

# Deep Learning for Natural Language Processing

## Basic Models for Sequence Labeling



UNIVERSITY OF  
GOTHENBURG

---

**CHALMERS**

**WASP** | WALLENBERG AI  
AUTONOMOUS SYSTEMS  
AND SOFTWARE PROGRAM

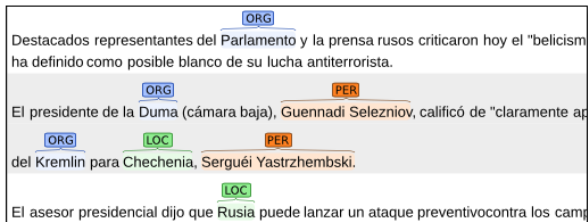
**Richard Johansson**

`richard.johansson@gu.se`

# recap: main type of sequence prediction tasks

- ▶ **sequence labeling** tasks
- ▶ **segmentation** tasks
- ▶ **bracketing** tasks

**Astronomers**    **announce**    **the**    **detection**    **of**    **water**    **in**    **the**    **atmosphere**  
**NNS**            **VBP**            **DT**            **NN**            **IN**            **NN**            **IN**            **DT**            **NN**



## recap: BIO coding

- ▶ for segmentation and bracketing tasks, we often convert into sequence labeling by applying BIO coding or similar

The	cases	of	metastatic cancer	of	the	gall	bladder
			Pathology				Organ
0	0	0	B-PAT	I-PAT	0	0	B-ORG I-ORG

# structured prediction: basic terminology

- ▶ sequence labeling is a **structured prediction** task
  - ▶ the output is a complex object, not just a category
  - ▶ later in the course: trees, graphs, ...
- ▶ input: a sequence  $\mathbf{x}$
- ▶ output: a sequence  $\mathbf{y}$  of the same length as  $\mathbf{x}$

$$\begin{array}{ccccccc} w_1 & w_2 & w_3 & \cdot & \cdot & \cdot & w_n \\ \downarrow & \downarrow & \downarrow & & & & \downarrow \\ l_1 & l_2 & l_3 & & & & l_n \end{array}$$

# Algorithmic approaches

- **Exhaustive search**

Cast structured prediction as a combinatorial optimisation problem over the set of target representations.

Viterbi algorithm, Eisner algorithm

- **Greedy search**

Cast structured prediction as a sequence of classification problems: at each point in time, predict one of several options.

window-based part-of-speech tagging, arc-standard algorithm


# Algorithmic approaches

- **Exhaustive search**

Cast structured prediction as a combinatorial optimisation problem over the set of target representations.

Viterbi algorithm, Eisner algorithm

- **Greedy search**

 Cast structured prediction as a sequence of classification problems: at each point in time, predict one of several options.

window-based part-of-speech tagging, arc-standard algorithm

## simplest approach: word-level classification

- ▶ what if we train a model that classifies each word independently of other words?

<b>She</b>	<b>plays</b>	<b>in</b>	<b>many</b>	<b>plays</b>
pronoun	verb	preposition	adverb	noun

<b>Manchester</b>	<b>United</b>	<b>will</b>	<b>return</b>	<b>to</b>	<b>the</b>	<b>United</b>	<b>States</b>
<b>B-ORG</b>	<b>I-ORG</b>	<b>O</b>	<b>O</b>	<b>O</b>	<b>O</b>	<b>B-LOC</b>	<b>I-LOC</b>

better: window-based classification

**Manchester    United    will**



better: window-based classification

Manchester United will



**I-ORG**

better: window-based classification

Manchester United will

the United States



**I-ORG**

better: window-based classification

Manchester United will



**I-ORG**

the United States

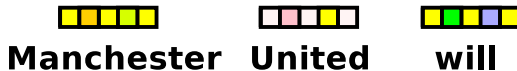


**B-LOC**

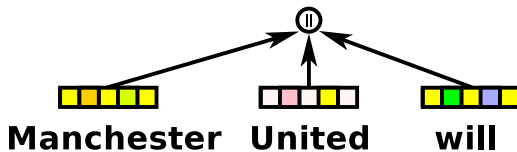
# window-based classification: implementation

**Manchester United will**

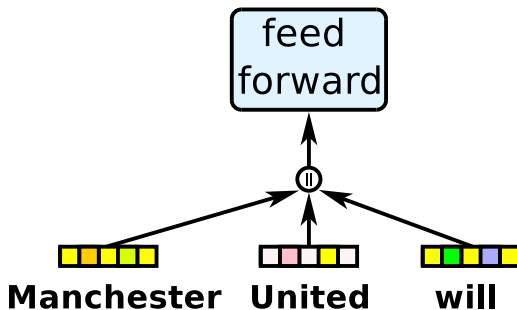
# window-based classification: implementation



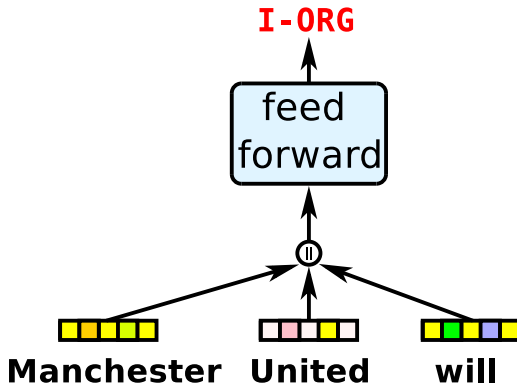
## window-based classification: implementation



## window-based classification: implementation



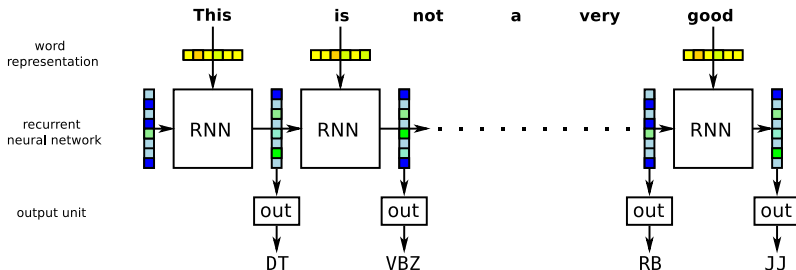
## window-based classification: implementation



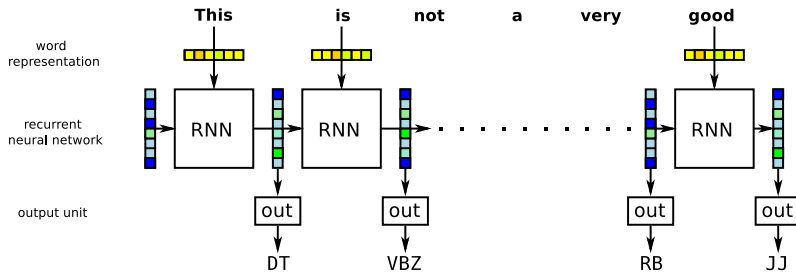


# RNN-based sequence labeling

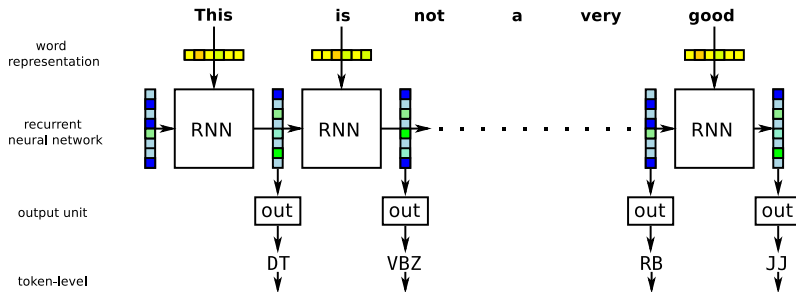
- ▶ we have already seen different types of **recurrent neural networks**, which operate on sequences
- ▶ it is straightforward to apply an output layer on top of an RNN



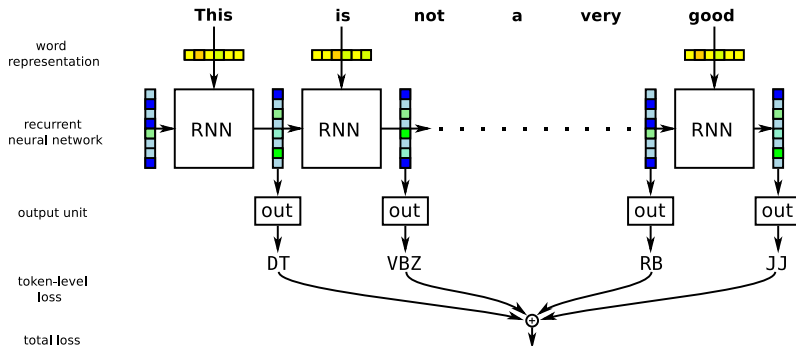
# computing the loss in an RNN-based sequence labeler



# computing the loss in an RNN-based sequence labeler

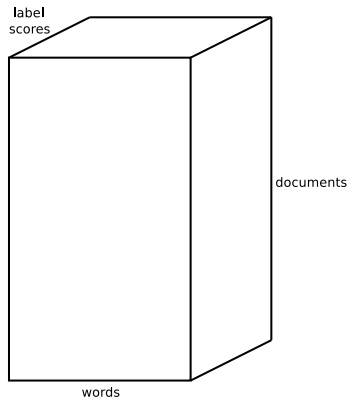


# computing the loss in an RNN-based sequence labeler



# a note about working with sequences in PyTorch

- ▶ the final output is a 3-dimensional tensor of shape  $(n\_docs, n\_words, n\_labels)$
- ▶ we use **padding** to make sure all documents in a batch have the same length
- ▶ if we don't want to compute the loss for the padded dummy tokens, we can use a **mask**



# padding and masking example

- ▶ label sequences:

0	0		
B-PER	0	0	
0	B-LOC	I-LOC	0

- ▶ encoded and padded:

1	1	0	0
2	1	1	0
1	4	5	1

- ▶ loss mask:

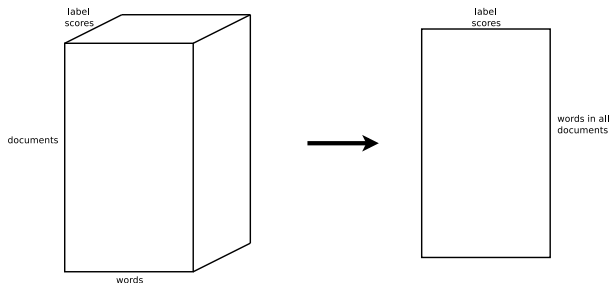
1	1	0	0
1	1	1	0
1	1	1	1

## easier solution

- ▶ PyTorch provides automatic masking for some loss functions:

```
loss = nn.CrossEntropyLoss(ignore_index=pad_label_id)
```

- ▶ note that most loss functions expect 2-dimensional tensors, so we need to reshape before computing the loss



# exercise 1

- ▶ simple models for **named entity recognition**

Manchester	United	will	return	to	the	United	States
<b>B-ORG</b>	<b>I-ORG</b>	<b>O</b>	<b>O</b>	<b>O</b>	<b>O</b>	<b>B-LOC</b>	<b>I-LOC</b>

- ▶ we will investigate more complex models in a second exercise