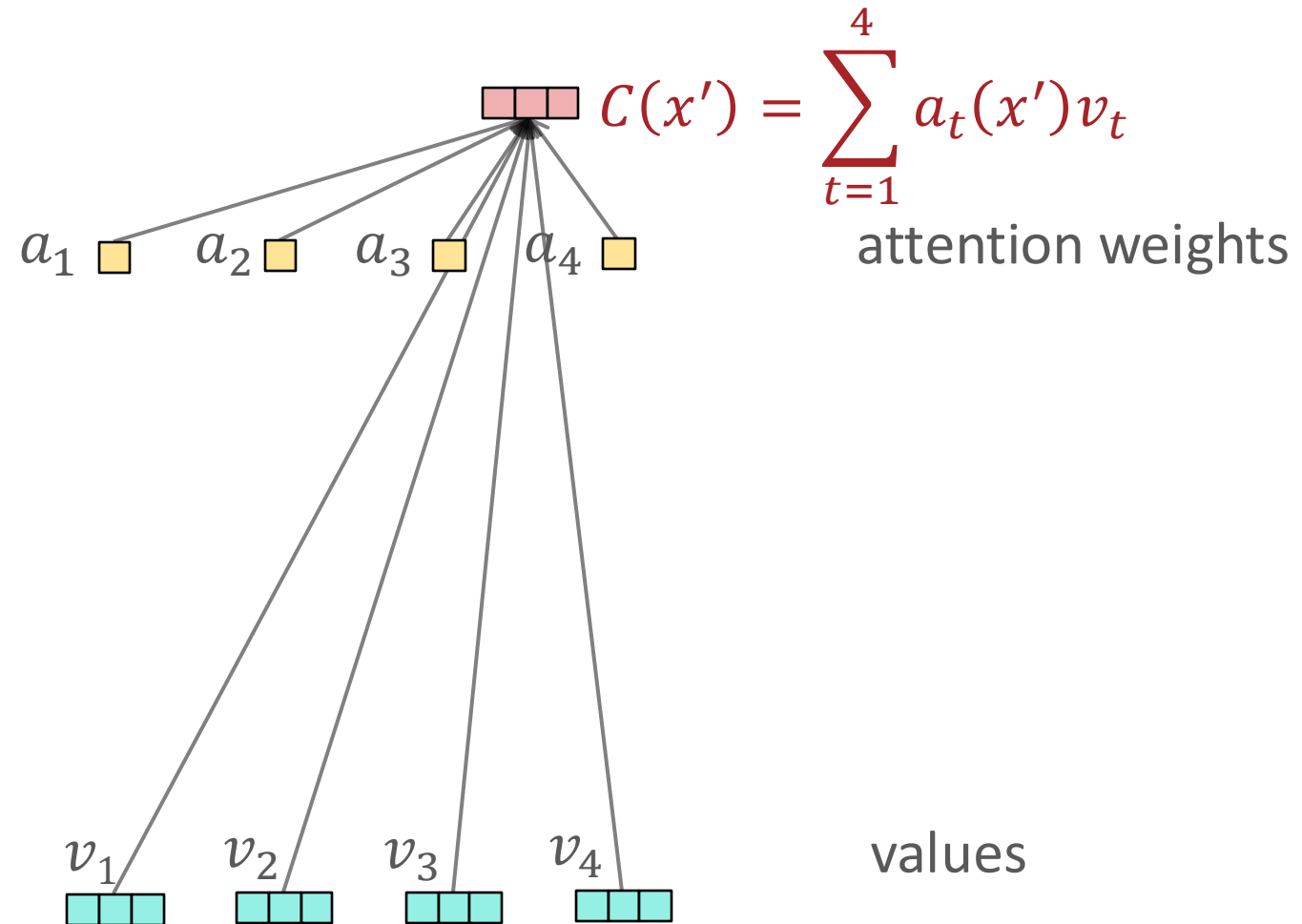
4. Those *scores* are turned into weights:  $\text{softmax}(s(x', x_t))$
5. The model returns a weighted sum of the **values** as the contextualized representation:

$$\sum_t \text{softmax}(s(x', x_t)) v(x_t)$$



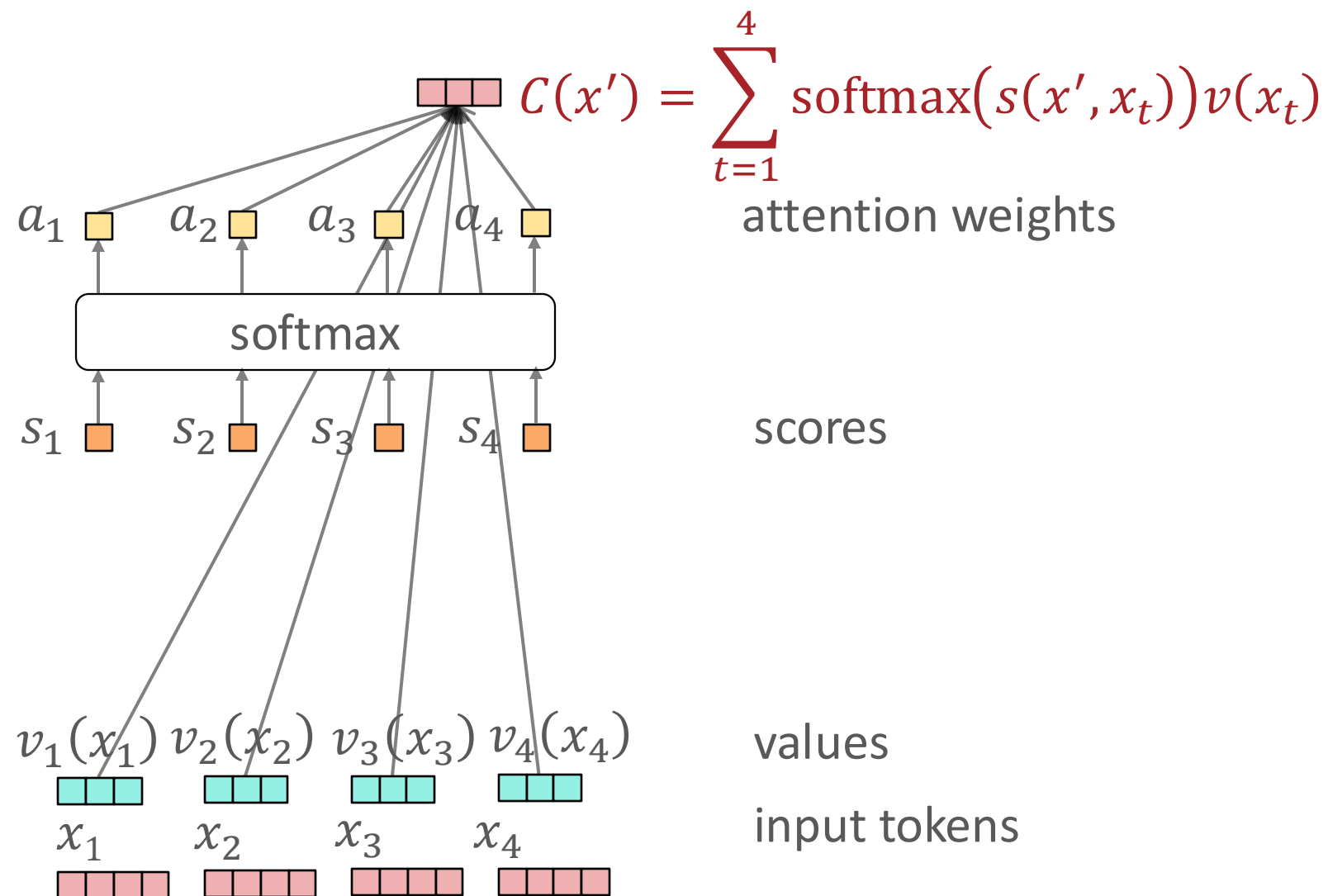
# Attention

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



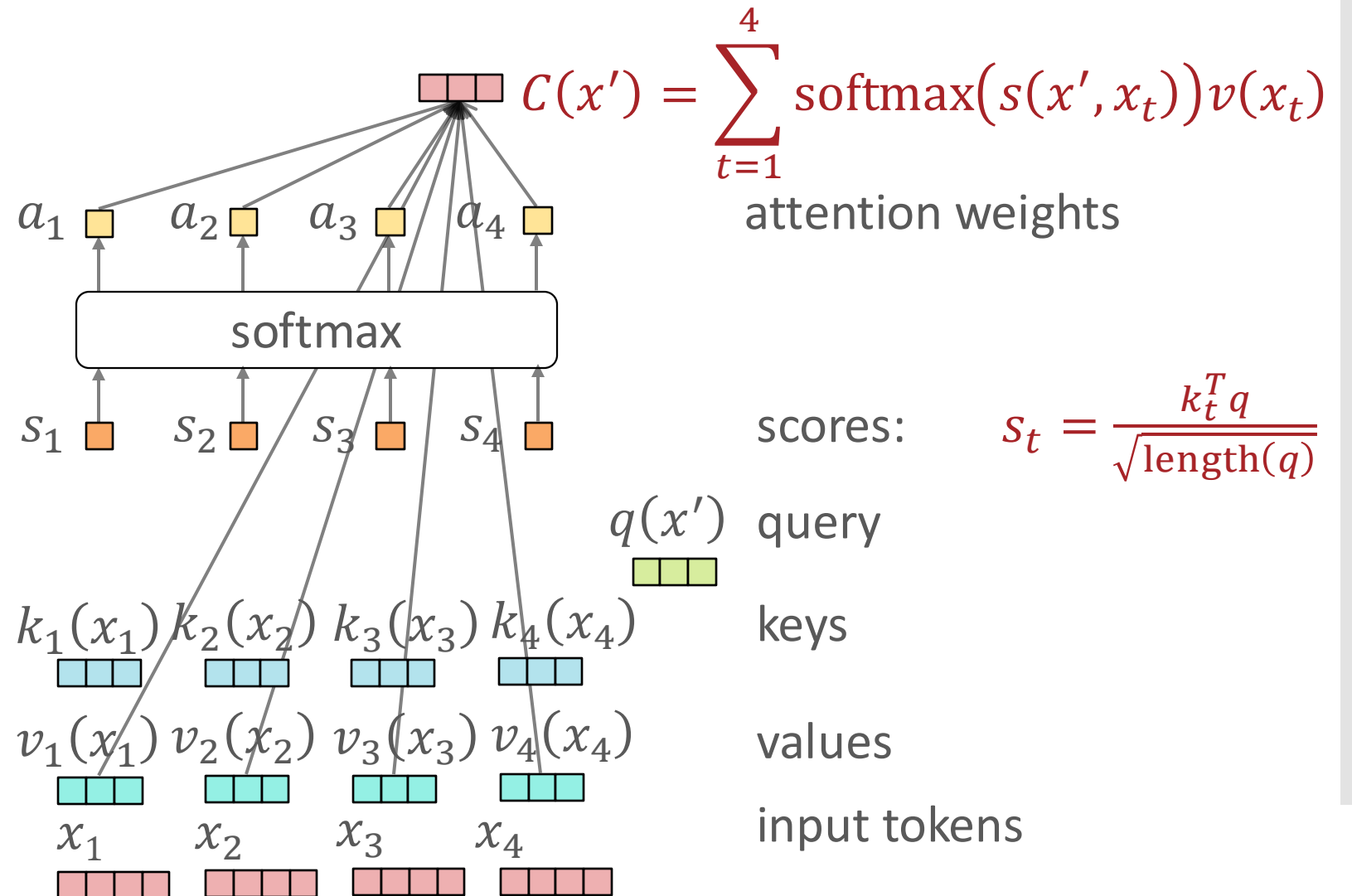
# Attention

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



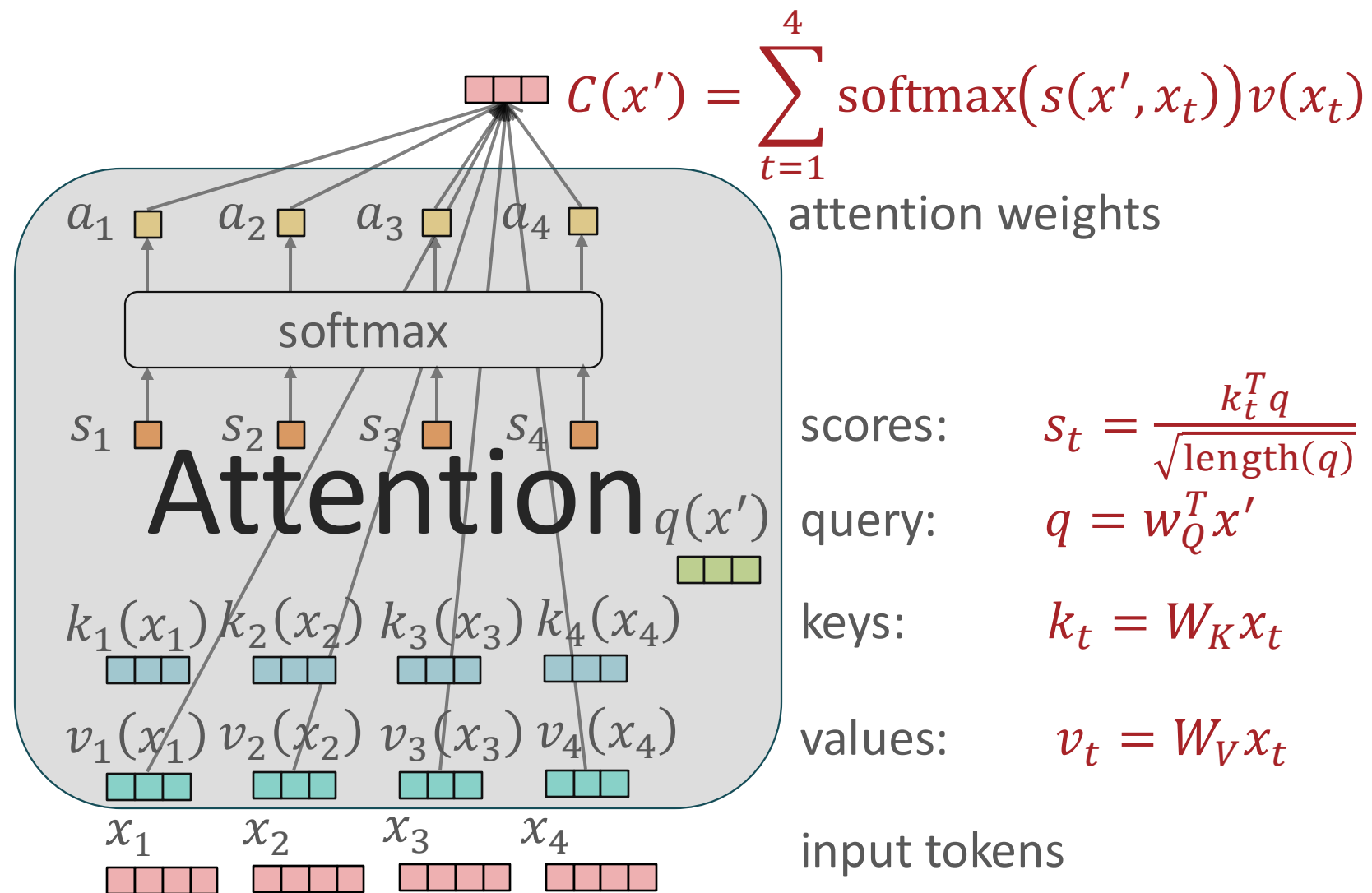
# Scaled Dot-product Attention

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



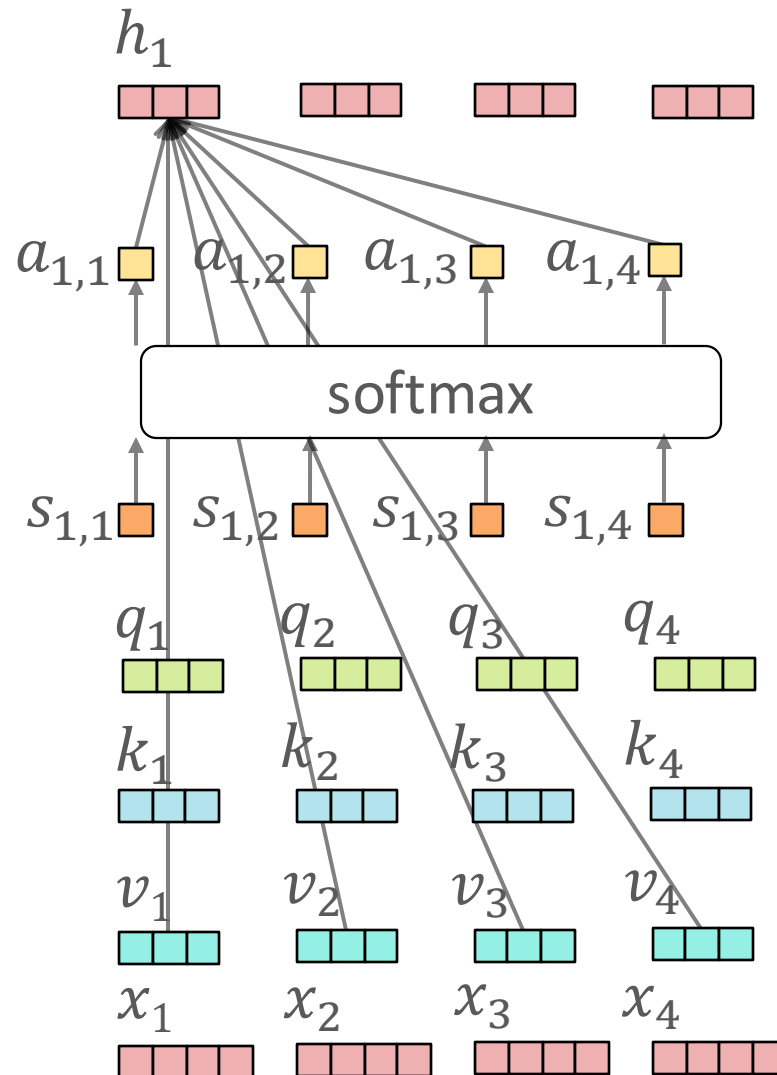
# Scaled Dot-product Attention

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



# 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

$$\text{scores: } s_{1,j} = \frac{k_j^T q_1}{\sqrt{\text{length}(k_j)}}$$

$$\text{queries: } q_t = W_Q x_t$$

$$\text{keys: } k_t = W_K x_t$$

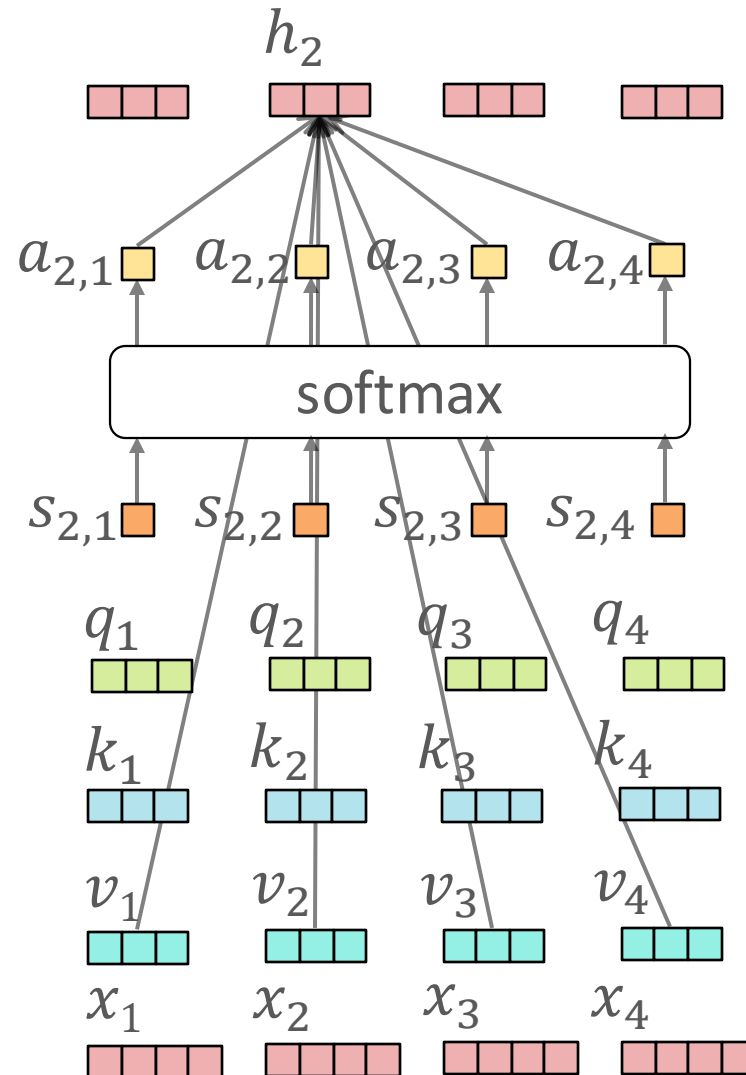
$$\text{values: } v_t = W_V x_t$$

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

$$\text{scores: } s_{2,j} = \frac{k_j^T q_2}{\sqrt{\text{length}(k_j)}}$$

$$\text{queries: } q_t = W_Q x_t$$

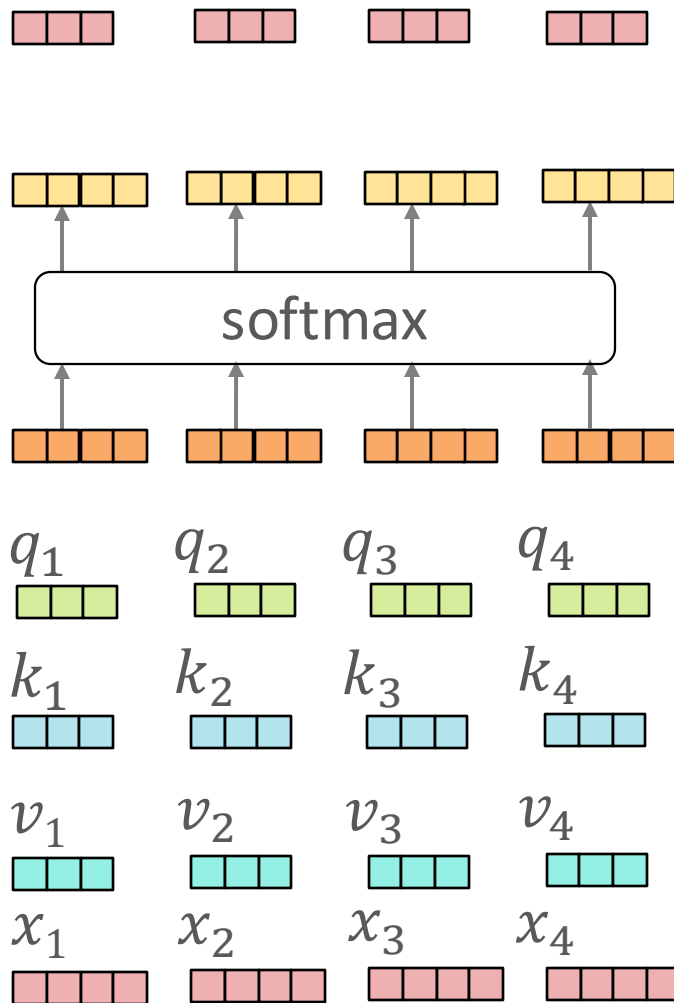
$$\text{keys: } k_t = W_K x_t$$

$$\text{values: } v_t = W_V x_t$$

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



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

attention weights

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

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

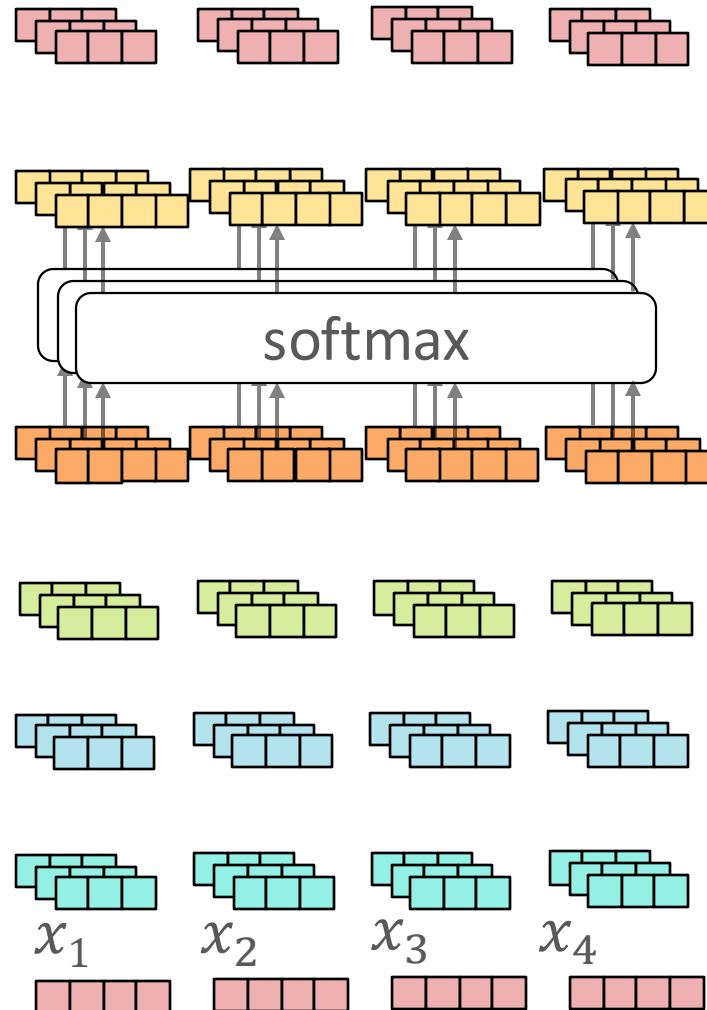
$$\text{keys: } K = XW_K \in \mathbb{R}^{N \times d_k}$$

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

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

# 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

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

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

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

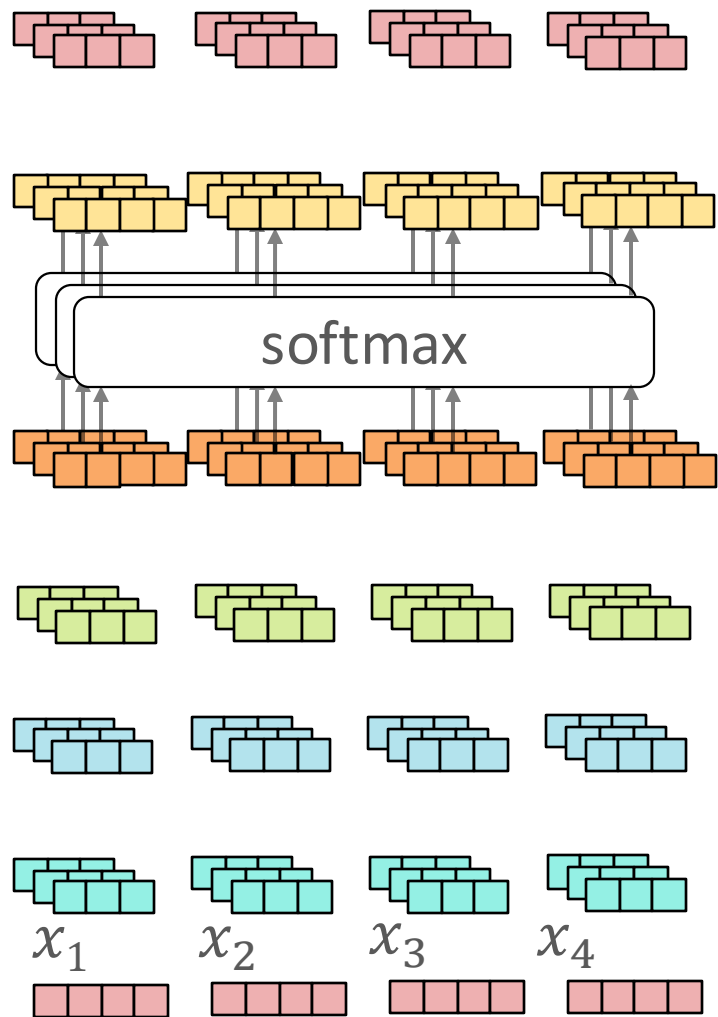
$$\text{values: } V^{(h)} = XW_V^{(h)}$$

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

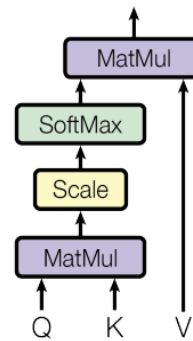
values: 
$$V^{(h)} = XW_V^{(h)}$$

design matrix:  $X$

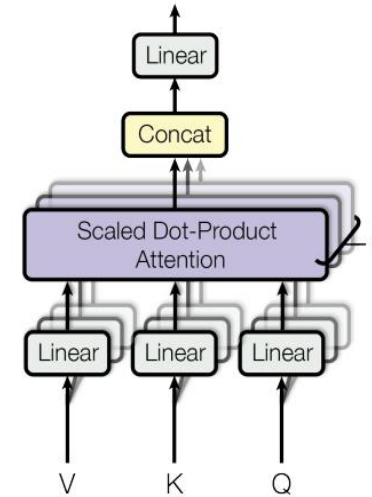
# Multi-head Scaled Dot-product Self-attention

- 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

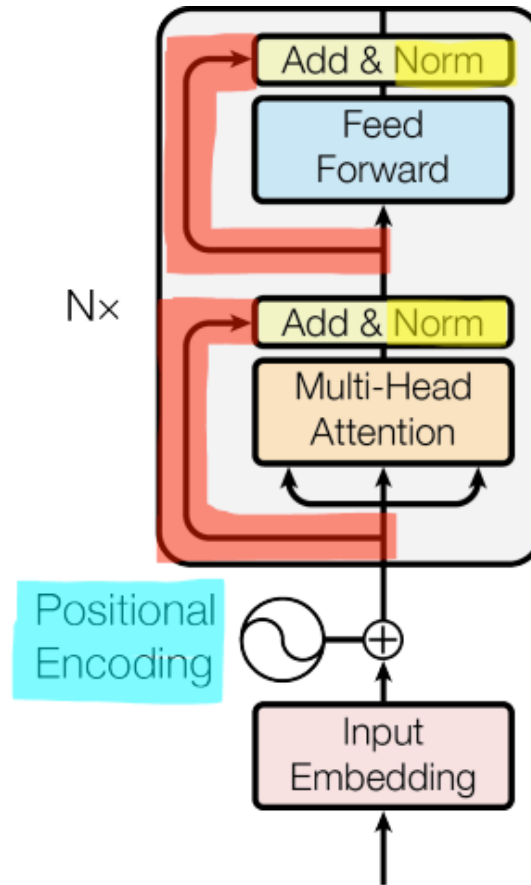


- The outputs from all the attention heads are concatenated together to get the final representation

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

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

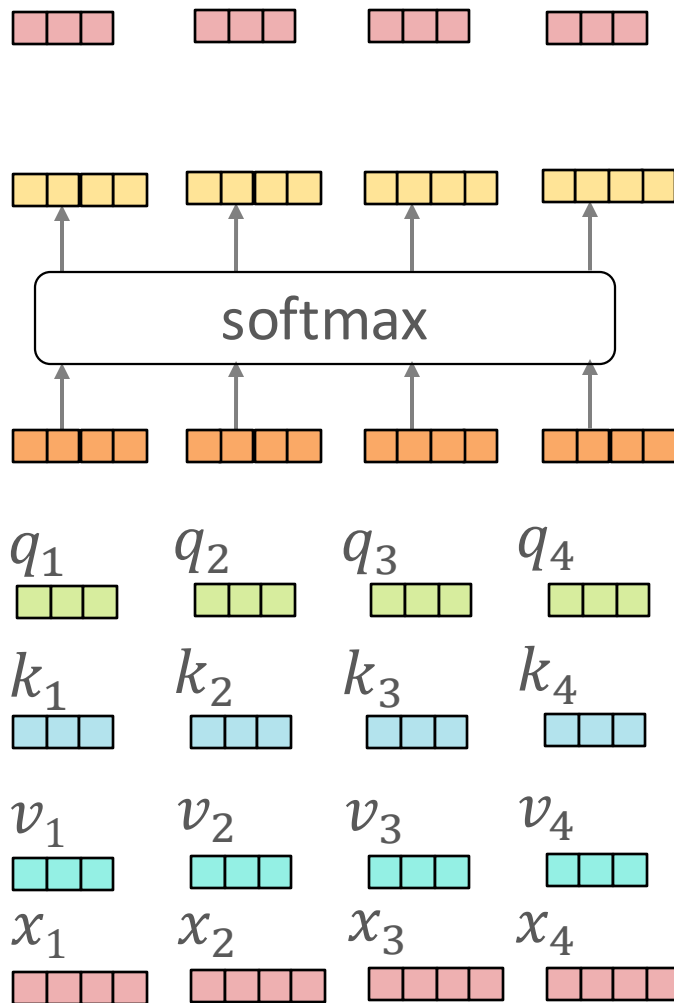
# Transformers



- 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



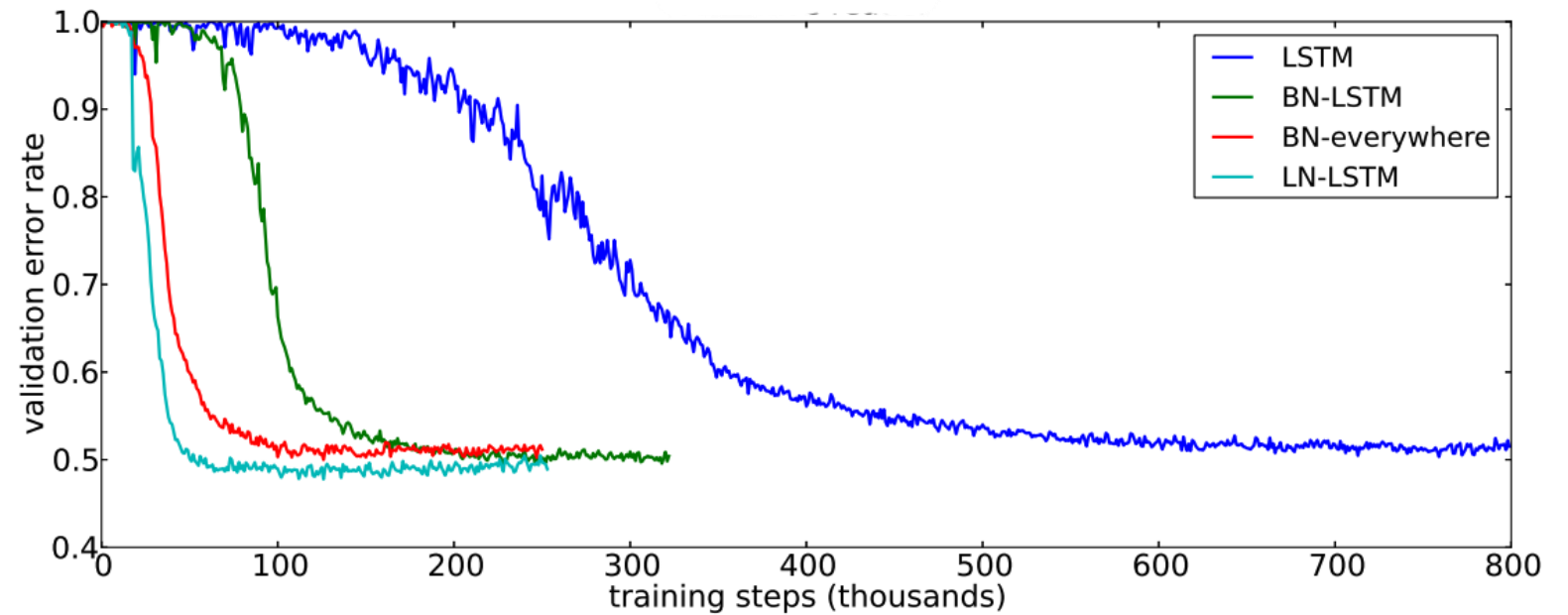
# Layer Normalization

- 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



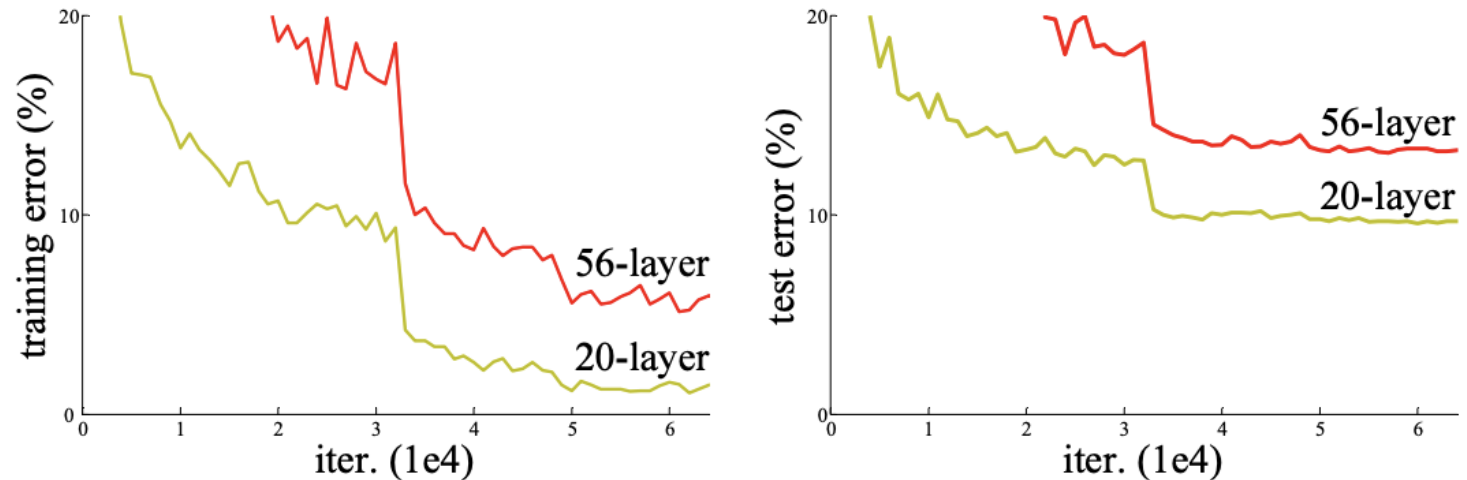
- Idea: normalize the output of a layer to always have the same (learnable) mean,  $\beta$ , and variance,  $\gamma^2$

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

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

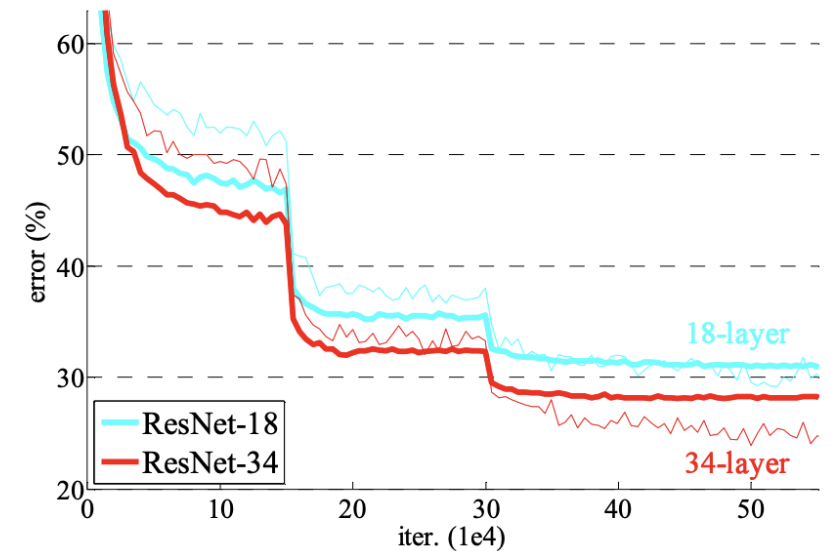
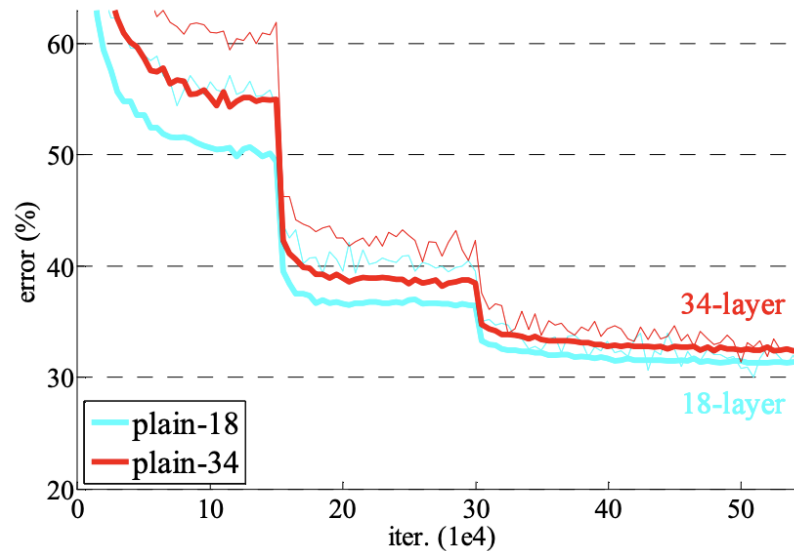
$$H' = H(x^{(i)}) + x^{(i)}$$

- 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