

# POS Tagging and Sequence Models

JURAFSKY AND MARTIN CHAPTER 8

# Recap: Sequence Models

A **sequence model** or **sequence classifier** is a model whose job is to assign a label or class to each unit in a sequence, thus mapping a sequence of observations to a sequence of labels.

A Hidden Markov Model (HMM) is a probabilistic sequence model: given a sequence of words, it computes a probability distribution over possible sequences of labels and chooses the best label sequence.

# Unknown words

One useful feature for distinguishing parts of speech is **word shape** (proper nouns start with a capital).

The strongest feature is **morphology**.

Words that end in

- -s tend to be **plural nouns (NNS)**
- -ed tend to be **past participles (VBN)**
- -able tend to be **adjectives (JJ)**
- and so on

# Learning suffix model

Store the final letter sequence (suffixes) for up to 10 letters.

For each such sequence, record the probability of the tag that it was associated with during training.

Use back-off to smooth these probabilities for successively shorter sequences.

Trigram HMM with unknown word handling:	96.7%
State of the art neural network POS tagging:	97%

# Maximum Entropy Markov Models

Could we add features like word shape and suffixes directly into the model in a clean way? We had this for classification with **logistic regression**. But logistic regression isn't a sequence model, since it assigns a class to a single observation.

We can turn logistic regression into a discriminative sequence model simply by running it on successive words, using the class assigned to the prior word as a feature in the classification of the next word. This is called a **MEMM**.

# MEMMs v HMMs

HMM:

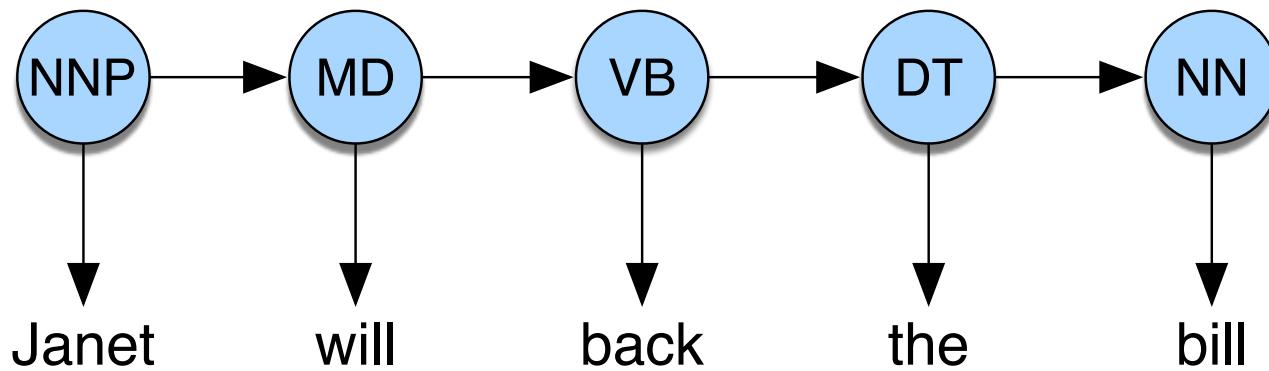
$$\begin{aligned}\hat{T} &= \underset{T}{\operatorname{argmax}} P(T|W) \\ &= \underset{T}{\operatorname{argmax}} P(W|T)P(T) \\ &= \underset{T}{\operatorname{argmax}} \prod_i P(\text{word}_i|\text{tag}_i) \prod_i P(\text{tag}_i|\text{tag}_{i-1})\end{aligned}$$

MEMM:

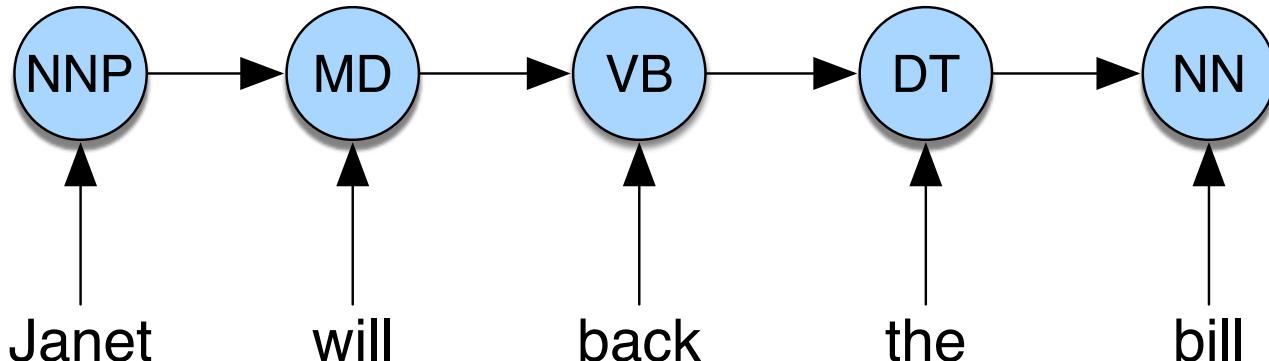
$$\begin{aligned}\hat{T} &= \underset{T}{\operatorname{argmax}} P(T|W) \\ &= \underset{T}{\operatorname{argmax}} \prod_i P(t_i|w_i, t_{i-1})\end{aligned}$$

# MEMMs v HMMs

HMM:

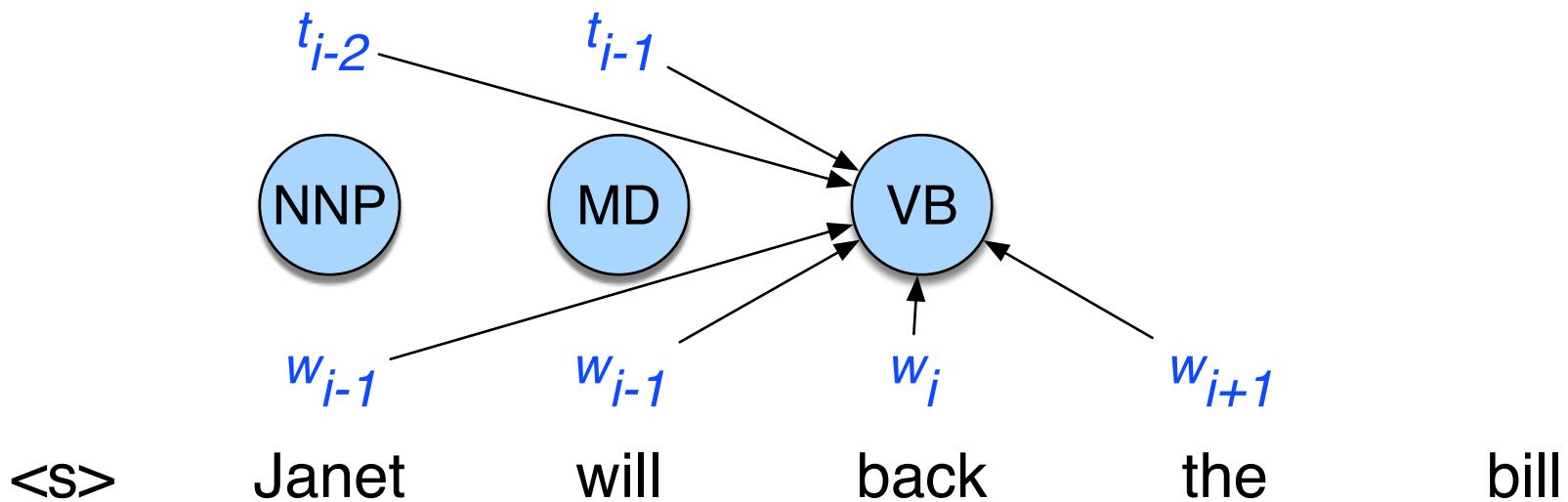


MEMM:



# Features in a MEMM

We can build MEMMs that don't just condition on  $w_i$  and  $t_{i-1}$ . It is easy to incorporate lots of features in a discriminative sequence model.



# Feature templates

A basic MEMM part-of-speech tagger conditions on the observation word it- self, neighboring words, and previous tags, and various combinations, using feature templates like the following

$$\begin{aligned} & \langle t_i, w_{i-2} \rangle, \langle t_i, w_{i-1} \rangle, \langle t_i, w_i \rangle, \langle t_i, w_{i+1} \rangle, \langle t_i, w_{i+2} \rangle \\ & \quad \langle t_i, t_{i-1} \rangle, \langle t_i, t_{i-2}, t_{i-1} \rangle, \\ & \quad \langle t_i, t_{i-1}, w_i \rangle, \langle t_i, w_{i-1}, w_i \rangle \langle t_i, w_i, w_{i+1} \rangle, \end{aligned}$$

*Janet/NNP will/MD back/VB the/DT bill/NN*, when  $w_i$  is the word *back*

$t_i = \text{VB}$  and  $w_{i-2} = \text{Janet}$

$t_i = \text{VB}$  and  $w_{i-1} = \text{will}$

$t_i = \text{VB}$  and  $w_i = \text{back}$

$t_i = \text{VB}$  and  $w_{i+1} = \text{the}$

$t_i = \text{VB}$  and  $w_{i+2} = \text{bill}$

$t_i = \text{VB}$  and  $t_{i-1} = \text{MD}$

$t_i = \text{VB}$  and  $t_{i-1} = \text{MD}$  and  $t_{i-2} = \text{NNP}$

$t_i = \text{VB}$  and  $w_i = \text{back}$  and  $w_{i+1} = \text{the}$

# Features for unknown words

- $w_i$  contains a particular prefix (from all prefixes of length  $\leq 4$ )
- $w_i$  contains a particular suffix (from all suffixes of length  $\leq 4$ )
- $w_i$  contains a number
- $w_i$  contains an upper-case letter
- $w_i$  contains a hyphen
- $w_i$  is all upper case
- $w_i$ 's word shape
- $w_i$ 's short word shape
- $w_i$  is upper case and has a digit and a dash (like *CFC-12*)
- $w_i$  is upper case and followed within 3 words by Co., Inc., etc.

# Features for *well-dressed*

$\text{prefix}(w_i) = \text{w}$

$\text{prefix}(w_i) = \text{we}$

$\text{prefix}(w_i) = \text{wel}$

$\text{prefix}(w_i) = \text{well}$

$\text{suffix}(w_i) = \text{ssed}$

$\text{suffix}(w_i) = \text{sed}$

$\text{suffix}(w_i) = \text{ed}$

$\text{suffix}(w_i) = \text{d}$

$\text{has-hyphen}(w_i)$

$\text{word-shape}(w_i) = \text{xxxx-xxxxxxxx}$

$\text{short-word-shape}(w_i) = \text{x-x}$

# Other sequence labeling tasks

A **sequence model** or **sequence classifier** is a model whose job is to assign a label or class to each unit in a sequence, thus mapping a sequence of observations to a sequence of labels.

What other sequence labeling tasks can you think of?

# Named Entity Recognition

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	The Mt. Sanitas loop is in Sunshine Canyon.
Geo-Political Entity	GPE	countries, states, provinces	Palo Alto is raising the fees for parking.
Facility	FAC	bridges, buildings, airports	Consider the Tappan Zee Bridge.
Vehicles	VEH	planes, trains, automobiles	It was a classic Ford Falcon.

# Category Ambiguity in NER

Name	Possible Categories
<i>Washington</i>	Person, Location, Political Entity, Organization, Vehicle
<i>Downing St.</i>	Location, Organization
<i>IRA</i>	Person, Organization, Monetary Instrument
<i>Louis Vuitton</i>	Person, Organization, Commercial Product

Examples of type ambiguities in the use of the name *Washington*:

[PER Washington] was born into slavery on the farm of James Burroughs.

[ORG Washington] went up 2 games to 1 in the four-game series.

Blair arrived in [LOC Washington] for what may well be his last state visit.

In June, [GPE Washington] passed a primary seatbelt law.

The [VEH Washington] had proved to be a leaky ship, every passage I made...

# Temporal Expression Extraction

Absolute	Relative	Durations
April 24, 1916	yesterday	four hours
The summer of '77	next semester	three weeks
10:15 AM	two weeks from yesterday	six days
The 3rd quarter of 2006	last quarter	the last three quarters

Lexical triggers for temporal expressions:

Category	Examples
Noun	<i>morning, noon, night, winter, dusk, dawn</i>
Proper Noun	<i>January, Monday, Ides, Easter, Rosh Hashana, Ramadan, Tet</i>
Adjective	<i>recent, past, annual, former</i>
Adverb	<i>hourly, daily, monthly, yearly</i>

# Fine grained entity recognition

<b>person</b>	doctor engineer monarch musician politician religious_leader soldier terrorist	<b>organization</b>	airline company educational_institution fraternity_sorority sports_league sports_team	terrorist_organization government_agency government political_party educational_department military news_agency
<b>location</b>	body_of_water island mountain glacier astral_body cemetery park	<b>product</b>	camera mobile_phone computer software game instrument weapon	<b>art</b> written_work film newspaper play music
				<b>event</b> military_conflict attack natural_disaster election sports_event protest terrorist_attack
<b>building</b>	time color award educational_degree title law ethnicity language religion god	chemical_thing biological_thing medical_treatment disease symptom drug body_part living_thing animal food	website broadcast_network broadcast_program tv_channel currency stock_exchange algorithm programming_language transit_system transit_line	

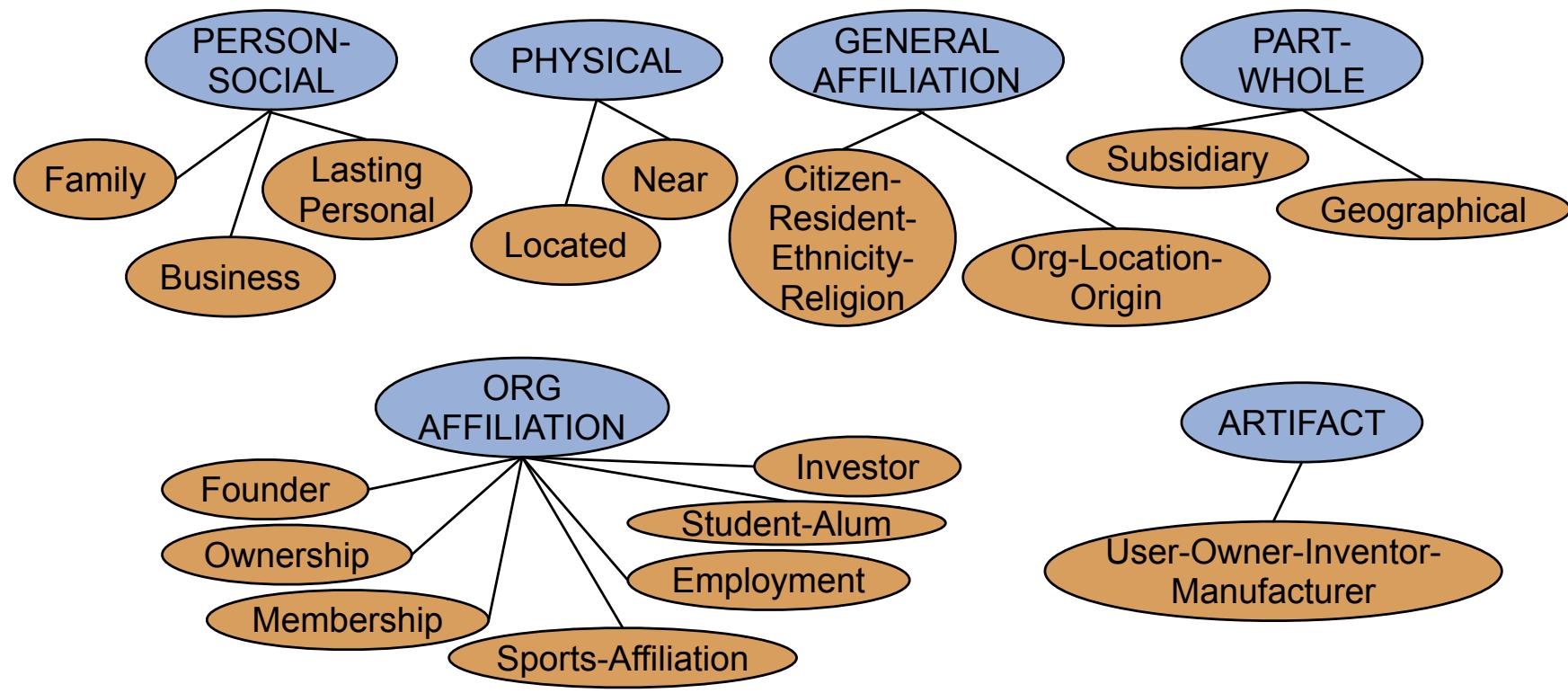
# Relation Extraction

Citing high fuel prices, [ORG United Airlines] said [TIME Friday] it has increased fares by [MONEY \$6] per round trip on flights to some cities also served by lower-cost carriers. [ORG American Airlines], a unit of [ORG AMR Corp.], immediately matched the move, spokesman [PER Tim Wagner] said. [ORG United], a unit of [ORG UAL Corp.], said the increase took effect [TIME Thursday] and applies to most routes where it competes against discount carriers, such as [LOC Chicago] to [LOC Dallas] and [LOC Denver] to [LOC San Francisco].

Tim Wagner **is a spokesman for** American Airlines

United **is a unit of** UAL Corp

American **is a unit of** AMR



Relations	Types	Examples
Physical-Located	PER-GPE	He was in <b>Tennessee</b>
Part-Whole-Subsidiary	ORG-ORG	<b>XYZ</b> , the parent company of <b>ABC</b>
Person-Social-Family	PER-PER	<b>Yoko</b> 's husband <b>John</b>
Org-AFF-Founder	PER-ORG	<b>Steve Jobs</b> , co-founder of <b>Apple</b> ...

# Unified Medical Language System

Entity	Relation	Entity
Injury	disrupts	Physiological Function
Bodily Location	location-of	Biologic Function
Anatomical Structure	part-of	Organism
Pharmacologic Substance	causes	Pathological Function
Pharmacologic Substance	treats	Pathologic Function

Given a medical sentence like:

Doppler echocardiography can be used to diagnose left anterior descending artery stenosis in patients with type 2 diabetes

Extract the UMLS relation:

*Echocardiography (Doppler) Diagnoses Acquired stenosis*

# Wikipedia info boxes

<b>Stephen Hawking</b> CH CBE FRS FRSA	
 A black and white photograph of Stephen Hawking, a theoretical physicist and cosmologist, sitting in his wheelchair. He is wearing glasses and a dark jacket over a light-colored shirt. He is smiling and looking towards the camera. In the background, there is a computer monitor and some papers on a desk.	
Hawking at <a href="#">NASA's StarChild Learning Center</a> , ca. 1999	
<b>Born</b>	Stephen William Hawking 8 January 1942 <a href="#">Oxford</a> , England
<b>Died</b>	14 March 2018 (aged 76) <a href="#">Cambridge</a> , England
<b>Education</b>	<a href="#">St Albans School</a> , <a href="#">Hertfordshire</a>
<b>Alma mater</b>	<a href="#">University of Oxford</a> (BA) <a href="#">University of Cambridge</a> (MA, PhD)
<b>Known for</b>	<a href="#">Hawking radiation</a> <a href="#">Penrose–Hawking theorems</a> <a href="#">Bekenstein–Hawking formula</a> <a href="#">Hawking energy</a> <a href="#">Gibbons–Hawking ansatz</a> <a href="#">Gibbons–Hawking effect</a> <a href="#">Gibbons–Hawking space</a> <a href="#">Gibbons–Hawking–York boundary term</a> <a href="#">Thorne–Hawking–Preskill bet</a>
<b>Spouse(s)</b>	<a href="#">Jane Wilde</a> (m. 1965; div. 1995) <a href="#">Elaine Mason</a> (m. 1995; div. 2006)
<b>Children</b>	3, including <a href="#">Lucy</a>
<b>Awards</b>	<a href="#">Adams Prize</a> (1966) <a href="#">Eddington Medal</a> (1975) <a href="#">Maxwell Medal and Prize</a>

# Your next HW assignment

Develop a NER system for Dutch and Spanish.

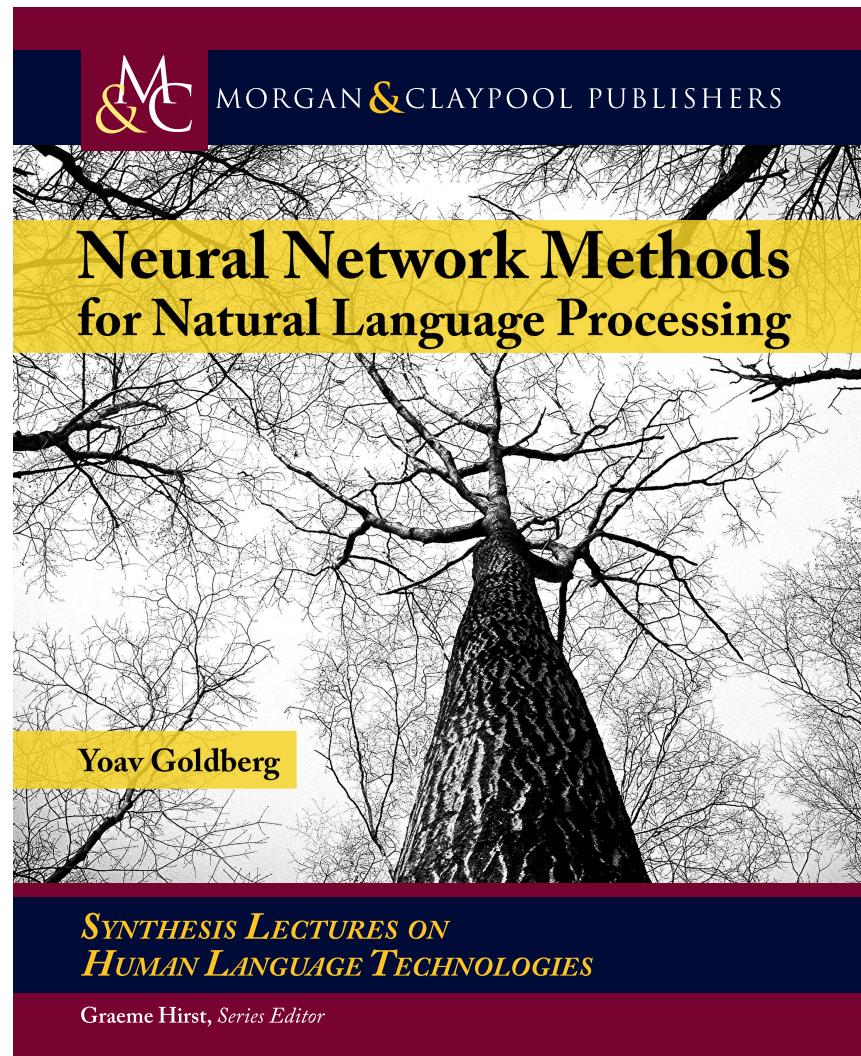
<http://computational-linguistics-class.org/assignment7.html>

# Feed forward vs Recurrent NNs

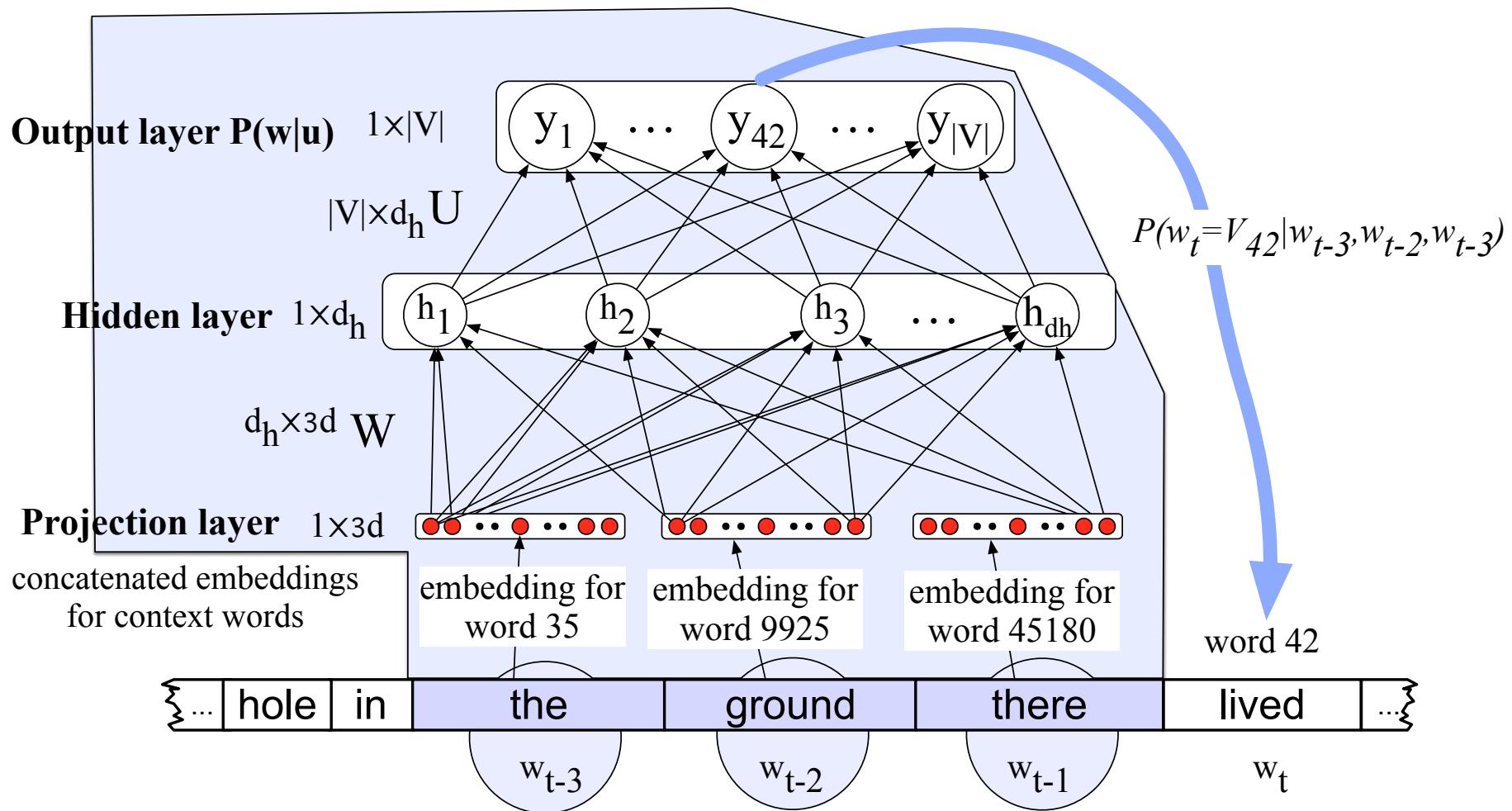
READ CHAPTER 8 AND  
CHAPTER 14 FROM

YOAV GOLDBER'S BOOK  
NEURAL NETWORKS  
METHODS FOR NLP

(IT'S FREE TO DOWNLOAD  
FROM PENN'S CAMPUS!)



# Feed forward neural network



# Fixed inputs

Feed-forward networks assume a fixed dimensional input. For feature-extraction (like in our MEMM), we can extract a fixed number of features, and represent each feature as a vector, and then concatenate all of the vectors. Each region of the resulting input vector corresponds to a different feature.

In some cases the number of features is not known in advance.

# Variable length inputs

If we want to use each word in the sentence as a feature, then we need to represent a variable number of features using a fixed size vector.

One way of achieving this is through a so-called continuous bag of words (CBOW) representation. We discard order information, and either sum or average the embeddings

$$\text{CBOW}(f_1, \dots, f_k) = \frac{1}{k} \sum_{i=1}^k v(f_i)$$

# Are feed forward networks good for sequences?

When dealing with language data, it is very common to work with sequences, such as words (sequences of letters), sentences (sequences of words), and documents.

Feed-forward networks can accommodate sequences through the use of vector concatenation and vector addition (CBOW). CBOW representations allows to encode arbitrary length sequences as fixed sized vectors. However, CBOW is limited, and it disregards order.

# Recurrent neural networks (RNNs)

RNNs can represent any size sequential inputs as fixed-size vectors. They model the structured properties of the input.

RNNs are very powerful at capturing statistical regularities in sequential inputs. They are arguably the strongest contribution of deep-learning to the statistical NLP tool-kit.

There are various types of RNNs: simple RNNs, Long-short-term Memory (LSTM) and the Gated Recurrent Unit (GRU).

# Recurrent neural networks (RNNs)

RNNs allow for language models that do not make the Markov assumption, and condition the next word on the entire sentence history (all the words preceding it).

This ability opens the way to *conditioned generation models*, where a language model that is used as a generator is conditioned on some other signal, such as a sentence in another language.

We'll talk about this application when we come to Neural Machine Translation next week.

# RNN Abstraction

$$y_{1:n} = \text{RNN}^\star(x_{1:n})$$

$$y_i = \text{RNN}(x_{1:i})$$

$$x_i \in \mathbb{R}^{d_{in}} \quad y_i \in \mathbb{R}^{d_{out}}$$

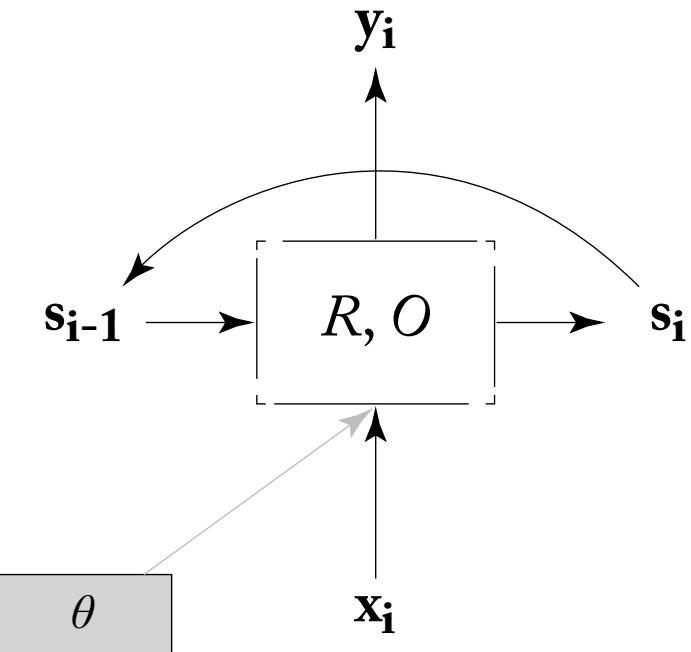
# RNN Abstraction

$$\text{RNN}^{\star}(x_{1:n}; s_0) = y_{1:n}$$

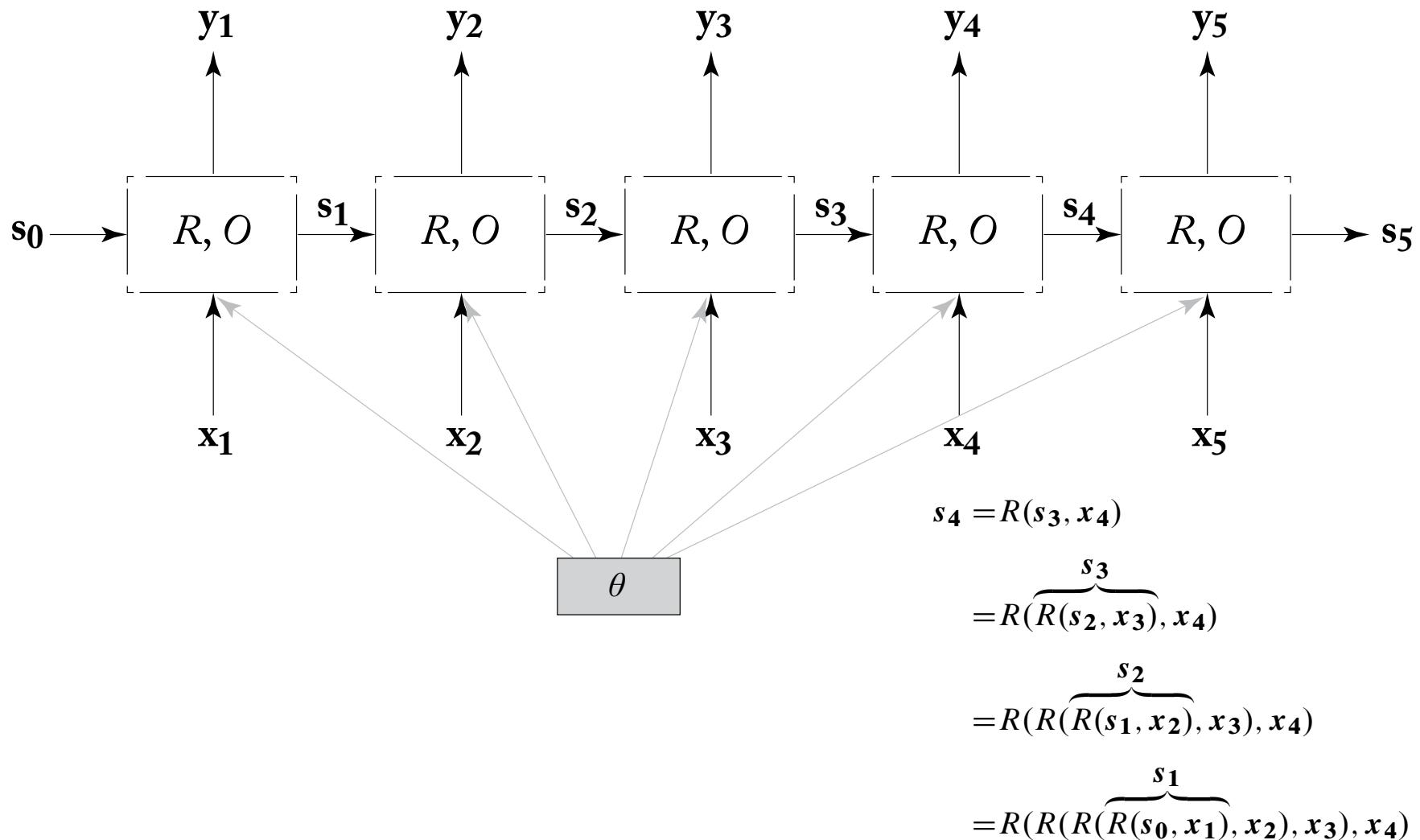
$$y_i = O(s_i)$$

$$s_i = R(s_{i-1}, x_i)$$

$$x_i \in \mathbb{R}^{d_{in}}, \quad y_i \in \mathbb{R}^{d_{out}}, \quad s_i \in \mathbb{R}^{f(d_{out})}$$



# RNN unrolled



# RNN training

To train an RNN network, all we need to do is to create the unrolled computation graph for a given input sequence, add a loss node to the unrolled graph, and then use the backpropogation algorithm to compute the gradients with respect to that loss. For RNNs, This is referred as *backpropagation through time* (BPTT)

To do this, we need a training objective.

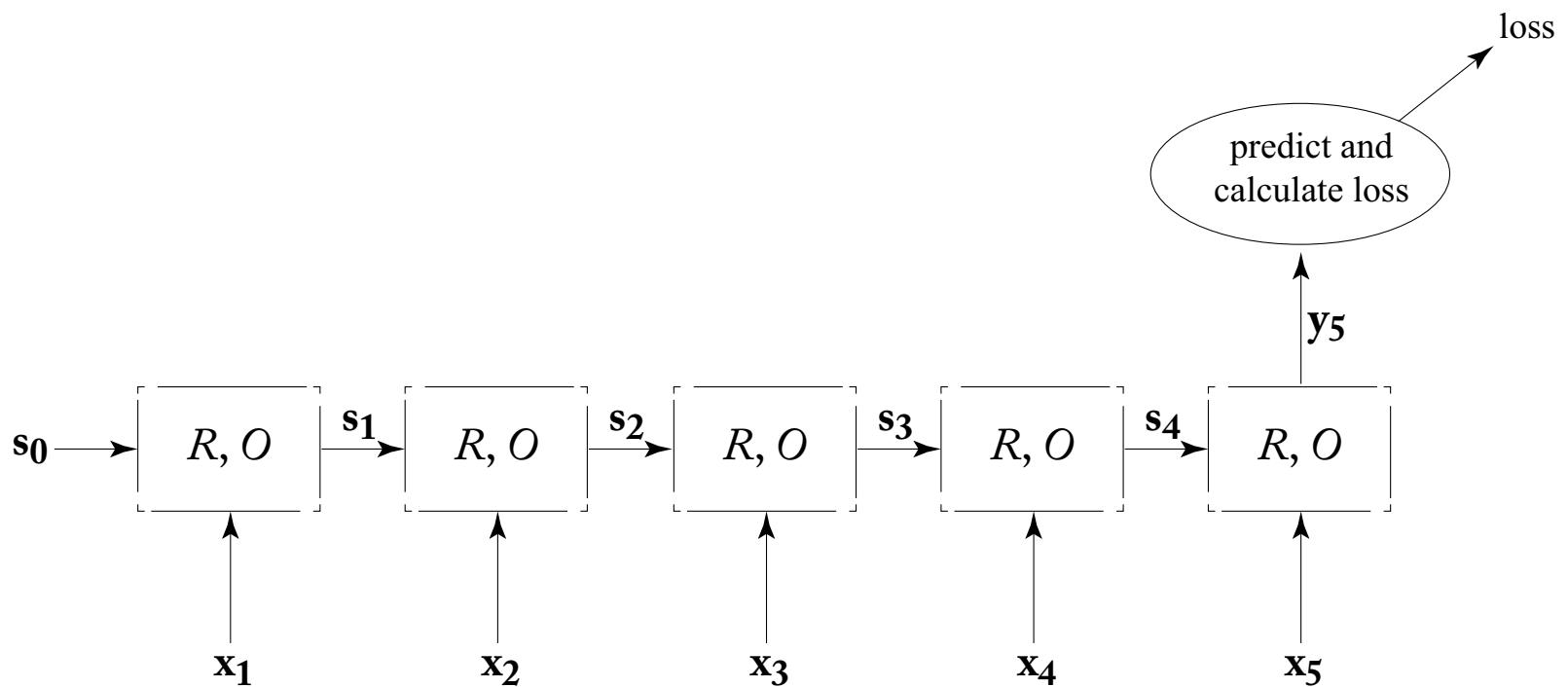
# What is the objective of the training?

The RNN does not do much on its own. Instead, it serves as a trainable component in a larger network.

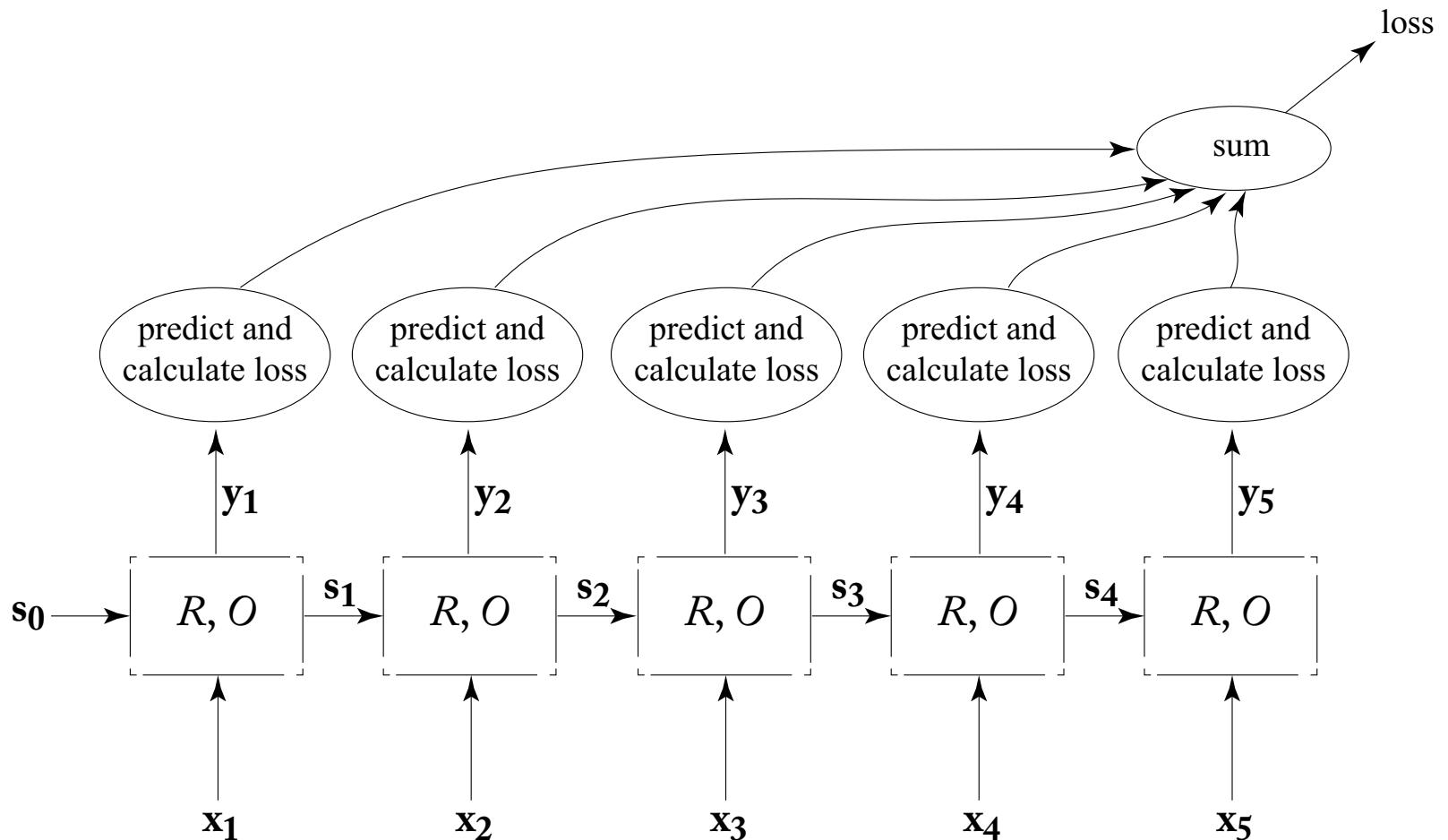
The final prediction and loss computation are performed by that larger network, and the error is back-propagated through the RNN.

This way, the RNN learns to encode properties of the input sequences that are useful for the further prediction task.

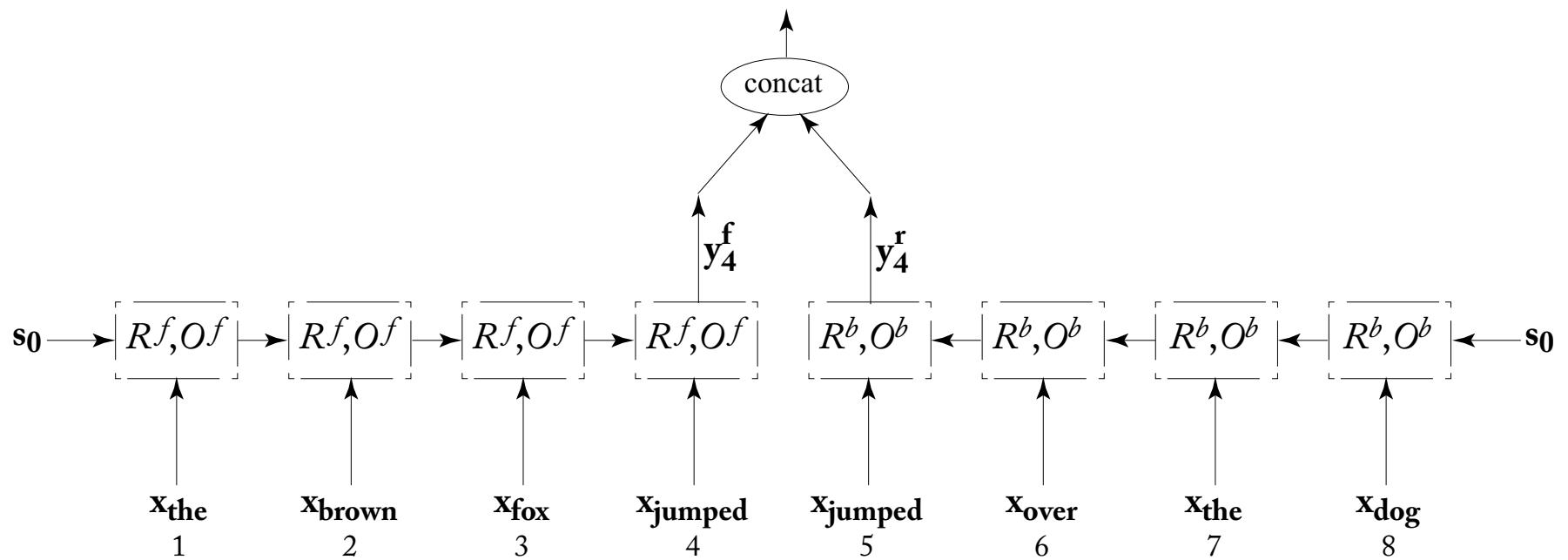
# Acceptor RNN training graph



# Transducer RNN training graph



# Bidirectional RNNs (biRNNs)



# Deep RNNs

