

# IN4080 – 2022 FALL

## NATURAL LANGUAGE PROCESSING

Jan Tore Lønning



## Lecture 13, Nov. 17 (&24)

# This lecture

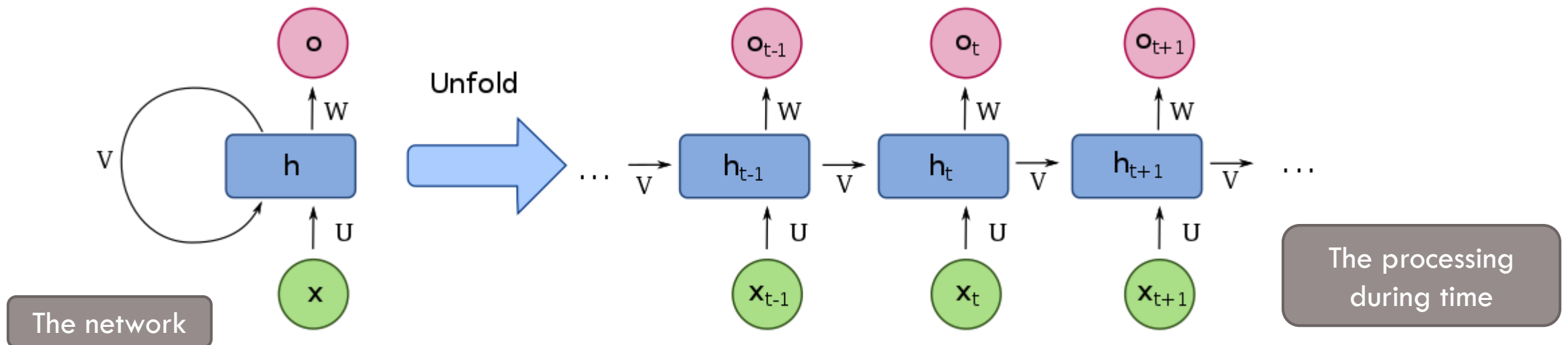
3

- **Recurrent Neural Networks**
  - ▣ RNN Language Models
  - ▣ Other applications of RNNs
  - ▣ Extended architectures and challenges
- Encoder-Decoder models and Machine Translation
- Transformers as Language Models
- Information extraction – some loose ends

# Recurrent neural nets

4

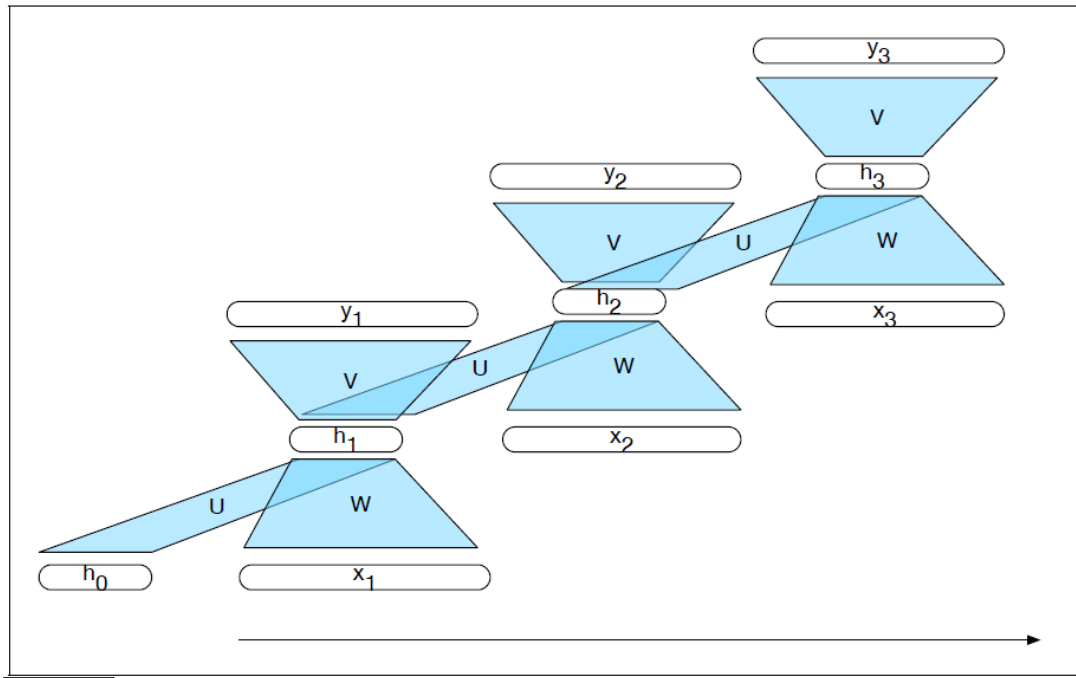
- Model sequences/temporal phenomena
- A cell may send a signal back to itself – at the next moment in time



[https://en.wikipedia.org/wiki/Recurrent\\_neural\\_network](https://en.wikipedia.org/wiki/Recurrent_neural_network)

# Forward

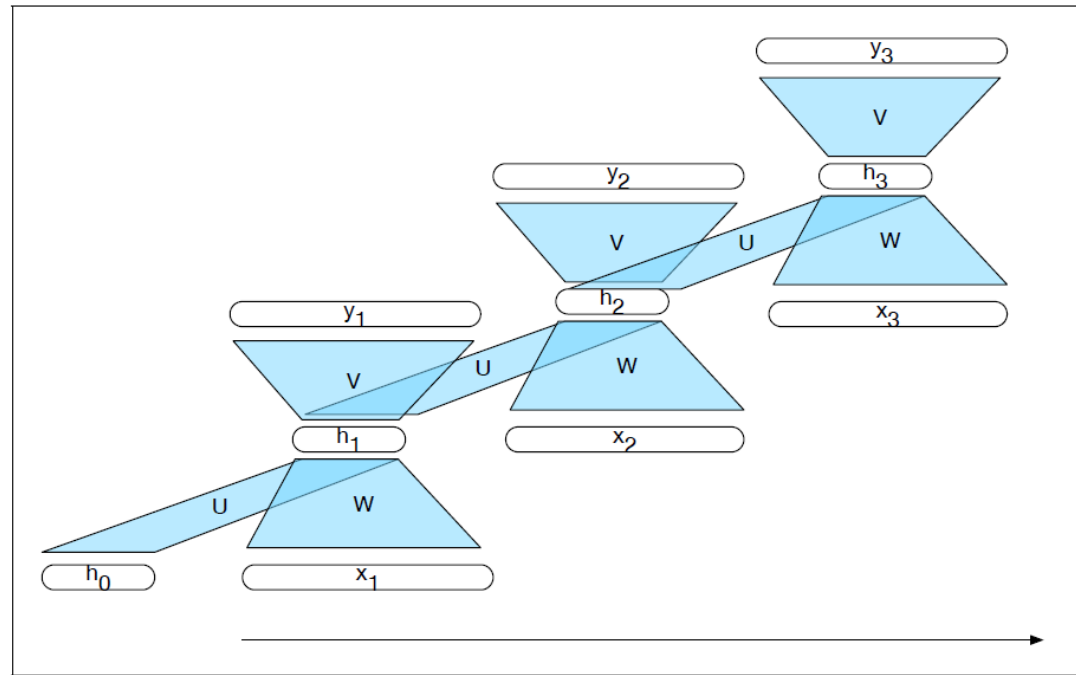
5



- Each  $U, V$  and  $W$  are edges with weights (matrices)
- $x_1, x_2, \dots, x_n$  is the input sequence
- $y_1, \dots, y_n$  is the output sequence
- Forward:
  - ▣ (Initialize  $h_0$ )
  - ▣ For  $i = 1$  to  $n$ :
    - Calculate  $h_i$  from  $h_{i-1}$  and  $x_i$ ,
    - calculate  $y_i$  from  $h_i$

# Forward

6

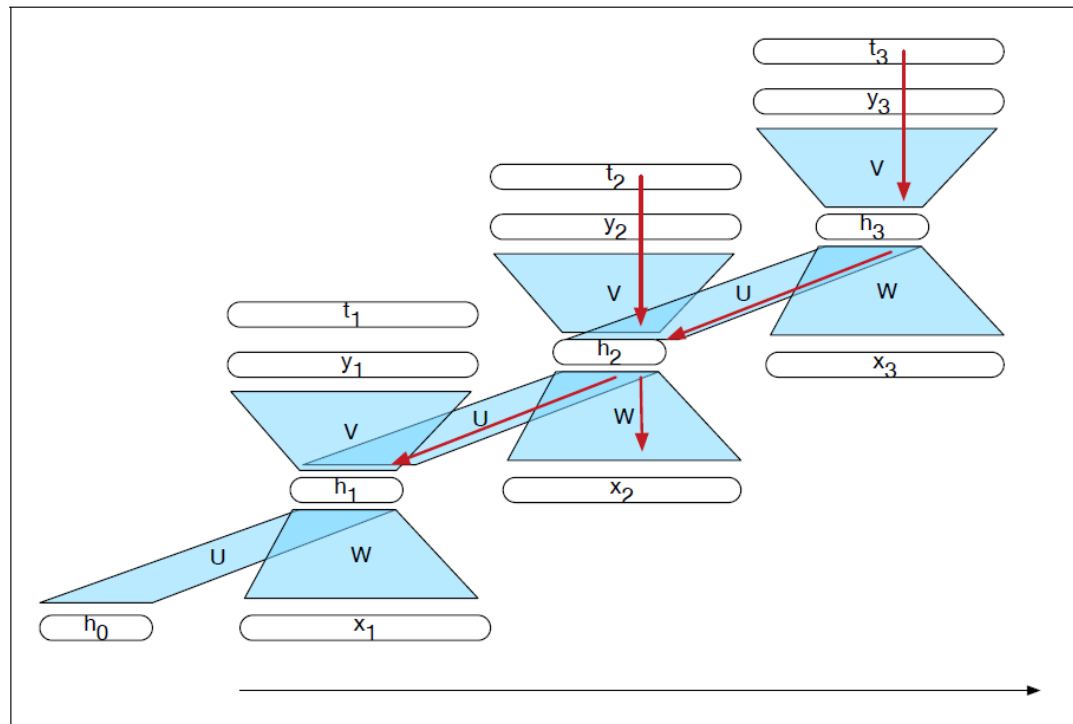


- $\mathbf{h}_t = g(U\mathbf{h}_{t-1} + W\mathbf{x}_t)$
- $\mathbf{y}_t = f(V\mathbf{h}_t)$
- $g$  and  $f$  are activation functions
- (There are also bias which we didn't include in the formulas)

From J&M, 3.ed., 2019

# Training

7



From J&M, 3.ed., 2019

Process one sequence:

- At each output node:
  - ▣ Calculate the loss and the
  - ▣  $\delta$ -term
- Back-propagate the error, e.g.
  - ▣ the  $\delta$ -term at  $h_2$  is calculated
    - ▣ from the  $\delta$ -term at  $h_3$  by U and
    - ▣ the  $\delta$ -term at  $y_2$  by V
- Update
  - ▣ V from the  $\delta$ -terms at the  $y_i$ -s and
  - ▣ U and W from the  $\delta$ -terms at the  $h_i$ -s

# This lecture

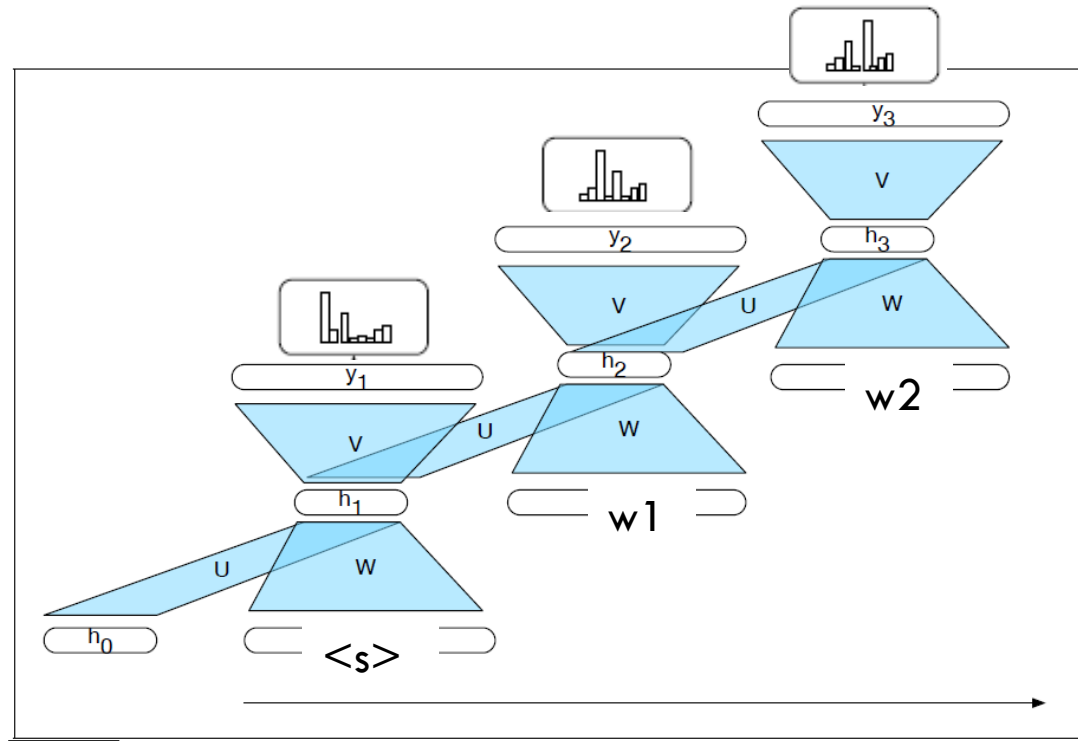
8

- Recurrent Neural Networks
  - ▣ RNN Language Models
  - ▣ Other applications of RNNs
  - ▣ Extended architectures and challenges
- Encoder-Decoder models and Machine Translation
- Transformers as Language Models
- Information extraction – some loose ends



# RNN Language model

9

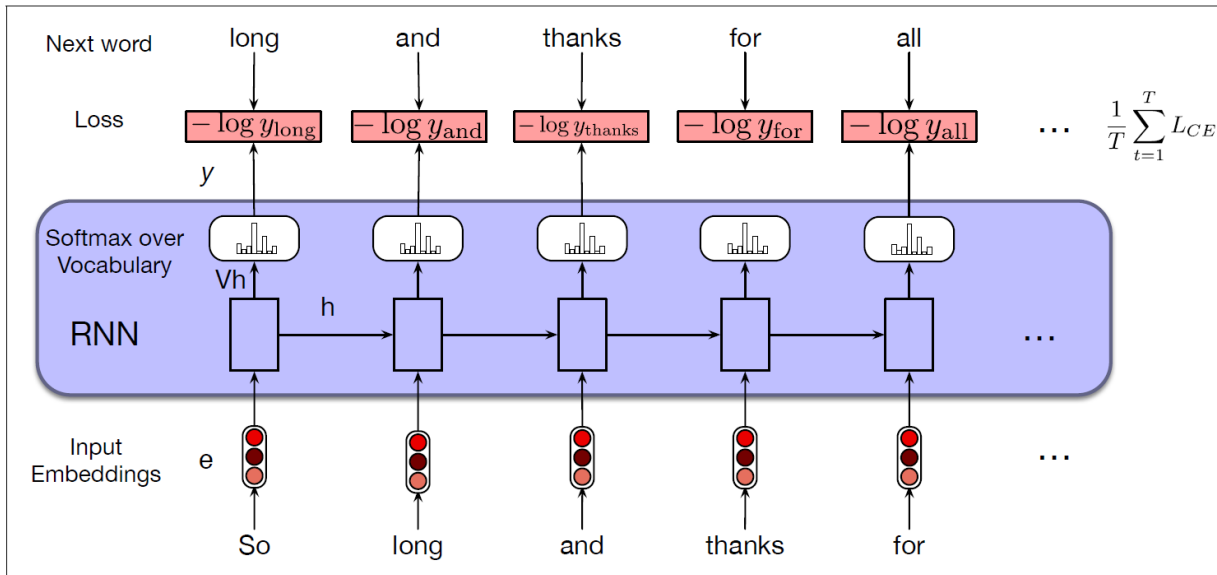


From J&M, 3.ed., 2019

- $\hat{y} = P(w_n | w_1^{n-1}) = \text{softmax}(V\mathbf{h}_n)$
- In principle:
  - ▣ unlimited history
  - ▣ a word depends on all preceding words
- The word  $w_i$  is represented by an embedding
  - ▣ or a one-hot and the embedding is made by the LM

# Training

10



**Figure 9.6** Training RNNs as language models.

From J&M, 3.ed., 2022

- The **predicted** output is a softmax prob.distribution over the vocabulary
- The **target** is the next word
- Cross-entropy loss compare the two
- The errors are back-propagated through the network
- The weights get updated.

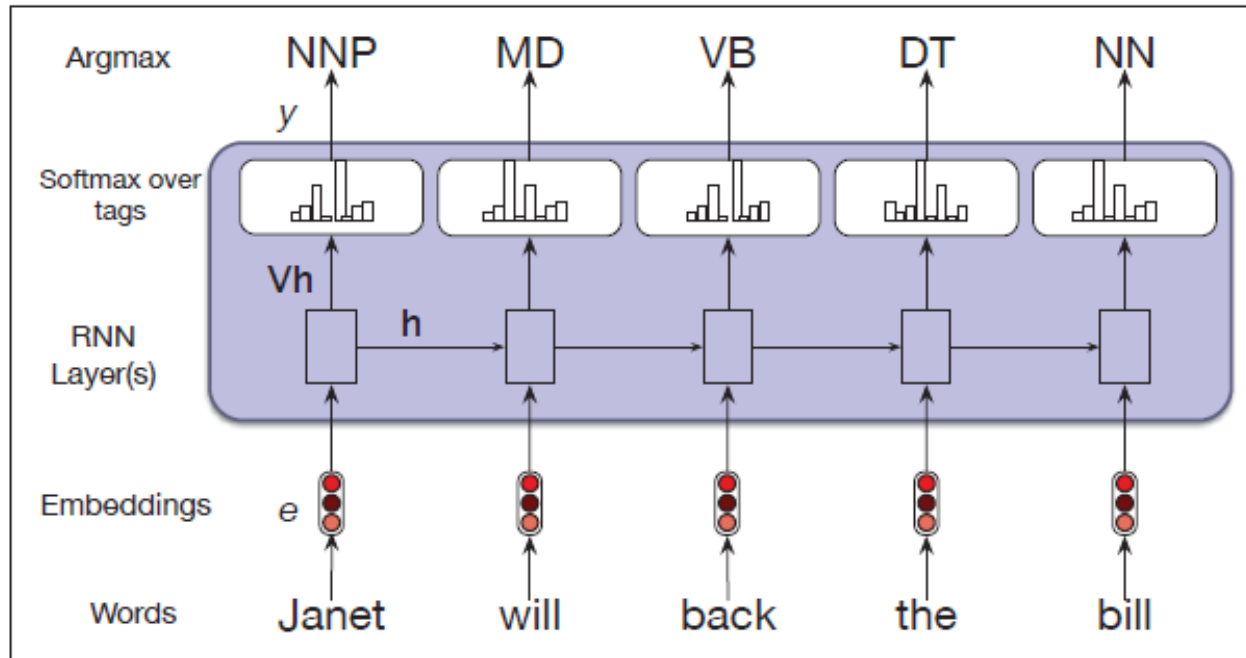
# This lecture

11

- Recurrent Neural Networks
  - RNN Language Models
  - Other applications of RNNs
  - Extended architectures and challenges
- Encoder-Decoder models and Machine Translation
- Transformers as Language Models
- Information extraction – some loose ends

# Neural sequence labeling: e.g., tagging

12

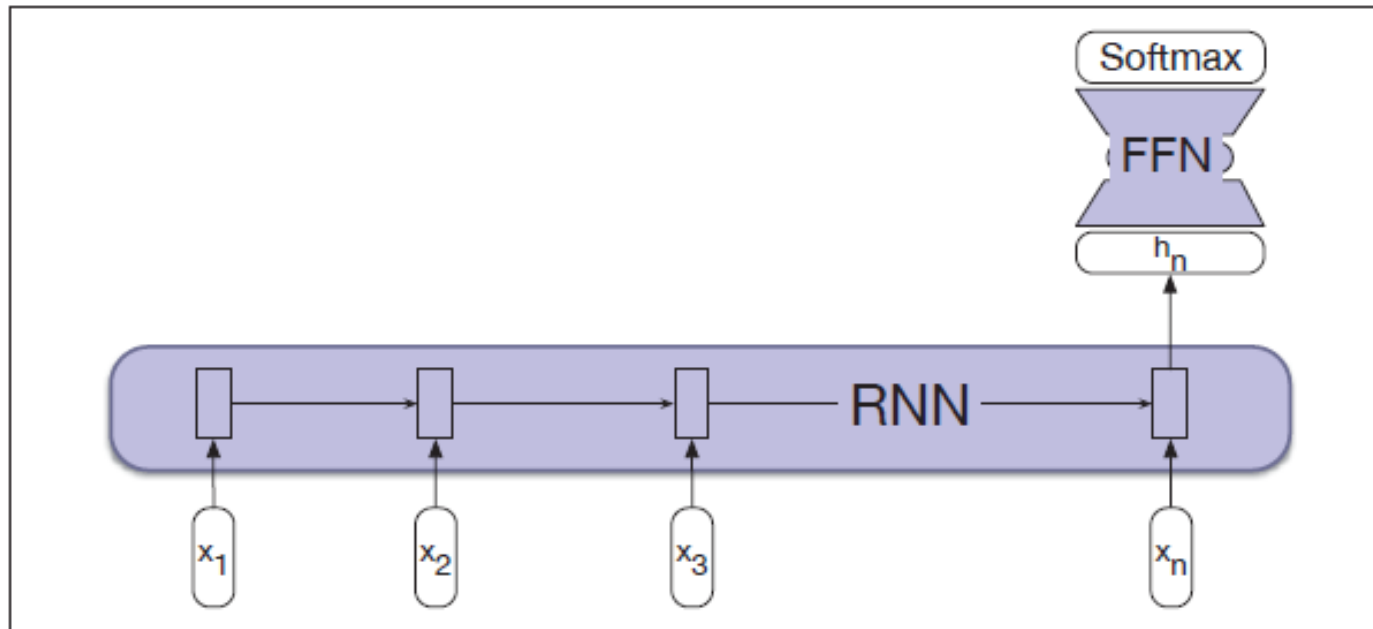


**Figure 9.7** Part-of-speech tagging as sequence labeling with a simple RNN. Pre-trained word embeddings serve as inputs and a softmax layer provides a probability distribution over the part-of-speech tags as output at each time step.

$$\square \hat{y} = P(t_n | w_1^n) = \text{softmax}(V\mathbf{h}_n)$$

# Neural sequence classification

13

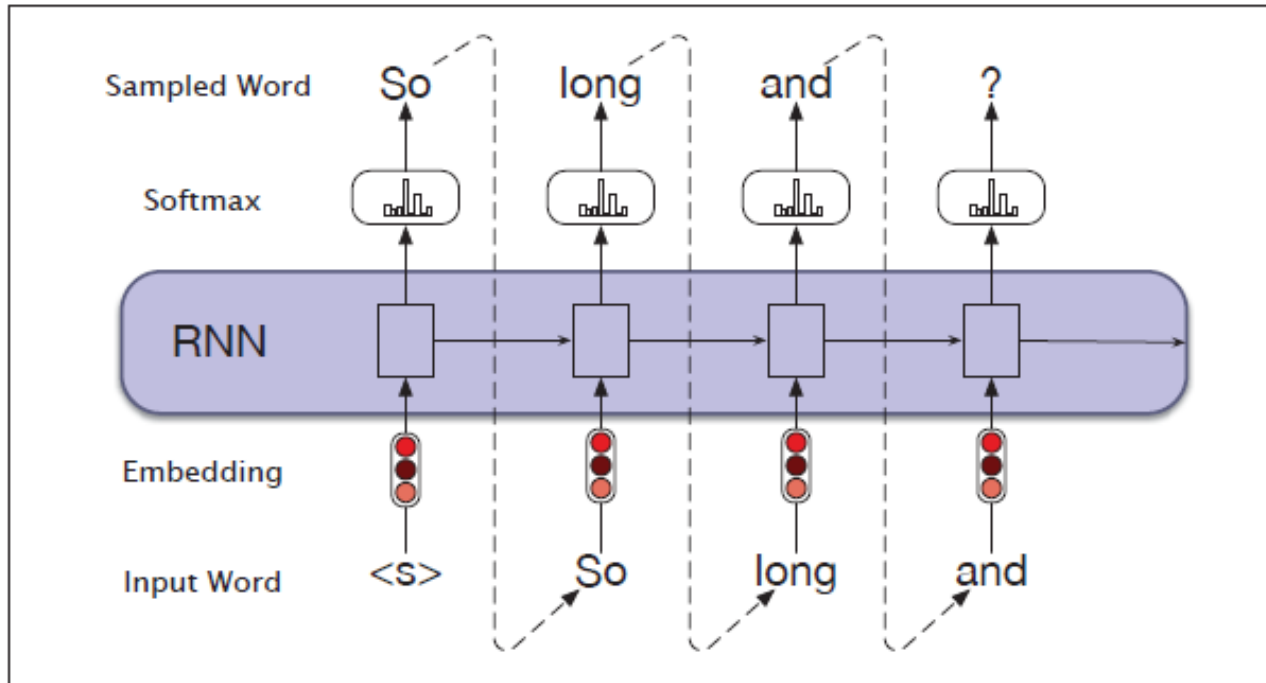


**Figure 9.8** Sequence classification using a simple RNN combined with a feedforward network. The final hidden state from the RNN is used as the input to a feedforward network that performs the classification.

- E.g., sentiment
- Consider only the final state
- Might have a feed-forward network before the softmax.

# Autoregressive generation

14



**Figure 9.9** Autoregressive generation with an RNN-based neural language model.

From J&M, 3.ed., 2021

- From a trained neural LM, generate text
- Initialize with some (or no) words
- The network will at each stage:
  - ▣ Alt. 1: pick the most likely word (argmax)
  - ▣ Alt. 2: sample a word randomly according to the softmax probability

# This lecture

15

- Recurrent Neural Networks
  - ▣ RNN Language Models
  - ▣ Other applications of RNNs
  - ▣ **Extended architectures and challenges**
- Encoder-Decoder models and Machine Translation
- Transformers as Language Models
- Information extraction – some loose ends

# Sequence labeling

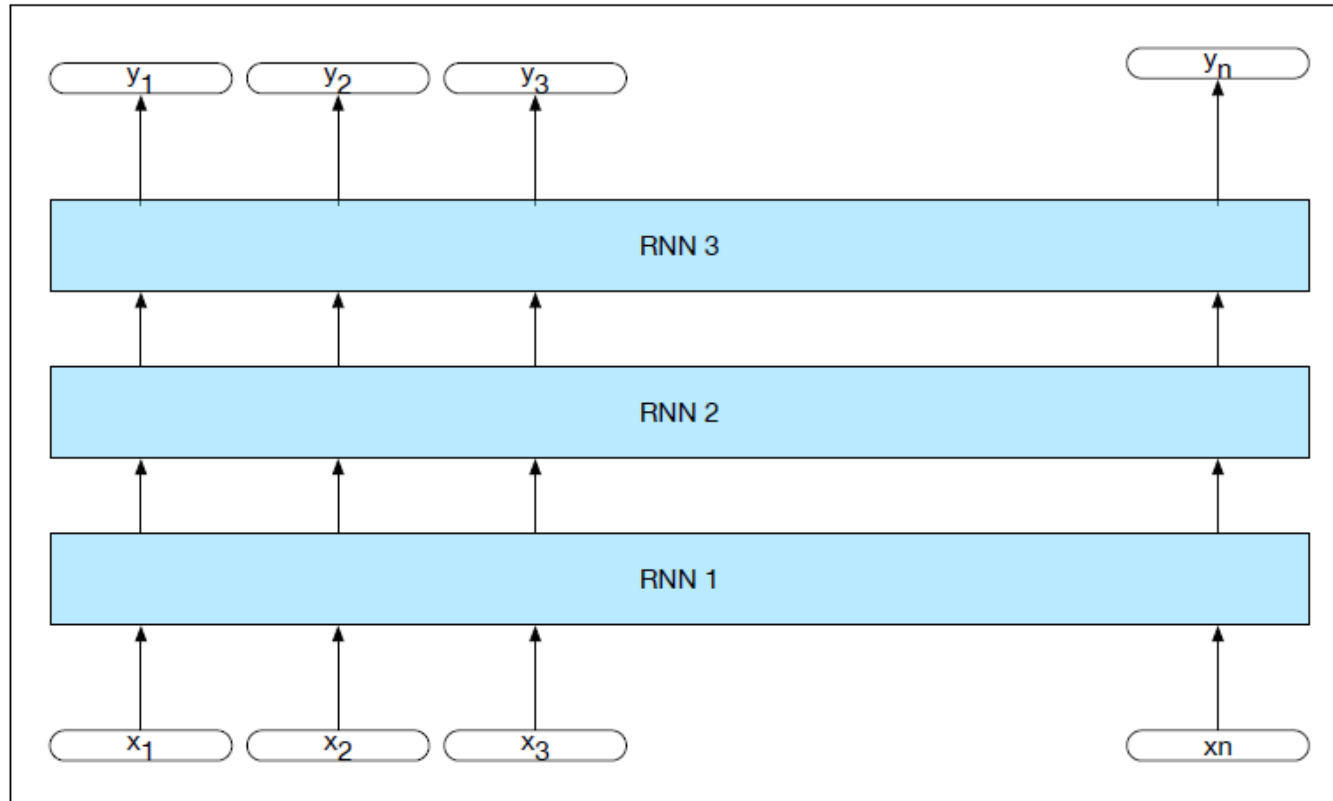
16

- Actual models for sequence labeling, e.g. tagging, are more complex
- For example, it may take words after the tag into consideration.



# Stacked RNN

17



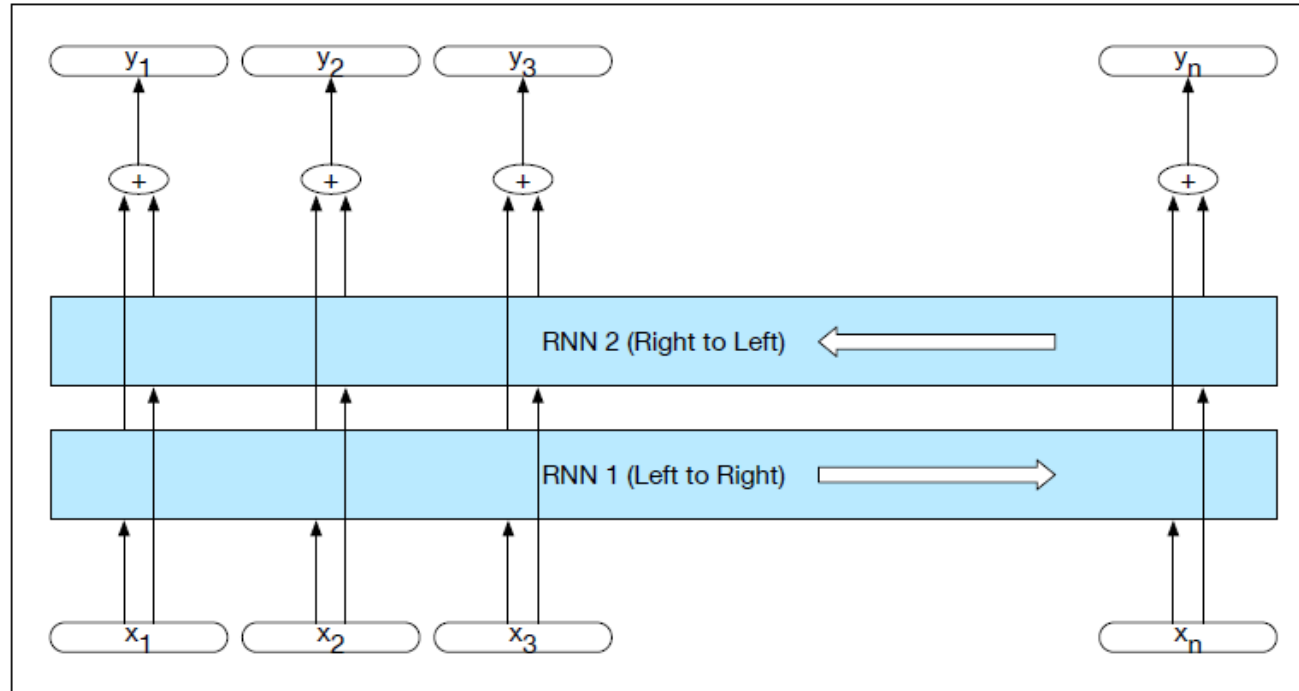
**Figure 9.10** Stacked recurrent networks. The output of a lower level serves as the input to higher levels with the output of the last network serving as the final output.

From J&M, 3.ed., 2019

- Can yield better results than single-layers
- Reason?
  - ▣ Higher-layers of abstraction
    - similar to image processing (convolutional nets)

# Bi-directional RNN

18



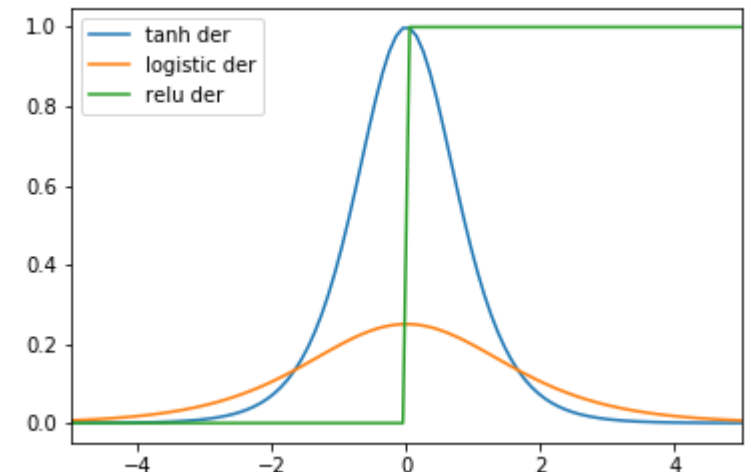
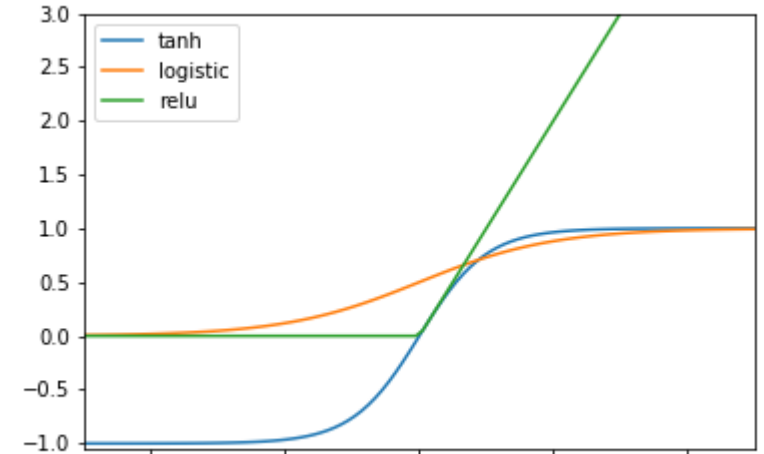
**Figure 9.11** A bidirectional RNN. Separate models are trained in the forward and backward directions with the output of each model at each time point concatenated to represent the state of affairs at that point in time. The box wrapped around the forward and backward network emphasizes the modular nature of this architecture.

- Example: Tagger
- Considers both preceding and following words

# Challenges: Vanishing gradient

19

- A problem for all deep neural networks
- When back-propagating through many layers
  - ▣ the gradient may approach 0
  - ▣ little update
- Partial help:
  - ▣ Other activation functions, e.g. ReLU
  - ▣ Various forms of data normalization
  - ▣ Adjusted architecture



# LSTM

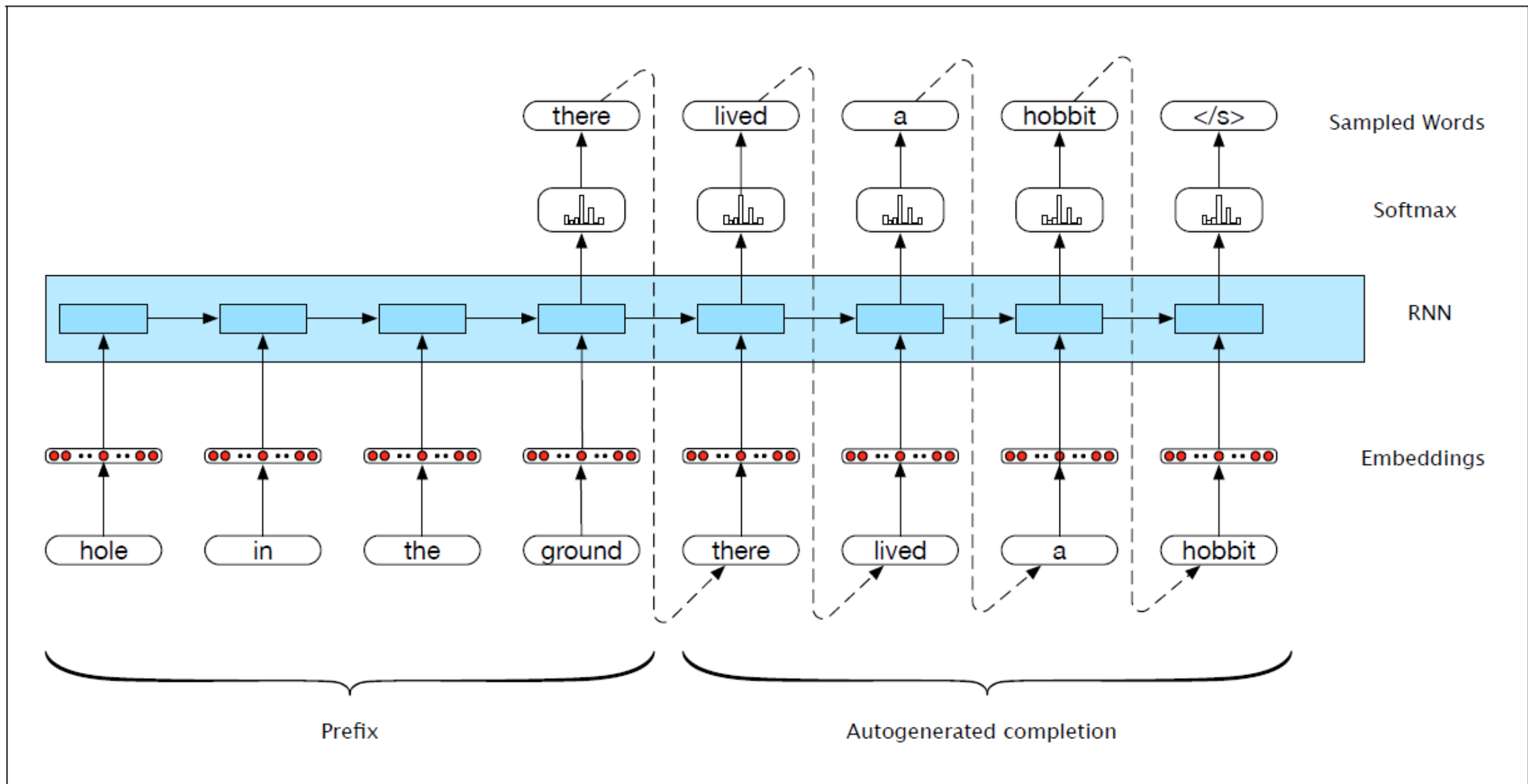
20

- Problems for RNN
  - ▣ Keep track of distant information
- Long Short-Term Memory
  - ▣ An advanced architecture with additional layers and weights
    - Not consider the details here
- Bi-LSTM (Binary LSTM)
  - ▣ Popular standard architecture in NLP

# This lecture

21

- Recurrent Neural Networks
- Encoder-Decoder models and Machine Translation
- Transformers as Language Models
- Information extraction – some loose ends

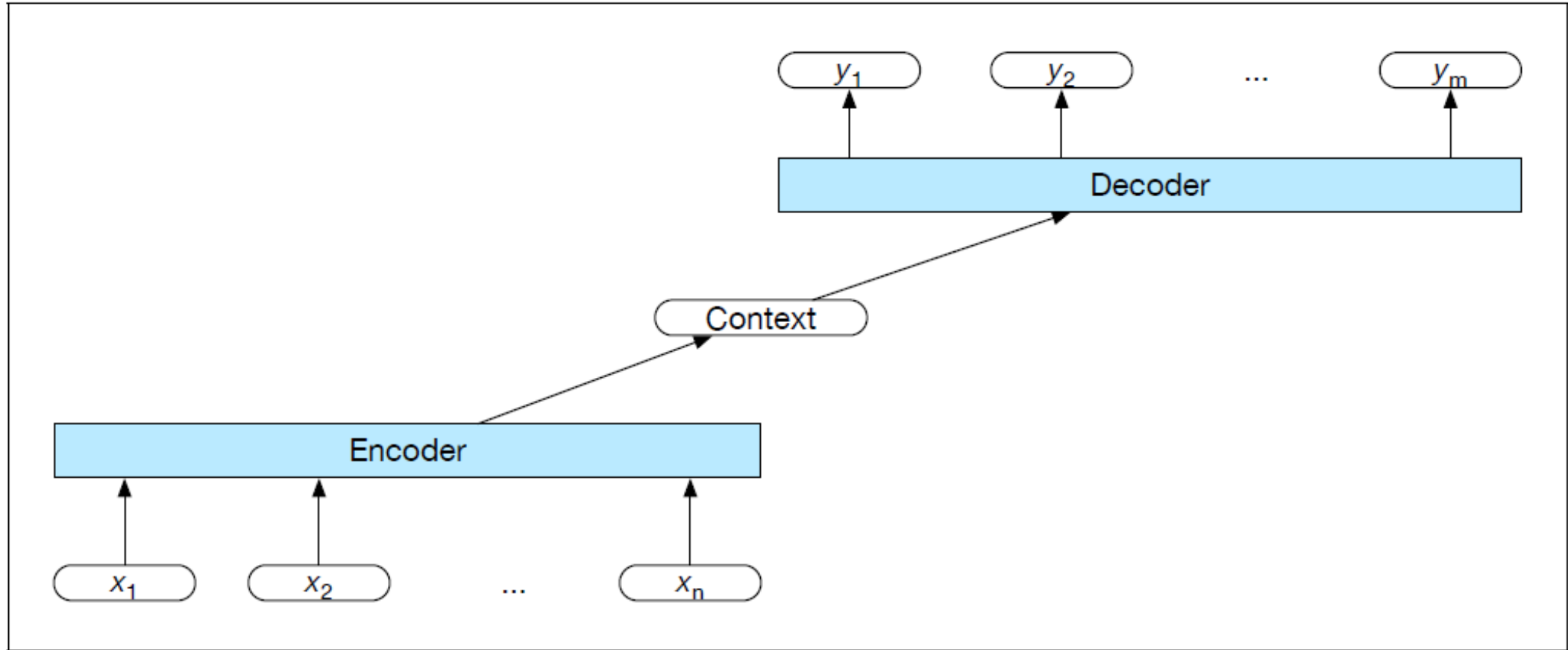


**Figure 10.1** Using an RNN to generate the completion of an input phrase.

# Idea

23

- C.f., the autoregressive generation:
  - ▣ Read-in the first part of the sentence, and
  - ▣ then predict the rest of the sentence
  - ▣ using an RNN trained on sentences
  
- Generalize to other tasks
  - ▣ e.g., Machine Translation



**Figure 10.4** Basic architecture for an abstract encoder-decoder network. The context is a function of the vector of contextualized input representations and may be used by the decoder in a variety of ways.



# Machine Learning-based Machine Translation

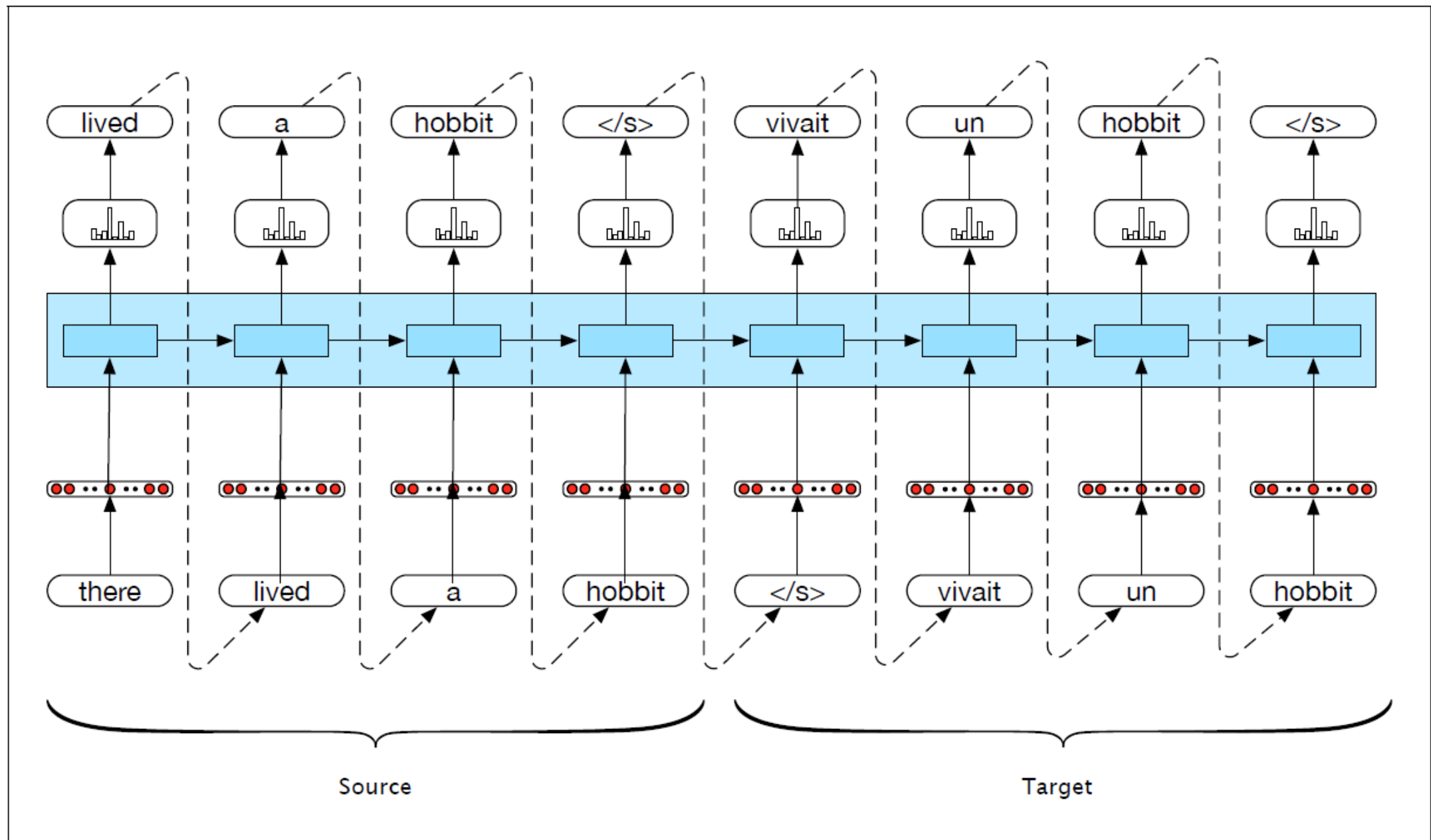
25

- Bi-text
  - ▣ Text translated between two languages
  - ▣ The translated sentences are aligned into sentence pairs
- Machine learning based translation systems are trained on large amounts of bi-text

# Encoder-decoder based translation

26

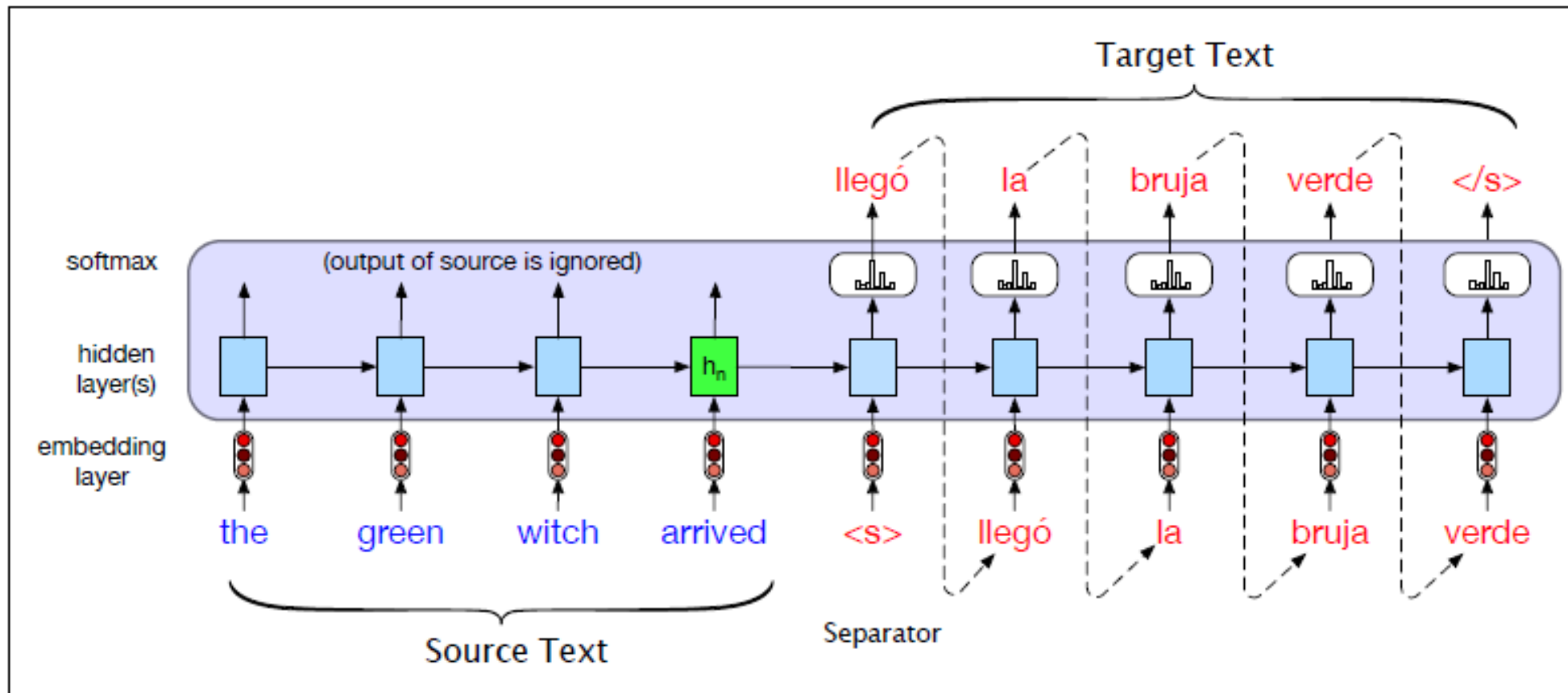
- Concatenate the two sentences in a pair:
  - ▣ source sentence\_<\s>\_target sentence
- Train an RNN on these concatenated pairs
- Apply by reading a source sentence and from there predict a target sentence



**Figure 10.2** Training setup for a neural language model approach to machine translation. Source-target bi-texts are concatenated and used to train a language model.

# Application

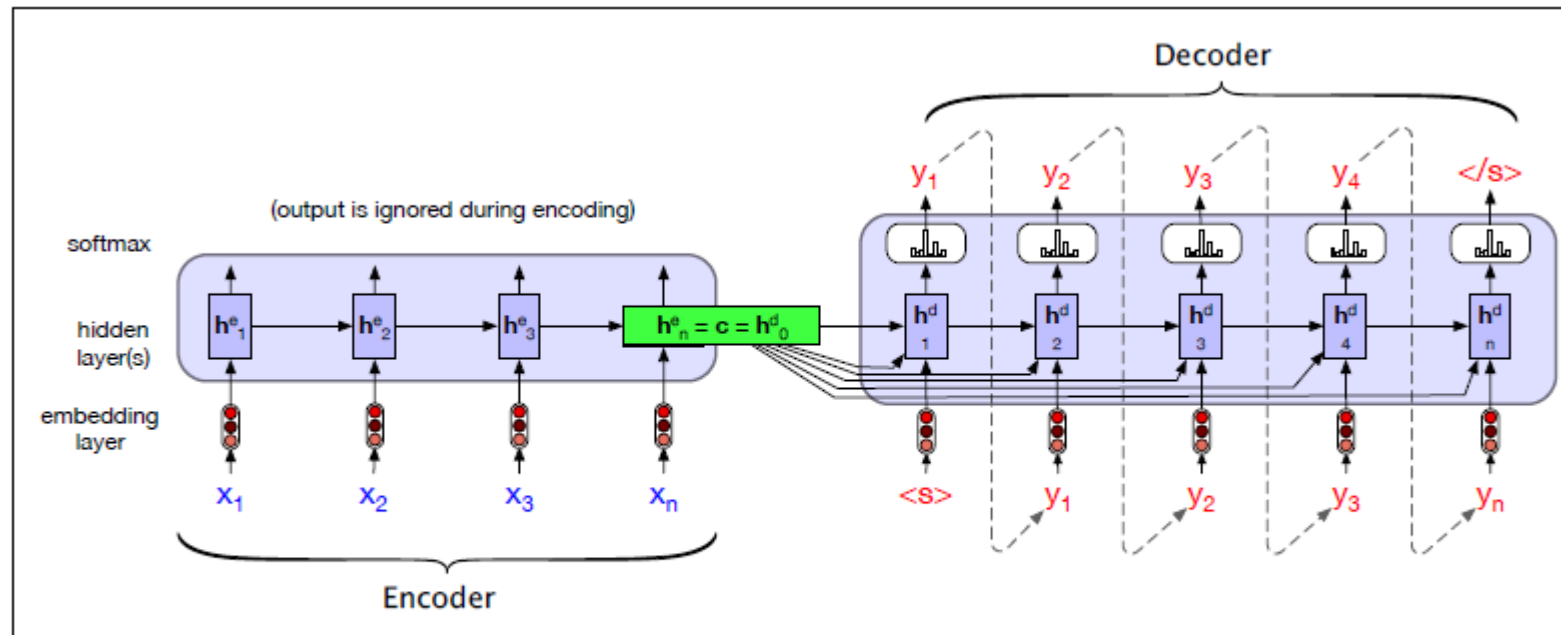
28



**Figure 10.4** Translating a single sentence (inference time) in the basic RNN version of encoder-decoder approach to machine translation. Source and target sentences are concatenated with a separator token in between, and the decoder uses context information from the encoder's last hidden state.

# Refinements

- The encoder can be more refined than a simple RNN,
  - ▣ e.g. bi-LSTM
- The decoder may take more information into consideration:
  - Each output state has access to  $c$



$$h_t^d = g(\hat{y}_{t-1}, h_{t-1}^d, c)$$

**Figure 10.5** A more formal version of translating a sentence at inference time in the basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN,  $h_n^e$ , serves as the context for the decoder in its role as  $h_0^d$  in the decoder RNN.

# This lecture

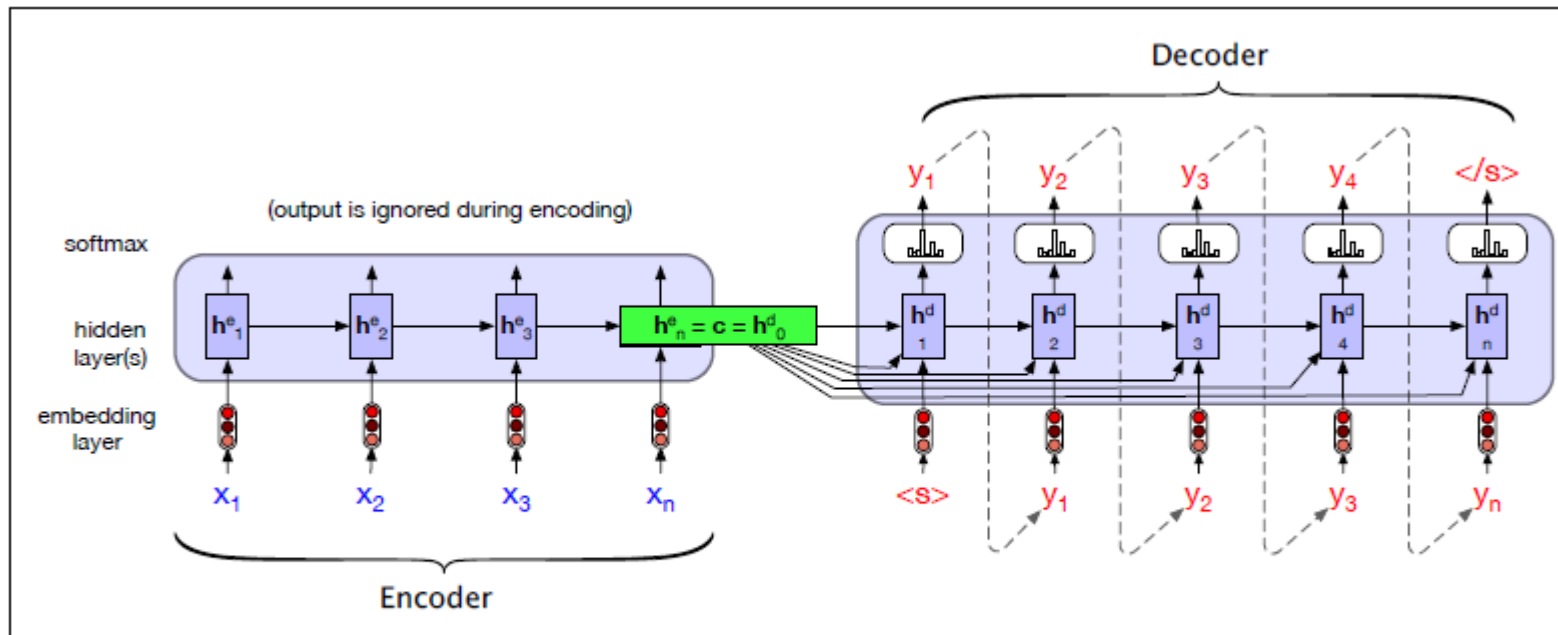
30

- Recurrent Neural Networks
- Encoder-Decoder models and Machine Translation
  - ▣ Attention
- Transformers as Language Models
- Information extraction – some loose ends

# Further refinement

31

- Challenge for one context vector to code a whole sentence.
  - ▣ In particular if the sentence is long

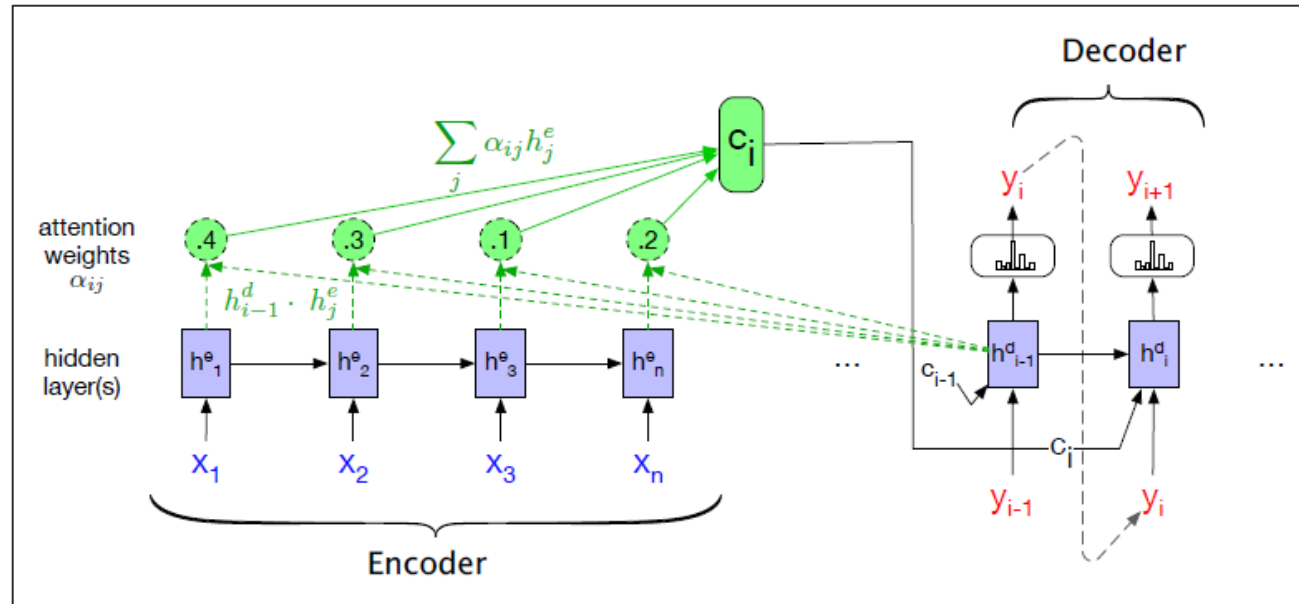


$$h_t^d = g(\hat{y}_{t-1}, h_{t-1}^d, c)$$

**Figure 10.5** A more formal version of translating a sentence at inference time in the basic RNN-based encoder-decoder architecture. The final hidden state of the encoder RNN,  $h_n^e$ , serves as the context for the decoder in its role as  $h_0^d$  in the decoder RNN.

# Attention - sketch

32



**Figure 10.10** A sketch of the encoder-decoder network with attention, focusing on the computation of  $c_i$ . The context value  $c_i$  is one of the inputs to the computation of  $h_i^d$ . It is computed by taking the weighted sum of all the encoder hidden states, each weighted by their dot product with the prior decoder hidden state  $h_{i-1}^d$ .

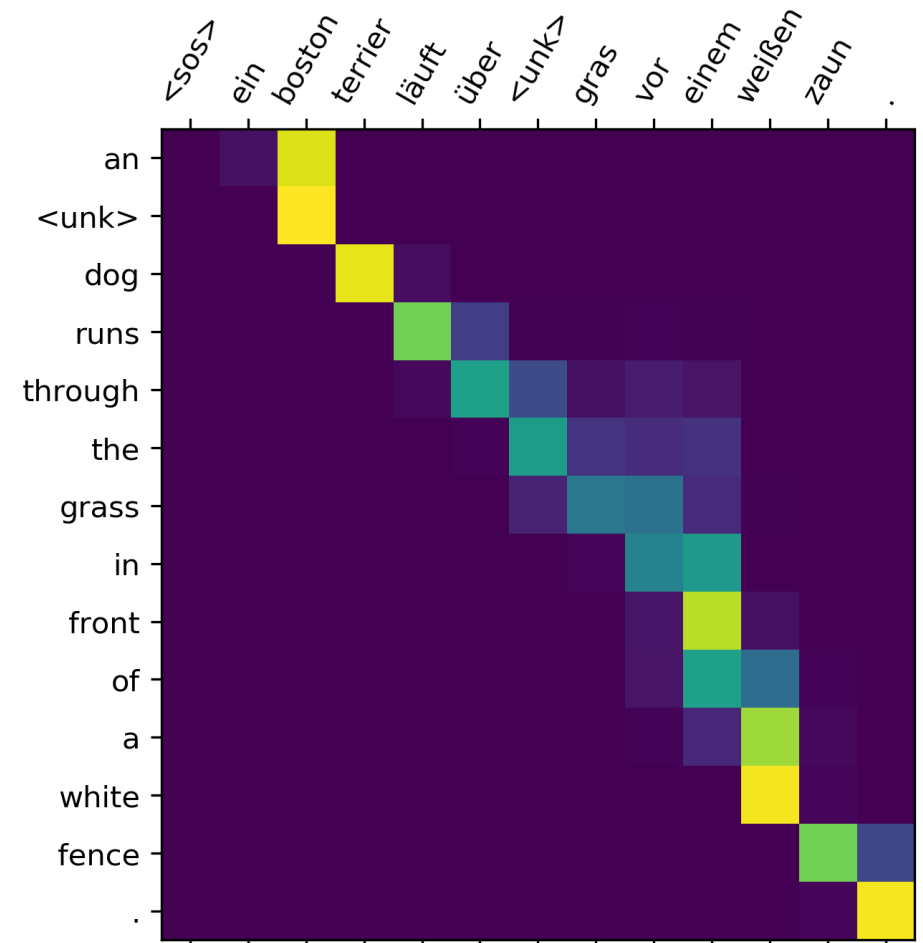
- The context vector is:
  - different for each step:
    - $h_t^d = g(\hat{y}_{t-1}, h_{t-1}^d, c_t)$
  - a weighted sum of input states:
    - $c_t = \sum_{j=1}^n \alpha_{i,j} h_j^e$
  - where the weight  $\alpha_{i,j}$  is determined by  $h_j^e$  and  $h_{t-1}^d$
- "How much attention shall  $h_t^d$  pay to each of the input words?"



# Attention for translation - sketch

33

- $\mathbf{h}_t^d = g(\hat{y}_{t-1}, \mathbf{h}_{t-1}^d, \mathbf{c}_t)$
- $\mathbf{c}_t = \sum_{j=1}^n \alpha_{i,j} \mathbf{h}_j^e$
- "How much attention shall  $\mathbf{h}_t^d$  pay to each of the input words?"
- Which words in the source sentence determine which words in the target sentence?
- $\alpha_{i,j}$  is learned from the examples sentence pairs (as the rest)



# This lecture

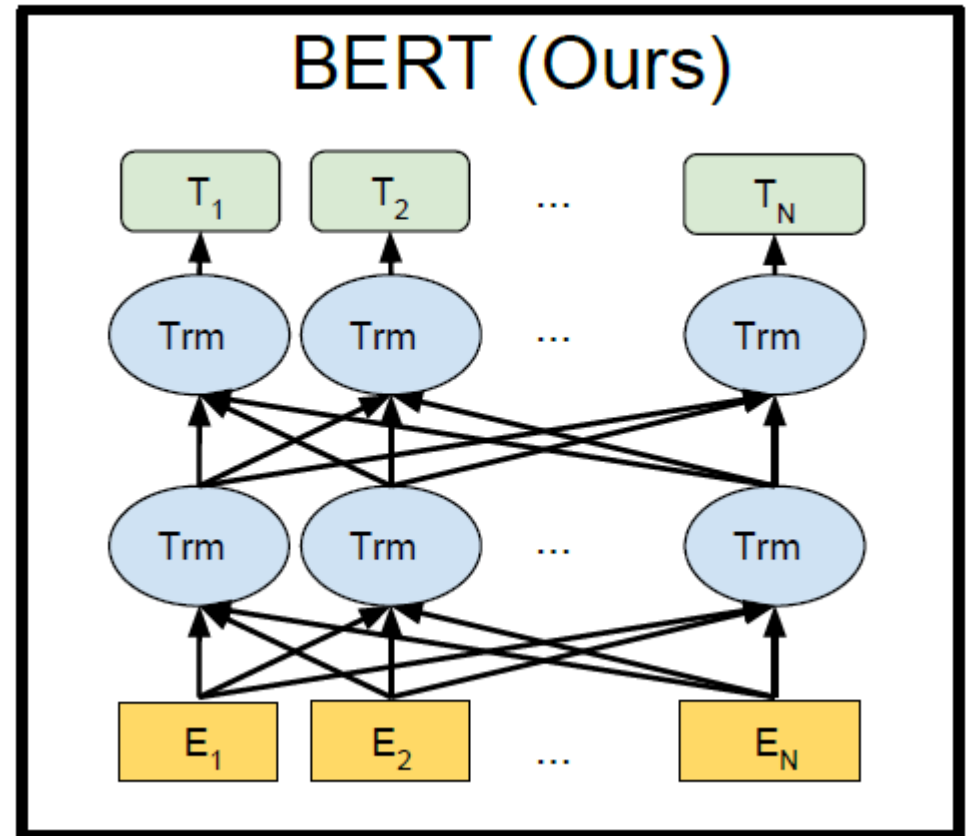
34

- Recurrent Neural Networks
- Encoder-Decoder models and Machine Translation
- **Transformers as Language Models**
- Information extraction – some loose ends

# Language transformers – self-attention (simplified)

35

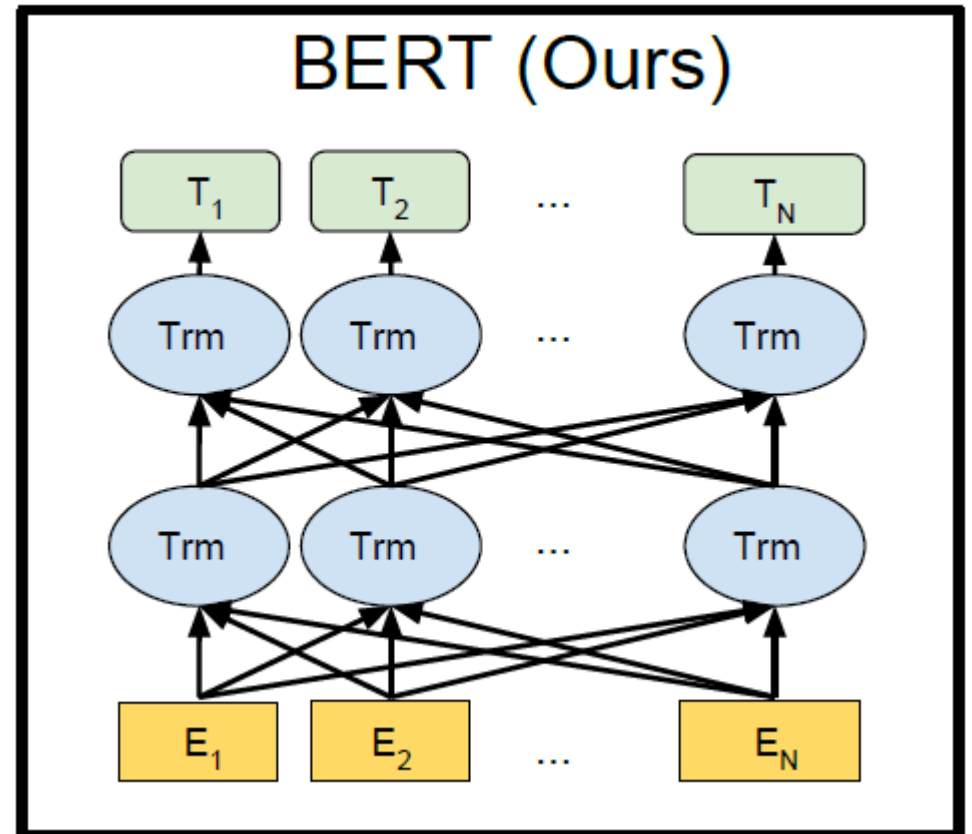
- **Self-attention:** Attention between one word and the other words in the sentence.
- **Transformer:**
  - ▣ A number of layers (e.g., 8)
  - ▣ Each layer:
    - A vector representation for each word in the input
    - Predicted from the representation of the same word in the preceding layer + the attention it receives from the other words



# Language transformers – self-attention (simplified)

36

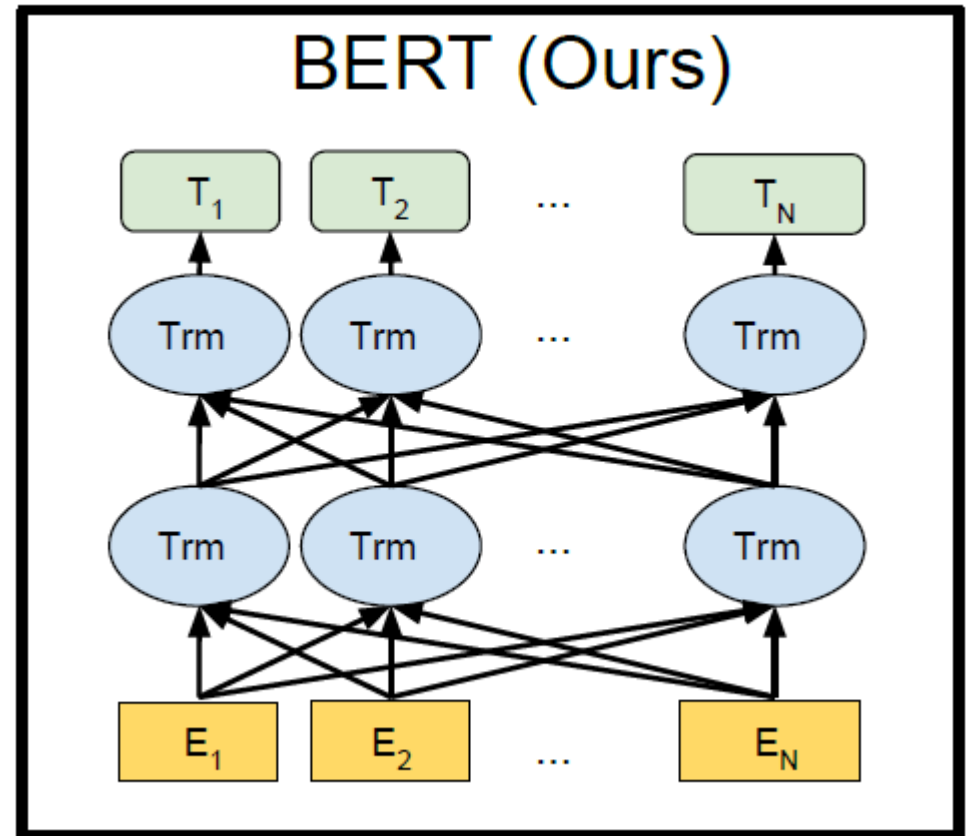
- This results in a vector (embedding) for each word which is context-dependent
- Can be applied to various down-street tasks
- In general, better results than with static embeddings (Word2Vec)



# Language transformers – self-attention (simplified)

37

- Pre-trained language models, e.g., BERT
- Trained on large amounts on data
- Trained on some prediction task, e.g., next word, masked word, next sentence, etc.
- Can then be applied to all kinds of jobs.
- The transformer can be tuned when training the downstream task



# Neural Methods in Natural Language Processing

38

- This completes the introduction to neural NLP
  - ▣ We did not go very far
- You will learn more in
  - ▣ IN5550 – Neural Methods in Natural Language Processing in January
  - ▣ In particular, implementations using HPCs

# This lecture

39

- Recurrent Neural Networks
- Encoder-Decoder models and Machine Translation
- Transformers as Language Models
- **Information extraction – some loose ends**
  - ▣ Repeat the basics of IE and NER
  - ▣ Evaluation of NER
  - ▣ Relation (extraction)
  - ▣ Comparison with the transformer approach

# IE basics

40

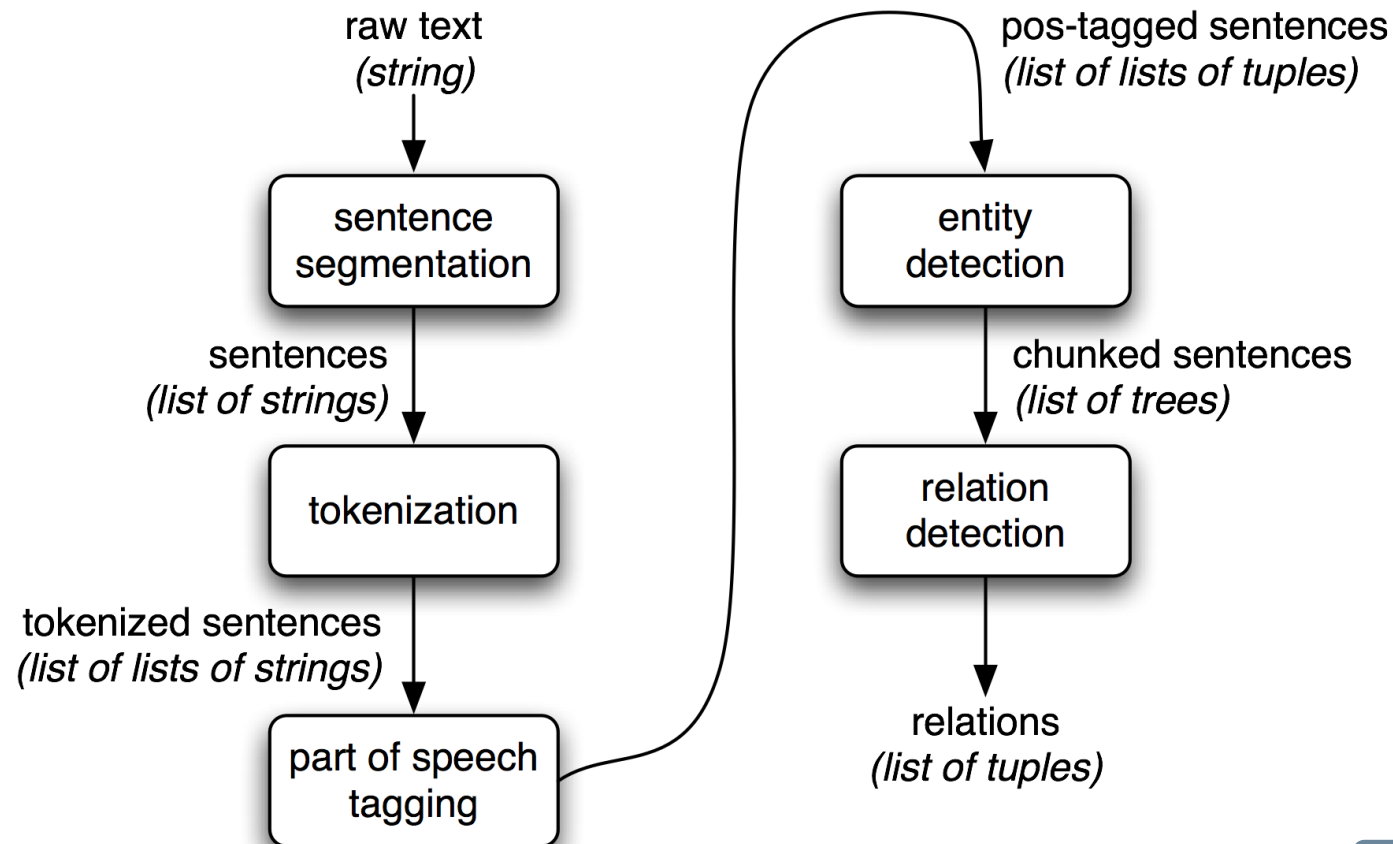
**Information extraction (IE)** is the task of automatically extracting structured information from **unstructured** and/or semi-structured **machine-readable** documents. (Wikipedia)

- ❑ Bottom-Up approach
- ❑ Start with unrestricted texts, and do the best you can
- ❑ The approach was in particular developed by the Message Understanding Conferences (MUC) in the 1990s
- ❑ Select a particular domain and task



# A typical pipeline

41



From NLTK

# Named entities

42

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].

- Named entity:
  - Anything you can refer to by a proper name
- NE Recognition
  - Find the phrases
  - Classify them

# BIO Labels (IOB)

43

- B-PER:
  - ▣ First word in this PER-NE
- I-NP:
  - ▣ Part of PER-NE
- O:
  - ▣ Not part of any NE

Words	BIO Label
Jane	B-PER
Villanueva	I-PER
of	O
United	B-ORG
Airlines	I-ORG
Holding	I-ORG
discussed	O
the	O
Chicago	B-LOC
route	O
.	O

- Can code where something begins and ends without altering the word sequence
- Applying "CONNL-format"
  - ▣ one word per line
  - ▣ info in columns
  - ▣ we may add more columns, e.g. for POS-tag

# Tag accuracy

44

In	IN	0	0
addition	NN	B-NP	B-NP
to	TO	0	0
his	PRP\$	B-NP	B-NP
previous	JJ	I-NP	I-NP
real-estate	NN	I-NP	I-NP
investment	NN	I-NP	I-NP
and	CC	I-NP	I-NP
asset-management	NN	I-NP	I-NP
duties	NNS	I-NP	I-NP
,	,	0	0
Mr.	NNP	B-NP	B-NP
Meador	NNP	I-NP	I-NP
takes	VBZ	0	0
responsibility	NN	B-NP	B-NP
for	IN	0	0
development	NN	B-NP	B-NP
and	CC	0	I-NP
property	NN	B-NP	I-NP
management	NN	I-NP	I-NP
.	.	0	0

- 2 out of 21 tags are incorrect
- Tag-accuracy: 19/21

In	IN	0	0
addition	NN	B-NP	B-NP
to	TO	0	0
his	PRP\$	B-NP	B-NP
previous	JJ	I-NP	I-NP
real-estate	NN	I-NP	I-NP
investment	NN	I-NP	I-NP
and	CC	I-NP	I-NP
asset-management	NN	I-NP	I-NP
duties	NNS	I-NP	I-NP
,	,	0	0
Mr.	NNP	B-NP	B-NP
Meador	NNP	T-NP	T-NP
takes	VBZ	0	0
responsibility	NN	B-NP	B-NP
for	IN	0	0
development	NN	B-NP	B-NP
and	CC	0	I-NP
property	NN	B-NP	I-NP
management	NN	I-NP	I-NP
.	.	0	0

# Counting chunks

- Left column: Gold
- Right column: Predicted

tp: 4   fp: 1   fn: 2

Precision: ?

Recall: ?

# Relation extraction

46

- Extract the relations that exist between the (named) entities in the text
- A fixed set of relations (normally)
  - ▣ Determined by application:
    - Jeopardy
    - Preventing terrorist attacks
    - Detecting illness from medical record
    - ...

- Born\_in
- Date\_of\_birth
- Parent\_of
  
- Author\_of
- Winner\_of
  
- Part\_of
- Located\_in
  
- Acquire
- Threaten
  
- Has\_symptom
- Has\_illness

# Examples

47

Relations	Examples	Types
Affiliations		
Personal	<i>married to, mother of</i>	PER → PER
Organizational	<i>spokesman for, president of</i>	PER → ORG
Artifactual	<i>owns, invented, produces</i>	(PER   ORG) → ART
Geospatial		
Proximity	<i>near, on outskirts</i>	LOC → LOC
Directional	<i>southeast of</i>	LOC → LOC
Part-Of		
Organizational	<i>a unit of, parent of</i>	ORG → ORG
Political	<i>annexed, acquired</i>	GPE → GPE

# Methods for relation extraction

48

1. **Hand-written patterns**
2. Machine Learning (Supervised classifiers)
3. Semi-supervised classifiers via bootstrapping
4. Semi-supervised classifiers via distant supervision
5. Unsupervised



# Different approaches

49

## Tradition pipeline

One step after another

1. Tokenize
2. Tag
3. Maybe some form of parsing
4. NER
5. Relations
6. Application

## Large pre-trained models

1. Tokenize
2. Pre-trained LM
3. Application:
  - ▣ Machine-learning
  - ▣ Fine-tune pre-trained model

# Some example systems

50

- Stanford core nlp (Java): <http://corenlp.run/>
- Stanza (Python): <https://stanfordnlp.github.io/stanza/>
  - ▣ with a wrapper for Stanford Core NLP
- SpaCy (Python): <https://spacy.io/docs/api/>