

Deep Learning for Natural Language Processing

Antoine Bosselut

EPFL



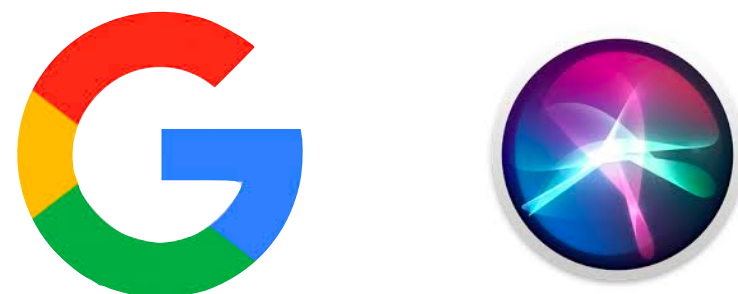
Natural Language Processing

Enabling Human-Machine Collaboration

Search Engines

Dialogue Agents

Text Generation

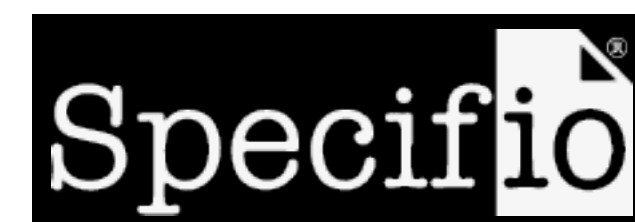


Accelerating Human-Human Communication

Machine Translation

Text Summarization

Information Extraction



Mining Human Insights

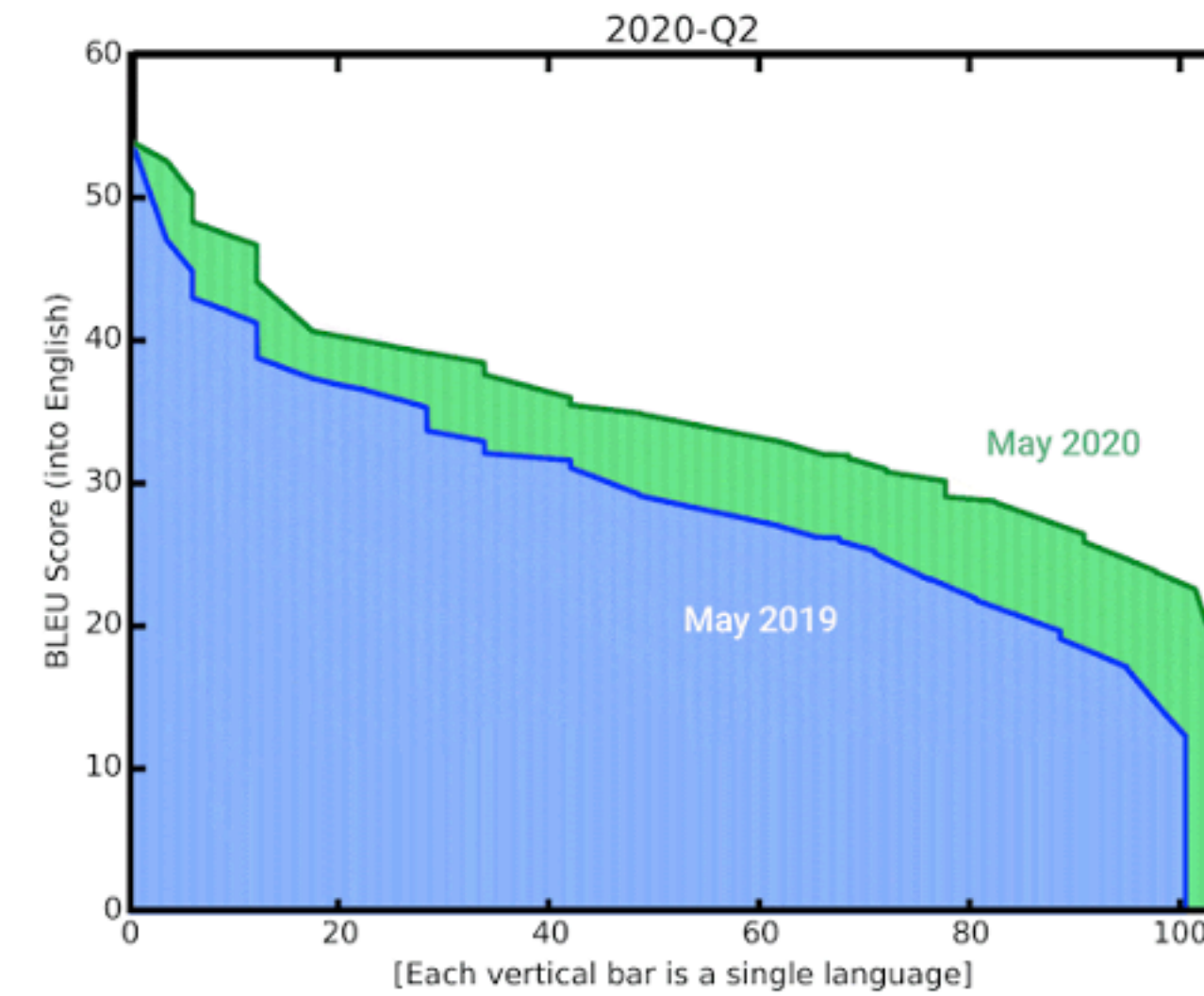
Sentiment Analysis

Motivation Analysis

Emotion Detection



Machine Translation



DETECT LANGUAGE **FRENCH** ENGLISH ▼ ↔ **ENGLISH** FRENCH SPANISH ▼

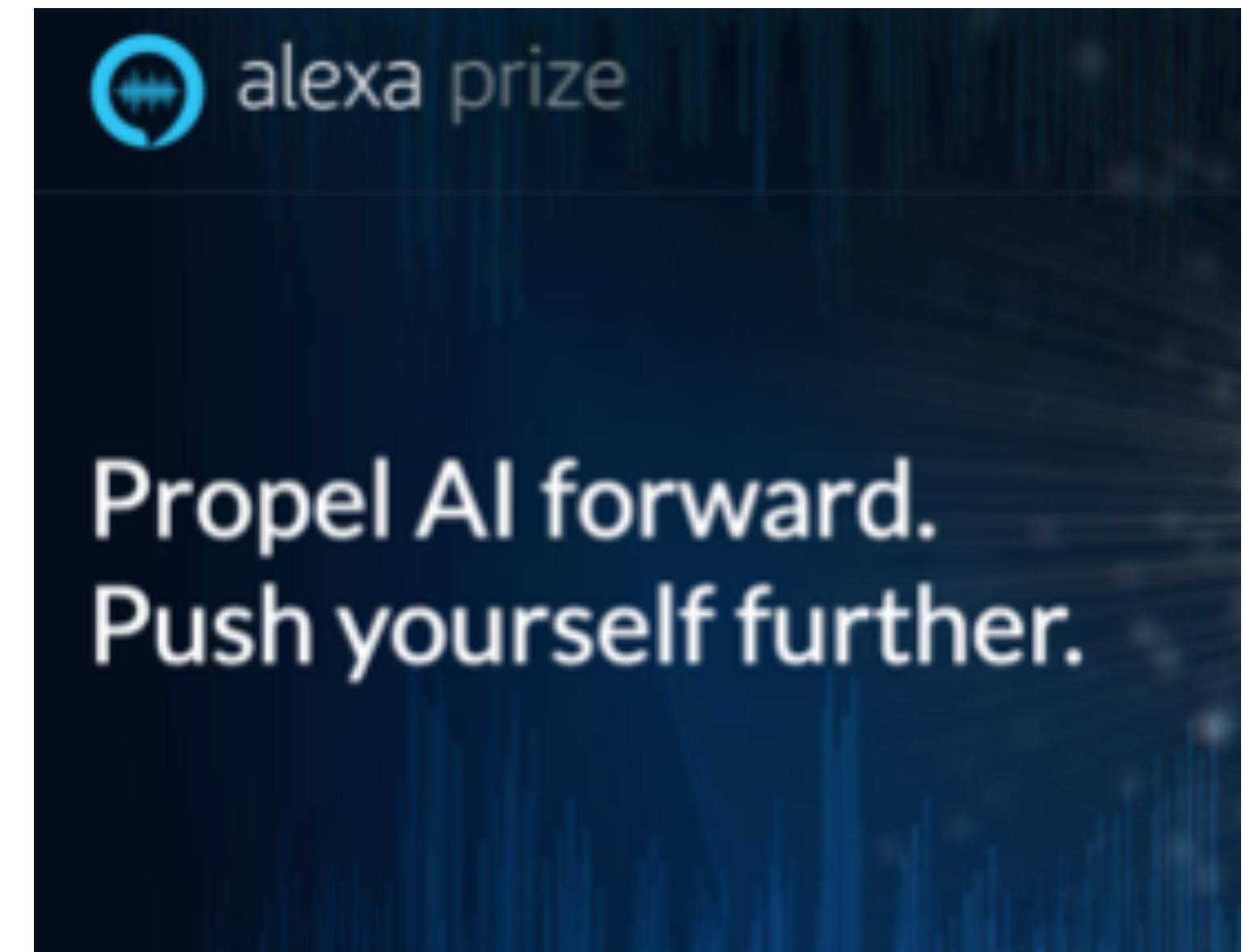
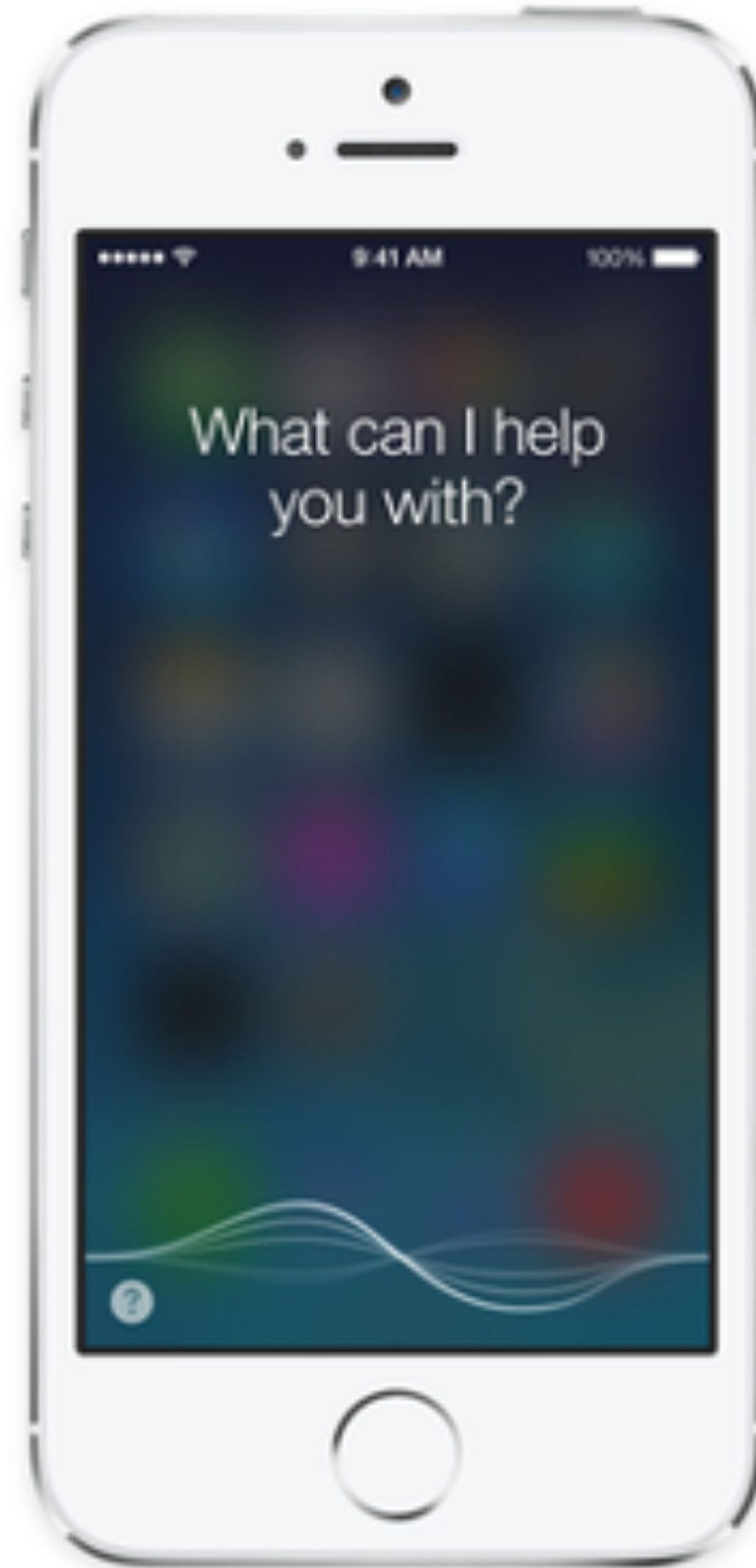
J'ai mangé avec mon avocat ×
aujourd'hui

I ate with my lawyer today ☆

38 / 5000

The screenshot shows the Google Translate interface. The source language is French and the target language is English. The input text is "J'ai mangé avec mon avocat aujourd'hui" and the output is "I ate with my lawyer today". The interface includes a microphone icon, a speaker icon, a character count of 38 / 5000, and a keyboard icon. The output text has a star icon and a speaker icon, along with copy, share, and refresh icons.

Conversational Systems



Question Answering

what is the tallest mountain in europe ?

All Images Maps News Videos More Tools

About 12'400'000 results (1.10 seconds)

Europe / Mountains / Maximal / Elevation

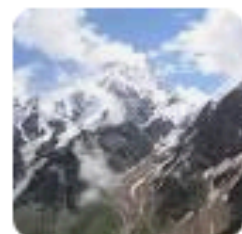
Mount Elbrus

5,642 m

Mountains



Mount Elbrus
5,642 m



Dykh-Tau
5,205 m



Shkhara
5,193 m



Koshtan-Tau
5,151 m



Mount Kazbek
5,033 m



Alps



Mont Blanc
4,810 m

who is the president of switzerland ?

All Images News Maps Videos More Tools

About 415'000'000 results (0.82 seconds)

Guy Parmelin

The President of the Swiss Confederation in 2021 is **Guy Parmelin** from the canton of Vaud. He was elected on 9 December 2020. The President's department in 2021 is the Federal Department of Economic Affairs, Education and Research EAER.

<https://www.admin.ch> > gov > start > federal-presidency

Presidential year 2021

About featured snippets Feedback

Text Generation

AN

Can you give me a two-day itinerary for a weekend trip to Paris ? Please don't include more than three activities per day.



Day 1:

1. Visit the Eiffel Tower - Start your day with the iconic symbol of Paris, the Eiffel Tower. Take a stroll around the Champ de Mars and enjoy the views from the tower.
2. Explore the Louvre Museum - After lunch, head to the Louvre Museum and admire the world-famous artwork such as the Mona Lisa and the Venus de Milo.
3. Stroll along the Seine River - In the evening, take a romantic stroll along the Seine River and enjoy the Parisian atmosphere.

Day 2:

1. Visit Notre-Dame Cathedral - Start your day by visiting one of the most famous Gothic cathedrals in the world, Notre-Dame Cathedral.
2. Explore the Marais neighborhood - After lunch, explore the Marais neighborhood, known for its historical architecture, trendy boutiques, and restaurants.
3. Enjoy a Cabaret Show - In the evening, enjoy a cabaret show at the famous Parisian venue, the Moulin Rouge, and end your trip with a memorable experience.

Next few weeks!

- **Today:** Deep Learning for Natural Language Processing
- **In two weeks:** Neural Text Generation
- **Final week:** Modern NLP & Ethical Implementation of NLP

Today's Outline

- **Introduction**
- **Section 1** - Neural NLP & Word Embeddings
- **Section 2** - Recurrent Neural Networks for Sequence Modeling
- **Section 3** - Attentive Neural Modeling with Transformers
- **Exercise Session:** Attention in Transformer Language Models

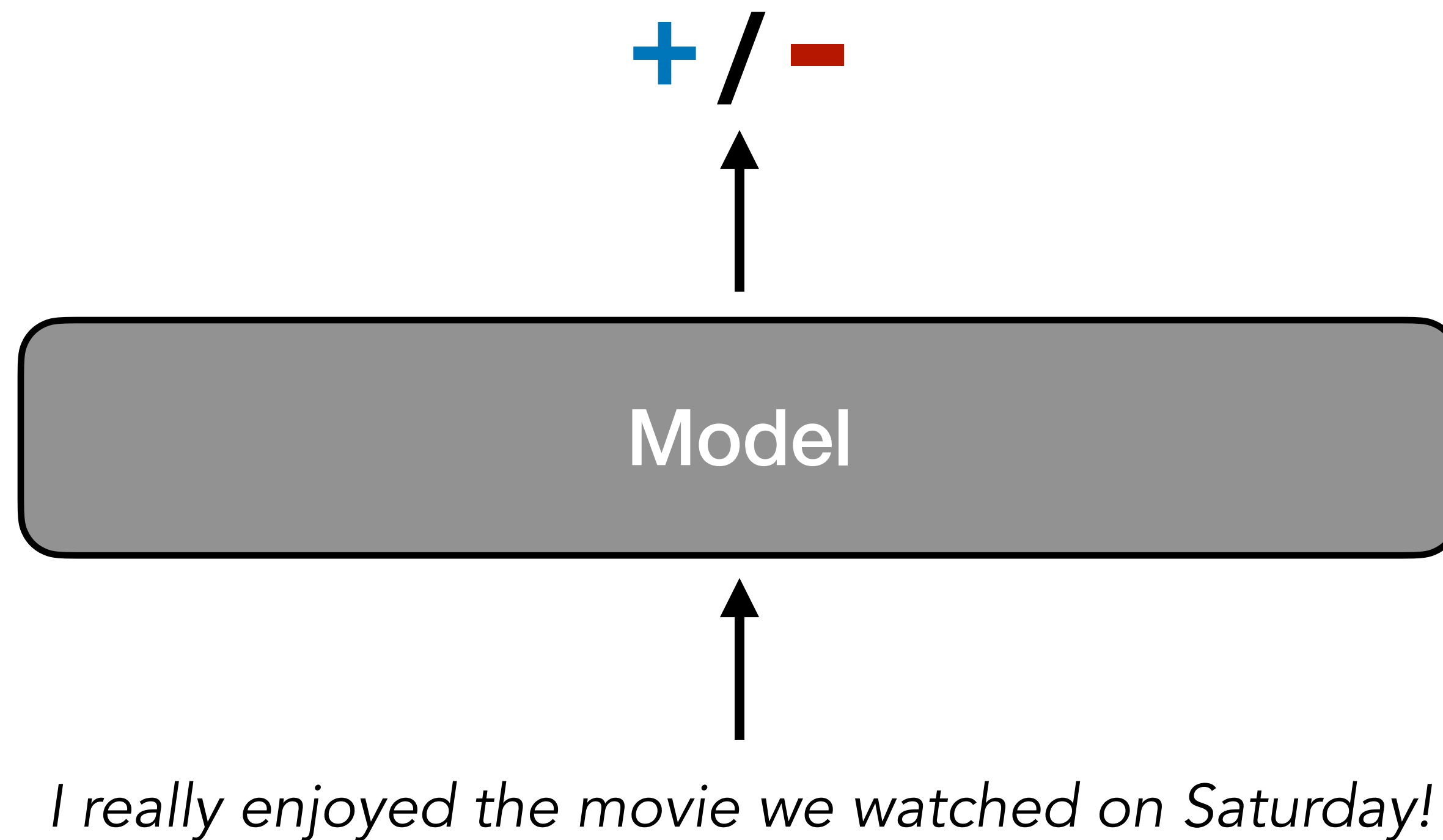
Part 1: Neural Embeddings

Section Outline

- **New:** Building our first neural classifier
- **Review:** sparse word vector representations
- **New:** Learning dense word vector representations - CBOW & Skipgram

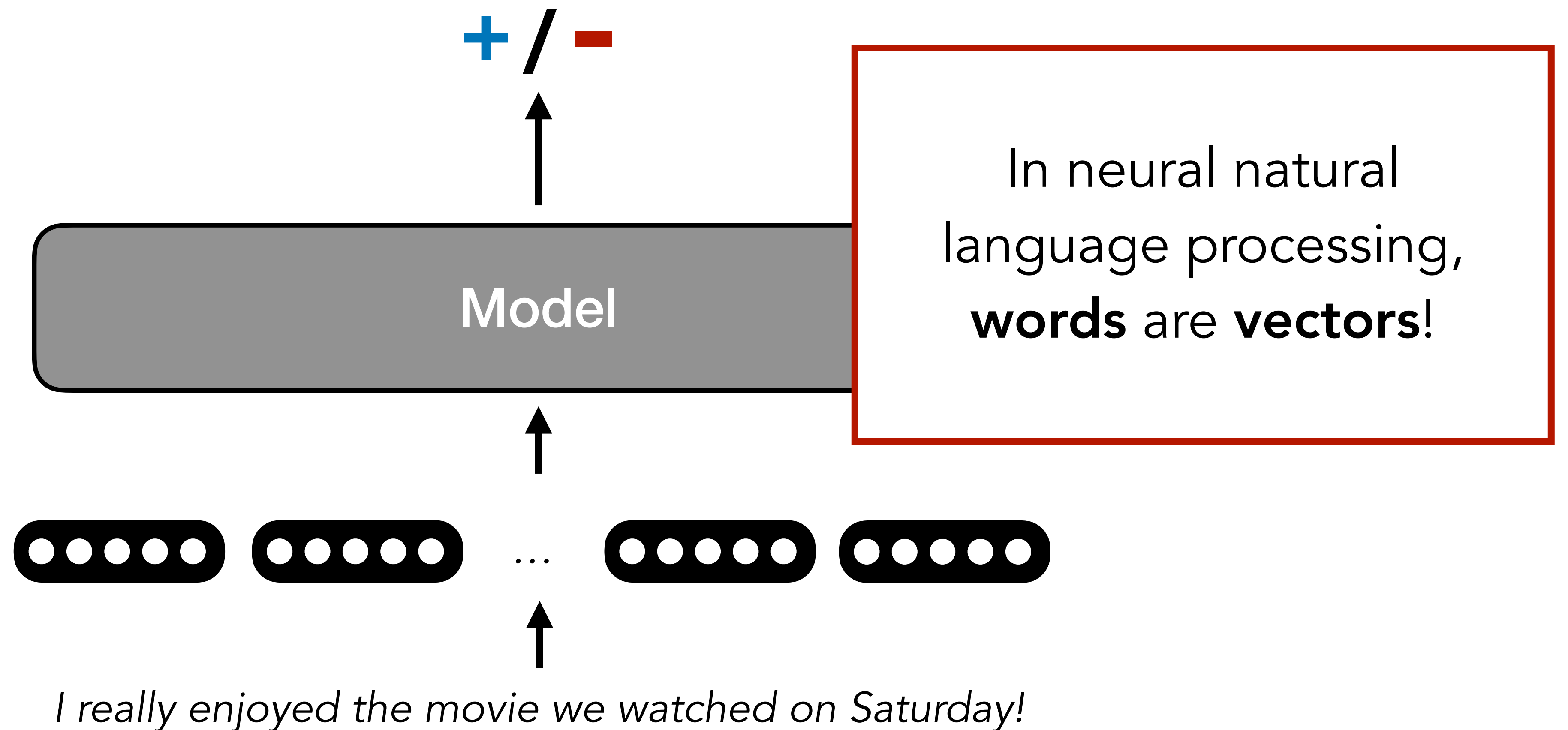
A simple NLP model

- How do we represent natural language sequences for NLP problems?



A simple NLP model

- How do we represent natural language sequences for NLP problems?



Question

What words should we model as vectors?

Choosing a vocabulary

- Language contains many words (e.g., ~600,000 in English)
 - **What about other tokens:** Capitalisation? Accents ? Typos!? Words in other languages!? In other scripts!? Emojis !? Unicode !?
 - **Millions of potential unique tokens!** Most rarely appear in our training data (Zipfian distribution)
 - Model has limited capacity

Choosing a vocabulary

- Language contains many words (e.g., ~600,000 in English)
 - **What about other tokens:** Capitalisation? Accents ? Typos!? Words in other languages!? In other scripts!? Emojis !? Unicode !?
 - **Millions of potential unique tokens!** Most rarely appear in our training data (Zipfian distribution)
 - Model has limited capacity
- How should we select which tokens we want our model to process?
 - CS-552: Modern NLP Week 13 - Tokenisation!
 - For now, initialize a vocabulary V of tokens that we can represent as a vector
 - Any token not in this vocabulary V is mapped to a special $\langle \text{UNK} \rangle$ token (e.g., unknown).

Question

How should we model a word as a vector?

Sparse Word Representations

$$w_i \in \{0,1\}^V$$

- Define a vocabulary V
- Each word in the vocabulary is represented by a sparse vector
- Dimensionality of sparse vector is size of vocabulary (e.g., thousands, possibly millions)

<i>I</i>	→	[0 ... 0 0 0 1 ... 0 0]
<i>really</i>	→	[0 ... 1 ... 0 0 0 0 0]
<i>enjoyed</i>	→	[0 ... 0 0 0 1 0 ... 0]
<i>the</i>	→	[0 ... 0 1 0 0 0 ... 0]
<i>movie</i>	→	[0 ... 0 0 0 0 0 ... 1]
<i>!</i>	→	[1 ... 0 0 0 0 0 0 0 0]

Word Vector Composition

- To represent sequences, beyond words, define a composition function over sparse vectors

I really enjoyed the movie ! → [1 ... 1 1 0 1 ... 0 1] **Simple Counts**

I really enjoyed the movie ! → [0.01 ... 0.1 0.1 0 0.001 ... 0 0.5] **Weighted by Corpus Statistics (e.g., TF-IDF)**

Many others...

Problem

With sparse vectors, similarity is a function of common words!

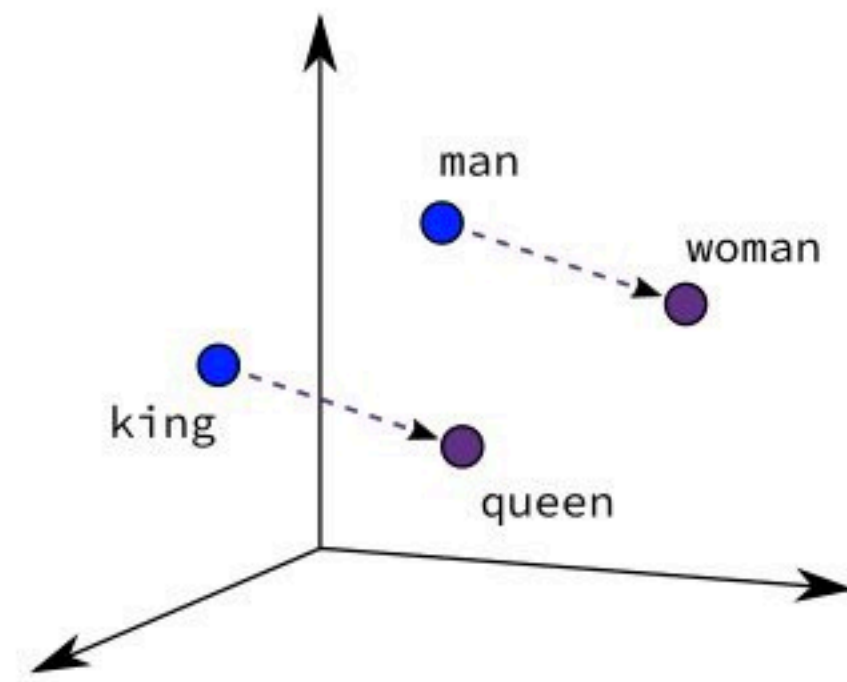
How do you learn learn similarity between words?

enjoyed → [0 ... 0 0 0 1 ... 0 0]

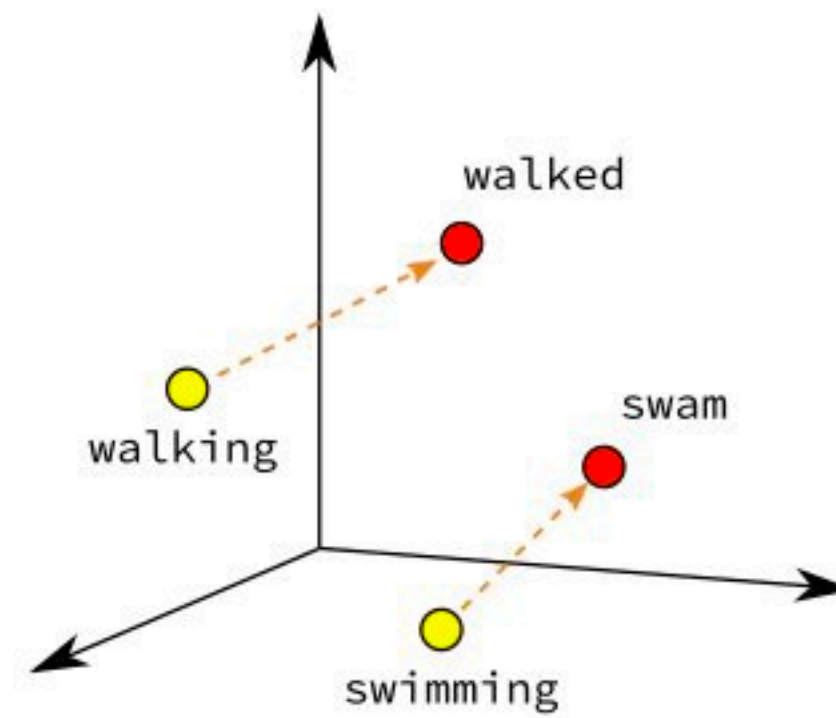
loved → [0 ... 1 ... 0 0 0 0 0]

$$\text{sim}(\textit{enjoyed}, \textit{loved}) = 0$$

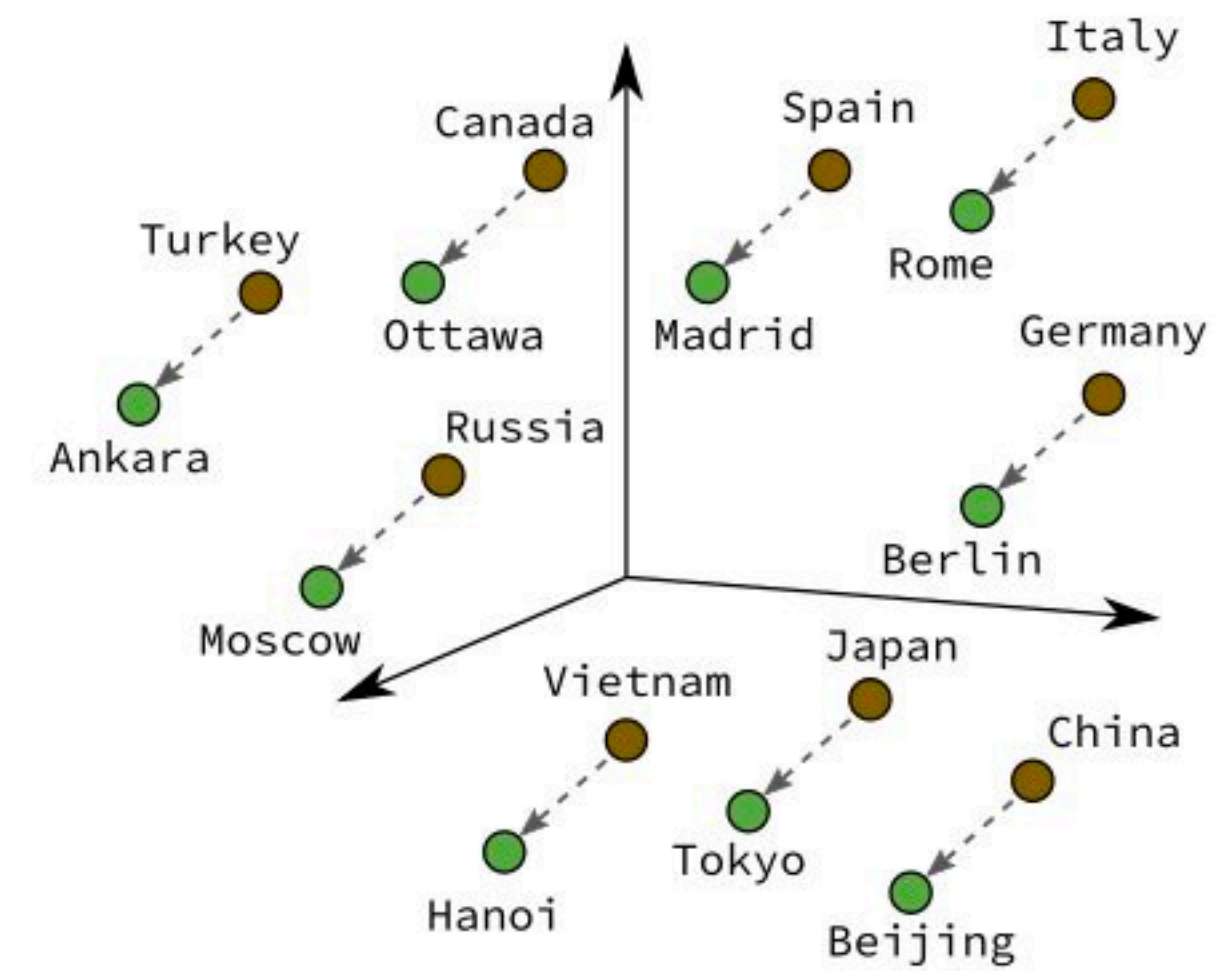
Embeddings Goal



Male-Female



Verb Tense



Country-Capital

How do we train semantics-encoding embeddings of words?

Dense Word Vectors

- Represent each word as a high-dimensional*, **real-valued** vector
 - *Low-dimensional compared to V-dimension sparse representations, but still usually $O(10^2 - 10^3)$

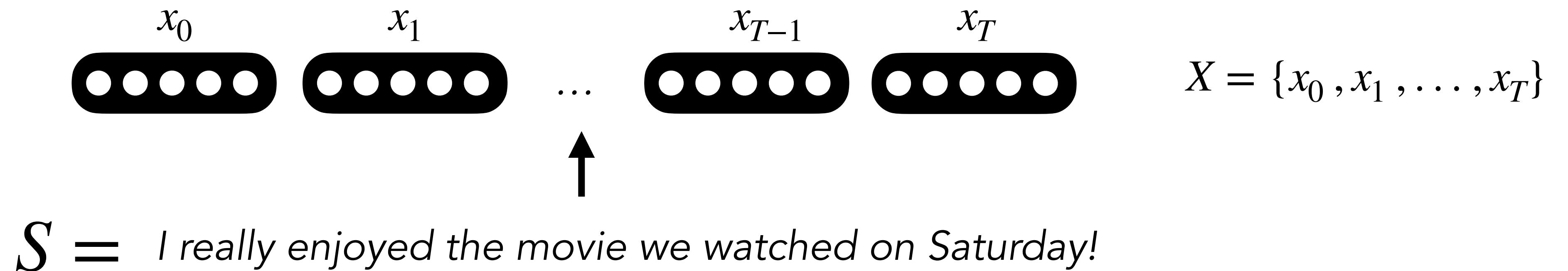
<i>I</i>	→	[0.113 -0.782 1.893 0.984 6.349 ...]
<i>really</i>	→	[0.906 0.661 -0.214 -0.894 -0.880 ...]
<i>enjoyed</i>	→	[-0.842 0.647 -0.882 0.045 0.029 ...]
<i>the</i>	→	[0.100 0.765 -0.333 -0.538 -0.150 ...]
<i>movie</i>	→	[0.104 -0.054 -0.268 -0.877 0.005 ...]
<i>!</i>	→	[0.439 -0.577 -0.727 0.261 0.699 ...]

word vectors
word embeddings
neural embeddings
dense embeddings
others...

- Similarity of vectors represents similarity of meaning for particular words

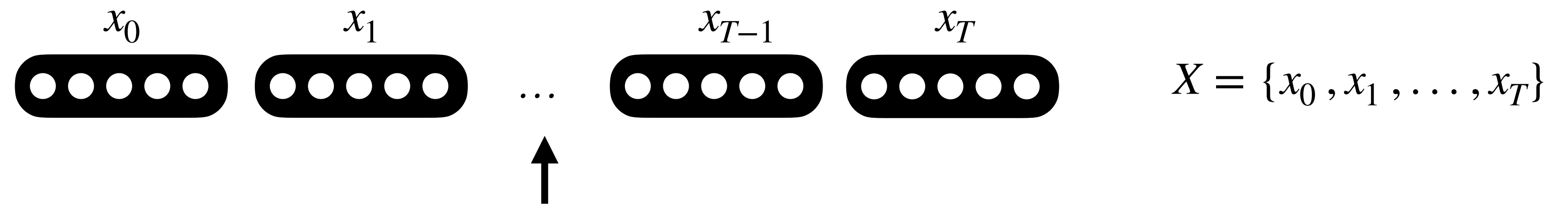
A simple NLP model

- For each sequence S , we have a corresponding sequence of embeddings X



A simple NLP model

- For each sequence S , we have a corresponding sequence of embeddings X



$S_1 =$ *I **really** enjoyed the movie **we** watched on Saturday !*

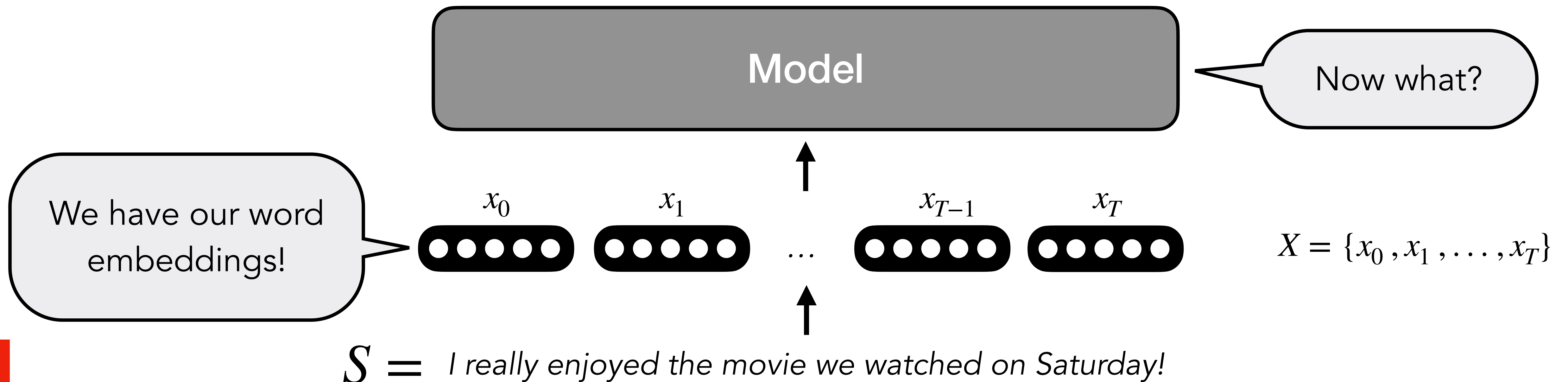
- Embeddings $x_t \in X$ are indexed from shared embedding dictionary \mathbb{E} for all items in vocabulary V

$S_2 =$ ***We** really loved a film **we** saw last Sunday !*

Bolded words would index the same embedding in \mathbb{E}

A simple NLP model

- For each sequence S , we have a corresponding sequence of embeddings X

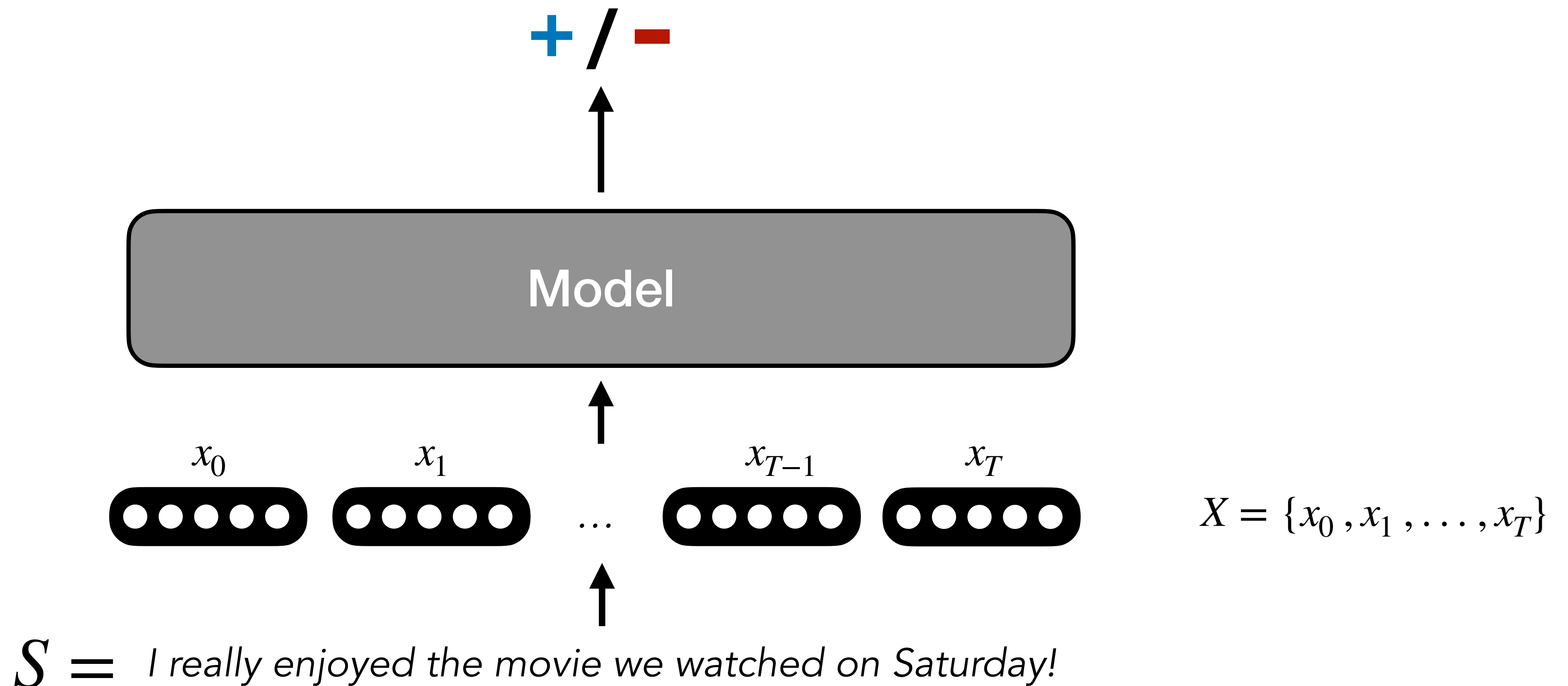


Question

What should we use as a model?

A simple NLP model

- Our model modifies and / or composes these word embeddings to formulate a representation that allows it to predict the correct label



A simple NLP model

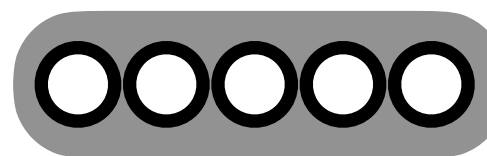
- Our model modifies and / or composes these word embeddings to formulate a representation that allows it to predict the correct label
 - Recurrent neural networks (RNNs) - Today!
 - RNN variants (LSTM, GRU, etc.) - Today!
 - Transformer - Today!

A simple NLP model

Notation: Typically, we represent the output of a model as h (or o).

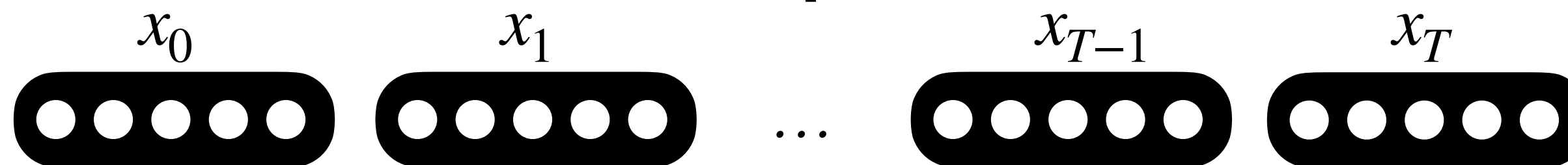
$$h_T = \sum_{t=0}^T x_t$$

+ / -



We composed our embeddings into a different representation!

Sum-pool



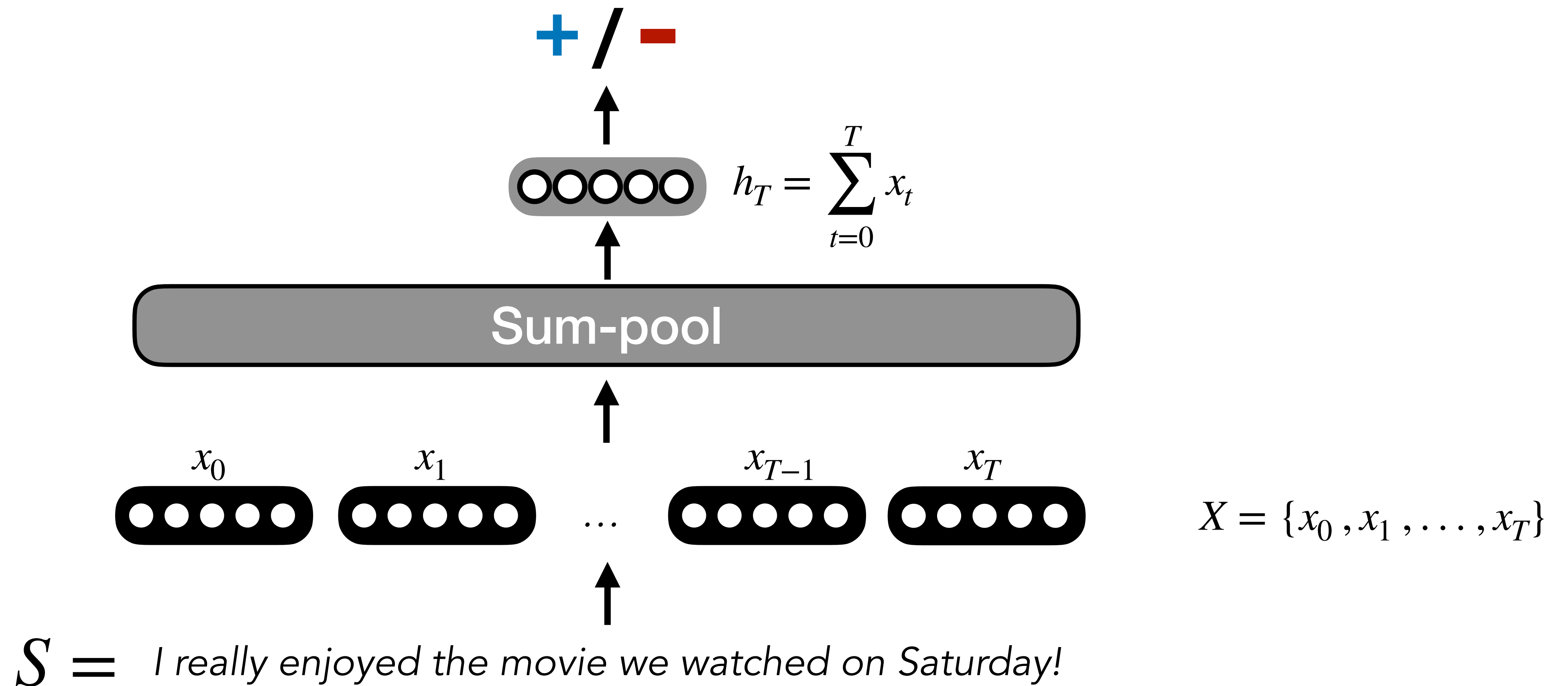
$$X = \{x_0, x_1, \dots, x_T\}$$

$S =$ *I really enjoyed the movie we watched on Saturday!*

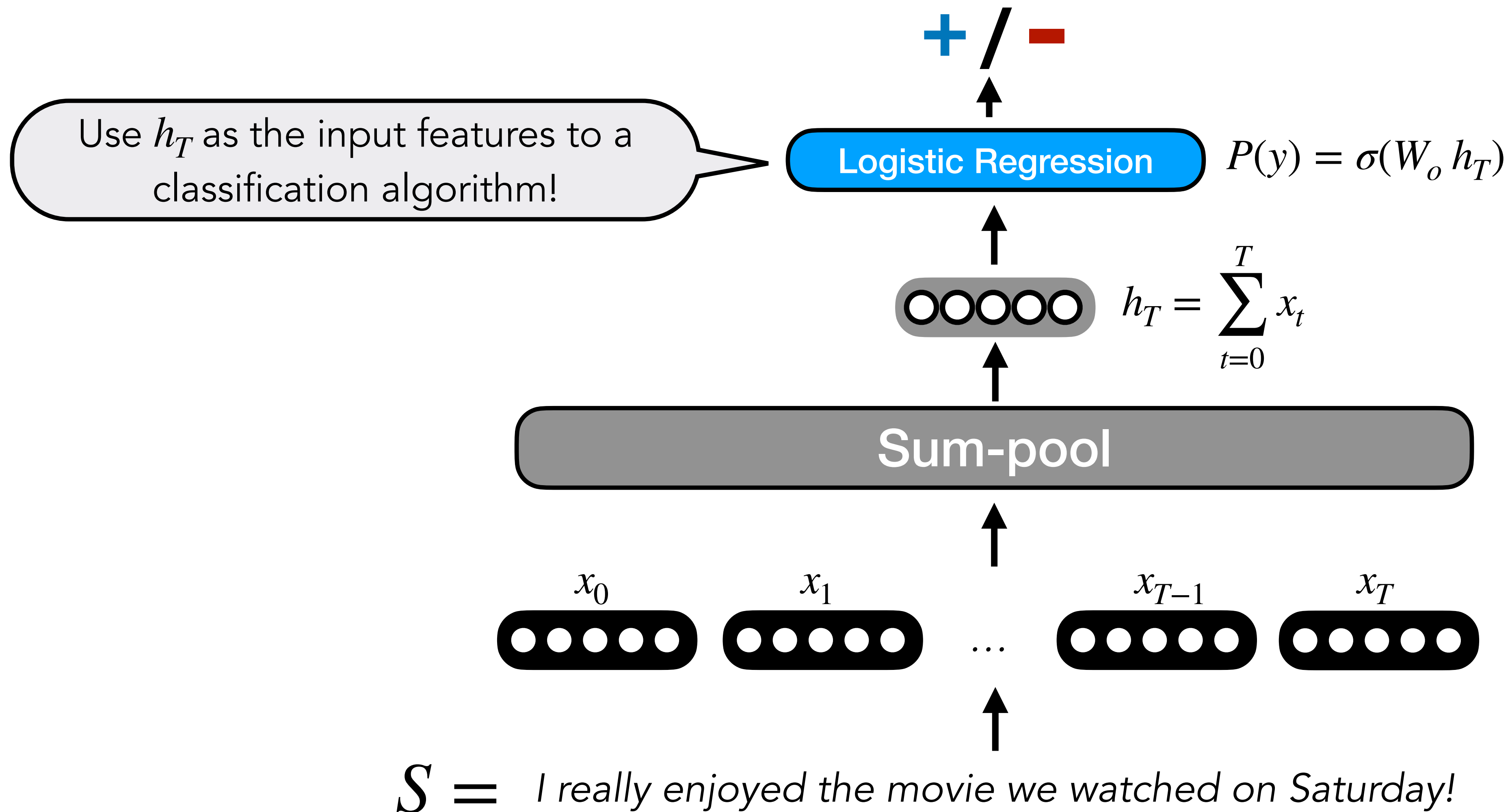
Question

How do we convert the output of our model to a prediction?

Predicting the label



Predicting the label



Learn using **backpropagation**: compute gradients of loss with respect to initial embeddings X

Learn embeddings that allow you to do the task successfully!

$$X = \{x_0, x_1, \dots, x_T\}$$

Question

What could be a better way to learn word embeddings?

“You shall know a word by the company it keeps”

–J.R. Firth, 1957

Context Representations

Solution:

Rely on the context in which words occur to learn their meaning

Context is the **set of words** that occur **nearby**

I really enjoyed the _____ we watched on Saturday!

The _____ growled at me, making me run away.

I need to go to the _____ to pick up some dinner.

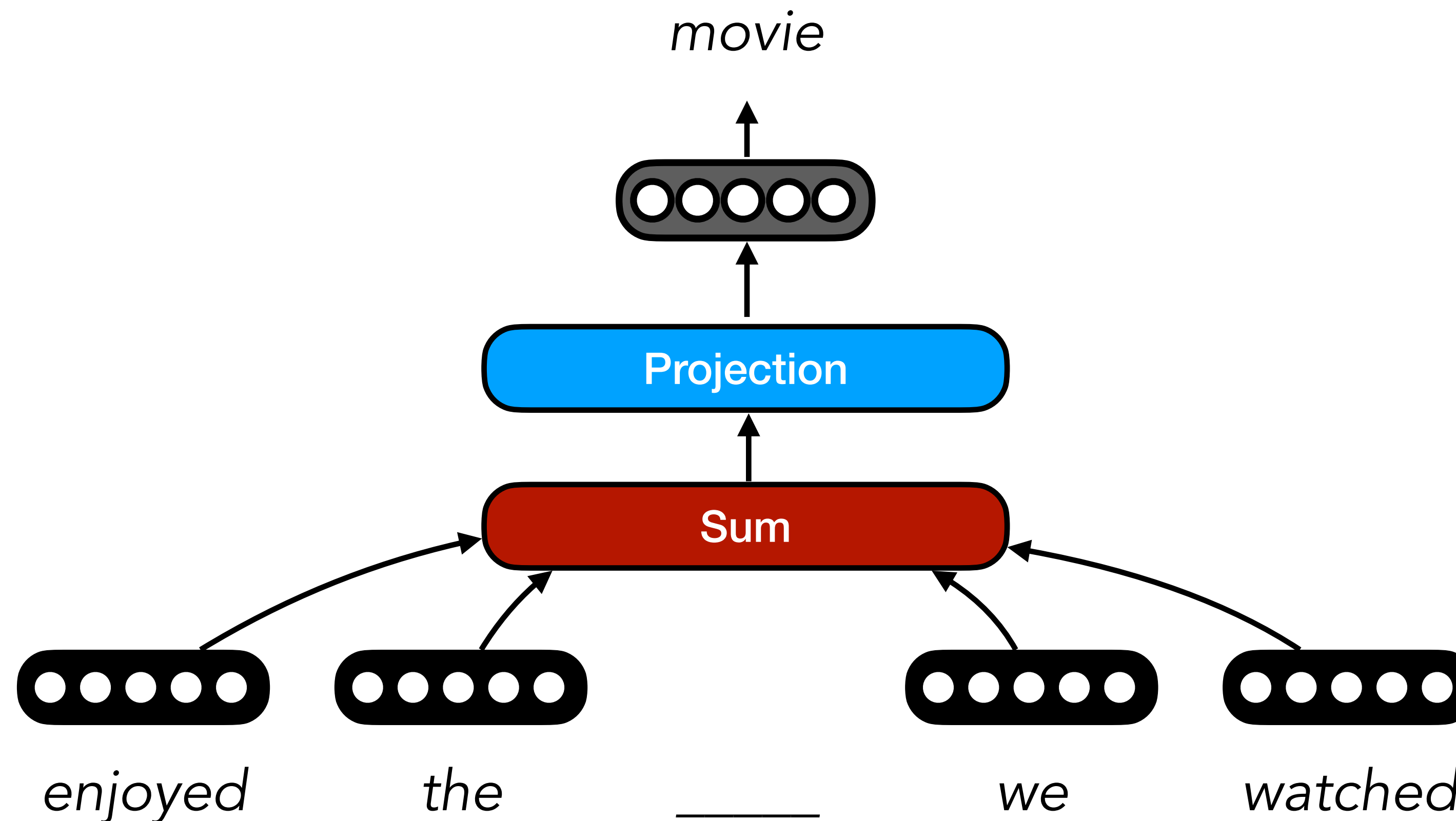
Foundation of **distributional semantics**

Learning Word Embeddings

- Many options, huge area of research, but three common approaches
- **Word2vec - Continuous Bag of Words (CBOW)**
 - Learn to predict missing word from surrounding window of words
- **Word2vec - Skip-gram**
 - Learn to predict surrounding window of words from given word
- **GloVe**
 - Not covered today

Continuous Bag of Words (CBOW)

- Predict the missing word from a window of surrounding words



Context:

enjoyed

the

we

watched

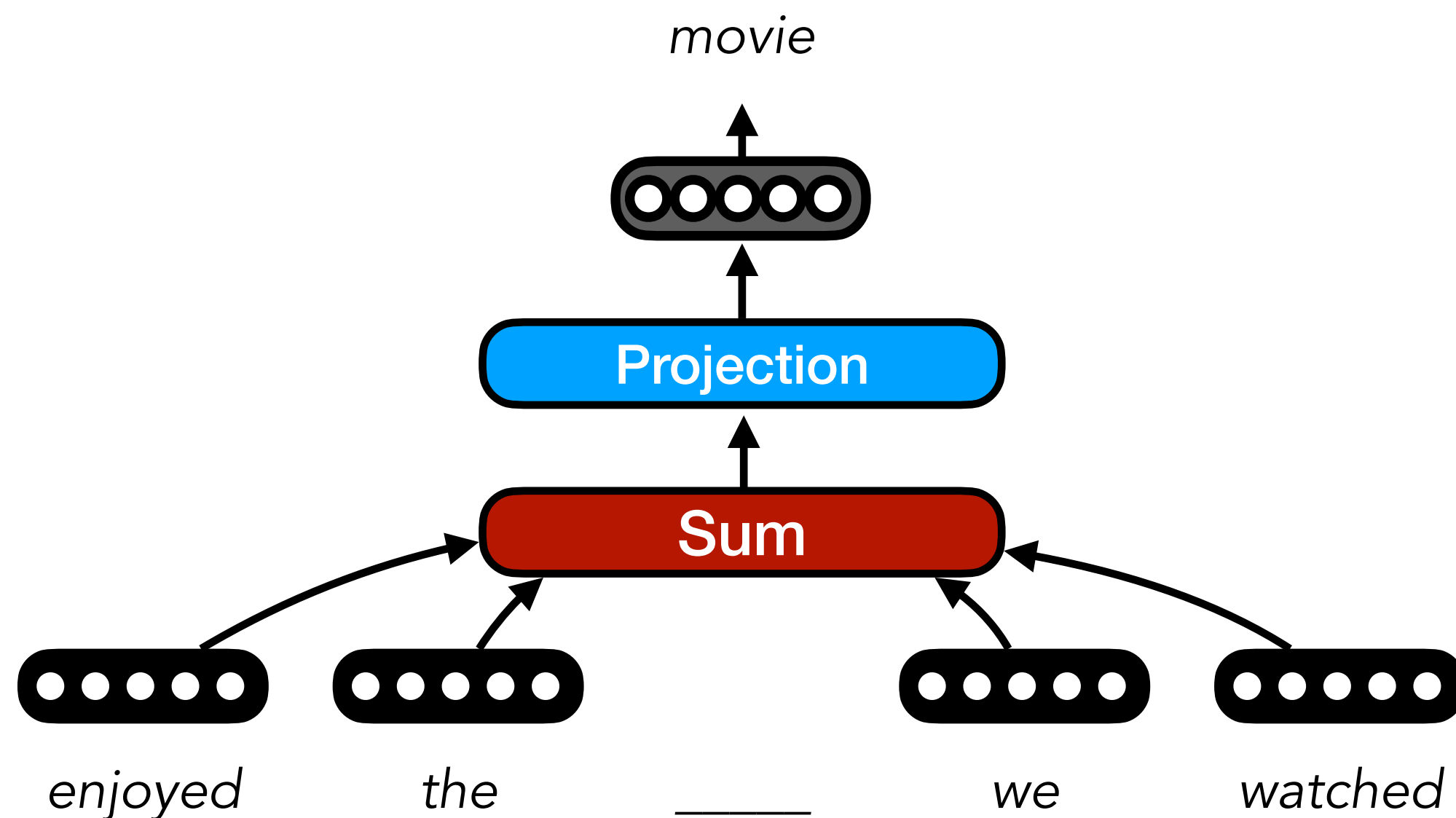
Continuous Bag of Words (CBOW)

- Predict the missing word from a window of surrounding words

$$\max P(\text{movie} \mid \text{enjoyed}, \text{the}, \text{we}, \text{watched})$$

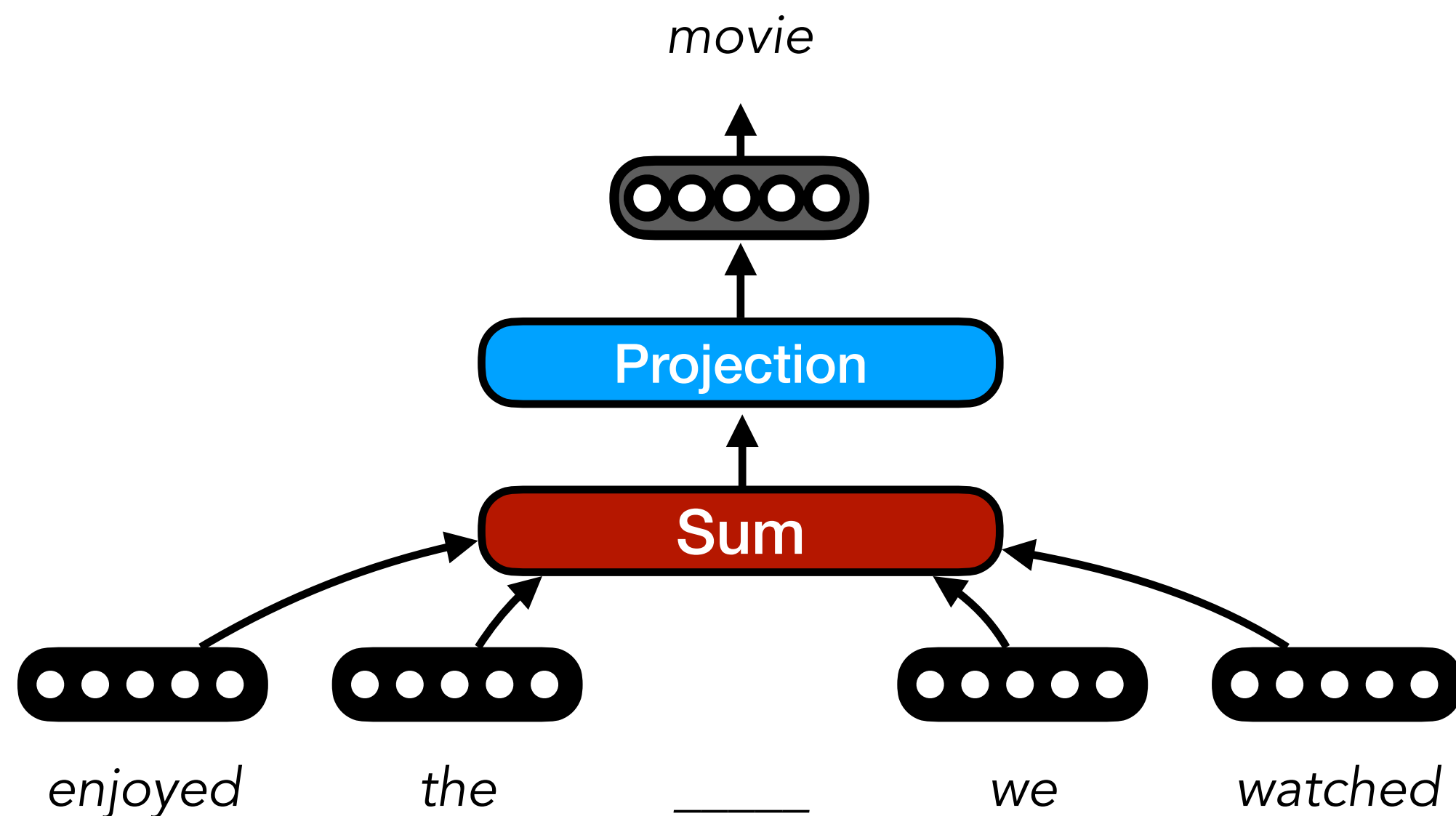
$$\max P(w_t \mid w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2})$$

$$\max P(w_t \mid \{w_x\}_{x=t-2}^{x=t+2})$$



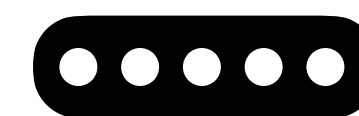
Continuous Bag of Words (CBOW)

- Predict the missing word from a window of surrounding words



$$P(w_t | \{w_x\}_{x=t-2}^{x=t+2}) = \mathbf{softmax} \left(\mathbf{U} \sum_{\substack{x=t-2 \\ x \neq t}}^{t+2} \mathbf{w}_x \right)$$

$$\mathbf{w}_x \in \mathbb{R}^{1 \times d}$$



$$\mathbf{U} \in \mathbb{R}^{d \times V}$$

Projection

Softmax Function

- The **softmax** function generates a probability distribution from the elements of the vector it is given

$$\mathbf{softmax}(\mathbf{a})_i = \frac{e^{a_i}}{\sum_{j=1}^{|\mathbf{a}|} e^{a_j}}$$

$\mathbf{V} = [0.790 \ -0.851 \ 0.506 \ 0.767 \ -0.788 \ 0.793 \ 0.887 \ 0.219 \ -0.052 \ 0.461]$

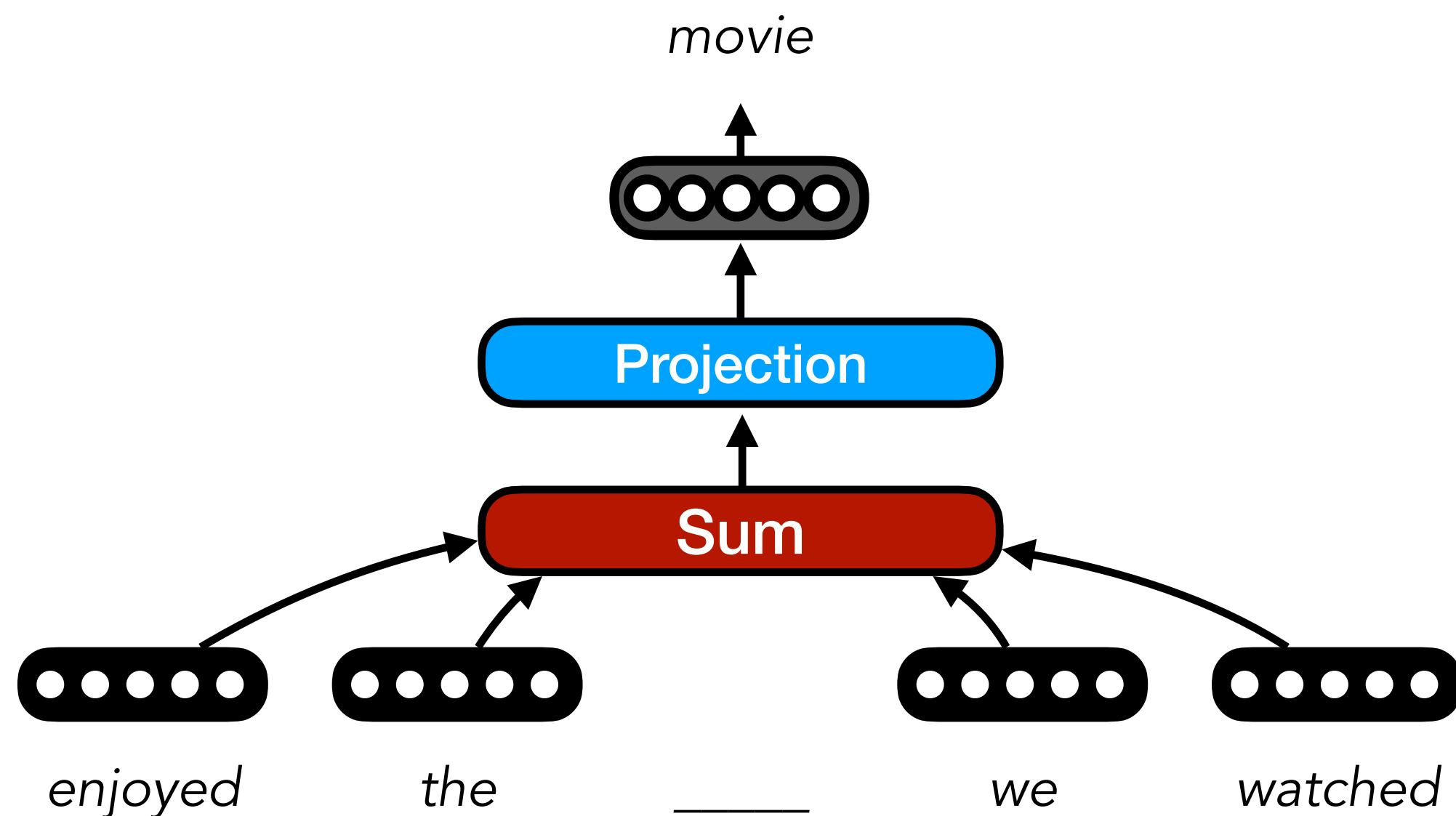


Softmax(V)

$\mathbf{P}(\mathbf{V}) = [0.144 \ 0.028 \ 0.108 \ 0.141 \ 0.030 \ 0.144 \ 0.159 \ 0.081 \ 0.062 \ 0.104]$

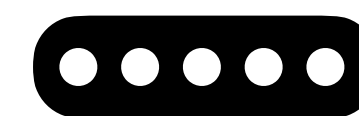
Continuous Bag of Words (CBOW)

- Predict the missing word from a window of surrounding words



$$P(w_t | \{w_x\}_{x=t-2}^{x=t+2}) = \mathbf{softmax} \left(\mathbf{U} \sum_{\substack{x=t-2 \\ x \neq t}}^{t+2} \mathbf{w}_x \right)$$

$$\mathbf{w}_x \in \mathbb{R}^{1 \times d}$$

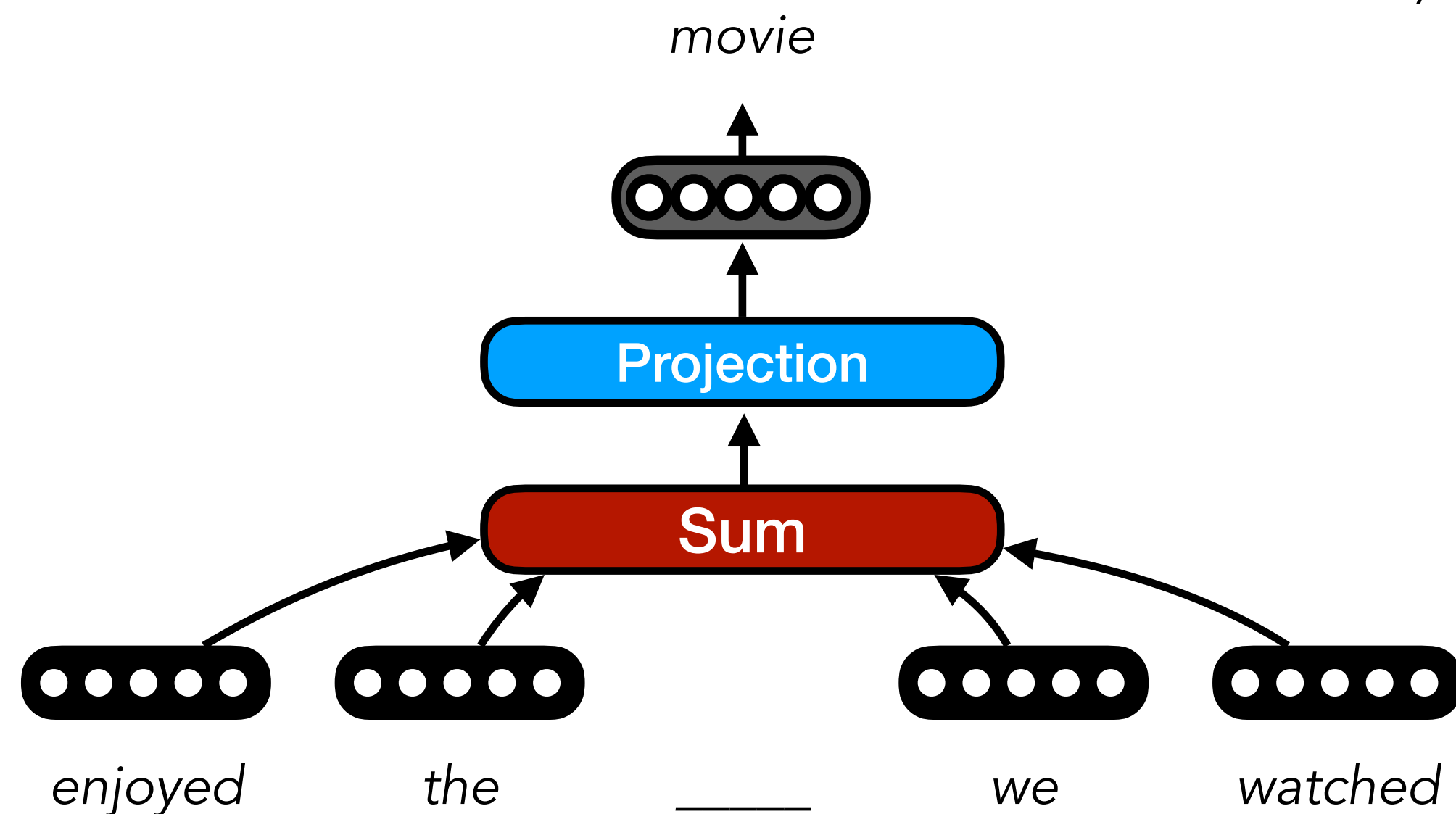


$$\mathbf{U} \in \mathbb{R}^{d \times V}$$

Projection

Continuous Bag of Words (CBOW)

$$P(w_t | \{w_x\}_{x=t-2}^{x=t+2}) = \mathbf{softmax} \left(\mathbf{U} \sum_{\substack{x=t-2 \\ x \neq t}}^{t+2} \mathbf{w}_x \right)$$

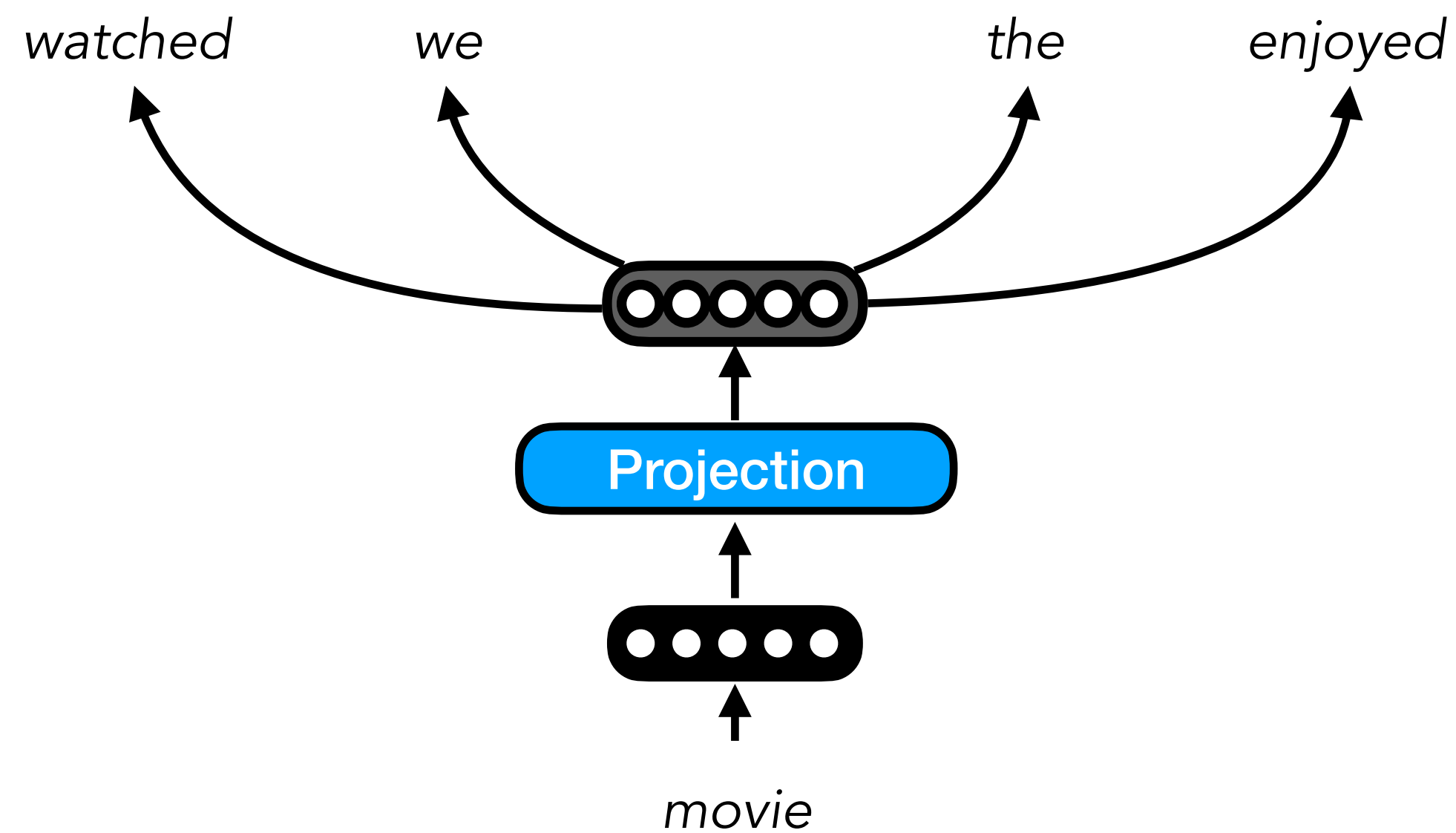


- Model is trained to **maximise** the **probability** of the missing word
 - For computational reasons, the model is typically trained to **minimise** the **negative log probability** of the missing word
- Here, we use a window of **N=2**, but the window size is a **hyperparameter**
- For computational reasons, a **hierarchical softmax** used to compute distribution

Skip-gram

- We can also learn embeddings by predicting the surrounding context from a single word

Context:



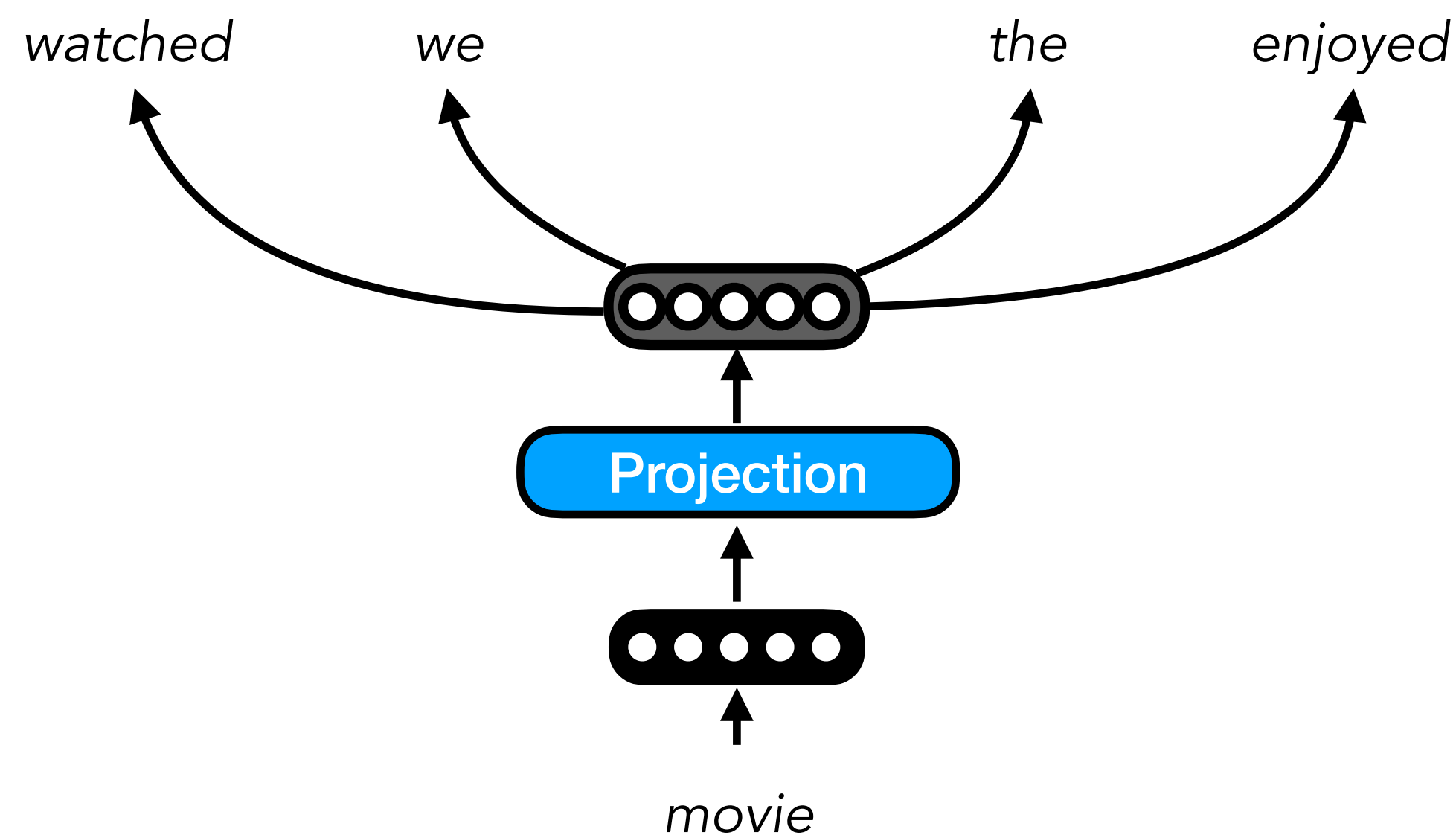
$$\max P(\textit{enjoyed}, \textit{the}, \textit{we}, \textit{watched} | \textit{movie})$$

$$\max P(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} | w_t)$$

Skip-gram

- We can also learn embeddings by predicting the surrounding context from a single word

Context:



$$\max P(\textit{enjoyed}, \textit{the}, \textit{we}, \textit{watched} | \textit{movie})$$

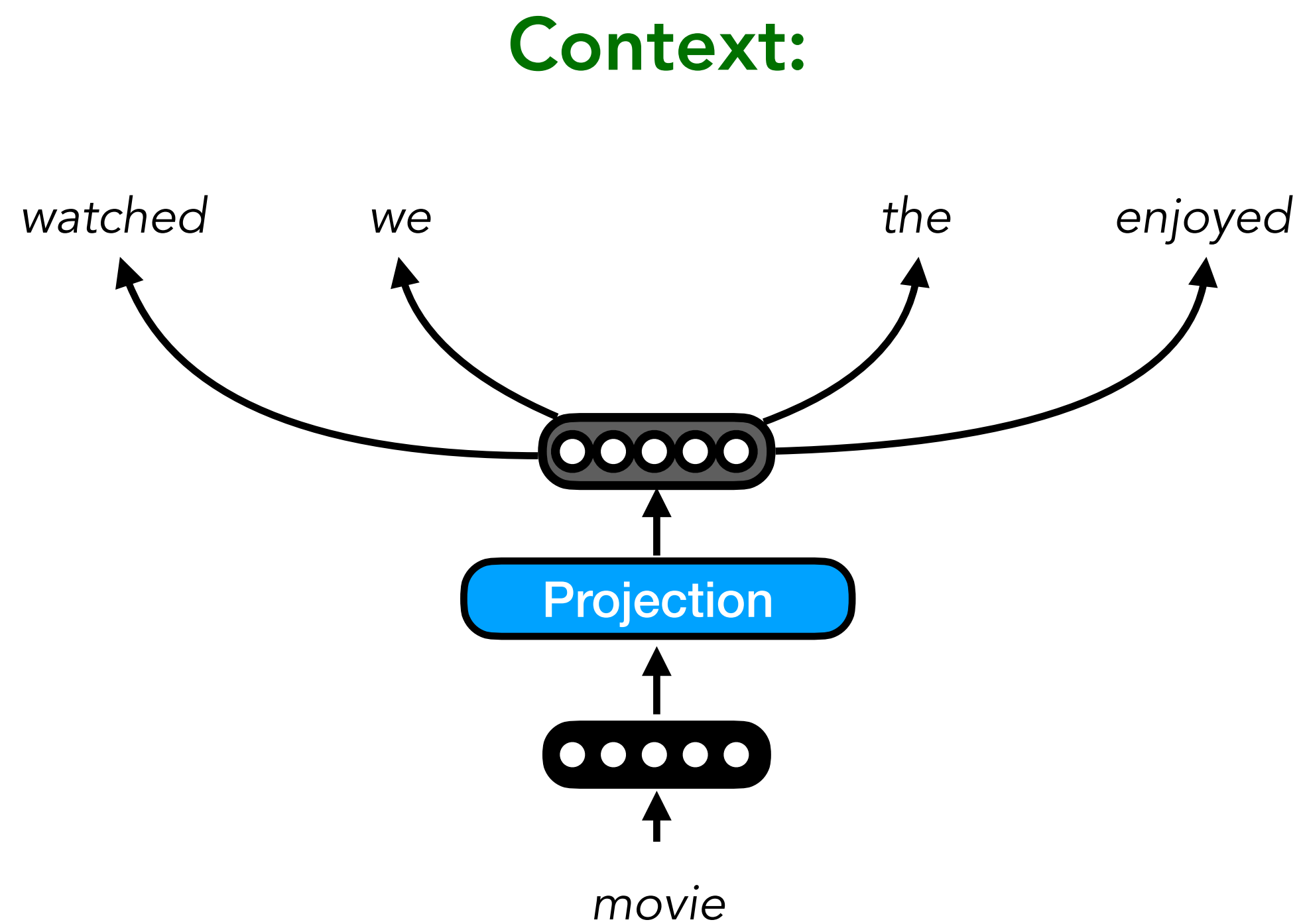
$$\max P(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} | w_t)$$

$$\max \log P(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} | w_t)$$

$$\max \left(\log P(w_{t-2} | w_t) + \log P(w_{t-1} | w_t) \right. \\ \left. + \log P(w_{t+1} | w_t) + \log P(w_{t+2} | w_t) \right)$$

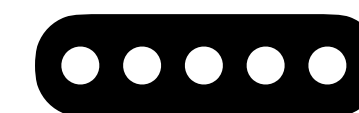
Skip-gram

- We can also learn embeddings by predicting the surrounding context from a single word



$$P(w_x | w_t) = \mathbf{softmax}(\mathbf{U}w_t)$$

$$w_t \in \mathbb{R}^{1 \times d}$$

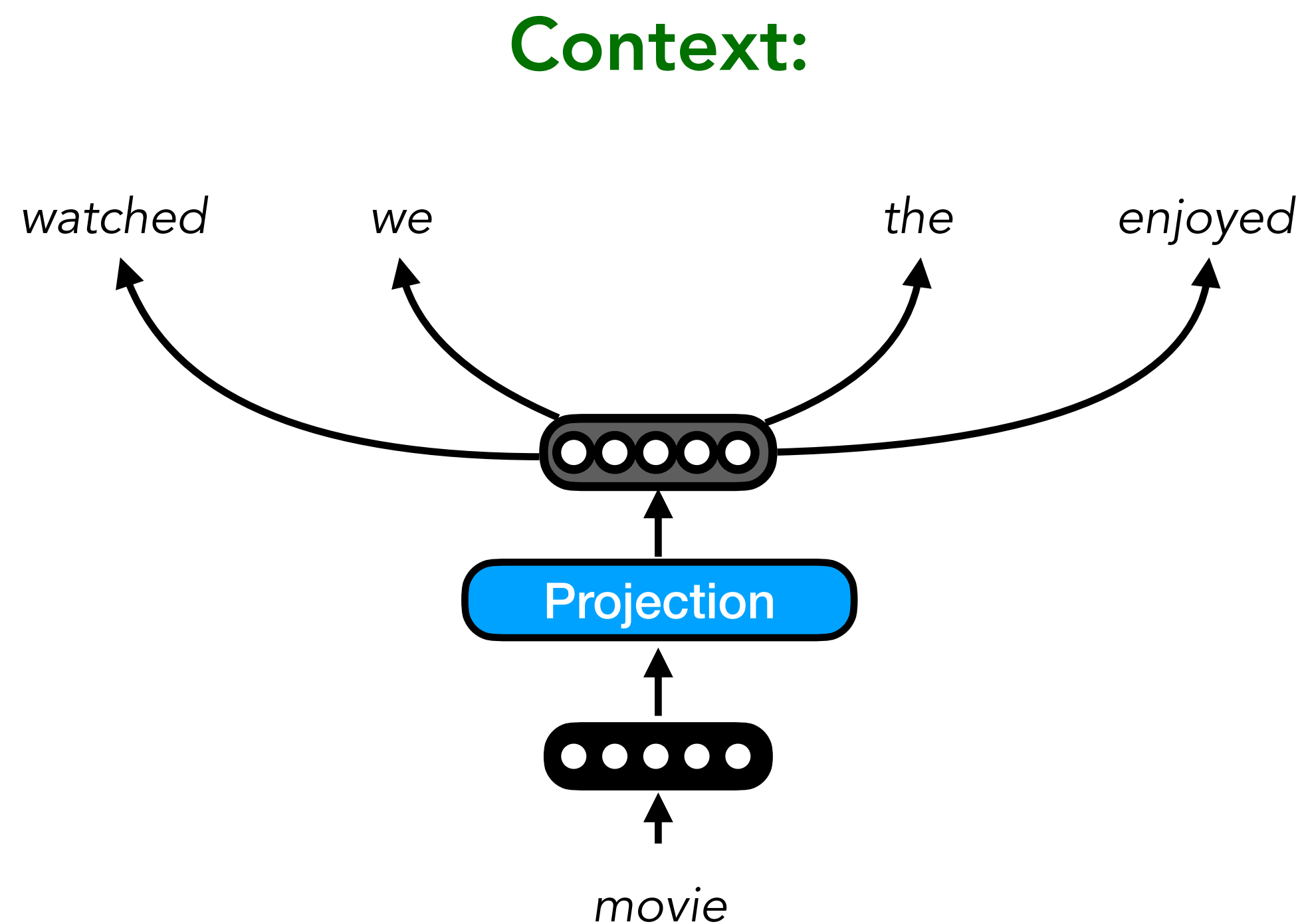


$$\mathbf{U} \in \mathbb{R}^{d \times V}$$

Projection

Skip-gram

- We can also learn embeddings by predicting the surrounding context from a single word



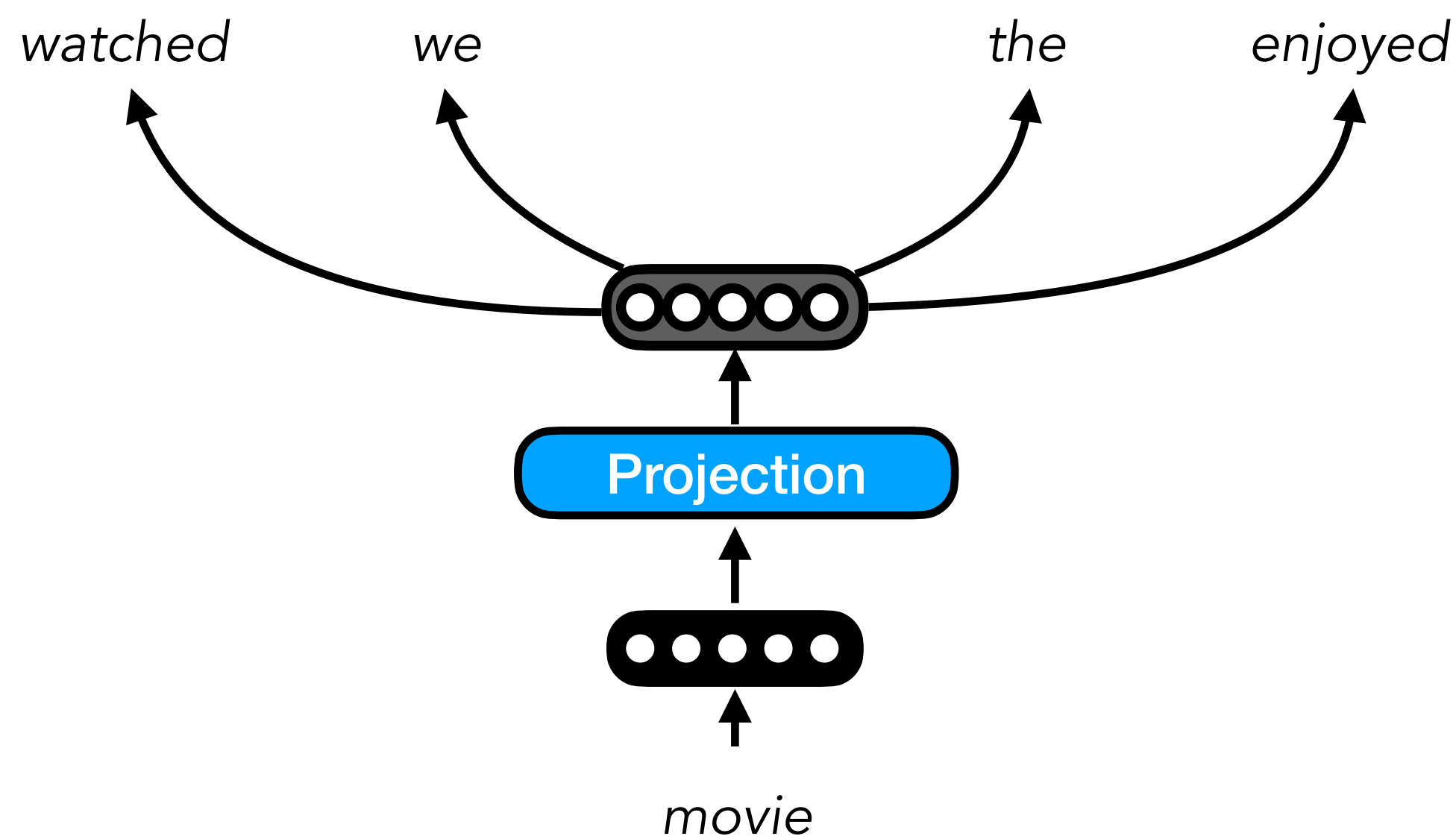
- Model is trained to **minimise** the **negative log probability** of the surrounding words
- Here, we use a window of **N=2**, but the window size is a **hyperparameter** to set
- Typically, set large window (**N=10**), but randomly select $i \in [1, N]$ as dynamic window size so that closer words contribute more to learning

Question

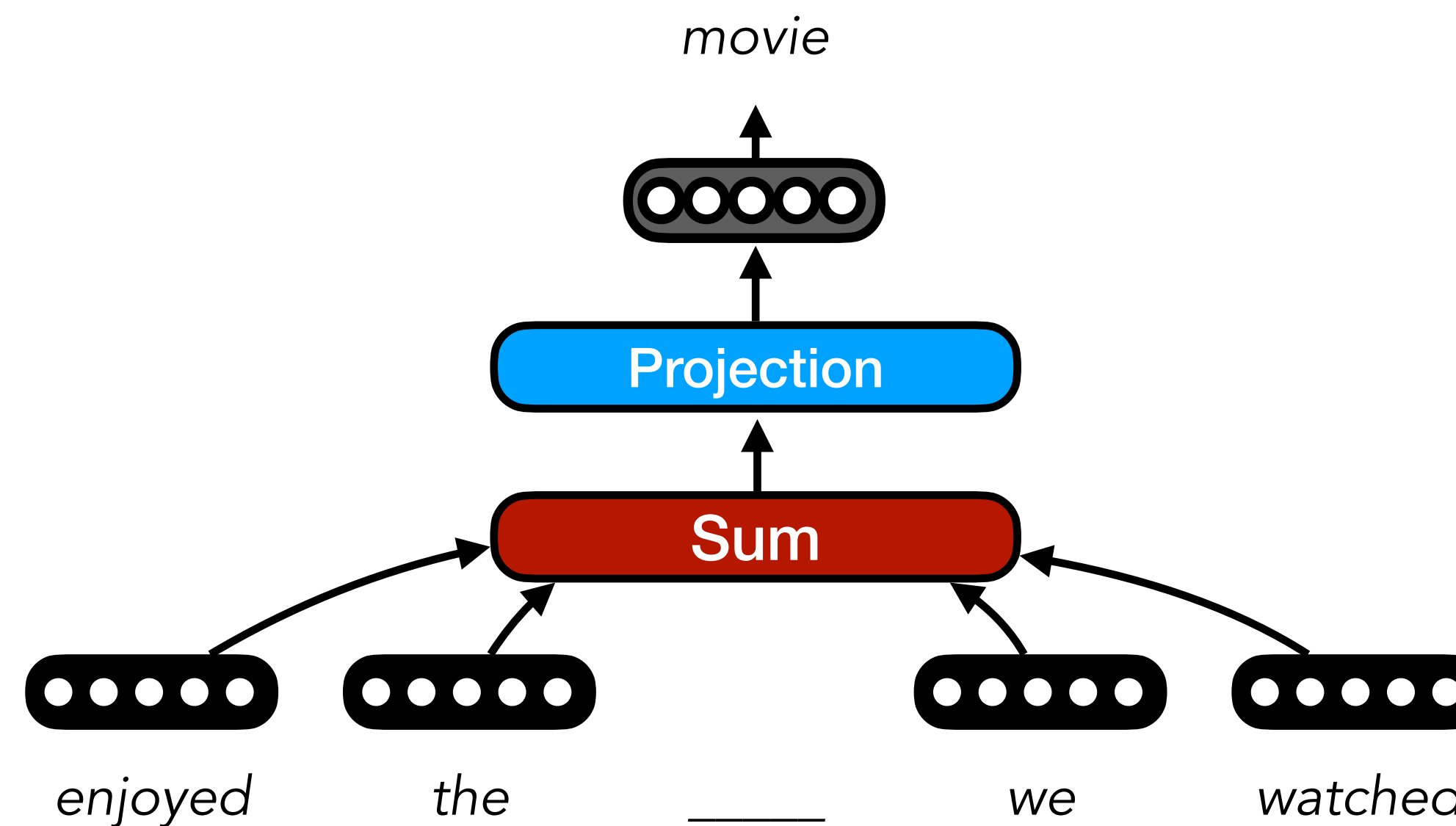
What is the major conceptual difference between the CBOW and Skipgram methods for training word embeddings?

Skip-gram vs. CBOW

- **Question:** Do you expect a difference between what is learned by CBOW and Skipgram methods?



(Mikolov et al., 2013b)



(Mikolov et al., 2013a)

Example

CBOW

```
[ ] top_cbow = cbow.wv.most_similar('cut', topn=10)
print(tabulate(top_cbow, headers=["Word", "Similarity"])
```

Word	Similarity
slice	0.662173
crosswise	0.650036
score	0.630569
tear	0.618827
dice	0.563946
lengthwise	0.557231
cutting	0.557228
break	0.551517
chop	0.541566
carve	0.537967

Skip-gram

```
[ ] top_sg = skipgram.wv.most_similar('cut', topn=10)
print(tabulate(top_sg, headers=["Word", "Similarity"])
```

Word	Similarity
crosswise	0.72921
score	0.702693
slice	0.696898
crossways	0.680091
1/2-inch-thick	0.678496
diamonds	0.671814
diagonally	0.670319
lengthwise	0.665378
cutting	0.66425
wise	0.656825

Recap

- Neural NLP: **Words are vectors!**
- Word embeddings can be learned in a self-supervised manner from large quantities of raw text
- **Two algorithms:** Continuous Bag of Words (CBOW) and Skip-gram

Resources

- **word2vec**: <https://code.google.com/archive/p/word2vec/>
- **GloVe**: <https://nlp.stanford.edu/projects/glove/>
- **FastText**: <https://fasttext.cc/>
- **Gensim**: <https://radimrehurek.com/gensim/>

Download pre-trained word vectors

- Pre-trained word vectors. This data is made available under the [Public Domain Dedication and License](http://www.opendatacommons.org/licenses/pddl/1.0/) v1.0 whose full text can be found at: <http://www.opendatacommons.org/licenses/pddl/1.0/>.
 - [Wikipedia 2014](#) + [Gigaword 5](#) (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download): [glove.6B.zip](#)
 - Common Crawl (42B tokens, 1.9M vocab, uncased, 300d vectors, 1.75 GB download): [glove.42B.300d.zip](#)
 - Common Crawl (840B tokens, 2.2M vocab, cased, 300d vectors, 2.03 GB download): [glove.840B.300d.zip](#)
 - Twitter (2B tweets, 27B tokens, 1.2M vocab, uncased, 25d, 50d, 100d, & 200d vectors, 1.42 GB download): [glove.twitter.27B.zip](#)
- Ruby [script](#) for preprocessing Twitter data

References

- Firth, J.R. (1957). *A Synopsis of Linguistic Theory, 1930-1955*.
- Mikolov, T., Chen, K., Corrado, G.S., & Dean, J. (2013a). Efficient Estimation of Word Representations in Vector Space. *International Conference on Learning Representations*.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G.S., & Dean, J. (2013b). Distributed Representations of Words and Phrases and their Compositionality. *ArXiv, abs/1310.4546*.
- Pennington, J., Socher, R., & Manning, C.D. (2014). GloVe: Global Vectors for Word Representation. *Conference on Empirical Methods in Natural Language Processing*.
- Bojanowski, P., Grave, E., Joulin, A., & Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the association for computational linguistics*.
- Mikolov, T., Grave, E., Bojanowski, P., Puhresch, C., & Joulin, A. (2018). Advances in pre-training distributed word representations. *International Conference on Language Resources and Evaluation*.

Deep Learning for Natural Language Processing

Antoine Bosselut

EPFL



Part 3: Attentive Neural Modeling with Transformers

Section Outline

- **Background:** Long-range Dependency Modeling
- **Content:** Attention, Self-Attention, Multi-headed Attention, Transformer Blocks, Transformers
- **Exercise Session:** Visualizing Transformer Attention

Issue with Recurrent Models

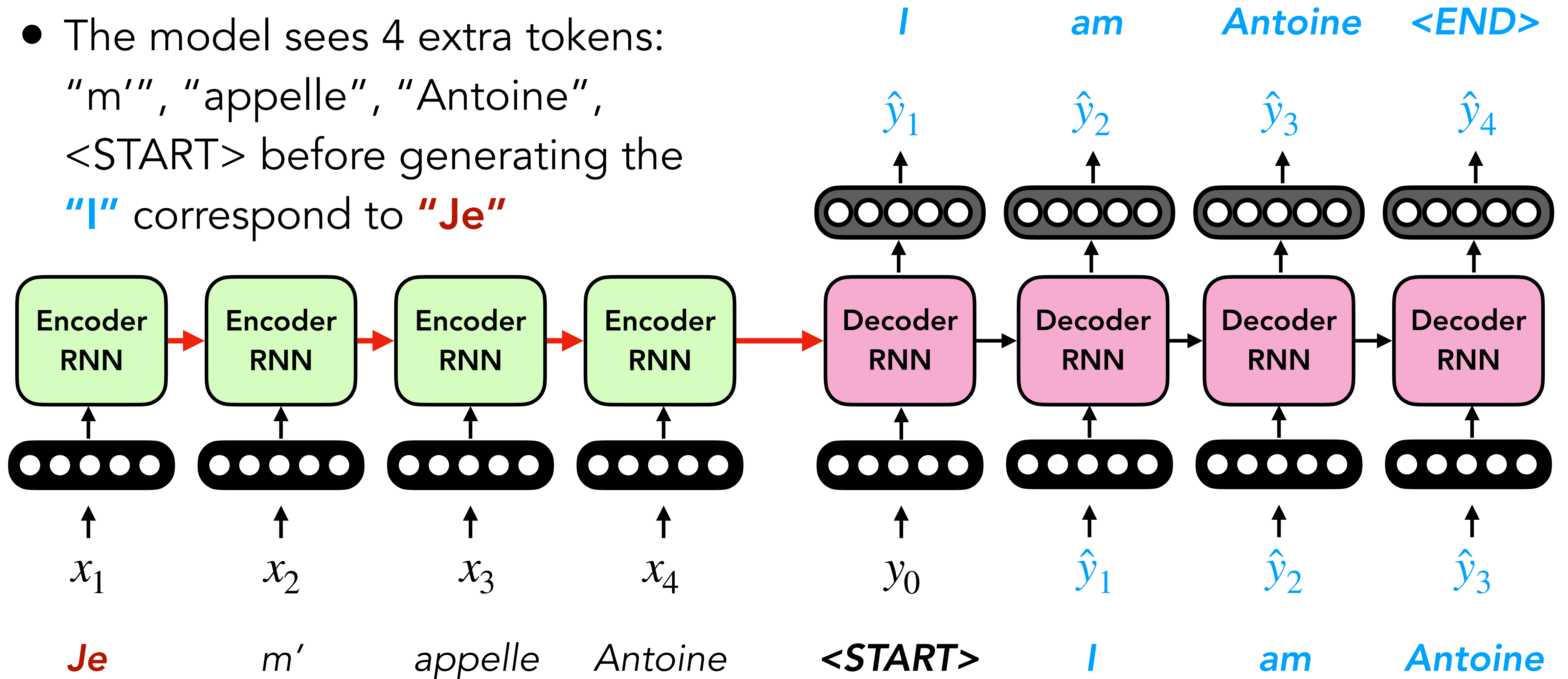
- Multiple steps of state overwriting makes it challenging to learn long-range dependencies.

*They tuned, discussed for a moment, then struck up a lively **jig**. Everyone joined in, turning the courtyard into an even more chaotic scene, people now **dancing** in circles, **swinging** and **spinning** in circles, everyone making up their own **dance steps**. I felt my feet tapping, my body wanting to move. Aside from writing, I 've always loved **dancing** .*

- Nearby words should affect each other more than farther ones, but RNNs make it challenging to learn **any** long-range interactions

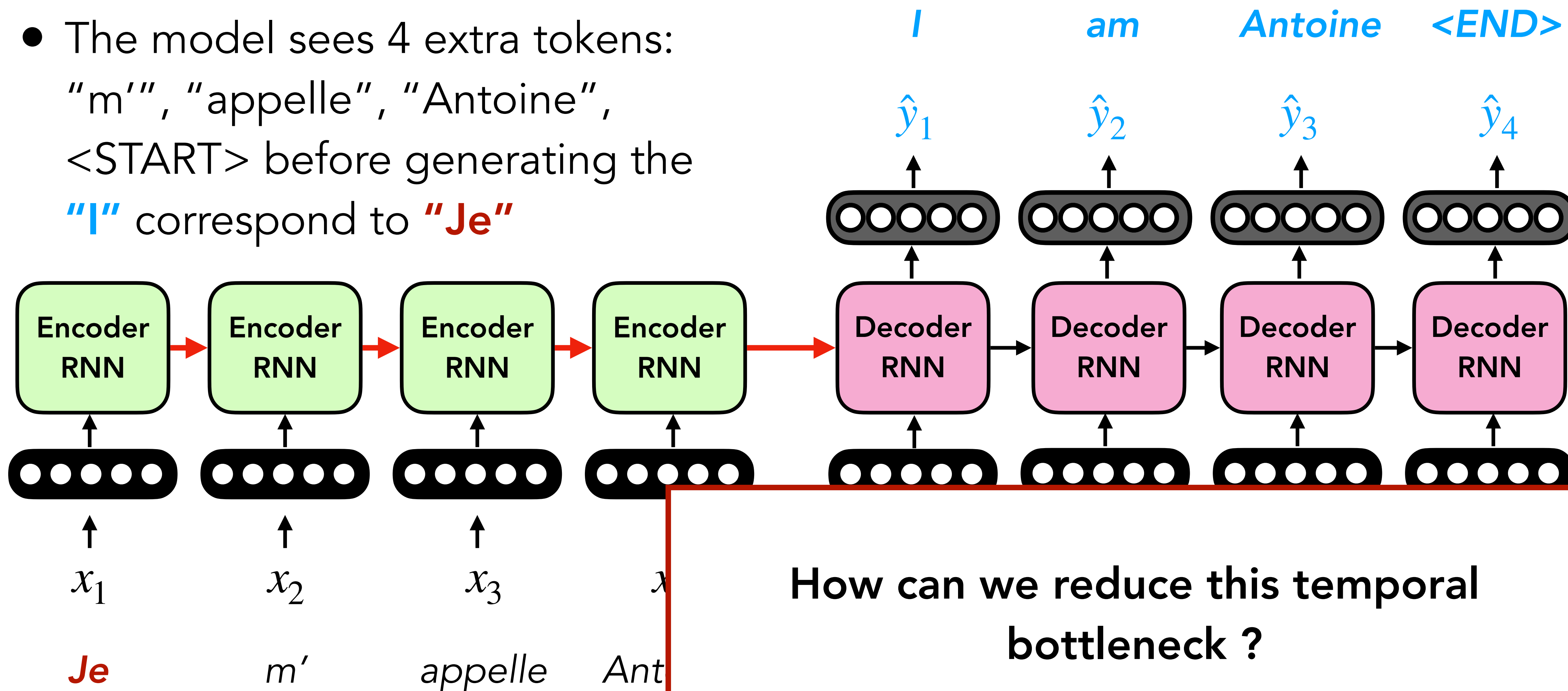
Toy Example

- The model sees 4 extra tokens: "m'", "appelle", "Antoine", <START> before generating the "I" correspond to "Je"



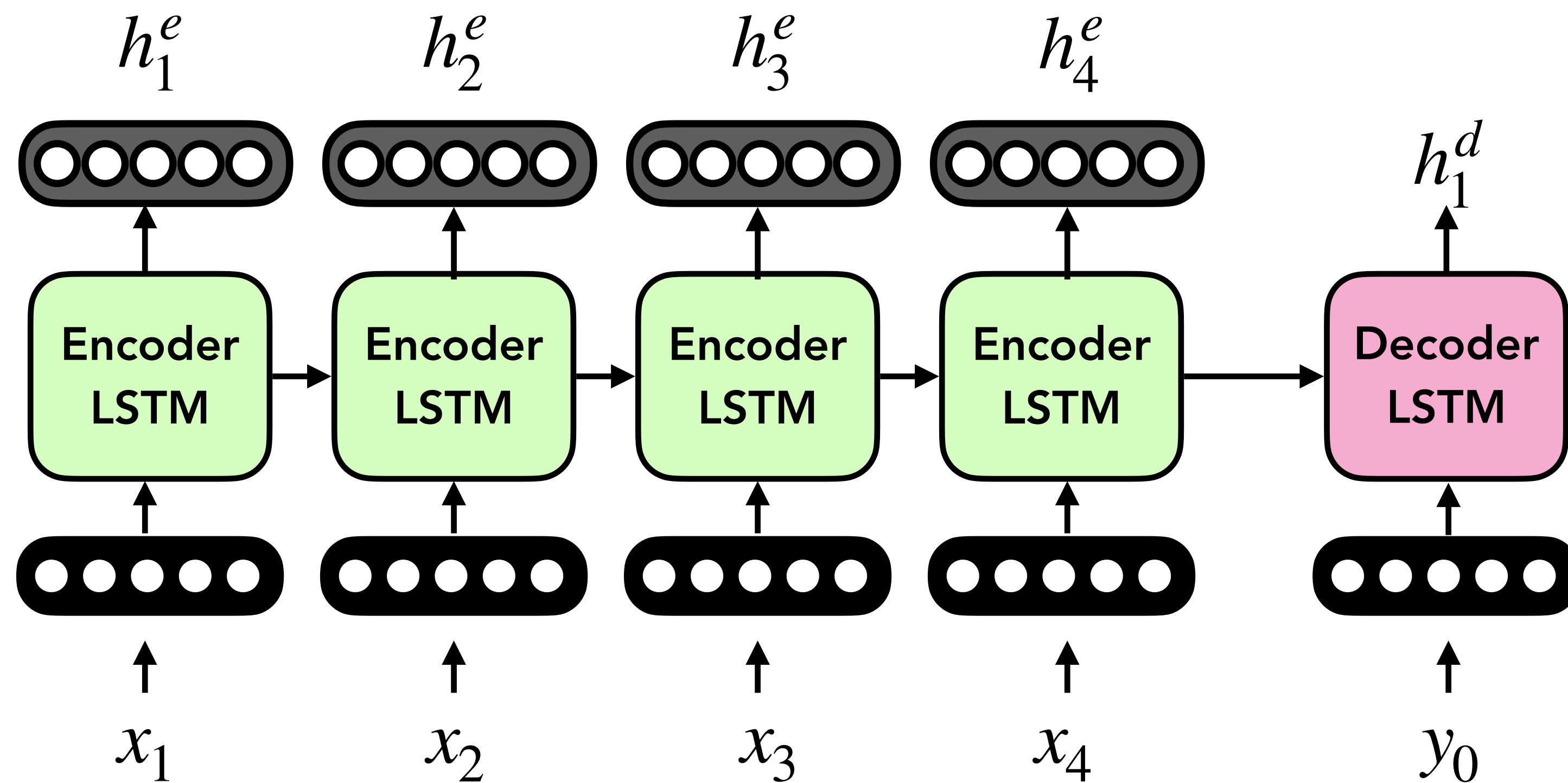
Toy Example

- The model sees 4 extra tokens: "m'", "appelle", "Antoine", <START> before generating the "!" correspond to "Je"

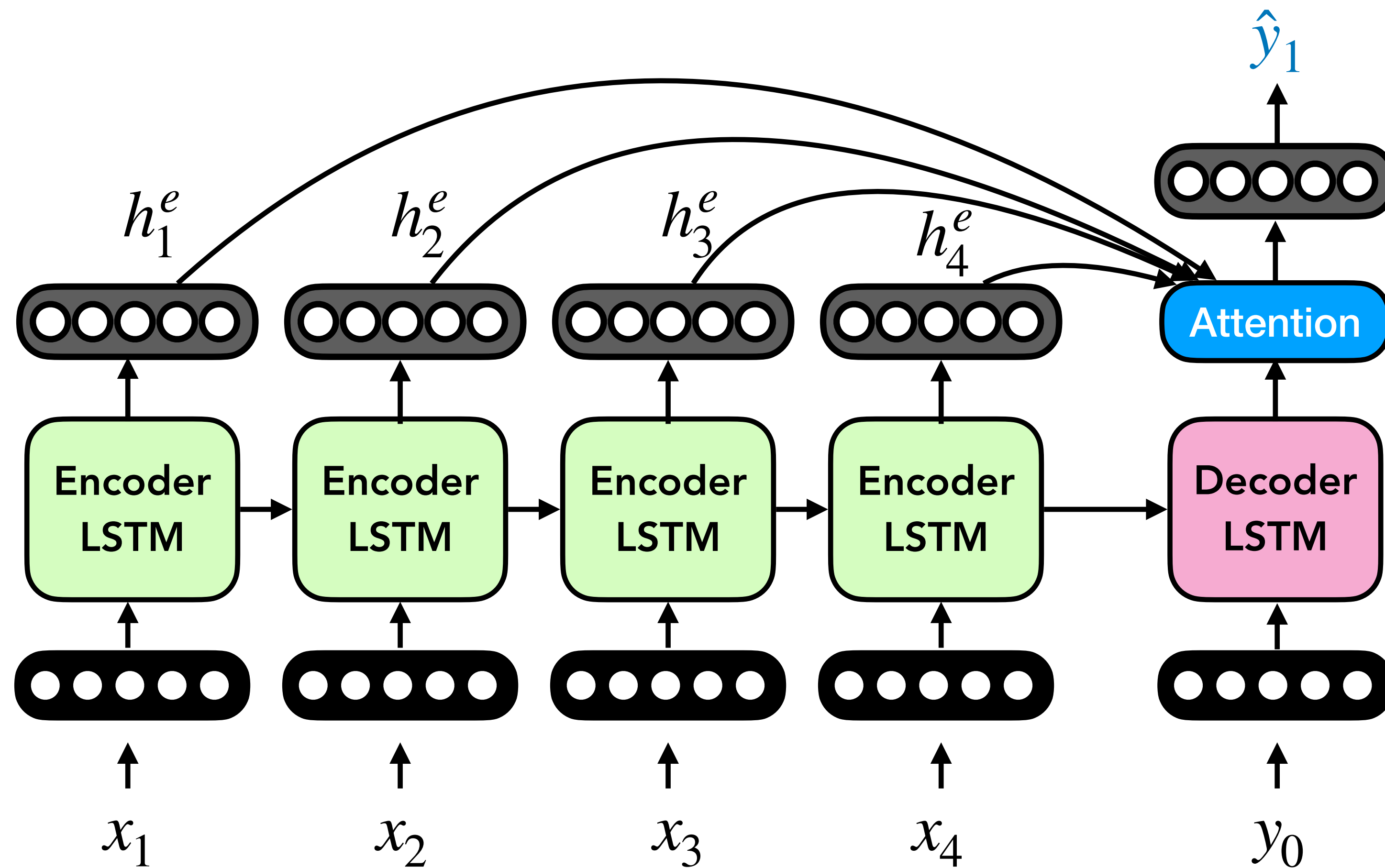


Attentive Encoder-Decoder Models

- **Recall:** At each encoder time step, there is an output of the RNN!



Attentive Encoder-Decoder Models



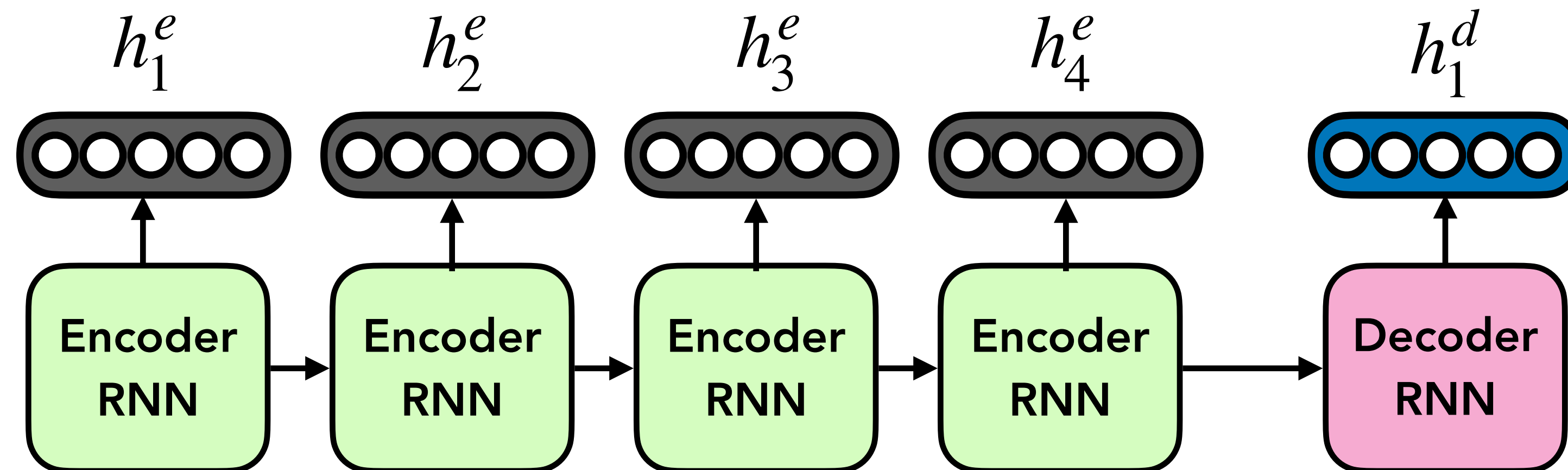
- **Recall:** At each encoder time step, there is an output of the RNN!
- **Idea:** Use the output of the Decoder LSTM to compute an **attention** (i.e., a mixture) over all the h_t^e outputs of the encoder LSTM
- **Intuition:** focus on different parts of the input at each time step

What is attention?

- Attention is a **weighted average over a set of inputs**

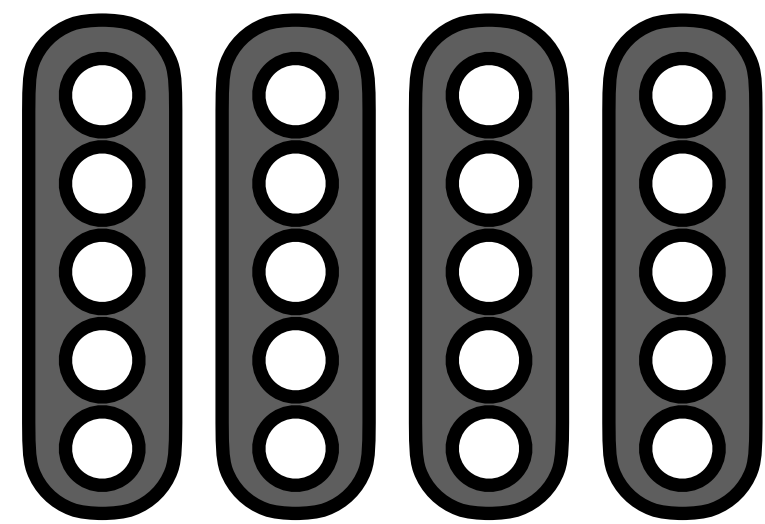
h_t^e = encoder output hidden states

- How should we compute this weighted average?



Attention Function

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



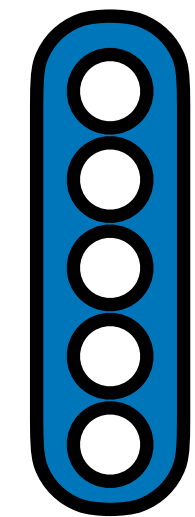
h_1^e h_2^e h_3^e h_4^e

h_t^e = encoder output hidden states

Also known as a “keys”

h_t^d = decoder output hidden state

Also known as a “query”



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

Also known as a "keys"

h_t^d = decoder output hidden state

Also known as a "query"

$$a_1 = f\left(h_1^e, h_1^d\right) \quad a_2 = f\left(h_2^e, h_1^d\right) \quad a_3 = f\left(h_3^e, h_1^d\right) \quad a_4 = f\left(h_4^e, h_1^d\right)$$

- We have a single query vector for multiple key vectors

Attention Function

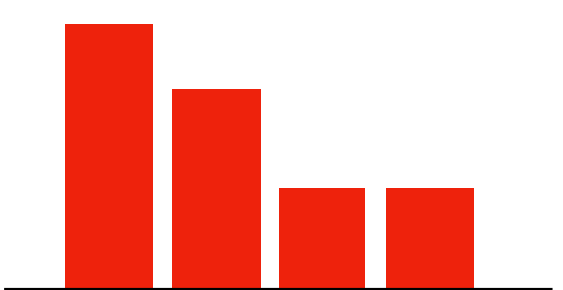
Attention Function	Formula
Multiplicative	$a = h^e \mathbf{W} h^d$
Linear	$a = v^T \phi(\mathbf{W}[h^e; h^d])$
Scaled Dot Product	$a = \frac{(\mathbf{W}h^e)^T (\mathbf{U}h^d)}{\sqrt{d}}$

Attention Function

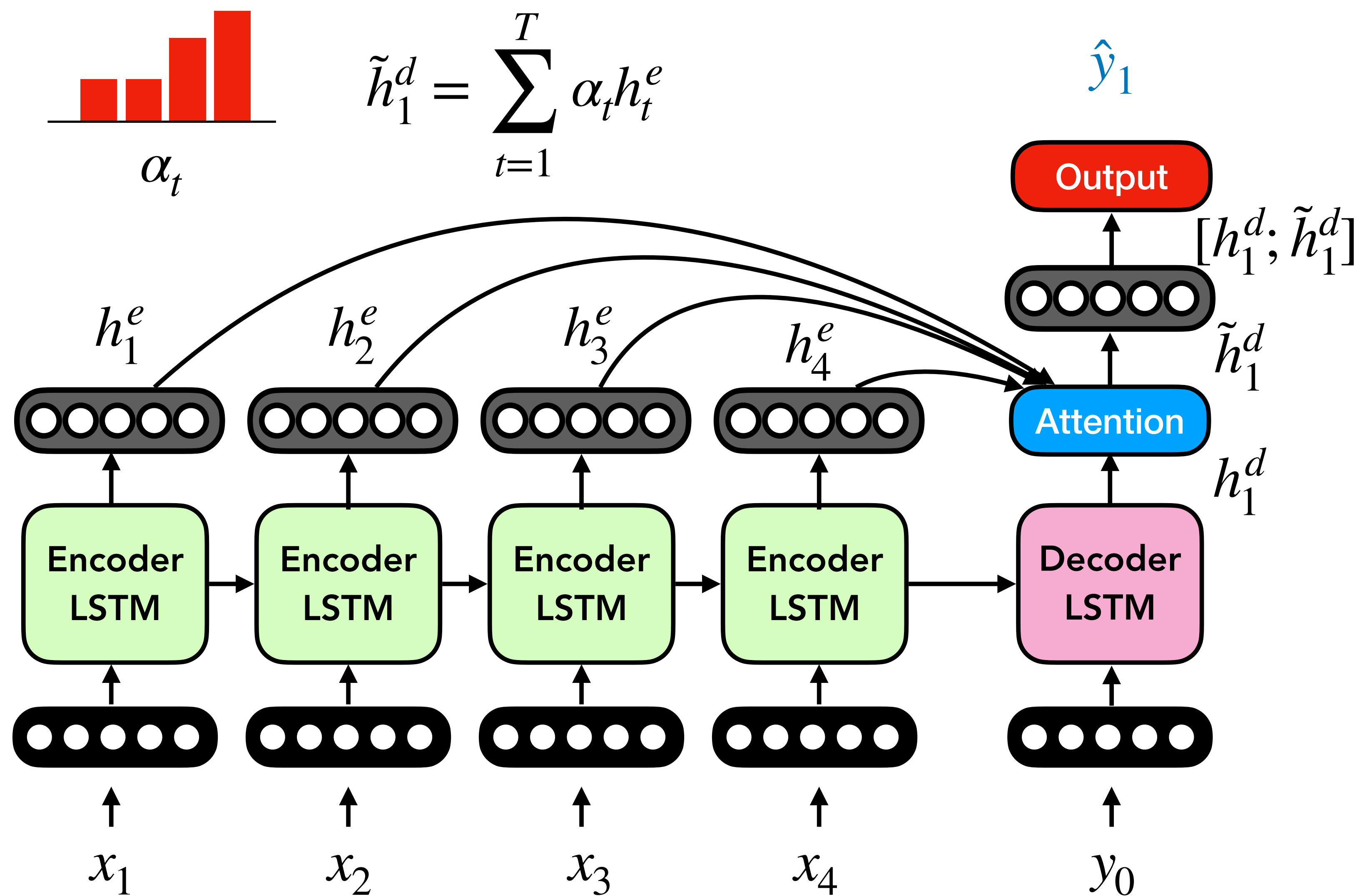
- **Compute** pairwise similarity between each encoder hidden state and decoder hidden state ("idea of what to decode")

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

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

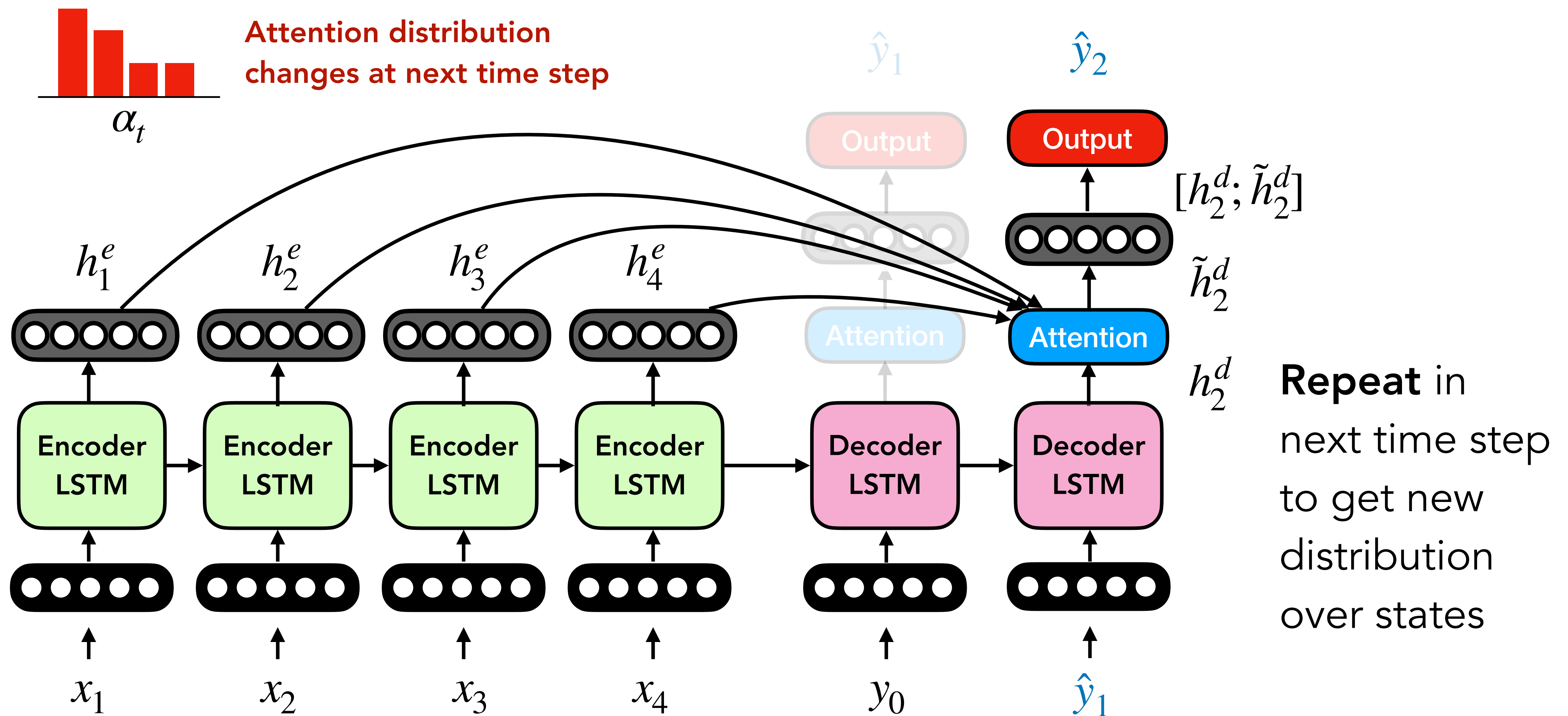
Softmax! $\alpha_t = \frac{e^{a_t}}{\sum_j e^{a_j}}$ \rightarrow  $\rightarrow \tilde{h}_1^d = \sum_{t=1}^T \alpha_t h_t^e$ **Here h_t^e is known as the "value"**

Attentive Encoder-Decoder Models

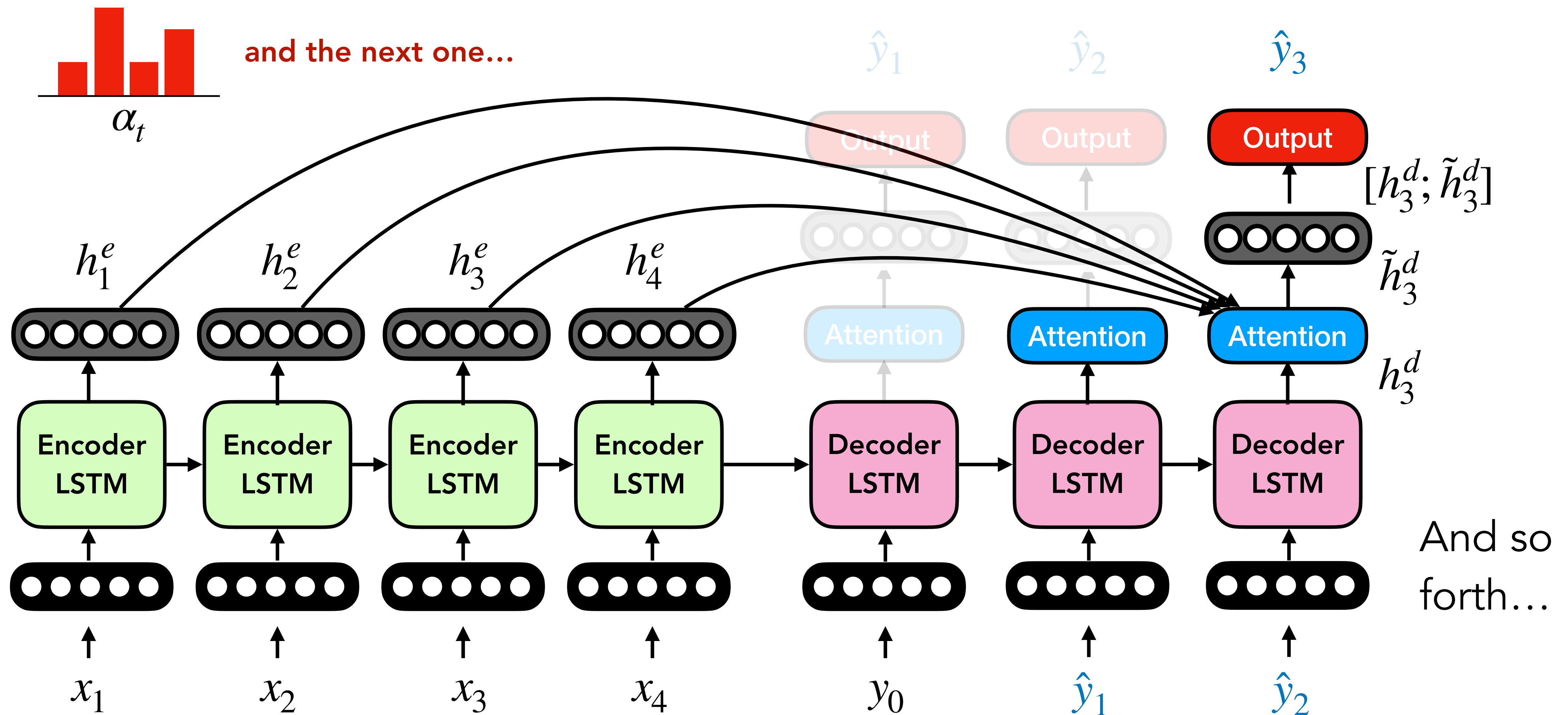


- **Intuition:** \tilde{h}_1^d contains information about hidden states that got **high** attention
- Typically, \tilde{h}_1^d is concatenated (or composed in some other manner) with h_1^d (the original decoder state) before being passed to the **output** layer
- **Output** layer predicts the most likely output token \hat{y}_1

Attentive Encoder-Decoder Models



Attentive Encoder-Decoder Models



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

h_t^e = encoder output hidden states

h_t^d = decoder initial hidden state

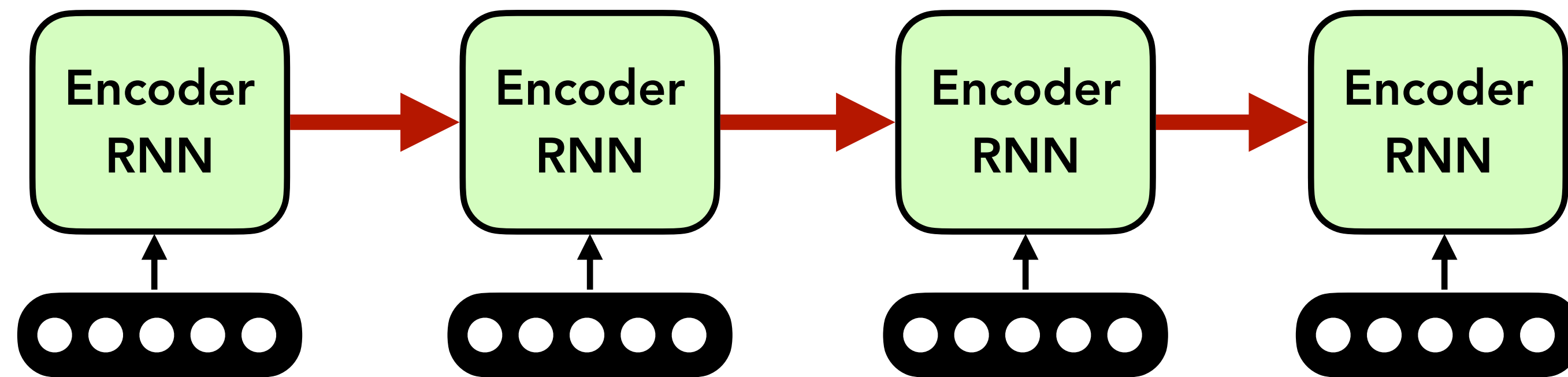
- Many possible functions for computing scores (dot product, bilinear, etc.)
- **Temporal Bottleneck Fixed! Direct connection** between decoder and ALL encoder states

Question

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

Encoder is still Recurrent

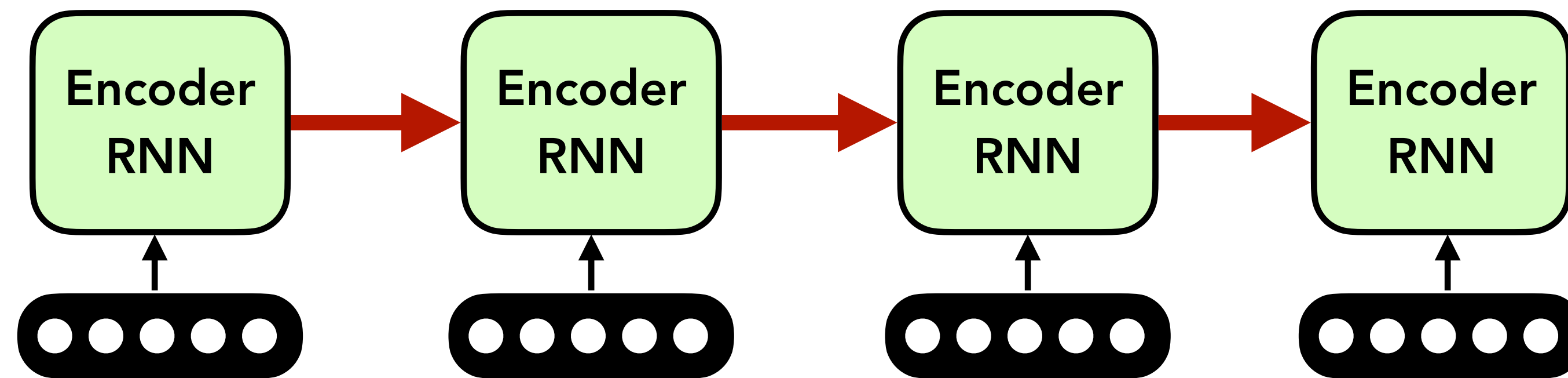
- **Encoder:** Recurrent functions can't be parallelized because previous state needs to be computed to encode next one



- **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**

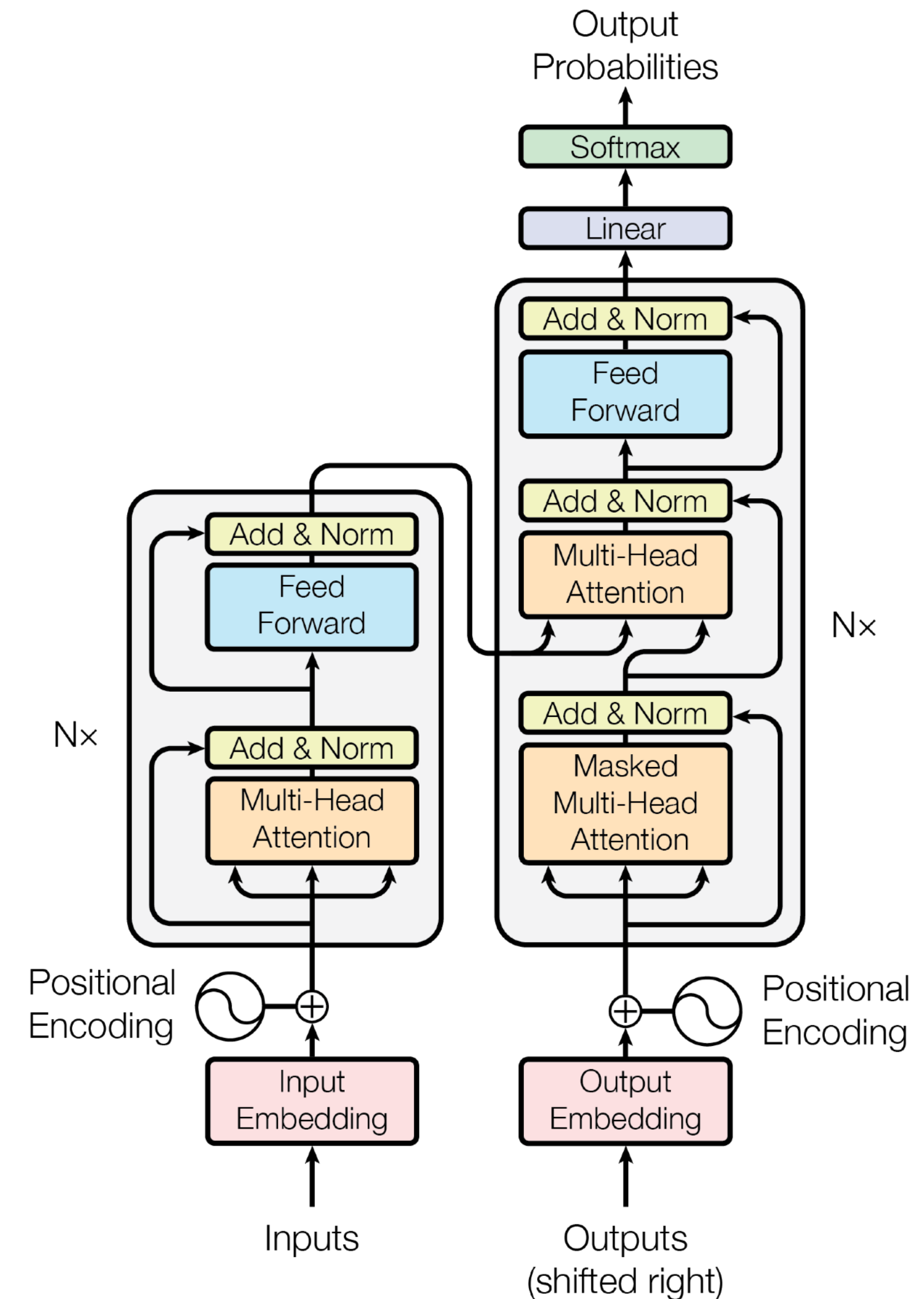
Who can think of a task where this might be a problem?

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!

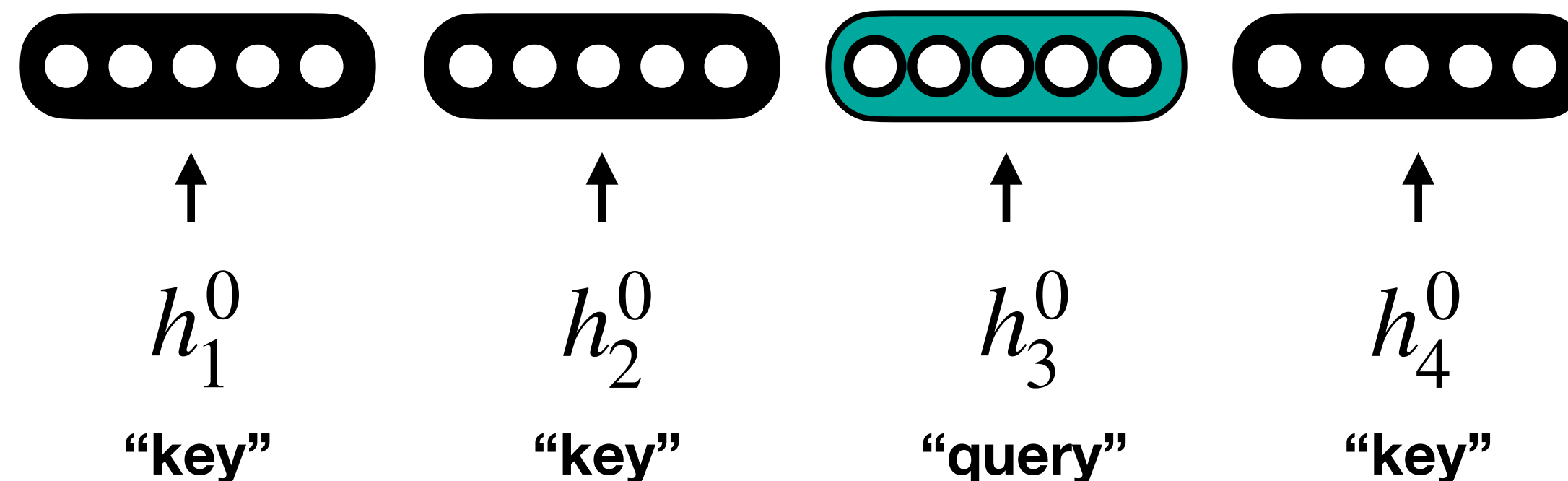
Encode sequences with self-attention



Self-Attention Toy Example

- **Original Idea:** Use decoder hidden state to compute attention distribution over encoder hidden states
- **New Idea:** Could we use encoder hidden states to compute attention distribution over themselves?
- **Ditch recurrence** and compute encoder state representations in parallel!

h_t^ℓ = encoder hidden state at time step t at layer ℓ

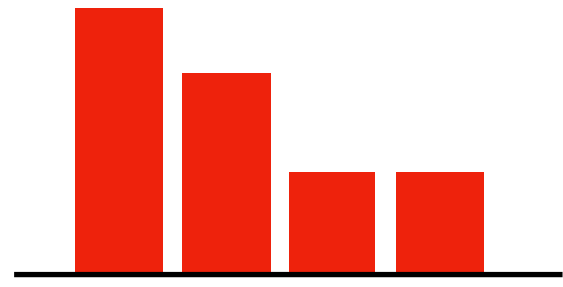


Recap: Attention with RNNs

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

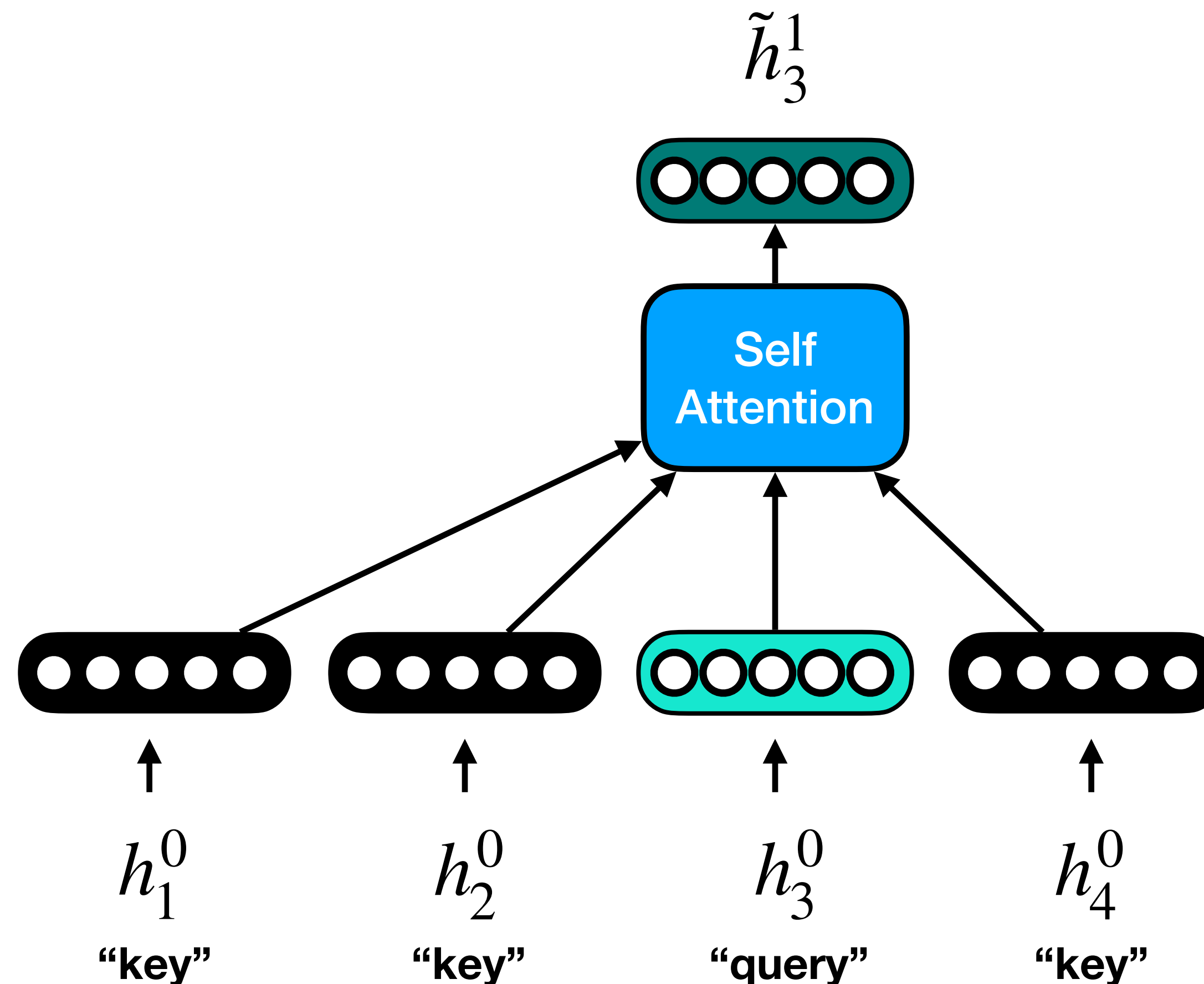
$$\begin{array}{ccc} a_1 = f\left(\begin{array}{c} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{array}, \begin{array}{c} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{array}\right) & a_2 = f\left(\begin{array}{c} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{array}, \begin{array}{c} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{array}\right) & a_3 = f\left(\begin{array}{c} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{array}, \begin{array}{c} \text{○} \\ \text{○} \\ \text{○} \\ \text{○} \end{array}\right) \\ \begin{array}{cc} h_1^e & h_1^d \\ \text{"key"} & \text{"query"} \end{array} & \begin{array}{cc} h_2^e & h_1^d \\ \text{"key"} & \text{"query"} \end{array} & \begin{array}{cc} h_3^e & h_1^d \\ \text{"key"} & \text{"query"} \end{array} \end{array}$$

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

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

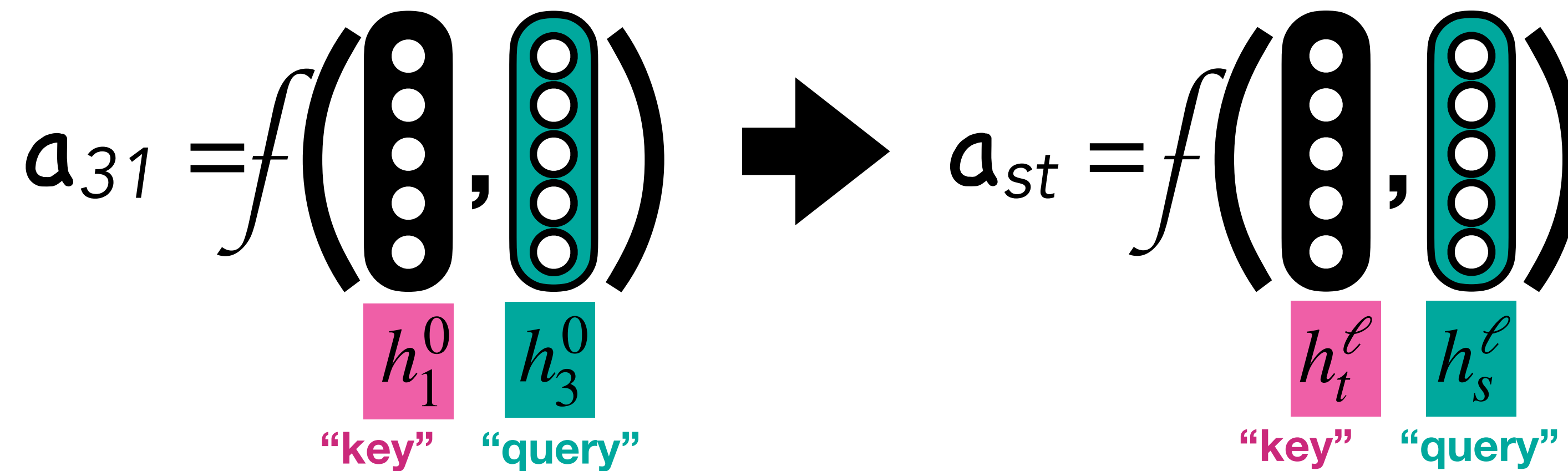
Self-Attention Toy Example

- For a particular encoder time step, compute pairwise score between this hidden state (the query) and the other encoder hidden states



Self-Attention Toy Example

h_t^ℓ = encoder hidden state at time step t at layer ℓ



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

Compute pairwise scores

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

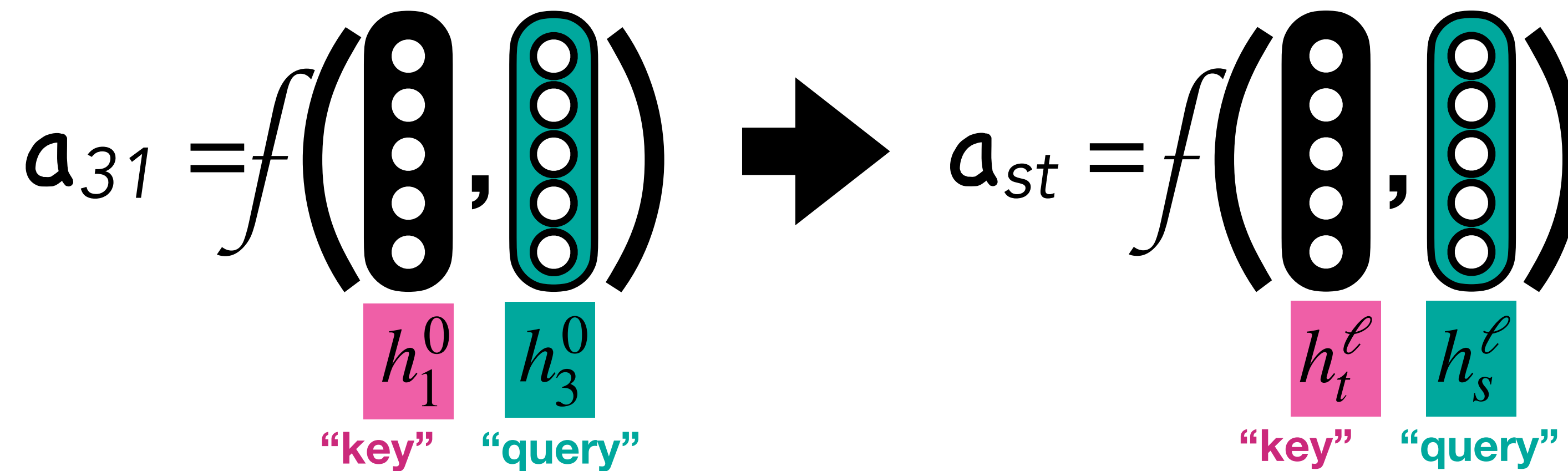
Get attention distribution

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

Attend to values to get weighted sum

Self-Attention Toy Example

h_t^ℓ = encoder hidden state at time step t at layer ℓ



$\{1, \dots, t, \dots, T\}$
includes s !

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

Compute pairwise scores

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

Get attention distribution

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

Attend to values to get weighted sum

Self-attention!

Self-Attention Toy Example

Compute pairwise scores

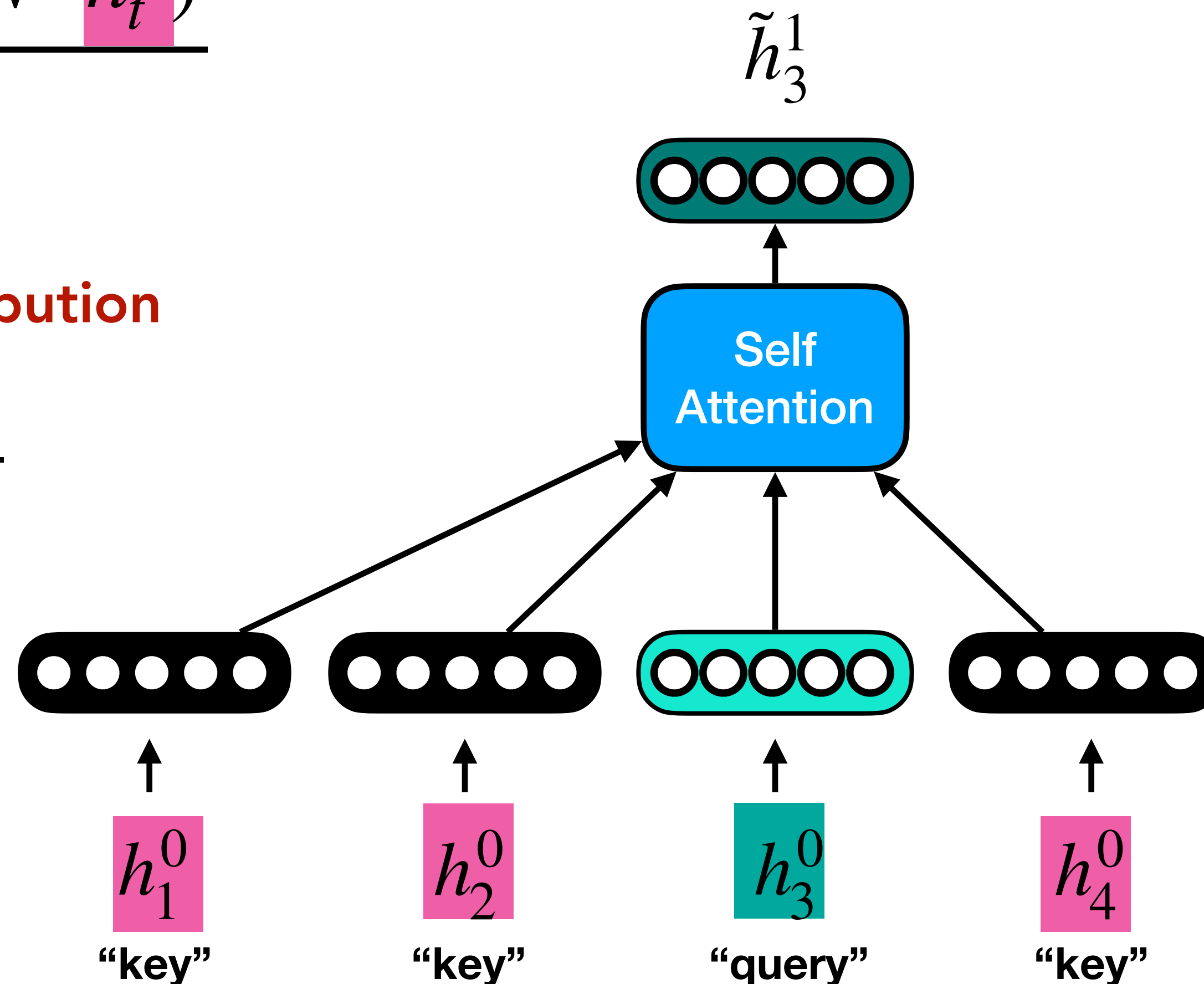
$$a_{st} = \frac{(\mathbf{W}^Q h_s^\ell)^T (\mathbf{W}^K h_t^\ell)}{\sqrt{d}}$$

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

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



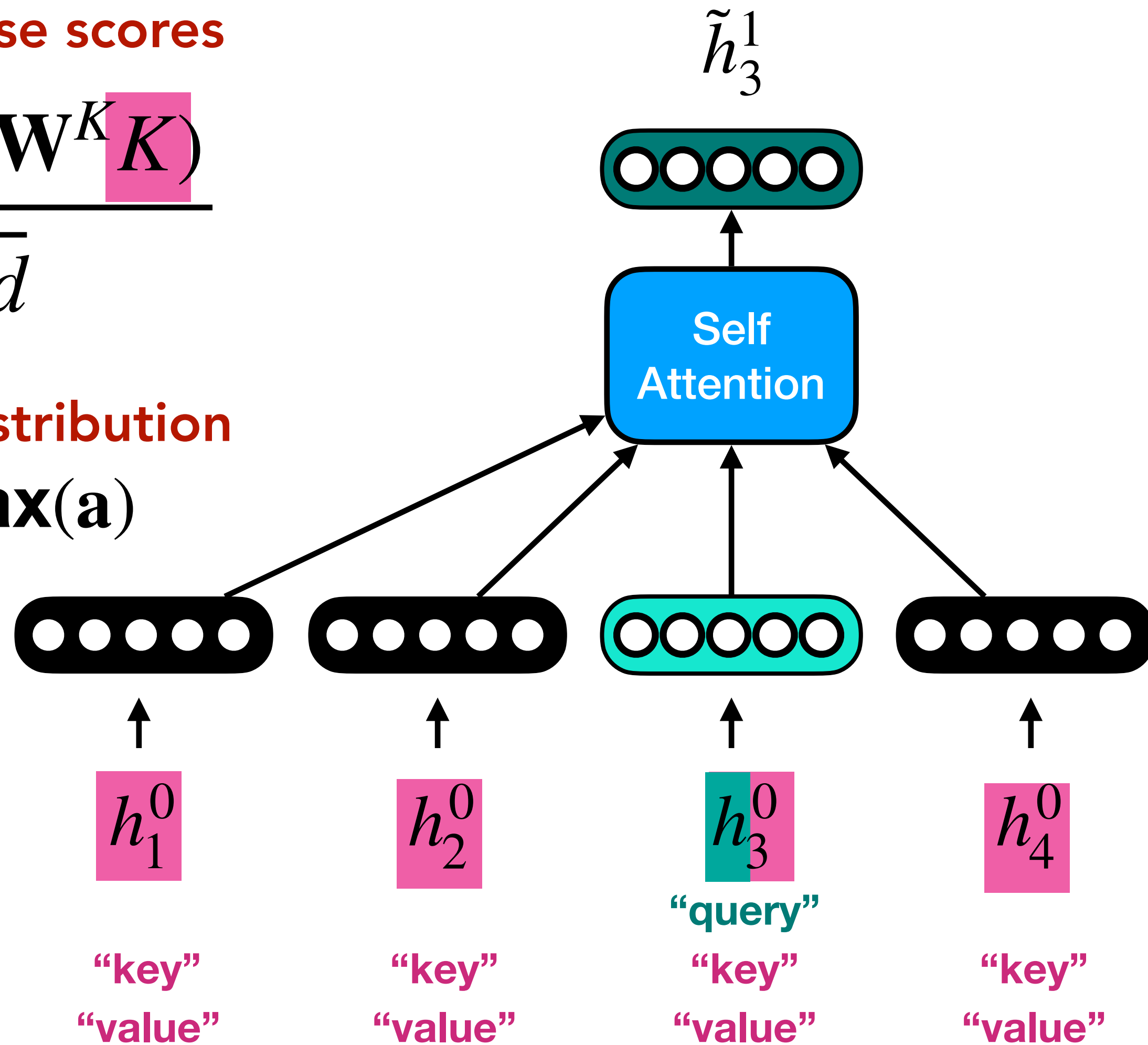
Self-Attention Toy Example

Compute pairwise scores

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

Get attention distribution

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



Attend to values to get weighted sum

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

“query” $\mathbf{q} = h_s^\ell$

“values”
“keys” $\mathbf{K} = \mathbf{V} = \{h_t^\ell\}_{t=0}^T$

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

Self-Attention Toy Example

Compute pairwise scores

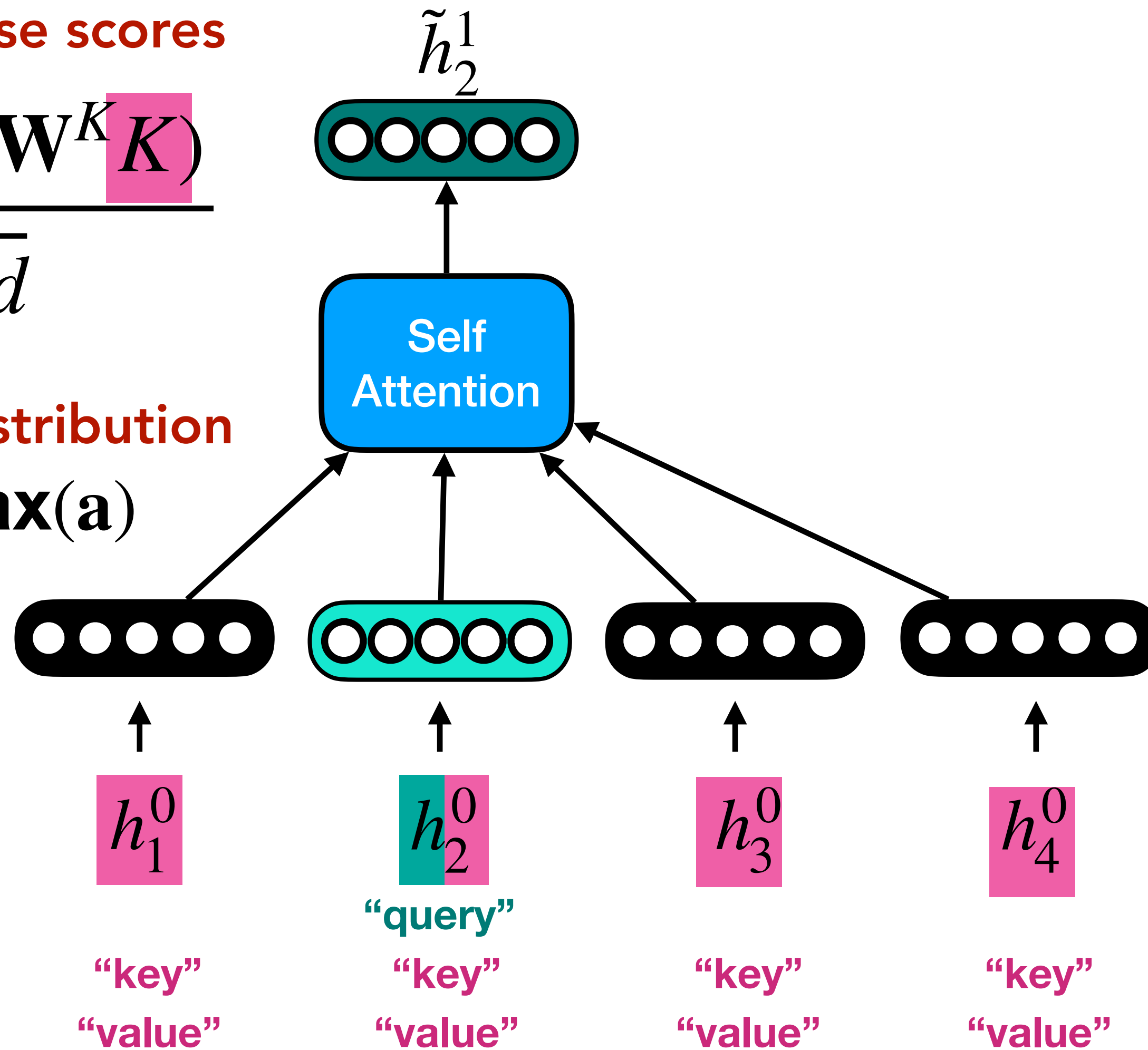
$$\mathbf{a} = \frac{(\mathbf{W}^Q \mathbf{q})(\mathbf{W}^K \mathbf{K})}{\sqrt{d}}$$

Attend to values to get weighted sum

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

Get attention distribution

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



“query” $\mathbf{q} = h_s^\ell$

“values” $\mathbf{K} = \mathbf{V} = \{h_t^\ell\}_{t=0}^T$
“keys”

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

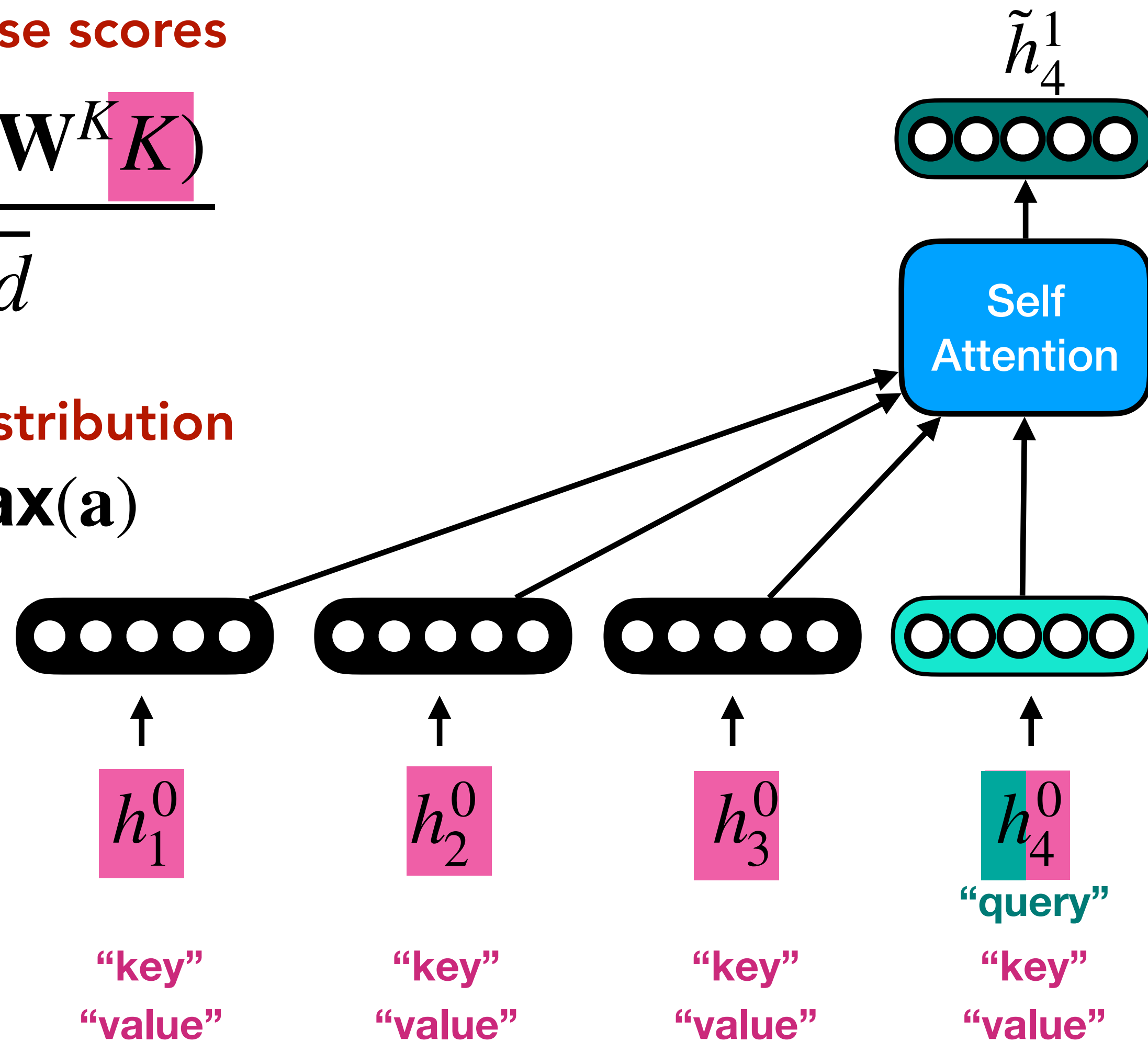
Self-Attention Toy Example

Compute pairwise scores

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

Get attention distribution

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



Attend to values to get weighted sum

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

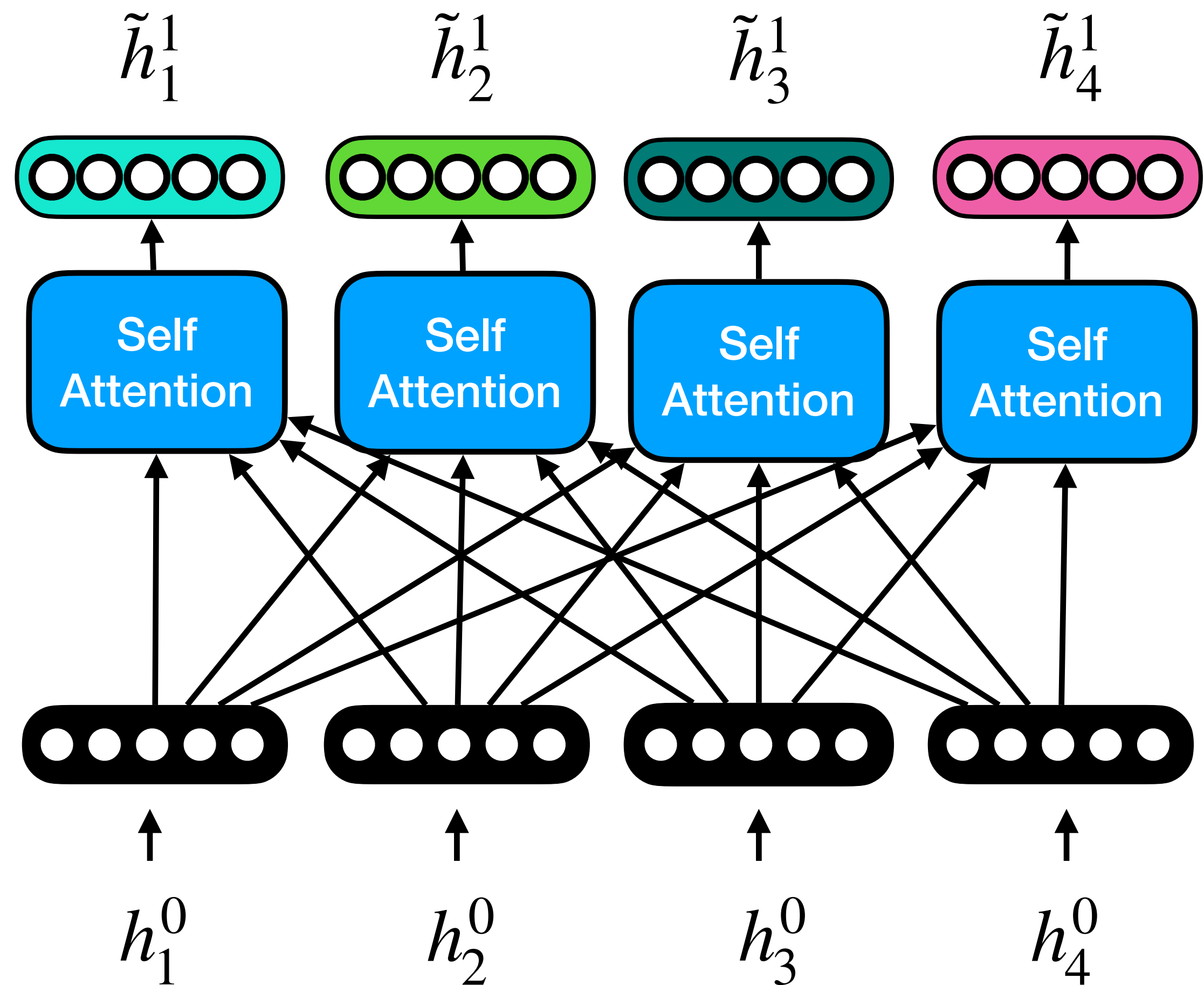
"query" $q = h_s^\ell$

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

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

Self-Attention Toy Example

- Every token is a query! Recompute self-attention value for each position in the sequence



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

$$\tilde{h}_2^1 = \text{Attention} \left(h_2^0, \{h_t^0\}_{t=0}^{t=3} \right)$$

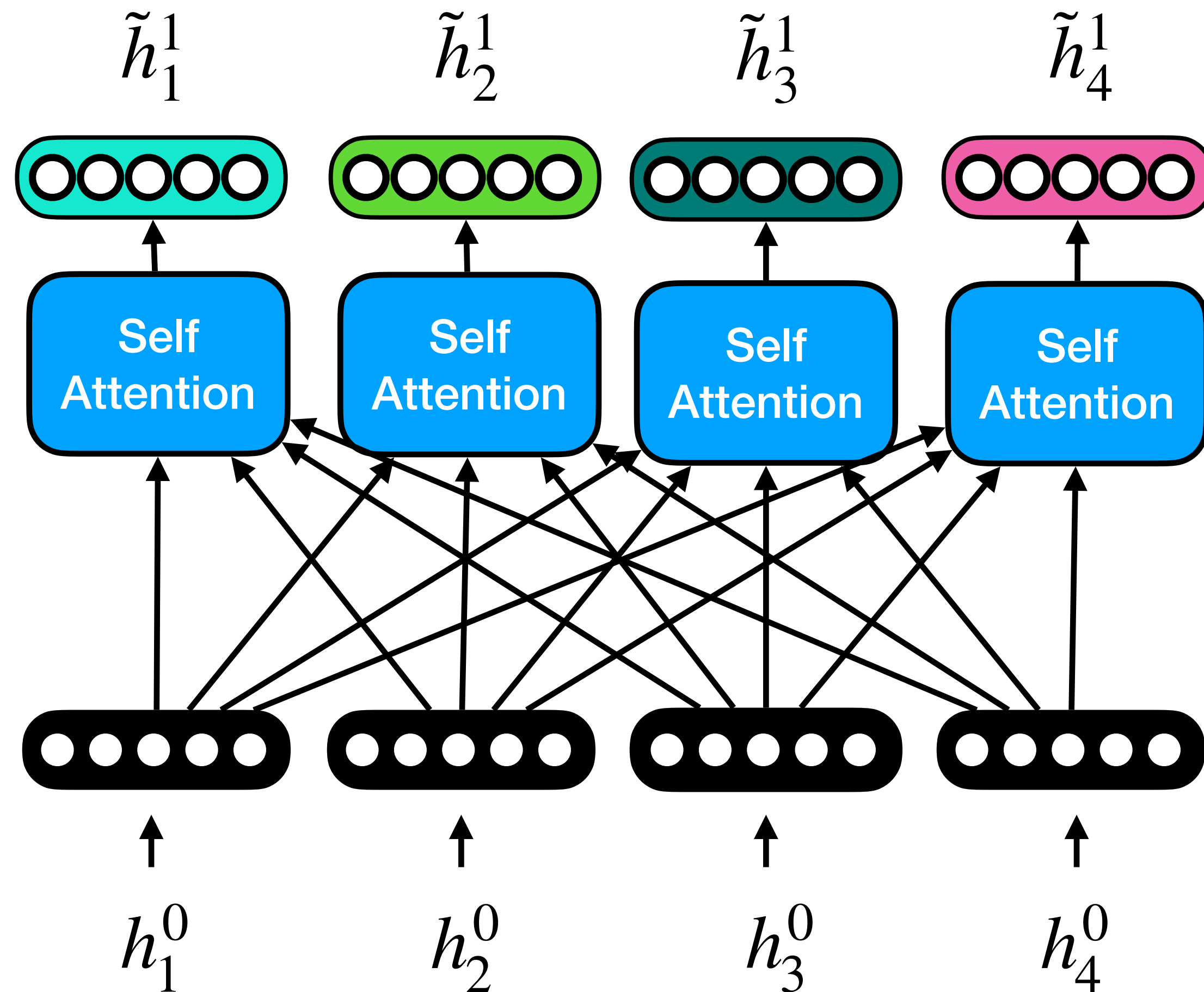
$$\tilde{h}_3^1 = \text{Attention} \left(h_3^0, \{h_t^0\}_{t=0}^{t=3} \right)$$

$$\tilde{h}_4^1 = \text{Attention} \left(h_4^0, \{h_t^0\}_{t=0}^{t=3} \right)$$

Question

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

Self-Attention Recap



- **Computed in parallel** — no previous time step computation needed for the next one
- **No long-term dependencies** — direct connection between all time-steps in sequence

Multi-Headed Self-Attention

- Project V, K, Q into H sub-vectors where H is the number of "heads"

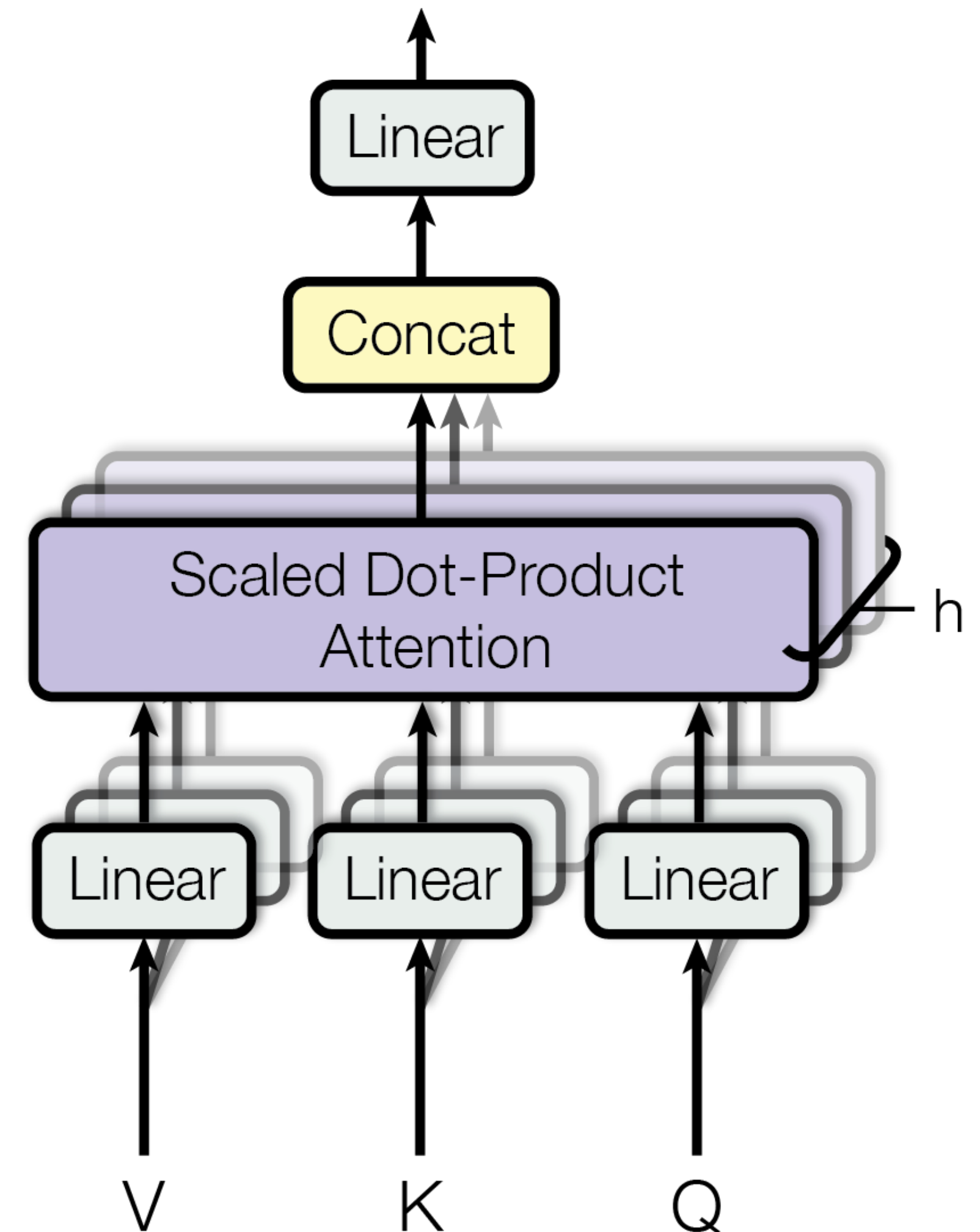
$$\mathbf{a}_i = \frac{(\mathbf{W}_i^Q \mathbf{q})(\mathbf{W}_i^K \mathbf{K})}{\sqrt{d/H}}$$

- Compute attention weights separately for each sub-vector

$$\alpha_i = \mathbf{softmax}(\mathbf{a}_i) \quad \tilde{h}_i^\ell = \alpha(V \mathbf{W}_i^V)$$

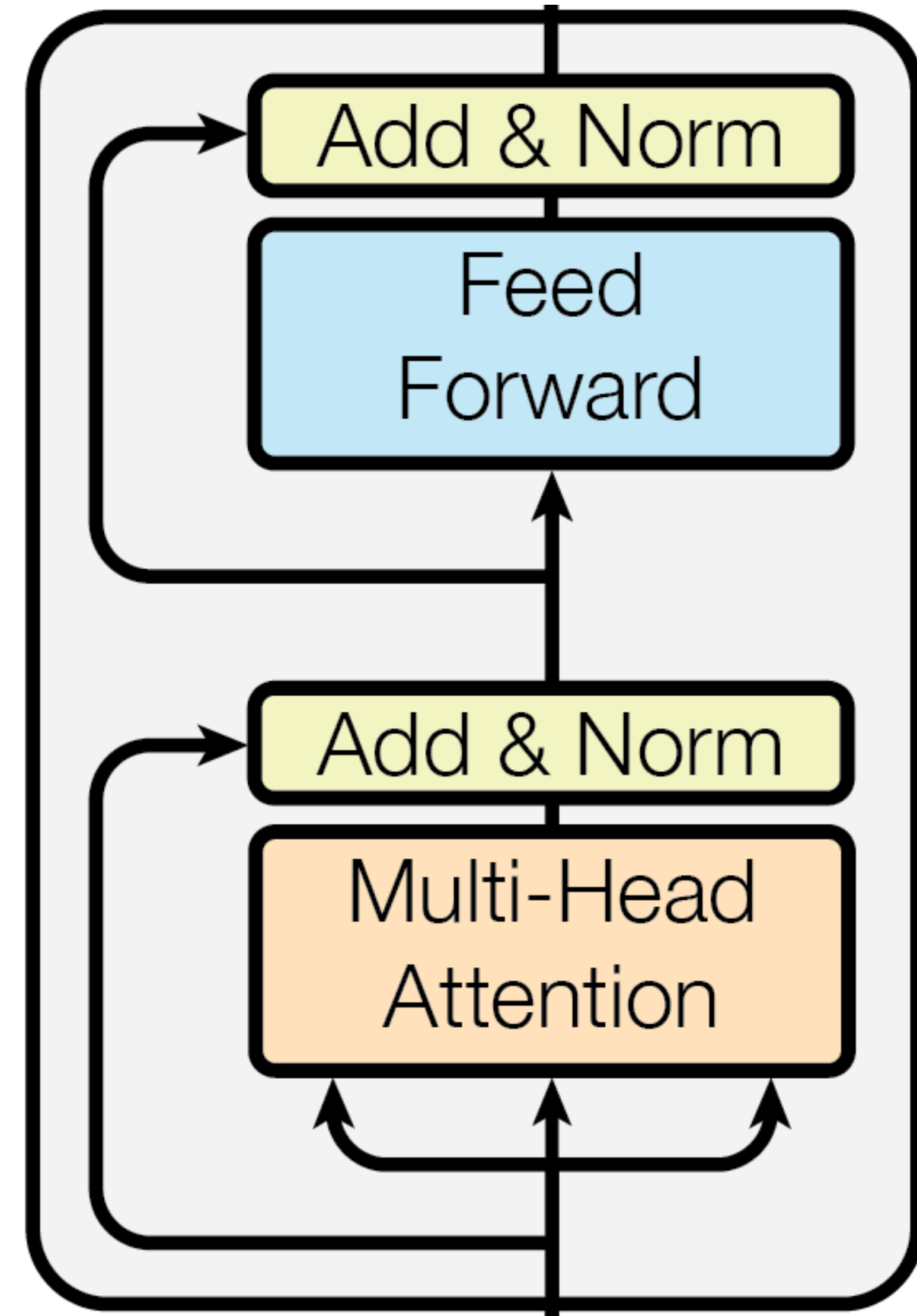
- Concatenate sub-vectors for each head and project

$$\tilde{h}^\ell = W^O [\tilde{h}_0^\ell; \dots; \tilde{h}_i^\ell; \dots; \tilde{h}_H^\ell]$$



Transformer Block

- Self-attention is the main innovation of the popular **transformer** model!
- Each transformer block receives as input the outputs of the previous layer at every time step
- Each block is composed of a multi-headed attention, a layer normalisation, a feedforward network, and another layer normalisation
- There are residual connections before every normalisation layer
- Layer normalisation + residual connections don't add capacity, but make training easier



Full Transformer

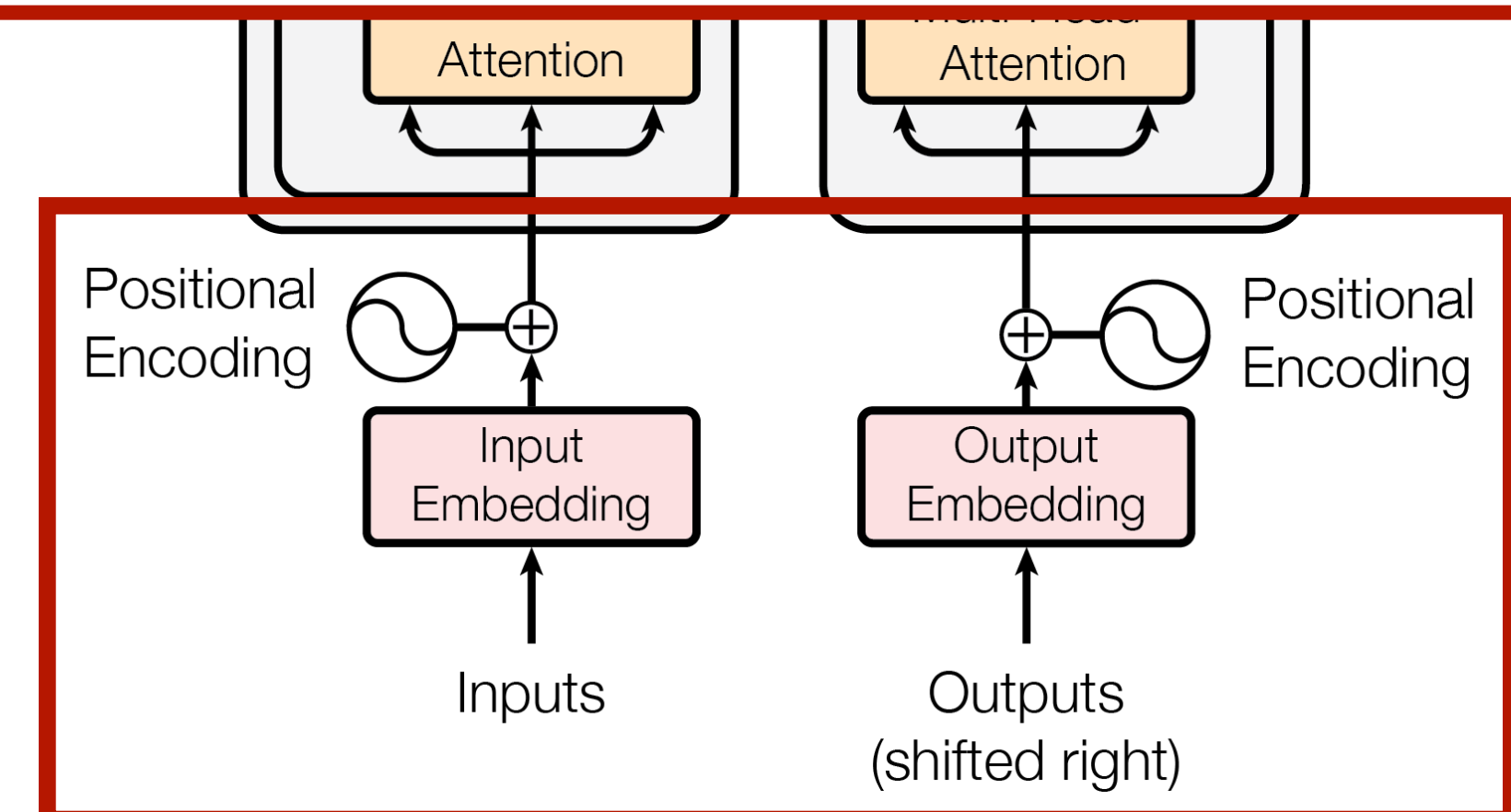
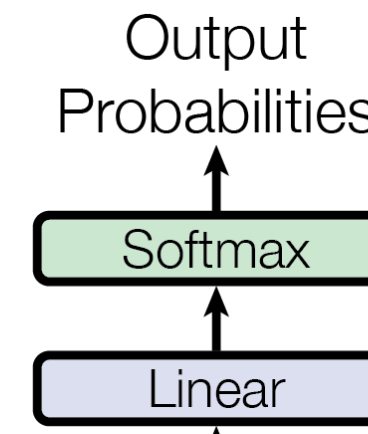
- Full transformer encoder is multiple cascaded transformer blocks

Recurrent models provided word order information

- N

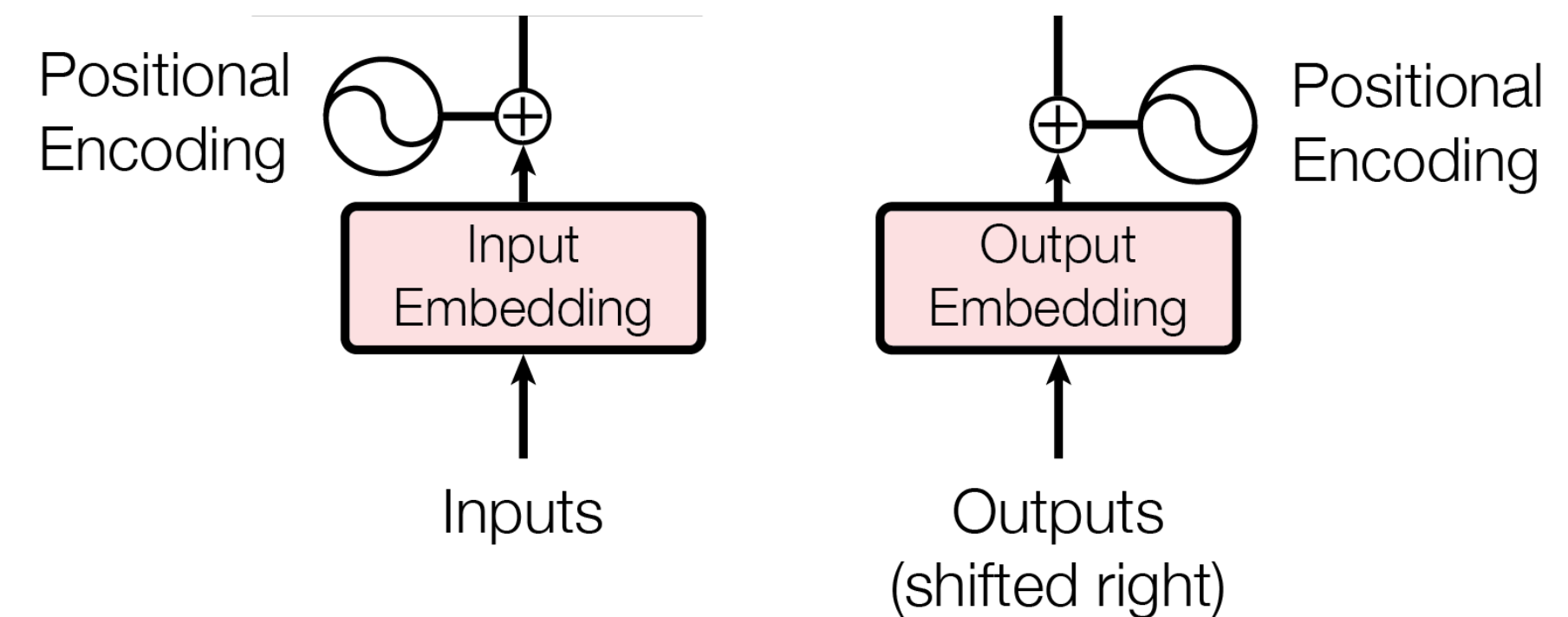
Does self-attention provide word order information?

- Transformer decoder (right) similar to encoder
 - second attention layer to compute weighted average of encoder states before FFN



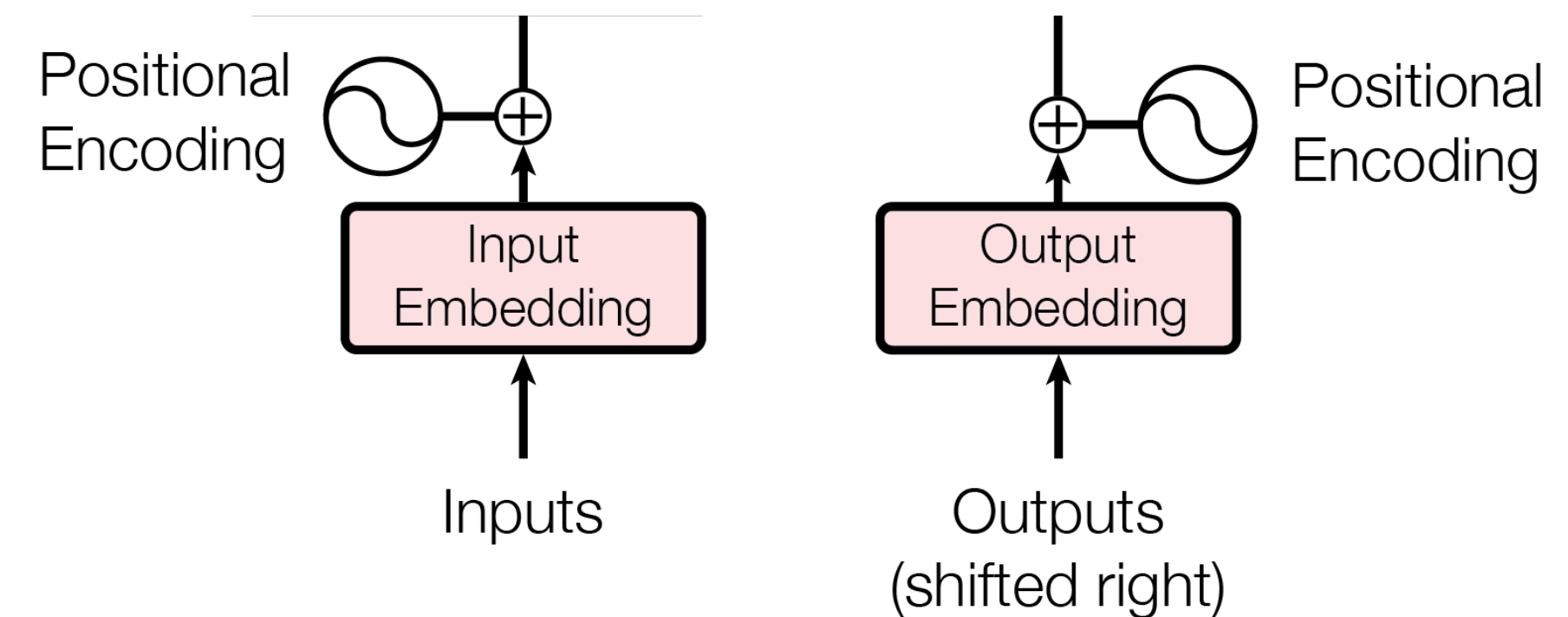
Position Embeddings

- Self-attention provides no word order information
 - Computes weighted average over set of vectors
- Word order is pretty crucial to understanding language
 - How do we fix this?
- Add an additional embedding to the input word that represents a position in the sequence



Position Embeddings

- Self-attention provides no word order information
 - Computes weighted average over set of vectors
- Word order is pretty crucial to understanding language
 - How do we fix this?
- Add an additional embedding to the input word that represents a position in the sequence



- Early position embeddings encoded a sinusoid function that was offset by a phase shift proportional to sequence position
- **In practice, position embeddings are learned scratch or more modern methods are used (e.g., Rotary position embeddings, AliBi)**

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 as primary function, but require position embeddings to capture sequence order

References

- Paperno, D., Kruszewski, G., Lazaridou, A., Pham, Q.N., Bernardi, R., Pezzelle, S., Baroni, M., Boleda, G., & Fernández, R. (2016). The LAMBADA dataset: Word prediction requiring a broad discourse context. *ArXiv, abs/1606.06031*.
- 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*.

Deep Learning for Natural Language Processing

Antoine Bosselut

EPFL



Part 2: Recurrent Neural Networks for Sequence Modeling

Section Outline

- **Background:** Language Modeling, Feedforward Neural Networks, Backpropagation
- **Content - Models:** Recurrent Neural Networks, Encoder-Decoders
- **Content - Algorithms:** Backpropagation through Time, Vanishing Gradients

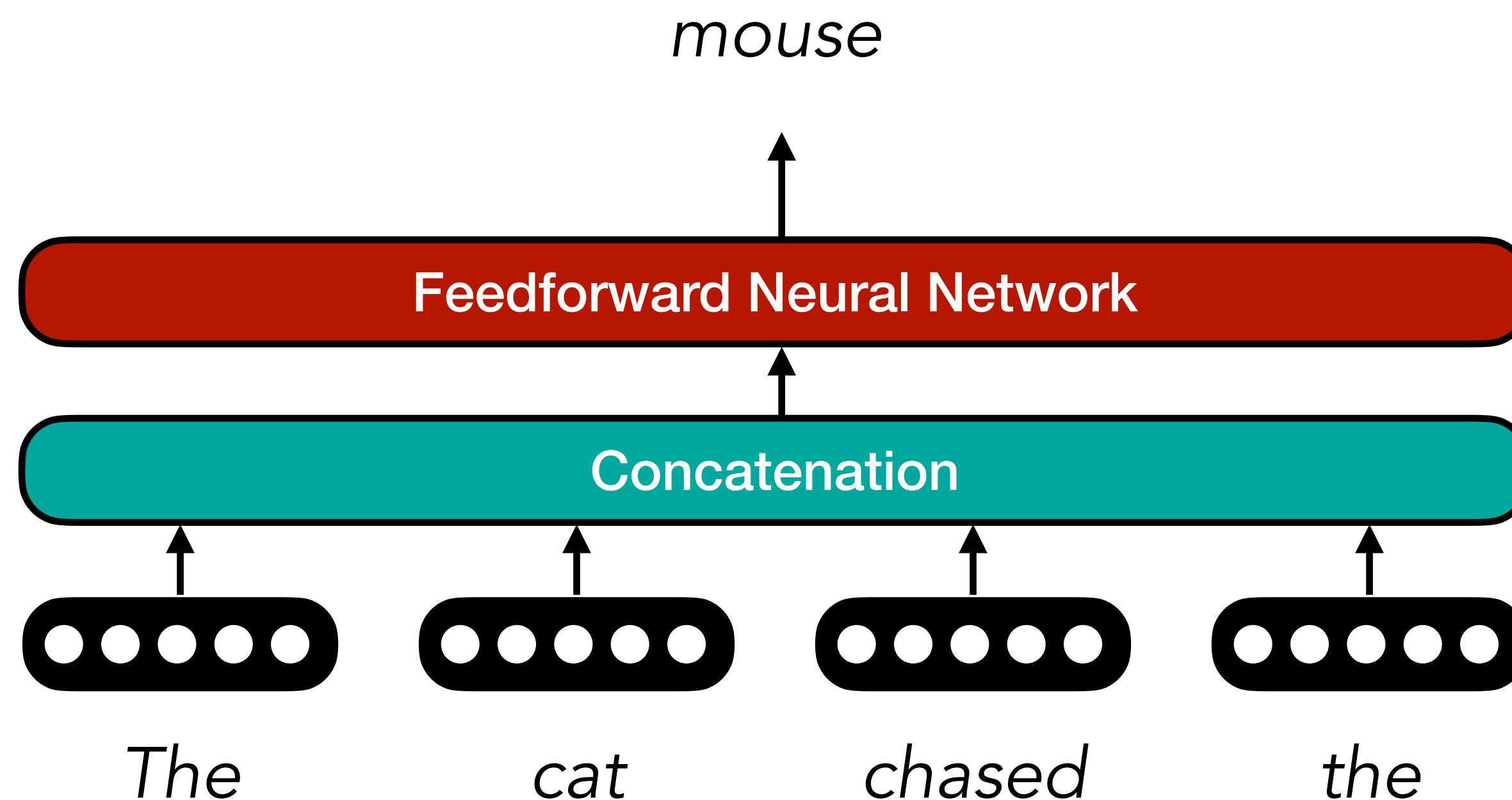
Language Modeling

- Given a subsequence, predict the next word: *The cat chased the _____*

Fixed Context Language Models

- Given a subsequence, predict the next word: *The cat chased the _____*

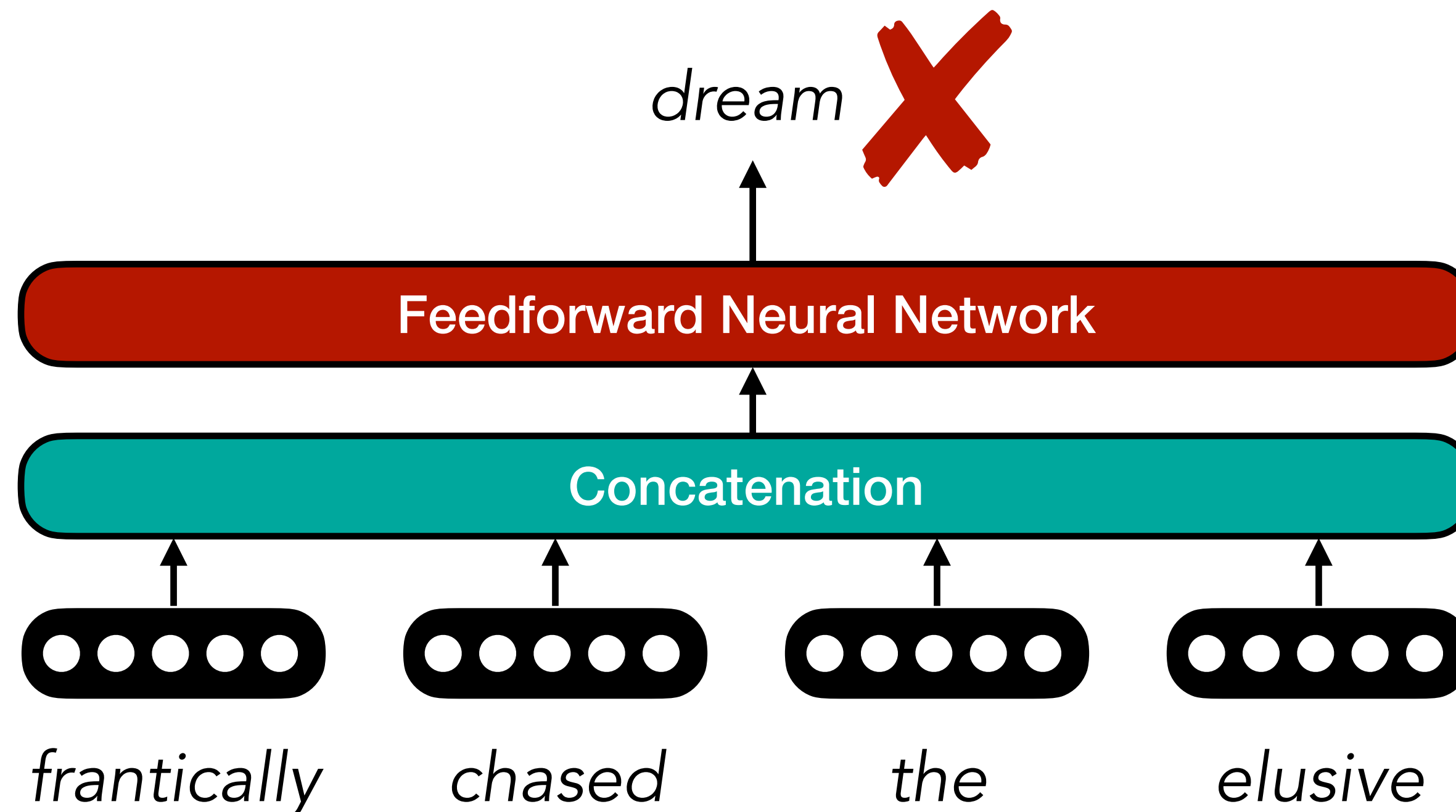
$$P(y) = \mathbf{softmax}(b_o + \mathbf{W}_o \mathbf{tanh}(b_h + \mathbf{W}_h x))$$



Fixed Context Language Models

- Given a subsequence, predict the next word:

The starving cat frantically chased the elusive _____



The starving cat

frantically

chased

the

elusive

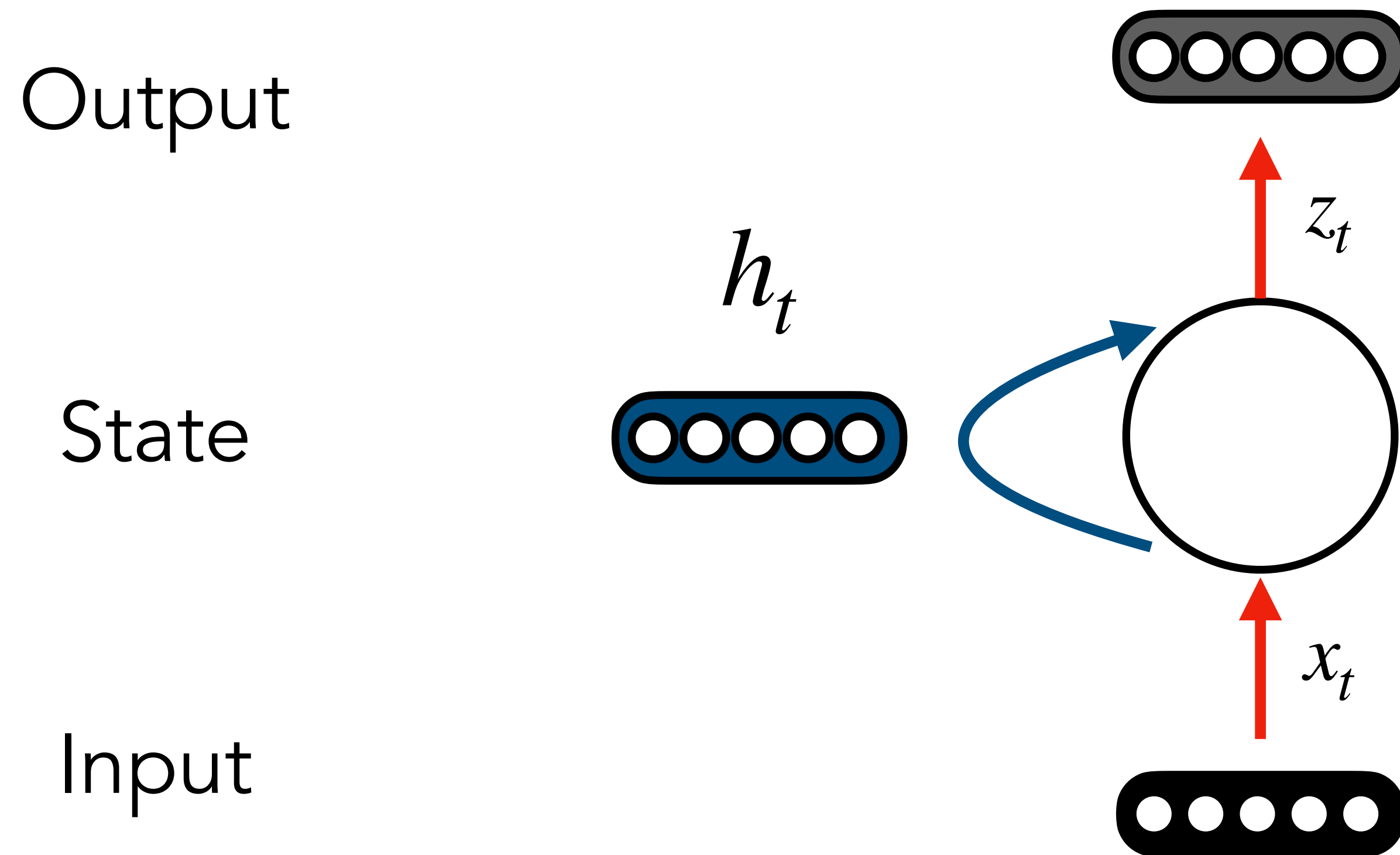
Problem

Fixed context windows limit language modelling capacity

How can we extend to arbitrary length sequences?

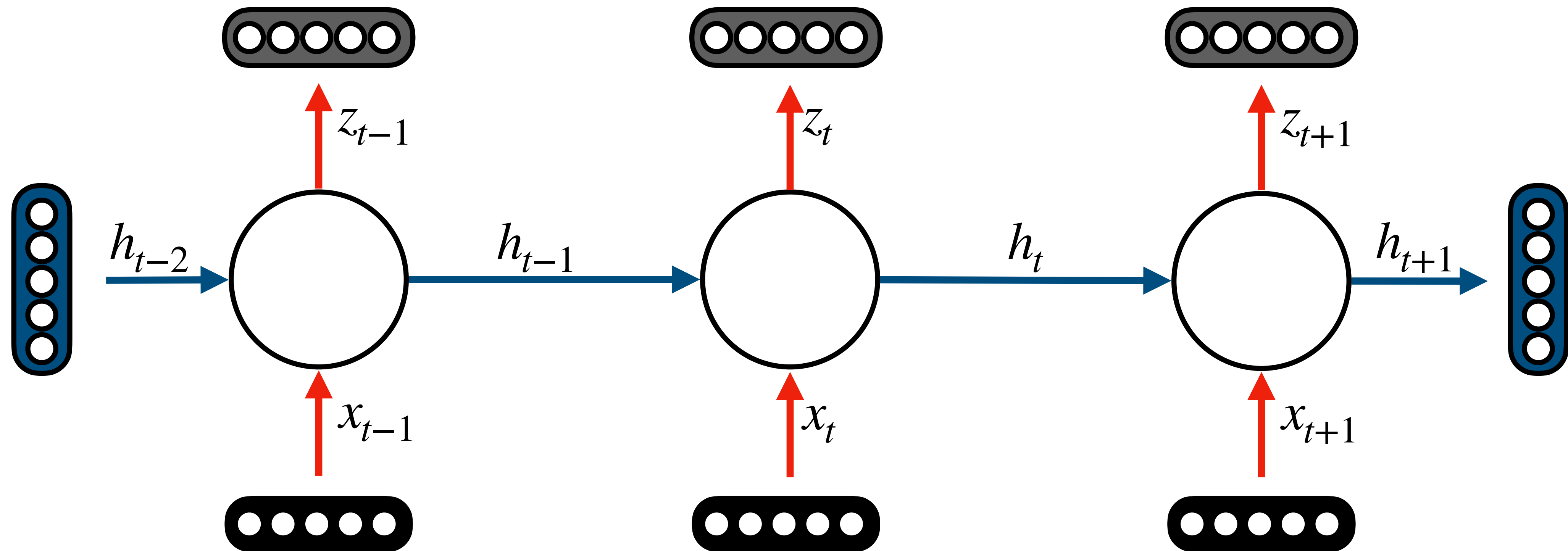
Recurrent Neural Networks

- **Solution:** Recurrent neural networks — NNs with feedback loops



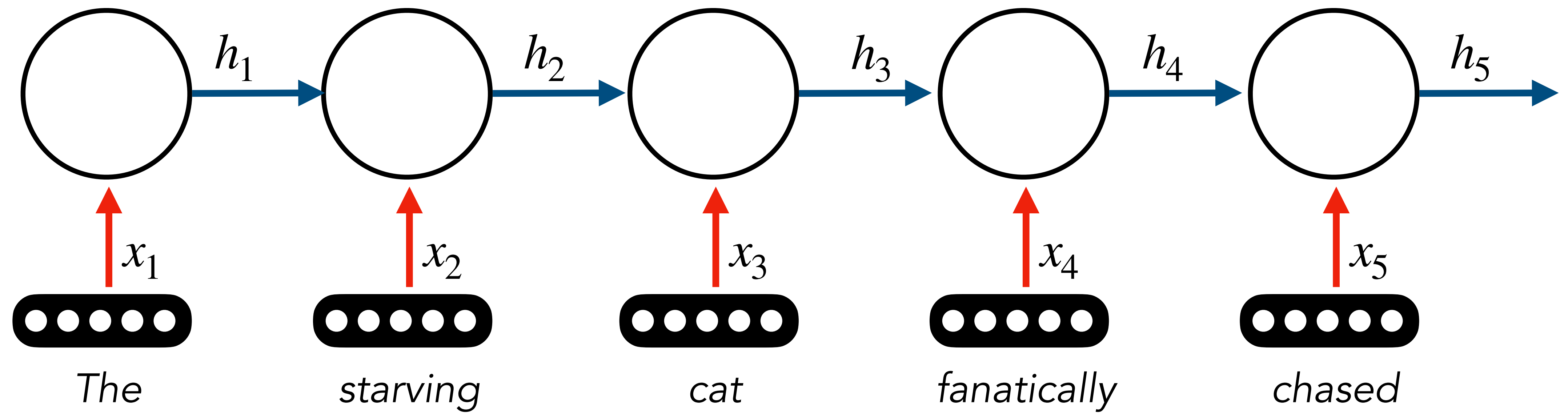
Unrolling the RNN

Unrolling the RNN across all time steps gives full computation graph

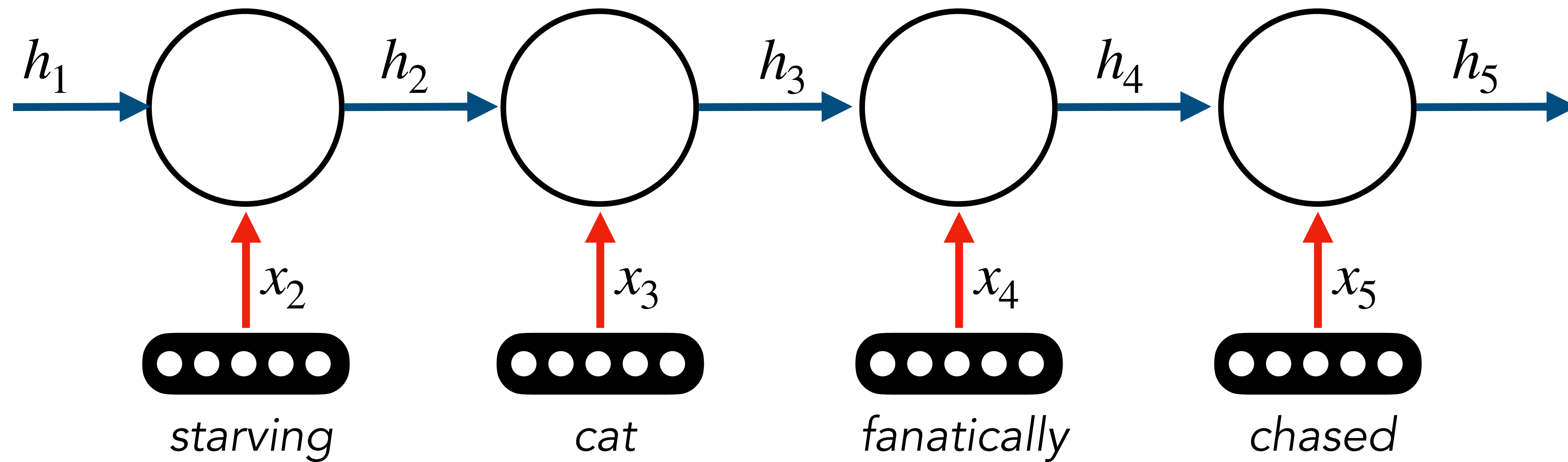


Allows for learning from entire sequence history, regardless of length

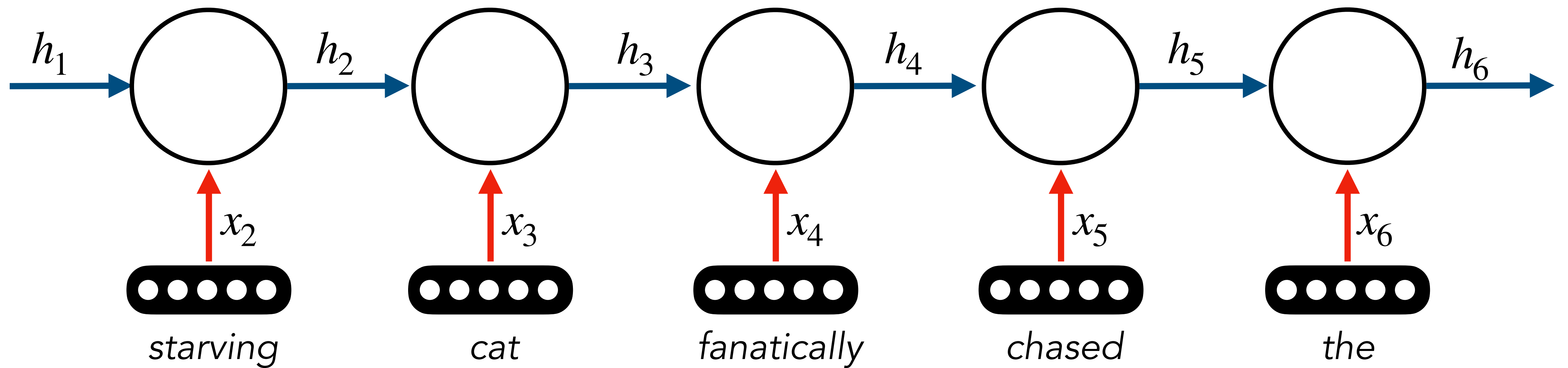
Unrolling the RNN



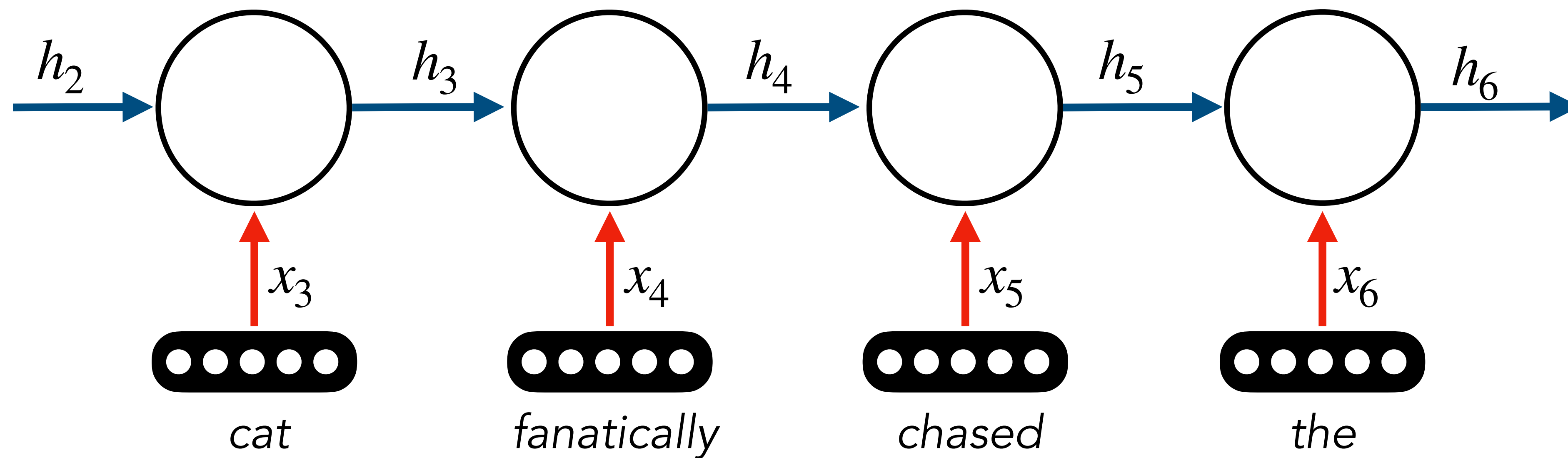
Unrolling the RNN



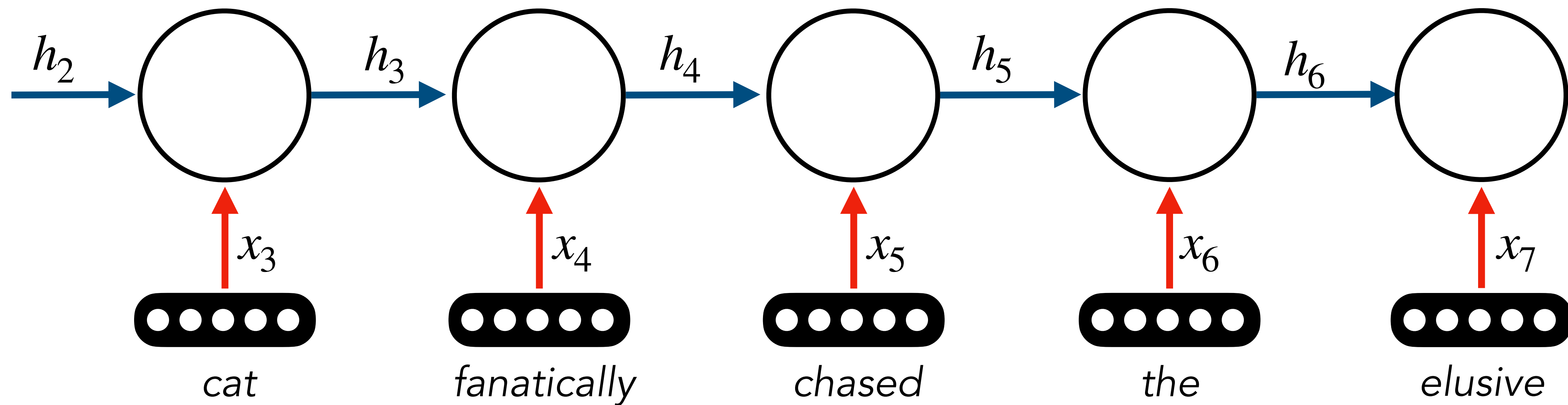
Unrolling the RNN



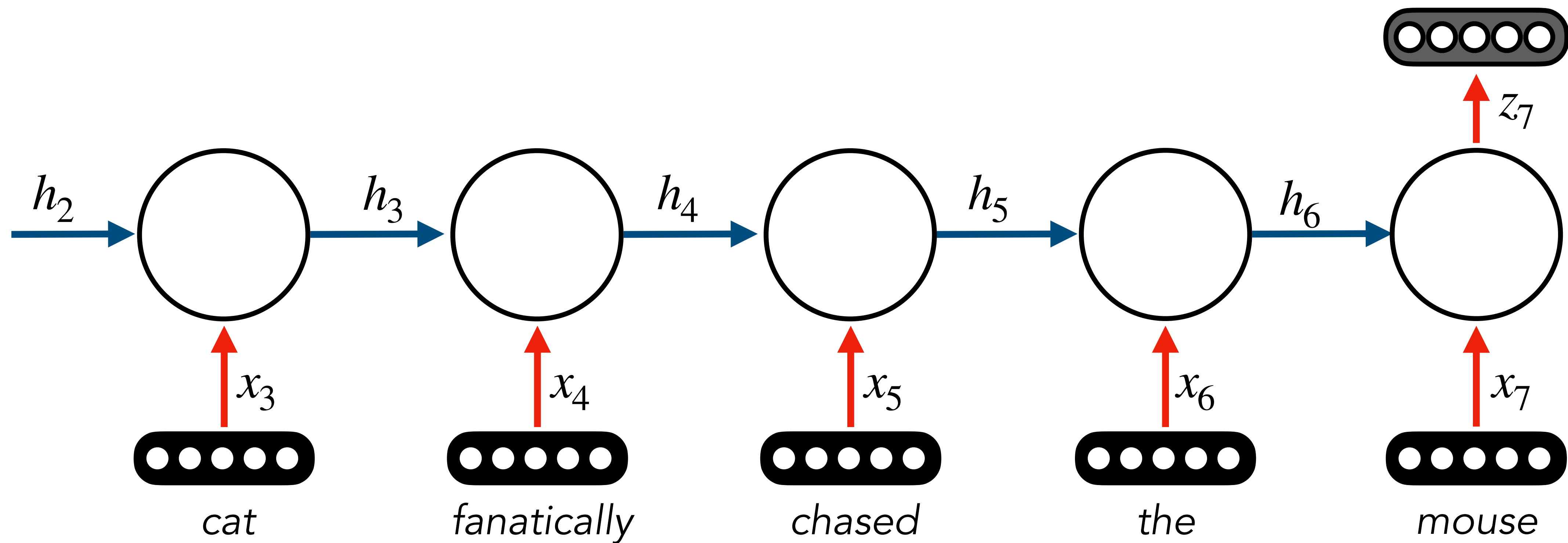
Unrolling the RNN



Unrolling the RNN



Unrolling the RNN



Classification

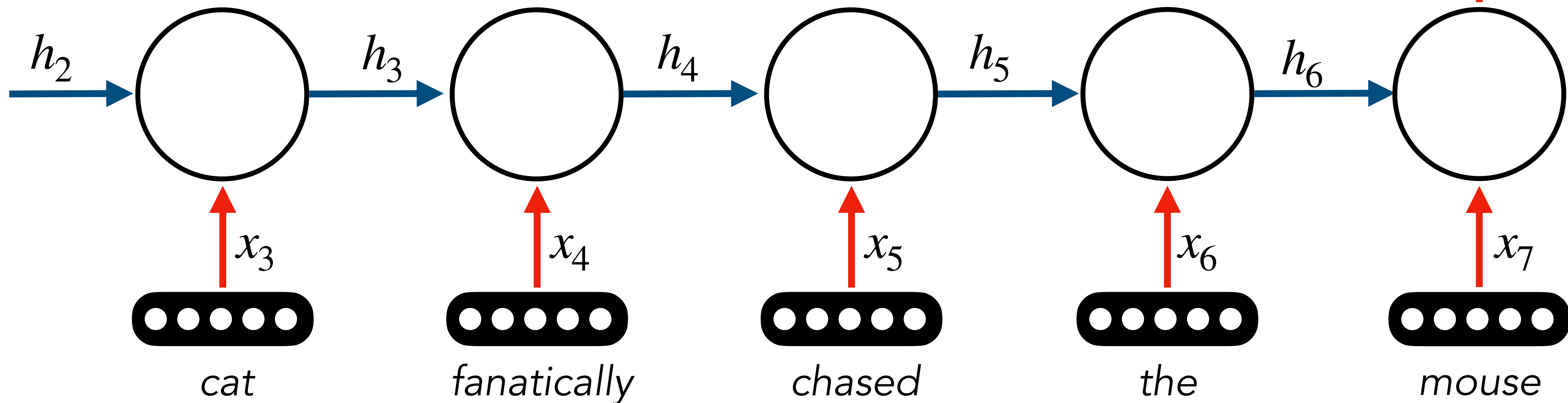
- Classifier is just an output projection followed by a softmax!

Binary

$$P(y) = \sigma(W_o z_T)$$

Multi-class

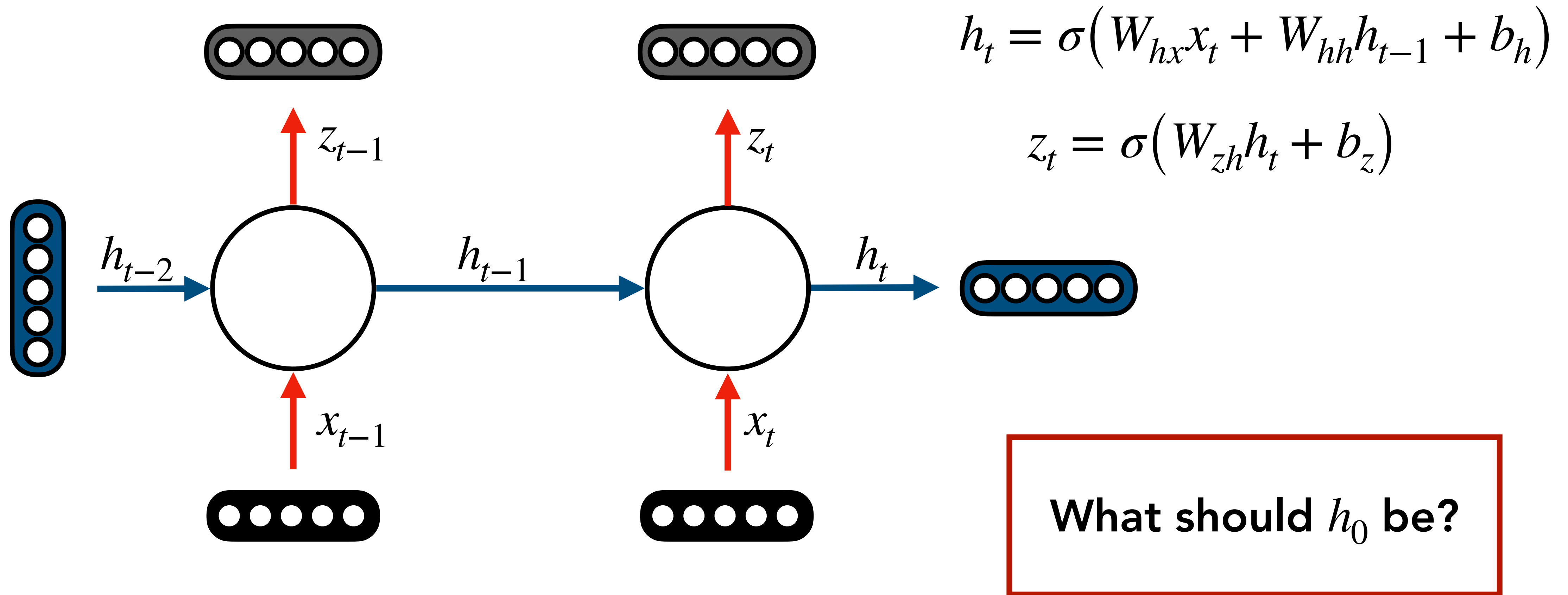
$$P(y) = \mathbf{softmax}(W_o z_T)$$



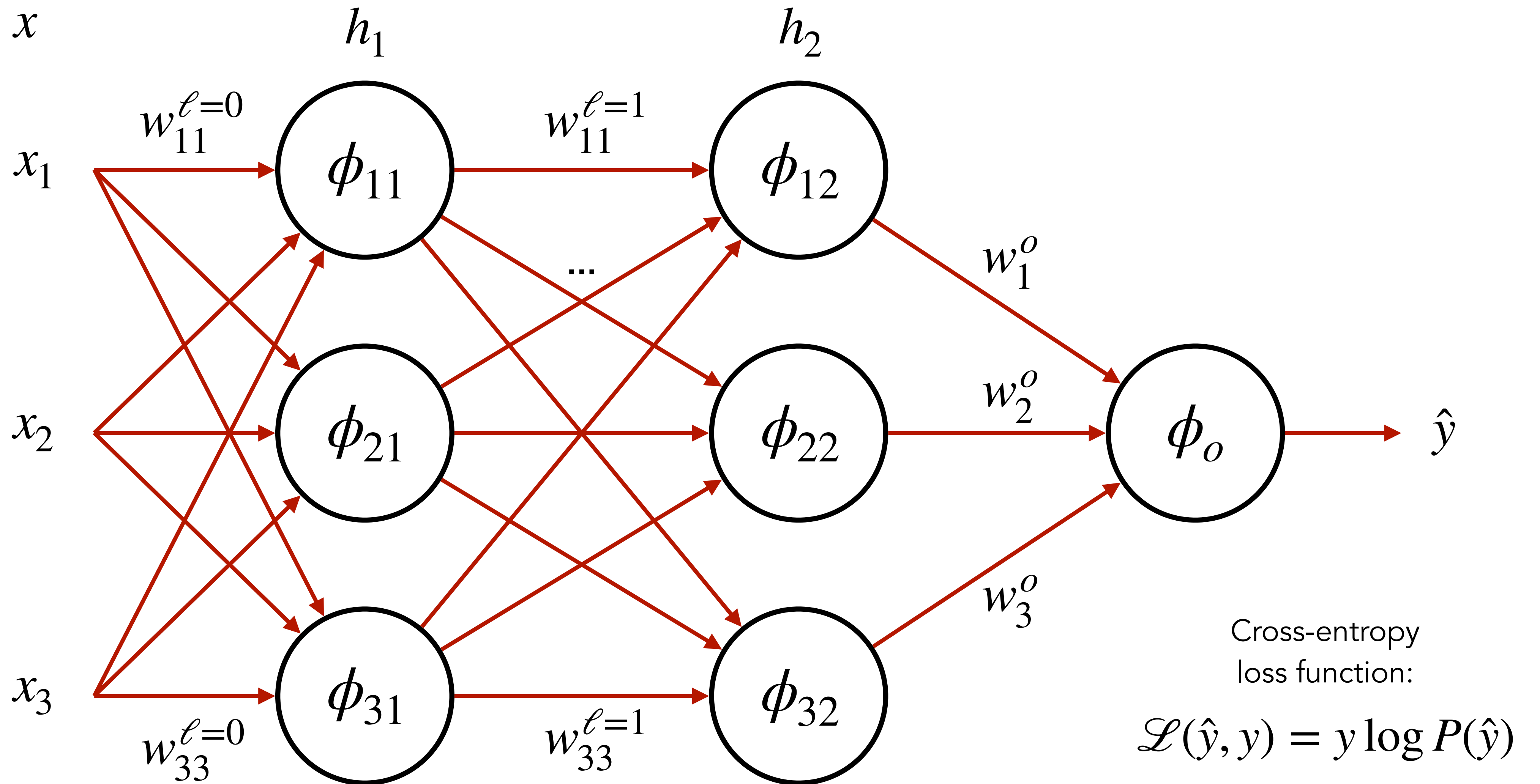
Question

Why would you use the output of the last recurrent unit as the one to predict a label?

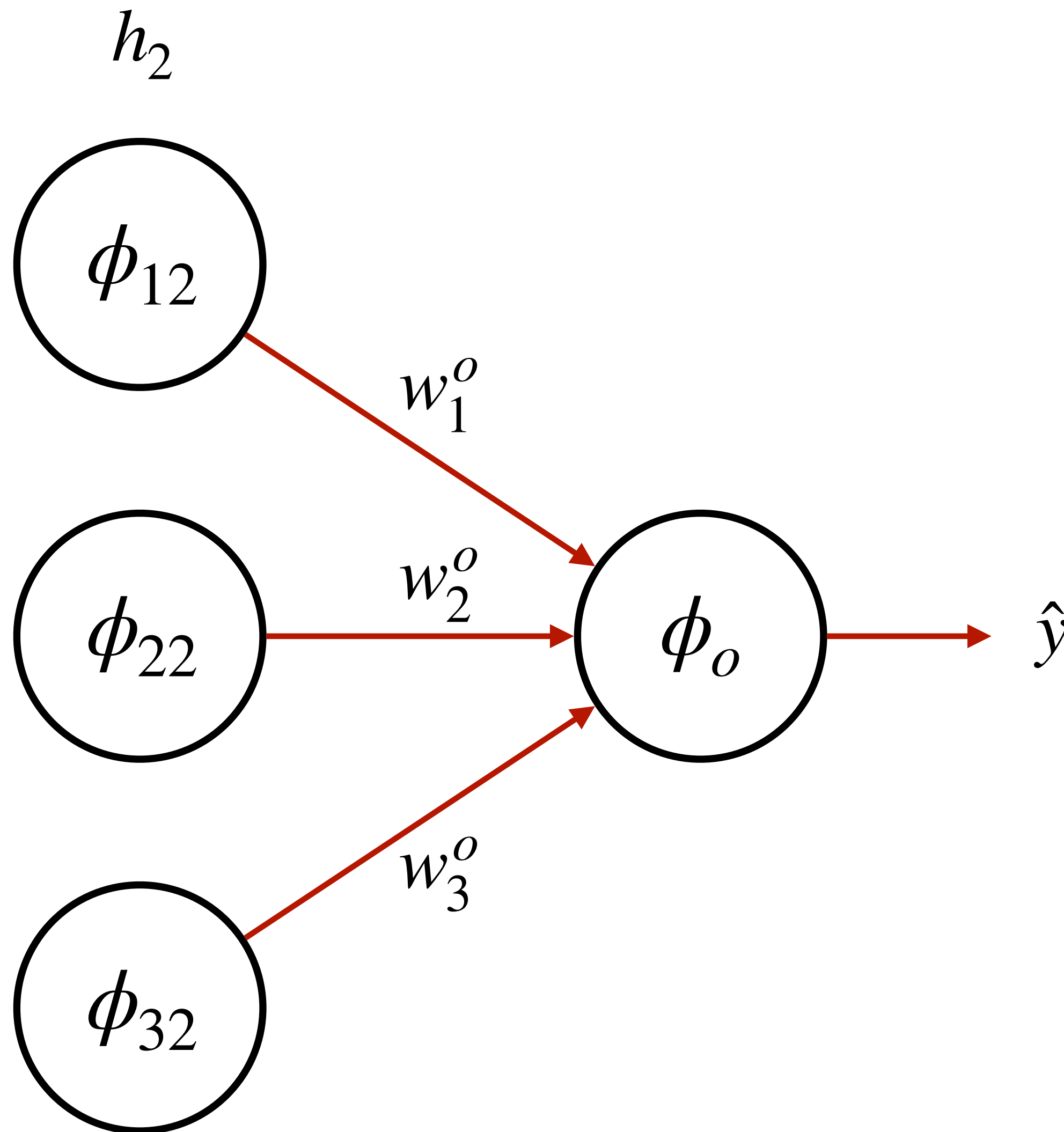
Classical RNN: Elman Network



Backpropagation Review: FFNs



Backpropagation Review: FFNs



$$\mathcal{L}(\hat{y}, y) = y \log P(\hat{y})$$

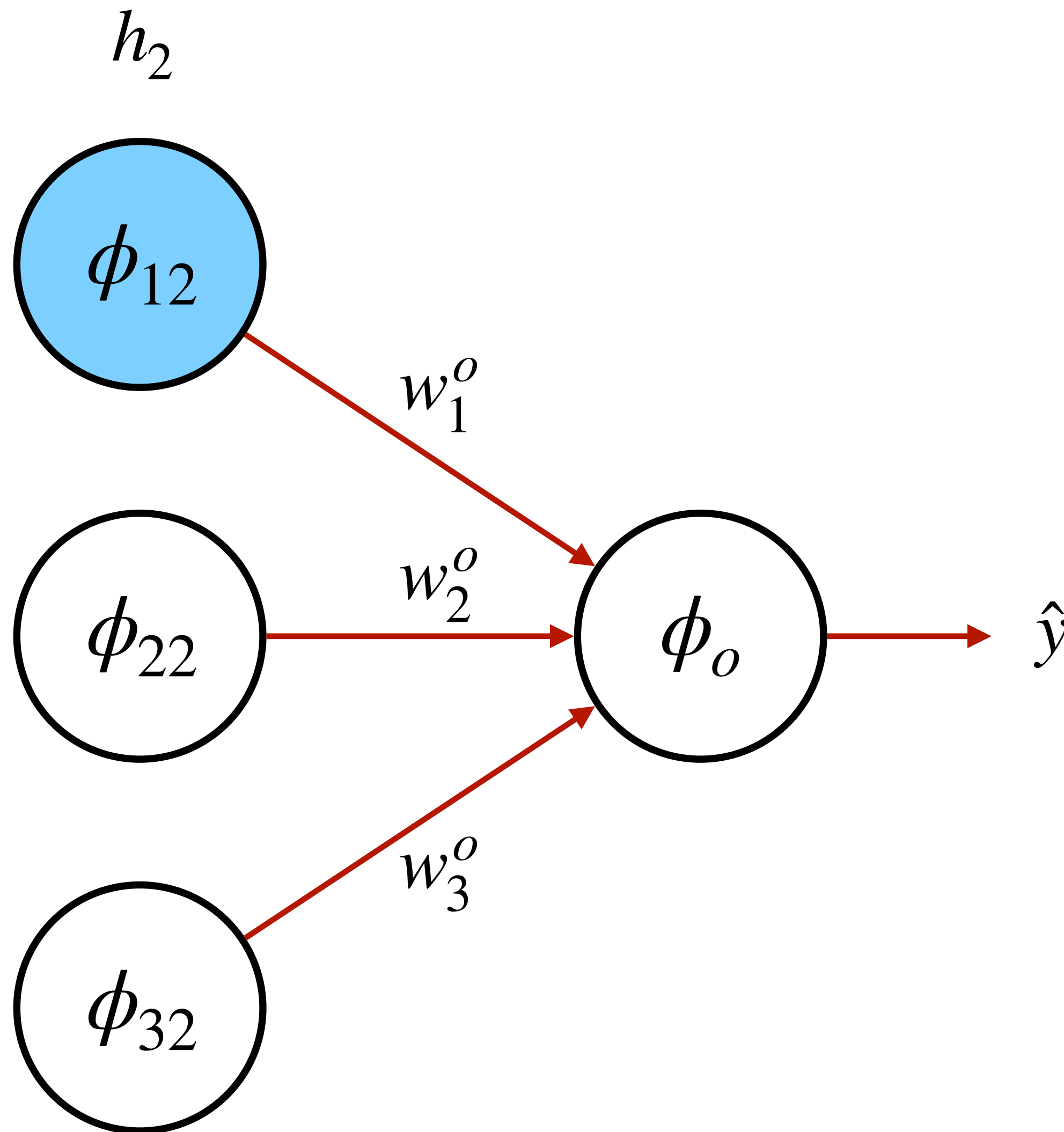
$$\hat{y} = \phi_o(u)$$

$$u = w_1^o \times \phi_{12}(\cdot) + w_2^o \times \phi_{22}(\cdot) + w_3^o \times \phi_{32}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

Backpropagation Review: FFNs



$$\mathcal{L}(\hat{y}, y) = y \log P(\hat{y})$$

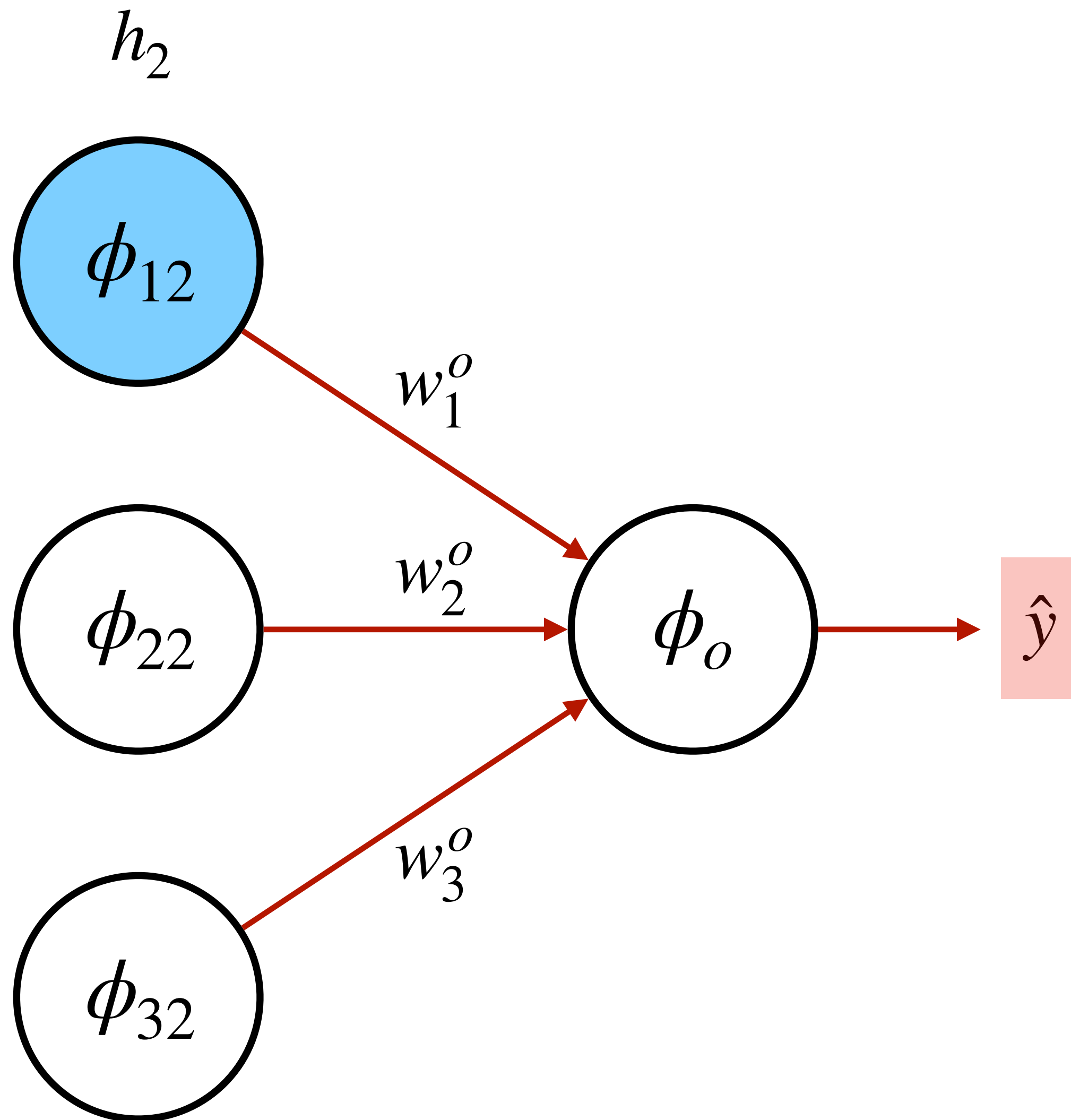
$$\hat{y} = \phi_o(u)$$

$$u = w_1^o \times \phi_{12}(\cdot) + w_2^o \times \phi_{22}(\cdot) + w_3^o \times \phi_{32}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

Backpropagation Review: FFNs



$$\mathcal{L}(\hat{y}, y) = y \log P(\hat{y})$$

$$\hat{y} = \phi_o(u)$$

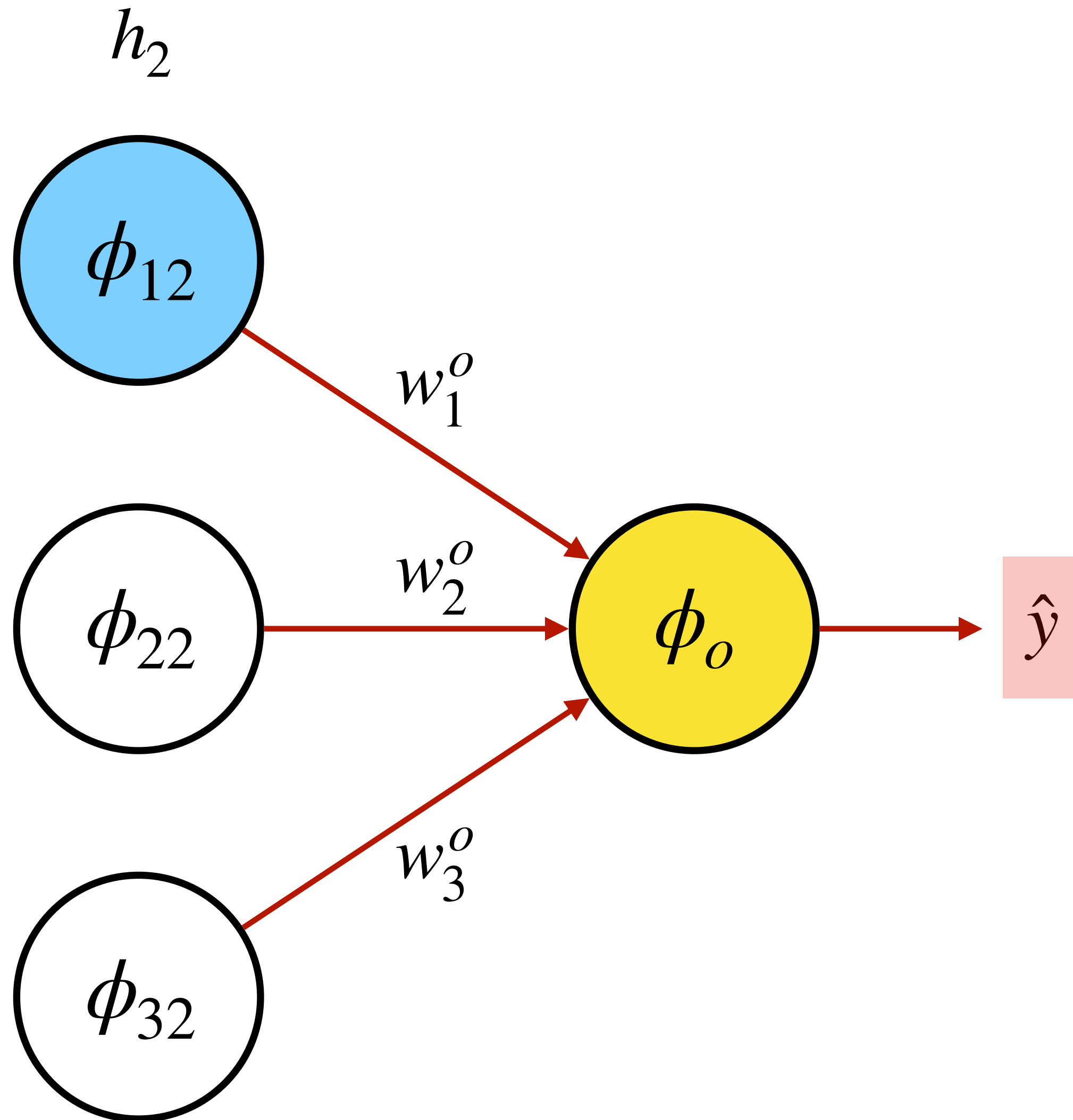
$$u = w_1^o \times \phi_{12}(\cdot) + w_2^o \times \phi_{22}(\cdot) + w_3^o \times \phi_{32}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

Depends on label y

Backpropagation Review: FFNs



$$\mathcal{L}(\hat{y}, y) = y \log P(\hat{y}) + (1 - y) \log P(1 - \hat{y})$$

$$\hat{y} = \phi_o(u)$$

$$u = w_1^o \times \phi_{12}(\cdot) + w_2^o \times \phi_{22}(\cdot) + w_3^o \times \phi_{32}(\cdot)$$

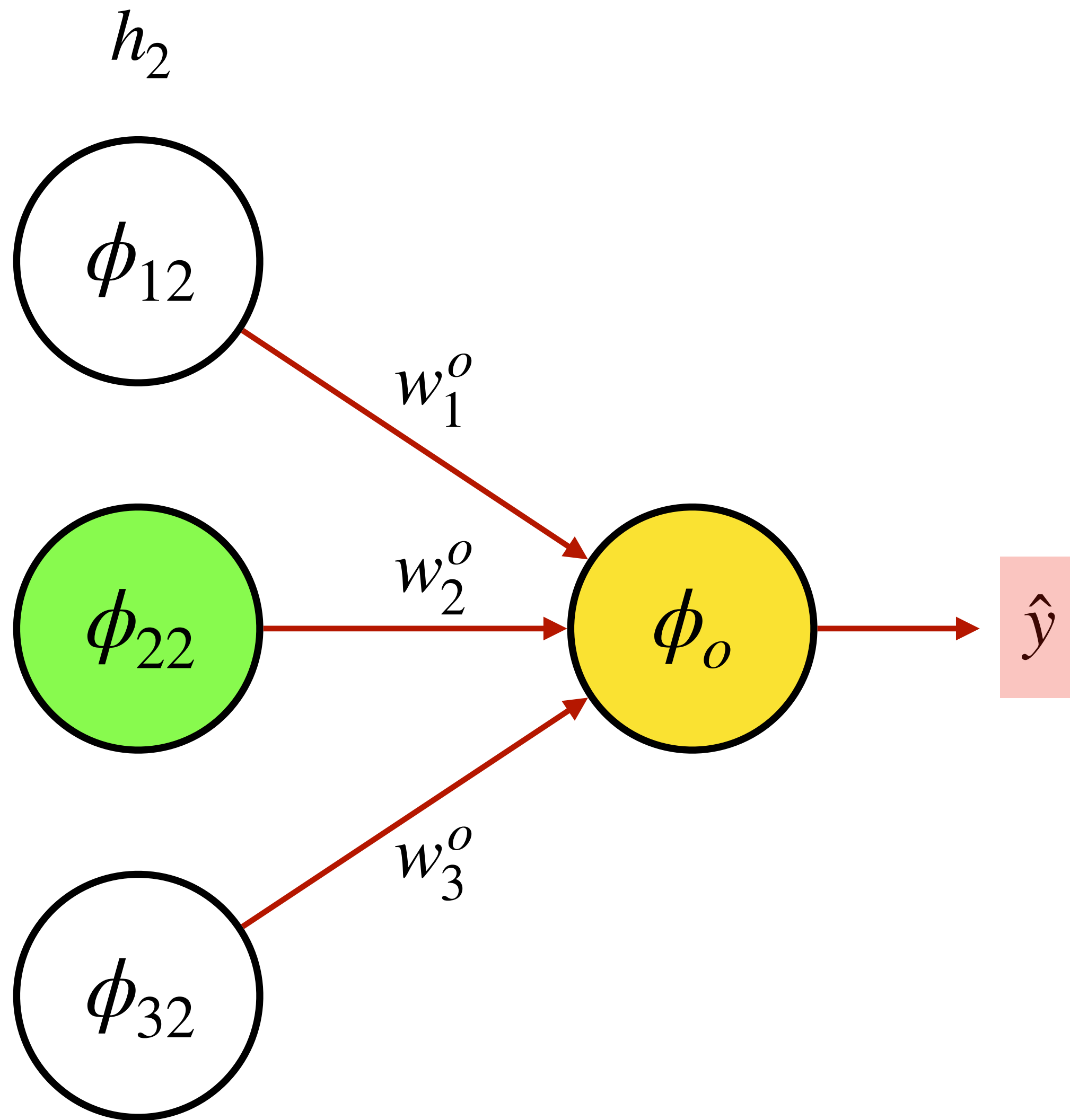
$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

Depends on label y

Depends on ϕ_o

Backpropagation Review: FFNs



$$\mathcal{L}(\hat{y}, y) = y \log P(\hat{y}) + (1 - y) \log P(1 - \hat{y})$$

$$\hat{y} = \phi_o(u)$$

$$u = w_1^o \times \phi_{12}(\cdot) + w_2^o \times \phi_{22}(\cdot) + w_3^o \times \phi_{32}(\cdot)$$

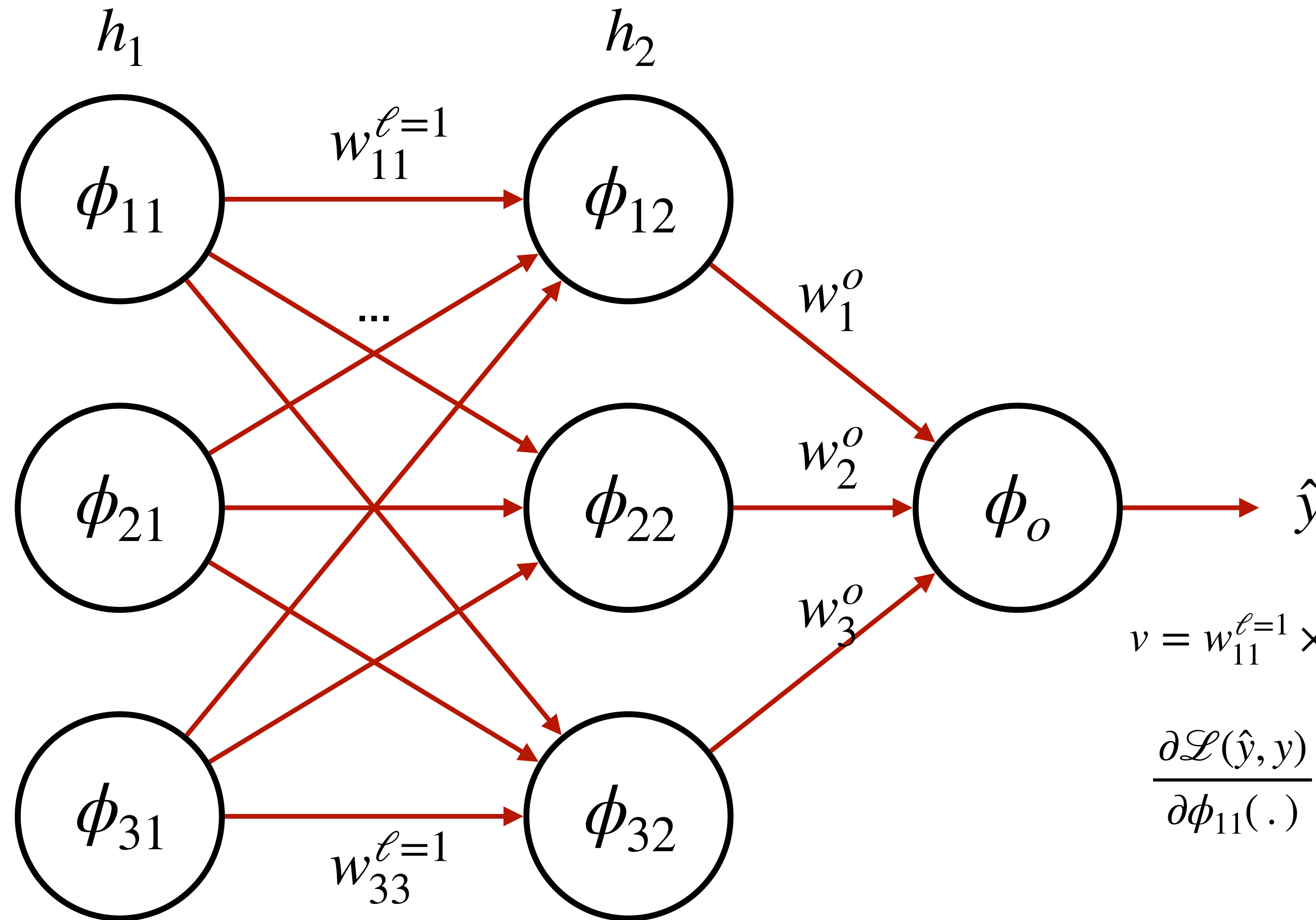
$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{22}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{22}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_2^o$$

Depends on label y

Depends on ϕ_o

Backpropagation Review: FFNs



$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

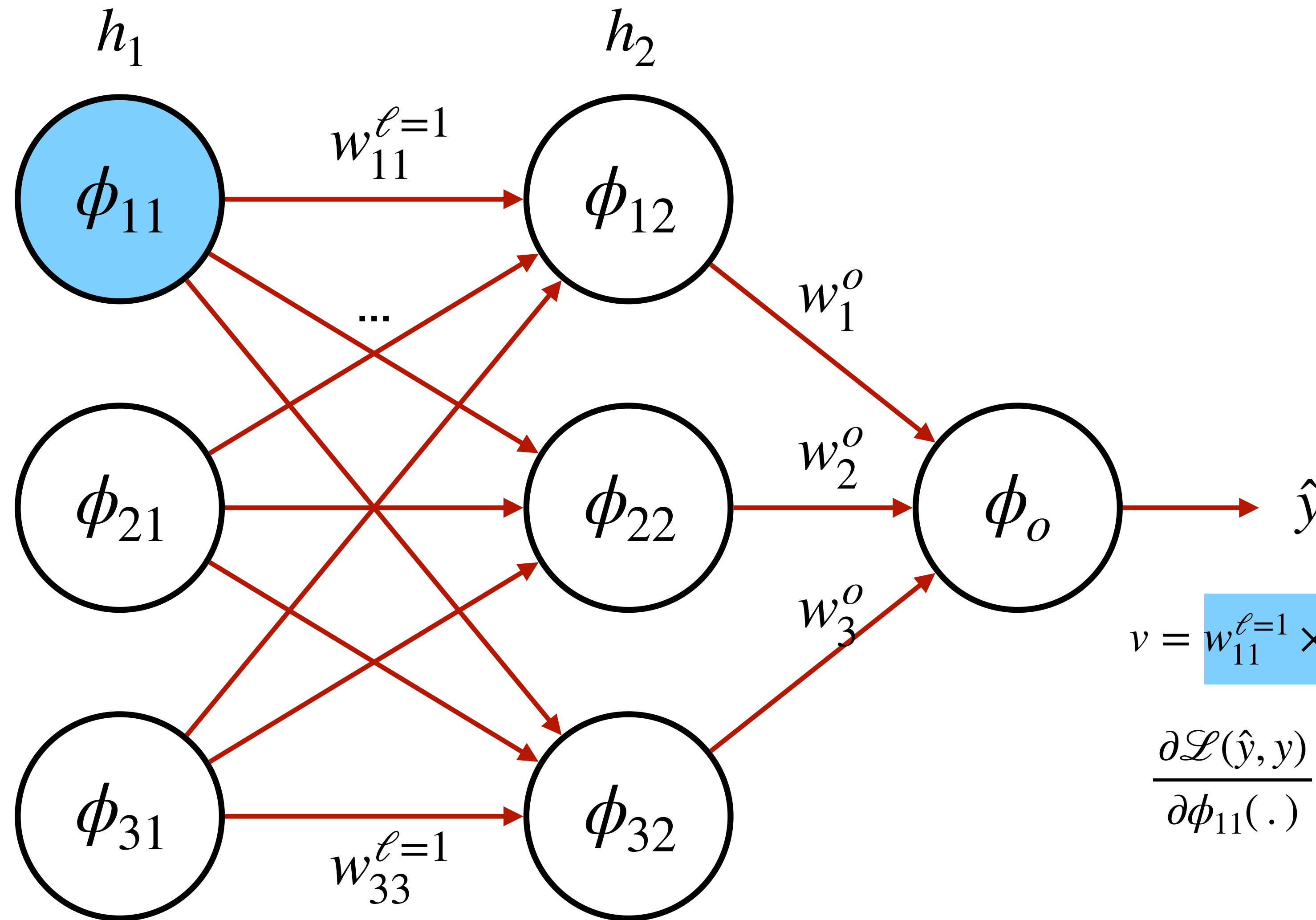
$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

$$v = w_{11}^{\ell=1} \times \phi_{11}(\cdot) + w_{21}^{\ell=1} \times \phi_{21}(\cdot) + w_{31}^{\ell=1} \times \phi_{31}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{11}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(v)} \frac{\partial \phi_{12}(v)}{\partial v} \frac{\partial v}{\partial \phi_{11}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o \frac{\partial \phi_{12}(v)}{\partial v} w_{11}^{\ell=1}$$

Backpropagation Review: FFNs



$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

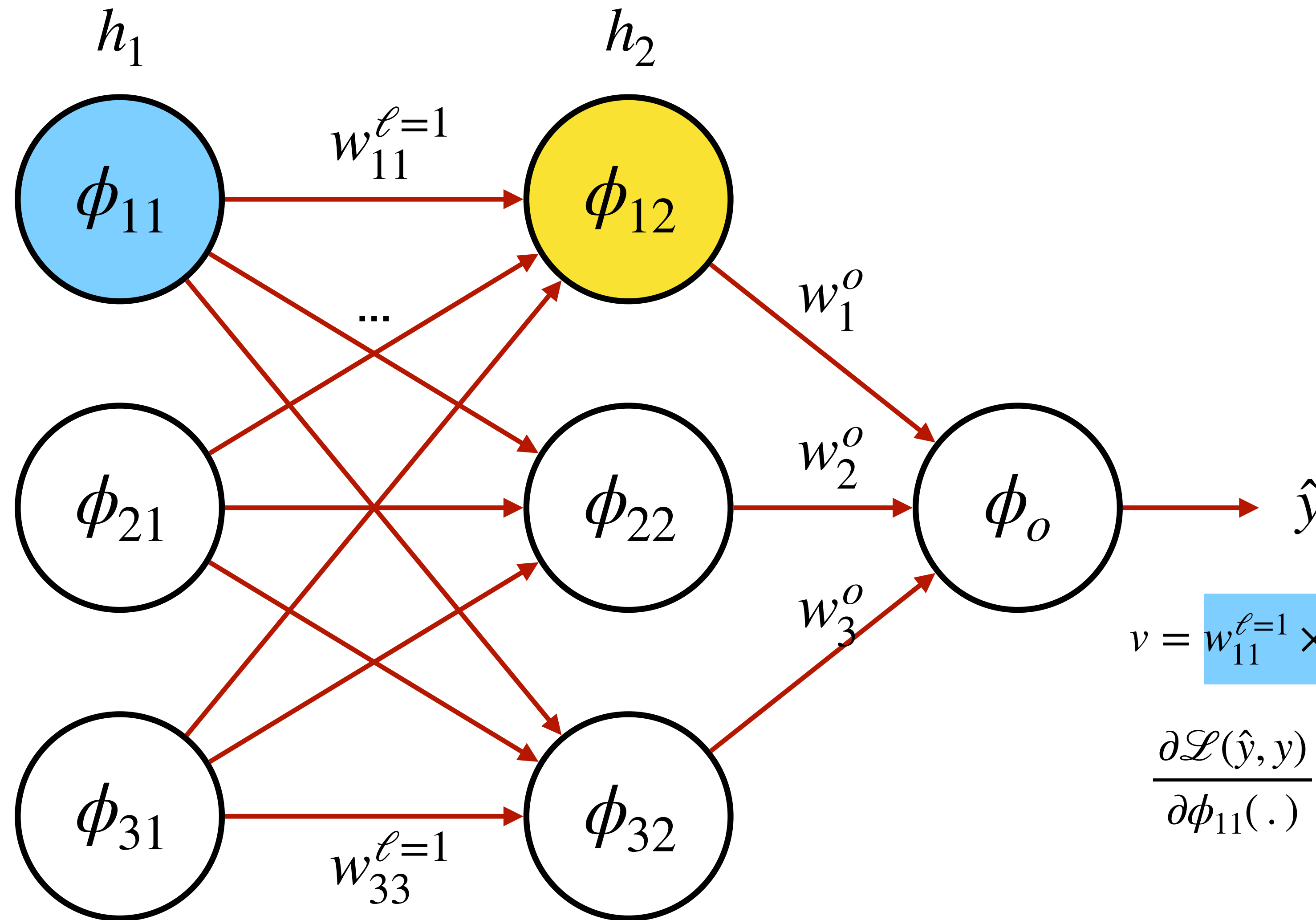
$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

$$v = w_{11}^{\ell=1} \times \phi_{11}(\cdot) + w_{21}^{\ell=1} \times \phi_{21}(\cdot) + w_{31}^{\ell=1} \times \phi_{31}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{11}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(v)} \frac{\partial \phi_{12}(v)}{\partial v} \frac{\partial v}{\partial \phi_{11}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o \frac{\partial \phi_{12}(v)}{\partial v} w_{11}^{\ell=1}$$

Backpropagation Review: FFNs



$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

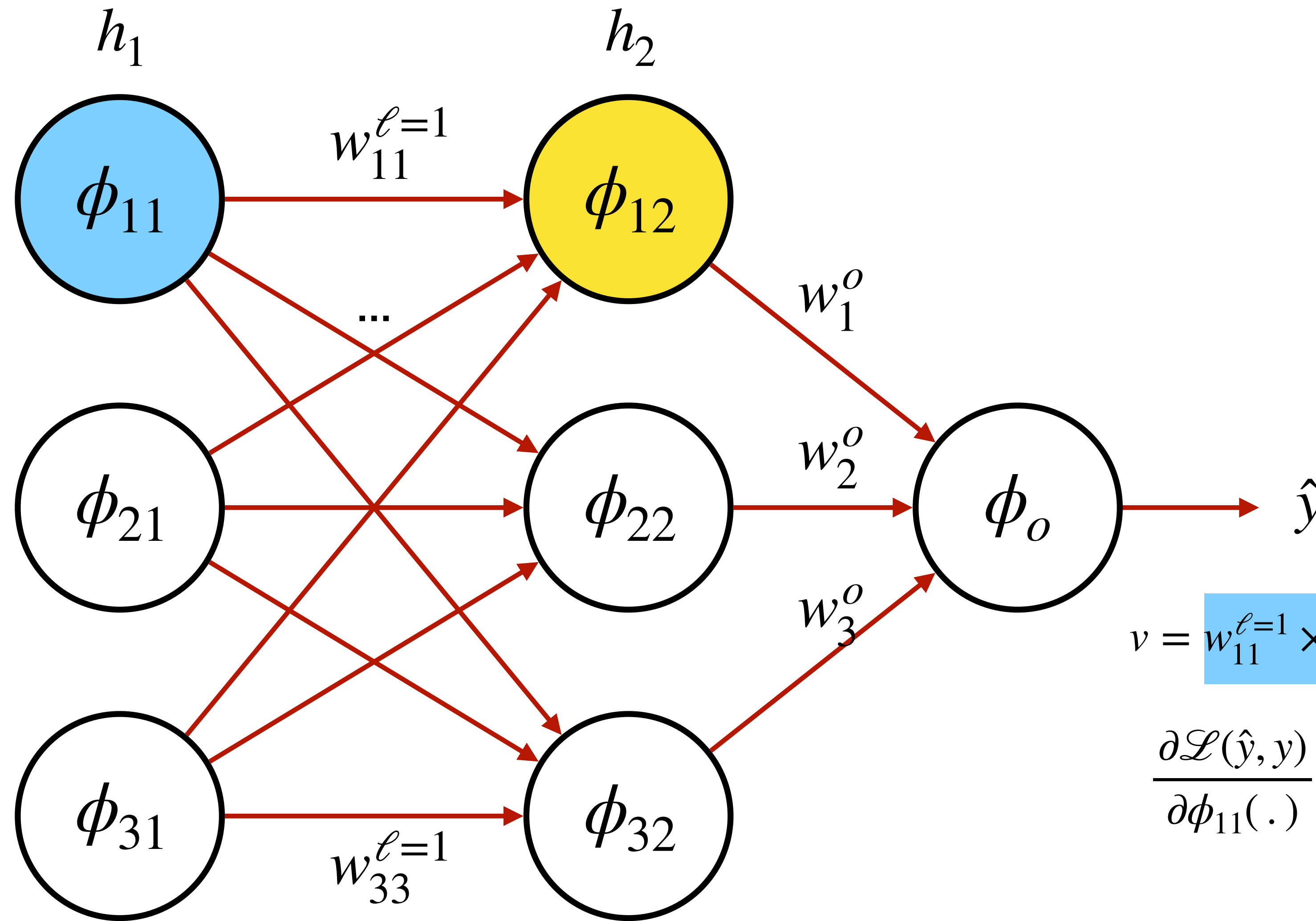
Depends on ϕ_{12}

$$v = w_{11}^{\ell=1} \times \phi_{11}(\cdot) + w_{21}^{\ell=1} \times \phi_{21}(\cdot) + w_{31}^{\ell=1} \times \phi_{31}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{11}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(v)} \frac{\partial \phi_{12}(v)}{\partial v} \frac{\partial v}{\partial \phi_{11}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o \frac{\partial \phi_{12}(v)}{\partial v} w_{11}^{\ell=1}$$

Backpropagation Review: FFNs



$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

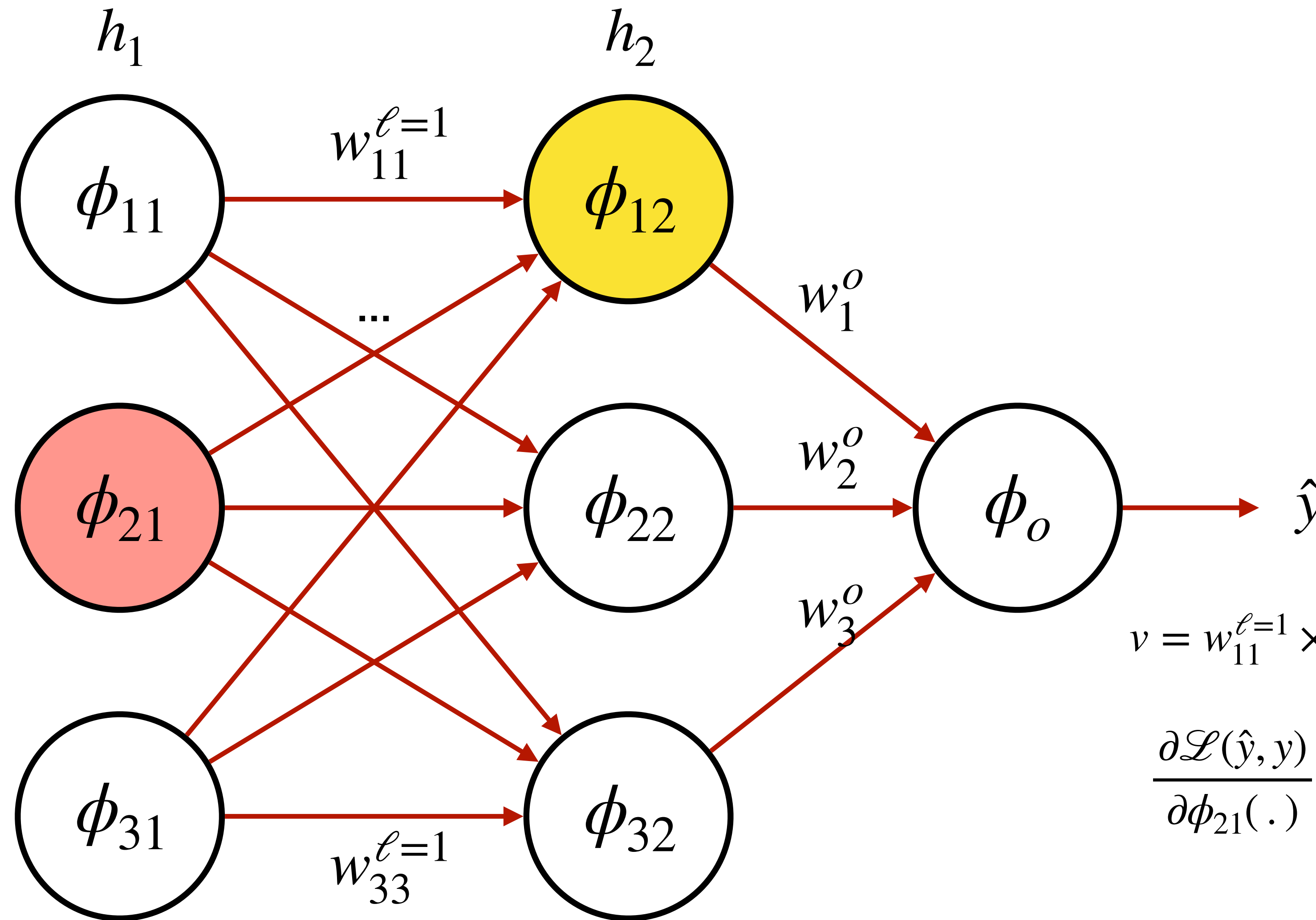
Depends on ϕ_{12}

$$v = w_{11}^{\ell=1} \times \phi_{11}(\cdot) + w_{21}^{\ell=1} \times \phi_{21}(\cdot) + w_{31}^{\ell=1} \times \phi_{31}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{11}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(v)} \frac{\partial \phi_{12}(v)}{\partial v} \frac{\partial v}{\partial \phi_{11}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o \frac{\partial \phi_{12}(v)}{\partial v} w_{11}^{\ell=1}$$

Backpropagation Review: FFNs



$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{12}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(\cdot)}$$

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o$$

Depends on ϕ_{12}

$$v = w_{11}^{\ell=1} \times \phi_{11}(\cdot) + w_{21}^{\ell=1} \times \phi_{21}(\cdot) + w_{31}^{\ell=1} \times \phi_{31}(\cdot)$$

$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{21}(\cdot)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{12}(v)} \frac{\partial \phi_{12}(v)}{\partial v} \frac{\partial v}{\partial \phi_{21}(\cdot)}$$

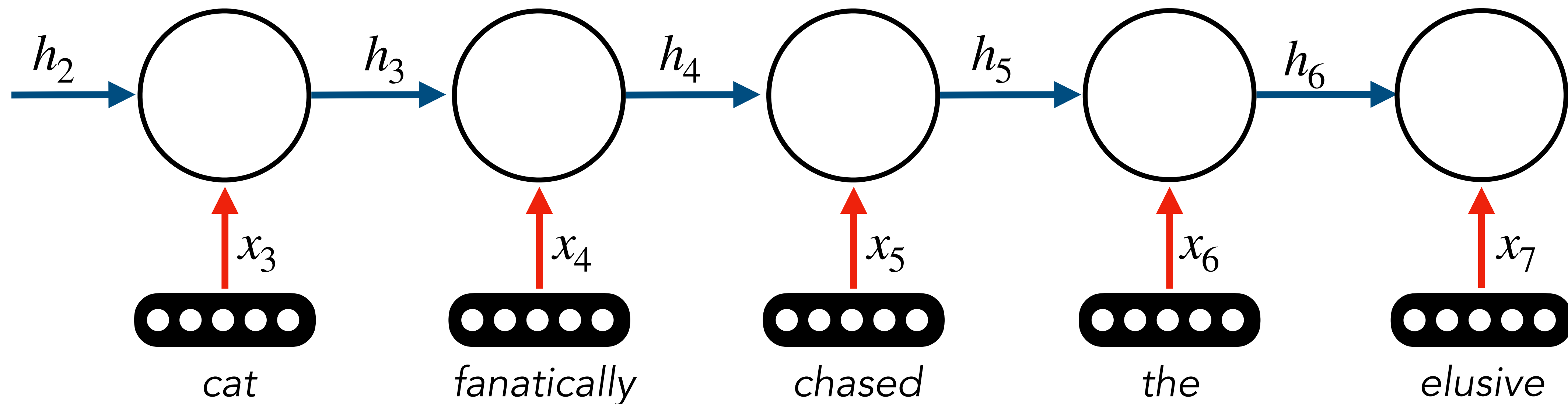
$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^o \frac{\partial \phi_{12}(v)}{\partial v} w_{21}^{\ell=1}$$

Question

**How would we extend backpropagation
to a recurrent neural network?**

Recall

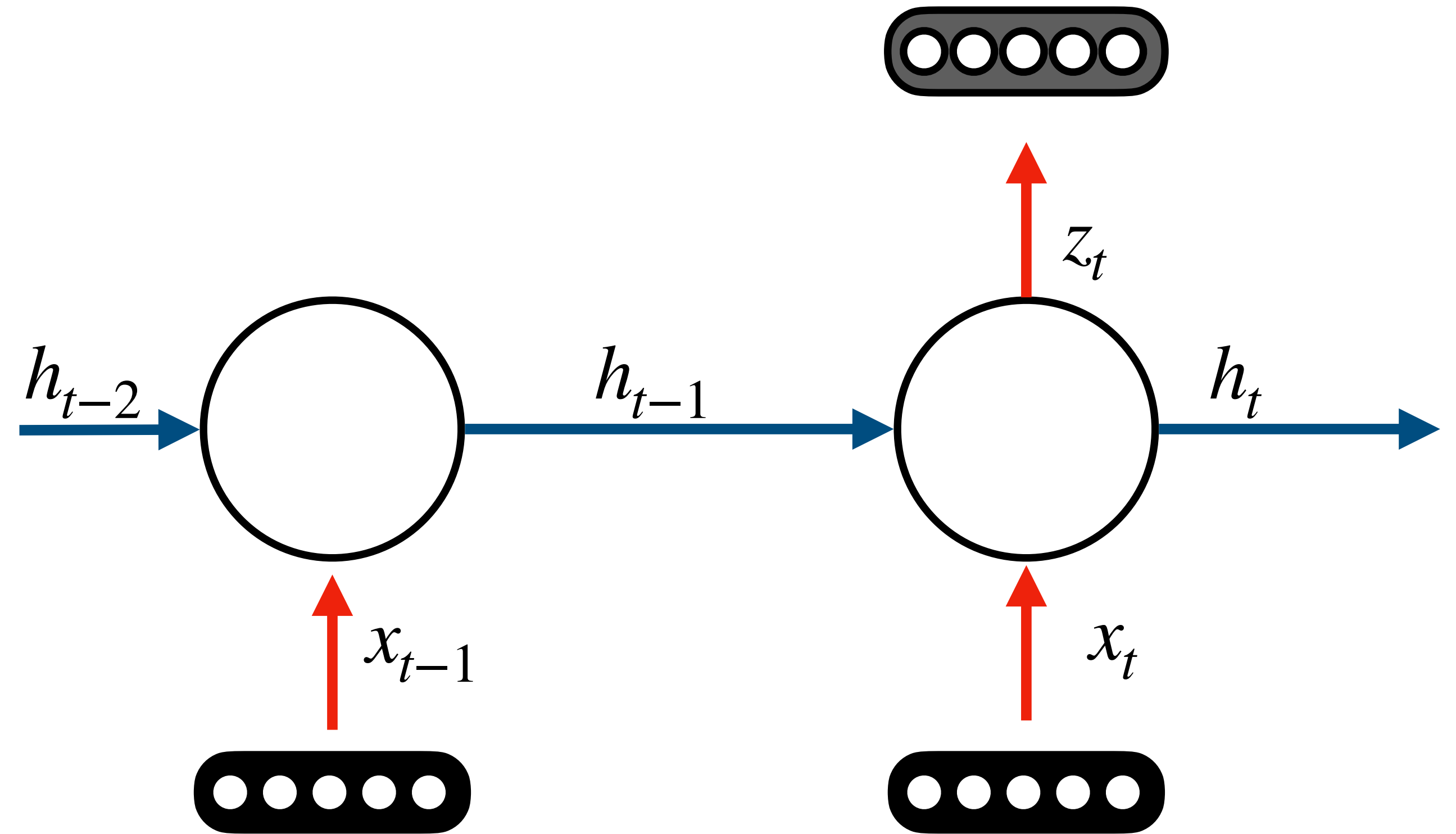
- RNN can be unrolled to a feedforward neural network
- Depth of feedforward neural network depends on length of the sequence



Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$



Backpropagation through Time

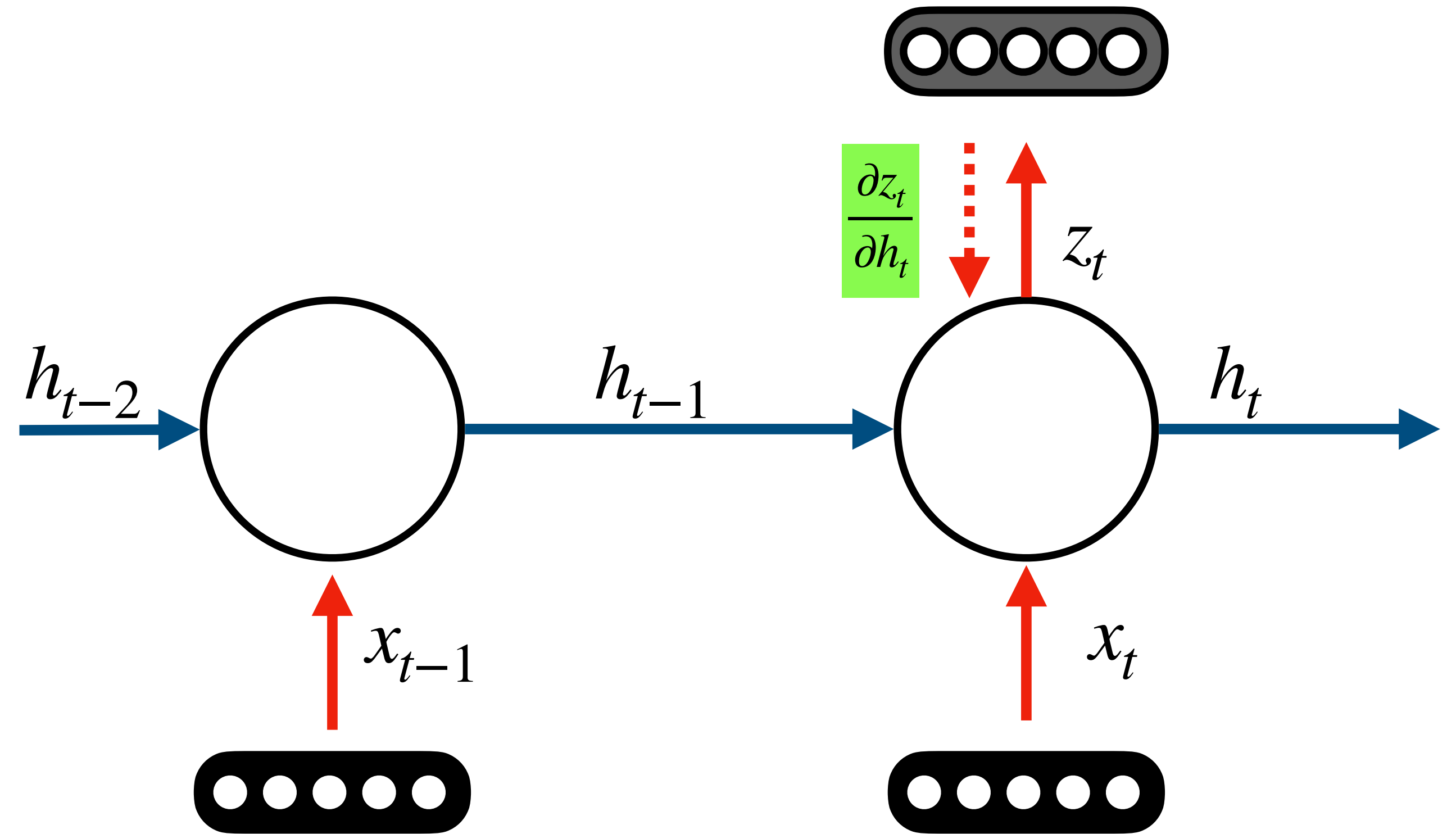
$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

$$v = W_{zh}h_t + b_z \quad z_t = \sigma(v)$$

$$u = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} \frac{\partial v}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} W_{zh}$$



Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

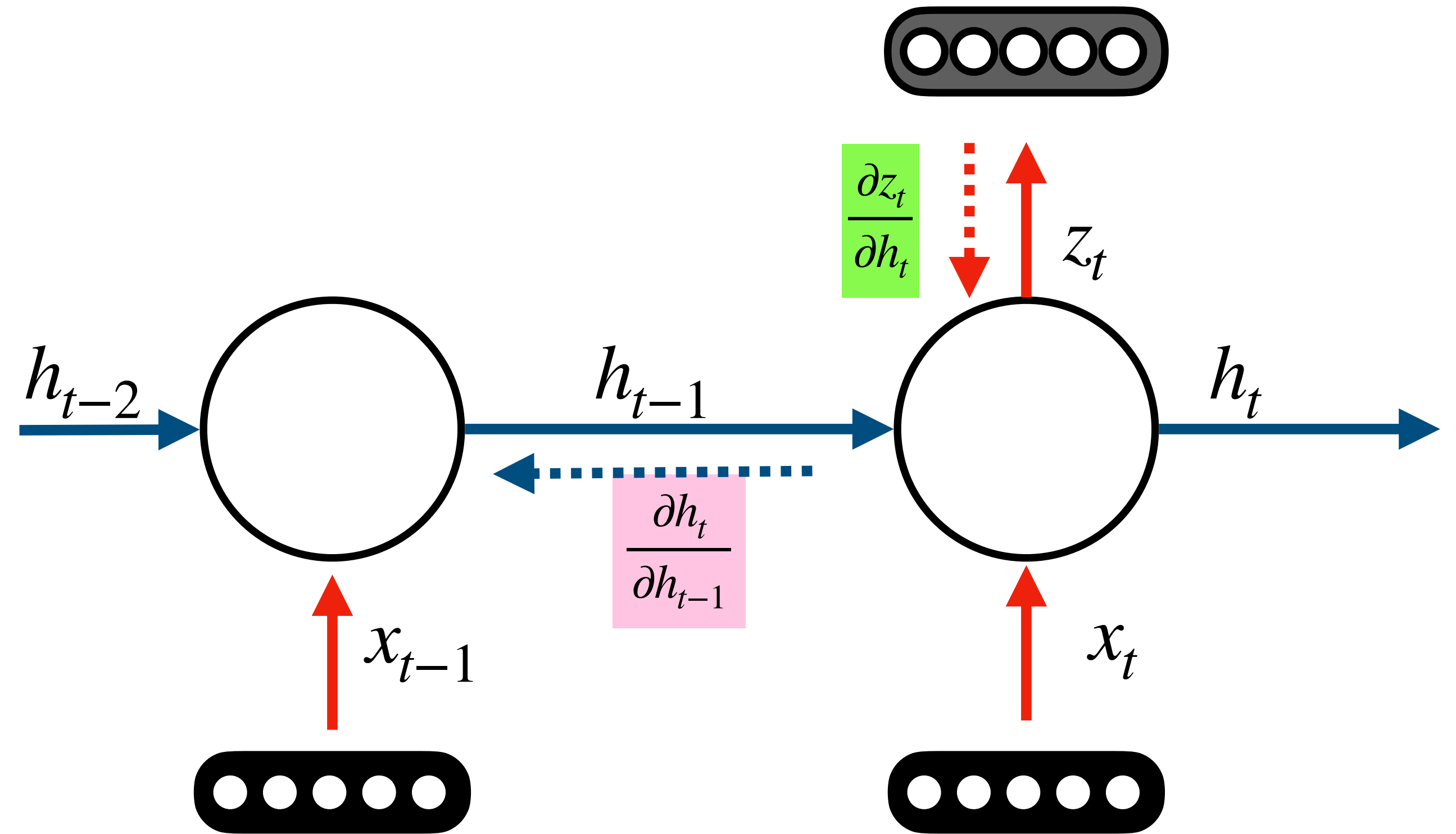
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

$$v = W_{zh}h_t + b_z \quad z_t = \sigma(v)$$

$$u = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} \frac{\partial v}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} W_{zh}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} W_{hh}$$



Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

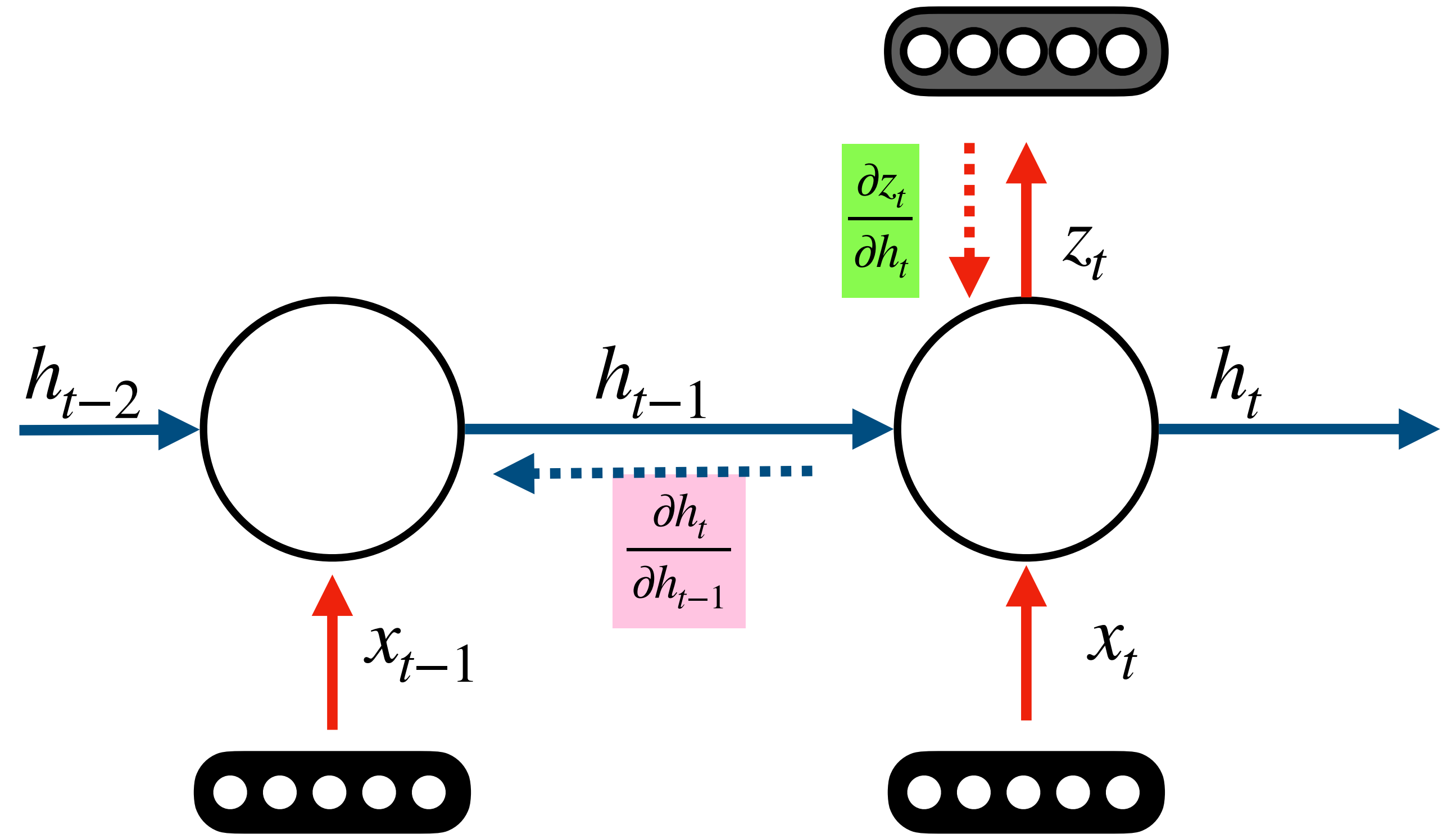
$$v = W_{zh}h_t + b_z \quad z_t = \sigma(v)$$

$$u = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} \frac{\partial v}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} W_{zh}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} W_{hh}$$

$$\frac{\partial z_t}{\partial h_{t-1}} = \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}}$$



Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

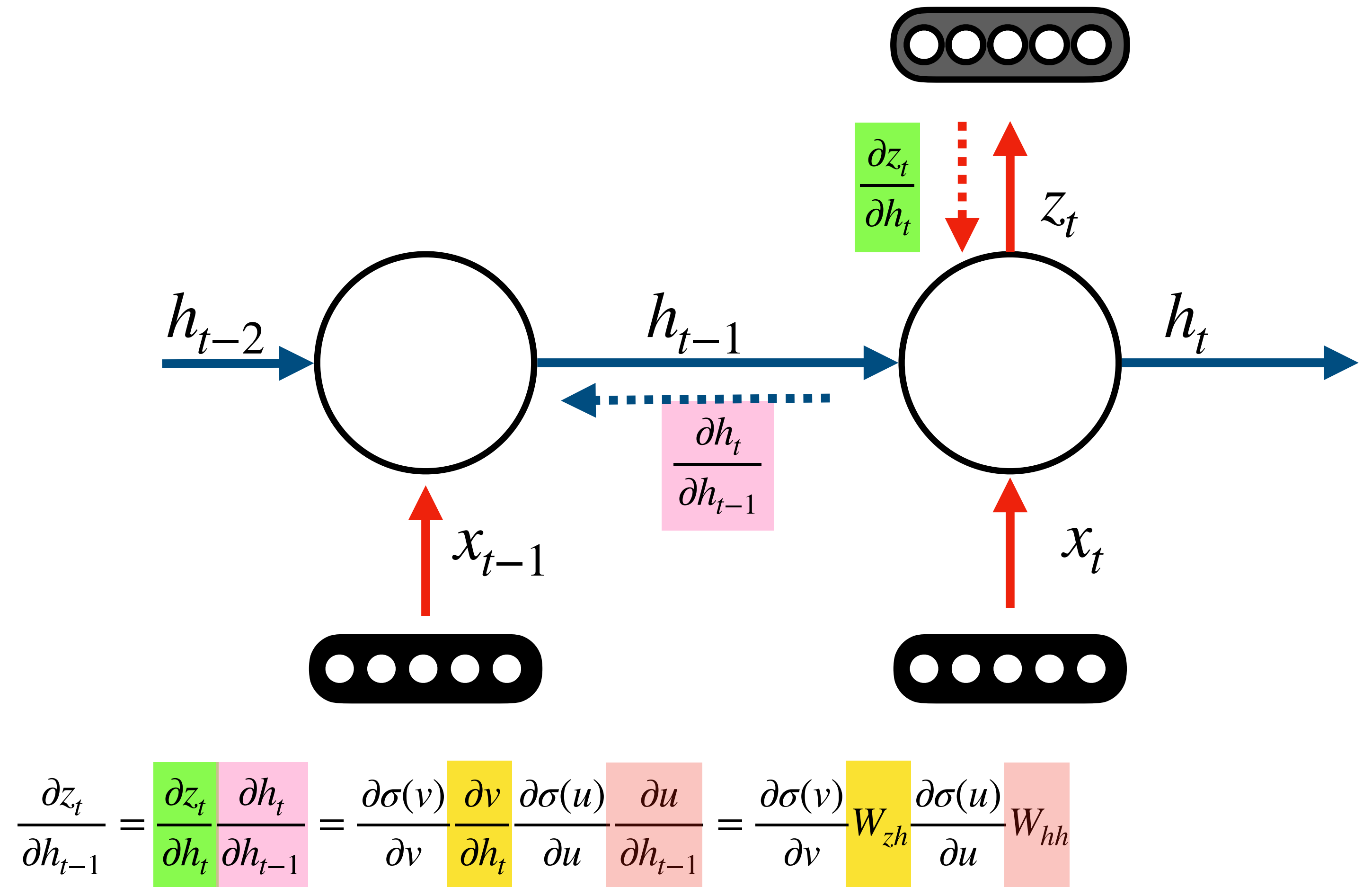
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

$$v = W_{zh}h_t + b_z \quad z_t = \sigma(v)$$

$$u = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} \frac{\partial v}{\partial h_t} = \frac{\partial \sigma(v)}{\partial v} W_{zh}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} W_{hh}$$



Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

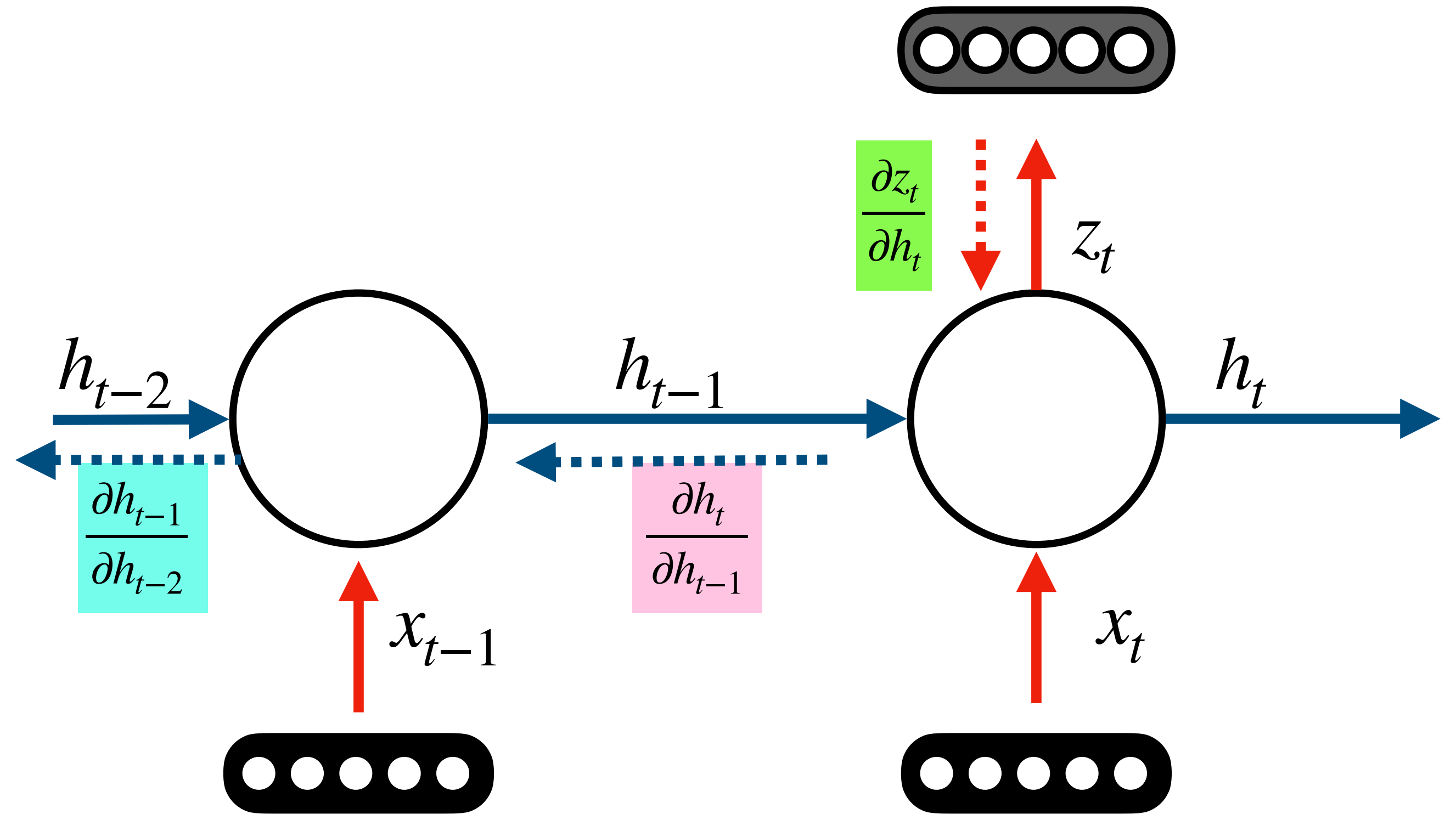
$$v_t = W_{zh}h_t + b_z \quad z_t = \sigma(v_t)$$

$$u_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u_t)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} \frac{\partial v_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} W_{zh}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \frac{\partial u_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh}$$

$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \frac{\partial u_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$



$$\frac{\partial z_t}{\partial h_{t-1}} = \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(v_t)}{\partial v_t} W_{zh} \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh} \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$

Backpropagation through Time

$$z_t = \sigma(W_{zh}h_t + b_z)$$

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

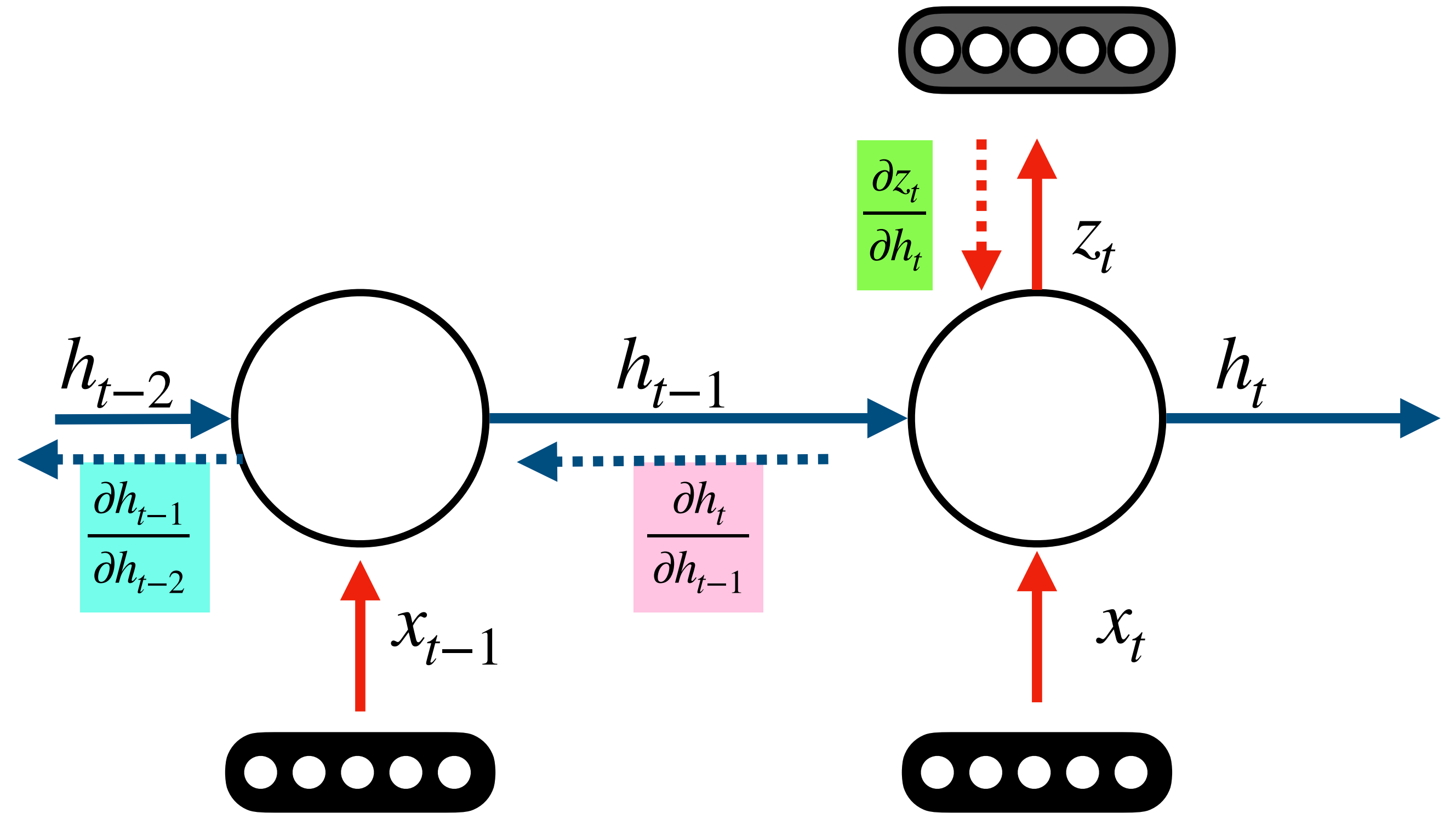
$$v_t = W_{zh}h_t + b_z \quad z_t = \sigma(v_t)$$

$$u_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h \quad h_t = \sigma(u_t)$$

$$\frac{\partial z_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} \frac{\partial v_t}{\partial h_t} = \frac{\partial \sigma(v_t)}{\partial v_t} W_{zh}$$

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \frac{\partial u_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh}$$

$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \frac{\partial u_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$

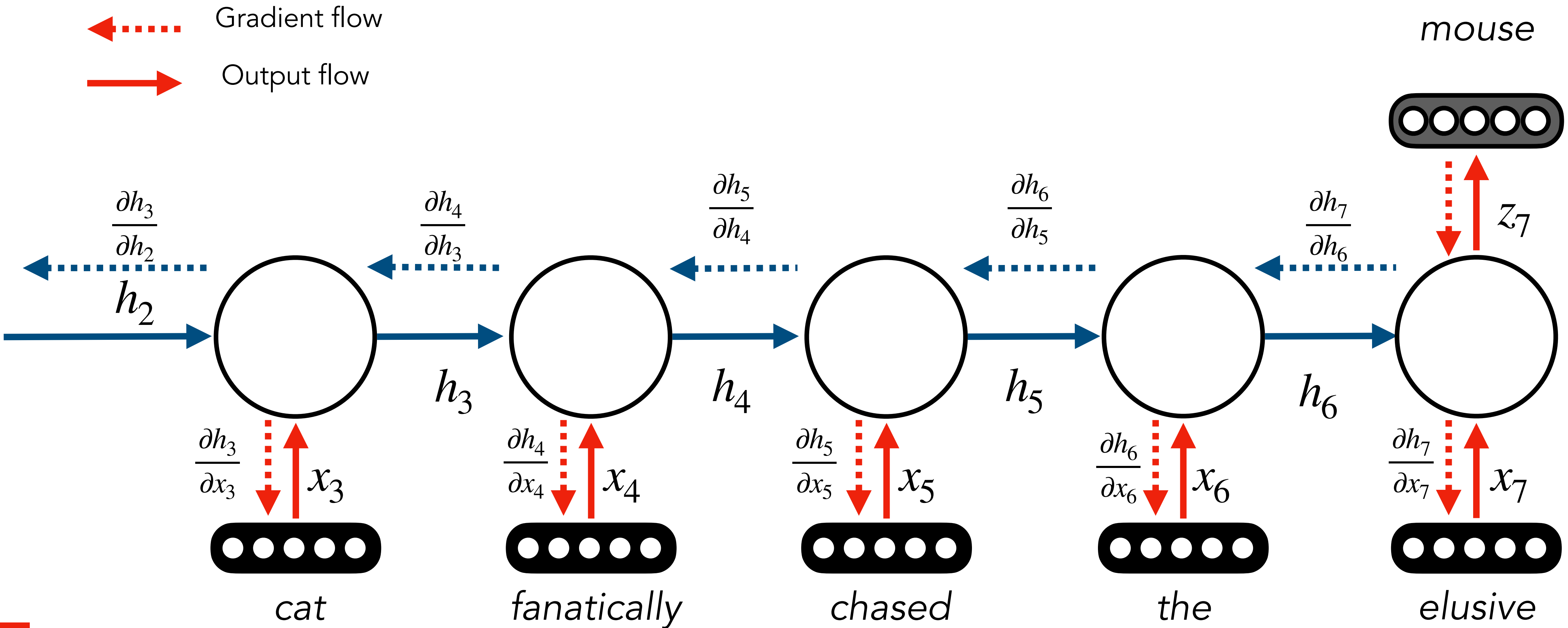


$$\frac{\partial z_t}{\partial h_{t-1}} = \frac{\partial z_t}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(v_t)}{\partial v_t} W_{zh} \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh} \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$

Note that these are actually the same matrix

Backpropagation through time

← Gradient flow
→ Output flow



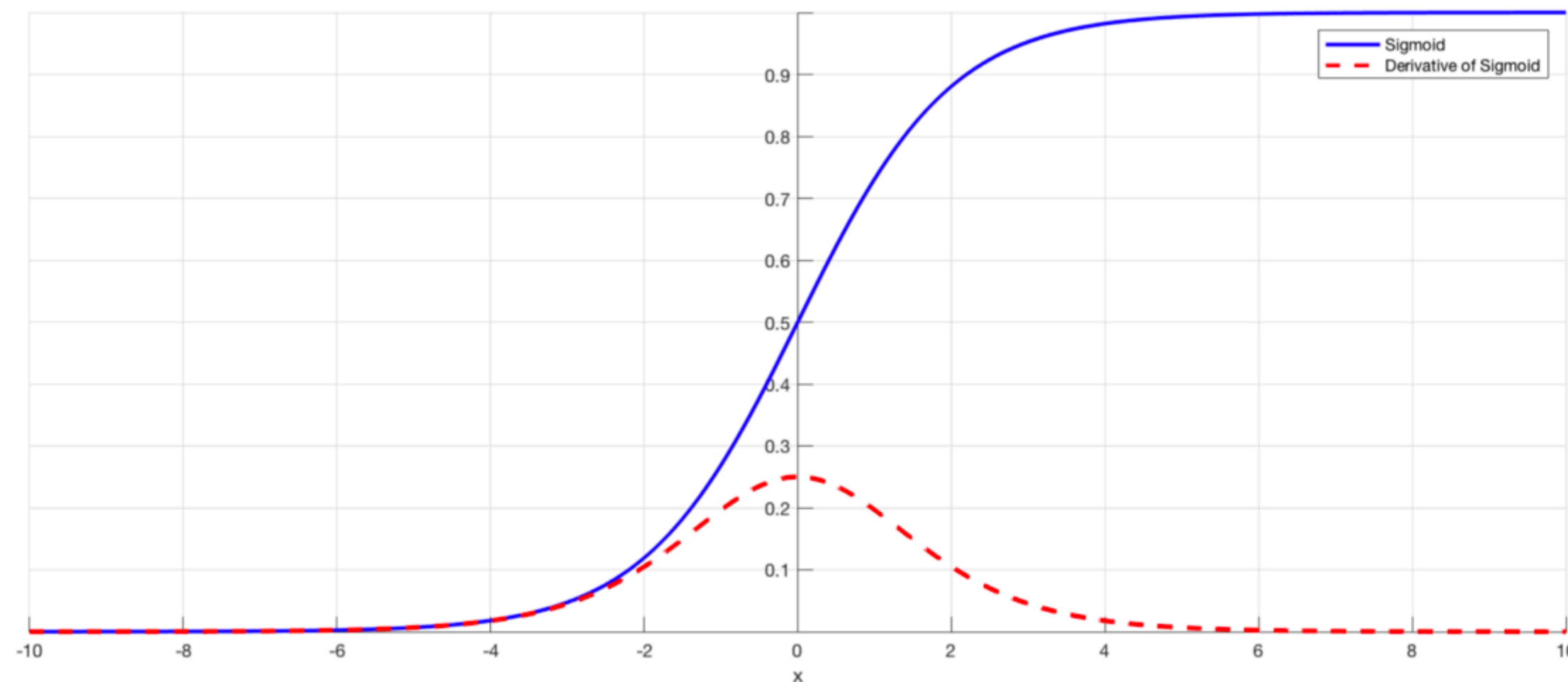
Summary

- Neural language models allow us to *share information* among similar sequences by learning neural representations that similarly represent them
- **Problem:** Fixed context language models can only process a limited window of the word history at a time
- **Solution:** recurrent neural networks can **theoretically** learn to model an **unbounded context length**

Vanishing Gradients

- **Learning Problem:** Long unrolled networks will crush gradients that backpropagate to earlier time steps

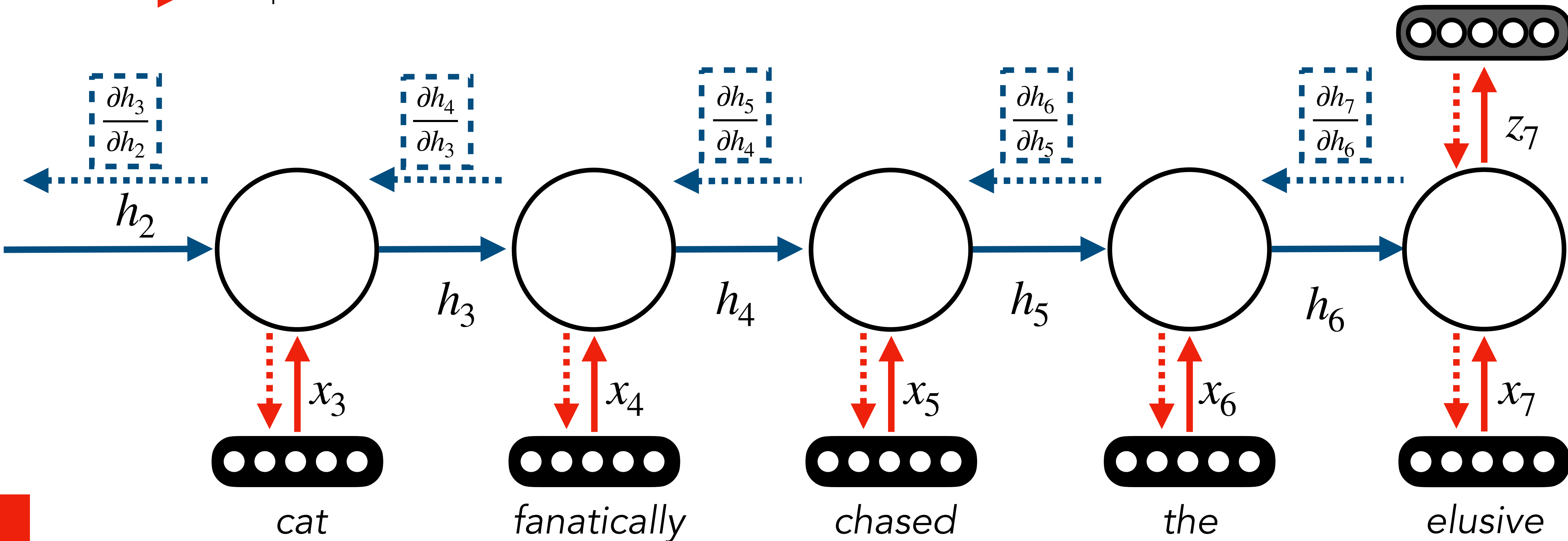
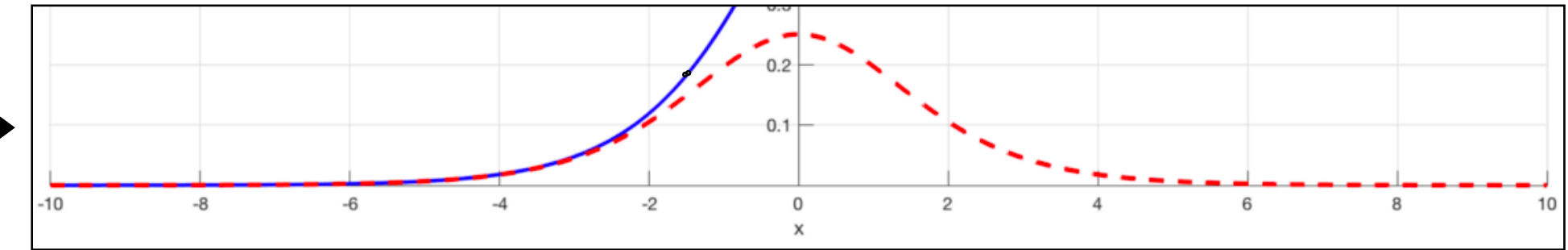
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$
$$u = W_{hx}x_t + W_{hh}h_{t-1} + b_h$$
$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = W_{hh} \frac{\partial \sigma(u)}{\partial u}$$



Backpropagation through time

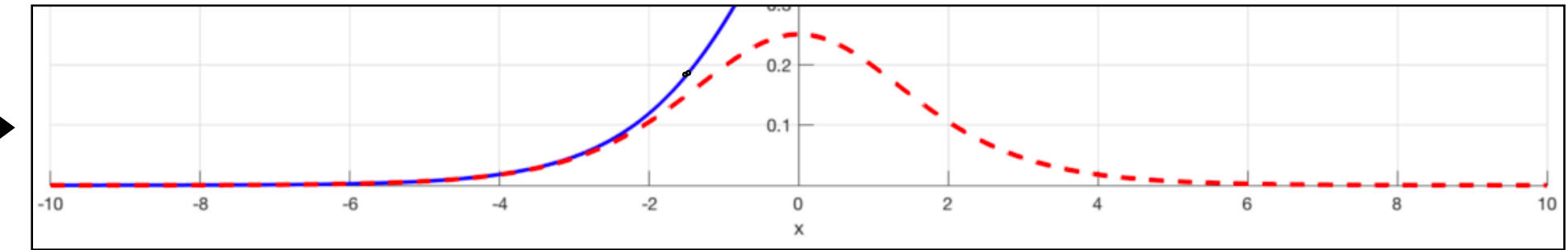
← Gradient flow
→ Output flow

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = W_{hh} \frac{\partial \sigma(u)}{\partial u}$$



Backpropagation through time

$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u)}{\partial u} \frac{\partial u}{\partial h_{t-1}} = W_{hh} \frac{\partial \sigma(u)}{\partial u}$$



Gradient flow



Output flow

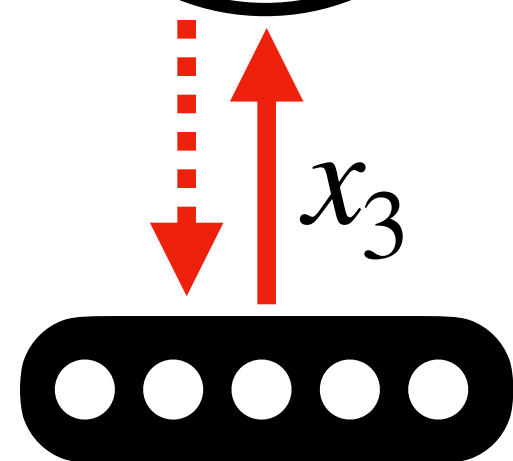


Problem in many recurrent neural networks,

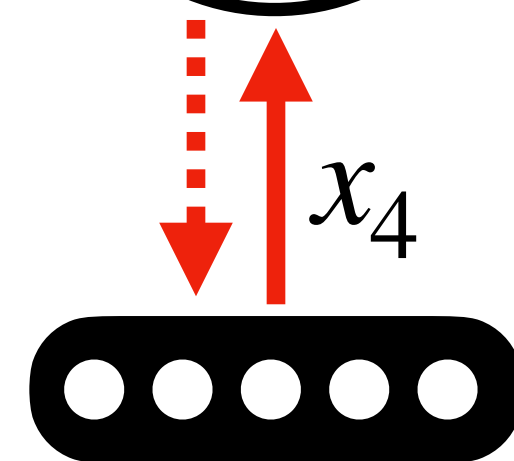
Especially pronounced in Elman networks (Vanilla RNNs) due to the sigmoid activation

$$\frac{\partial h_3}{\partial h_2}$$

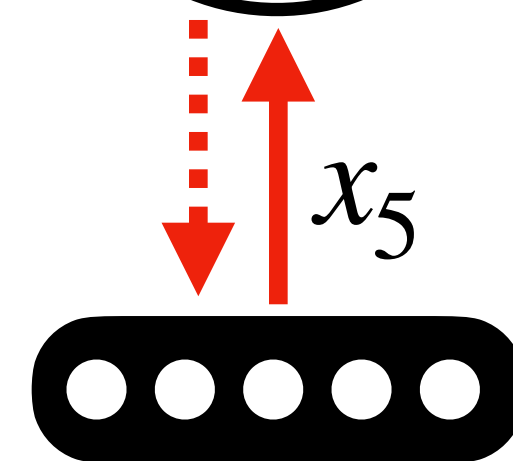
h_2



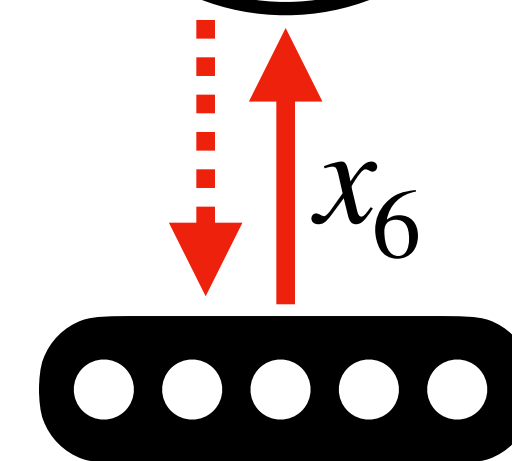
cat



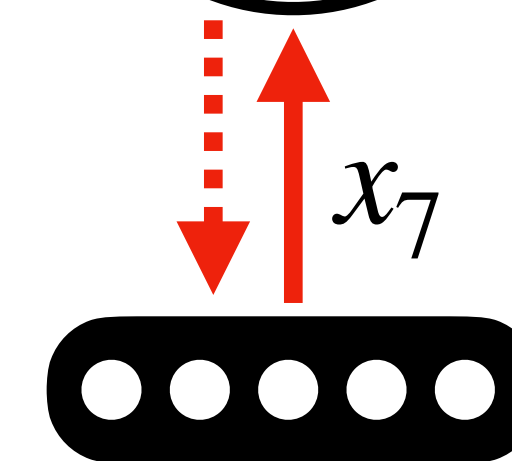
fanatically



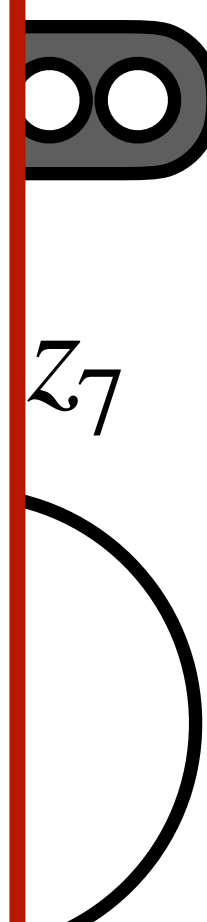
chased



the



elusive



Issue with Recurrent Models

- Multiple steps of state overwriting makes it challenging to learn long-range dependencies.

*They tuned, discussed for a moment, then struck up a lively **jig**. Everyone joined in, turning the courtyard into an even more chaotic scene, people now **dancing** in circles, **swinging** and **spinning** in circles, everyone making up their own **dance steps**. I felt my feet tapping, my body wanting to move. Aside from writing, I 've always loved **dancing** .*

- Nearby words should affect each other more than farther ones, but RNNs make it challenging to learn **any** long-range interactions

Gated Recurrent Neural Networks

- Use gates to avoid dampening gradient signal every time step

$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

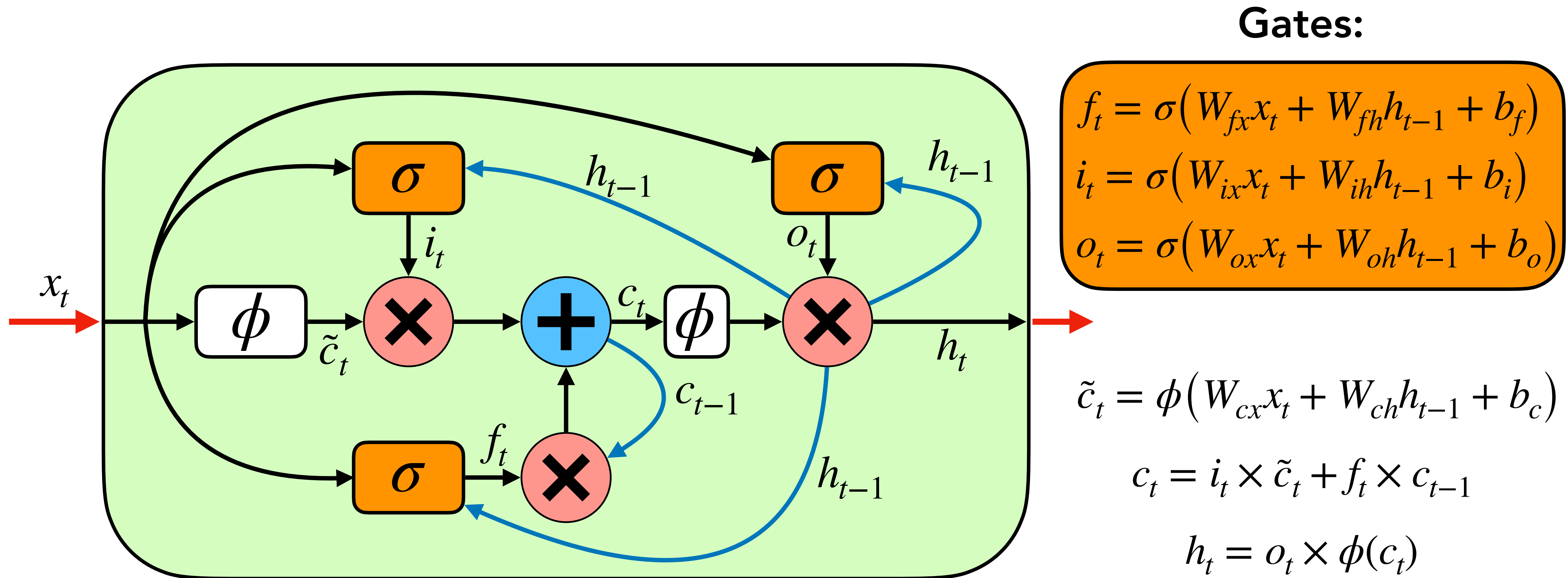
Elman Network

$$h_t = h_{t-1} \odot \mathbf{f} + \mathbf{func}(x_t)$$

Gated Network Abstraction

- Gate value \mathbf{f} computes how much information from previous hidden state moves to the next time step $\rightarrow 0 < \mathbf{f} < 1$
- Because h_{t-1} is no longer inside the activation function, it is not automatically constrained, reducing vanishing gradients!

Long Short Term Memory (LSTM)



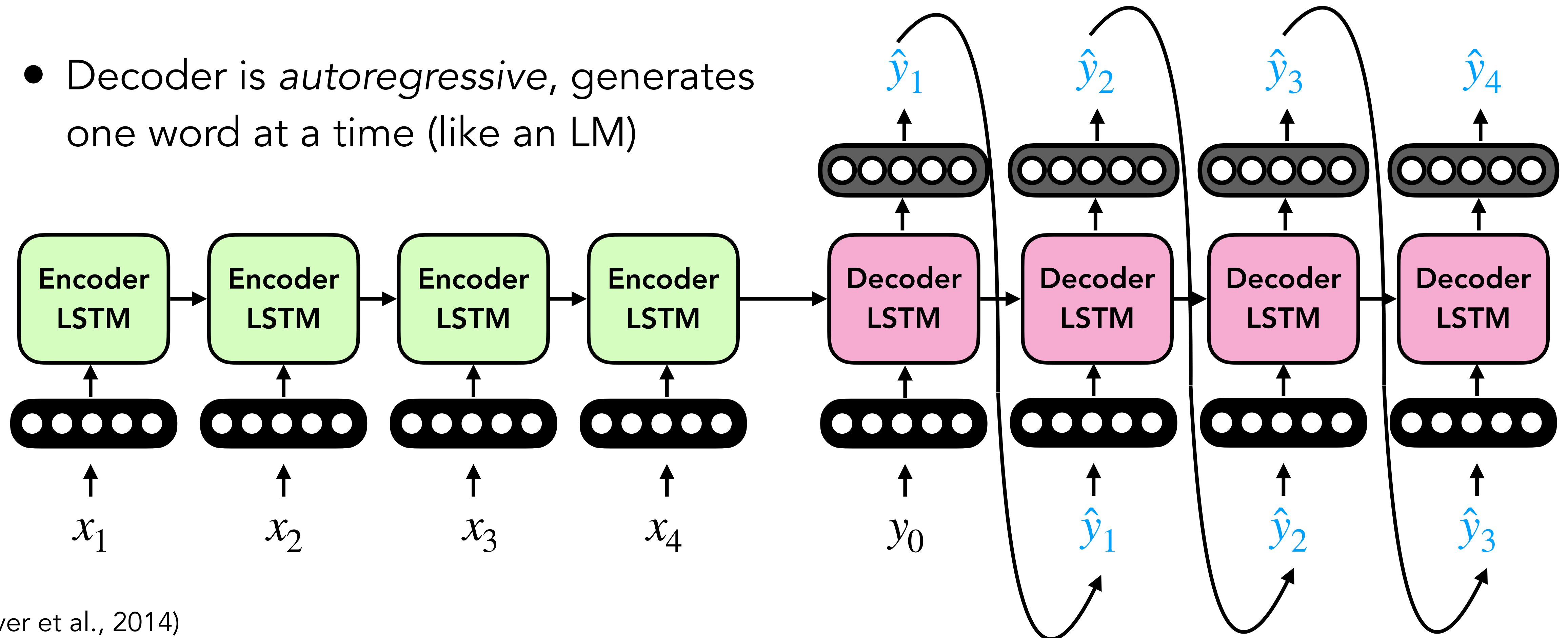
Question

How can we use recurrent neural networks in practice?

Machine Translation involves more than estimating the probability next word; requires generating a full translation of a given context into another language

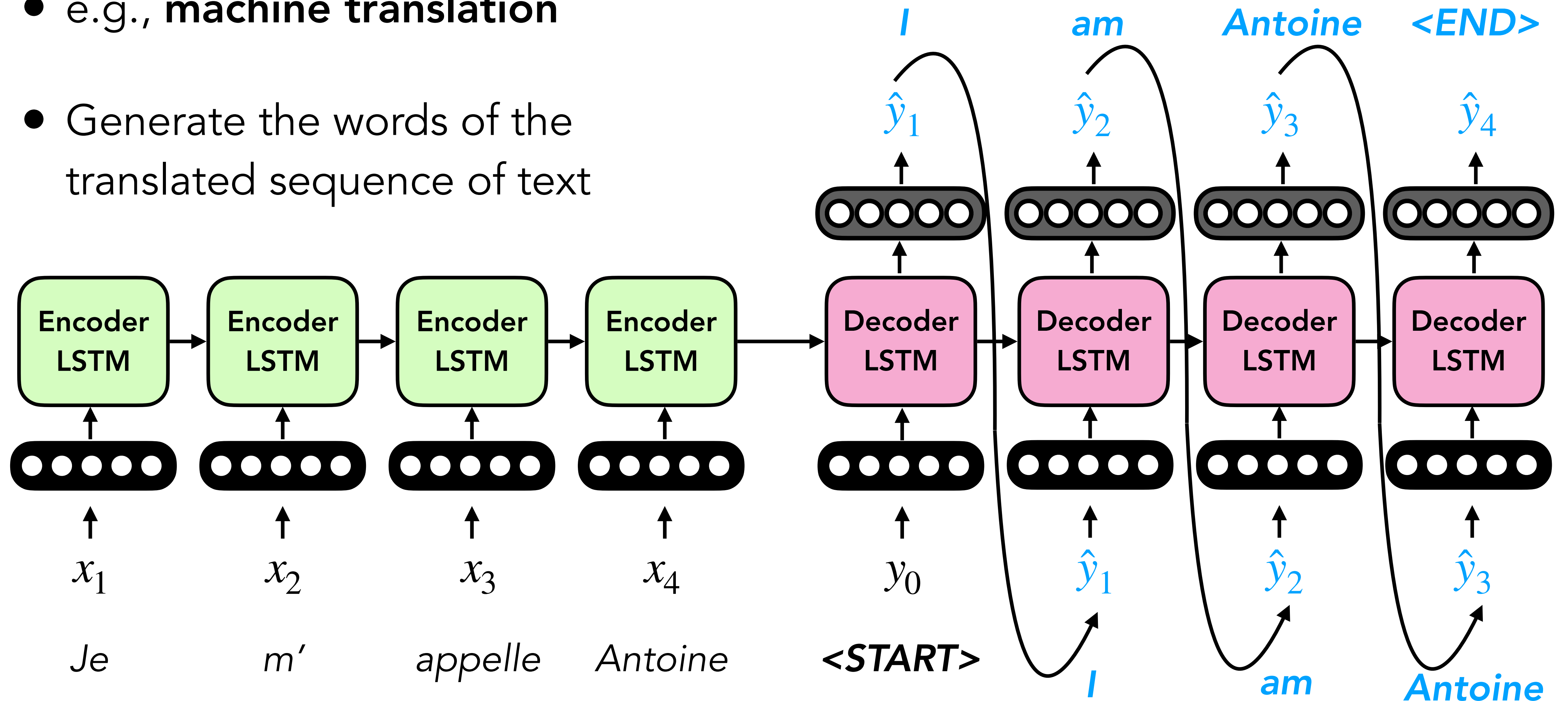
Encoder-Decoder Models

- 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)



Encoder-Decoder Models

- e.g., machine translation
- Generate the words of the translated sequence of text



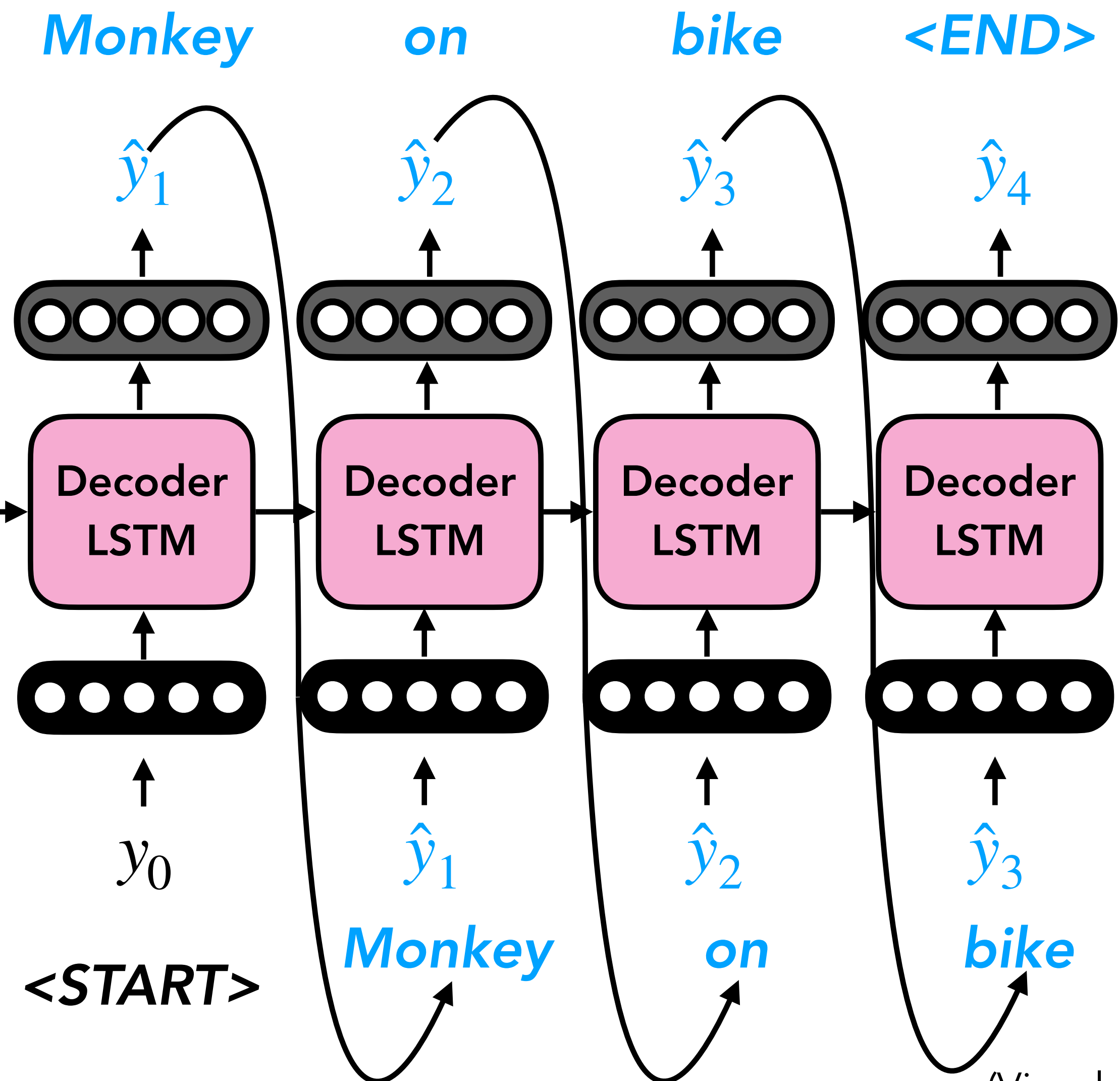
Encoder-Decoder Models

- Input doesn't need to be text
- e.g., image captioning



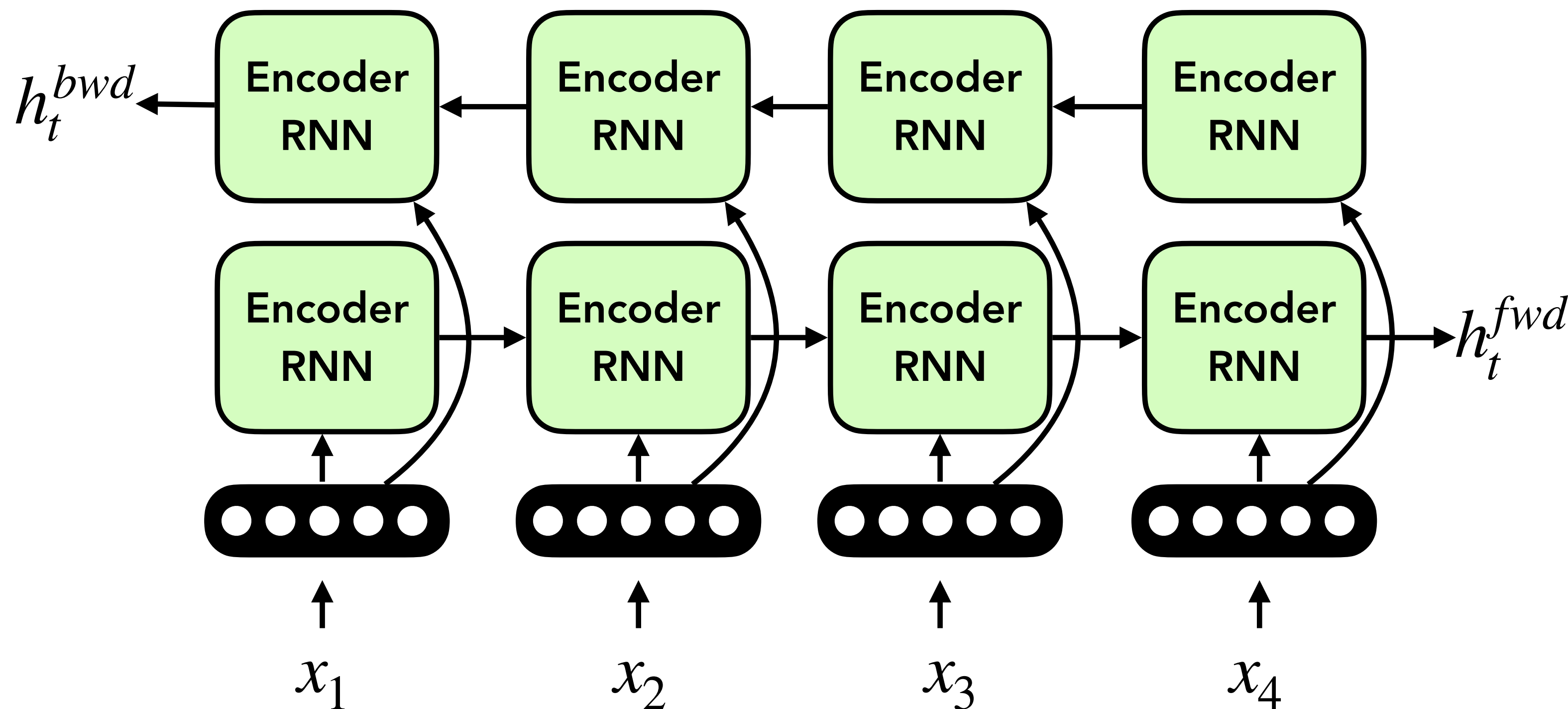
Photo credit: J Hovenstine Studios

- Generate words of image description



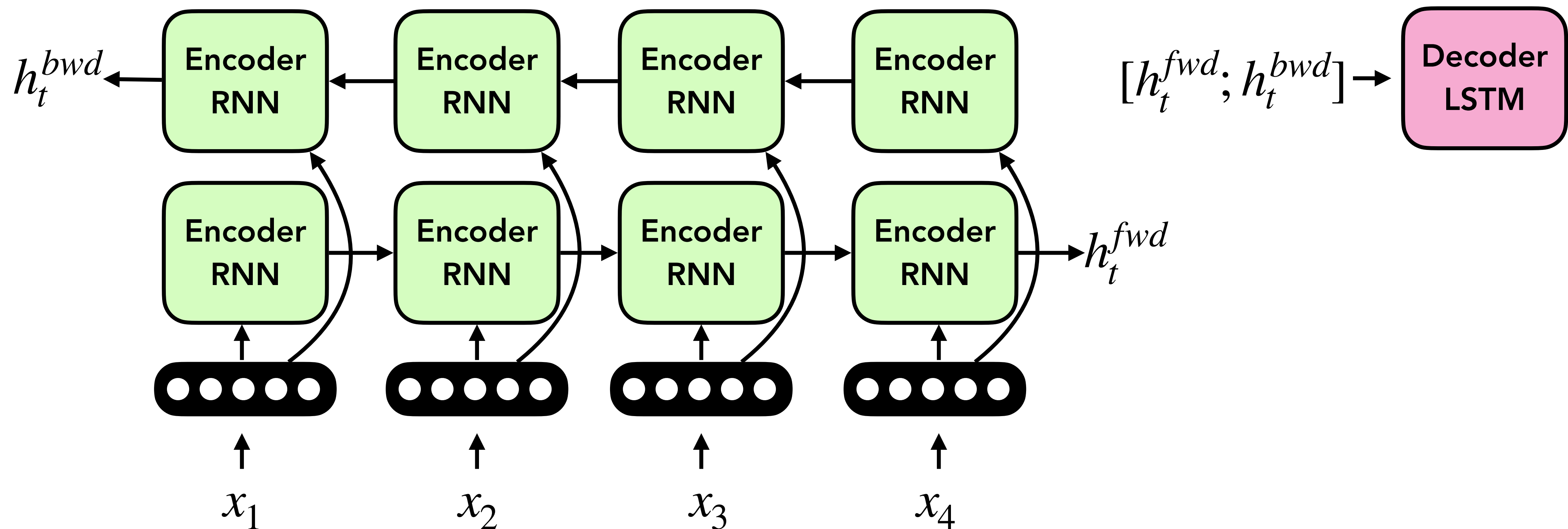
Bidirectional Encoders

- Decoder needs to be unidirectional (can't know the future...)
- Encoder sequence representation augmented by encoding in both directions



Bidirectional Encoders

- Decoder needs to be unidirectional (can't know the future...)
- Encoder sequence representation augmented by encoding in both directions



Other Resources of Interest

- Approaches for maintaining state and avoiding vanishing gradients
 - Long Short-Term Memory (Hochreiter and Schmidhuber, 1997):
 - Gated Recurrent Units (Cho et al., 2014):
- LSTM: A Search Space Odyssey (Greff et al., 2015)
 - Examine 5000 different modifications to LSTMs — none significantly better than original architecture
- Only basics presented here today! Many offshoots of these techniques!

Recap

- Early neural language models (and n-gram models) suffer from **fixed context windows**
- Recurrent neural networks can **theoretically** learn to model an **unbounded context length** using back propagation through time (BPTT)
- Practically, however, **vanishing gradients** stop many RNN architectures from learning **long-range dependencies**
- RNNs (and modern variants) remain useful for **many sequence-to-sequence tasks**

References

- Bengio, Y., Ducharme, R., Vincent, P., & Janvin, C. (2003). A Neural Probabilistic Language Model. *Journal of machine learning research*.
- Elman, J.L. (1990). Finding Structure in Time. *Cogn. Sci.*, 14, 179-211.
- Schuster, M., & Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal Processing*, 45(11), 2673–2681.
- Hochreiter, S., & Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation*, 9, 1735-1780.
- Cho, K., Merriënboer, B.V., Gülçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., & Bengio, Y. (2014). Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation. *Conference on Empirical Methods in Natural Language Processing*.
- Sutskever, I., Vinyals, O., & Le, Q.V. (2014). Sequence to Sequence Learning with Neural Networks. *NIPS*.
- Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2014). Show and tell: A neural image caption generator. *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3156-3164.
- Greff, K., Srivastava, R.K., Koutník, J., Steunebrink, B.R., & Schmidhuber, J. (2015). LSTM: A Search Space Odyssey. *IEEE Transactions on Neural Networks and Learning Systems*, 28, 2222-2232.