

# Natural Language Processing with Deep Learning

RUHR  
UNIVERSITÄT  
BOCHUM

RUB

## Lecture 7 — Recurrent neural networks and encoder-decoder architectures

---

Prof. Dr. Ivan Habernal

December 5, 2024

[www.trusthlt.org](http://www.trusthlt.org)

Trustworthy Human Language Technologies Group (TrustHLT)

Ruhr University Bochum & Research Center Trustworthy Data Science and Security



CENTER FOR TRUSTWORTHY  
DATA SCIENCE AND SECURITY

# Motivation

Language data – working with sequences (of tokens, characters, etc.)

MLP – fixed input vector size

# Motivation

Language data – working with sequences (of tokens, characters, etc.)

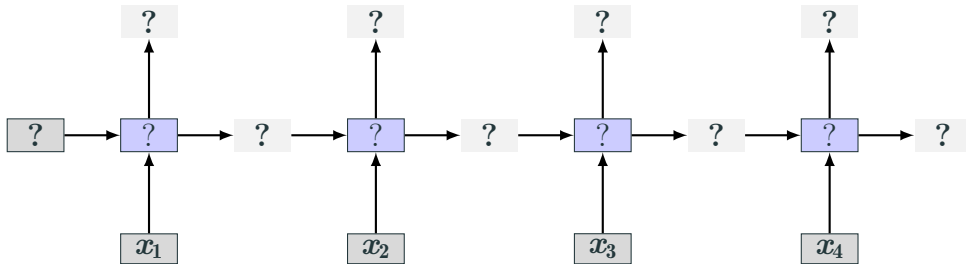
MLP – fixed input vector size

How we dealt with it

- Vector concatenation
- Vector addition/averaging (CBOW)
- Limiting context (e.g., Markov property)

What we want to really work with: Sequence of inputs, fixed-size output(s)

# Our goal would be to build something like this



Example for 4 input tokens

# Recurrent Neural Networks (RNN) abstraction

---

- 1 Recurrent Neural Networks (RNN) abstraction
- 2 RNN architectures
- 3 Encoder-decoder architectures

# RNN abstraction

We have a sequence of  $n$  **input** vectors  $\mathbf{x}_{1:n} = \mathbf{x}_1, \dots, \mathbf{x}_n$

Each input vector has the same dimension  $d_{in} : \mathbf{x}_i \in \mathbb{R}^{d_{in}}$

What might  $\mathbf{x}_i$  contain?

- Typically a word embedding of token  $i$ , but could be any arbitrary input, e.g., one-hot encoding of token  $i$

# RNN abstraction

We have a sequence of  $n$  **input** vectors  $\mathbf{x}_{1:n} = \mathbf{x}_1, \dots, \mathbf{x}_n$

Each input vector has the same dimension  $d_{in} : \mathbf{x}_i \in \mathbb{R}^{d_{in}}$

What might  $\mathbf{x}_i$  contain?

- Typically a word embedding of token  $i$ , but could be any arbitrary input, e.g., one-hot encoding of token  $i$

We have a single **output**  $d_{out}$ -dimensional vector  $\mathbf{y}_n \in \mathbb{R}^{d_{out}}$

RNN is a function from input to output

$$\mathbf{y}_n = \text{RNN}(\mathbf{x}_{1:n})$$

# RNN in fact returns a sequence of outputs

RNN definition:  $y_n = \text{RNN}(x_{1:n})$

Let's have  $n = 3$ , so our input sequence is  $x_1, x_2, x_3$ :

$$y_2 = \text{RNN}(x_1, x_2, x_3)$$



# RNN in fact returns a sequence of outputs

RNN definition:  $\mathbf{y}_n = \text{RNN}(\mathbf{x}_{1:n})$

Let's have  $n = 3$ , so our input sequence is  $\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3$ :

$$\mathbf{y}_2 = \text{RNN}(\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$$

But our input sequence also contains  $\mathbf{x}_1, \mathbf{x}_2$ , so:

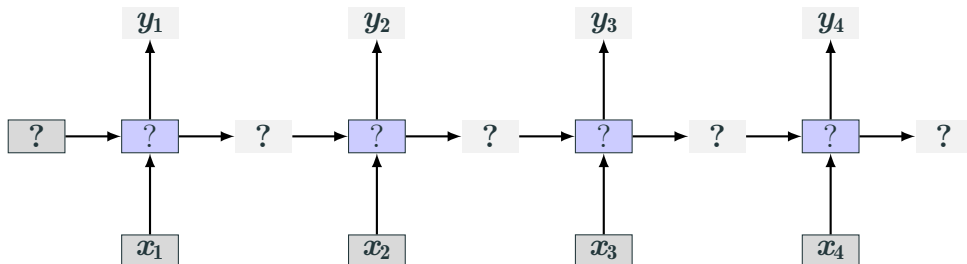
$$\mathbf{y}_2 = \text{RNN}(\mathbf{x}_1, \mathbf{x}_2)$$

Which makes RNN outputting a vector  $\mathbf{y}_i \in \mathbb{R}^{d_{out}}$  at each position  $i \in (1, \dots, n)$

Let's call this sequence-outputting function  $\text{RNN}^*$ :

$$\mathbf{y}_{1:n} = \text{RNN}^*(\mathbf{x}_{1:n})$$

# Adding outputs to our sketch



# Recap

For a sequence of input vectors  $\mathbf{x}_{1:i}$

$$\mathbf{y}_i = \text{RNN}(\mathbf{x}_{1:i}) \quad \mathbf{y}_{1:n} = \text{RNN}^*(\mathbf{x}_{1:n})$$

# Recap

For a sequence of input vectors  $\mathbf{x}_{1:i}$

$$\mathbf{y}_i = \text{RNN}(\mathbf{x}_{1:i}) \quad \mathbf{y}_{1:n} = \text{RNN}^*(\mathbf{x}_{1:n})$$

Without knowing what RNN actually is, what are the advantages?

# Recap

For a sequence of input vectors  $\mathbf{x}_{1:i}$

$$\mathbf{y}_i = \text{RNN}(\mathbf{x}_{1:i}) \quad \mathbf{y}_{1:n} = \text{RNN}^*(\mathbf{x}_{1:n})$$

Without knowing what RNN actually is, what are the advantages?

- Each output  $\mathbf{y}_i$  takes into account the entire history  $\mathbf{x}_{1:i}$  without Markov property

What to do with  $\mathbf{y}_n$  or  $\mathbf{y}_{1:n}$ ?

# Recap

For a sequence of input vectors  $\mathbf{x}_{1:i}$

$$\mathbf{y}_i = \text{RNN}(\mathbf{x}_{1:i}) \quad \mathbf{y}_{1:n} = \text{RNN}^*(\mathbf{x}_{1:n})$$

Without knowing what RNN actually is, what are the advantages?

- Each output  $\mathbf{y}_i$  takes into account the entire history  $\mathbf{x}_{1:i}$  without Markov property

What to do with  $\mathbf{y}_n$  or  $\mathbf{y}_{1:n}$ ?

- Use for further prediction, e.g., plug into softmax, MLP, etc.

# Underlying mechanism of RNNs — states

For “passing information” from one position to the next, i.e. from

$$\mathbf{y}_i = \text{RNN}(\mathbf{x}_{1:i})$$

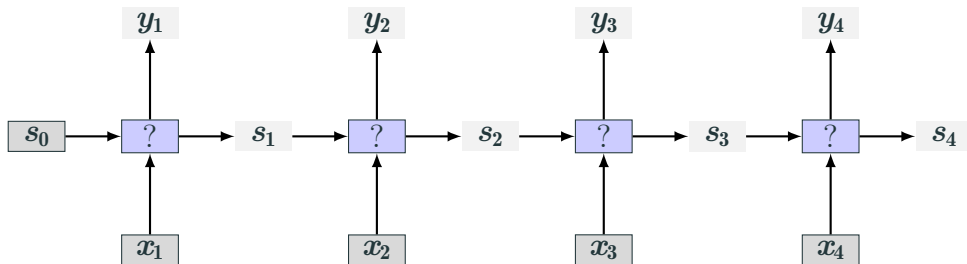
to

$$\mathbf{y}_{i+1} = \text{RNN}(\mathbf{x}_{1:i+1})$$

we use a “state” vector

$$\mathbf{s}_i \in \mathbb{R}^{d_{state}}$$

# Adding state vectors





# Define RNN recursively — Computing current state

At each step  $i \in (1, \dots, n)$  we have

- Current input vector  $\mathbf{x}_i$
- Vector of the previous state  $\mathbf{s}_{i-1}$ <sup>1</sup>

and compute

- Current state  $\mathbf{s}_i$

$$\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i) \quad (\text{we will specify } R \text{ later})$$

---

<sup>1</sup>Initial state vector  $\mathbf{s}_0$  — often omitted, assumed to be zero-filled

# Define RNN recursively — Computing current output

At each step  $i \in (1, \dots, n)$  we have

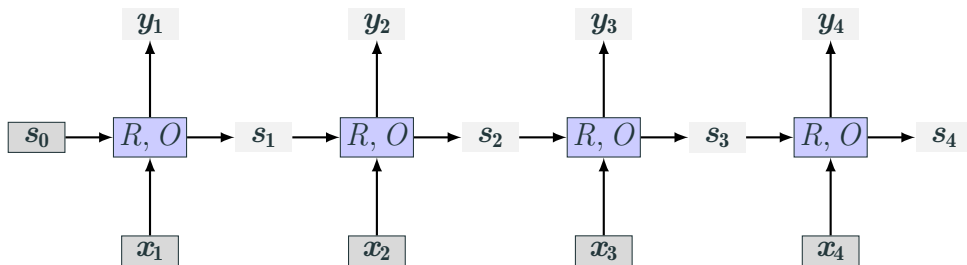
- Current input vector  $\mathbf{x}_i$
- Vector of the previous state  $\mathbf{s}_{i-1}$

and compute

- Current state  $\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$
- Current output  $\mathbf{y}_i$

$$\mathbf{y}_i = O(\mathbf{s}_i) \quad (\text{we will specify } O \text{ later})$$

## Adding $R$ and $O$



# Summary

At each step  $i \in (1, \dots, n)$  we have

- Current input  $\mathbf{x}_i$  and previous state  $\mathbf{s}_{i-1}$

and compute

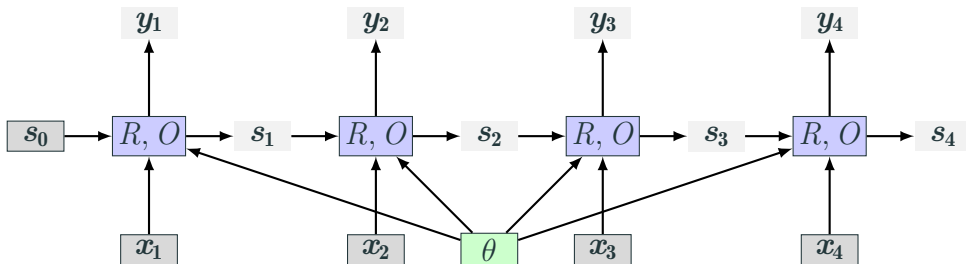
- $\mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i)$  and  $\mathbf{y}_i = O(\mathbf{s}_i)$

The functions  $R$  and  $O$  are the same for each position  $i$

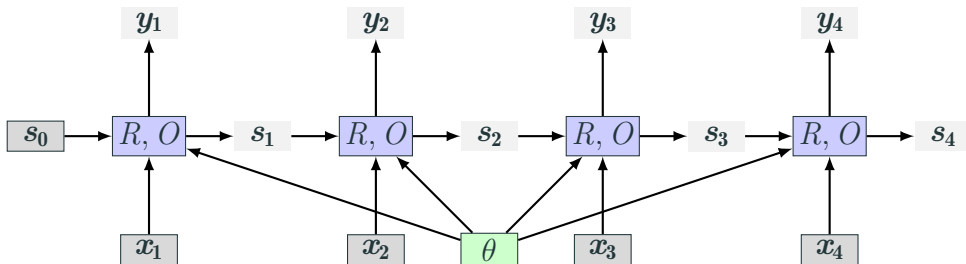
## RNN

$$\mathbf{y}_{1:n} = \text{RNN}^*(\mathbf{x}_{1:n}, \mathbf{s}_0) \quad \mathbf{s}_i = R(\mathbf{s}_{i-1}, \mathbf{x}_i) \quad \mathbf{y}_i = O(\mathbf{s}_i)$$

# Graphical visualization of abstract RNN (unrolled)

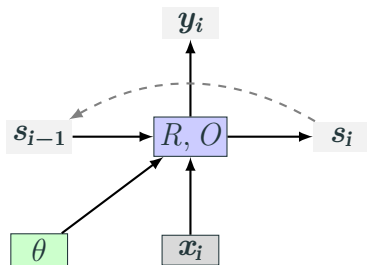


# Graphical visualization of abstract RNN (unrolled)



Note that  $\theta$  (parameters) are “shared” (the same) for all positions

# Graphical visualization of abstract RNN (recursive)



# Recurrent Neural Networks (RNN)

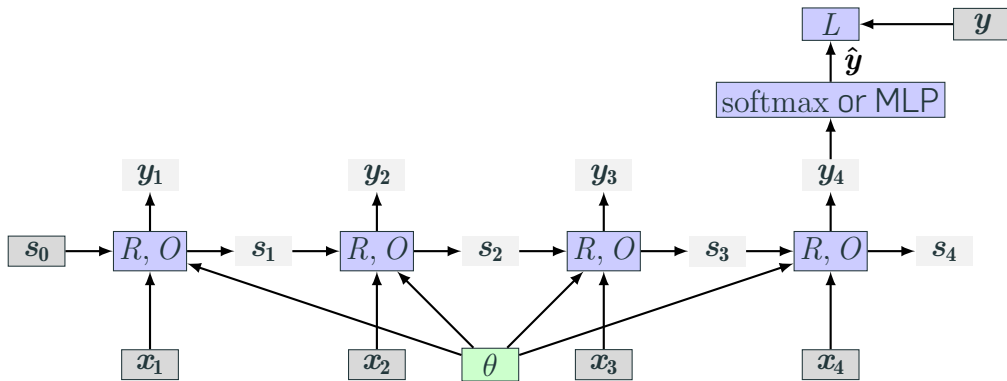
## abstraction

---

RNN as 'acceptor' or 'encoder'



# Supervision on the last output



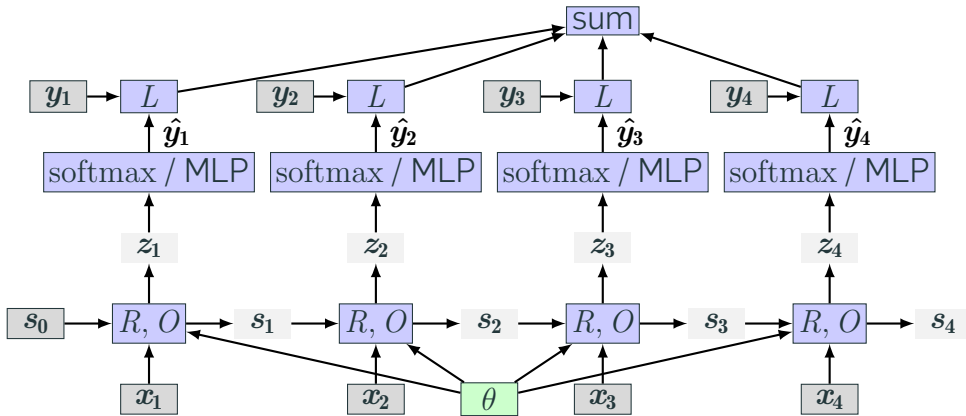
The loss is computed on the final output (e.g., directly on  $y_n$  or by putting  $y_n$  through MLP)

# Recurrent Neural Networks (RNN) abstraction

---

RNN as 'transducer'

# Supervision on each output



For sequence tagging — loss on each position, overall network's loss simply as a sum of losses

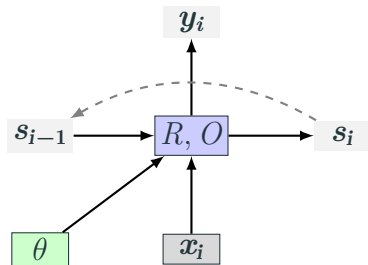
# Bi-directional RNNs

Simple idea: Run one RNN from left-to-right (forward,  $f$ ) and another RNN from right-to-left (backward,  $b$ ), and concatenate

$$\text{biRNN}(\mathbf{x}_{1:i}, i) = \mathbf{y}_i = [\text{RNN}_f(\mathbf{x}_{1:i}); \text{RNN}_b(\mathbf{x}_{n:i})]$$

Both for encoder (concatenate the last outputs) and transducer (concatenate each step's output)

## But what is happening ‘inside’ $R$ and $O$ ?



# RNN architectures

---

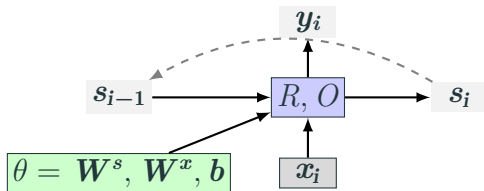
- 1 Recurrent Neural Networks (RNN) abstraction
- 2 RNN architectures**
- 3 Encoder-decoder architectures

# RNN architectures

---

## Simple RNN

# Elman Network or Simple-RNN (S-RNN)

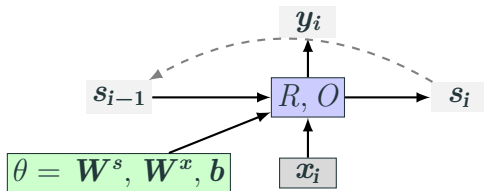


$$s_i = R(x_i, s_{i-1}) = g(s_{i-1} W^s + x_i W^x + b)$$

$$y_i = O(s_i) = s_i$$



# Elman Network or Simple-RNN (S-RNN)



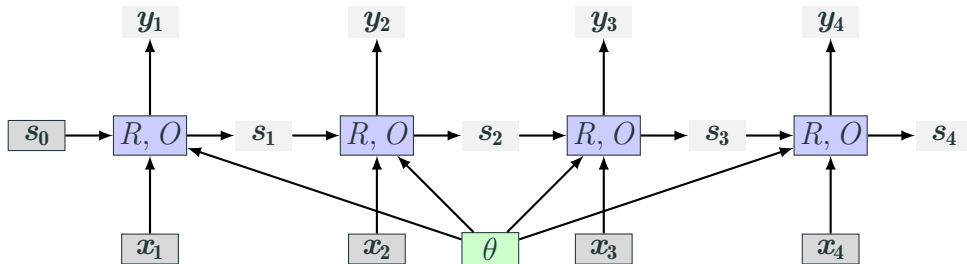
$$s_i = R(x_i, s_{i-1}) = g(s_{i-1} W^s + x_i W^x + b)$$

$$y_i = O(s_i) = s_i$$

$$s_i, y_i \in \mathbb{R}_s^d \quad x_i \in \mathbb{R}_{in}^d \quad W^x \in \mathbb{R}^{d_{in} \times d_s} \quad W^s \in \mathbb{R}^{d_s \times d_s} \quad b \in \mathbb{R}^{d_s}$$

$g$  — commonly tanh or ReLU

# Elman Network and vanishing gradient



Gradients might vanish ( $\rightarrow 0$ ) as they propagate back through the computation graph

- Severe in deeper nets, especially in recurrent networks
- Hard for the S-RNN to capture long-range dependencies

# RNN architectures

---

## Gated architectures

# RNN as a general purpose computing device

State  $s_i$  represents a finite memory

## Recall: Simple RNN

$$s_i = R(x_i, s_{i-1}) = g(s_{i-1} W^s + x_i W^x + b)$$

# RNN as a general purpose computing device

State  $s_i$  represents a finite memory

## Recall: Simple RNN

$$s_i = R(x_i, s_{i-1}) = g(s_{i-1} \mathbf{W}^s + x_i \mathbf{W}^x + \mathbf{b})$$

Each application of function  $R$

- Reads the current memory  $s_{i-1}$
- Reads the current input  $x_i$
- Operates on them in some way
- Writes the result to the memory  $s_i$

# RNN as a general purpose computing device

State  $s_i$  represents a finite memory

## Recall: Simple RNN

$$s_i = R(x_i, s_{i-1}) = g(s_{i-1} \mathbf{W}^s + x_i \mathbf{W}^x + \mathbf{b})$$

Each application of function  $R$

- Reads the current memory  $s_{i-1}$
- Reads the current input  $x_i$
- Operates on them in some way
- Writes the result to the memory  $s_i$

Memory access not controlled: At each step, entire memory state is read, and entire memory state is written

# How to provide more controlled memory access?

Memory vector  $\mathbf{s} \in \mathbb{R}^d$  and input vector  $\mathbf{x} \in \mathbb{R}^d$

# How to provide more controlled memory access?

Memory vector  $\mathbf{s} \in \mathbb{R}^d$  and input vector  $\mathbf{x} \in \mathbb{R}^d$

Let's have a binary vector ("gate")  $\mathbf{g} \in \{0, 1\}^d$



# How to provide more controlled memory access?

Memory vector  $\mathbf{s} \in \mathbb{R}^d$  and input vector  $\mathbf{x} \in \mathbb{R}^d$

Let's have a binary vector ("gate")  $\mathbf{g} \in \{0, 1\}^d$

**Hadamard-product**  $\mathbf{z} = \mathbf{u} \odot \mathbf{v}$

Fancy name for element-wise multiplication  $z_{[i]} = u_{[i]} \cdot v_{[i]}$

$$\mathbf{s}' \leftarrow \mathbf{g} \odot \mathbf{x} + (\mathbf{1} + \mathbf{g}) \odot \mathbf{s}$$

# How to provide more controlled memory access?

Memory vector  $\mathbf{s} \in \mathbb{R}^d$  and input vector  $\mathbf{x} \in \mathbb{R}^d$

Let's have a binary vector ("gate")  $\mathbf{g} \in \{0, 1\}^d$

**Hadamard-product**  $\mathbf{z} = \mathbf{u} \odot \mathbf{v}$

Fancy name for element-wise multiplication  $z_{[i]} = u_{[i]} \cdot v_{[i]}$

$$\mathbf{s}' \leftarrow \mathbf{g} \odot \mathbf{x} + (\mathbf{1} + \mathbf{g}) \odot \mathbf{s}$$

- Reads the entries in  $\mathbf{x}$  corresponding to ones in the gate, writes them to the memory

# How to provide more controlled memory access?

Memory vector  $\mathbf{s} \in \mathbb{R}^d$  and input vector  $\mathbf{x} \in \mathbb{R}^d$

Let's have a binary vector ("gate")  $\mathbf{g} \in \{0, 1\}^d$

**Hadamard-product**  $\mathbf{z} = \mathbf{u} \odot \mathbf{v}$

Fancy name for element-wise multiplication  $z_{[i]} = u_{[i]} \cdot v_{[i]}$

$$\mathbf{s}' \leftarrow \mathbf{g} \odot \mathbf{x} + (\mathbf{1} + \mathbf{g}) \odot \mathbf{s}$$

- Reads the entries in  $\mathbf{x}$  corresponding to ones in the gate, writes them to the memory
- Remaining locations are copied from the memory
- Note that the operation  $+$  here is modulo 2

## Gate example

Updating memory position 2

$$\begin{pmatrix} 8 \\ 11 \\ 3 \end{pmatrix} \leftarrow \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \odot \begin{pmatrix} 10 \\ 11 \\ 12 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \odot \begin{pmatrix} 8 \\ 9 \\ 3 \end{pmatrix}$$
$$\mathbf{s}' \leftarrow \mathbf{g} \odot \mathbf{x} + (\mathbf{1} + \mathbf{g}) \odot \mathbf{s}$$

# Gate example

Updating memory position 2

$$\begin{pmatrix} 8 \\ 11 \\ 3 \end{pmatrix} \leftarrow \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \odot \begin{pmatrix} 10 \\ 11 \\ 12 \end{pmatrix} + \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} \odot \begin{pmatrix} 8 \\ 9 \\ 3 \end{pmatrix}$$
$$\mathbf{s}' \leftarrow \mathbf{g} \odot \mathbf{x} + (\mathbf{1} + \mathbf{g}) \odot \mathbf{s}$$

Could be used for gates in RNNs! But:

- Our gates are not learnable
- Our hard-gates are not differentiable

Solution: Replace with 'soft' gates

# RNN architectures

---

## LSTM

# Long Short-Term Memory (LSTM)

Designed to solve the vanishing gradients problem, first to introduce the gating mechanism

# Long Short-Term Memory (LSTM)

Designed to solve the vanishing gradients problem, first to introduce the gating mechanism

LSTM splits the state vector  $s_i$  exactly in two halves

- One half is treated as 'memory cells'
- The other half is 'working memory'



# Long Short-Term Memory (LSTM)

Designed to solve the vanishing gradients problem, first to introduce the gating mechanism

LSTM splits the state vector  $s_i$  exactly in two halves

- One half is treated as 'memory cells'
- The other half is 'working memory'

## Memory cells

- Designed to preserve the memory, and also the error gradients, across time
- Controlled through *differentiable gating components* — smooth functions that simulate logical gates

# Long Short-Term Memory (LSTM)

The state at time  $j$  is composed of two vectors:

- $c_j$  — the memory component
- $h_j$  — the hidden state component

# Long Short-Term Memory (LSTM)

The state at time  $j$  is composed of two vectors:

- $c_j$  — the memory component
- $h_j$  — the hidden state component

At each input state  $j$ , a gate decides how much of the new input should be written to the memory cell, and how much of the memory cell should be forgotten

# Long Short-Term Memory (LSTM)

The state at time  $j$  is composed of two vectors:

- $c_j$  — the memory component
- $h_j$  — the hidden state component

At each input state  $j$ , a gate decides how much of the new input should be written to the memory cell, and how much of the memory cell should be forgotten

There are three gates

- $i$  — input gate
- $f$  — forget gate
- $o$  — output gate

# LSTM architecture

$c_{j-1}$

$h_{j-1}$

$x_j$

# LSTM architecture

$c_{j-1}$

$h_{j-1}$

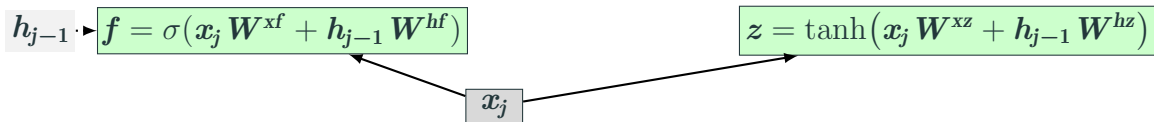
$x_j$

$$z = \tanh(x_j W^{xz} + h_{j-1} W^{hz})$$

$z$  — update candidate

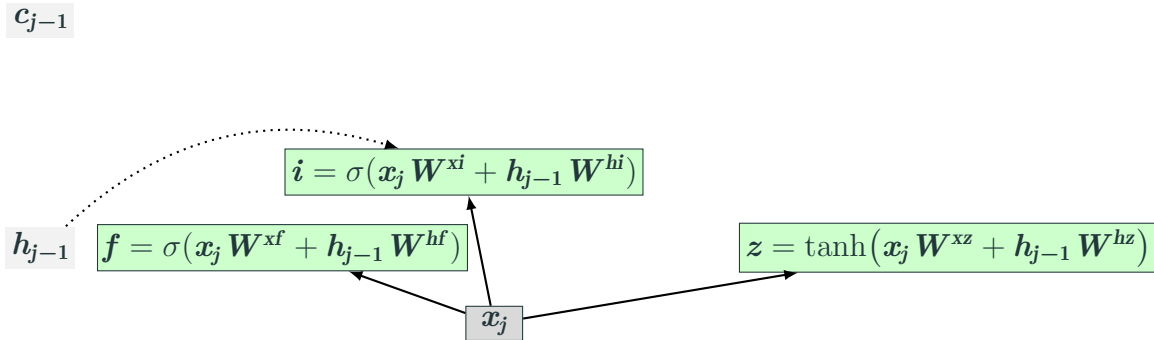
# LSTM architecture

$c_{j-1}$



$z$  — update candidate

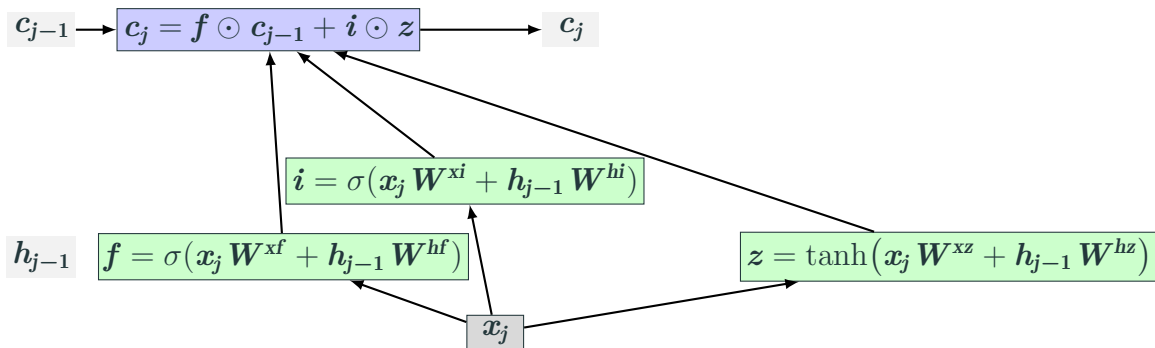
# LSTM architecture



$z$  — update candidate

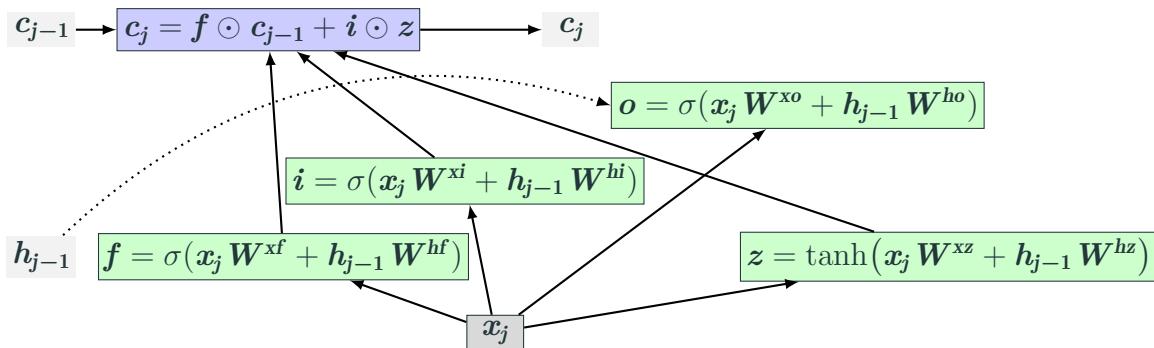


# LSTM architecture



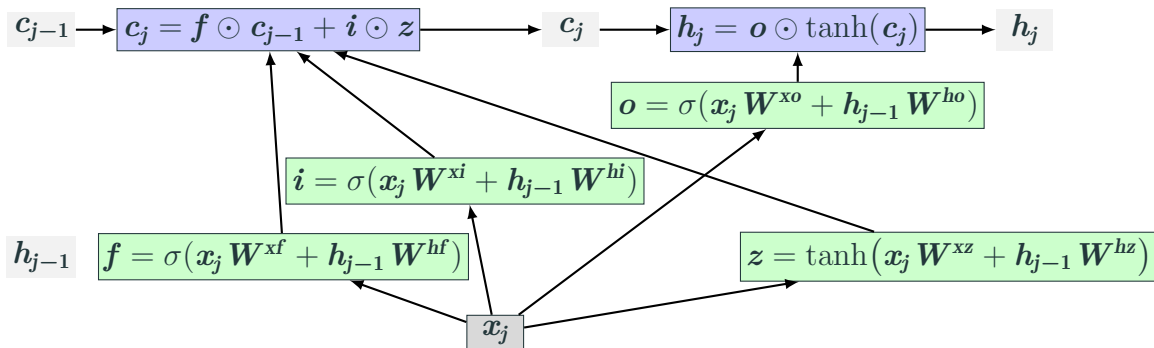
$z$  — update candidate

# LSTM architecture



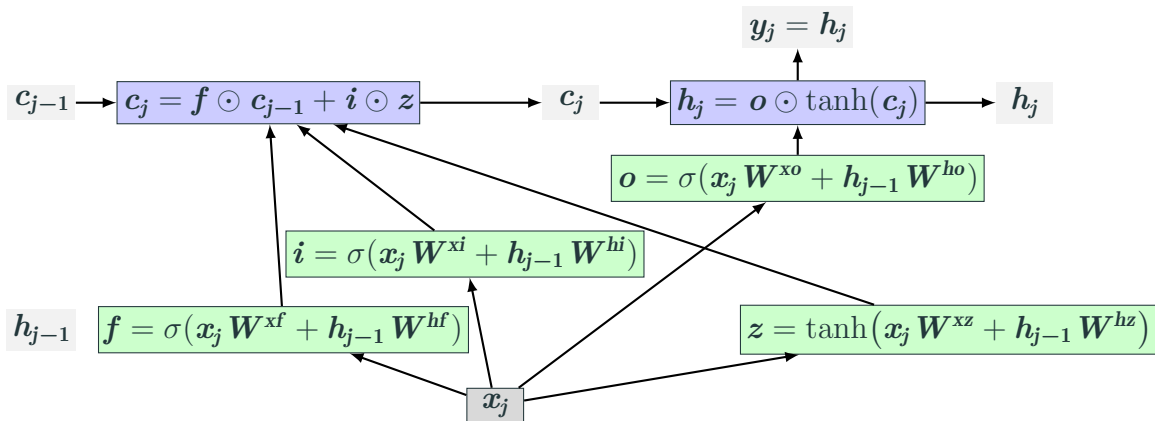
$z$  — update candidate

# LSTM architecture



$z$  — update candidate

# LSTM architecture

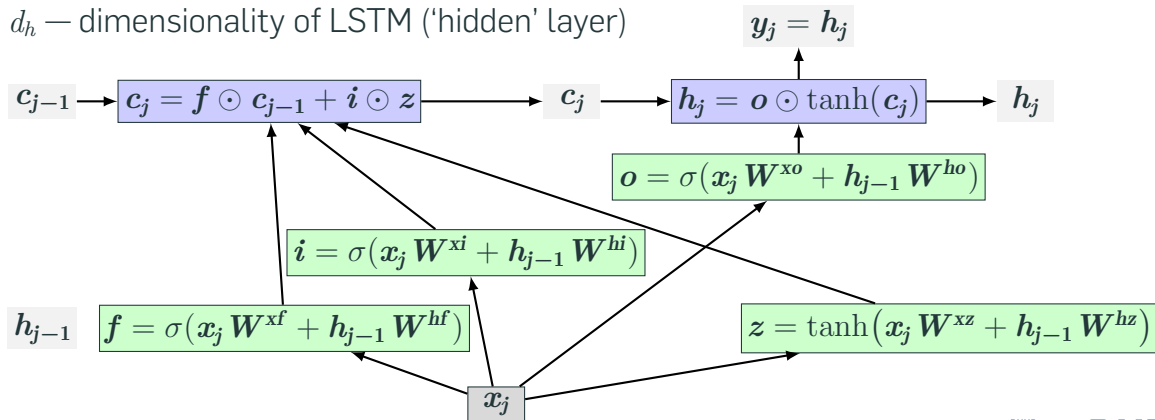


$z$  — update candidate

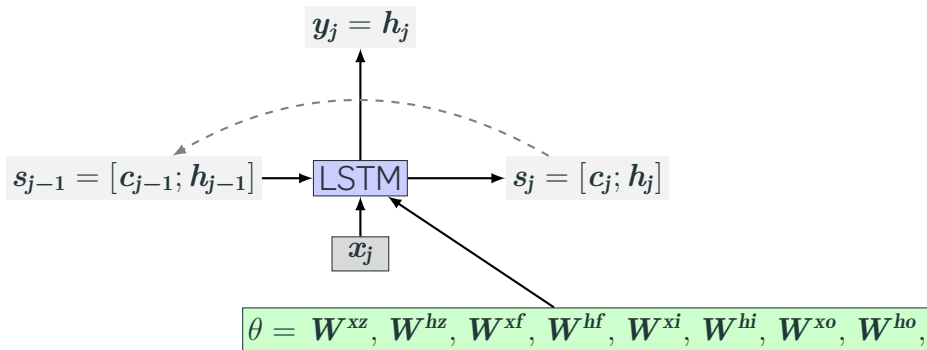
# LSTM parameters and dimensions

$$\mathbf{x}_j \in \mathbb{R}^{d_{in}} \quad \mathbf{c}_j, \mathbf{h}_j, \mathbf{y}_j, \mathbf{i}, \mathbf{f}, \mathbf{o}, \mathbf{z} \in \mathbb{R}^{d_h} \quad \mathbf{W}^{x\star} \in \mathbb{R}^{d_{in} \times d_h} \quad \mathbf{W}^{h\star} \in \mathbb{R}^{d_h \times d_h}$$

$d_h$  — dimensionality of LSTM ('hidden' layer)



# LSTM as a 'layer'



We also ignored bias terms for each gate

# Recap

---

- 1 Recurrent Neural Networks (RNN) abstraction
- 2 RNN architectures
- 3 Encoder-decoder architectures

# Encoder-decoder architectures

---

- 1 Recurrent Neural Networks (RNN) abstraction
- 2 RNN architectures
- 3 Encoder-decoder architectures**



# The problem of variable output sequence length

We have a sequence of  $n$  **input** vectors  $\mathbf{x}_{1:n} = \mathbf{x}_1, \dots, \mathbf{x}_n$

Each input vector has the same dimension  $d_{in} : \mathbf{x}_i \in \mathbb{R}^{d_{in}}$

# The problem of variable output sequence length

We have a sequence of  $n$  **input** vectors  $\mathbf{x}_{1:n} = \mathbf{x}_1, \dots, \mathbf{x}_n$

Each input vector has the same dimension  $d_{in} : \mathbf{x}_i \in \mathbb{R}^{d_{in}}$

We also have a **sequence** of  $d_{out}$ -dimensional vector

$\mathbf{y}_{1:\hat{n}} \in \mathbb{R}^{\hat{n} \times d_{out}}$  **outputs**

# The problem of variable output sequence length

We have a sequence of  $n$  **input** vectors  $\mathbf{x}_{1:n} = \mathbf{x}_1, \dots, \mathbf{x}_n$

Each input vector has the same dimension  $d_{in}$  :  $\mathbf{x}_i \in \mathbb{R}^{d_{in}}$

We also have a **sequence** of  $d_{out}$ -dimensional vector

$\mathbf{y}_{1:\hat{n}} \in \mathbb{R}^{\hat{n} \times d_{out}}$  **outputs**

RNNs produce a sequence of outputs

$$\mathbf{y}_{1:n} = \text{RNN}(\mathbf{x}_{1:n})$$

What are we missing?

# The problem of variable output sequence length

We have a sequence of  $n$  **input** vectors  $\mathbf{x}_{1:n} = \mathbf{x}_1, \dots, \mathbf{x}_n$

Each input vector has the same dimension  $d_{in}$  :  $\mathbf{x}_i \in \mathbb{R}^{d_{in}}$

We also have a **sequence** of  $d_{out}$ -dimensional vector

$\mathbf{y}_{1:\hat{n}} \in \mathbb{R}^{\hat{n} \times d_{out}}$  **outputs**

RNNs produce a sequence of outputs

$$\mathbf{y}_{1:n} = \text{RNN}(\mathbf{x}_{1:n})$$

What are we missing?

- The input and output sequence: rarely of same length

# Generating a variable length sequence

**Translate to German:** *I like attending deep learning lectures*

## Generating a variable length sequence

**Translate to German:** *I like attending deep learning lectures*

**Output:** *Ich besuche gerne Deep-Learning-Vorlesungen*

# Generating a variable length sequence

**Translate to German:** *I like attending deep learning lectures*

**Output:** *Ich besuche gerne Deep-Learning-Vorlesungen*

Current approach:

- 1 Tokenize input sequence
- 2 Obtain a word embedding (e.g. word2vec) for each token
- 3 Use a RNN (e.g. LSTM) to encode sequence of tokens
- 4 Generate token sequence in target language

# Generating a variable length sequence

**Translate to German:** *I like attending deep learning lectures*

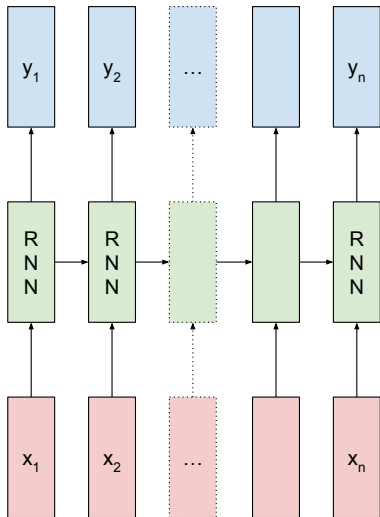
**Output:** *Ich besuche gerne Deep-Learning-Vorlesungen*

Current approach:

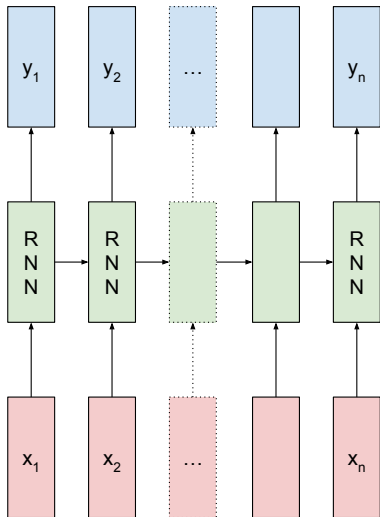
- 1 Tokenize input sequence
- 2 Obtain a word embedding (e.g. word2vec) for each token
- 3 Use a RNN (e.g. LSTM) to encode sequence of tokens
- 4 Generate token sequence in target language
  - Multi-class classification over target vocabulary



# Generating a variable length sequence

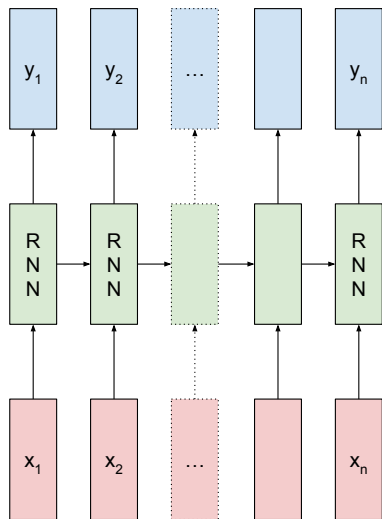


# Generating a variable length sequence



How to solve the issue of  
varying input/output  
lengths?

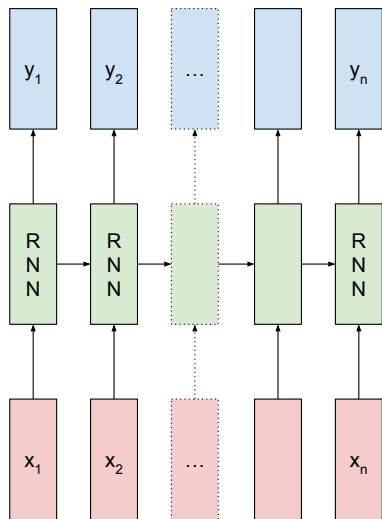
# Generating a variable length sequence



How to solve the issue of varying input/output lengths?

- 1 We **don't have to** stop generating after the last input
- 2 We can only consider outputs up to a special **"end token"**

# Generating a variable length sequence

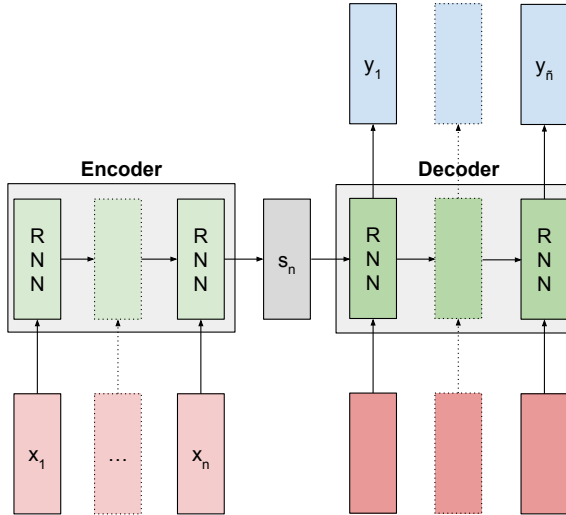


How to solve the issue of varying input/output lengths?

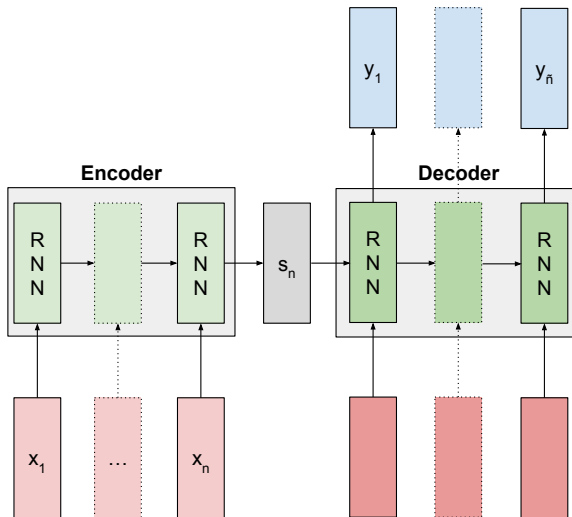
- 1 We **don't have to** stop generating after the last input
- 2 We can only consider outputs up to a special **"end token"**

Neither ideal

# Sequence-to-sequence models



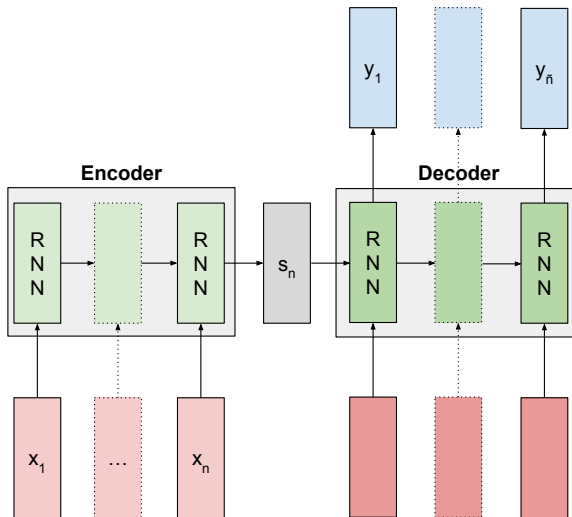
# Sequence-to-sequence models



Two networks

- **Encoder** (reader)  
RNN
- **Decoder** (writer)  
RNN

# Sequence-to-sequence models



Two networks

- **Encoder** (reader)  
RNN
- **Decoder** (writer)  
RNN

Note:

- Encoder and decoder have **separate** params

# The encoder-decoder architecture specifics

1 How to **initialize** decoder hidden **state**?



# The encoder-decoder architecture specifics

## 1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

# The encoder-decoder architecture specifics

1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

2 When do we **stop generating** with the decoder?

# The encoder-decoder architecture specifics

## 1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

## 2 When do we **stop generating** with the decoder?

- We use a **special token** (<EOS>, \n) to indicate the end-of-sequence
- When the **maximum generation length** is exceeded

# The encoder-decoder architecture specifics

## 1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

## 2 When do we **stop generating** with the decoder?

- We use a **special token** (<EOS>, \n) to indicate the end-of-sequence
- When the **maximum generation length** is exceeded

## 3 What are the **inputs** of the decoder?

# The encoder-decoder architecture specifics

## 1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

## 2 When do we **stop generating** with the decoder?

- We use a **special token** (<EOS>, \n) to indicate the end-of-sequence
- When the **maximum generation length** is exceeded

## 3 What are the **inputs** of the decoder?

- The **previous output** of the decoder
  - Teacher forcing (with probability  $p$ ): use the **correct output**

# The encoder-decoder architecture specifics

## 1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

## 2 When do we **stop generating** with the decoder?

- We use a **special token** (<EOS>, \n) to indicate the end-of-sequence
- When the **maximum generation length** is exceeded

## 3 What are the **inputs** of the decoder?

- The **previous output** of the decoder
  - Teacher forcing (with probability  $p$ ): use the **correct output**
- What is the **initial input**  $x_0^{dec}$ ?

# The encoder-decoder architecture specifics

## 1 How to **initialize** decoder hidden **state**?

- $h_0^{dec} = h_n^{enc}$ : simply copy the last encoder state
- $h_0^{dec} = \text{NN}_\theta(h_n^{enc})$ : transform the last encoder state

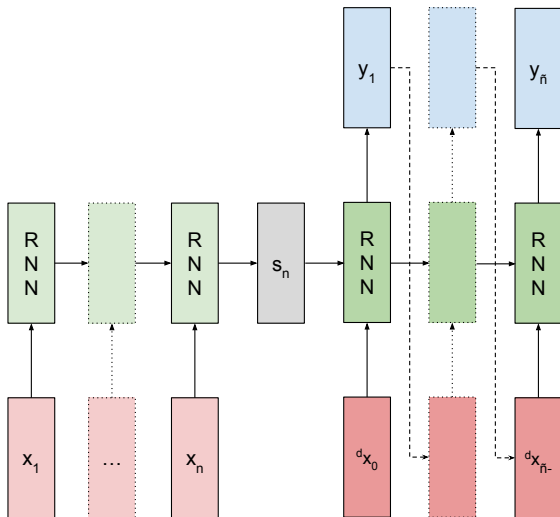
## 2 When do we **stop generating** with the decoder?

- We use a **special token** (<EOS>, \n) to indicate the end-of-sequence
- When the **maximum generation length** is exceeded

## 3 What are the **inputs** of the decoder?

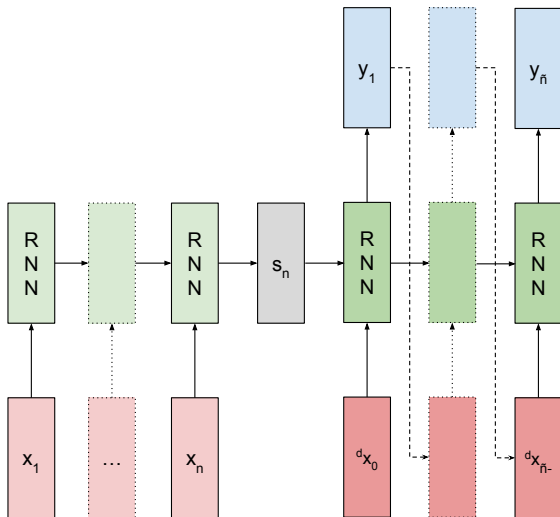
- The **previous output** of the decoder
  - Teacher forcing (with probability  $p$ ): use the **correct output**
- What is the **initial input**  $x_0^{dec}$ ?
  - A beginning-of-sequence **special token** (<BOS>)

# The encoder-decoder architecture





# The encoder-decoder architecture



## Decoder inputs

- $x_0^{dec} = \langle \text{BOS} \rangle$
- $x_i^{dec} = y_i^{dec}$  **if** no teacher forcing
- $x_i^{dec} = \hat{y}_i$  **if** we use teacher forcing

# Take aways

- RNNs for arbitrary long input
- Encoding the entire sequence and/or each step
- Modeling freedom with bi-directional RNNs
- Vanishing gradients in deep nets — gating mechanism, memory cells
- LSTM a particularly powerful RNN
- Encoder-decoder RNNs for text-to-text tasks

# License and credits

Licensed under Creative Commons  
Attribution-ShareAlike 4.0 International  
(CC BY-SA 4.0)



## Credits

Ivan Habernal, Martin Tutek

Content from ACL Anthology papers licensed under CC-BY  
<https://www.aclweb.org/anthology>