Deep Learning for Natural Language Processing

Antoine Bosselut





Natural Language Processing

Enabling
Human-Machine
Collaboration

Accelerating
Human-Human
Communication

Mining
Human Insights

Search Engines

Machine Translation

Sentiment Analysis

Dialogue Agents

Text Summarization

Motivation Analysis

Text Generation

Information Extraction

Emotion Detection











Machine Translation



Question Answering



ChatBots!



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

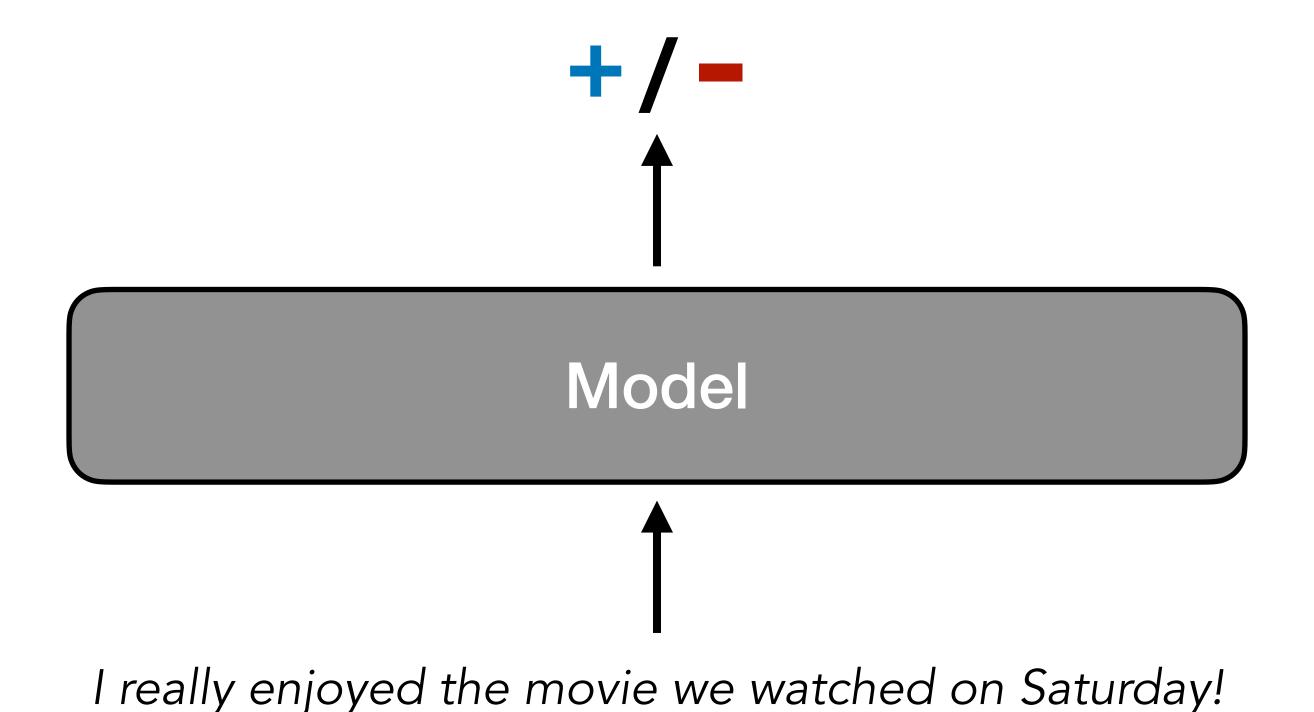
- Introduction
- Section 1 Neural NLP & Word Embeddings
- Section 2 Recurrent Neural Networks for Sequence Modeling
- Section 3 Attentive Neural Modeling with Transformers
- Exercise Session: Attention in Transformer Language Models

Part 1: Neural Embeddings

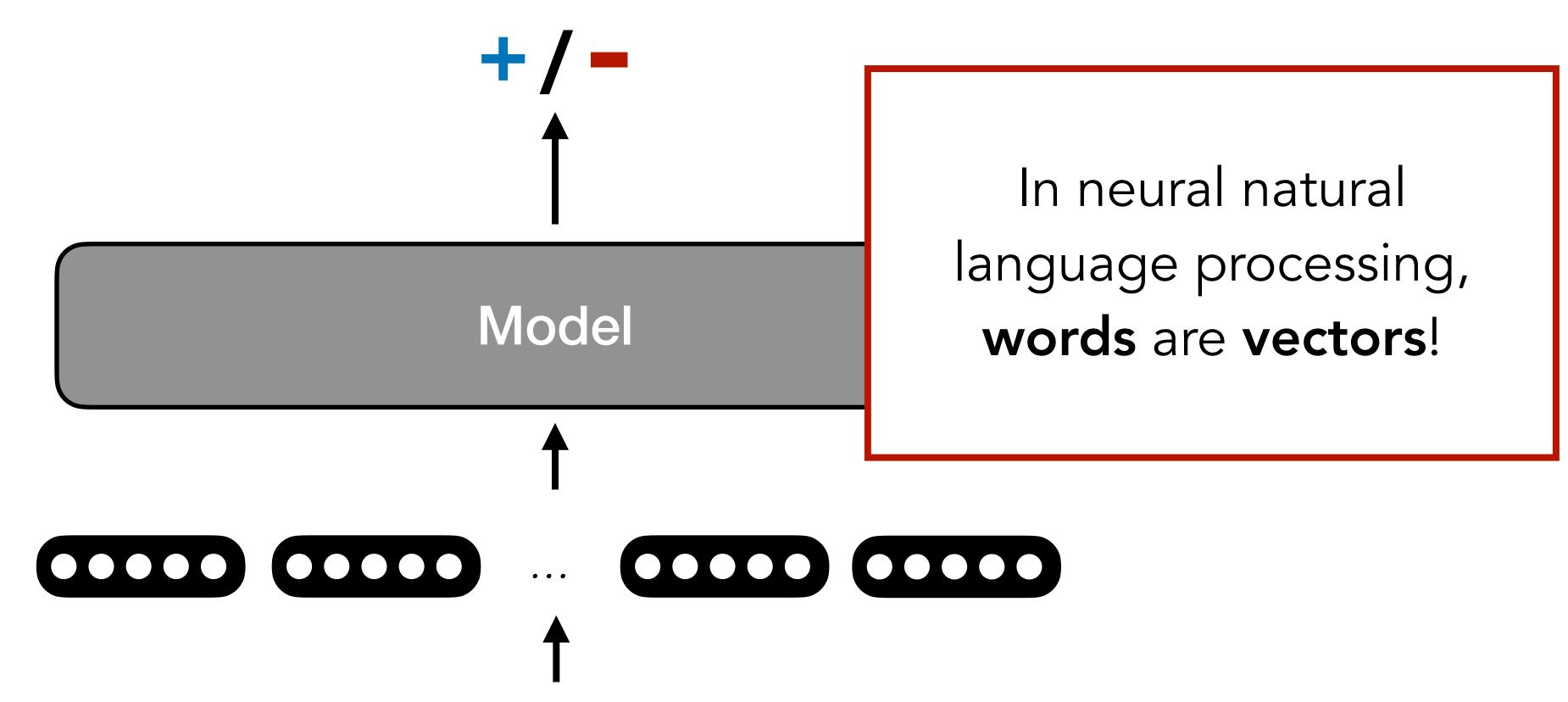
Section Outline

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

How do we represent natural language sequences for NLP problems?



How do we represent natural language sequences for NLP problems?



I really enjoyed the movie we watched on Saturday!

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 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 <UNK> token (e.g., unknown).

Question

How should we model a word as a vector?

Sparse Word Representations

- 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)

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

Word Vector Composition

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

```
I really enjoyed the movie! \longrightarrow [1...1101...01] 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!

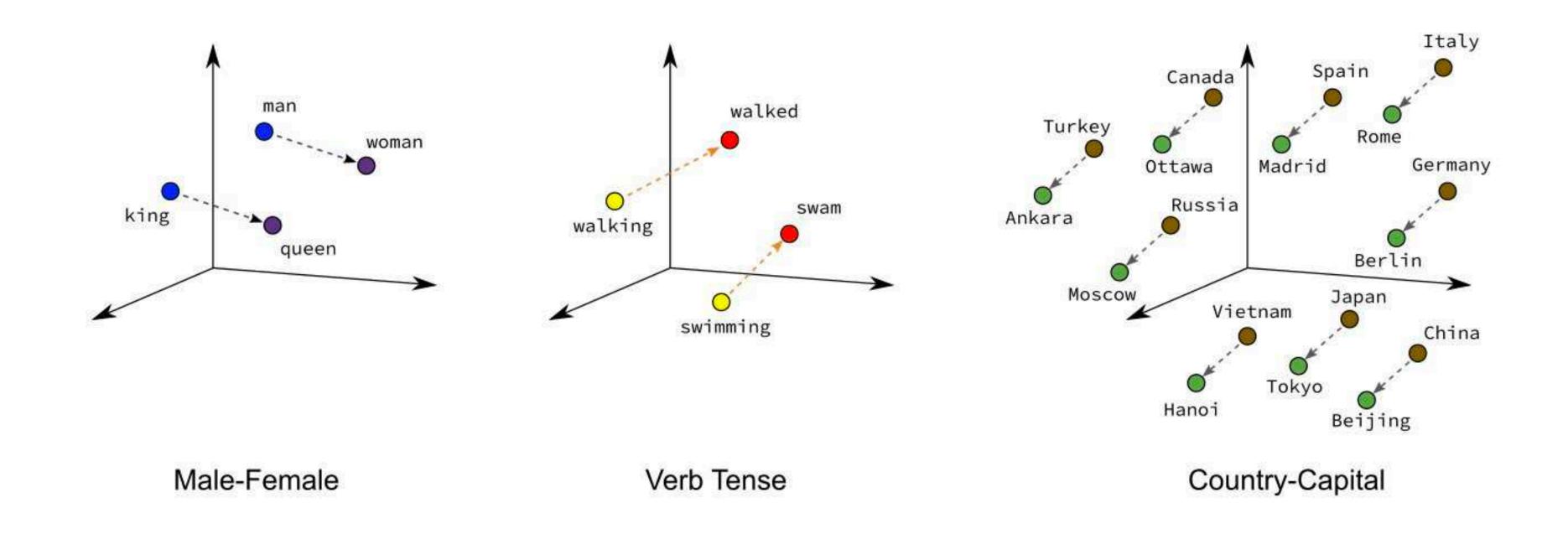
How do you learn learn similarity between words?

```
enjoyed → [0...0001...00]

loved → [0...1...0000]
```

sim(enjoyed, loved) = 0

Embeddings Goal



How do we train semantics-encoding embeddings of words?

Dense Word Vectors

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

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

word vectors

word embeddings

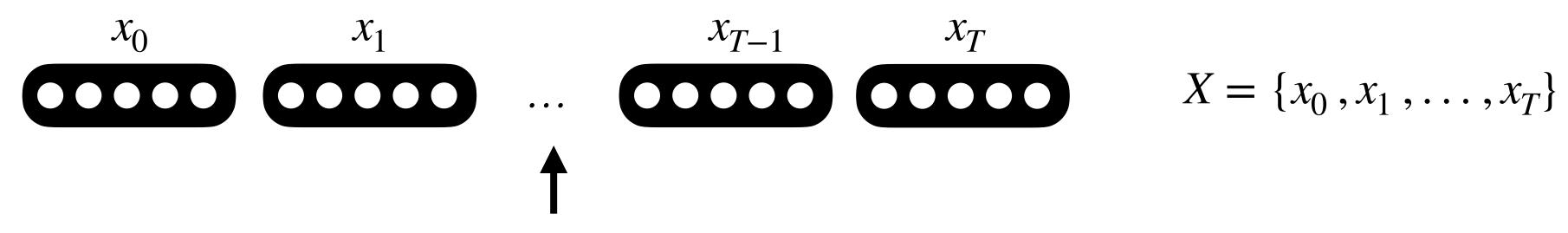
neural embeddings

dense embeddings

others...

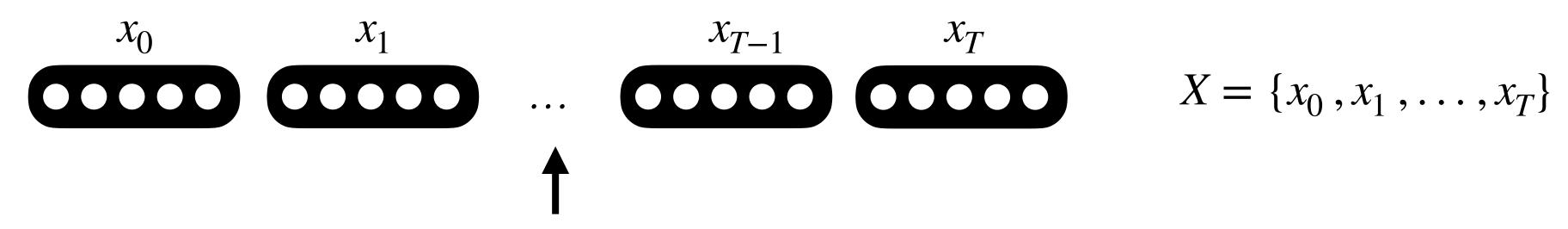
• Similarity of vectors represents similarity of meaning for particular words

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



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

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



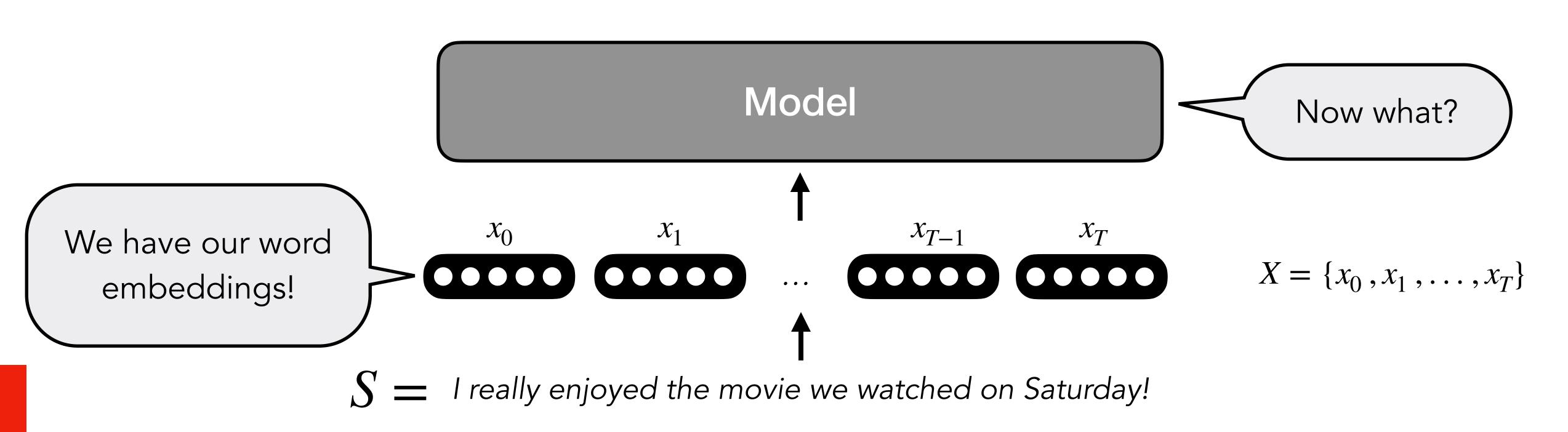
 $S_1 = \text{I really}$ enjoyed the movie we watched on Saturday!

• 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 **E**

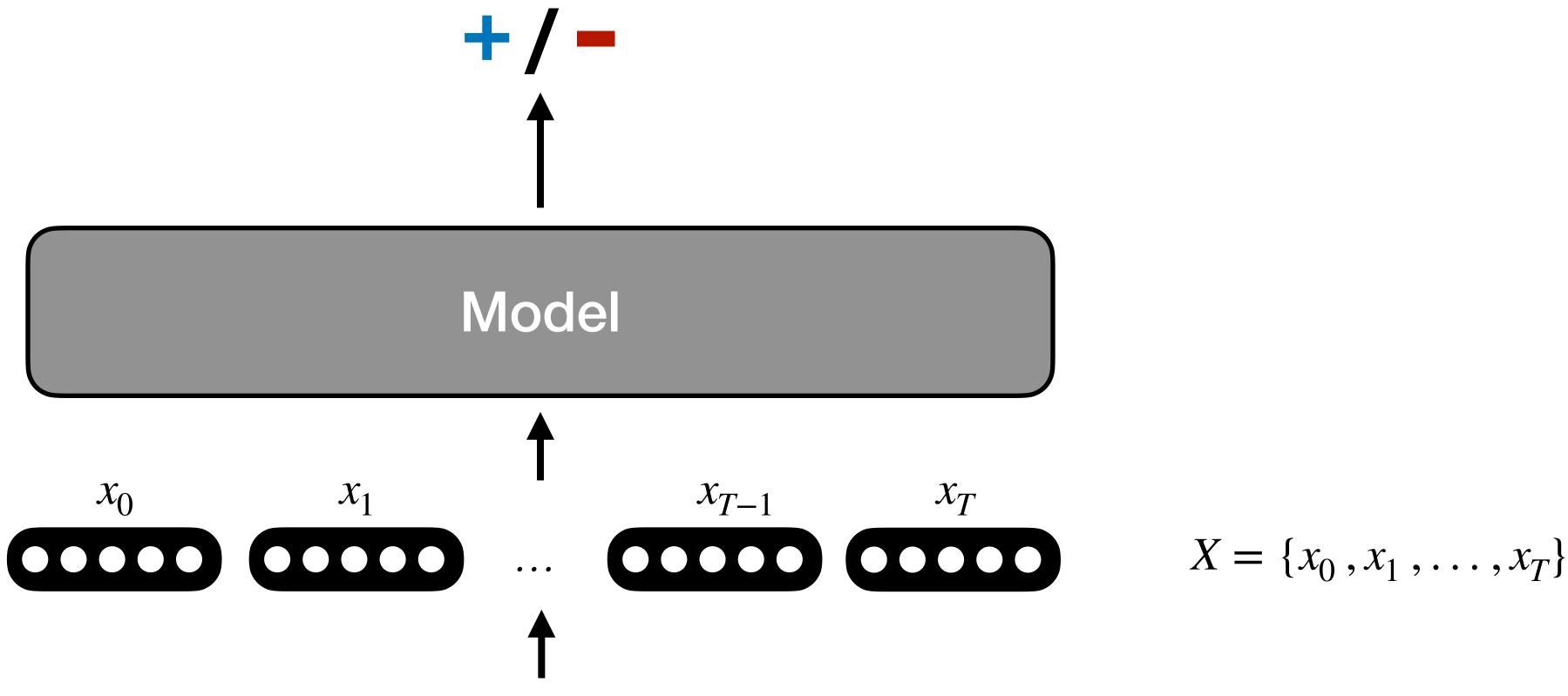
ullet For each sequence S, we have a corresponding sequence of embeddings X



Question

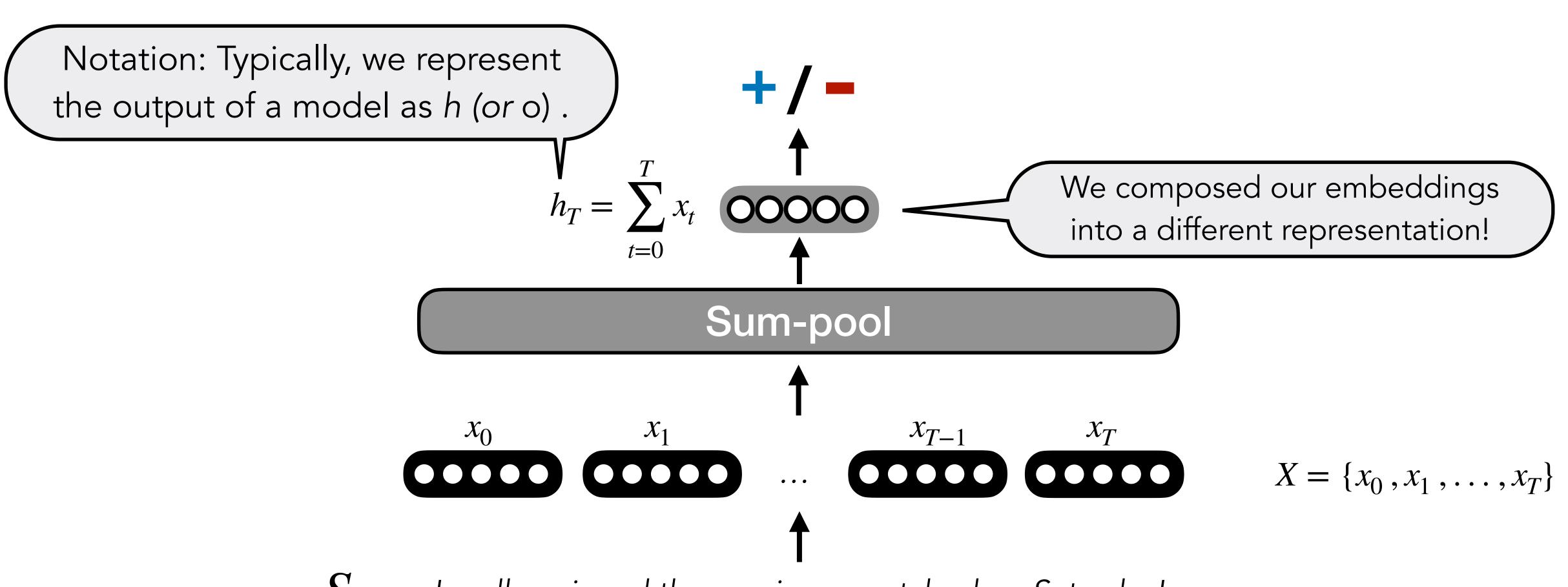
What should we use as a model?

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



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

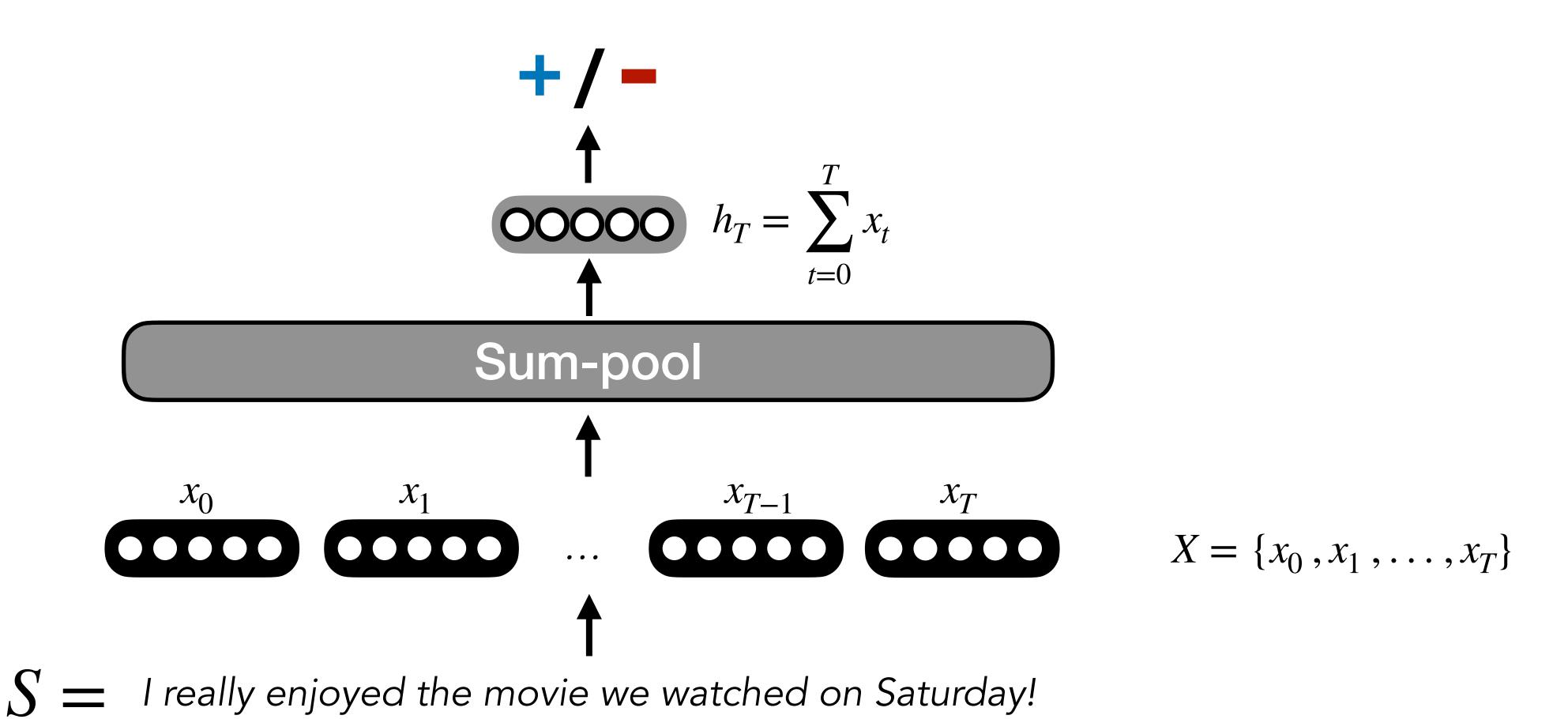
- Our model modifies and / or composes these word embeddings to formulate a representation that allows it to predict the correct label
 - Recurrent neural networks (RNNs) Today!
 - RNN variants (LSTM, GRU, etc.) Today!
 - Transformer Today!



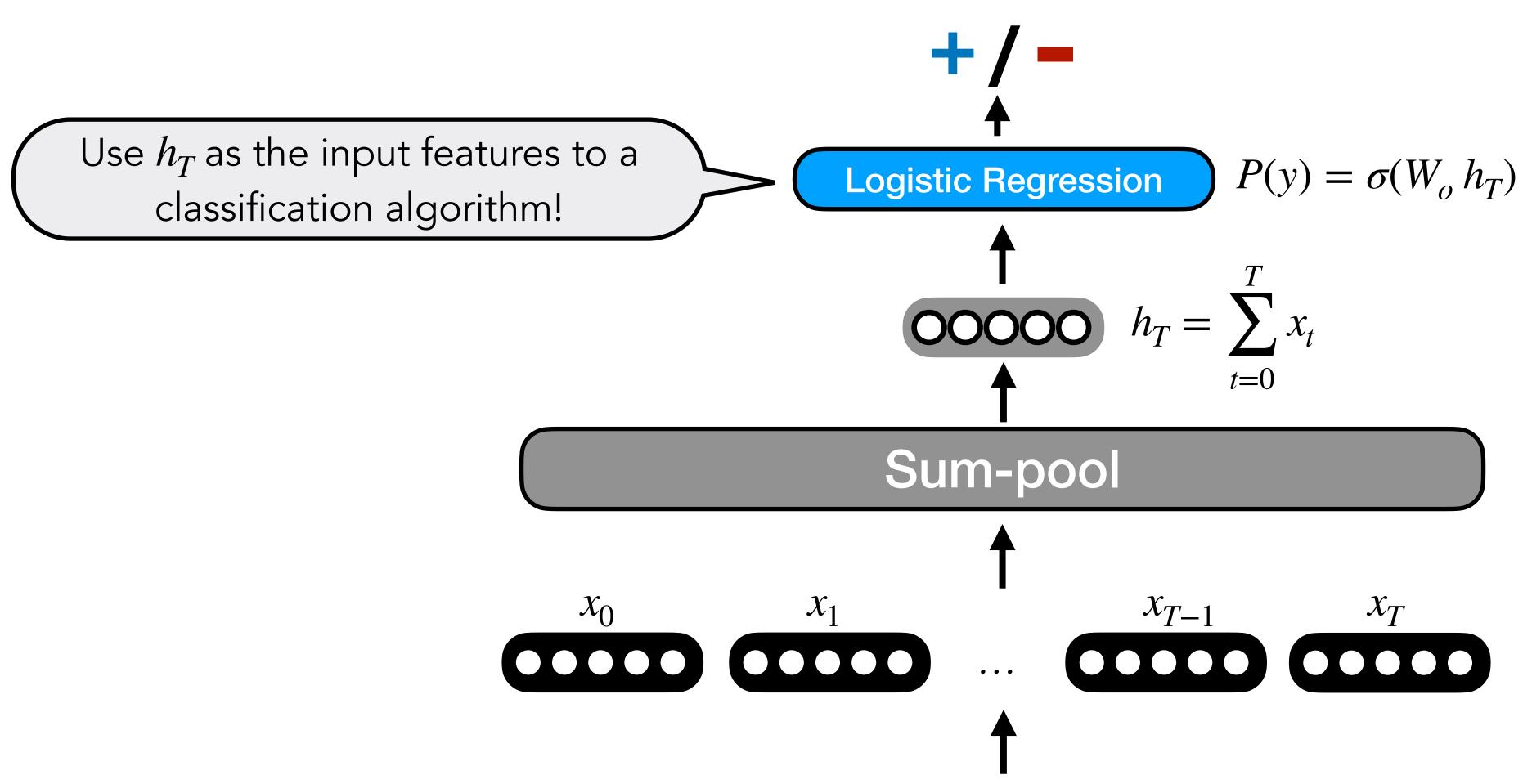
Question

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

Predicting the label



Predicting the label



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

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.

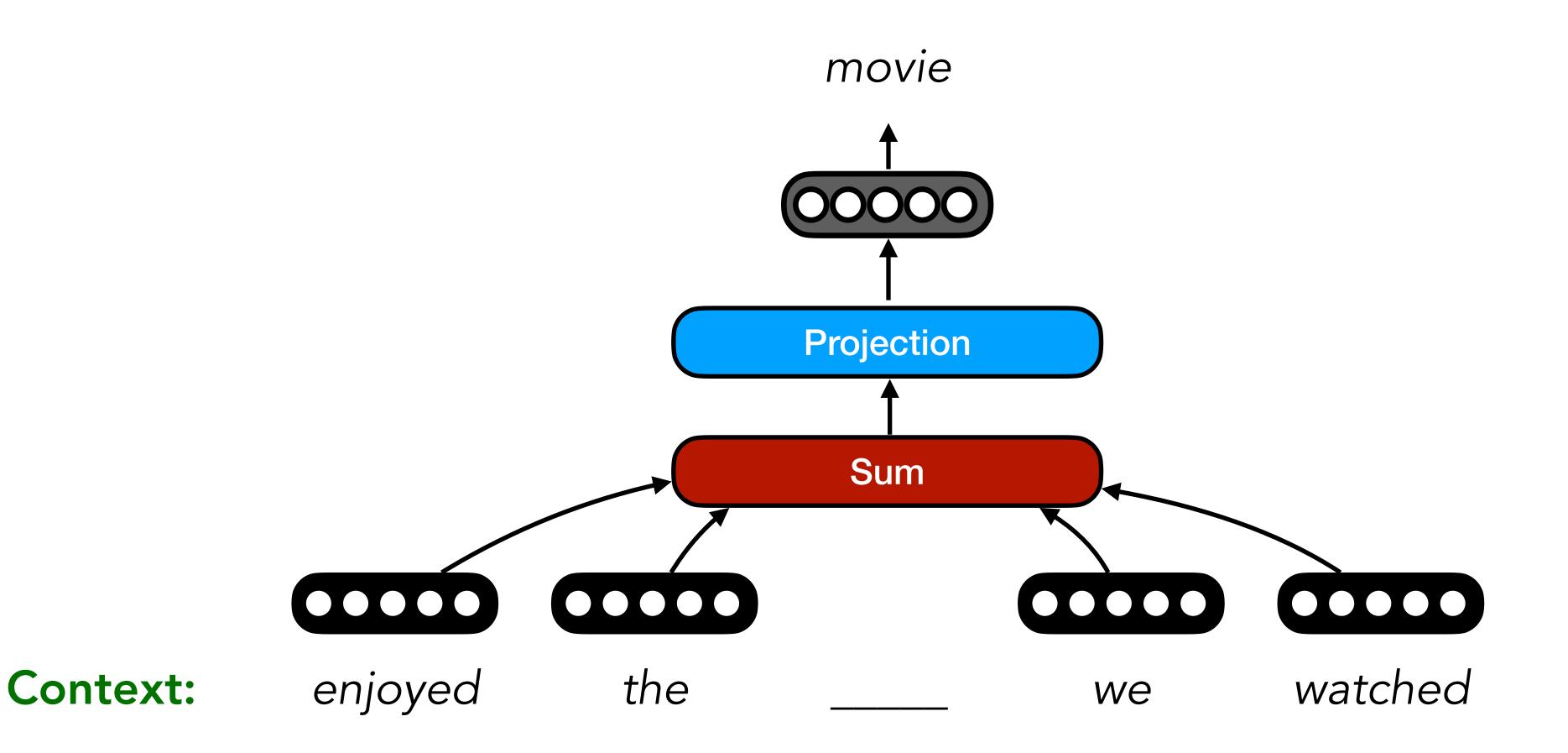
Foundation of distributional semantics

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

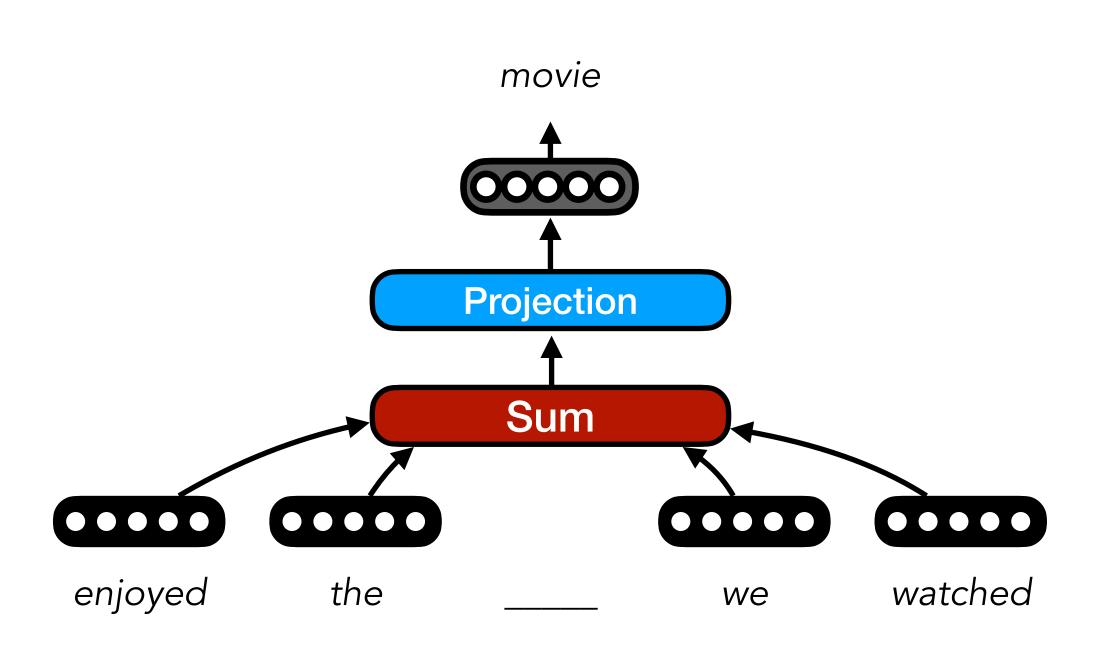
Continuous Bag of Words (CBOW)

Predict the missing word from a window of surrounding words



Continuous Bag of Words (CBOW)

Predict the missing word from a window of surrounding words



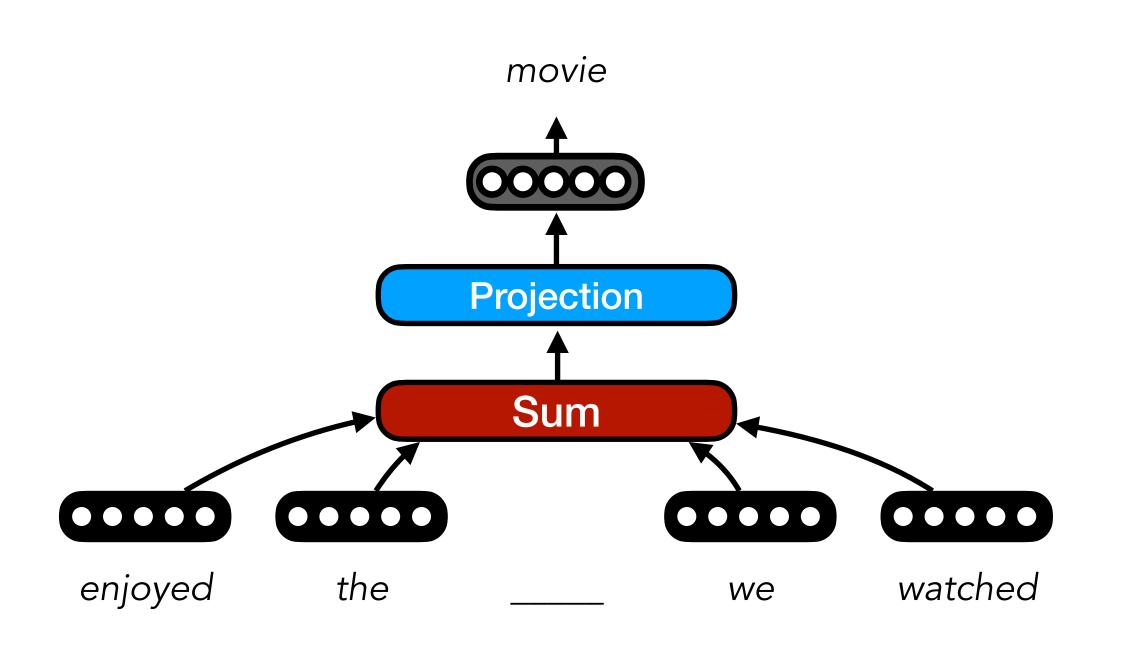
max *P*(*movie* | *enjoyed*, *the*, *we*, *watched*)

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

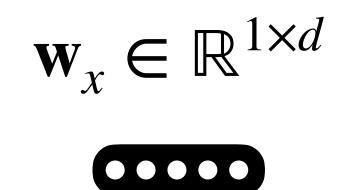
$$\max P(w_t | \{w_x\}_{x=t-2}^{x=t+2})$$

Continuous Bag of Words (CBOW)

Predict the missing word from a window of surrounding words



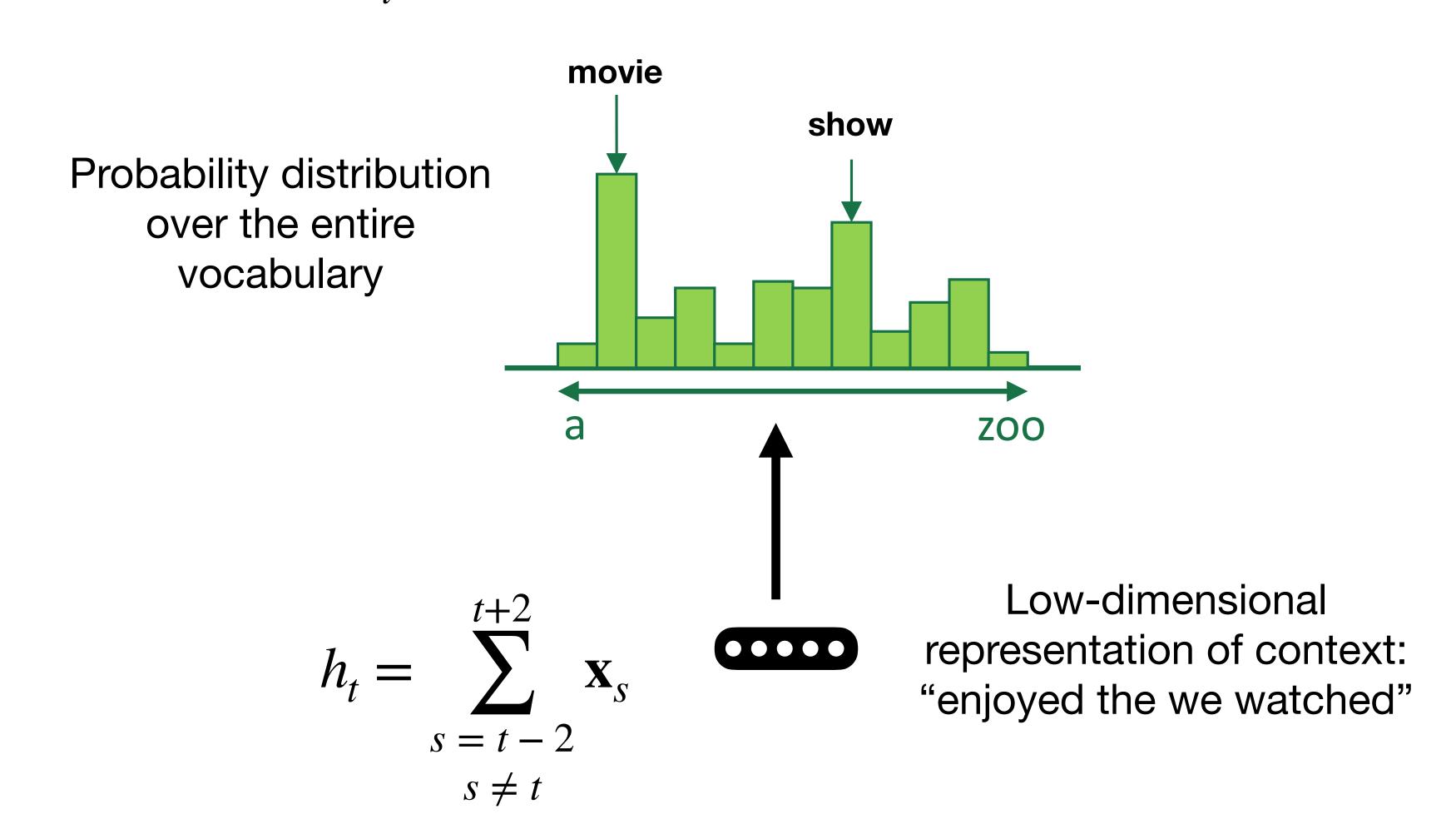
$$P(w_t | \{w_x\}_{x=t-2}^{x=t+2}) = \mathbf{softmax} \left(\mathbf{U} \sum_{\substack{x=t-2 \ x \neq t}}^{t+2} \mathbf{w}_x \right)$$





Vocabulary Space Projection

 $P(w_i | \text{vector for "enjoyed the we watched"})$



Let's say our output vocabulary consists of just four words: "movie", "show", "book", and "shelf".

$$h_t = \sum_{s=t-2}^{t+2} \mathbf{x}_s$$

$$s = t-2$$

$$s \neq t$$

Low-dimensional representation of context: "enjoyed the we watched"

Let's say our output vocabulary consists of just four words: "movie", "show", "book", and "shelf".

movie show book shelf <0.6, 0.2, 0.1, 0.1>

We want to get a probability distribution over these four words



Low-dimensional representation of context: "enjoyed the we watched"

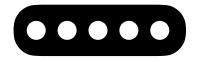
Let's say our output vocabulary consists of just four words: "movie", "show", "book", and "shelf".

$$\mathbf{U} = \left\{ \begin{array}{l} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{array} \right\}$$

first, we'll project our 3-d context representation to 4-d with a matrix-vector product

$$h_t = \langle -2.3, 0.9, 5.4 \rangle$$

Here's an example 3-d prefix vector



How do we get there?

$$\mathbf{U} = \left\{ \begin{array}{l} 1.2, -0.3, 0.9 \\ 0.2, 0.4, -2.2 \\ 8.9, -1.9, 6.5 \\ 4.5, 2.2, -0.1 \end{array} \right\}$$

$$h_t = \langle -2.3, 0.9, 5.4 \rangle$$

intuition: each dimension of h_t corresponds to a feature of the context

How do we get there?

$$h_t = \langle -2.3, 0.9, 5.4 \rangle$$

intuition: each dimension of h_t corresponds to a feature of the context

$$\mathbf{U}h_t = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

How did we compute this? It's just the dot product of each row of ${\bf U}$ with h_t

$$h_t = \langle -2.3, 0.9, 5.4 \rangle$$

$$\mathbf{U}h_t = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

How did we compute this? It's just the dot product of each row of ${\bf U}$ with h_t

$$U = \begin{cases} 1.2, & -0.3, & 0.9 \\ 0.2, & 0.4, & -2.2 \\ 8.9, & -1.9, & 6.5 \\ 4.5, & 2.2, & -0.1 \end{cases}$$

$$h_t = \langle -2.3, 0.9, 5.4 \rangle$$

$$\mathbf{U}h_t = \langle 1.8, -11.9, 12.9, -8.9 \rangle$$

How did we compute this? It's just the dot product of each row of ${\bf U}$ with h_t

$$\mathbf{U} = \begin{cases} 1.2, & -0.3, & 0.9 \\ 0.2, & 0.4, & -2.2 \\ 8.9, & -1.9, & 6.5 \\ 4.5, & 2.2, & -0.1 \end{cases}$$

$$h_t = \langle -2.3, 0.9, 5.4 \rangle$$

Softmax

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

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

- a is a vector
- a_i is dimension *i* of **a**
- each dimension *i* of the softmaxed output represents the probability of class *i*

Softmax

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

$$\mathbf{softmax(a)}_{i} = \frac{e^{a_{i}}}{\sum_{j=1}^{|\mathbf{a}|} e^{a_{j}}}$$

- a is a vector
- a_i is dimension *i* of **a**
- each dimension *i* of the softmaxed output represents the probability of class *i*

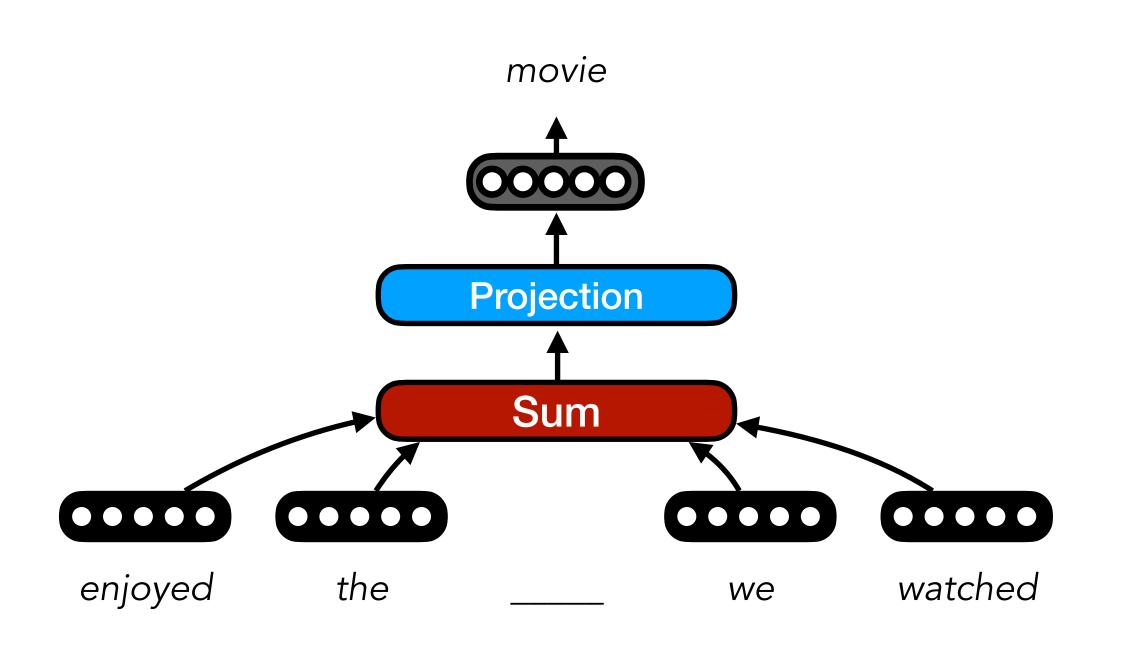
$$Uh_t = \langle 1.8, -1.9, 2.9, -0.9 \rangle$$

softmax(U h_t **)** = $\langle 0.24, 0.006, 0.73, 0.02 \rangle$

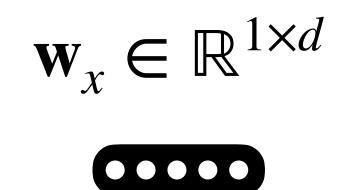
Softmax will keep popping up, so be sure to understand it!

Continuous Bag of Words (CBOW)

Predict the missing word from a window of surrounding words



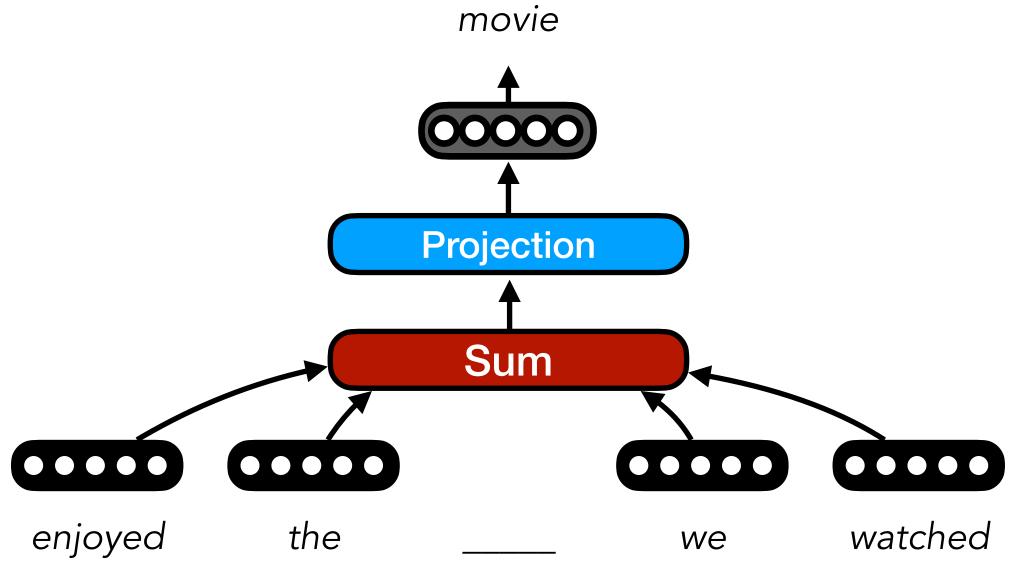
$$P(w_t | \{w_x\}_{x=t-2}^{x=t+2}) = \mathbf{softmax} \left(\mathbf{U} \sum_{\substack{x=t-2 \ x \neq t}}^{t+2} \mathbf{w}_x \right)$$





Continuous Bag of Words (CBOW)

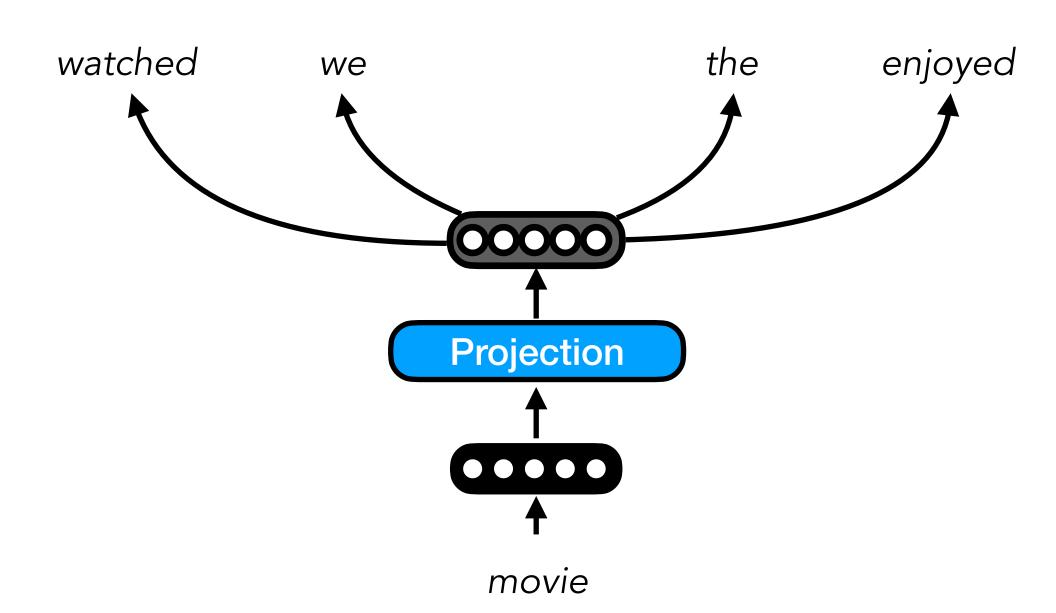
$$P(w_t | \{w_x\}_{x=t-2}^{x=t+2}) = \mathbf{softmax} \left(\mathbf{U} \sum_{\substack{x=t-2 \ x \neq t}}^{t+2} \mathbf{w}_x \right)$$



- 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

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

Context:

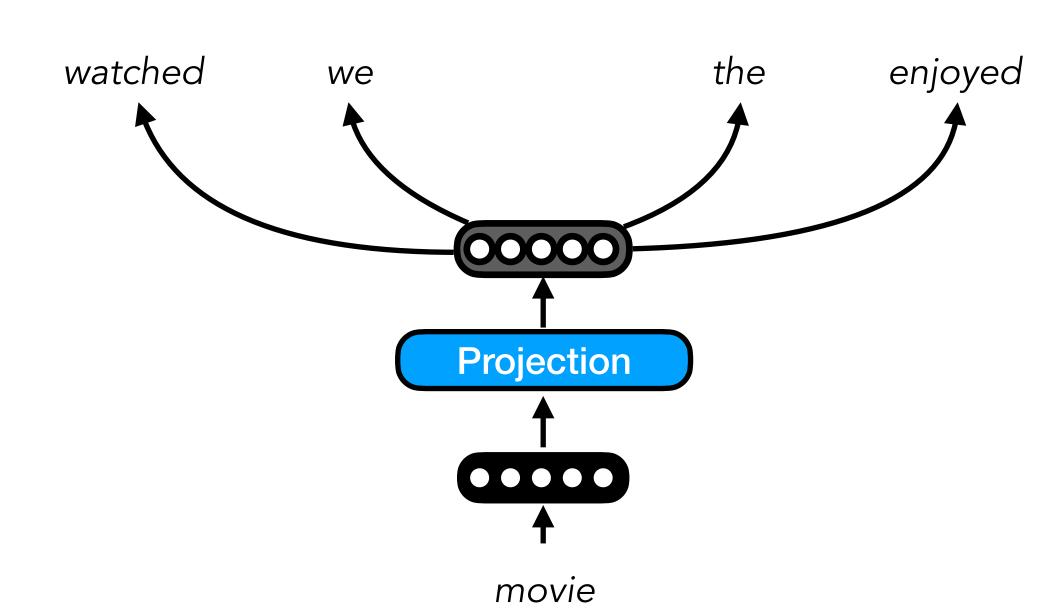


max *P*(*enjoyed*, *the*, *we*, *watched* | *movie*)

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

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

Context:



max *P*(*enjoyed*, *the*, *we*, *watched* | *movie*)

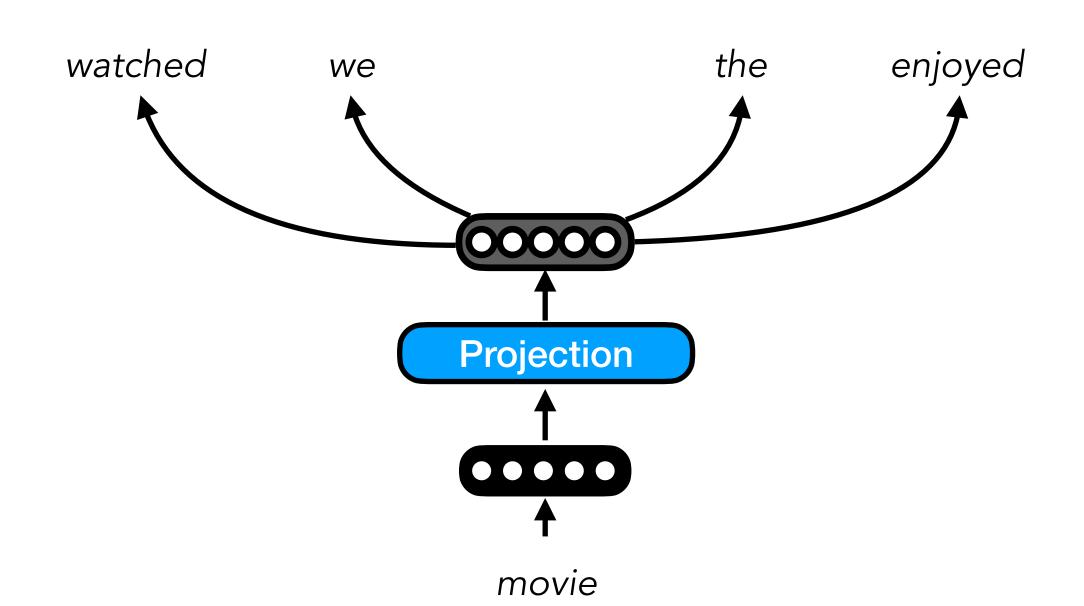
$$\max P(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2} | w_t)$$

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

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

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

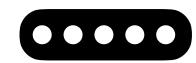
Context:



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

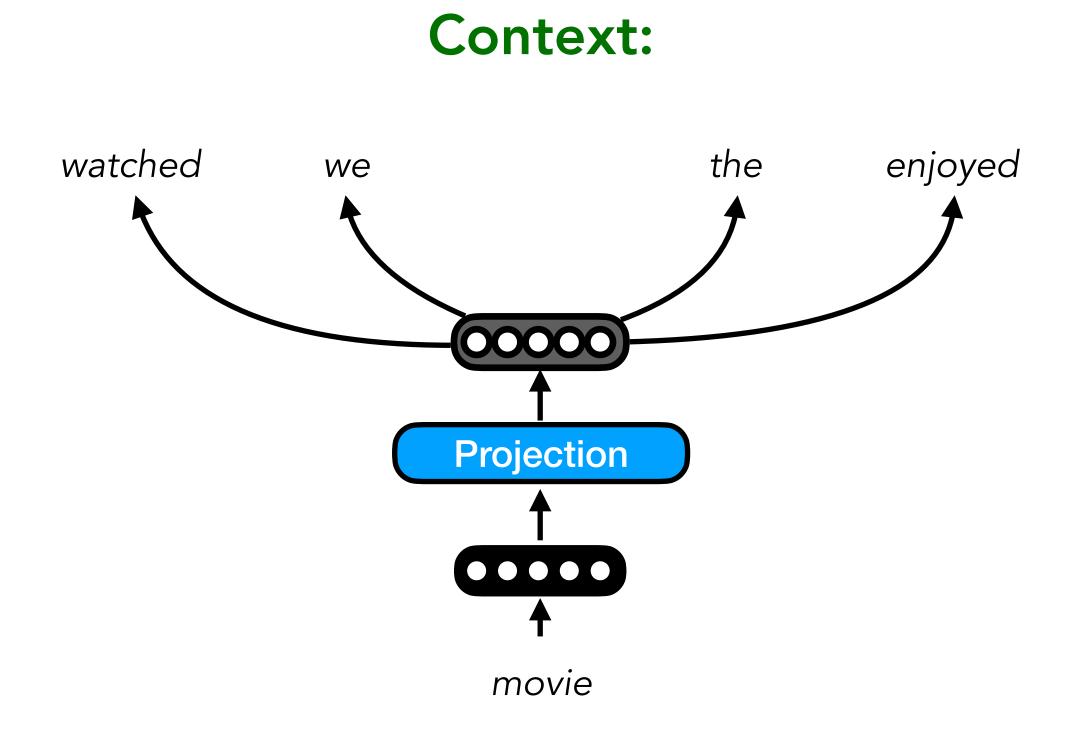
$$\mathbf{w}_t \in \mathbb{R}^{1 \times d}$$





Projection

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



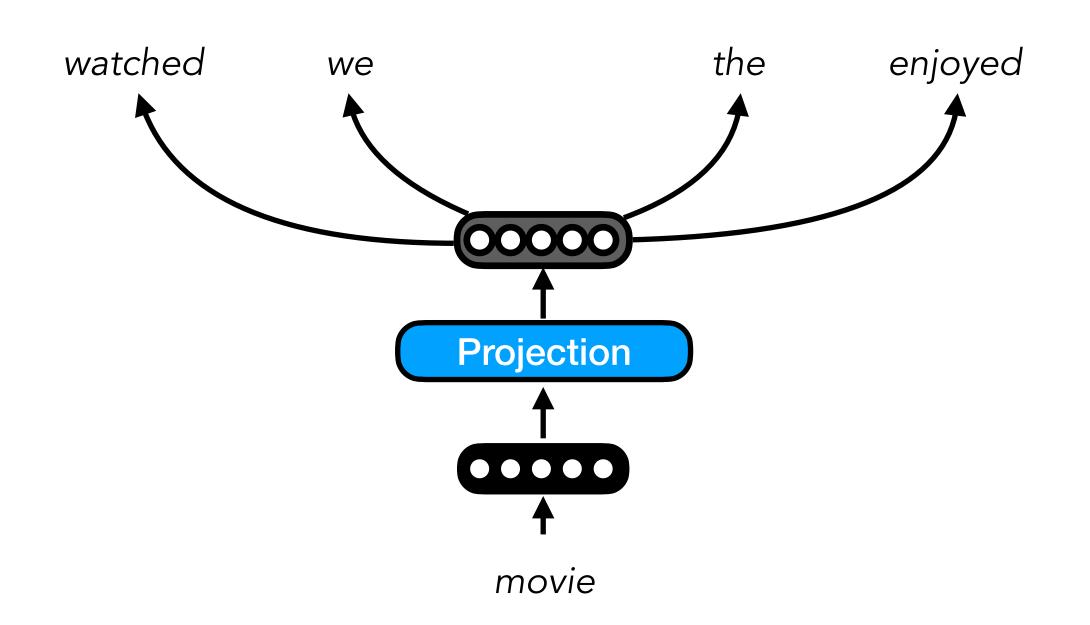
- Model is trained to minimise the negative log probability of the surrounding words
- Here, we use a window of N=2, but the window size is a hyperparameter to set
- Typically, set large window (N=10), but randomly select $i \in [1,N]$ as dynamic window size so that closer words contribute more to learning

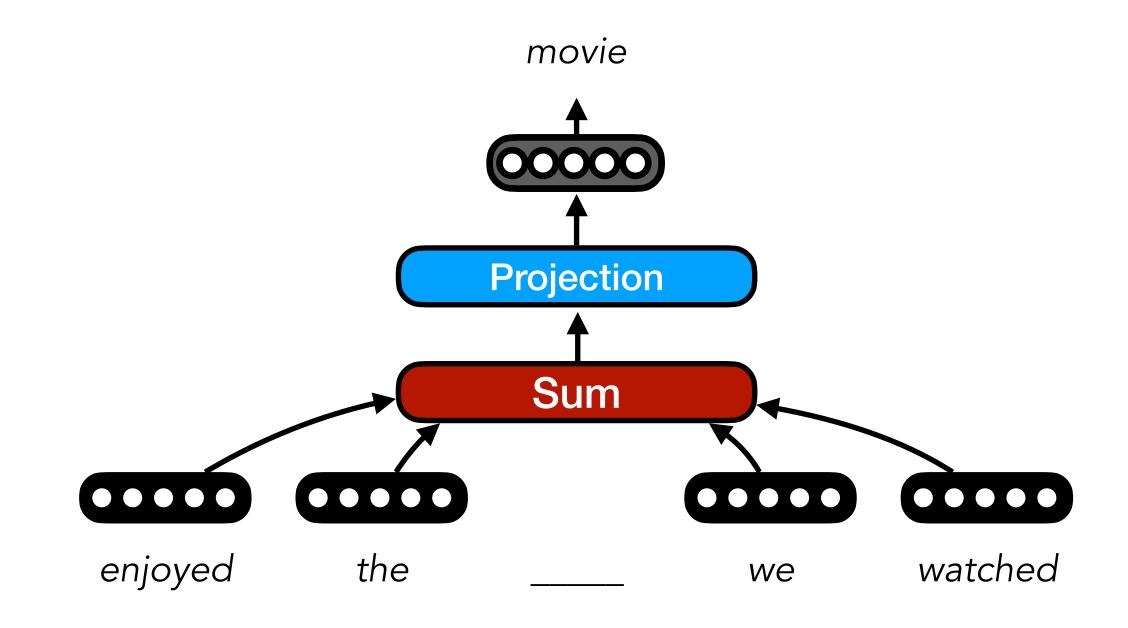
Question

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?





Example

CBOW

```
top_cbow = cbow.wv.most_similar('cut', topn=10)
print(tabulate(top_cbow, headers=["Word", "Simi
              Similarity
Word
slice
                0.662173
                0.650036
crosswise
                0.630569
score
                0.618827
tear
                0.563946
dice
lengthwise
                0.557231
                0.557228
cutting
break
                0.551517
                0.541566
chop
                0.537967
carve
```

Skip-gram

```
top_sg = skipgram.wv.most_similar('cut', topn=10)
print(tabulate(top_sg, headers=["Word", "Similarit
Word
                  Similarity
                    0.72921
crosswise
                    0.702693
score
slice
                    0.696898
                    0.680091
crossways
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 <u>Public Domain Dedication and License</u> 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 <u>script</u> 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., Puhrsch, 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





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

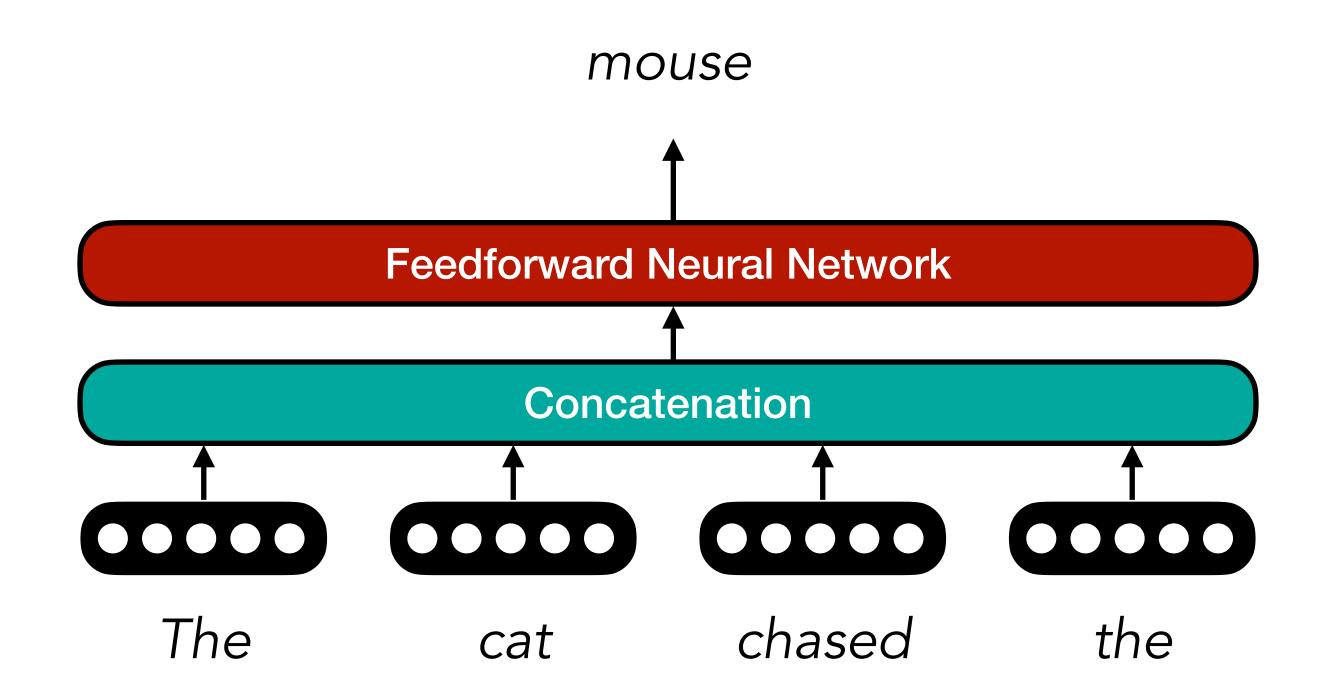
Language Modeling

• 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 _____

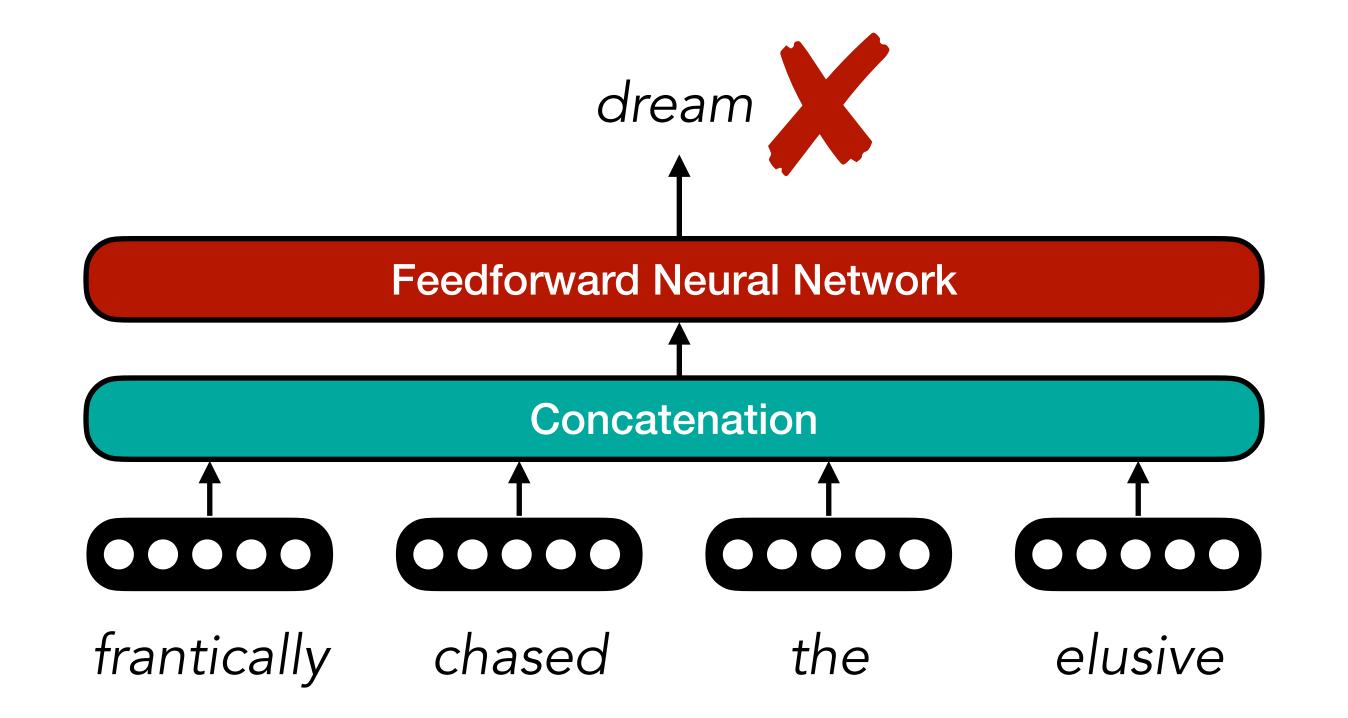
$$P(y) = \mathbf{softmax}(b_o + \mathbf{W}_o \ \mathbf{tanh}(b_h + \mathbf{W}_h x))$$



Fixed Context Language Models

• Given a subsequence, predict the next word:

The starving cat frantically chased the elusive _____



The starving cat

(Bengio et al., 2003)

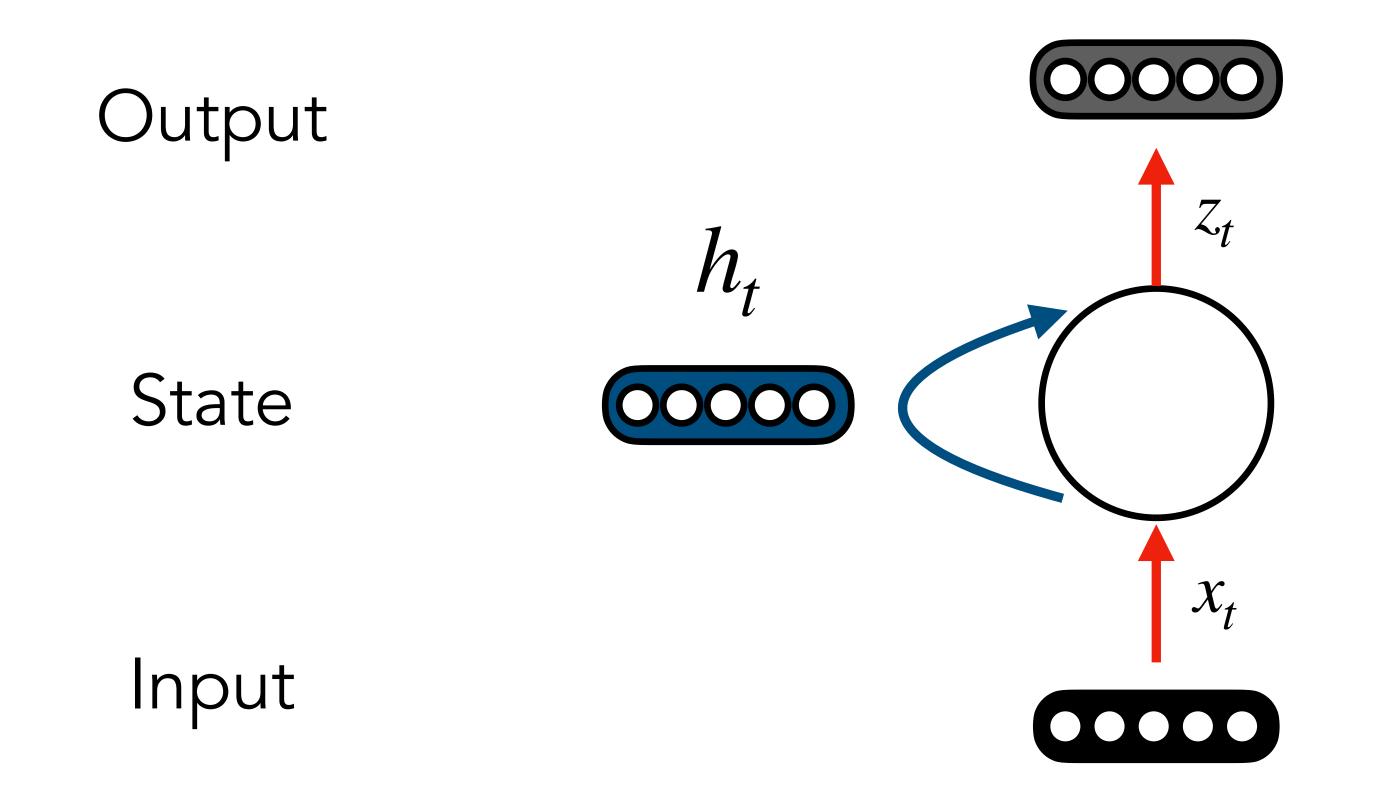
Problem

Fixed context windows limit language modelling capacity

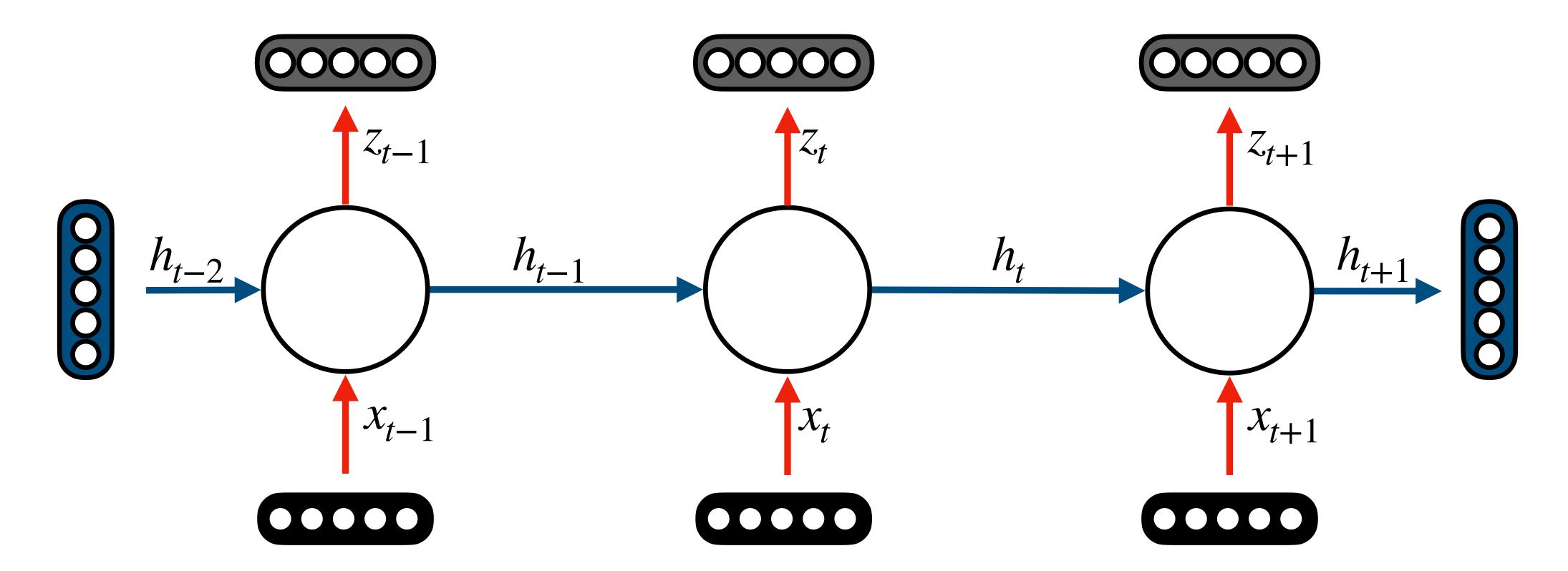
How can we extend to arbitrary length sequences?

Recurrent Neural Networks

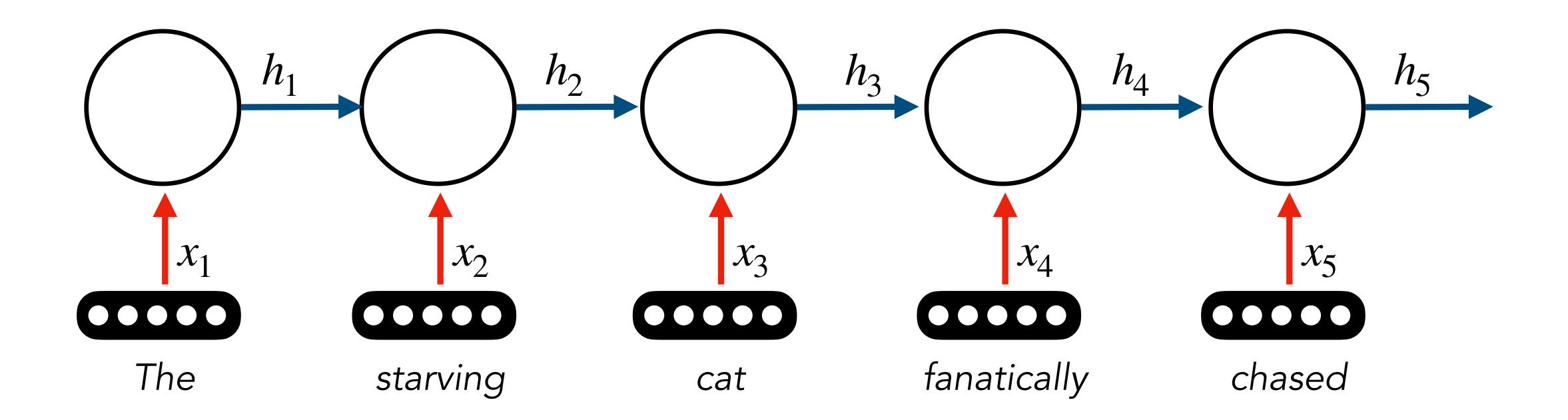
• Solution: Recurrent neural networks — NNs with feedback loops

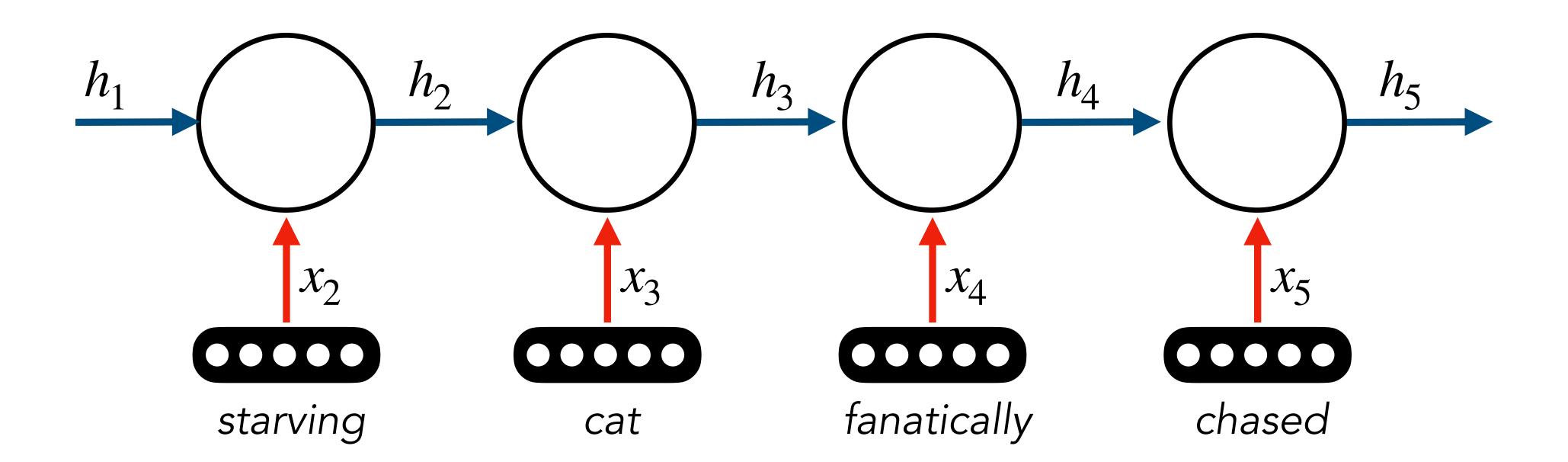


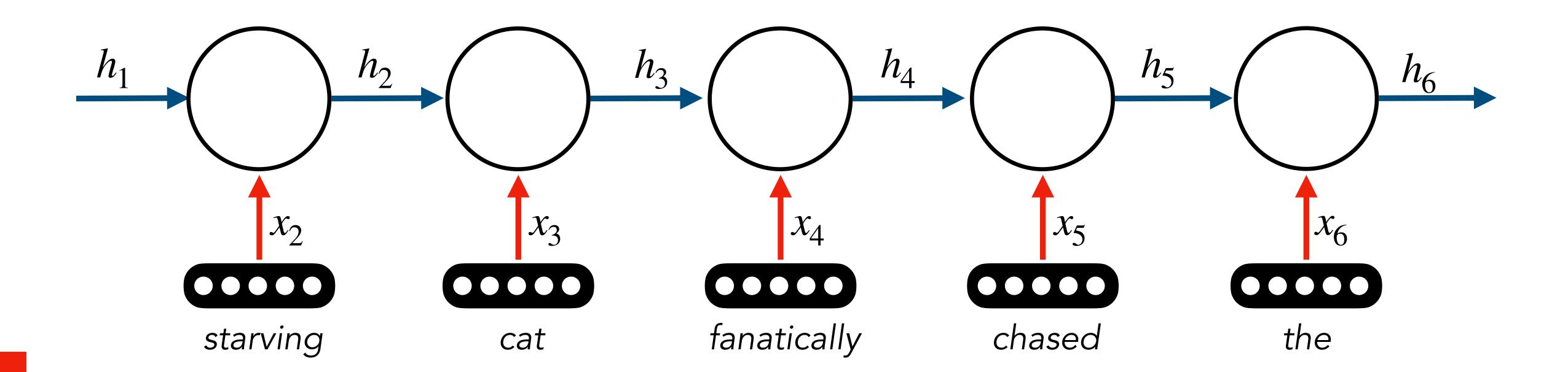
Unrolling the RNN across all time steps gives full computation graph



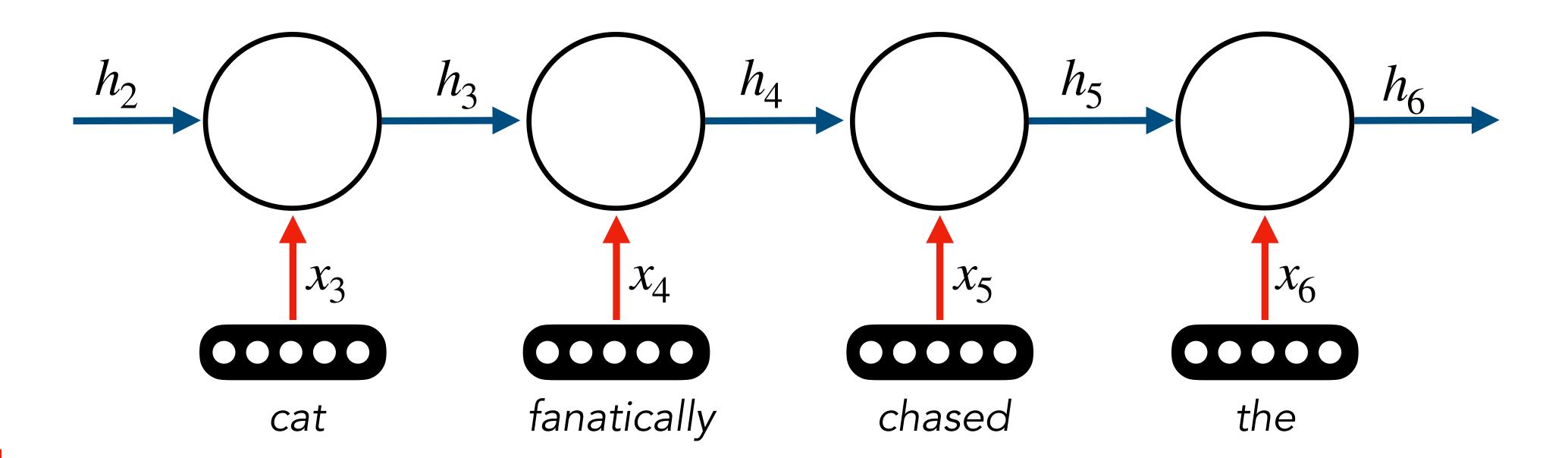
Allows for learning from entire sequence history, regardless of length



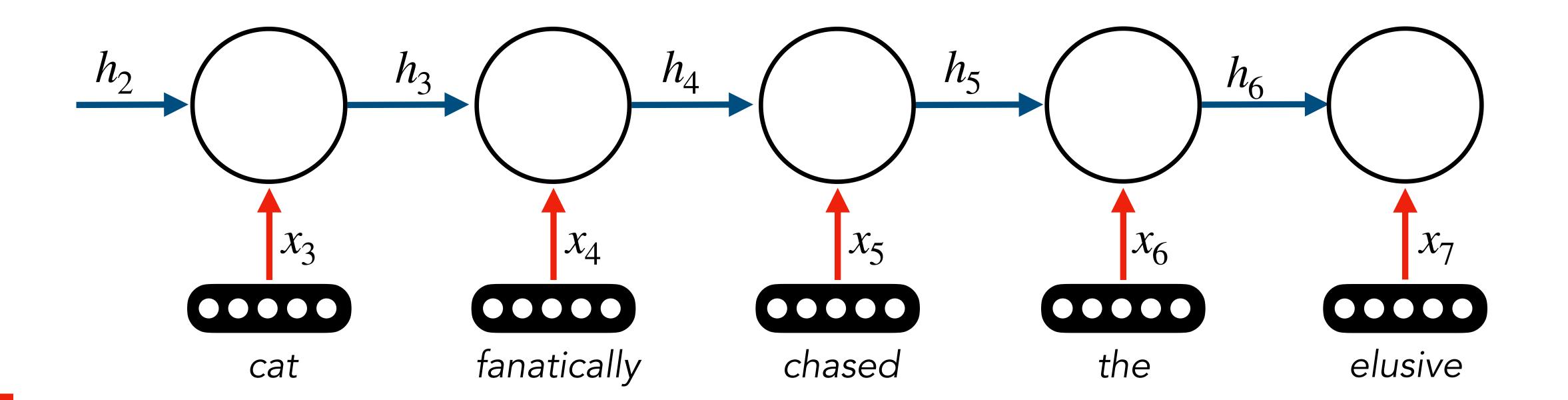




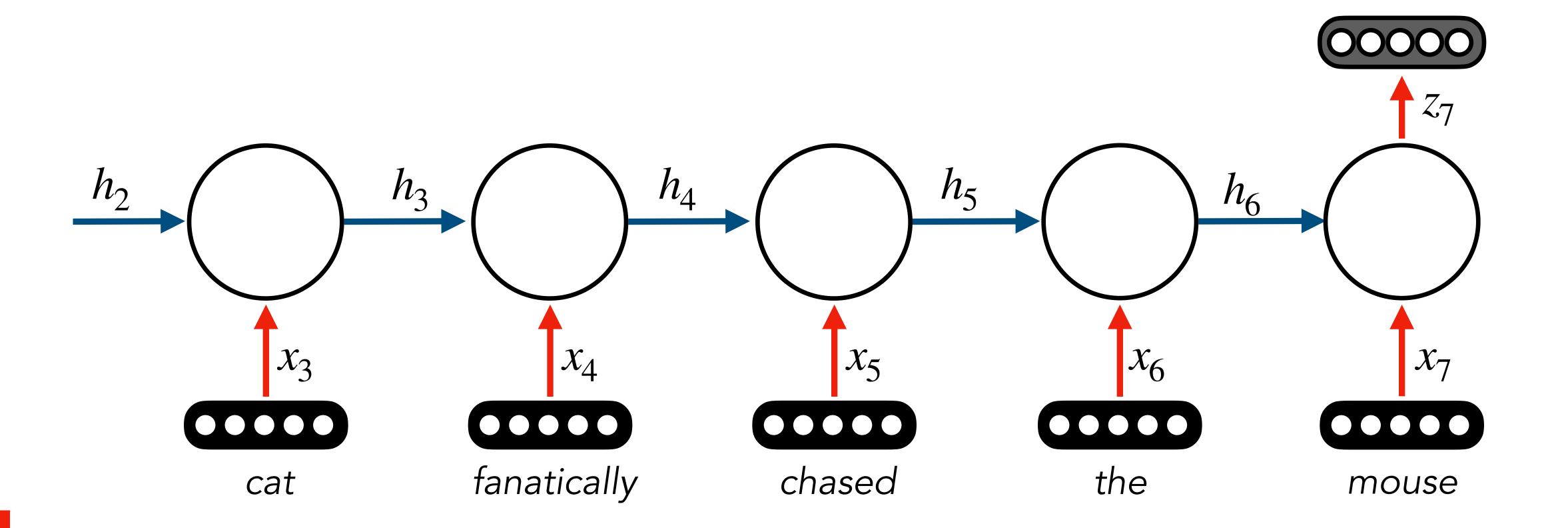
Unrolling the RNN



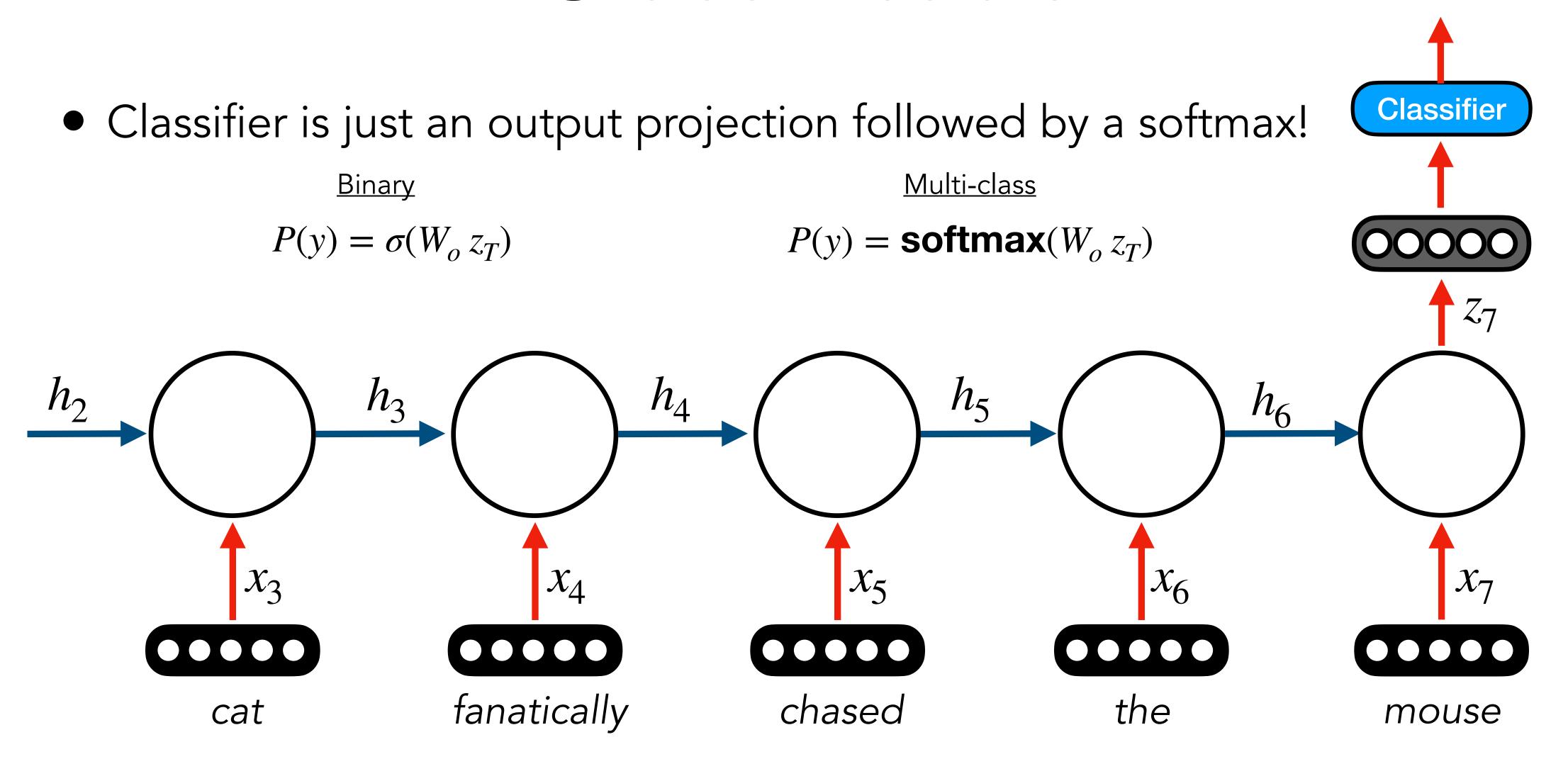
Unrolling the RNN



Unrolling the RNN



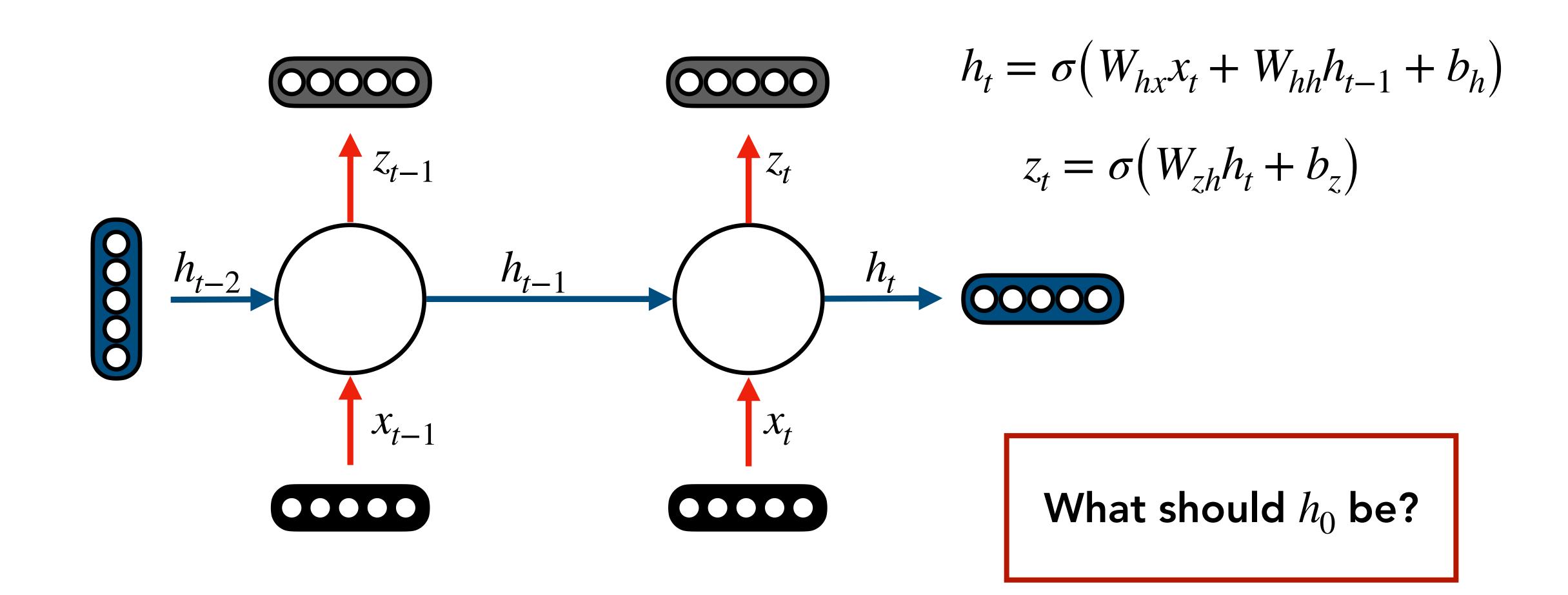
Classification



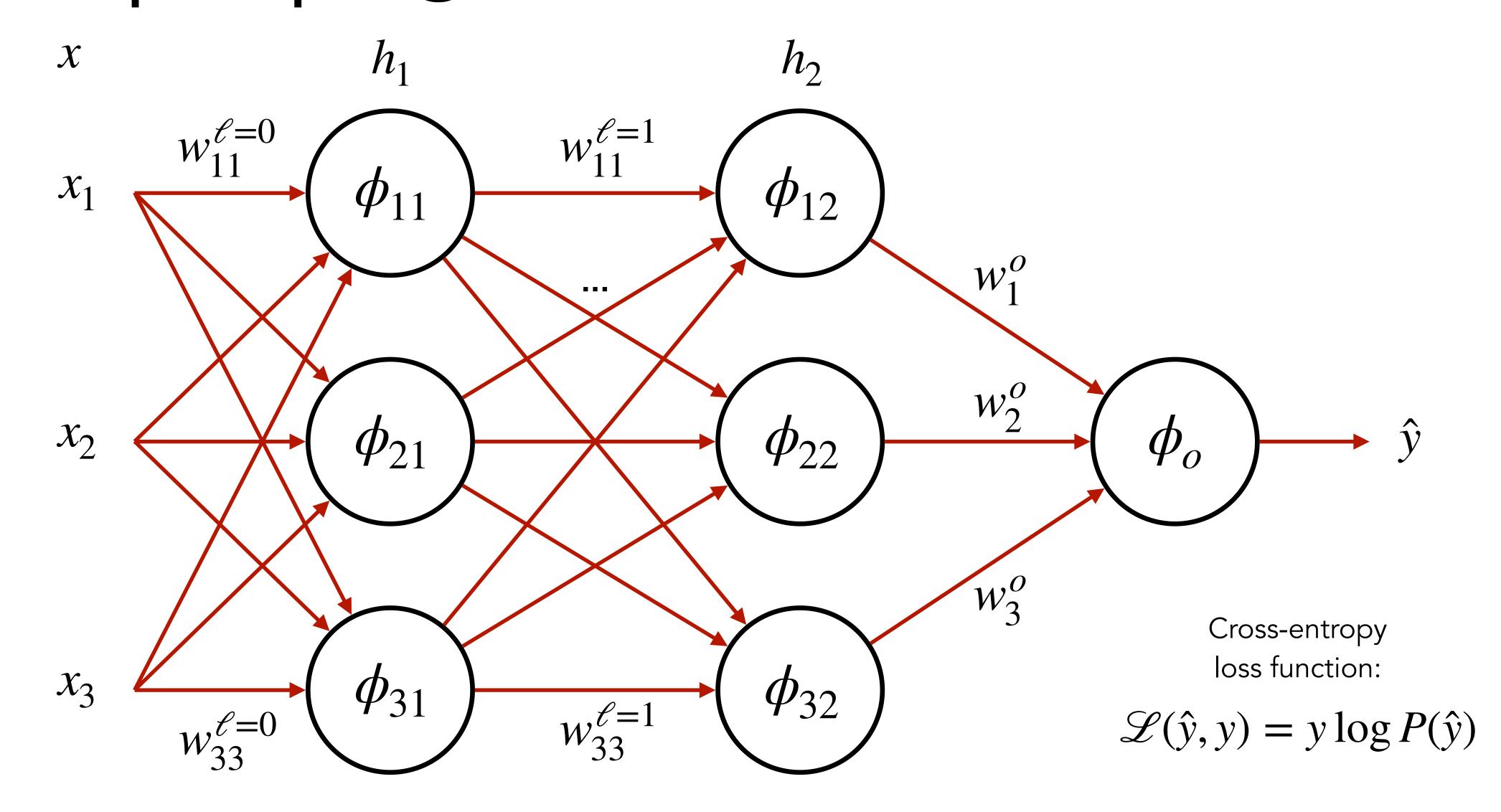
Question

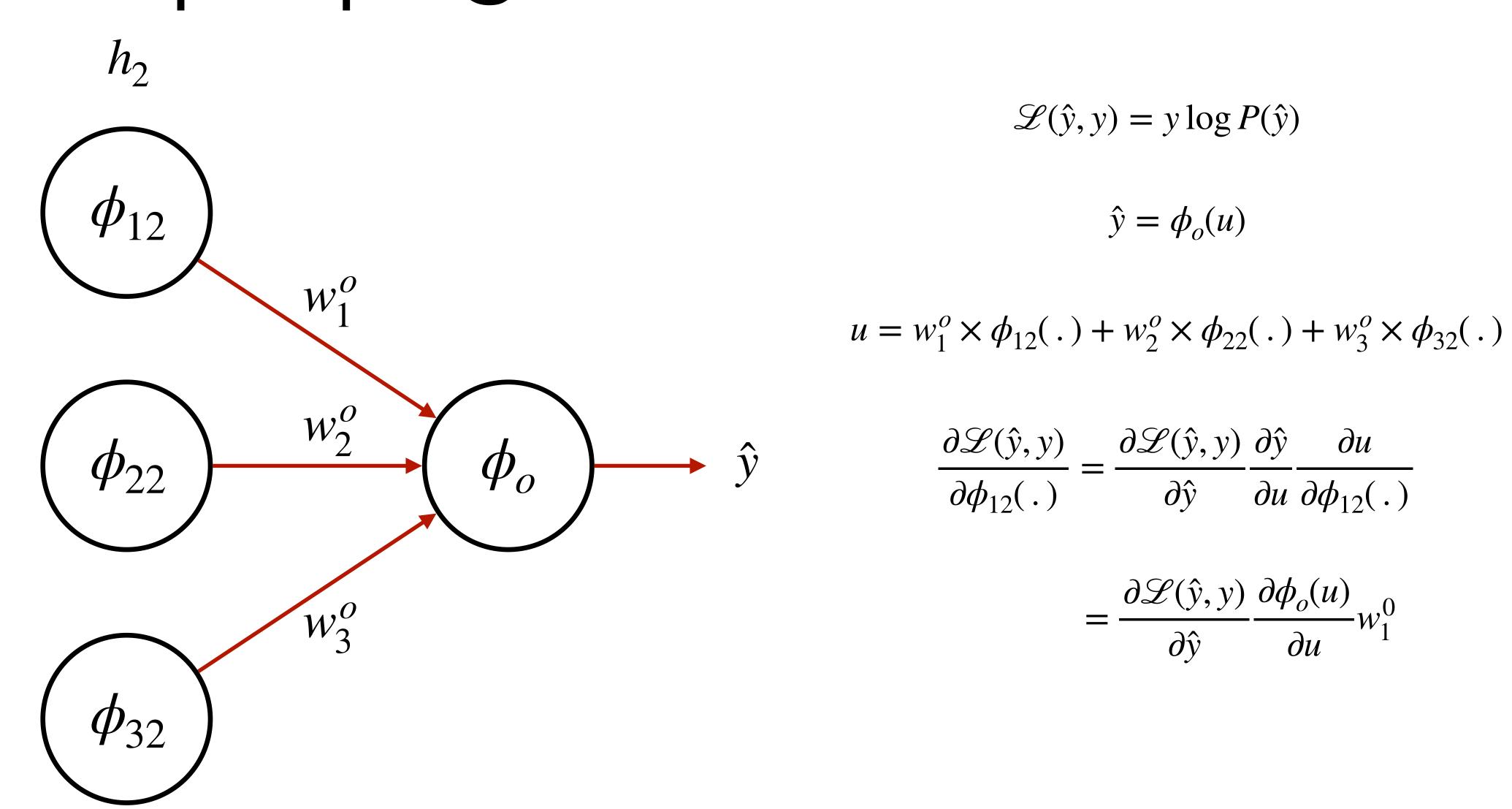
Why would you use the output of the last recurrent unit as the one to predict a label?

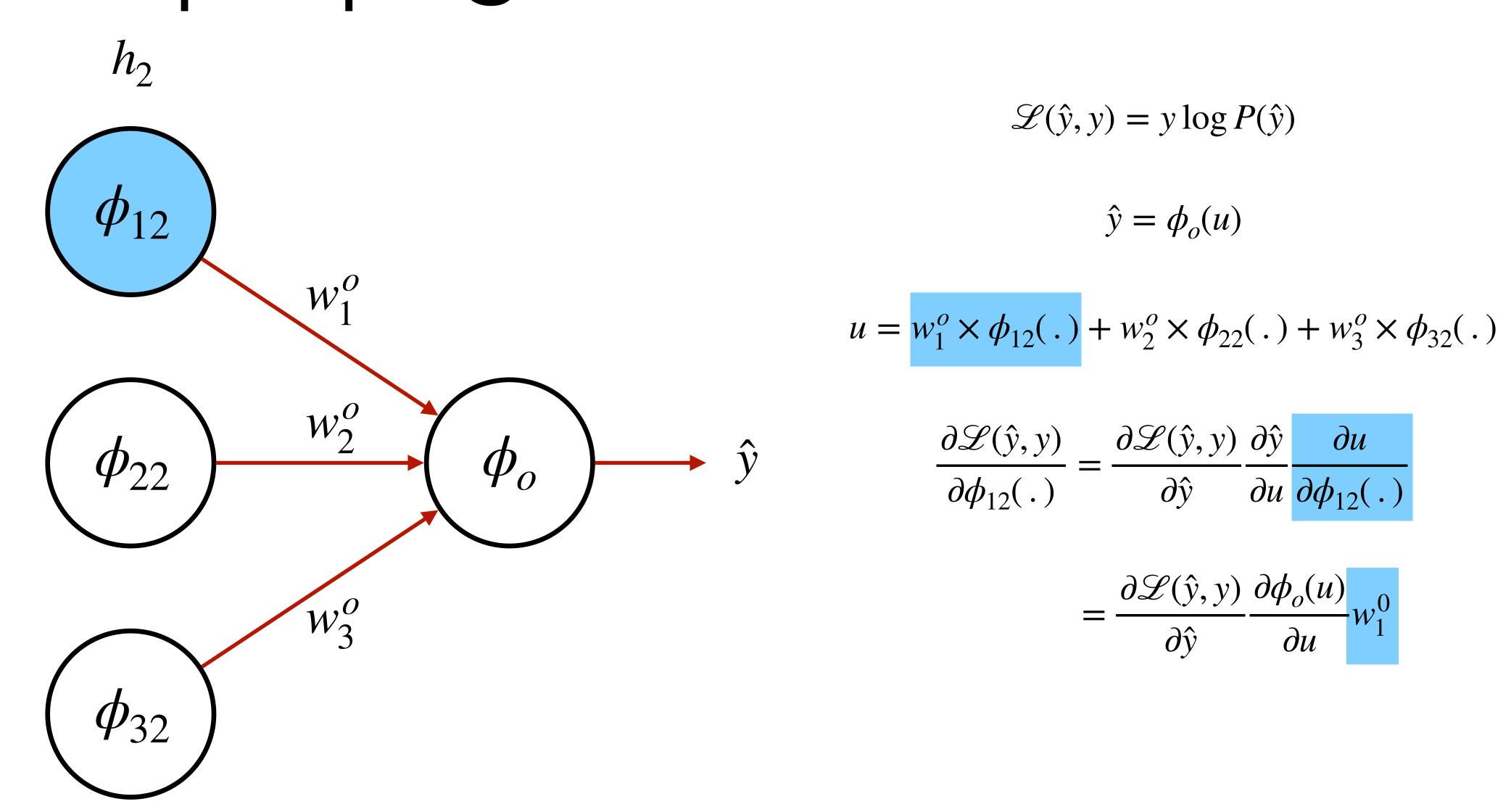
Classical RNN: Elman Network

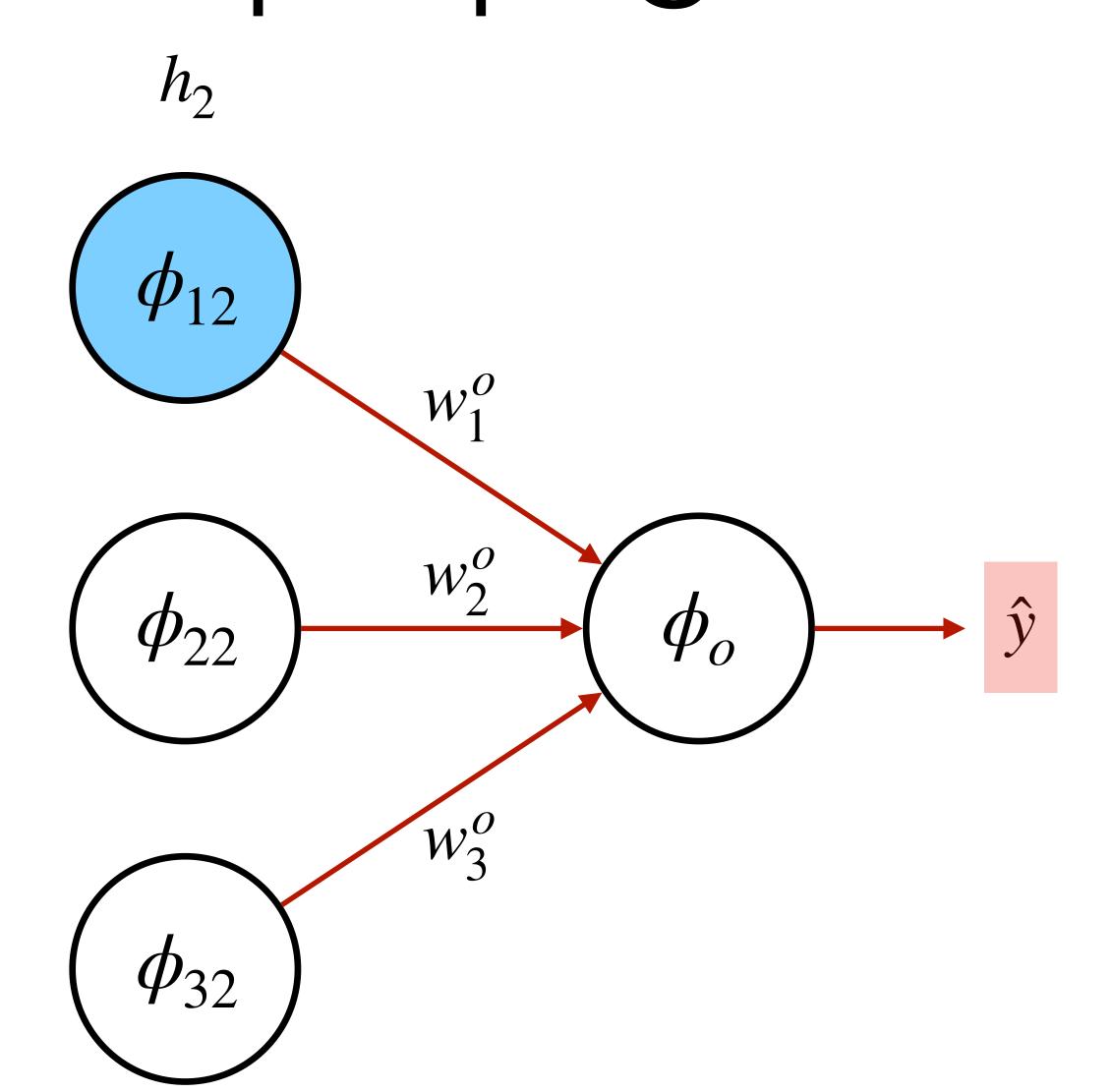


(<mark>E</mark>lman, 1990)









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

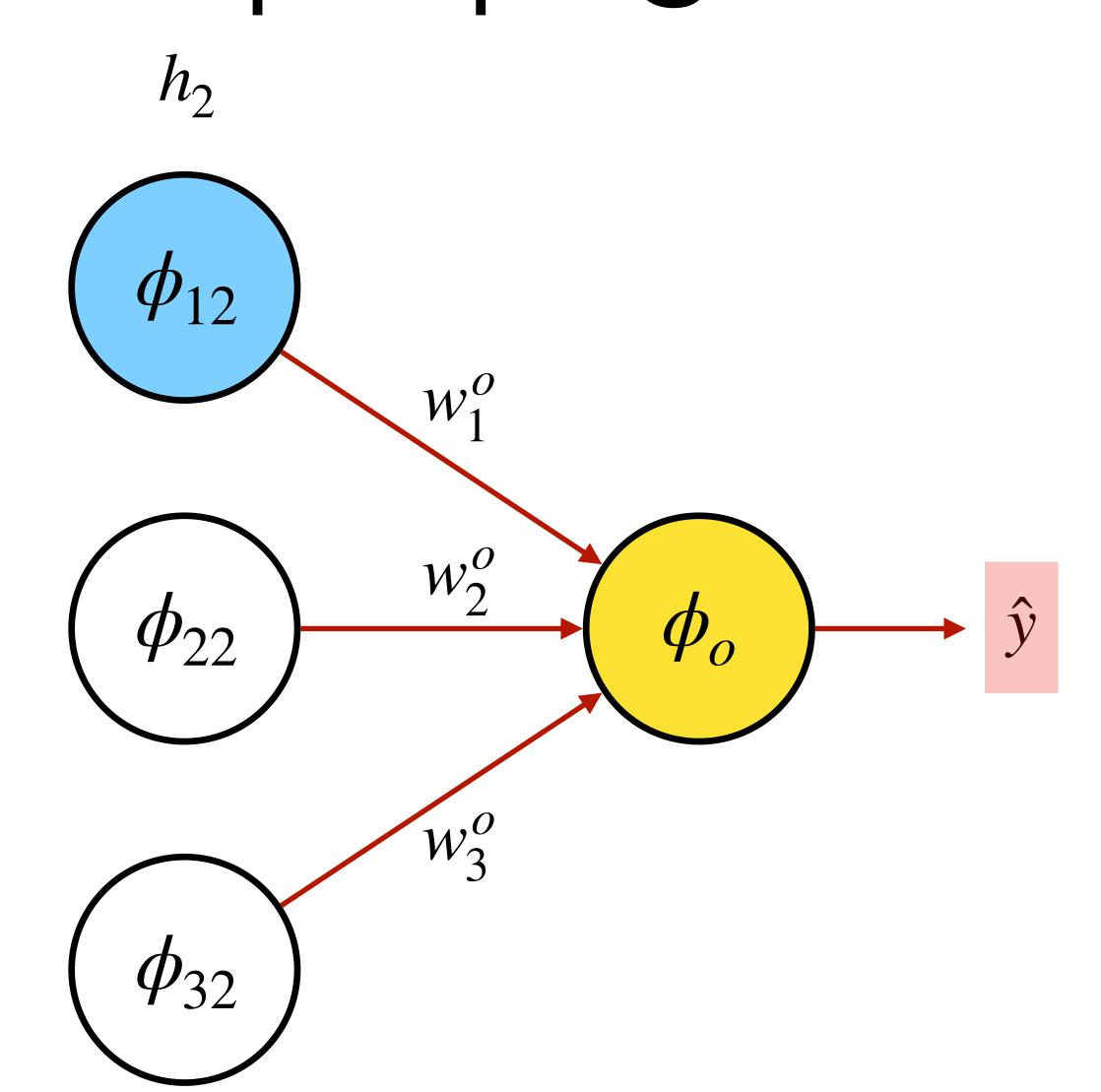
$$\hat{y} = \phi_o(u)$$

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

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

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^0$$

Depends on label y



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

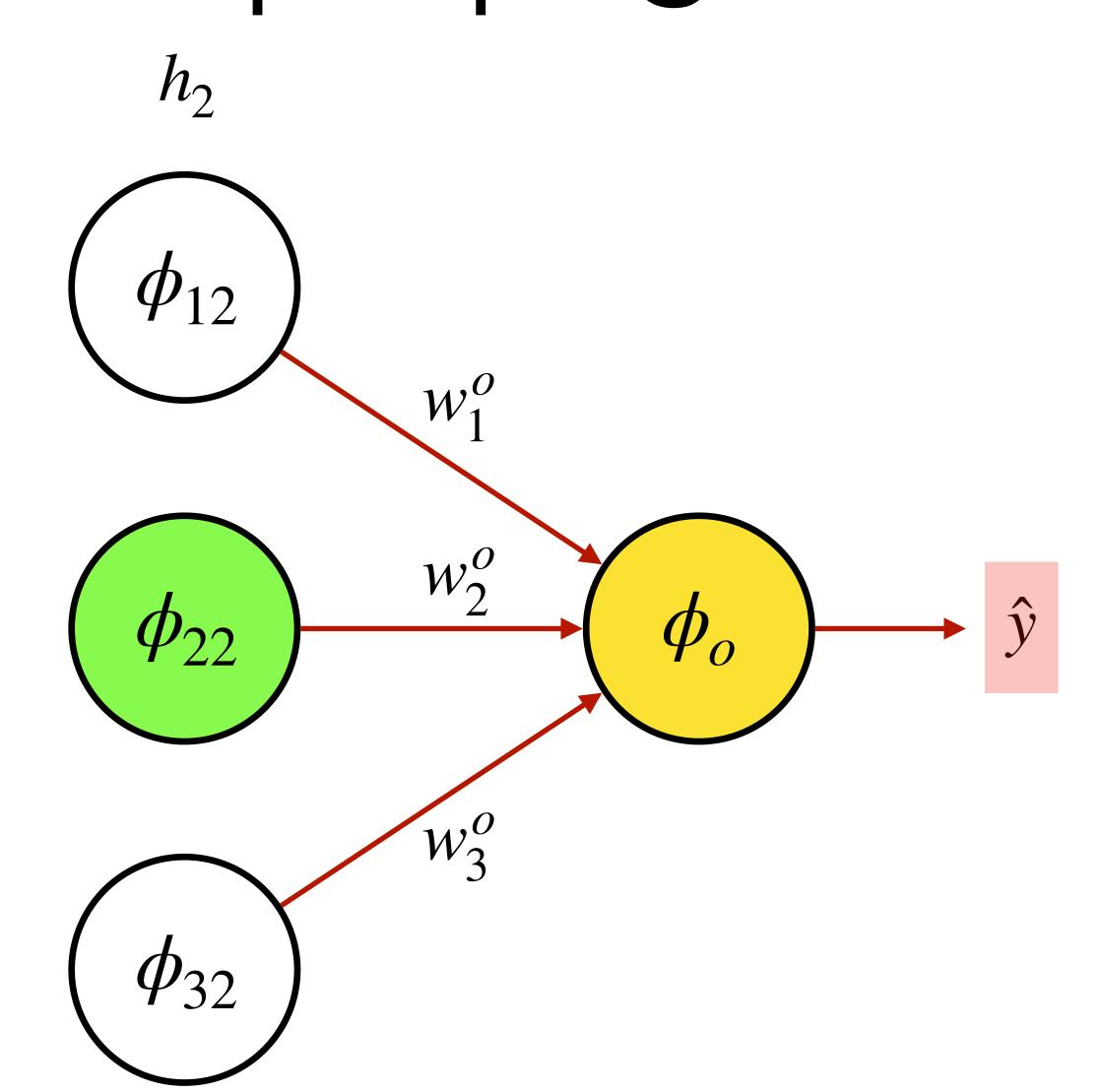
$$\hat{y} = \phi_o(u)$$

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

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

$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_1^0$$

Depends on label y



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

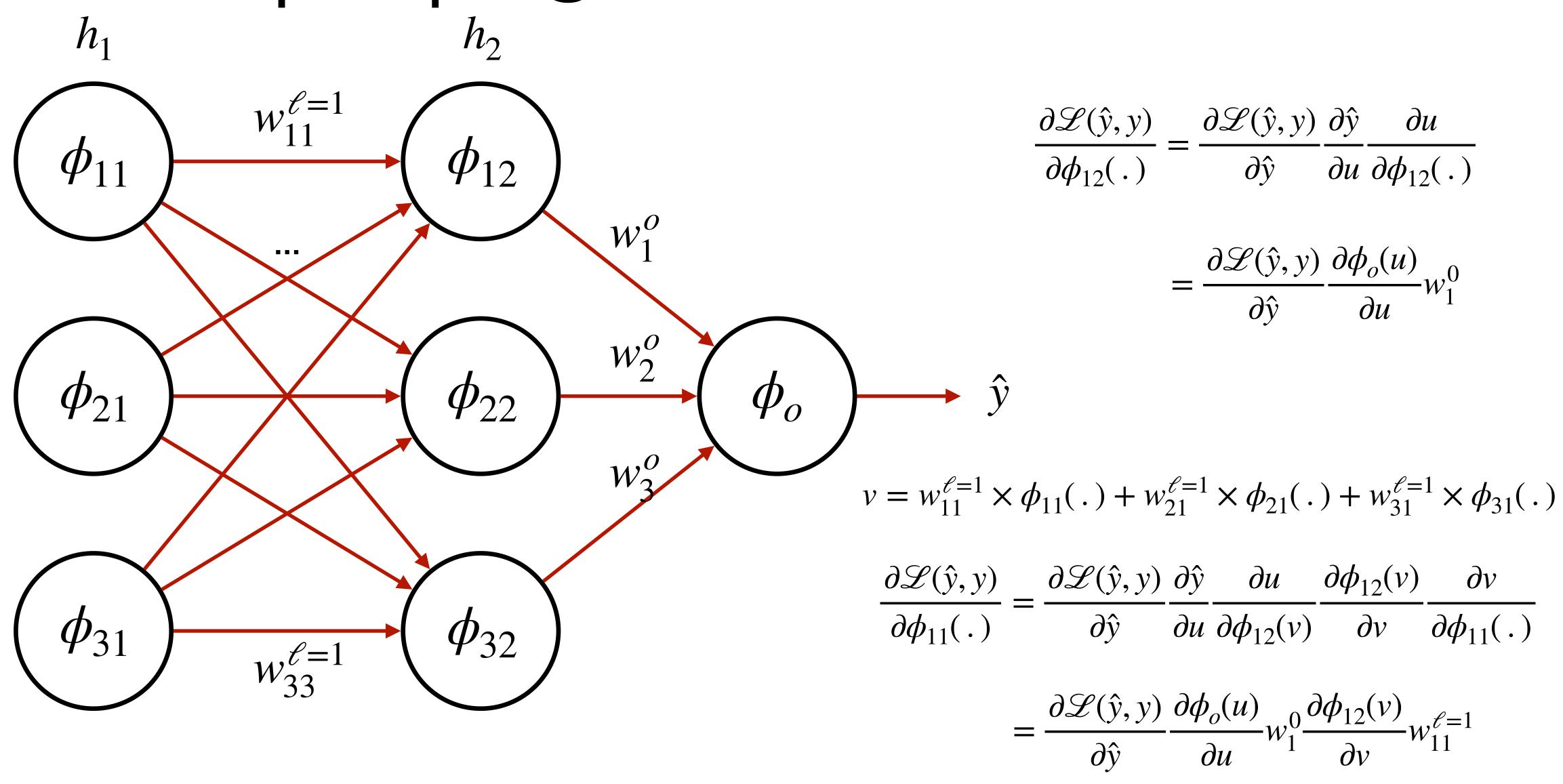
$$\hat{y} = \phi_o(u)$$

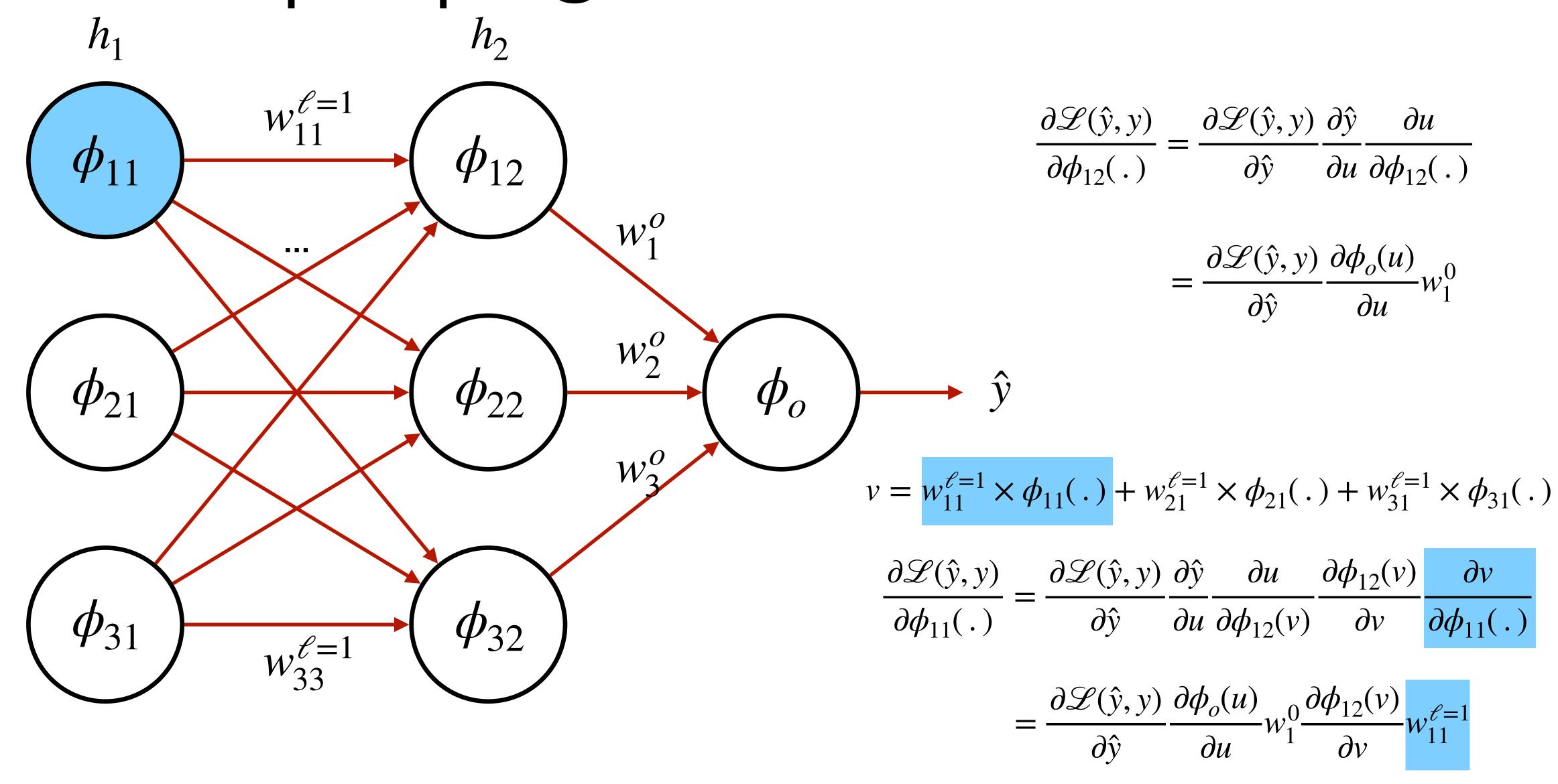
$$u = w_1^o \times \phi_{12}(.) + w_2^o \times \phi_{22}(.) + w_3^o \times \phi_{32}(.)$$

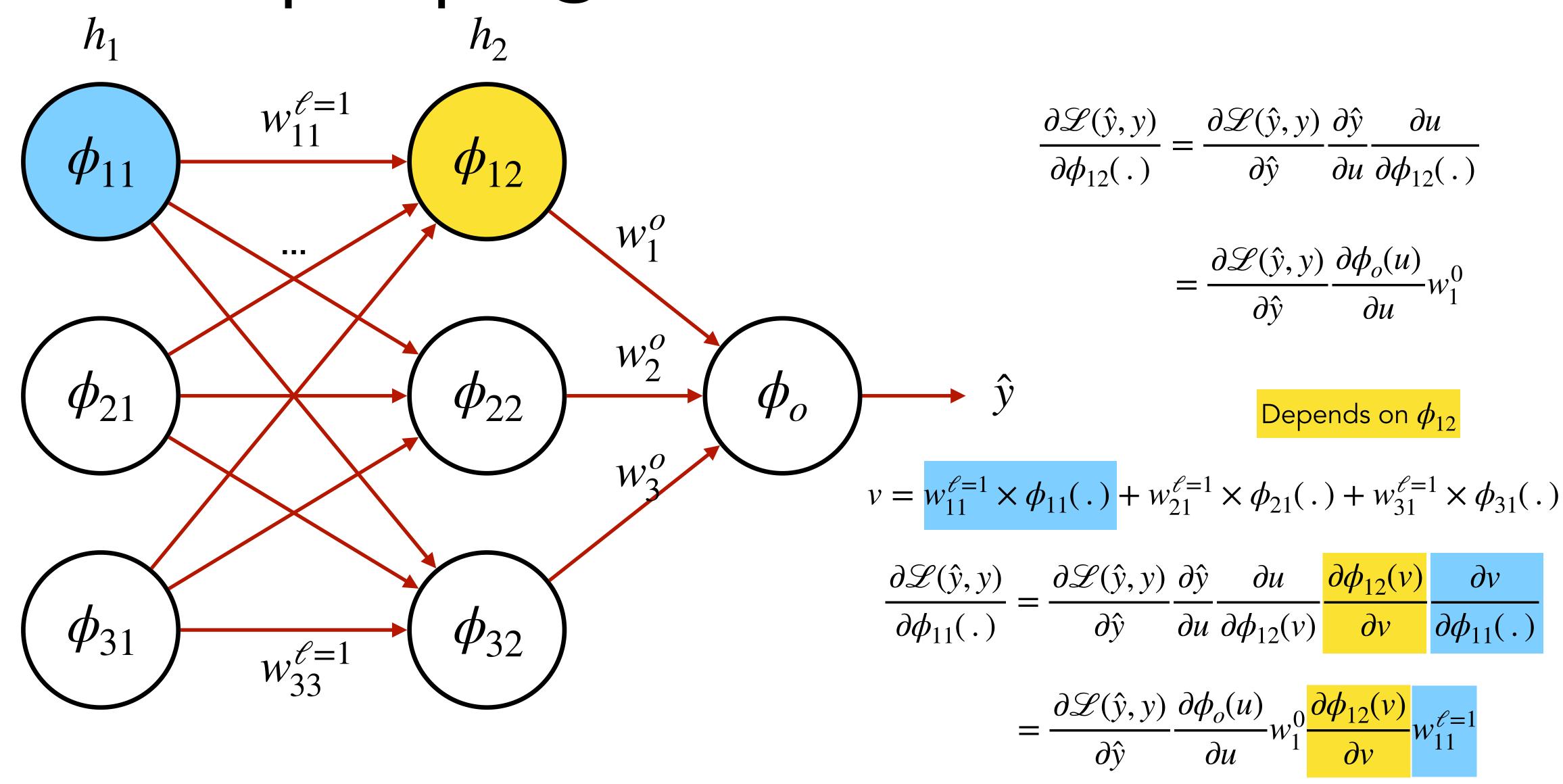
$$\frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \phi_{22}(.)} = \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial u} \frac{\partial u}{\partial \phi_{22}(.)}$$

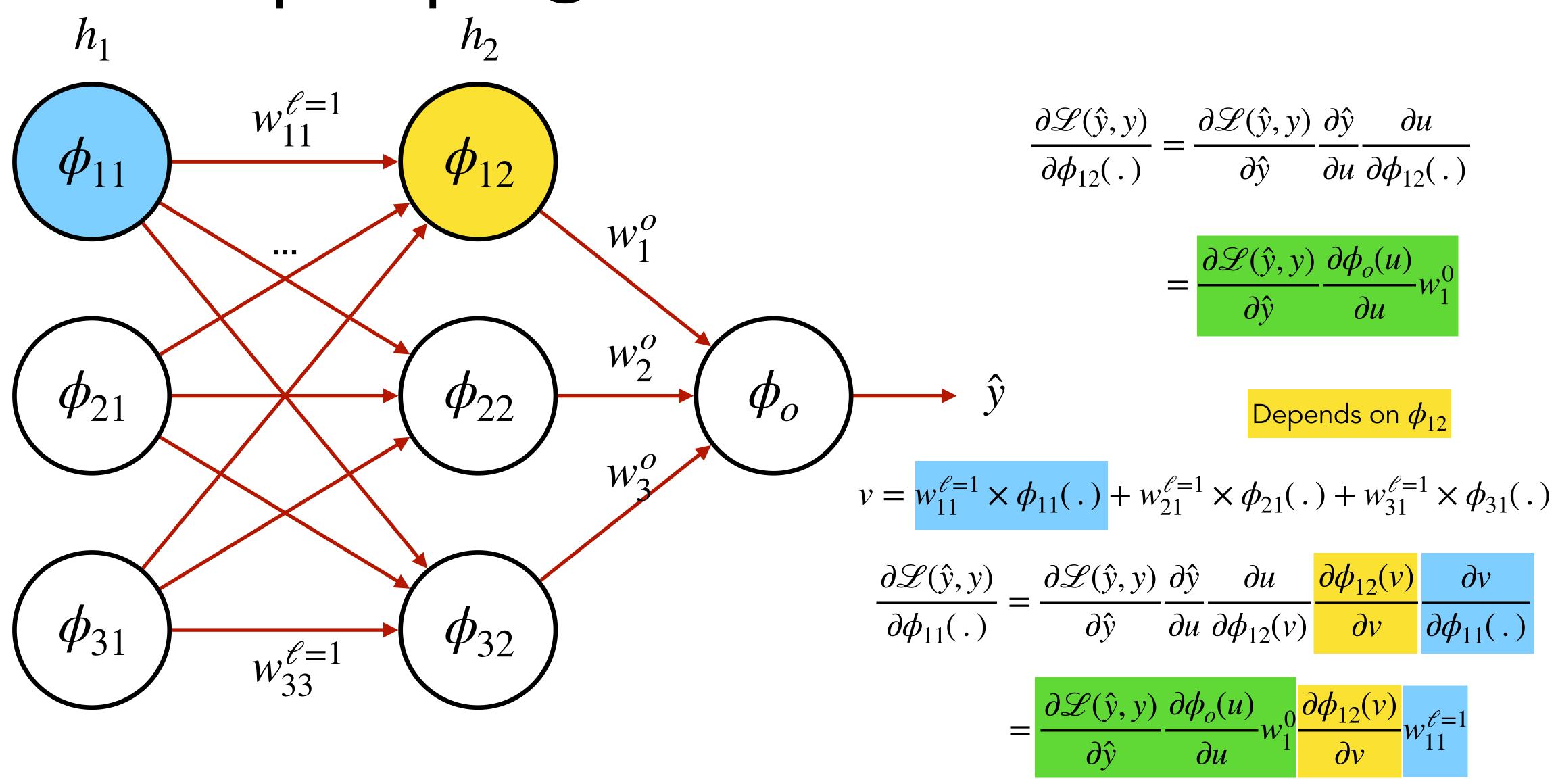
$$= \frac{\partial \mathcal{L}(\hat{y}, y)}{\partial \hat{y}} \frac{\partial \phi_o(u)}{\partial u} w_2^0$$

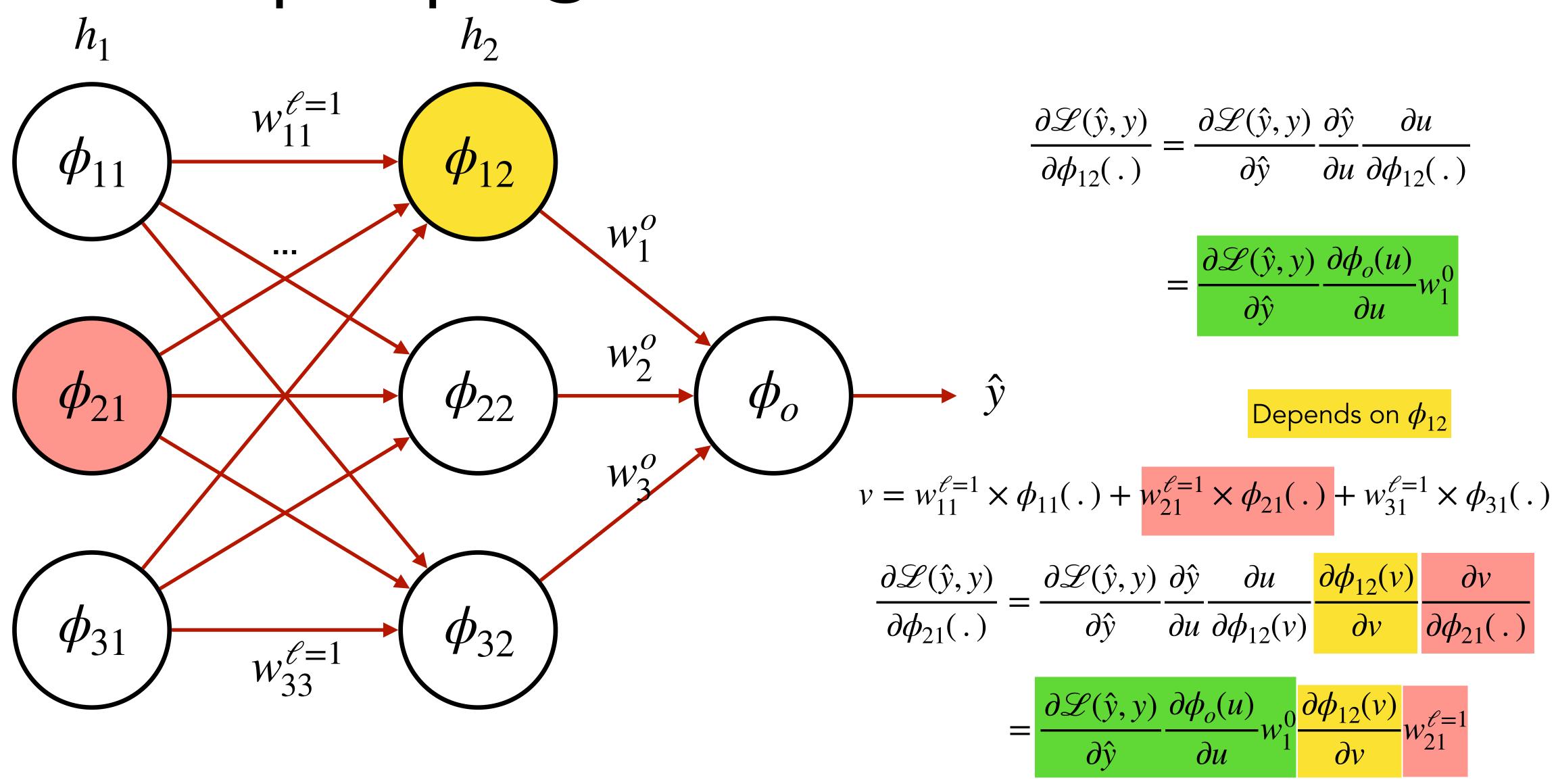
Depends on label y









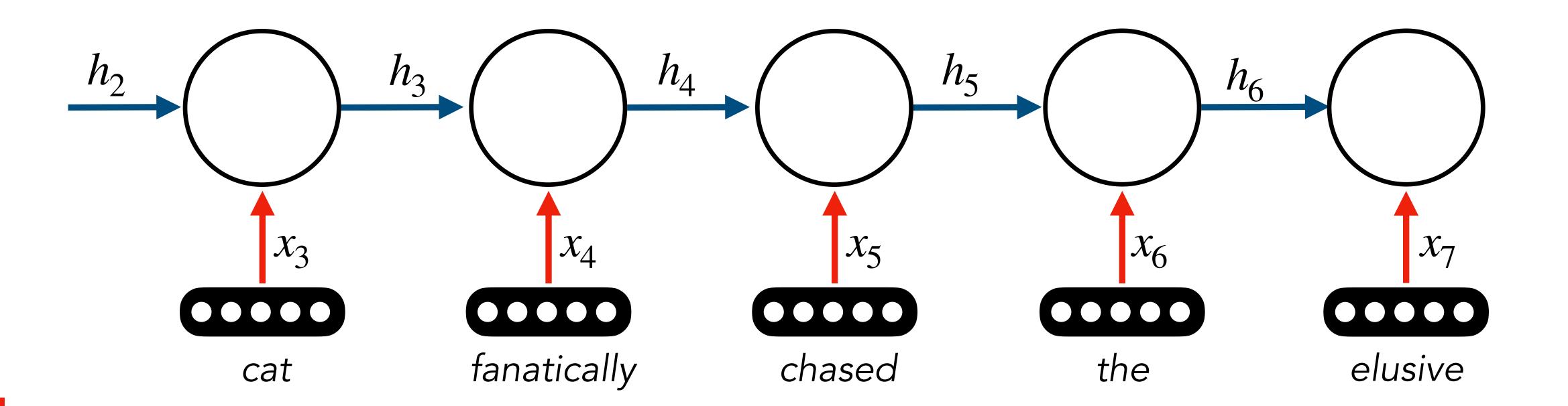


Question

How would we extend backpropagation to a recurrent neural network?

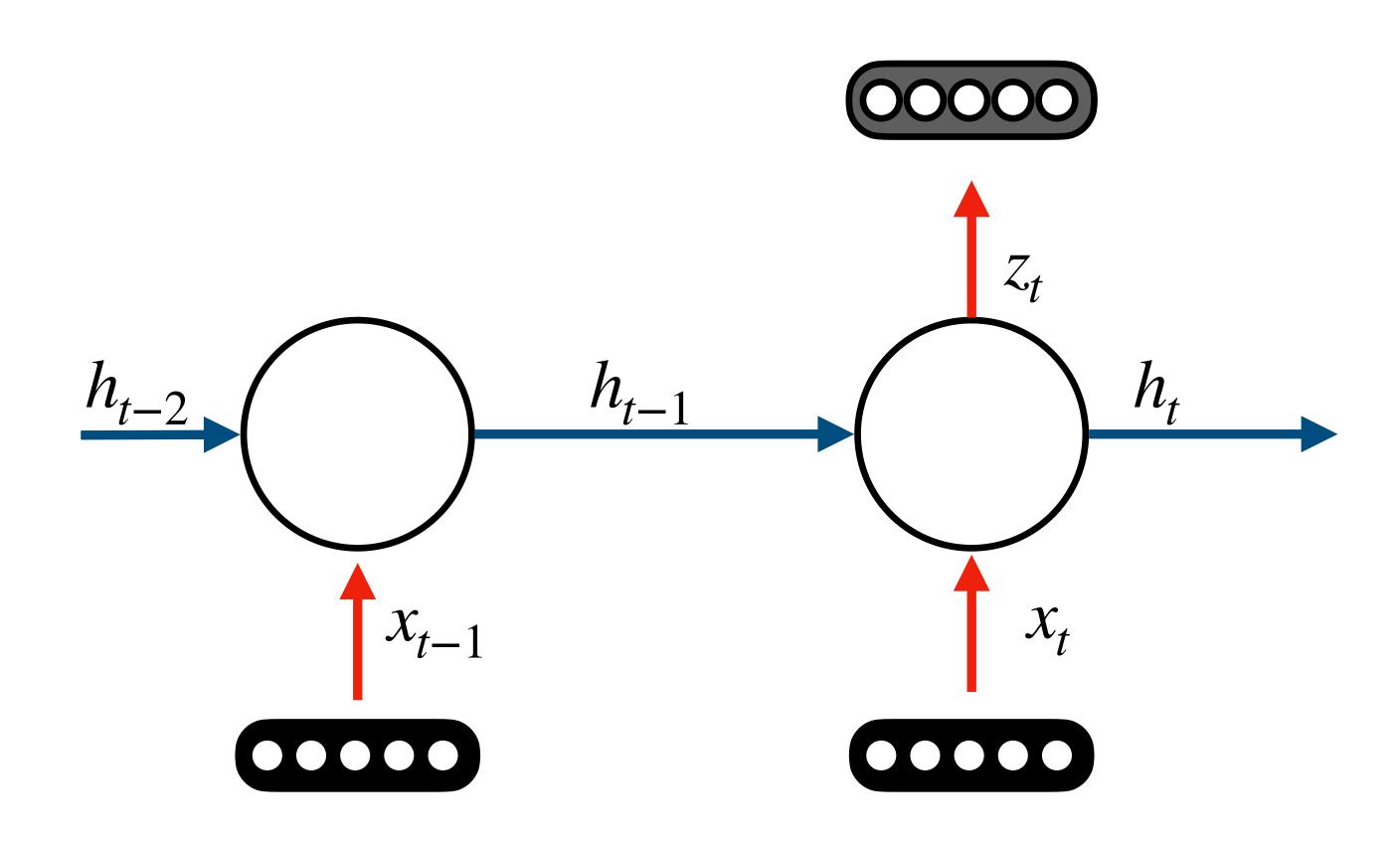
Recall

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



$$z_{t} = \sigma(W_{zh}h_{t} + b_{z})$$

$$h_{t} = \sigma(W_{hx}x_{t} + W_{hh}h_{t-1} + b_{h})$$



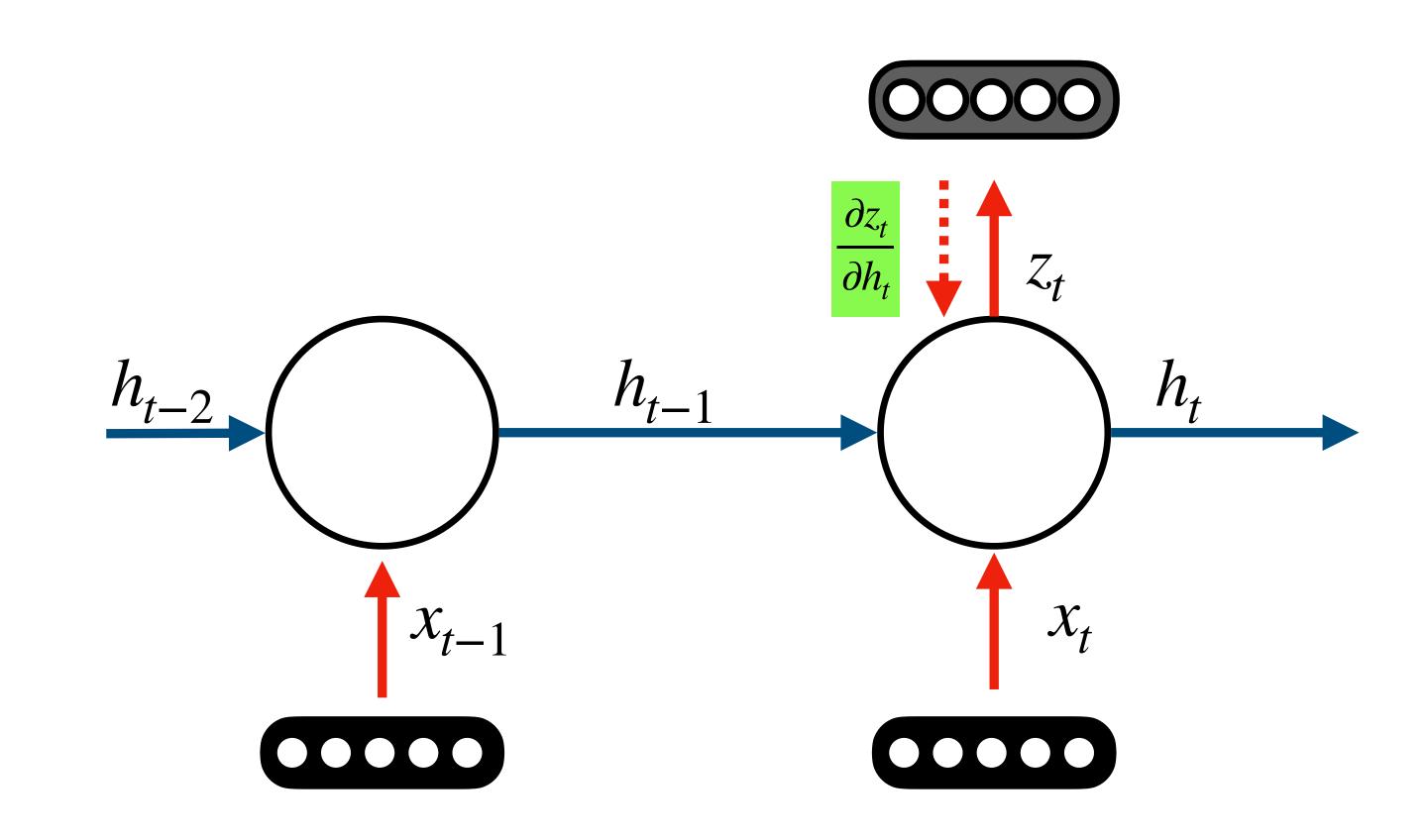
$$z_t = \sigma(W_{zh}h_t + b_z)$$

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

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

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

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



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

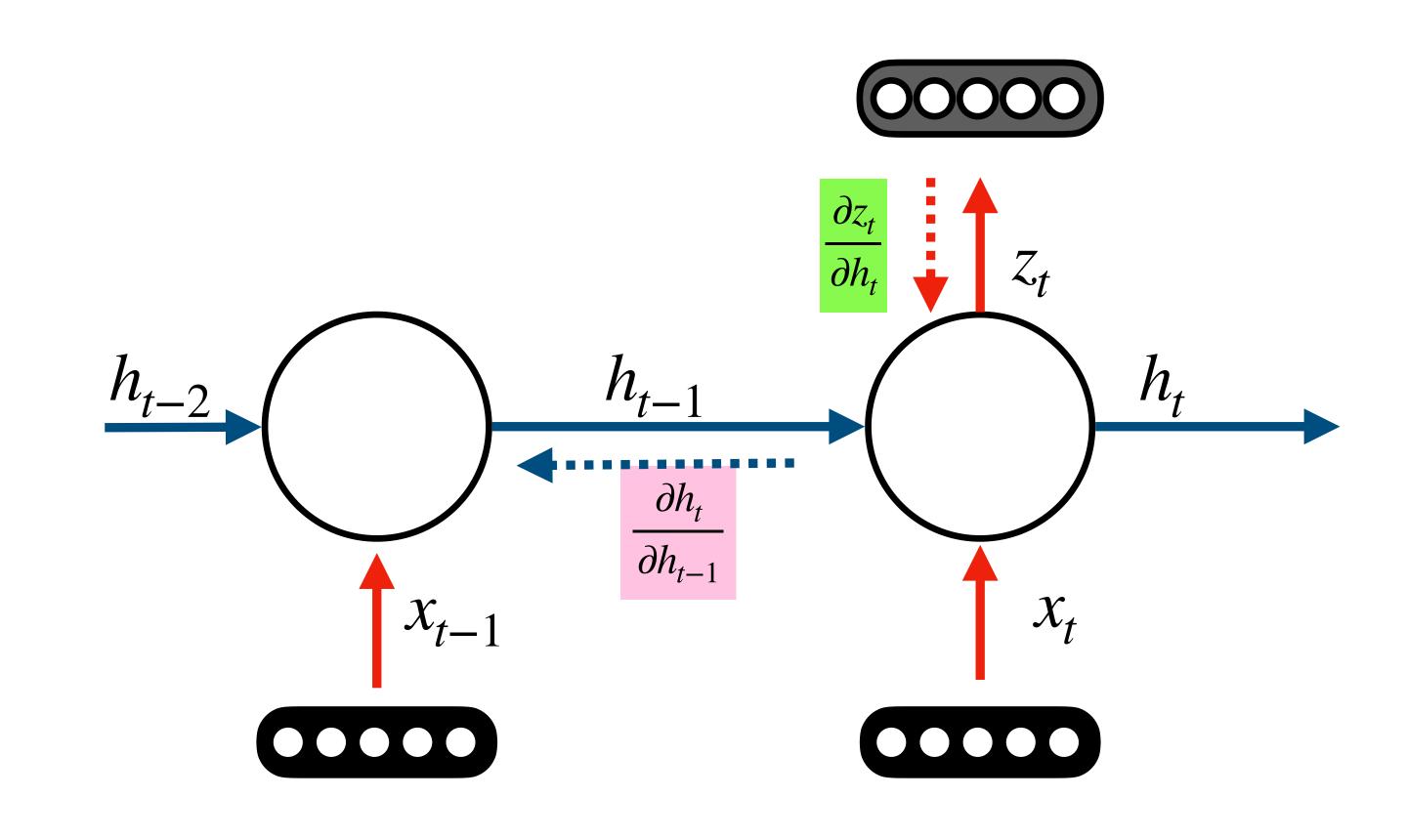
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

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

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

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

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



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

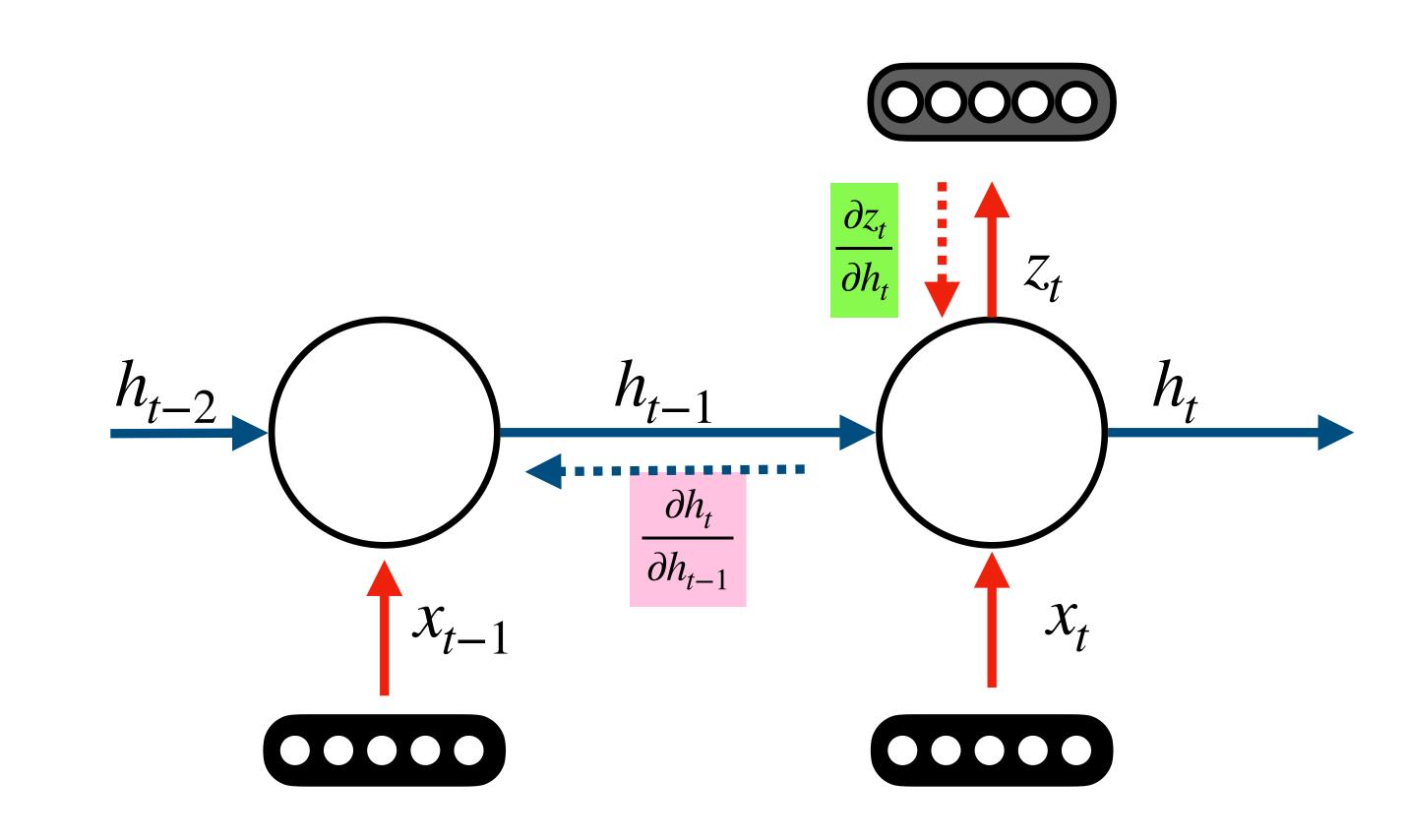
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

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

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

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

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



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

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

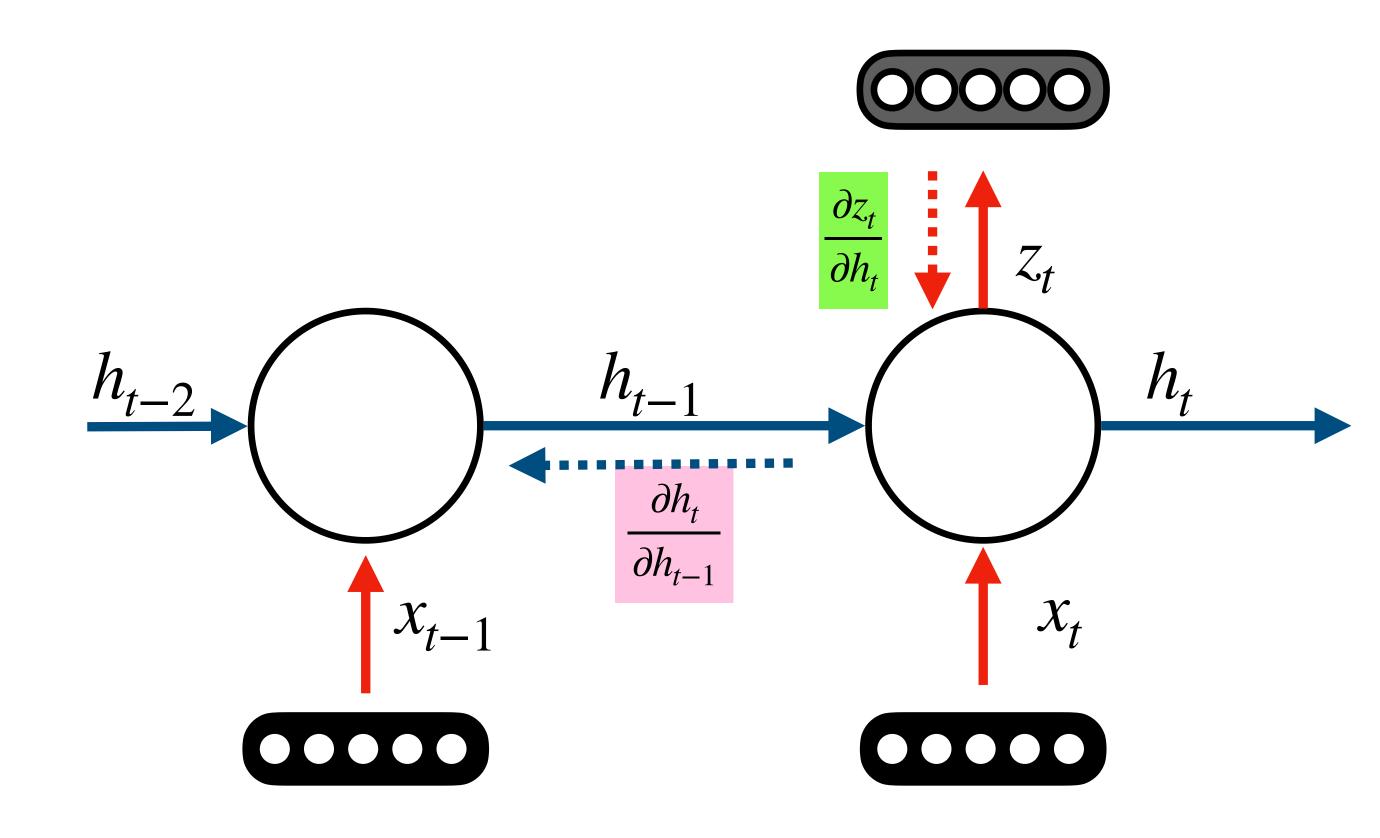
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$

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

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

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

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



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

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

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

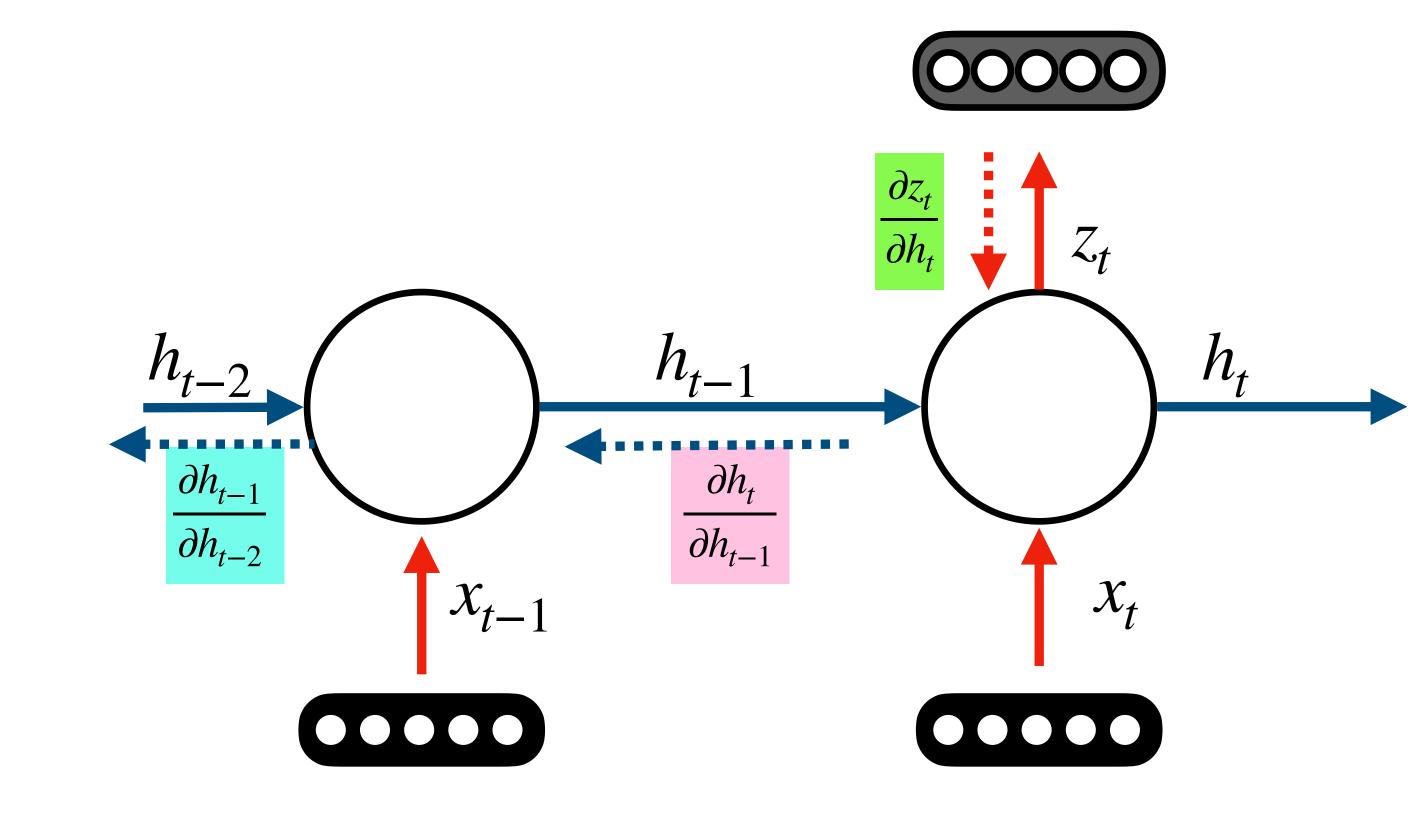
$$v_t = W_{zh}h_t + b_z \qquad z_t = \sigma(v_t)$$

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

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

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

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



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

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

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

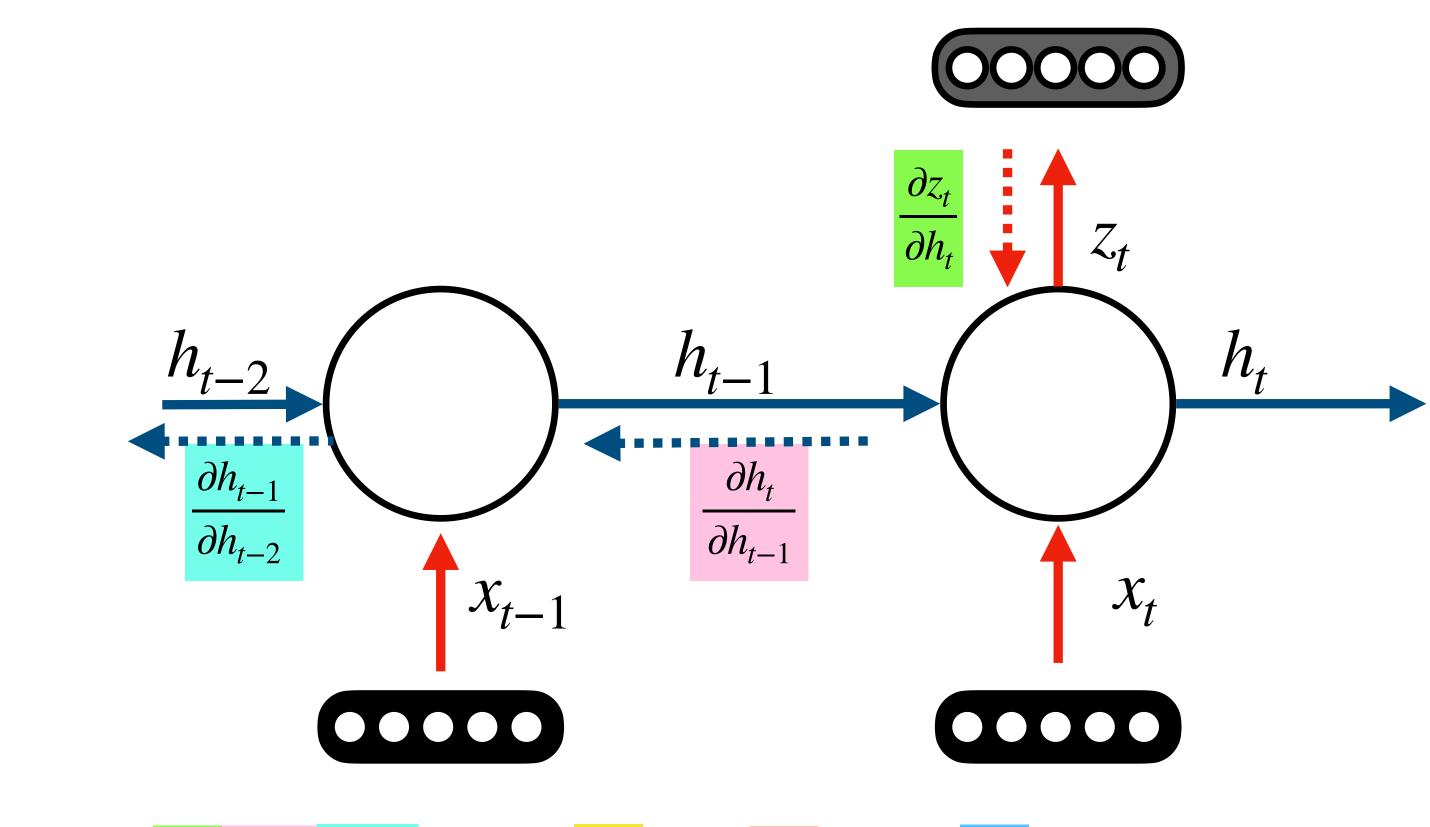
$$v_t = W_{zh}h_t + b_z \qquad z_t = \sigma(v_t)$$

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

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

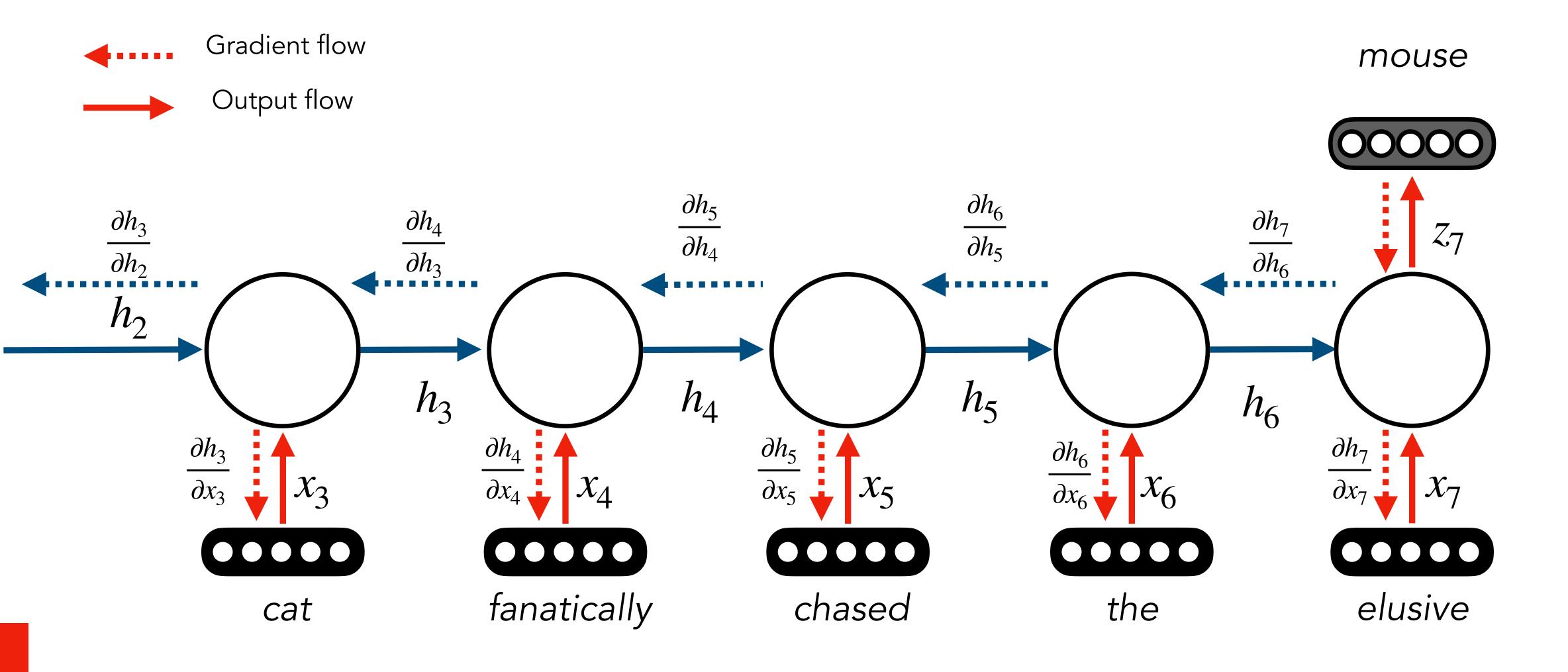
$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \frac{\partial u_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh}$$

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



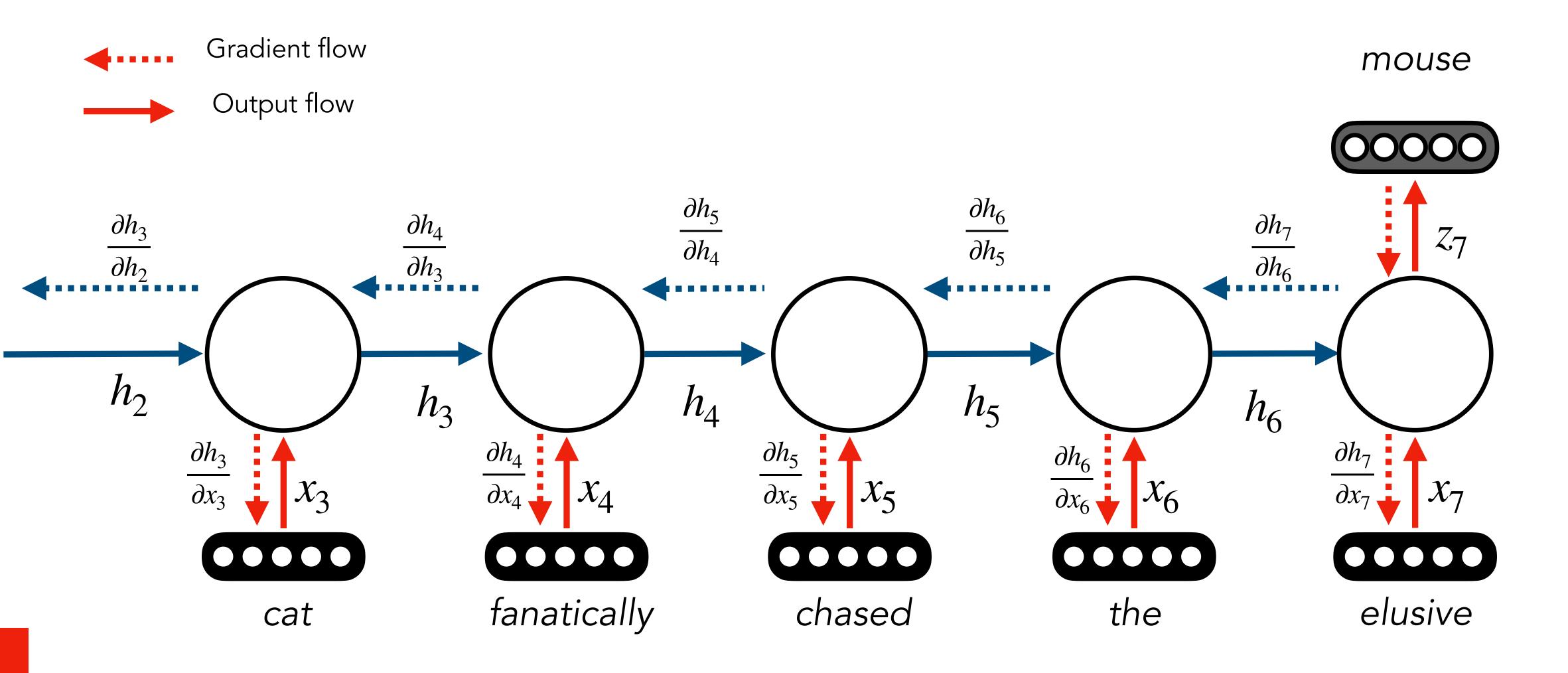
$$\frac{\partial z_{t}}{\partial h_{t-1}} = \frac{\partial z_{t}}{\partial h_{t}} \frac{\partial h_{t}}{\partial h_{t-1}} \frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(v_{t})}{\partial v_{t}} \frac{\partial \sigma(u_{t})}{\partial u_{t}} \frac{\partial \sigma(u_{t})}{\partial u_{t}}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

Note that these are actually the same matrix



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

$$z_{t} = \sigma(W_{zh}h_{t} + b_{z})$$

$$h_{t} = \sigma(W_{hx}x_{t} + W_{hh}h_{t-1} + b_{h})$$

$$v_t = W_{zh}h_t + b_z$$

$$z_t = \sigma(v_t)$$

$$u_t = W_{hx}x_t + W_{hh}h_{t-1} + b_h$$

$$h_t = \sigma(u_t)$$

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

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

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

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

Generalising this:

$$\frac{\partial h_t}{\partial h_{t-T}} = \prod_{i=t-T}^{i=t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t-T}^{i=t} \frac{\partial \sigma(u_i)}{\partial u_i} W_{hh}$$

Vanishing Gradients

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

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

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

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

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

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

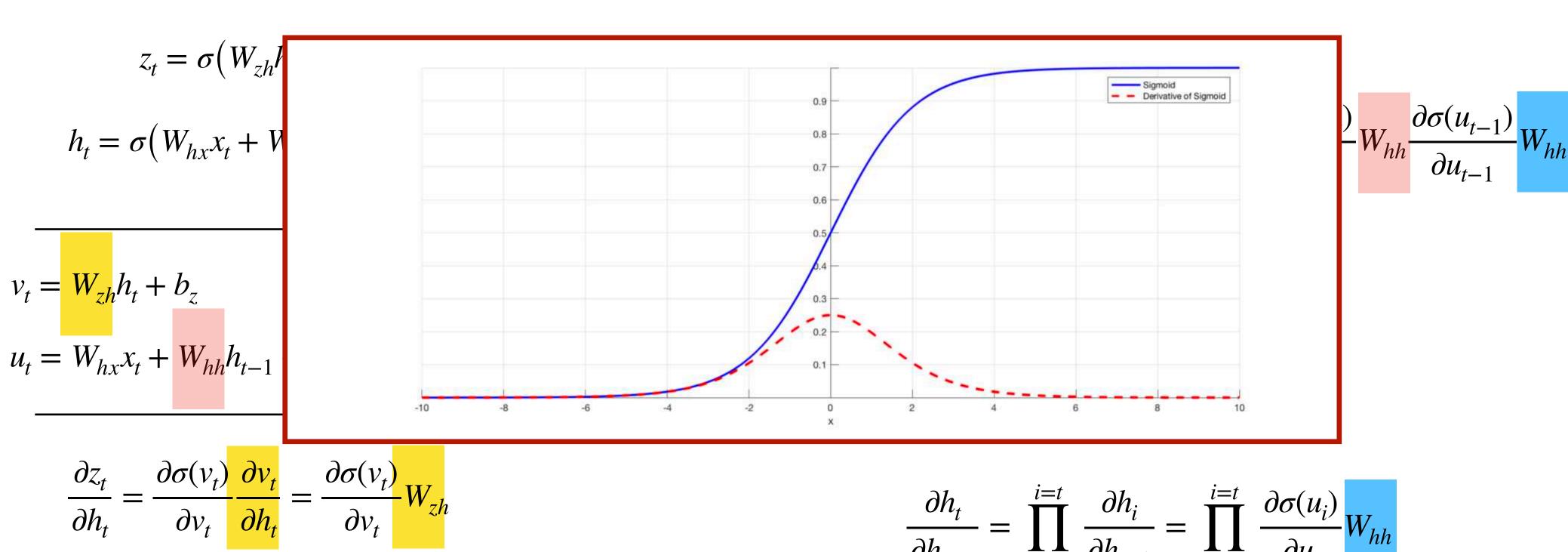
$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \frac{\partial u_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$

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

Generalising this:

$$\frac{\partial h_t}{\partial h_{t-T}} = \prod_{i=t-T}^{i=t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t-T}^{i=t} \frac{\partial \sigma(u_i)}{\partial u_i} W_{hh}$$
< 1 for many
Typically small activation fxns
(Regularisation)

Vanishing Gradients



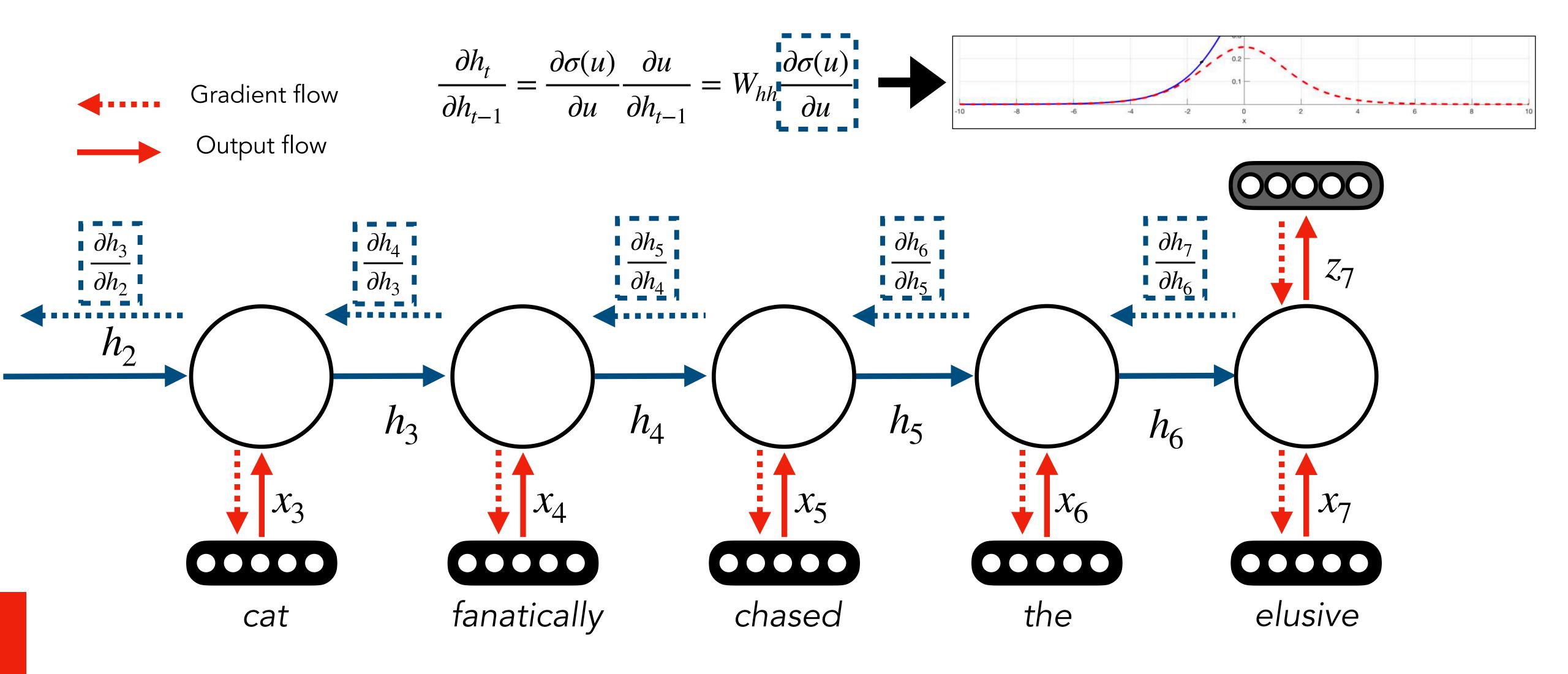
$$\frac{\partial h_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} \frac{\partial u_t}{\partial h_{t-1}} = \frac{\partial \sigma(u_t)}{\partial u_t} W_{hh}$$

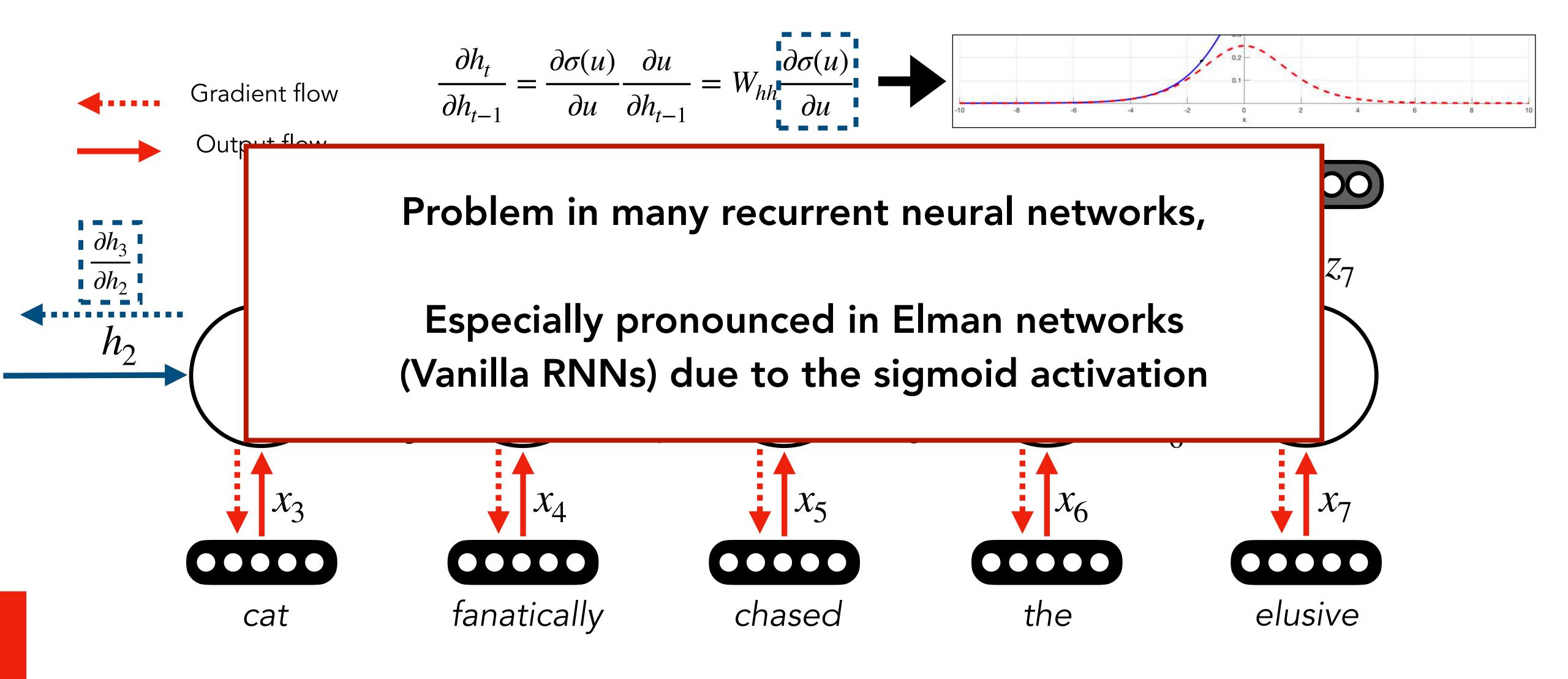
$$\frac{\partial h_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} \frac{\partial u_{t-1}}{\partial h_{t-2}} = \frac{\partial \sigma(u_{t-1})}{\partial u_{t-1}} W_{hh}$$

$$\frac{\partial h_t}{\partial h_{t-T}} = \prod_{i=t-T}^{i=t} \frac{\partial h_i}{\partial h_{i-1}} = \prod_{i=t-T}^{i=t} \frac{\partial \sigma(u_i)}{\partial u_i} W_{hh}$$

< 1 for many activation fxns</p>

Typically small (Regularisation)





Issue with Recurrent Models

 Multiple steps of state overwriting makes it challenging to learn longrange 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 <u>any</u> long-range interactions

Gated Recurrent Neural Networks

Use gates to avoid dampening gradient signal every time step

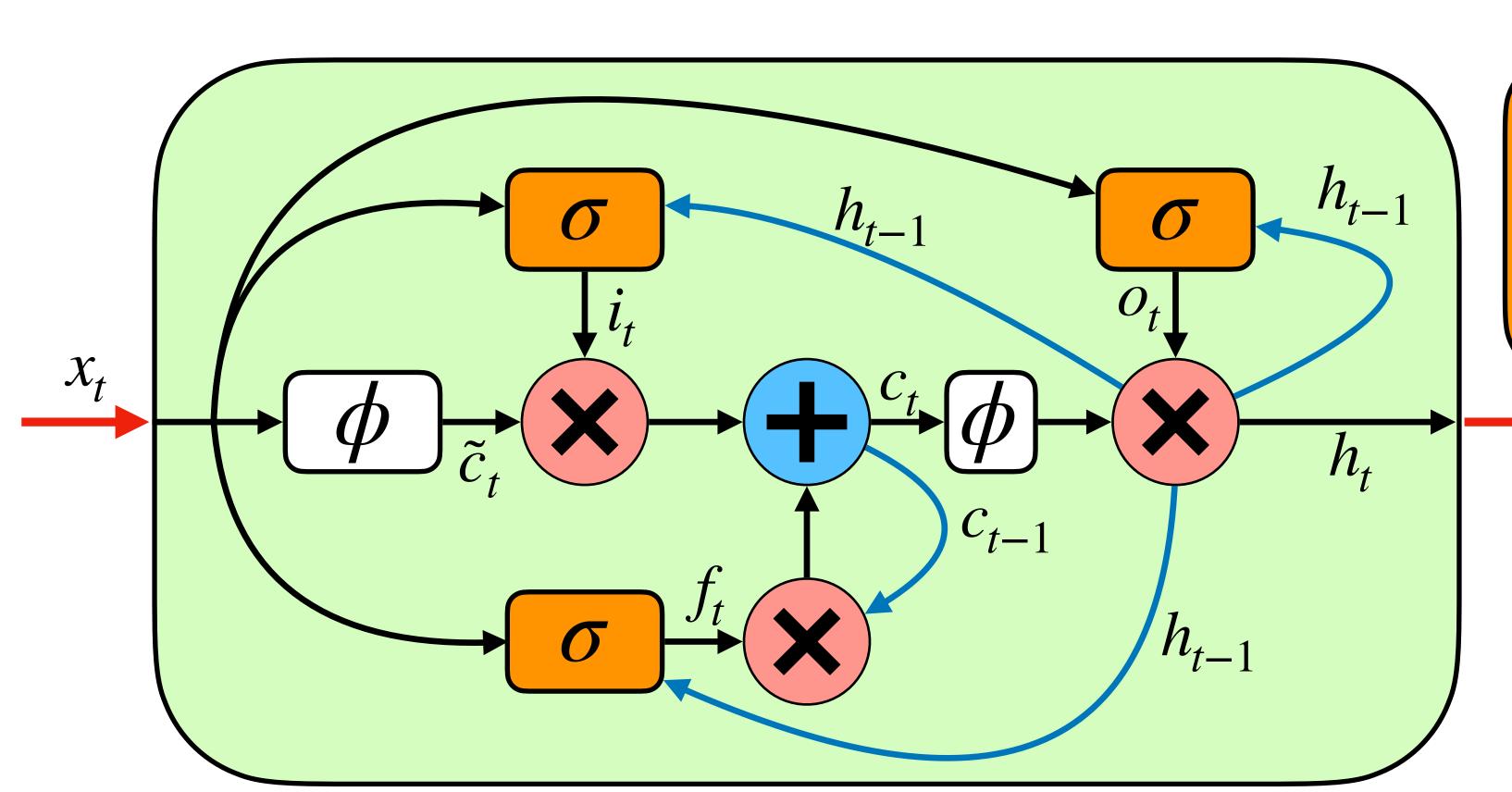
$$h_t = \sigma(W_{hx}x_t + W_{hh}h_{t-1} + b_h)$$
 $h_t = h_{t-1} \odot \mathbf{f} + \mathbf{func}(x_t)$

Elman Network

Gated Network Abstraction

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

Long Short Term Memory (LSTM)



Gates:

$$f_{t} = \sigma(W_{fx}x_{t} + W_{fh}h_{t-1} + b_{f})$$

$$i_{t} = \sigma(W_{ix}x_{t} + W_{ih}h_{t-1} + b_{i})$$

$$o_{t} = \sigma(W_{ox}x_{t} + W_{oh}h_{t-1} + b_{o})$$

$$\tilde{c}_t = \phi \left(W_{cx} x_t + W_{ch} h_{t-1} + b_c \right)$$

$$c_t = i_t \times \tilde{c}_t + f_t \times c_{t-1}$$

$$h_t = o_t \times \phi(c_t)$$

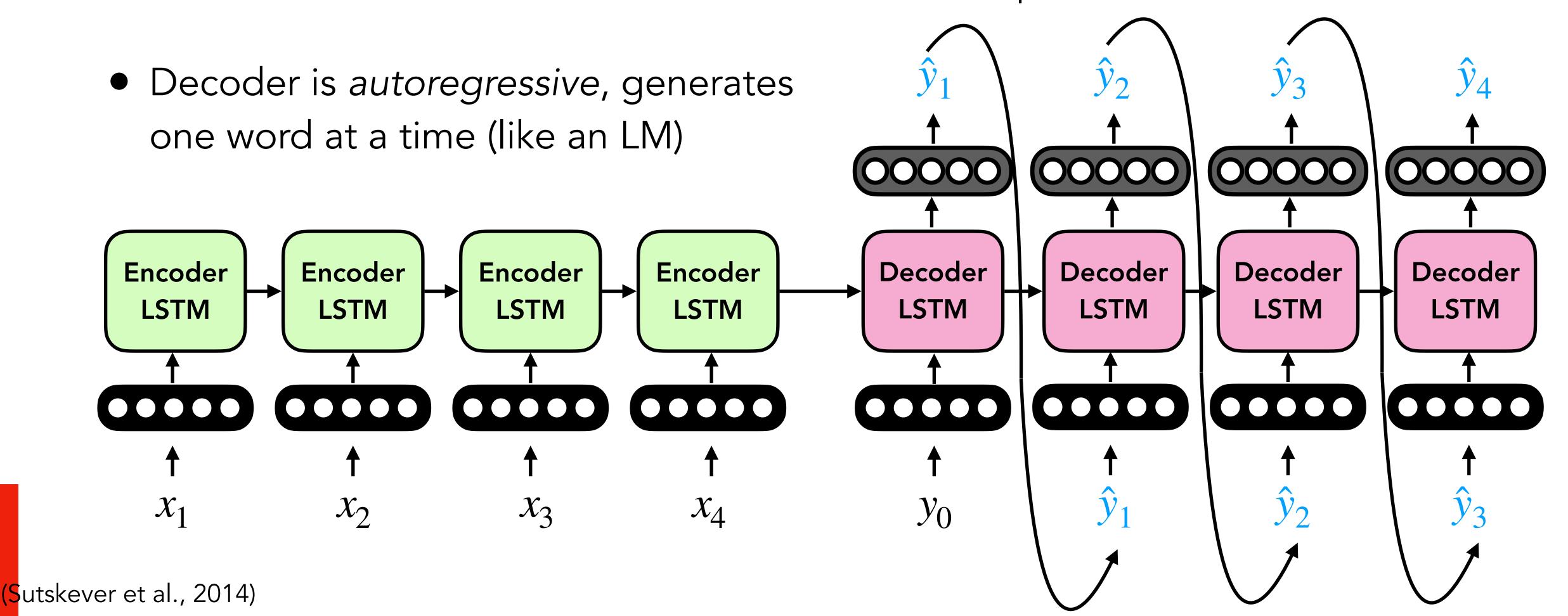
Question

How can we use recurrent neural networks in practice?

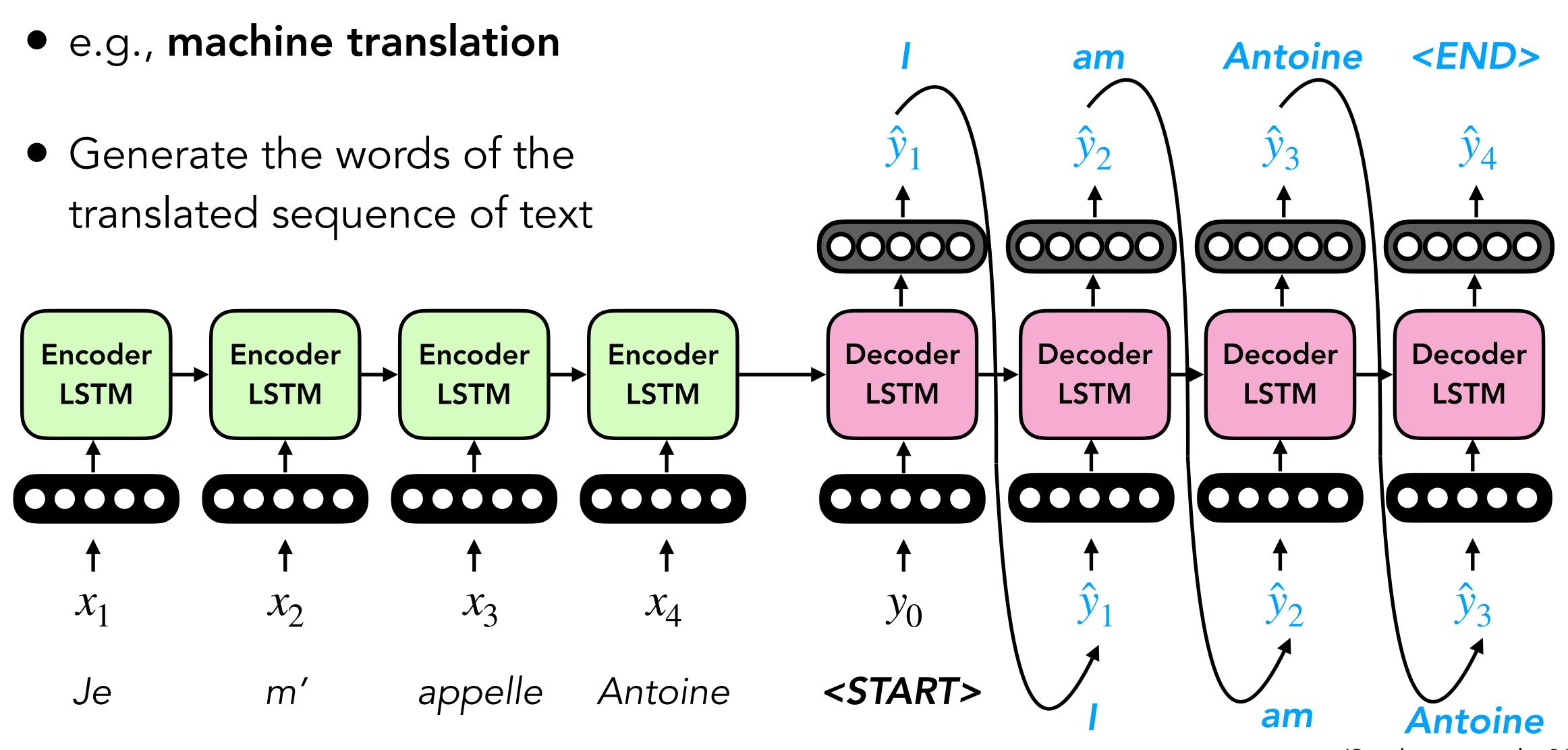
Machine Translation involves more than estimating the probability next word; requires generating a full translation of a given context into another language

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**)



Encoder-Decoder Models

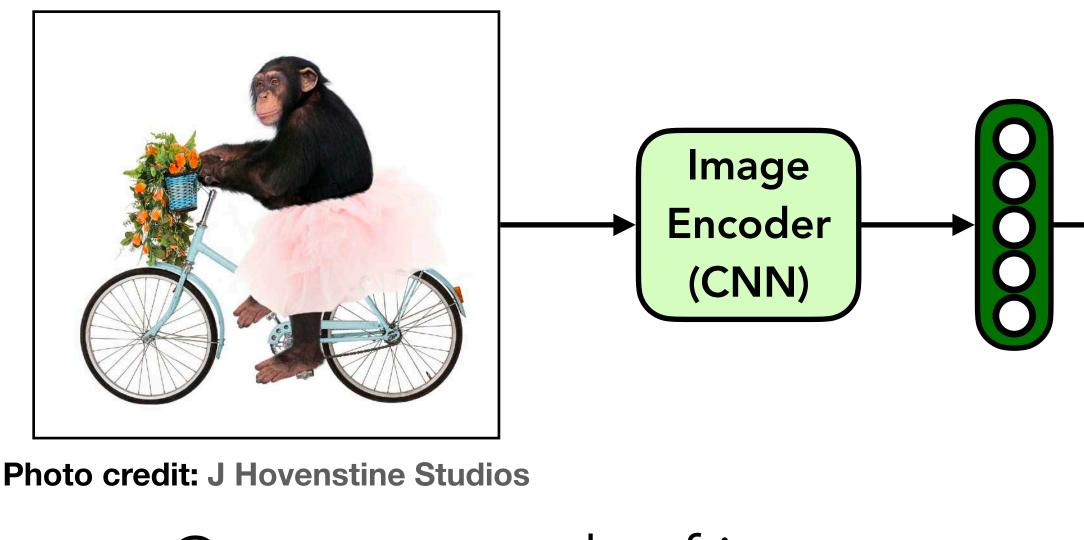


(Sutskever et al., 2014)

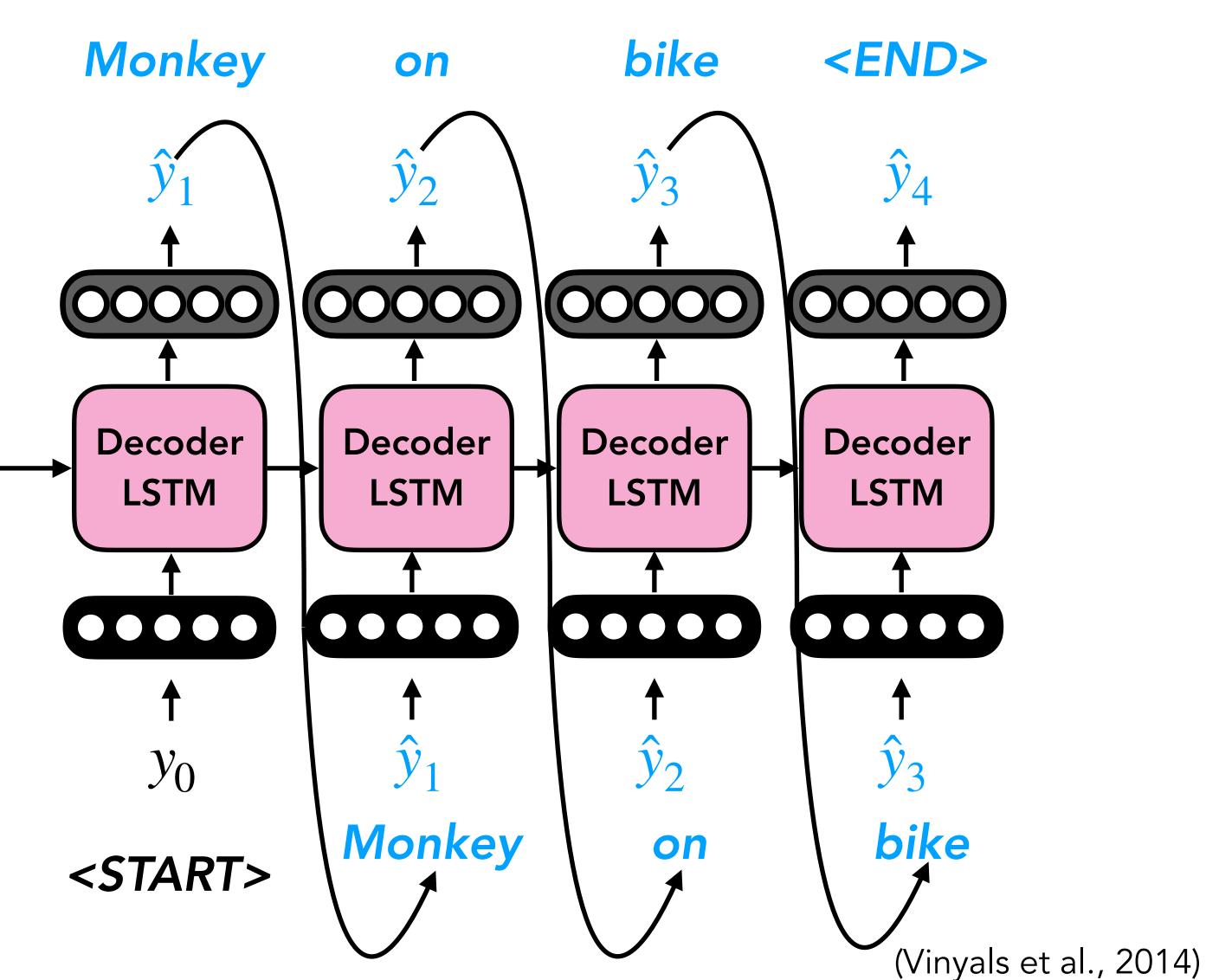
Encoder-Decoder Models

Input doesn't need to be text

• e.g., image captioning

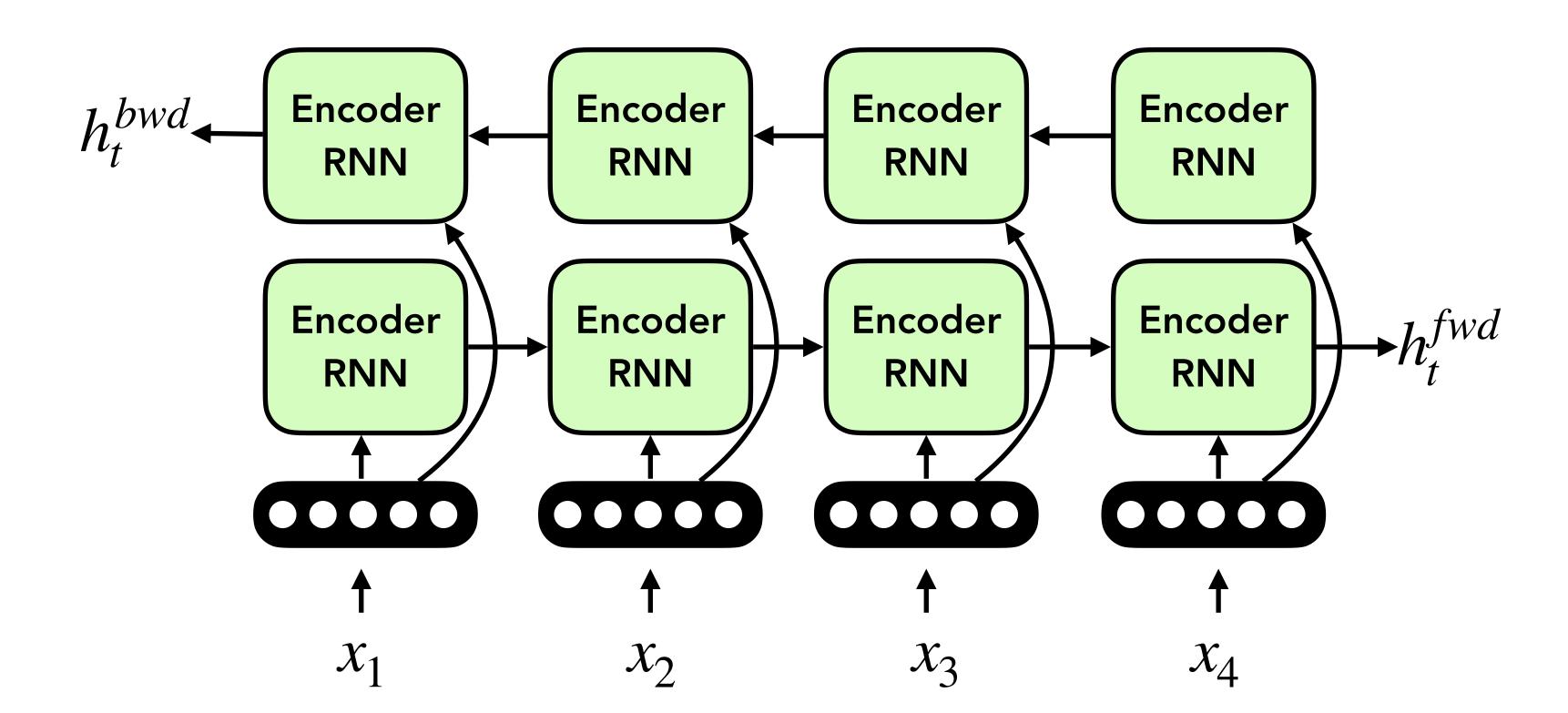


Generate words of image description



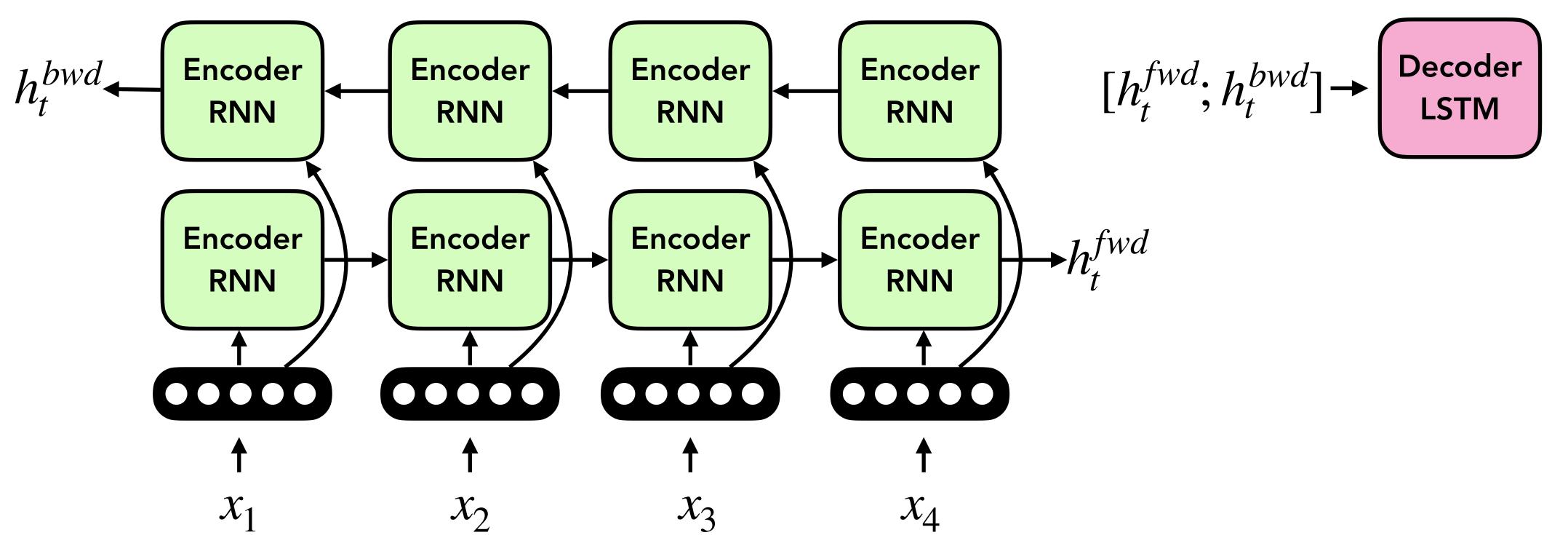
Bidirectional Encoders

- 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-tosequence 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., Merrienboer, 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.

Deep Learning for Natural Language Processing

Antoine Bosselut





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

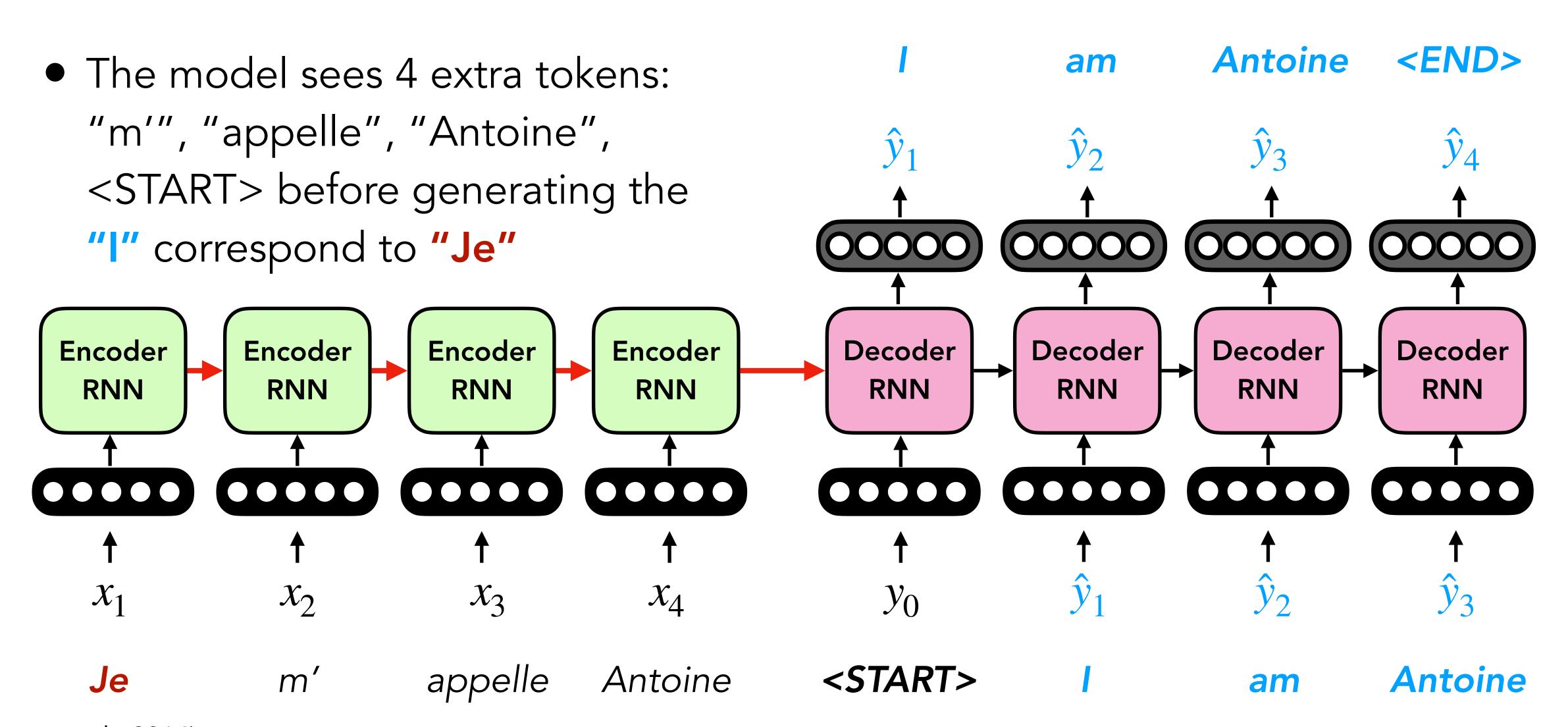
Issue with Recurrent Models

 Multiple steps of state overwriting makes it challenging to learn longrange dependencies.

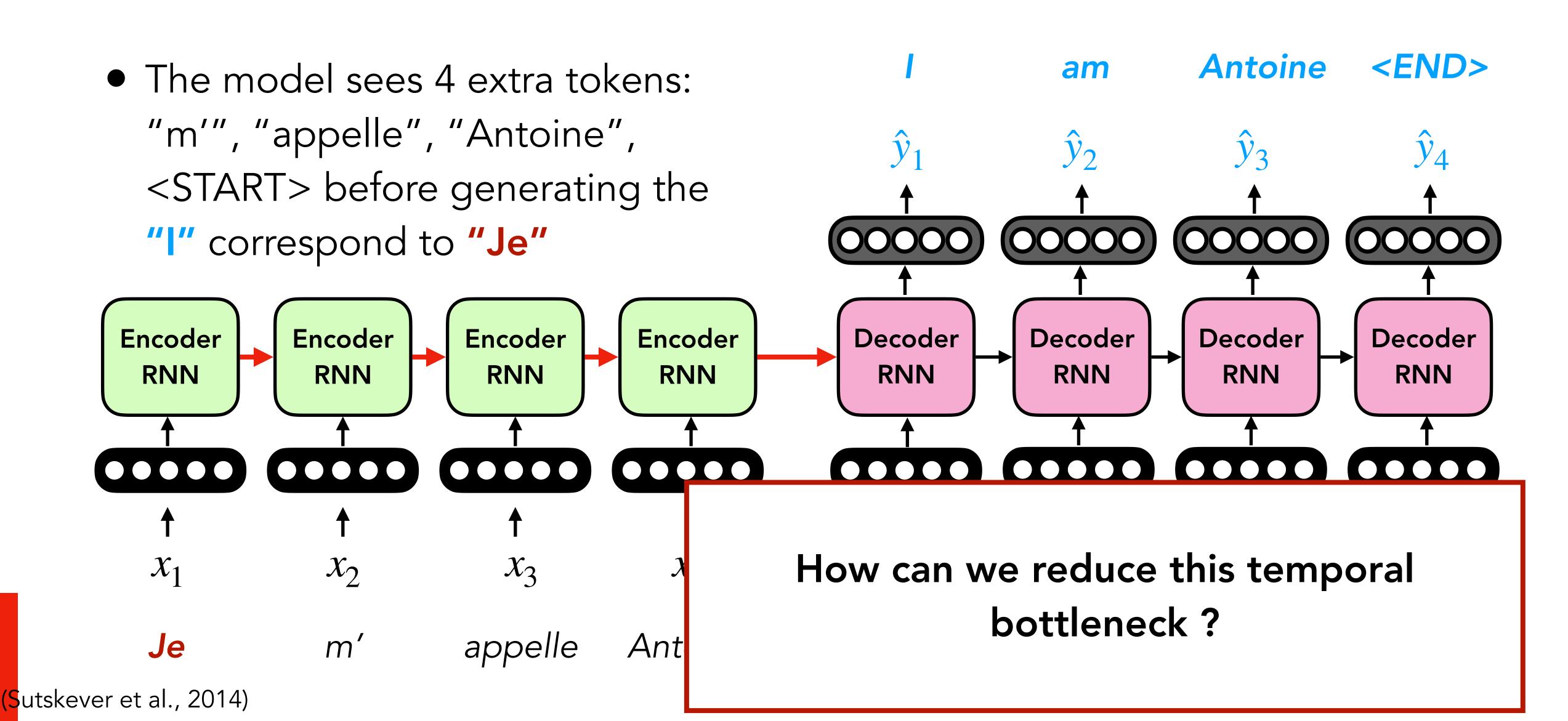
They tuned, discussed for a moment, then struck up a lively jig. Everyone joined in, turning the courtyard into an even more chaotic scene, people now dancing in circles, swinging and spinning in circles, everyone making up their own dance steps. I felt my feet tapping, my body wanting to move. Aside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

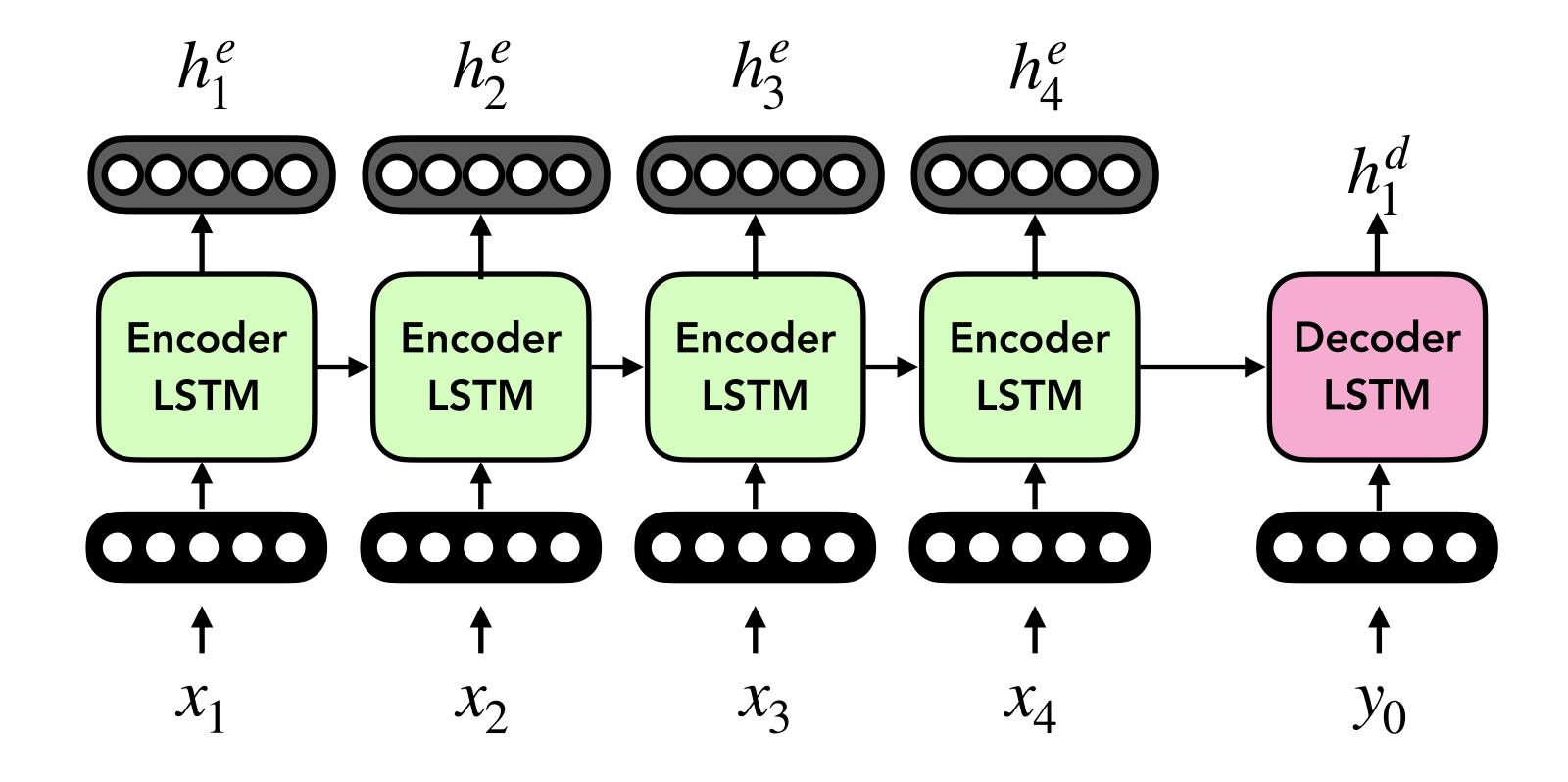
 Nearby words should affect each other more than farther ones, but RNNs make it challenging to learn <u>any</u> long-range interactions

Toy Example

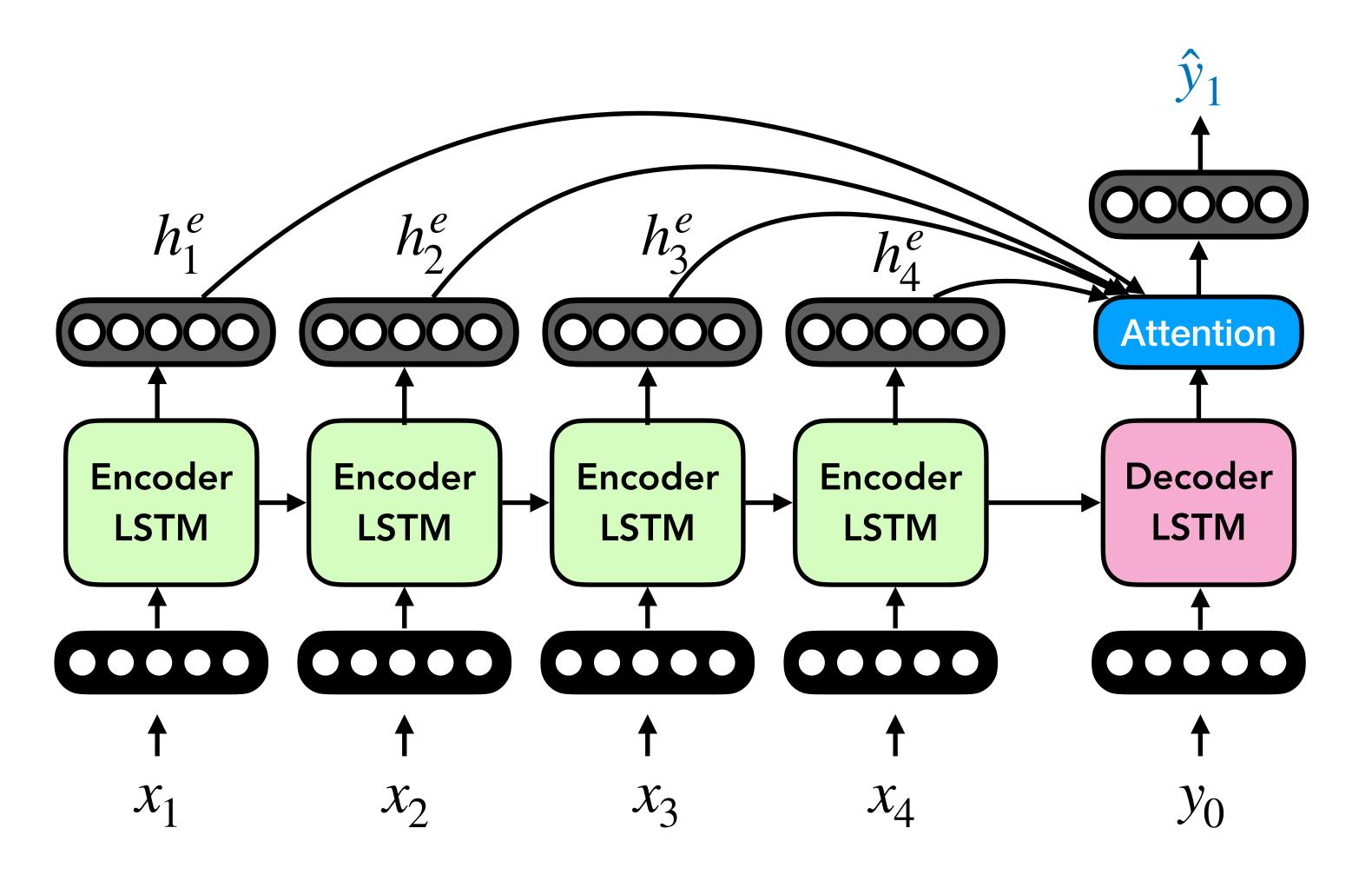


Toy Example





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



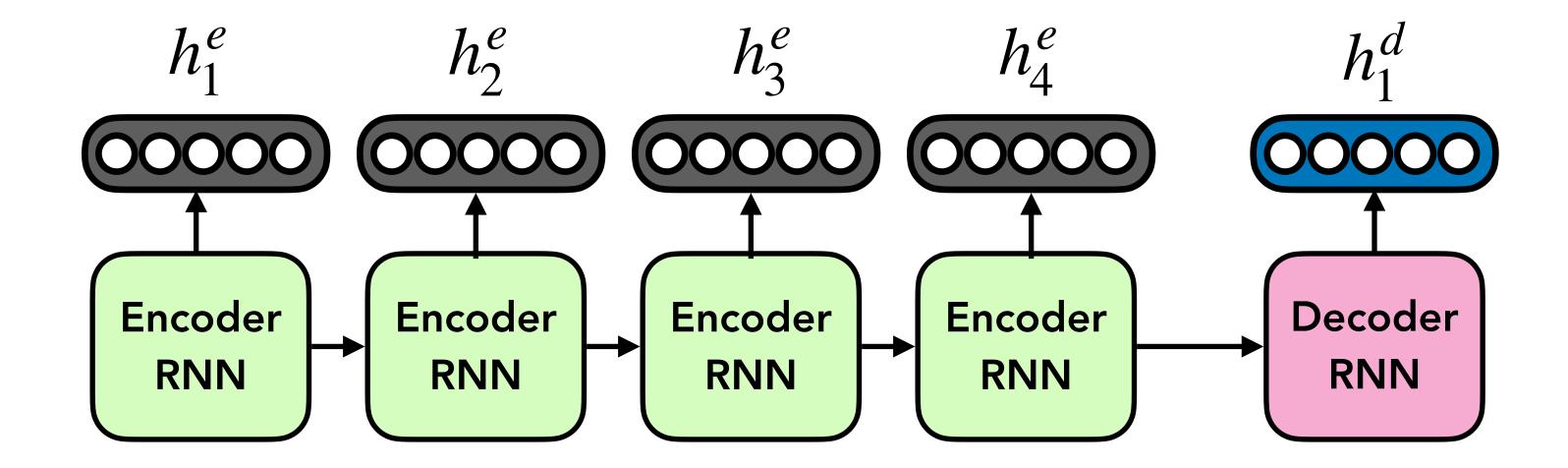
- Recall: At each encoder time step, there is an output of the RNN!
- Idea: Use the output of the Decoder LSTM to compute an attention (i.e., a mixture) over all the h_t^e outputs of the encoder LSTM
- Intuition: focus on different parts of the input at each time step

What is attention?

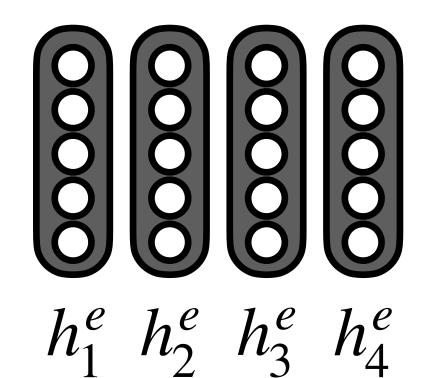
• Attention is a weighted average over a set of inputs

$$h_t^e$$
 = encoder output hidden states

How should we compute this weighted average?



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

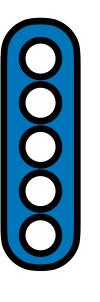


 h_t^e = encoder output hidden states

Also known as a "keys"

 h_t^d = decoder output hidden state

Also known as a "query"



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

$$h_t^e$$
 = encoder output hidden states

 h_t^d = decoder output hidden state

Also known as a "keys"

Also known as a "query"

$$a_1 = f(0), a_2 = f(0), a_3 = f(0), a_4 = f(0), a_4$$

We have a single query vector for multiple key vectors

	•		•
Atter	ition	Fund	ction

Formula

Multiplicative

Linear

$$a = h^e W h^d$$

$$a = v^T \phi(\mathbf{W}[h^e; h^d])$$

Scaled Dot Product

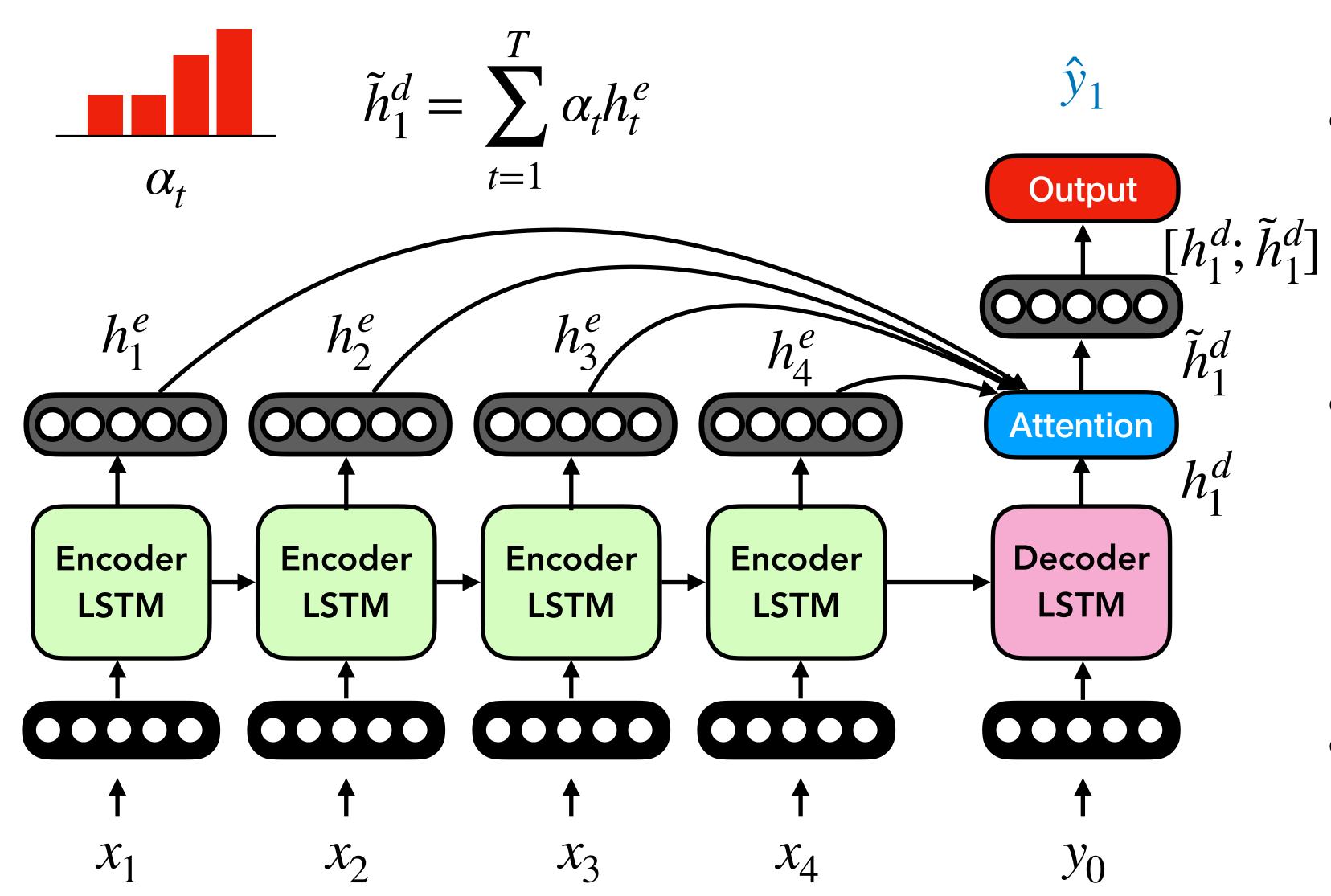
$$a = \frac{(\mathbf{W}h^e)^T(\mathbf{U}h^d)}{\sqrt{d}}$$

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

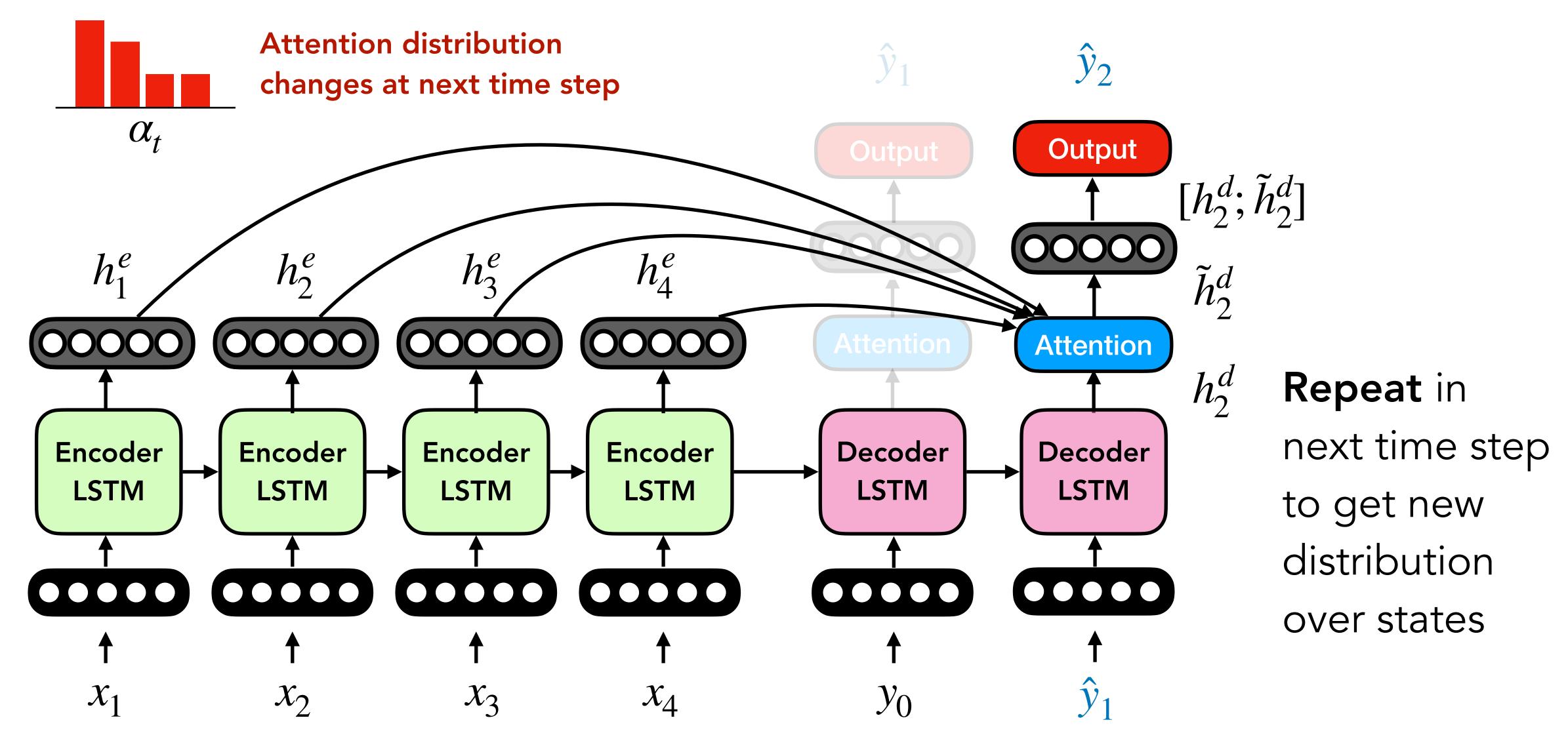
$$a_1 = f(0), 0 \qquad a_2 = f(0), 0 \qquad a_3 = f(0), 0 \qquad a_3 = f(0), 0 \qquad a_4 = f(0), 0 \qquad a_4 = f(0), 0 \qquad a_5 = f(0), 0 \qquad a_6 = f(0),$$

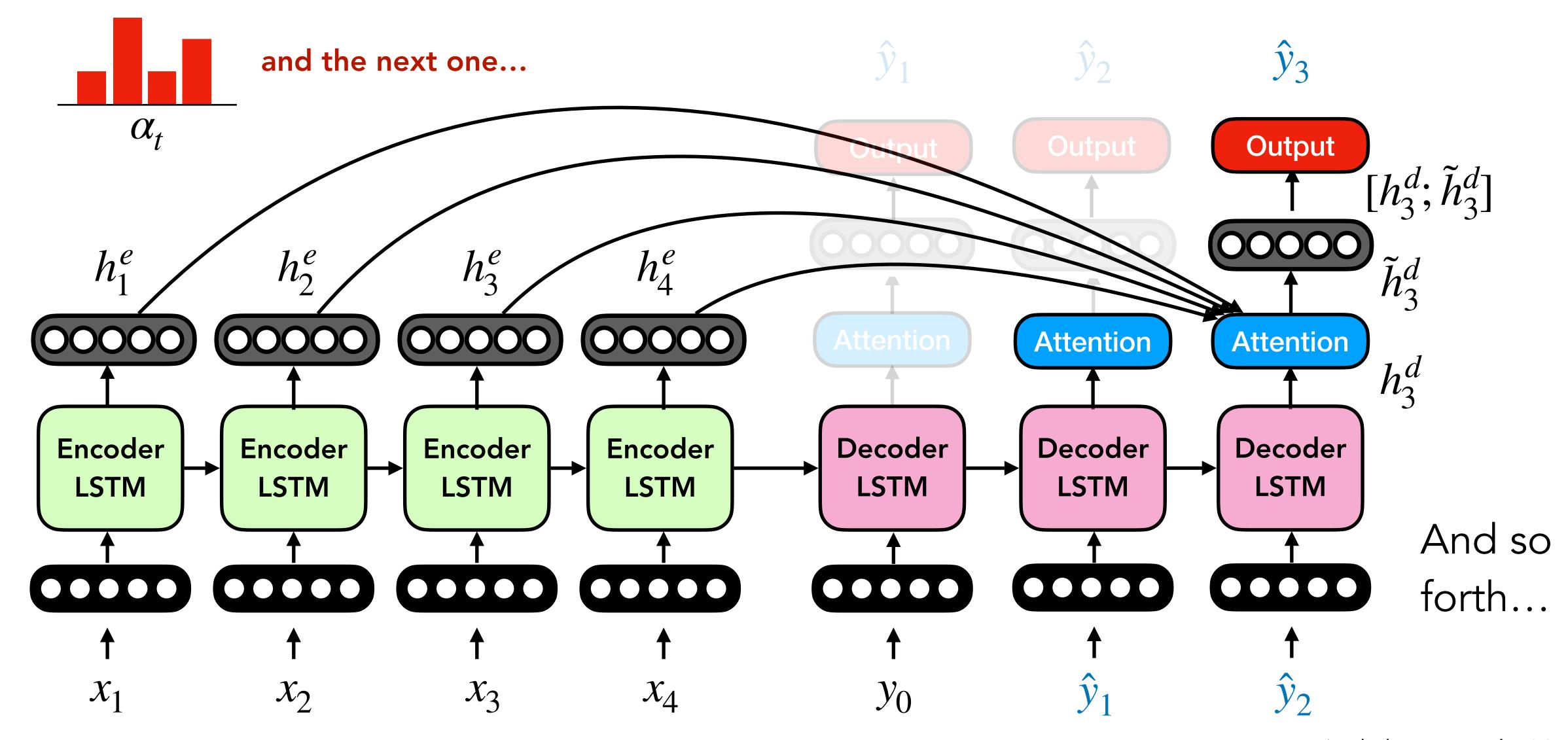
• Convert pairwise similarity scores to probability distribution (using softmax!) over encoder hidden states and compute weighted average:

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



- 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





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

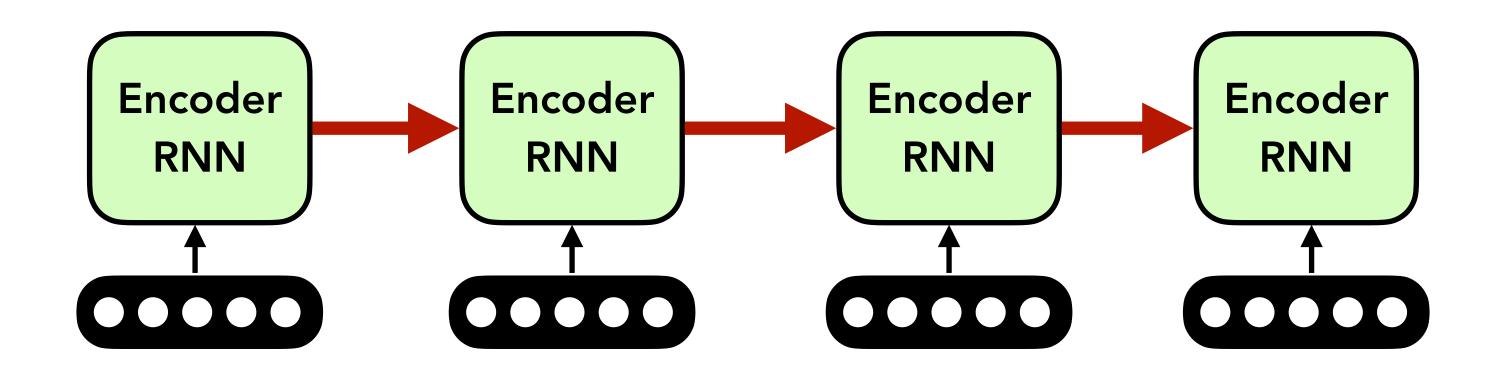
- Many possible functions for computing scores (dot product, bilinear, etc.)
- Temporal Bottleneck Fixed! Direct connection between decoder and ALL encoder states

Question

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

Encoder is still Recurrent

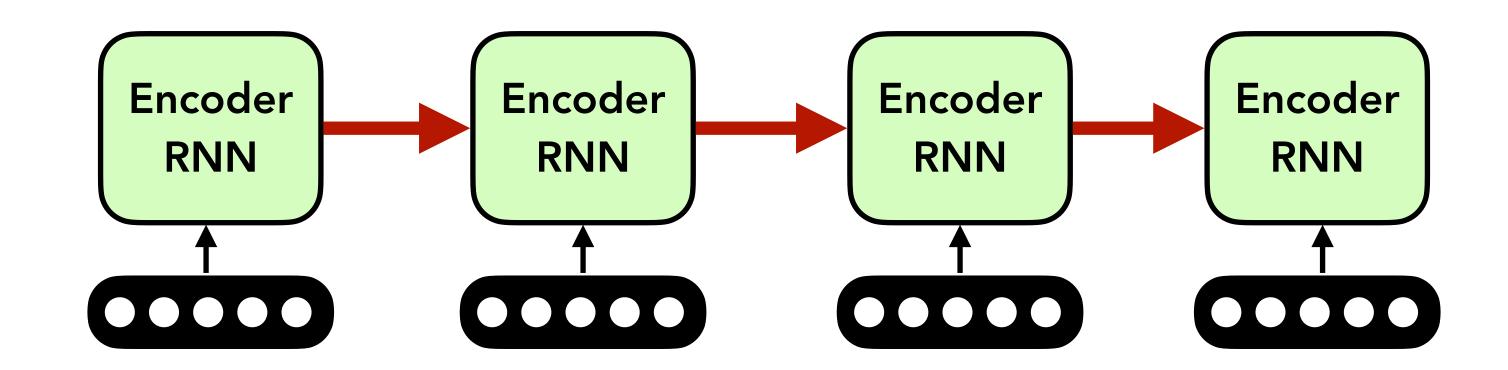
• **Encoder:** Recurrent functions can't be parallelized because previous state needs to be computed to encode next one



• 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

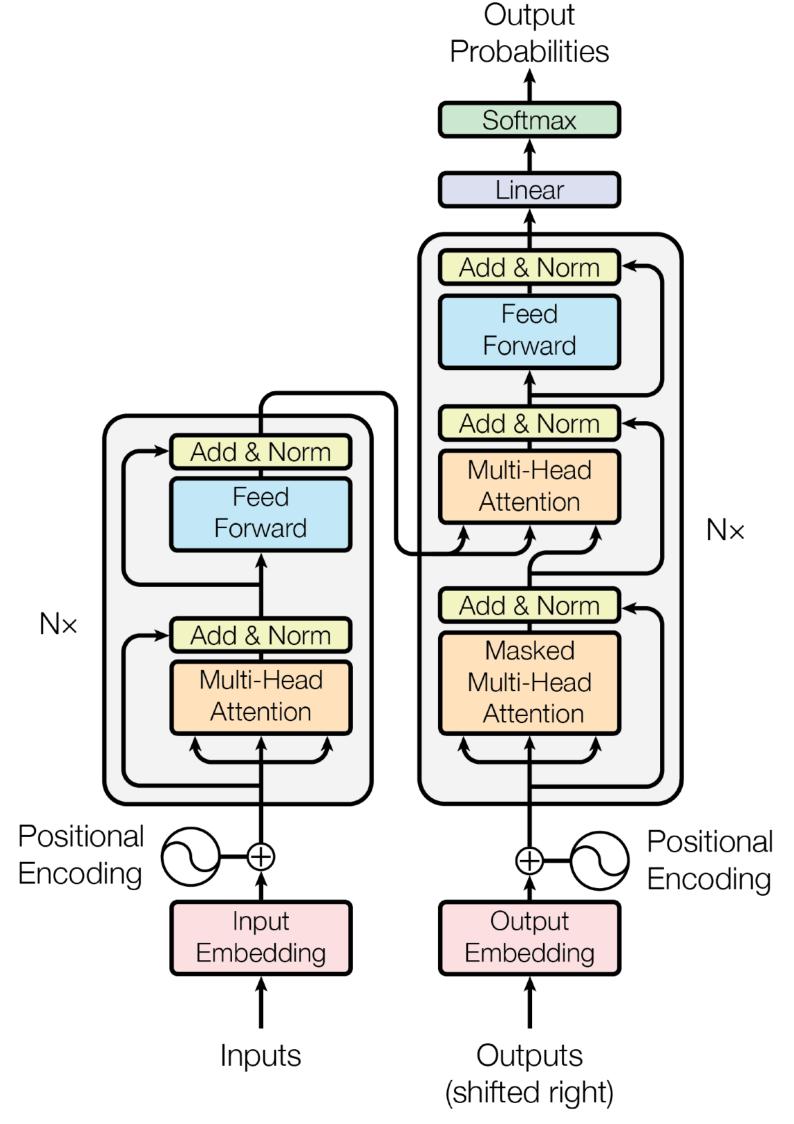
Who can think of a task where this might be a problem?

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!

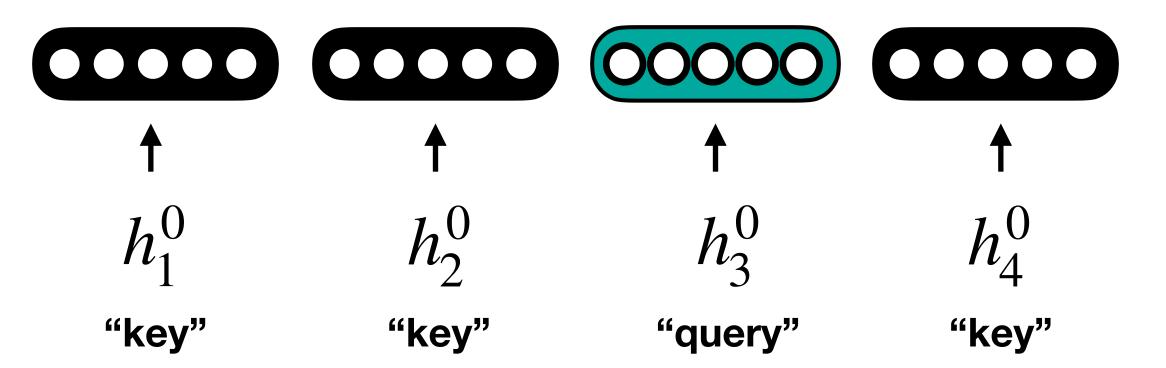
Encode sequences with self-attention



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!

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



Recap: Attention with RNNs

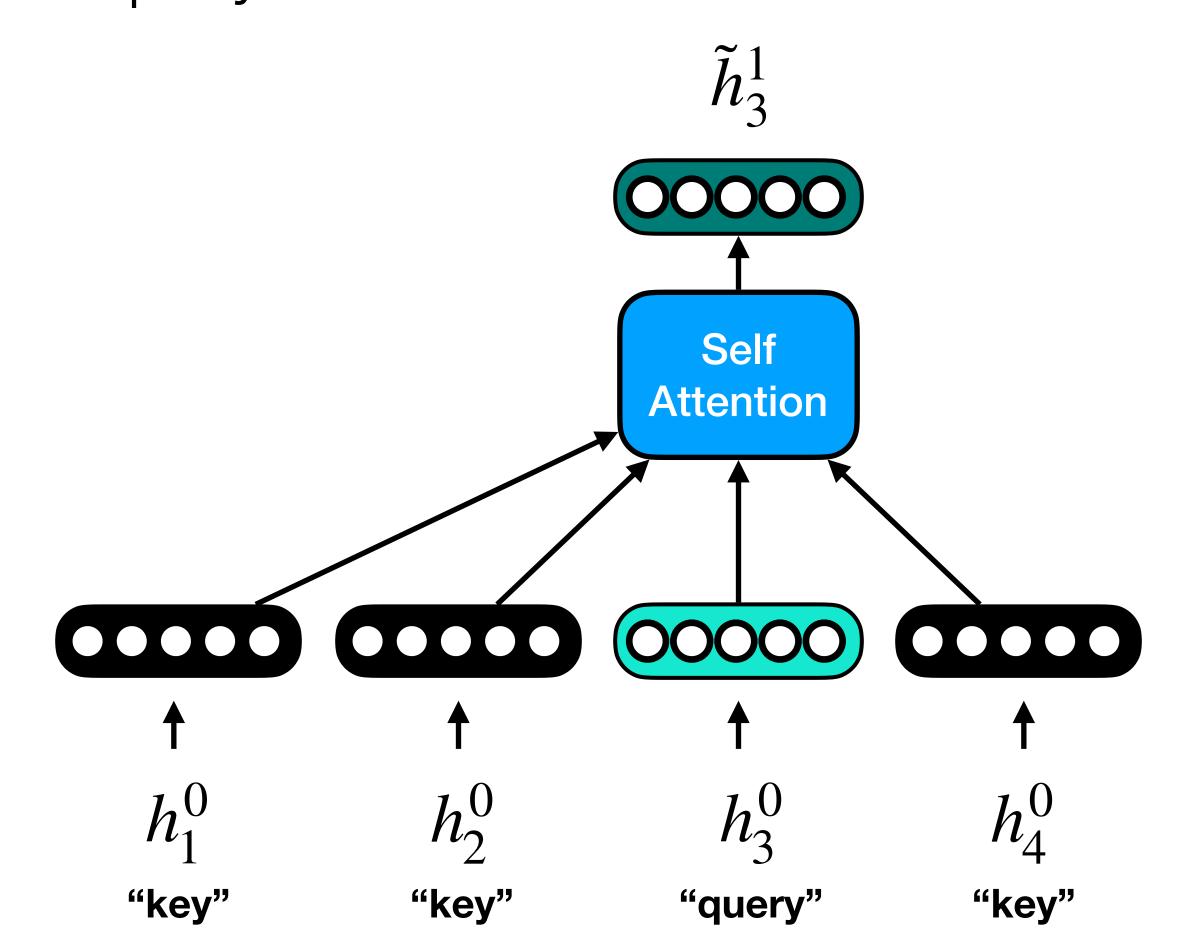
• Compute pairwise similarity between each encoder hidden state and decoder hidden state ("idea of what to decode")

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

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

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

$$a_{31} = f(\mathbf{B}, \mathbf{B}) \rightarrow a_{st} = f(\mathbf{B}, \mathbf{B})$$
 $h_1^0 h_3^0 h_3^0 h_t^\ell h_s^\ell h$

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

Compute pairwise scores

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

Get attention distribution

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

Attend to values to get weighted sum

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

$$a_{31} = f(\mathbf{B}, \mathbf{B}) \rightarrow a_{st} = f(\mathbf{B}, \mathbf{B})$$
 $h_1^0 h_3^0 h_s^0$

"key" "query"

$$a_{st} = \frac{(\mathbf{W}^{Q} \mathbf{h}_{s}^{\ell})^{T} (\mathbf{W}^{K} \mathbf{h}_{t}^{\ell})}{\sqrt{d}}$$

Compute pairwise scores

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

Get attention distribution

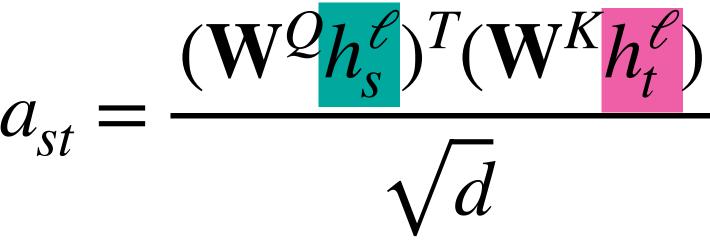
$$\tilde{h}_{s}^{\ell} = \sum_{t=1}^{T} \alpha_{st}(\mathbf{W}^{V} \mathbf{h}_{t}^{\ell})$$

Attend to values to get weighted sum

{1, ..., t, ..., T} includes *s!*

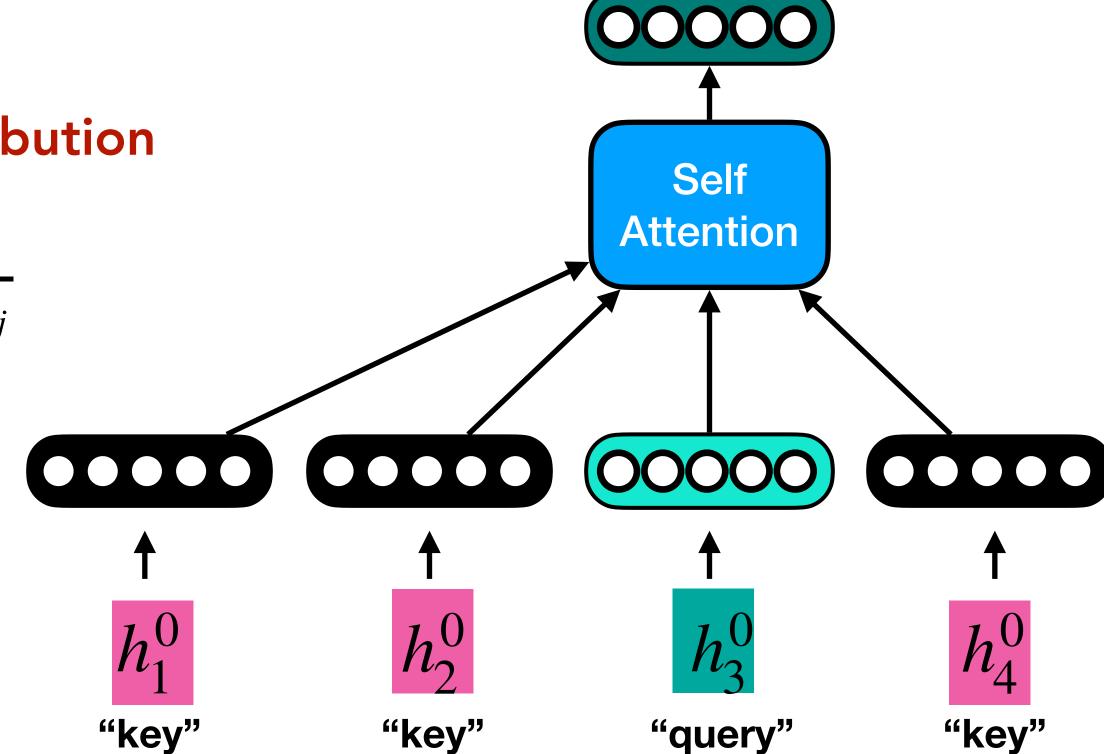
Self-attention

Compute pairwise scores



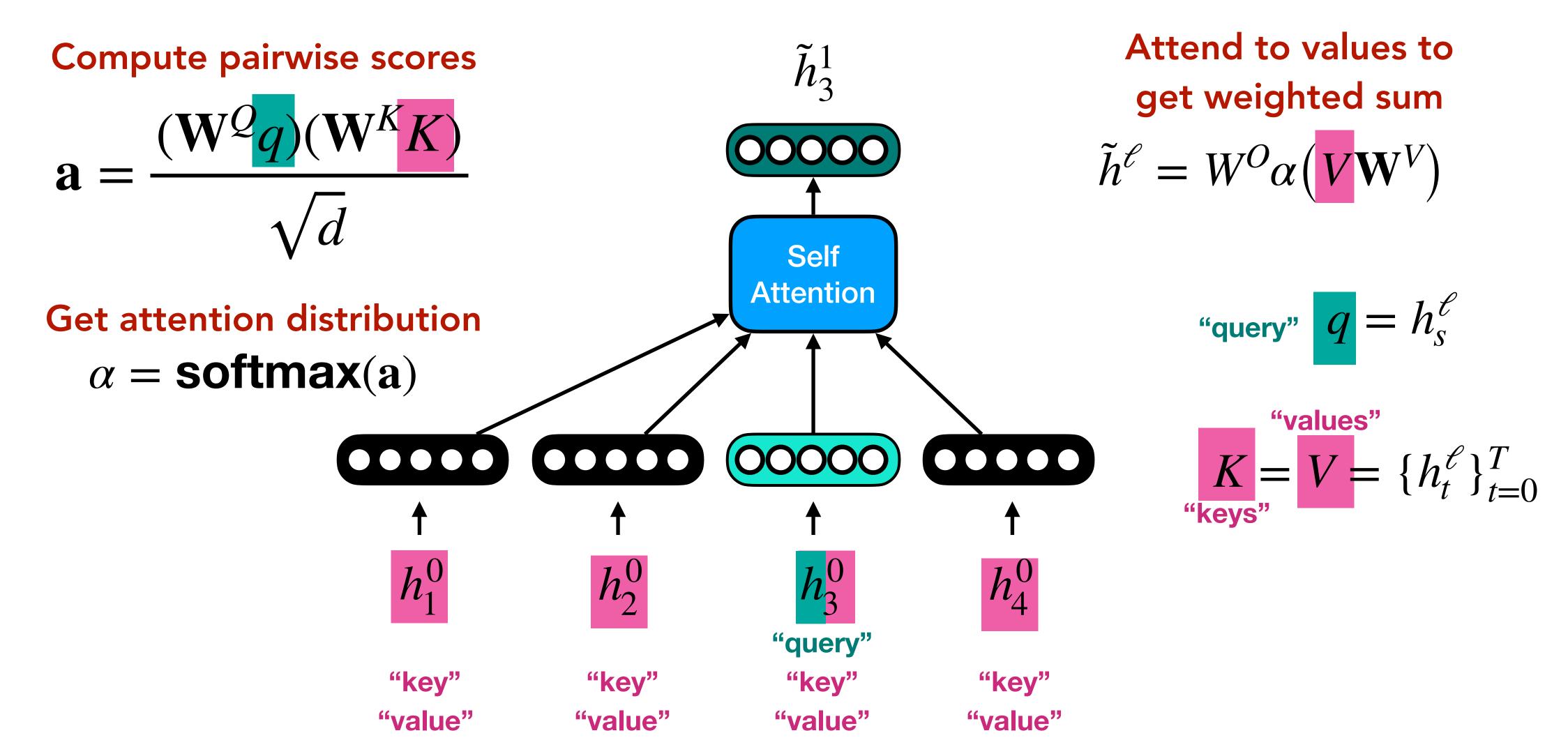
Get attention distribution

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

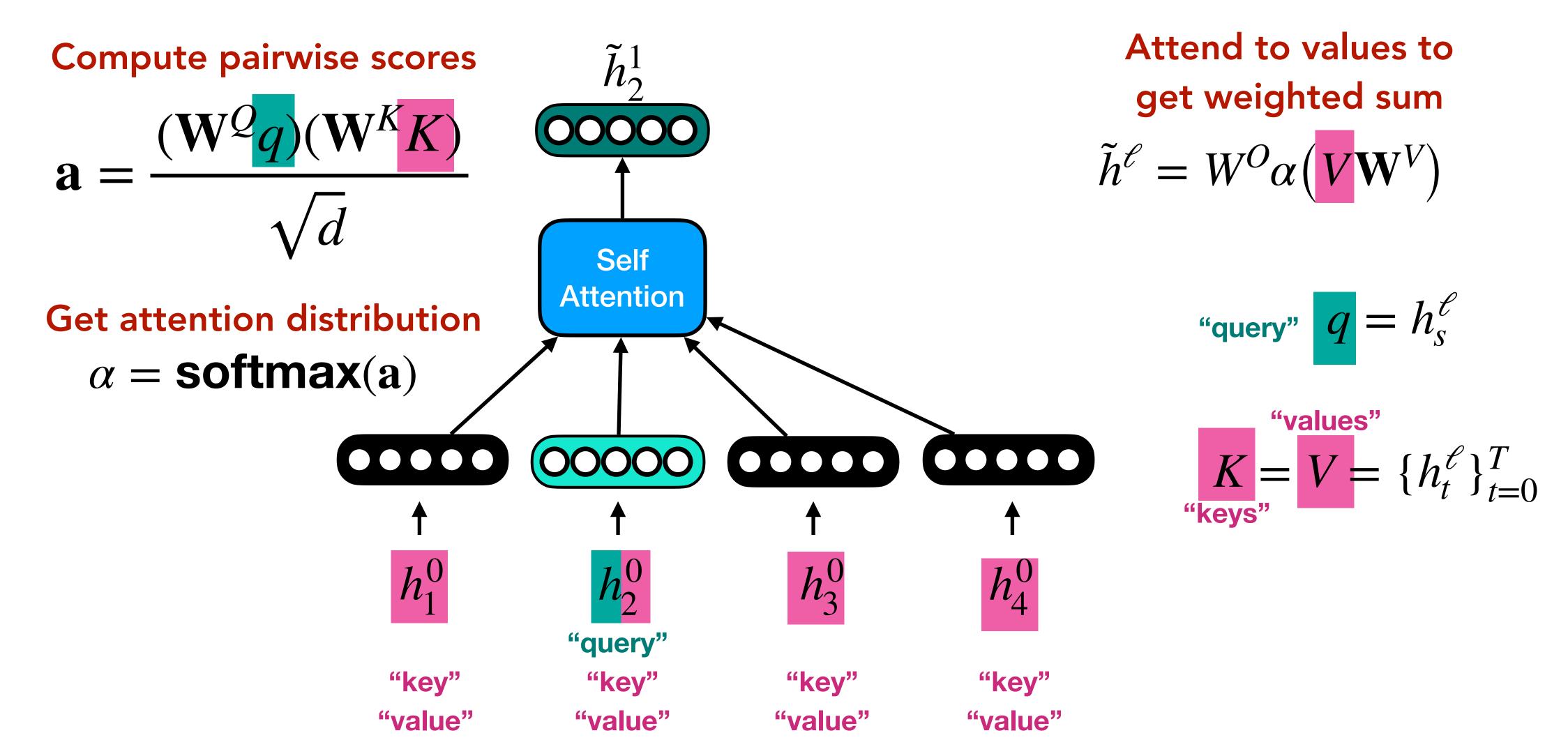


Attend to values to get weighted sum

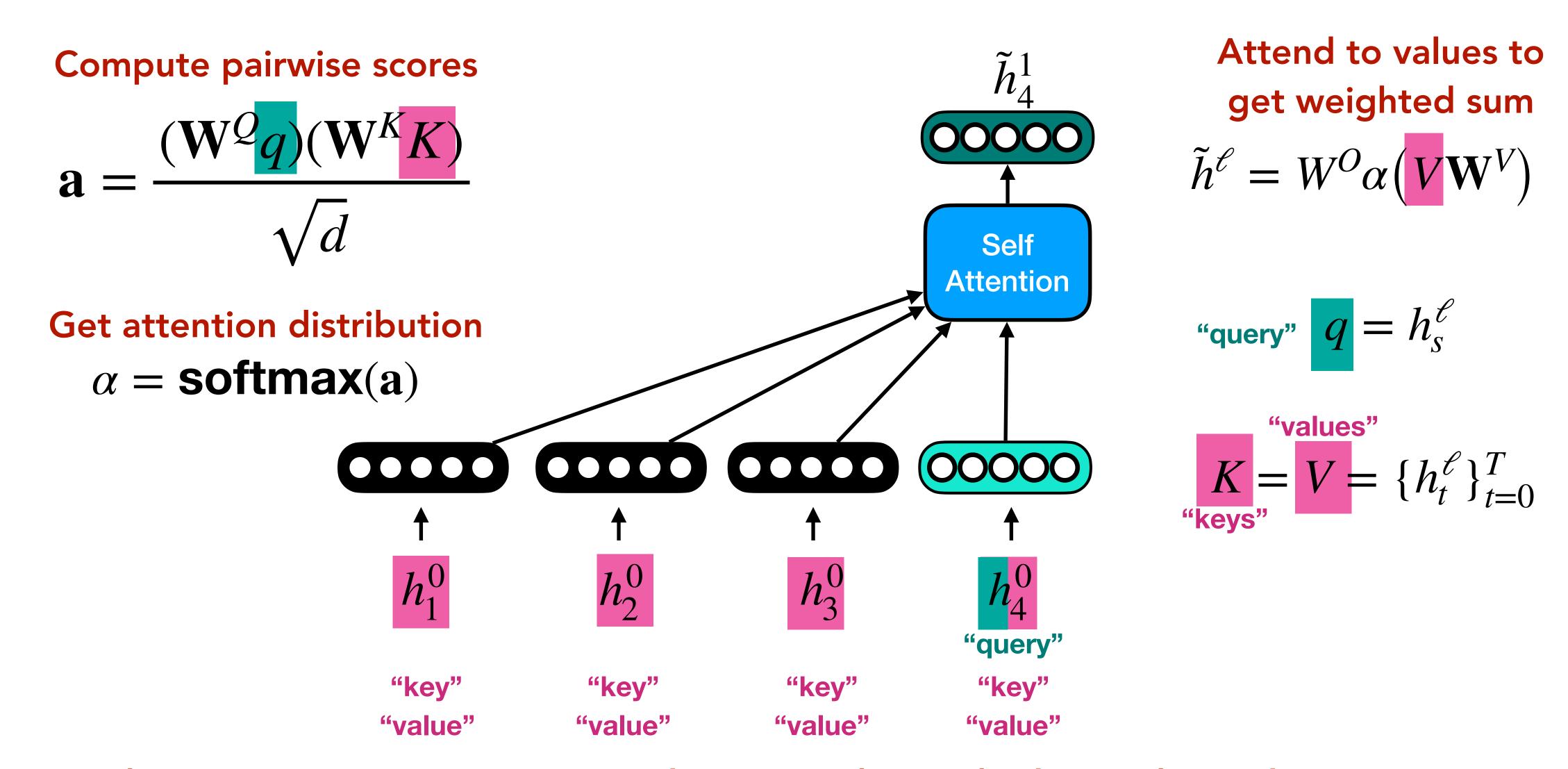
$$\tilde{h}_{s}^{\ell} = \sum_{t=1}^{T} \alpha_{st}(\mathbf{W}^{V} \mathbf{h}_{t}^{\ell})$$



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

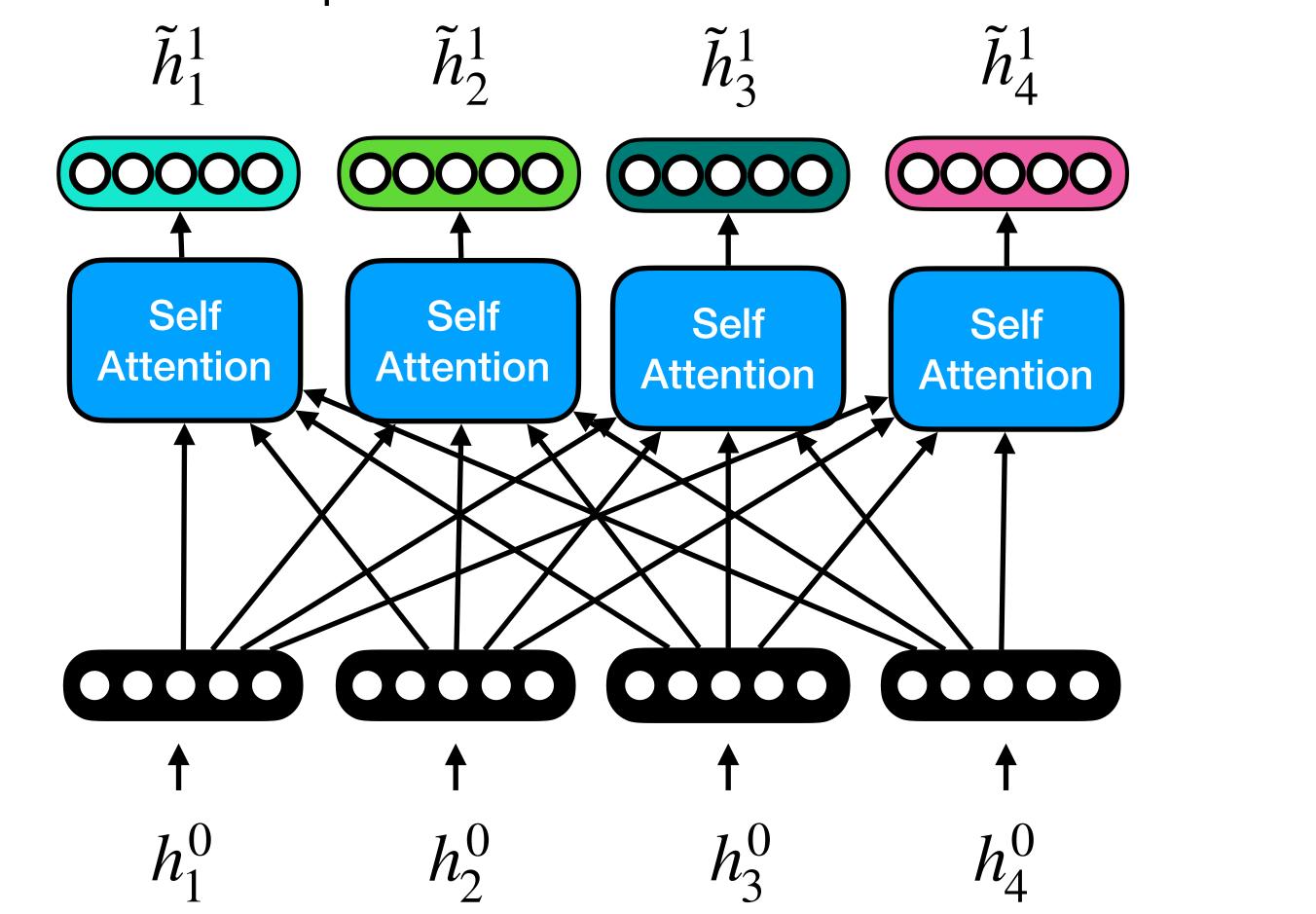


For each attention computation, every element is a key and value, and one element is a query



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

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



$$\tilde{h}_{1}^{1} = Attention(h_{1}^{0}, \{h_{t}^{0}\}_{t=0}^{t=3})$$

$$\tilde{h}_{1}^{2} = Attention(h_{2}^{0}, \{h_{t}^{0}\}_{t=0}^{t=3})$$

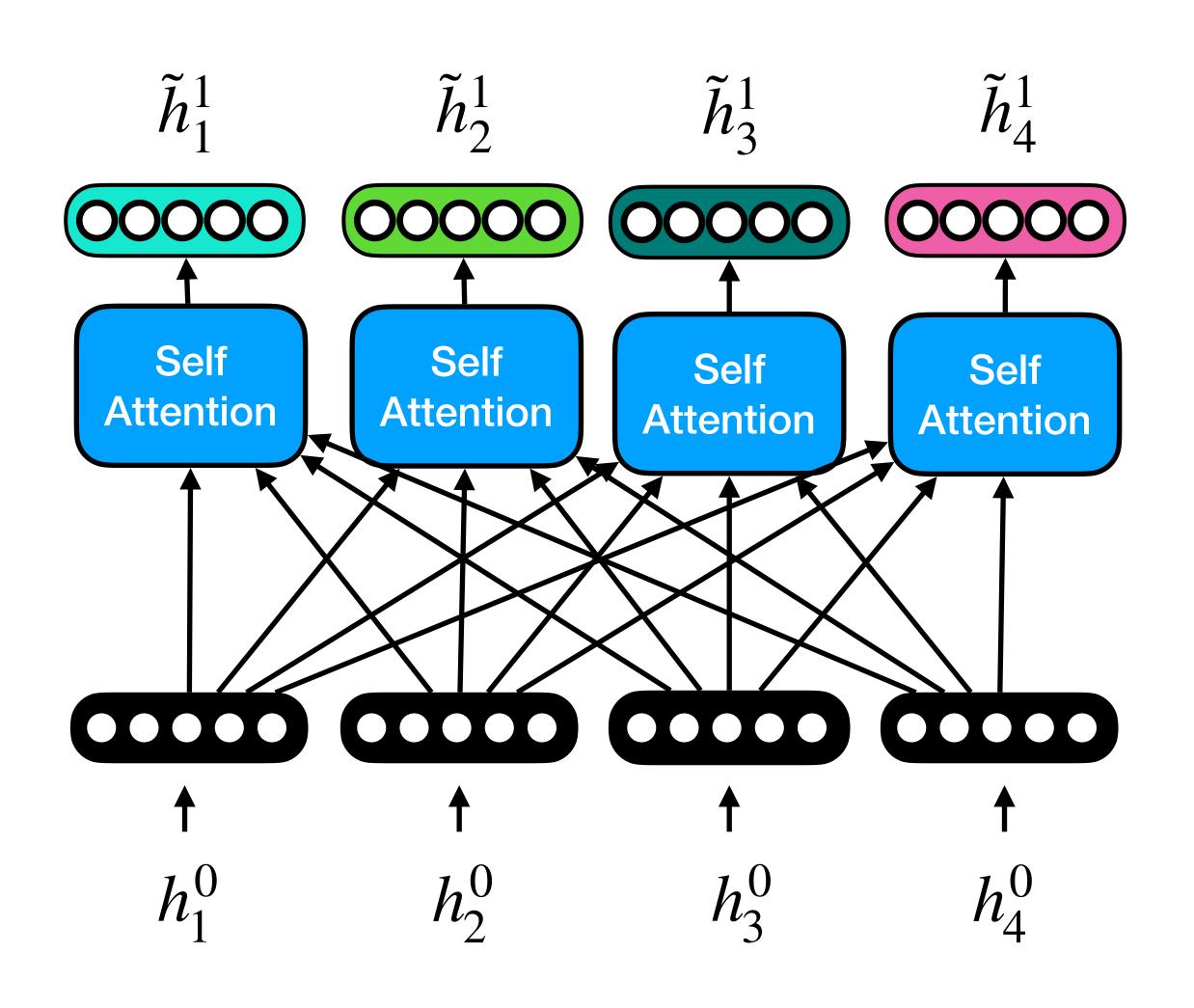
$$\tilde{h}_{1}^{3}$$
 = Attention $\left(h_{3}^{0}, \{h_{t}^{0}\}_{t=0}^{t=3}\right)$

$$\tilde{h}_{1}^{4}$$
 = 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 betweenall time-steps in sequence

Multi-Headed Self-Attention

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

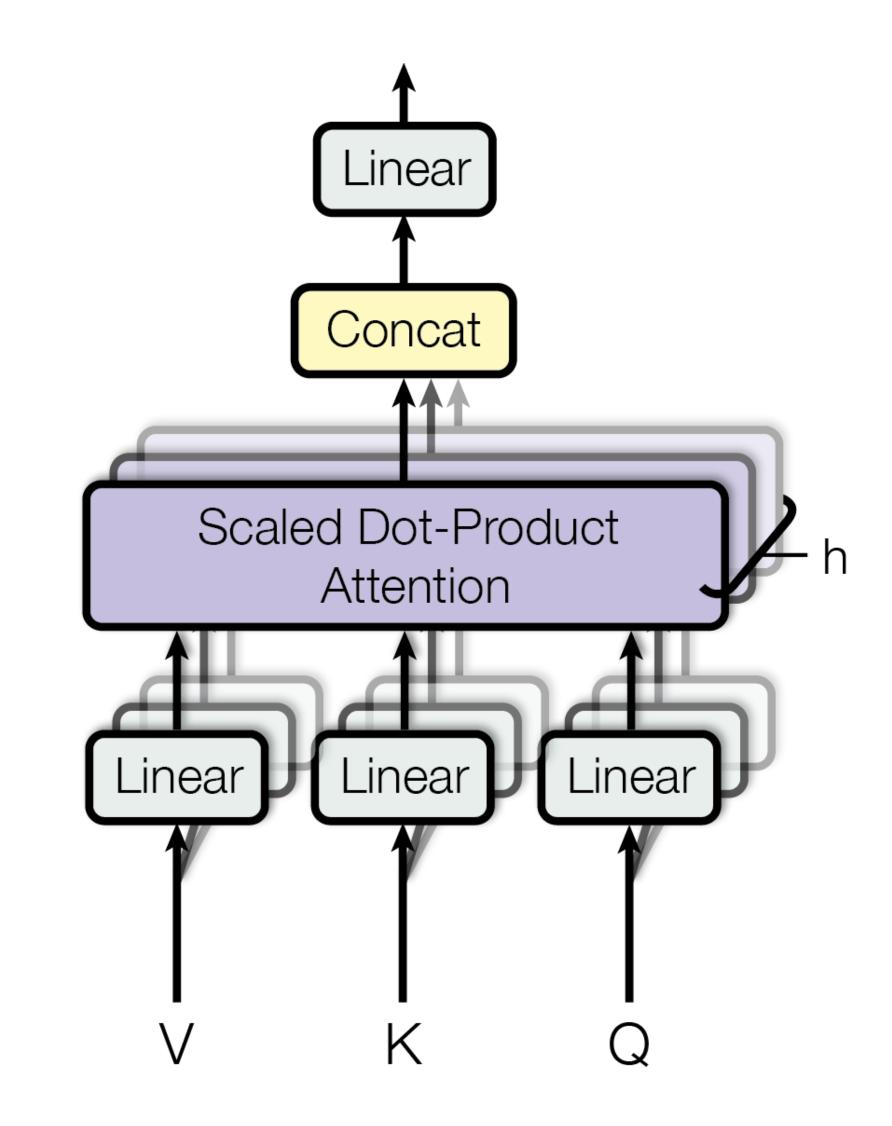
"heads"
$$\mathbf{a}_i = \frac{(\mathbf{W}_i^Q q)(\mathbf{W}_i^K K)}{\sqrt{d/H}}$$

Compute attention weights separately for each sub-vector

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

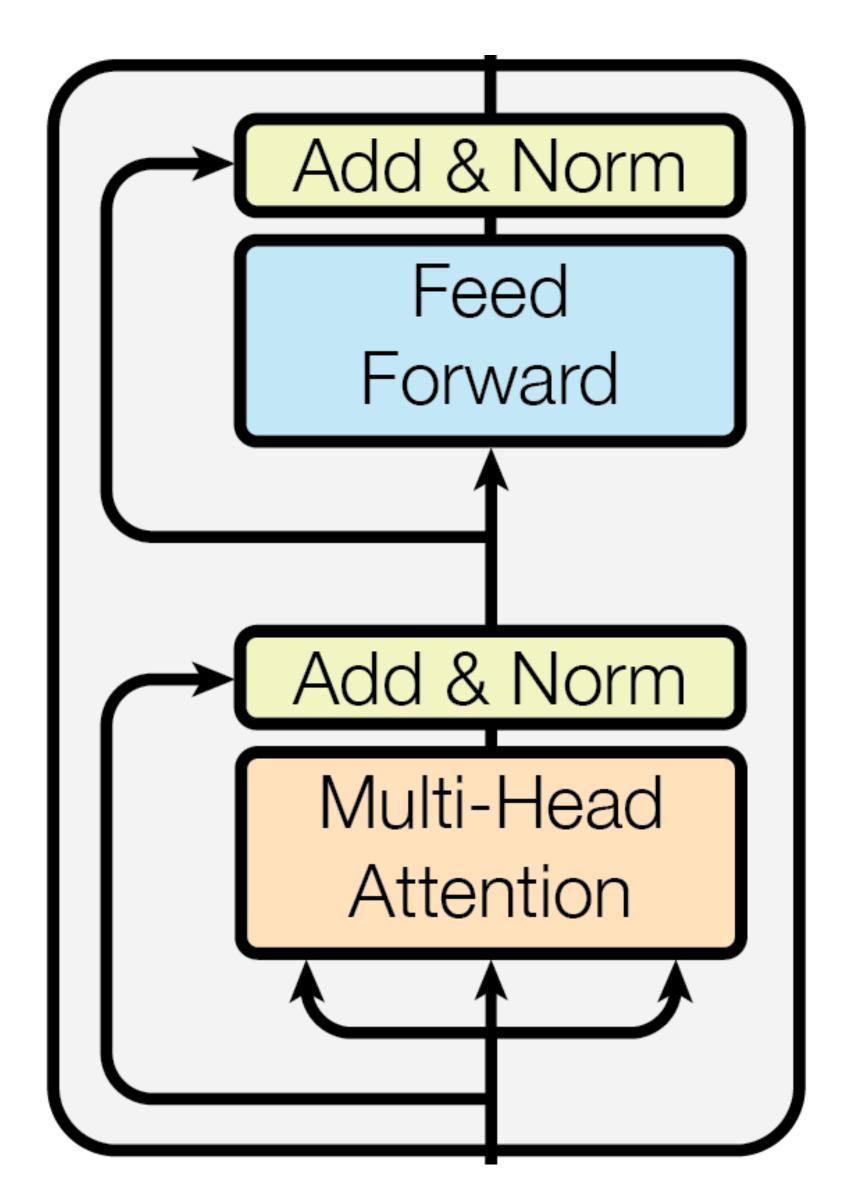
Concatenate sub-vectors for each head and project

$$\tilde{h}^{\ell} = W^{O}[\tilde{h}_{0}^{\ell}; \ldots; \tilde{h}_{i}^{\ell}; \ldots; \tilde{h}_{H}^{\ell}]$$



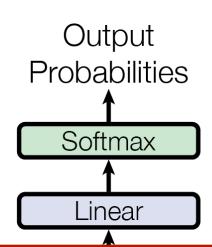
Transformer Block

- 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

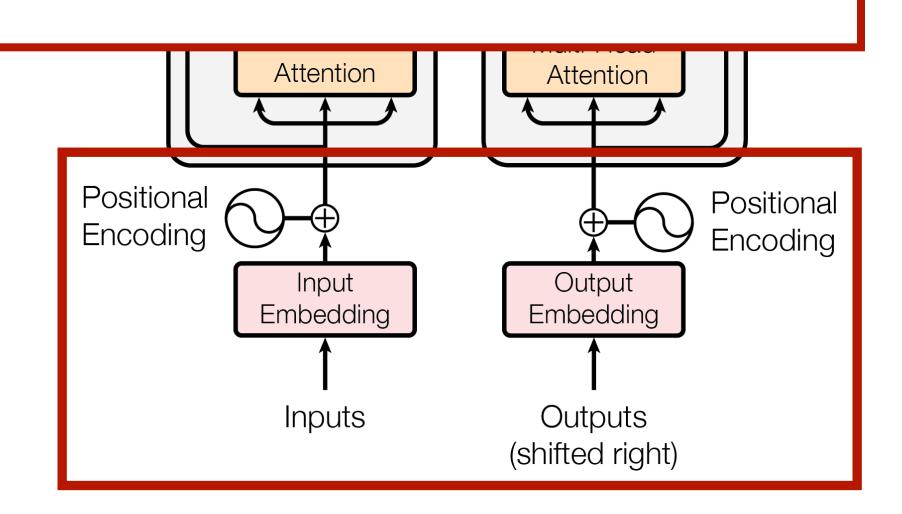


Recurrent models provided word order information

Does self-attention provide word order information?

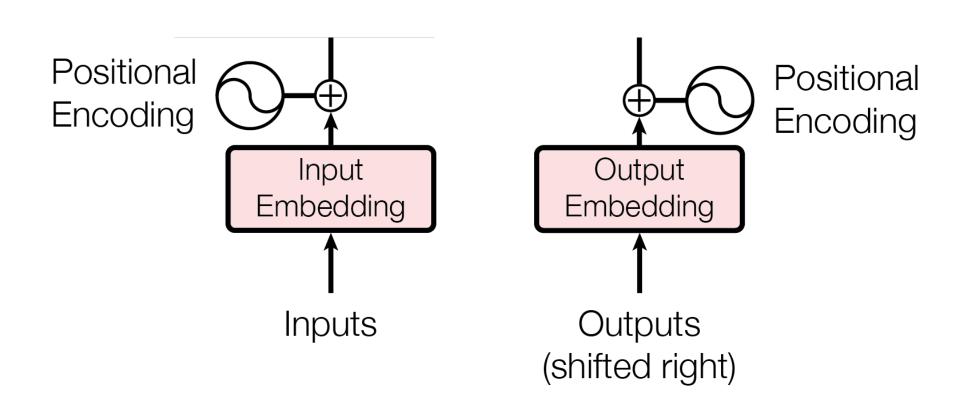
paraner.

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



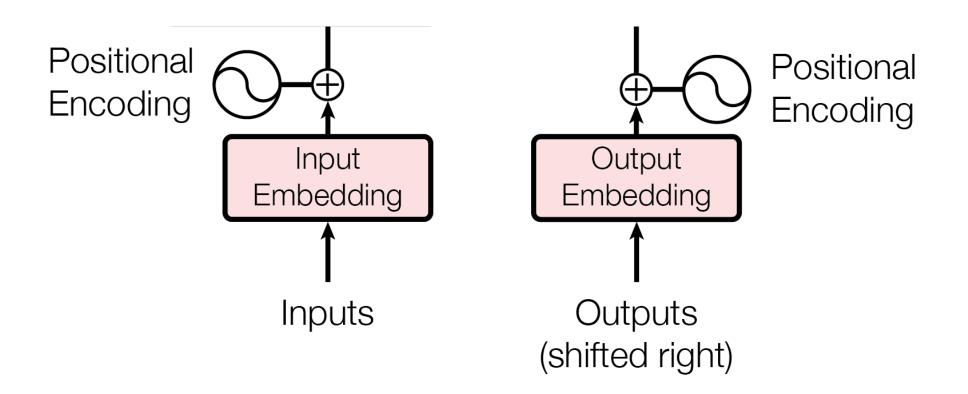
Position Embeddings

- 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.