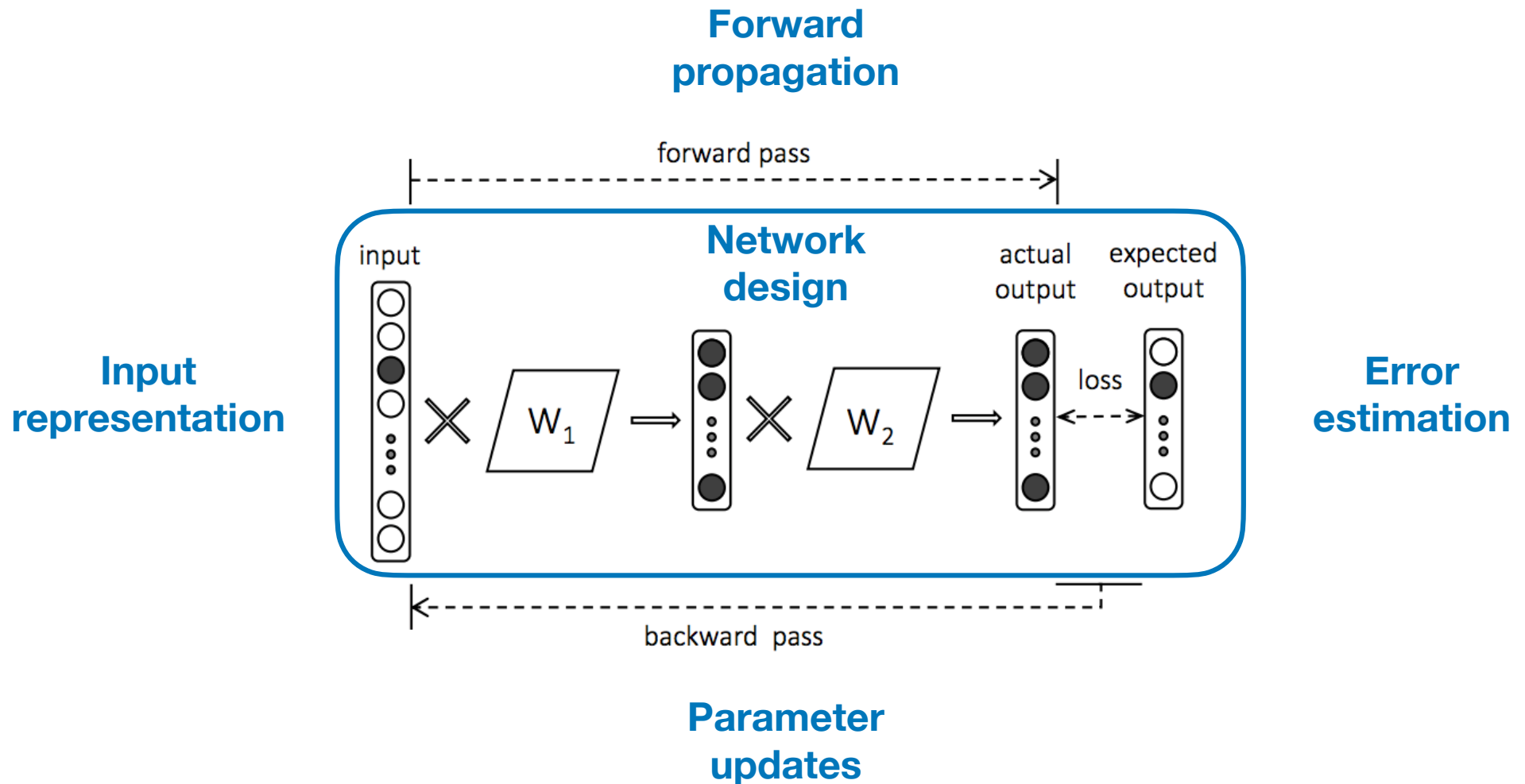


Advanced Topics in Deep Learning: Recurrent Neural Networks

Anders Søgaard

Previous Lectures



Week	Lecturer	Subject	Literature	Assignment
1	Stefan	Introduction to Neural Networks.	d2l 2.1-2.5, 2.7, 11.5.1, slides	
2	Stefan	CNNs; FCNs; U-Nets. Data augmentation; invariance; regularization e.g. dropout	d2l 6, 7, slides	Assignment 1 (May 10)
3	Anders/Phillip	RNNs (May 9) and Transformers (May 11)	Transformers: Vaswani et al. (2017) ↗	Assignment 2 (May 17)
4	Phillip/Anders	May 16-18: Representation Learning	Autoencoders (1/2): tba Self-supervised learning: blog.post ↗ Contrastive learning: Dor et al. (2018) ↗	
5	Anders	May 23-25: Adversarial Learning	Autoencoders (2/2): Chandar et al. (2011) ↗ GANs: Lample et al. (2018) ↗ Applications: Goodfellow et al. (2015) ↗	Assignment 3 [<i>MC on Representation Learning/1p Report on Adversarial Learning</i>] (May 31)
6	Anders	May 30-June 1: Interpretability	Literature: Søgaard (2022) ↗	
7	Anders	June 8: Interpretability (Note: June 6 off)	Literature: Feng and Boyd-Graber (2018) ↗ ; Jiang and Senge (2021) ↗	
8	Anders	June 13-15: Best Practices	Literature: Dodge et al. (2019) ↗ and Raji et al. (2021) ↗	Assignment 4 [<i>MC on Interpretability; 1p Report on Best Practices</i>] (June 21)

This Lecture

- Word Embeddings (Recap)
- Recurrent Neural Networks
- RNNs for Embeddings

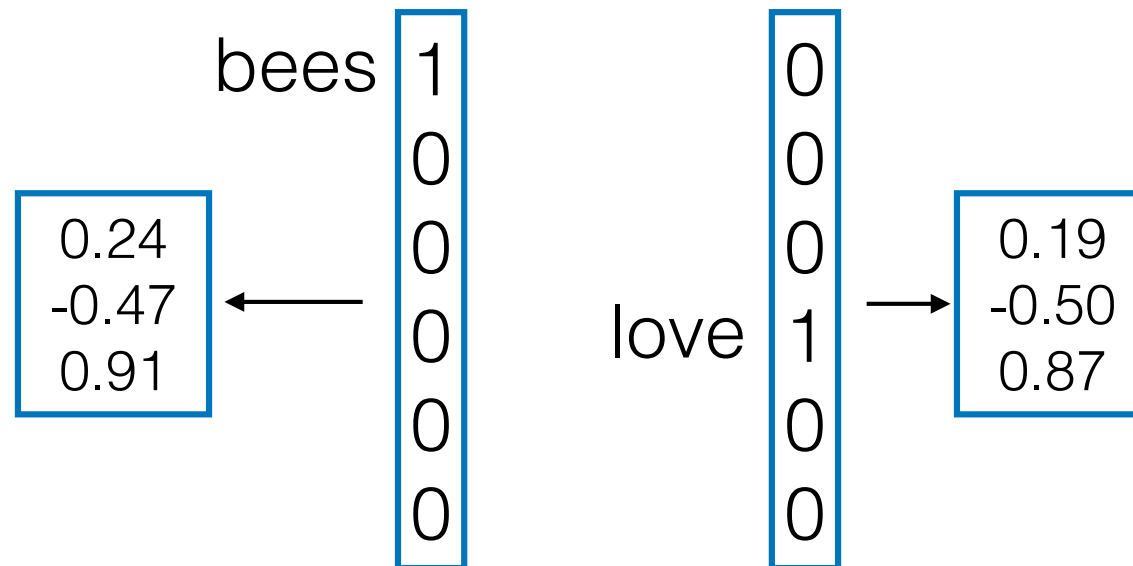
Word Embeddings

Representing Text

- One-hot vector over a vocabulary of V word types $w \in \mathcal{R}^{|V|}$
- Continuous dense vector with E dimensions (word embedding)

$$x \in \mathcal{R}^{|E|}$$

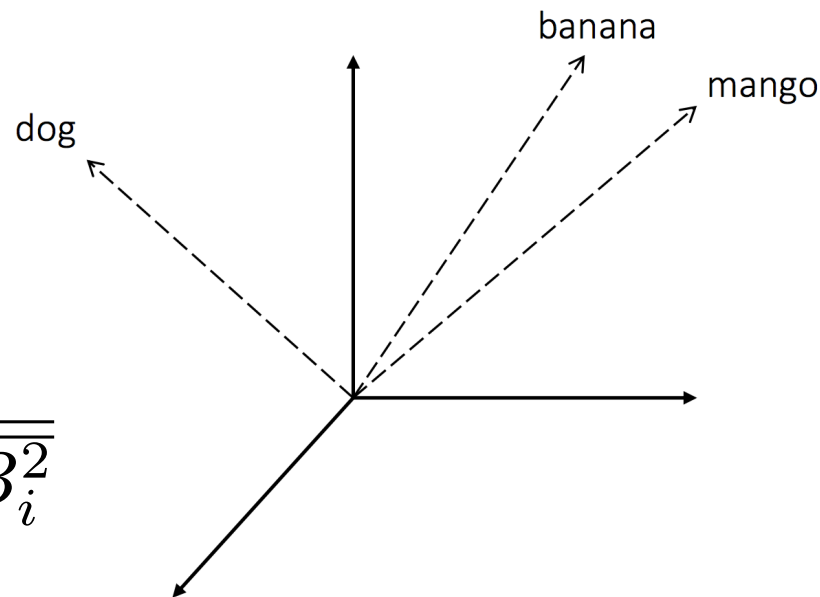
- One-hot -to- word embedding is a lookup function



Why learn Embeddings?

- Overcome sparsity of one-hot vectors
- Cosine similarity will tell us the **angle formed between vectors \mathbf{A} and \mathbf{B}** , where A and B are representations of the different terms

$$\begin{aligned}\cos(\theta) &= \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} \\ &= \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}\end{aligned}$$



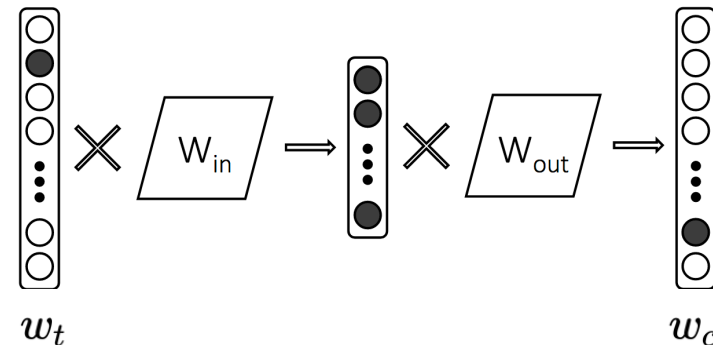
Type-based Skip-gram Embeddings

- Skip-gram embeddings are learned in the context of **predicting the surrounding words** $w_{i-|c|}, w_{i+|c|}$, given the target word w_t

$$\mathcal{L} = \sum_{i=1}^T \sum_{c \in \mathcal{C}_t} \log p(w_c | w_t)$$

Over all surrounding windows in the corpus

Estimate probability of predicting each word in the surrounding window



$$\log p(w_c | w_t) := \text{softmax}(W_{out} W_{in} w_t)$$

- Softmax doesn't scale so you can use noise-contrastive estimation

$$\sum_{t=1}^T \left[\sum_{c \in \mathcal{C}_t} \ell(s(w_t, w_c)) + \sum_{n \in \mathcal{N}_{t,c}} \ell(-s(w_t, n)) \right] \quad s(w_t, w_c) = \mathbf{u}_{w_t}^\top \mathbf{v}_{w_c}$$

Type-based FastText Embeddings

- The main shortcoming of type-based embeddings is **how should we deal with out-of-vocabulary tokens?**
- FastText proposes to split words into all sequences of 3—6 characters and learn embeddings \mathbf{z}_g for each sequence

subwords (**where**)

<wh	<whe	<wher	<where
whe	wher	where	where>
her	here	here>	
ere	ere>		
re>			

FastText

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c$$

Type-based

$$s(w_t, w_c) = \mathbf{u}_{w_t}^\top \mathbf{v}_{w_c}$$

Recurrent Neural Networks

Recurrent Neural Networks

- The **order of the words** in a sentence is not random. The order often **conforms to rules of syntax and semantics**.

The girl saw the dog with a telescope

Semantic nonsense → The telescope saw the girl with a dog

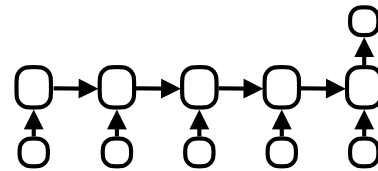
The dog with the a girl telescope saw ← Syntactic gibberish

- An averaged embedding representation of a sentence **does not capture the difference in these sentences**.

Applications of RNNs

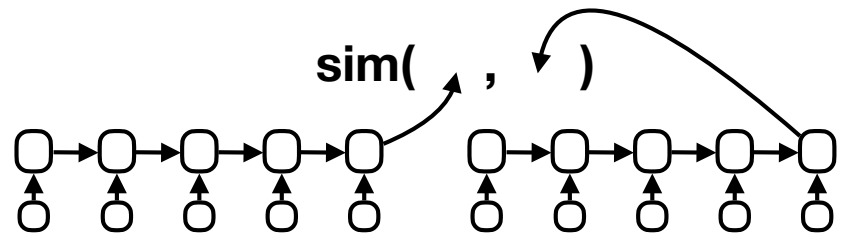
- **Classification**

- Text Categorization



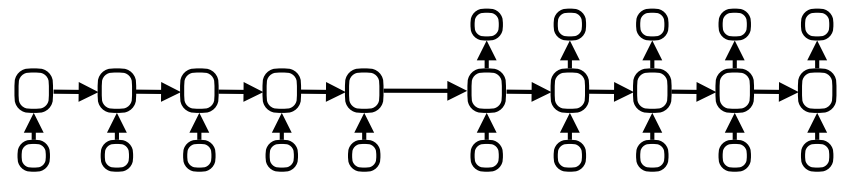
- **Sequence encoding**

- Sentence similarity



- **Autoregressive modelling**

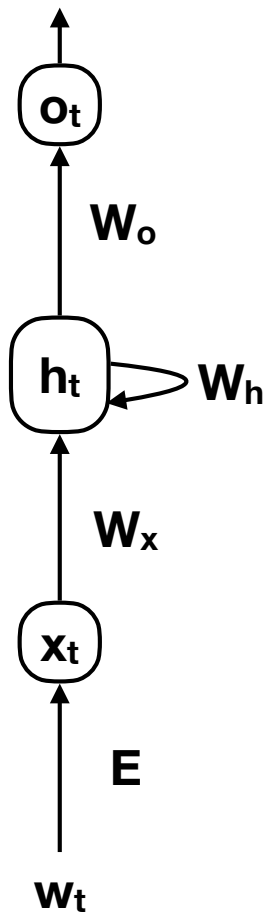
- Machine translation



Simple RNN: Definition

w: Explaining the Simple RNN model
t:

$$\mathbf{w}_{t+1} = \text{softmax}(\mathbf{o}_t)$$



Output vector

$$\mathbf{o}_t = \mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o$$

Hidden state

$$\mathbf{h}_t = f(\mathbf{W}_x \mathbf{x}_t + \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{b}_h)$$

Embedded word

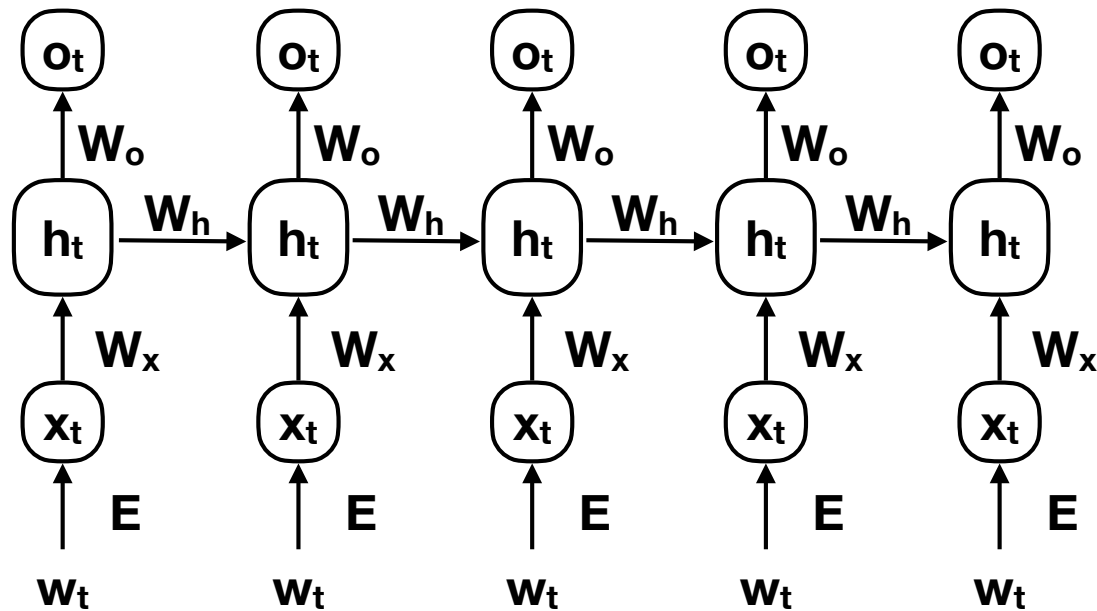
$$\mathbf{x}_t = \mathbf{E}[\mathbf{w}_t]$$

1-hot vector

$$\mathbf{w}_t = [0 \ 0 \ 0 \ 0 \ 1]$$

Unrolled Simple RNN

- This is a feed-forward neural network that has depth-in-time.



```

1  class SingleRNN(nn.Module):
2      def __init__(self, n_inputs, n_neurons):
3          super(SingleRNN, self).__init__()
4
5          self.Wx = torch.randn(n_inputs, n_neurons) # 4 X 1
6          self.Wy = torch.randn(n_neurons, n_neurons) # 1 X 1
7
8          self.b = torch.zeros(1, n_neurons) # 1 X 4
9
10     def forward(self, X0, X1):
11         self.Y0 = torch.tanh(torch.mm(X0, self.Wx) + self.b) # 4 X 1
12
13         self.Y1 = torch.tanh(torch.mm(self.Y0, self.Wy) +
14                                 torch.mm(X1, self.Wx) + self.b) # 4 X 1
15
16     return self.Y0, self.Y1

```

Exercise: What is this?

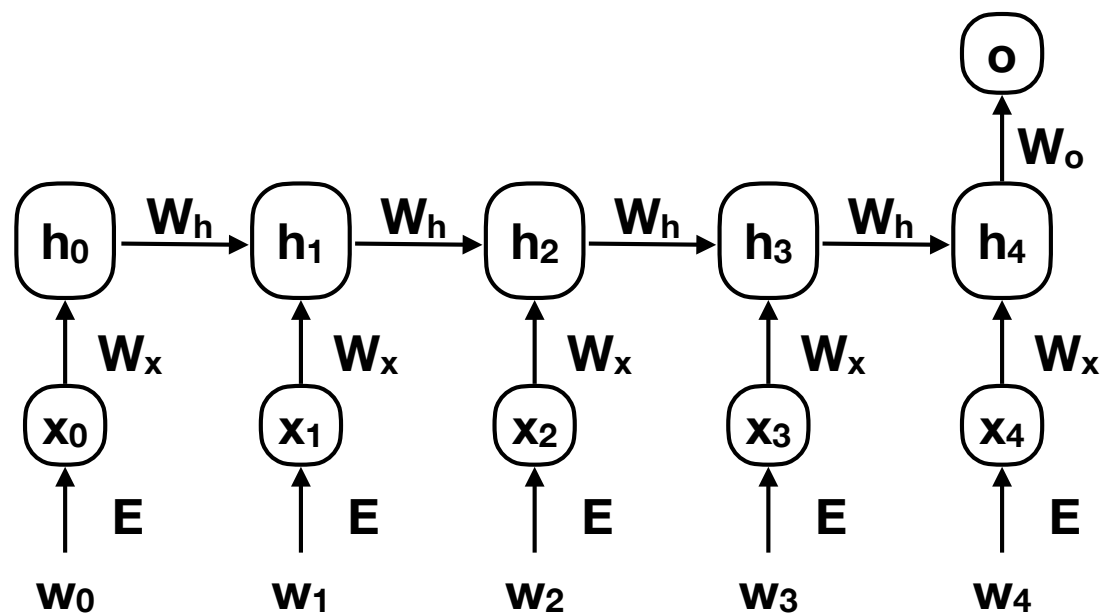
Padding data in RNNs

X

The	small	car		
The	car			
The	car	with	red	wheels
My	favourite	car		

Padded X

The	small	car	<p>	<p>
The	car	<p>	<p>	<p>
The	car	with	red	wheels
My	favourite	car	<p>	<p>



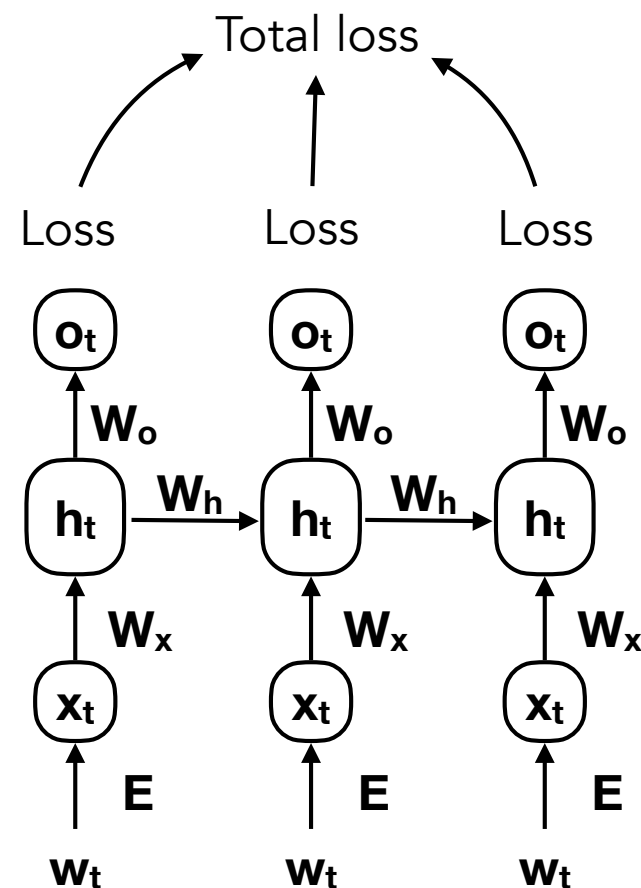
Loss masking in RNNs

Padded X

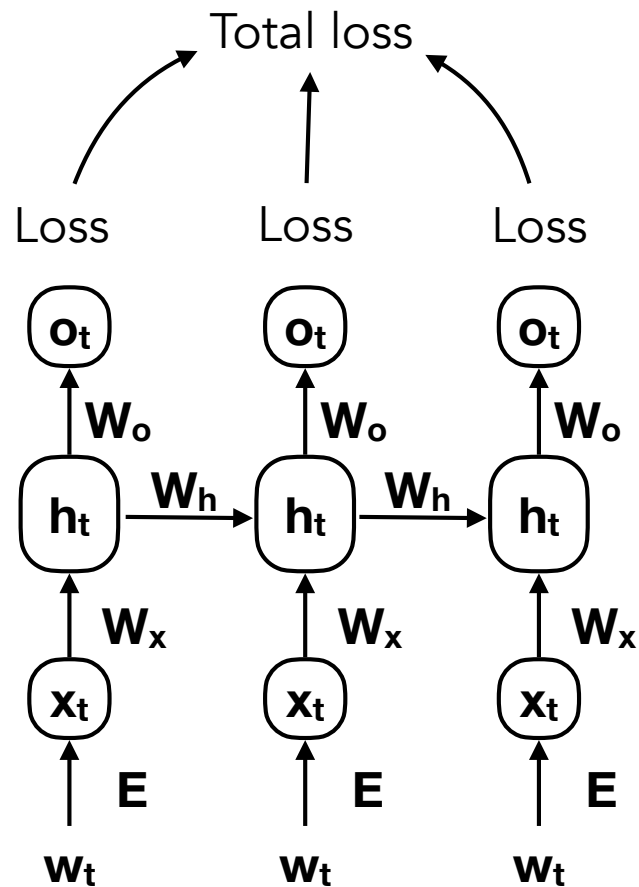
The	small	car	<p>	<p>
The	car	<p>	<p>	<p>
The	car	with	red	wheels
My	favourite	car	<p>	<p>

Mask

1	1	1	0	0
1	1	0	0	0
1	1	1	1	1
1	1	1	0	0



Training Autoregressive RNNs

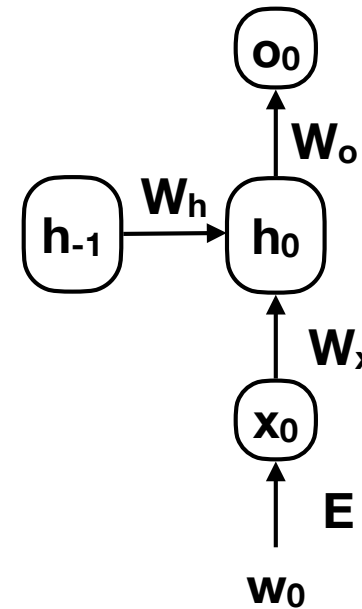


- Objective

$$- \sum_{i=1}^n \sum_{t=1}^T \log p(x_t^i | x_{<t}^i)$$

Simple RNN: Extra Notes

- What about the initial hidden state?

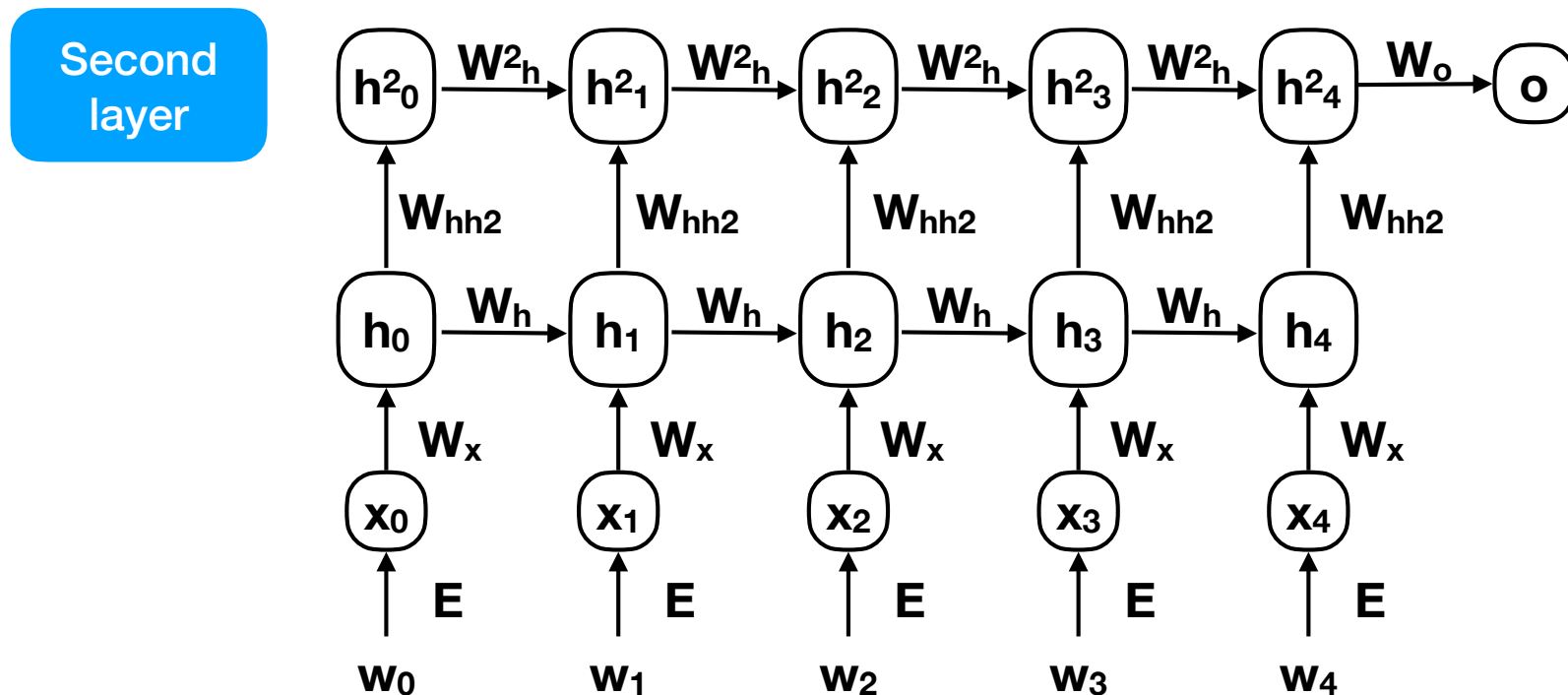


- Why are we reusing the transition vectors?

w_h w_x w_o

Multi-layer RNN

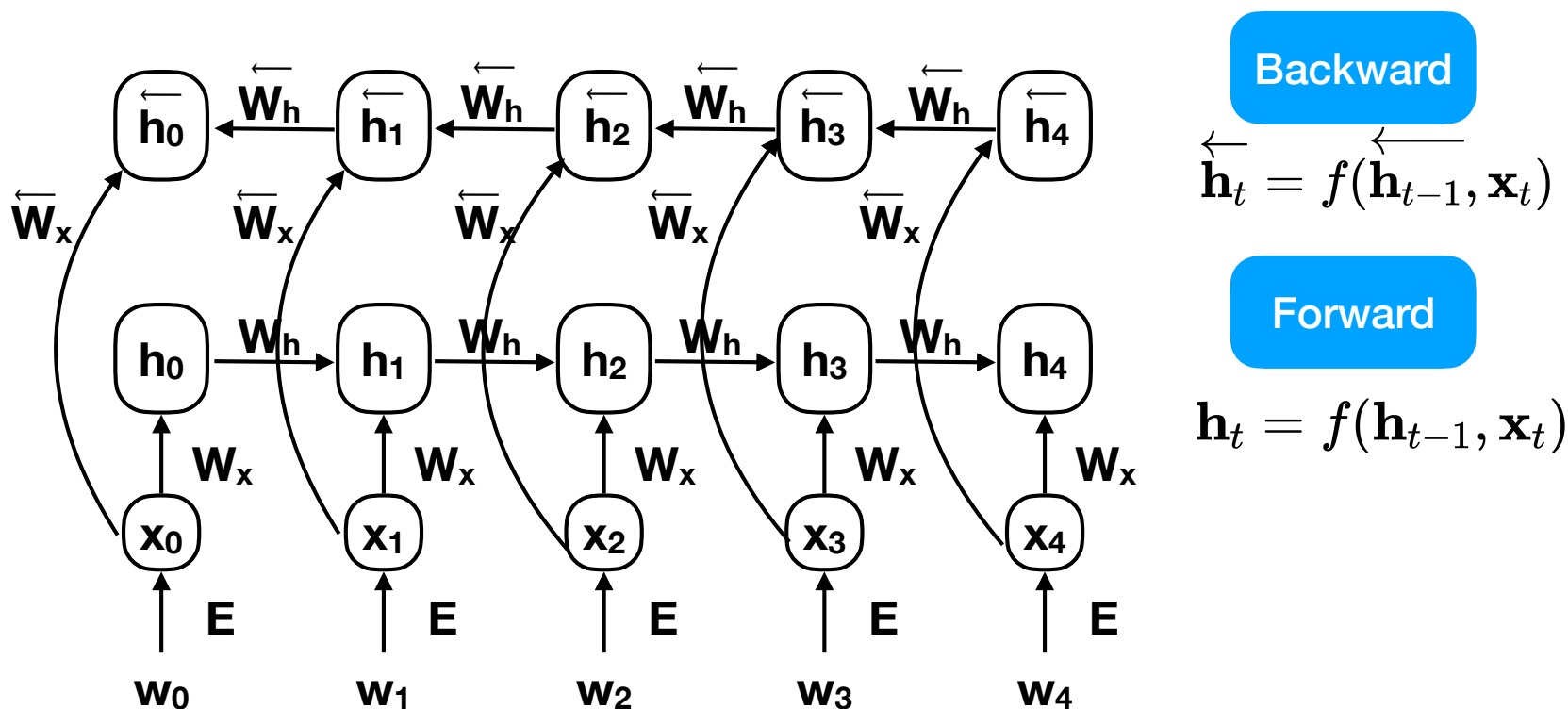
- Recurrent Neural Networks can also have layer-wise depth



$$h_t^2 = f(W_{hh2}h_t + W_h^2h_{t-1}^2 + b_h^2)$$

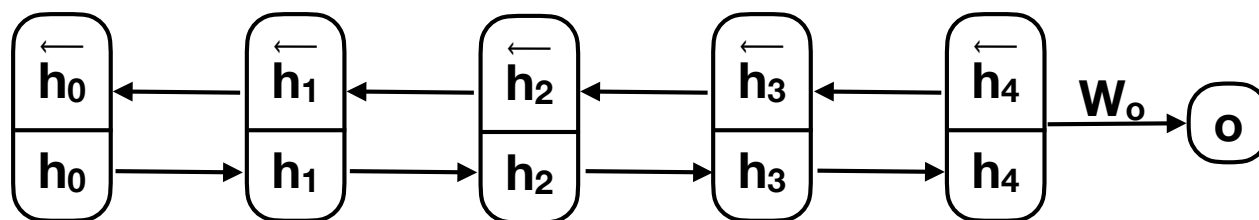
Bidirectional Networks

- Alternatively, we can encode the sequence in reverse.
- This “bidirectional” model provides a richer representation.



Bi-Directional Networks

- Now we can train our classifiers on a concatenation of the forward and backward representations of the sequence



$$\mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t) \quad \overleftarrow{\mathbf{h}}_t = f(\overleftarrow{\mathbf{h}}_{t-1}, \mathbf{x}_t)$$

$$\mathbf{h}_t = [\mathbf{h}_t \oplus \overleftarrow{\mathbf{h}}_t]$$

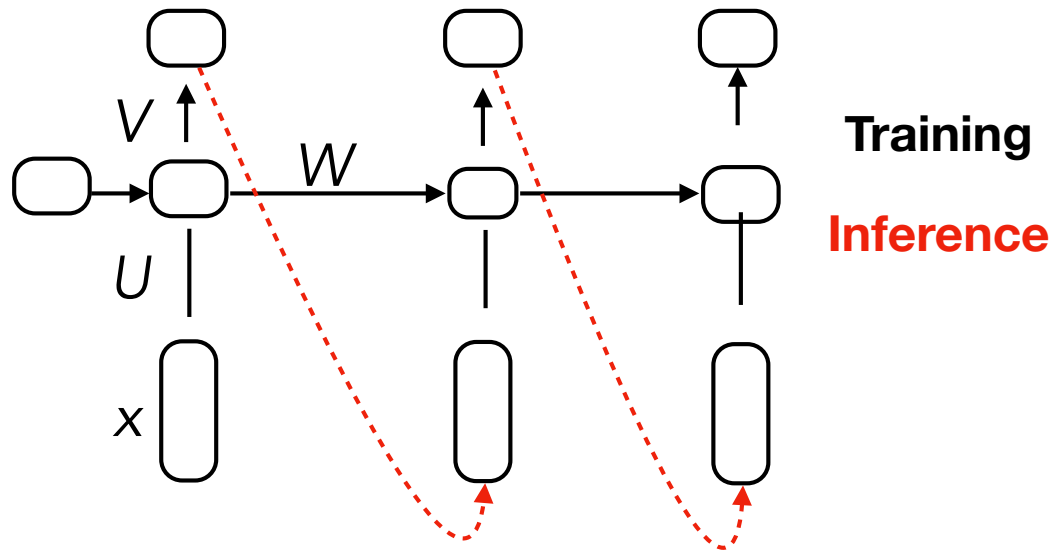
RNNs for Embeddings

Autoregressive Recurrent Neural Networks

- Elman-style model (Cognitive Science, 1990)
 - LSTM (Hochreiter and Schmidhuber, 1997)
 - Gated Recurrent Unit (Cho et al. 2014)

$$o_t = \text{softmax}(V s_t)$$

$$s_t = f(Ux_t + Ws_{t-1})$$



- Create a **contextualized representation** of a word in the sequence of hidden states s_t

Contextual Embeddings from Language Models

- Represent each token using multi-layer bidirectional LSTMs.
- Final representation is a **weighted sum of the representations** at the different layers in the model

$$\mathbf{ELMo}_k^{task} = E(R_k; \Theta^{task}) = \gamma^{task} \sum_{j=0}^L s_j^{task} \mathbf{h}_{k,j}^{LM}$$

Sum over each Layer Softmax-normalized task weights Feature vector for token k at layer j

