

Deep Learning

*Ms Computer Science,
Data Science, Cybersecurity
2nd semester - 6 CFU*

Part 3 by Alessandro Sperduti

References:

[**Deep Learning Book \(main book\)**](#)

University of Padova, A.A. 2022/23

Sequence Modeling: Recurrent and Recursive Nets (chapter 10) ...and Transformers

University of Padova, A.A. 2022/23

Learning in Sequential Domains

Large amounts of data in sequential form more and more common:

- sensor-equipped devices (e.g., drones);
 - sensor networks (such as *Internet of Things*);
 - development of inexpensive human-machine interaction interfaces;
 - ...

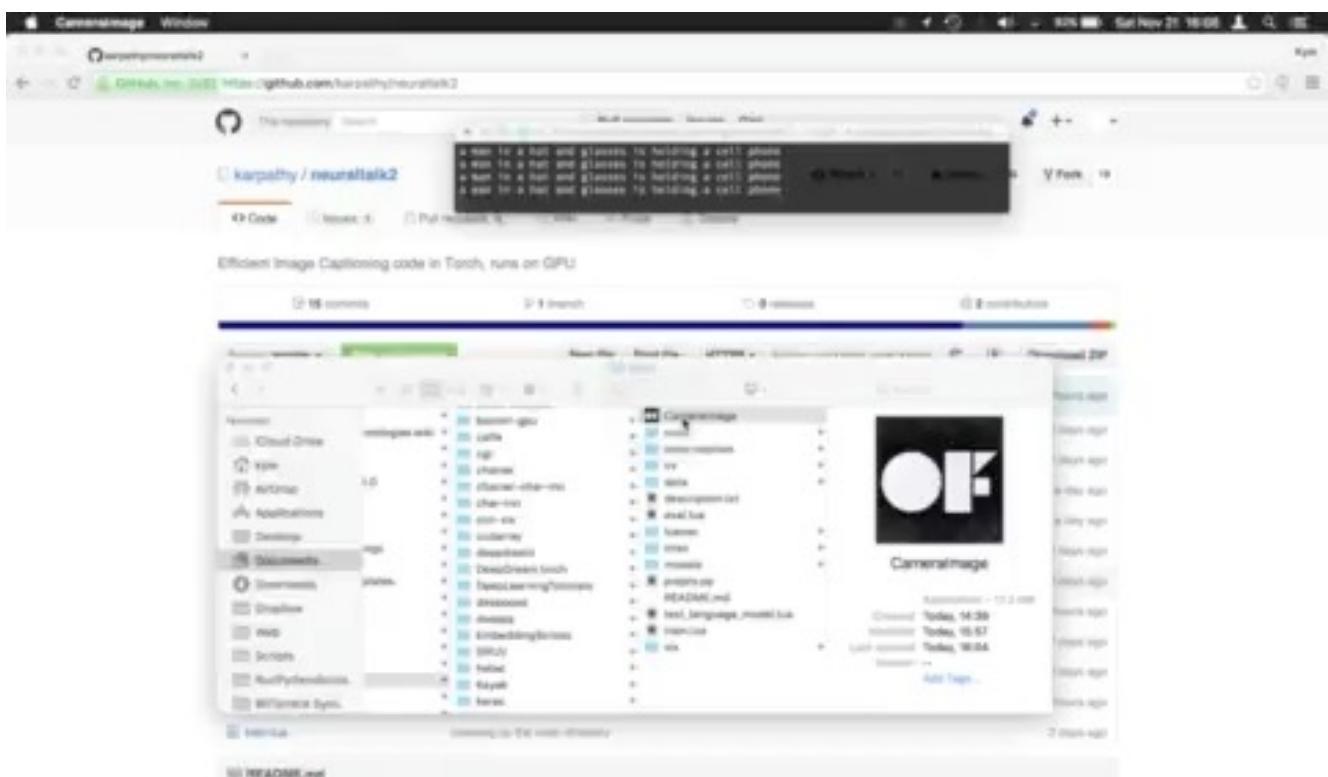
Many are the applicative tasks of interest:

- monitoring and classification of human behaviour ([sequence → category](#))
 - content video captioning ([image → sequence of text](#))
 - sentiment classification ([sequence of text → sentiment](#))
 - language machine translation ([sequence of text → sequence of text](#))
 - ...

Not only classification or transduction, also generation of sequences from static data (e.g., image), or from text (e.g., ChatGPT)

A.A. 2022/23: Deep Learning

Video Captioning



A.A. 2022/23: Deep Learning

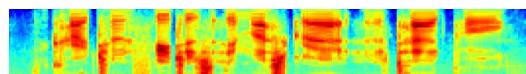
Learning in Sequential Domains

Why learning in sequential domains is different wrt static domains ?

Because successive points in sequential data are strongly correlated

Machine learning models and algorithms for sequence learning have to

- consider that data points are not independent
- deal with sequential distortions and/or variations (e.g. In speech, variations in speaking rate)



- make use of contextual information



Static Data

Learn $P(o|x)$

- $\text{typeof}(x)$:
 - a *fixed-size tuple* of predictive attributes
When using neural networks categorical attributes are encoded using numerical values \Rightarrow a fixed-size array of real numbers
- $\text{typeof}(o)$:
 - Classification: a categorical variable
 - Regression: a numerical (multivariate) variable

Sequential Data

Learn $P(o|x)$

- $\text{typeof}(x)$: a **sequence** $x^{(1)}, \dots, x^{(t)}, \dots$ where each $x^{(t)}$ has a static type

Examples:

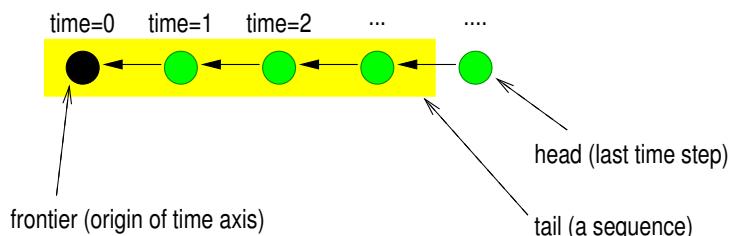
- Speech recognition: $x^{(t)}$ is an array of k real values representing acoustic features measured in a 10 ms time interval
 - Molecular biology: $x^{(t)}$ is a symbol in a finite alphabet of 20 amino-acids
 - Image captioning: only one item per sequence, i.e. $x^{(1)}$ is an image
-
- $\text{typeof}(o)$ may be either static (e.g., sequence classification) or a sequence

A.A. 2022/23: Deep Learning

Sequences

Using mathematical induction, a sequence is either

- an external vertex, or
- an ordered pair (t, h) where the *head* h is a vertex and the *tail* t is a sequence



A.A. 2022/23: Deep Learning

Sequential Transductions

- Let \mathcal{X} and \mathcal{O} be the **input** and **output** label spaces
- We denote by \mathcal{X}^* the set of all sequences with labels in \mathcal{X}
- Cf. the classical definition in language theory: $\mathcal{X} = \Sigma$ (a set of terminal symbols) and Σ^* is the free monoid over Σ
- A general transduction \mathcal{T} is a subset of $\mathcal{X}^* \times \mathcal{O}^*$
- We shall limit our discussion to functions:

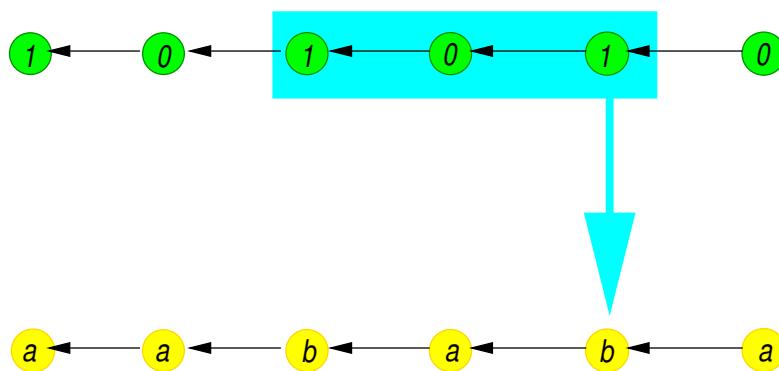
$$\mathcal{T} : \mathcal{X}^* \rightarrow \mathcal{O}^*$$

A.A. 2022/23: Deep Learning

Algebraic and Finite-Memory Transductions

$\mathcal{T}(\cdot)$ has **finite memory** $k \in \mathbb{N}$ if

- $\forall x \in \mathcal{X}^*$ and $\forall t$, $\mathcal{T}(x^{(t)})$ only depends on $\{x^{(t)}, x^{(t-1)}, \dots, x^{(t-k)}\}$



$\mathcal{T}(\cdot)$ is **algebraic** if it has 0 finite memory (i.e., no memory at all)

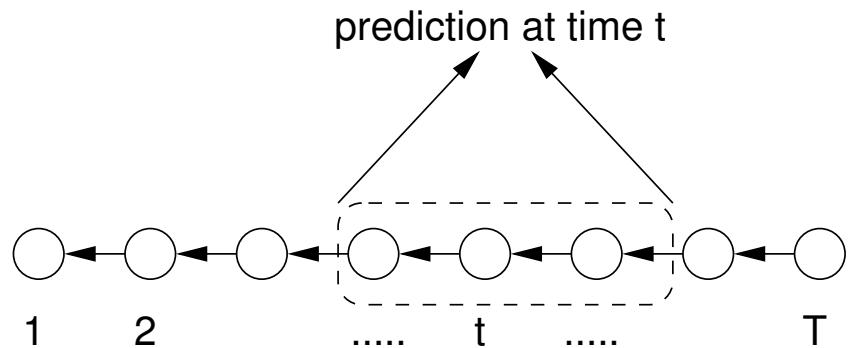
A.A. 2022/23: Deep Learning

Learning Sequences

Sequences have variable length but typical *machine learning models* have a fixed number of inputs

Solutions:

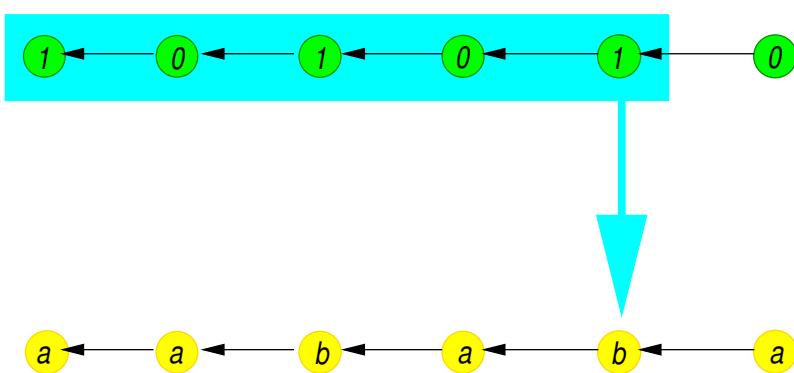
- Limit context to a **fixed-size window**



- Use **recurrent** models (causality)
- Use **Transformers** for non-causal sequences (e.g., text)

Causality

A transduction $\mathcal{T}(\cdot)$ is **causal** if the output at time t does not depend on **future** inputs (at time $t+1, t+2, \dots$)



Recursive State Representation

Intuition:

outputs depend on **hidden state variables (label space \mathcal{H})**

For each time t :

$$\begin{aligned}\mathbf{h}^{(t)} &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}, t) \\ \mathbf{o}^{(t)} &= g(\mathbf{h}^{(t)}, \mathbf{x}^{(t)}, t)\end{aligned}$$

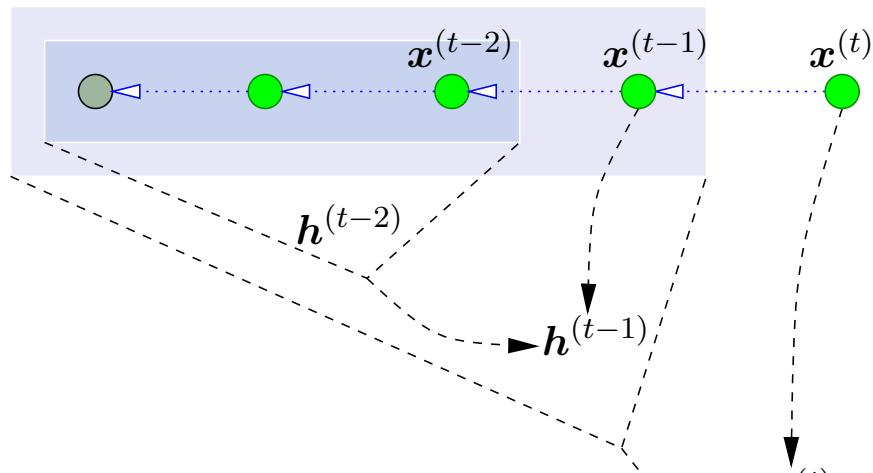
where

$$\begin{aligned}f : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{H} &\quad \text{state transition function} \\ g : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{O} &\quad \text{output function}\end{aligned}$$

A recursive state representation exists only if $\mathcal{T}(\cdot)$ is *causal*

$\mathcal{T}(\cdot)$ is **stationary** if $f(\cdot)$ and $g(\cdot)$ do not depend on t

Recursive State Update



$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)})$$

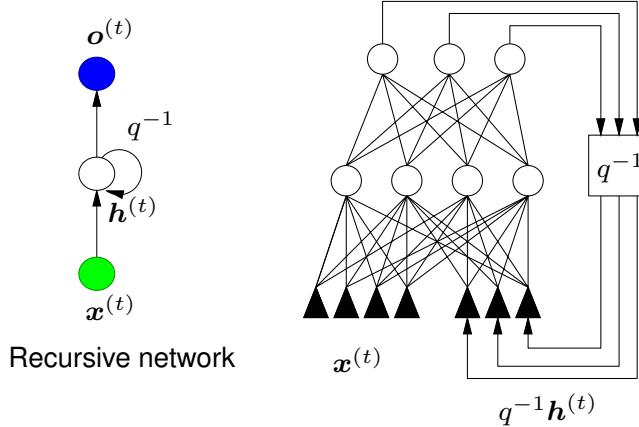
for $t = 0$, $\mathbf{h} = \mathbf{h}^{(0)}$ is the **initial state**

The initial state is associated with the external vertex (frontier)

Graphical Description

Graphical representation of the state transition function using the time shift operator q^{-1} :

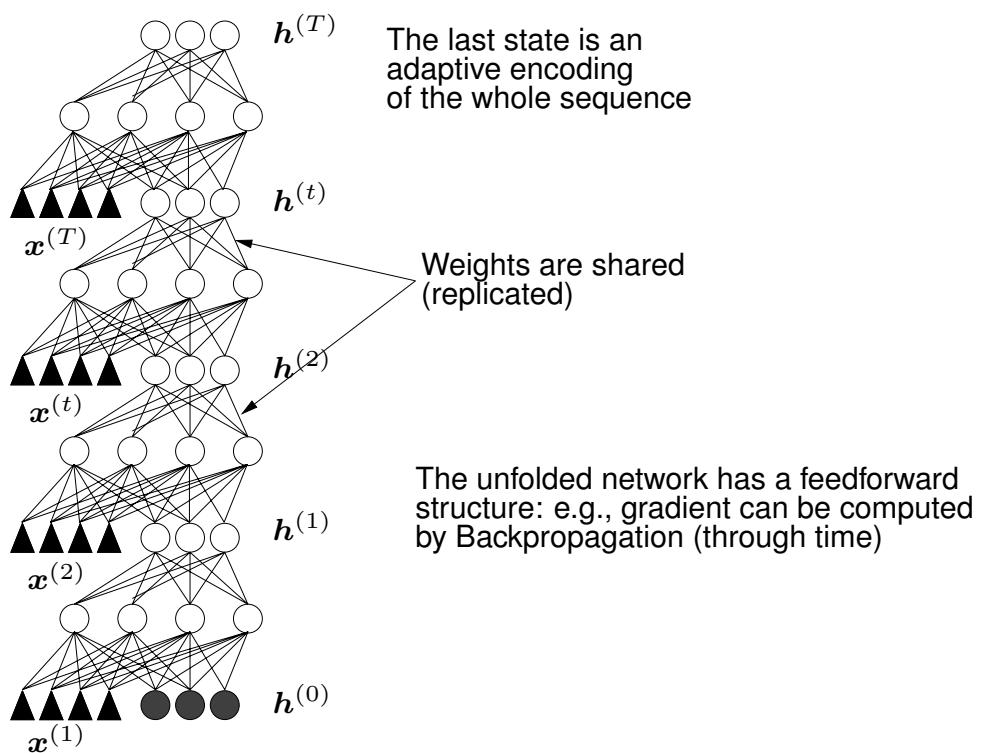
$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)})$$



$$q^{-1} \mathbf{h}^{(t)} = \mathbf{h}^{(t-1)} \text{ (unitary time delay)}$$

A.A. 2022/23: Deep Learning

Time Unfolding

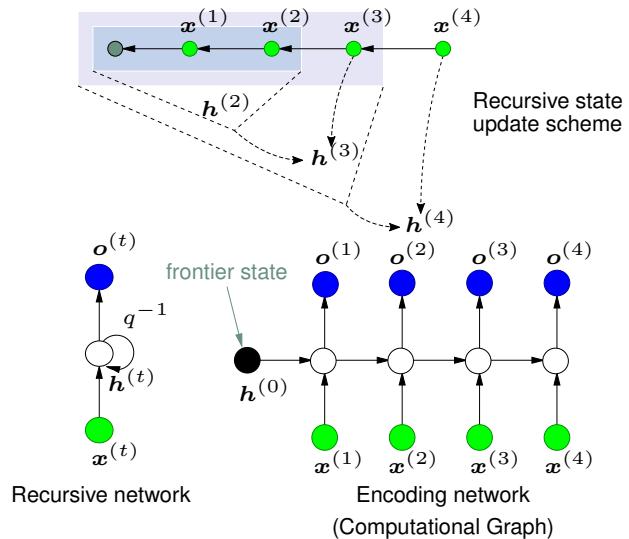


A.A. 2022/23: Deep Learning

Encoding Networks

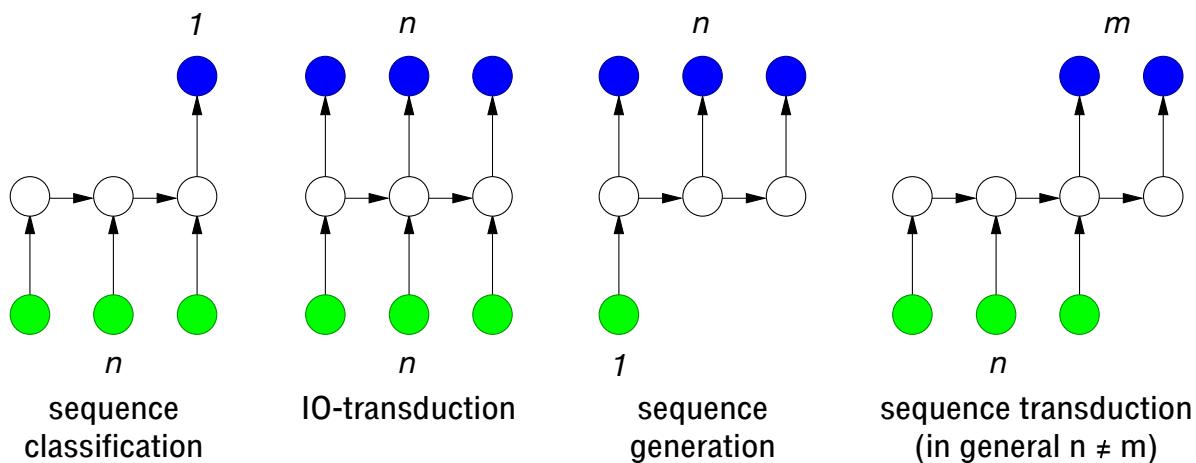
(Computational Graph)

Given a sequence $s \in \mathcal{X}^*$ and a recursive transduction \mathcal{T} , the *encoding network* associated to s and \mathcal{T} is formed by unrolling (*time unfolding*) the recursive network of \mathcal{T} through the input sequence s



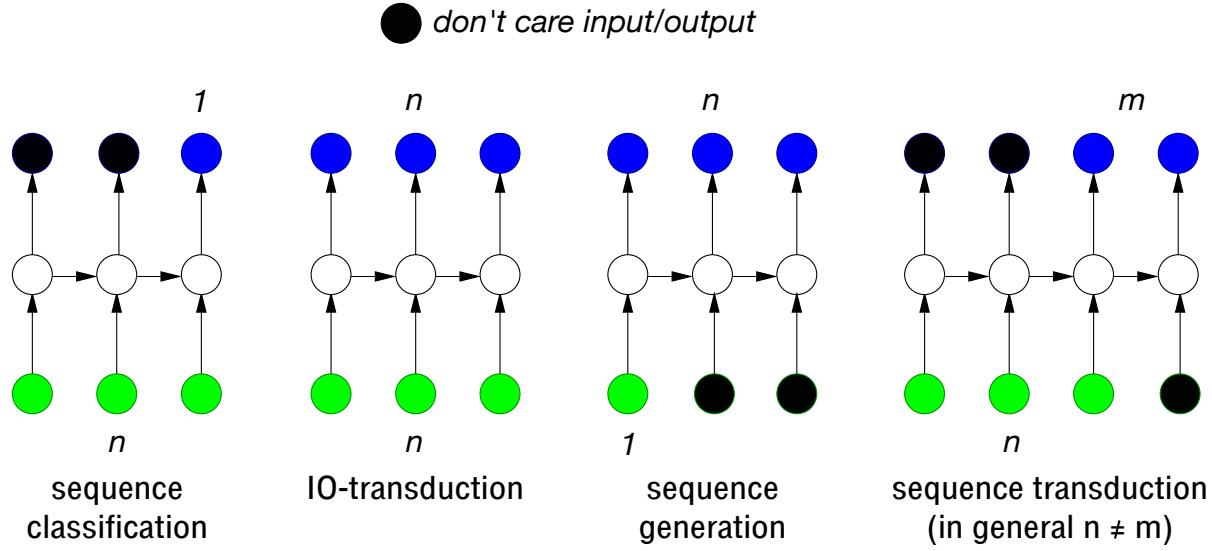
A.A. 2022/23: Deep Learning

Examples of Sequential Transductions



A.A. 2022/23: Deep Learning

A Unifying Perspective as IO-transductions



A.A. 2022/23: Deep Learning

Recursive State Representation

Intuition:

outputs depend on **hidden state variables (label space \mathcal{H})**

For each time t :

$$\begin{aligned}\mathbf{h}^{(t)} &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}, t) \\ \mathbf{o}^{(t)} &= g(\mathbf{h}^{(t)}, \mathbf{x}^{(t)}, t)\end{aligned}$$

where

$$\begin{aligned}f : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{H} &\quad \text{state transition function} \\ g : \mathcal{H} \times \mathcal{X} \rightarrow \mathcal{O} &\quad \text{output function}\end{aligned}$$

A recursive state representation exists only if $\mathcal{T}(\cdot)$ is *causal*

$\mathcal{T}(\cdot)$ is **stationary** if $f(\cdot)$ and $g(\cdot)$ do not depend on t

A.A. 2022/23: Deep Learning

Model Types

How to implement $f(\cdot)$ and $g(\cdot)$?

Two general families of models:

- Linear

- Kalman Filter
- Hidden Markov Models
- Linear Dynamical Systems
 - Linear Autoencoders for Sequences

- ...

- Nonlinear

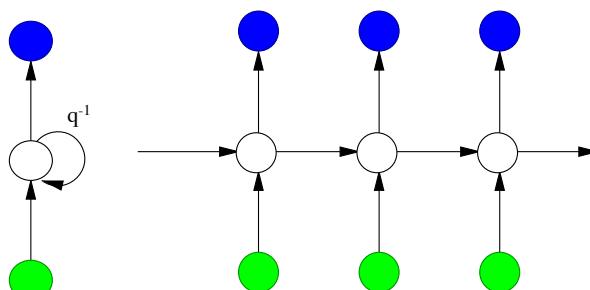
- Recurrent Neural Networks
- ...

Shallow Recurrent Neural Networks

Non-linear Dynamical System

$$\begin{aligned}\mathbf{h}^{(t)} &= f(\mathbf{Ux}^{(t)} + \mathbf{Wh}^{(t-1)} + \mathbf{b}), \\ \mathbf{o}^{(t)} &= g(\mathbf{Vh}^{(t)} + \mathbf{c}),\end{aligned}$$

where $f()$ and $g()$ are non-linear functions (e.g. $\text{tanh}()$), and $\mathbf{h}^{(0)} = \mathbf{0}$ (or can be learned jointly with the other parameters)



Shallow Recurrent Neural Networks

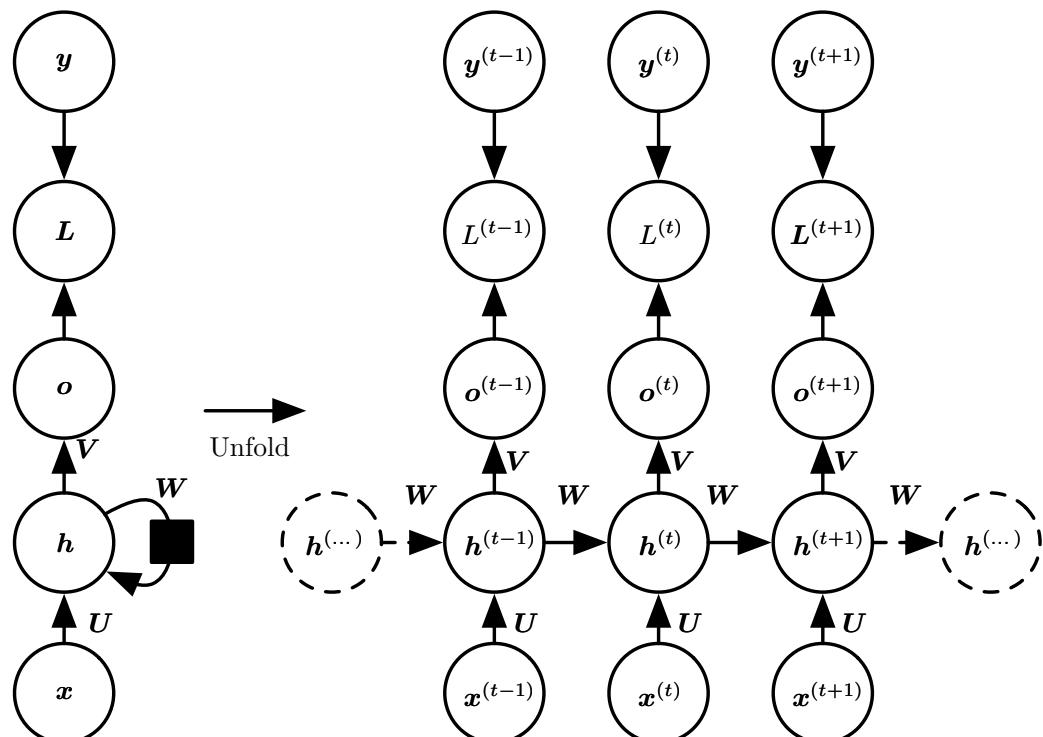
Example: IO-transduction RNN with discrete outputs

$$\begin{aligned} \mathbf{h}^{(t)} &= \tanh(\mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{b}), \\ \mathbf{o}^{(t)} &= \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c}, \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}), \\ (\text{Loss}) \quad L &= -\sum_t \log p_{\text{model}} \left(\hat{\mathbf{y}}^{(t)} | \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\} \right) \end{aligned}$$

where

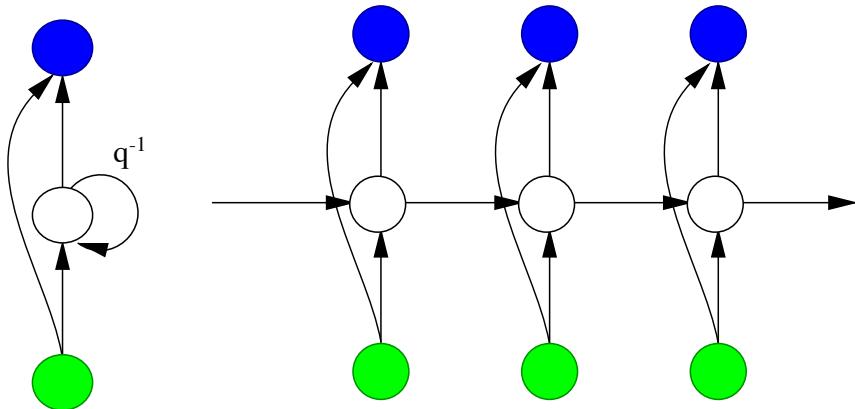
- $\mathbf{o}^{(t)}$ is the unnormalized log probabilities a time t
- $\mathbf{y}^{(t)}$ is the target vector a time t
- $p_{\text{model}} \left(\mathbf{y}^{(t)} | \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\} \right)$ is given by reading the entry for $y^{(t)}$ from the model's output vector $\hat{\mathbf{y}}^{(t)}$, i.e. $\hat{\mathbf{y}}^{(t)}$ is computed internally to L

Recurrent and Computational Graph



Some Additional Architectural Features for RNN

short-cut connections

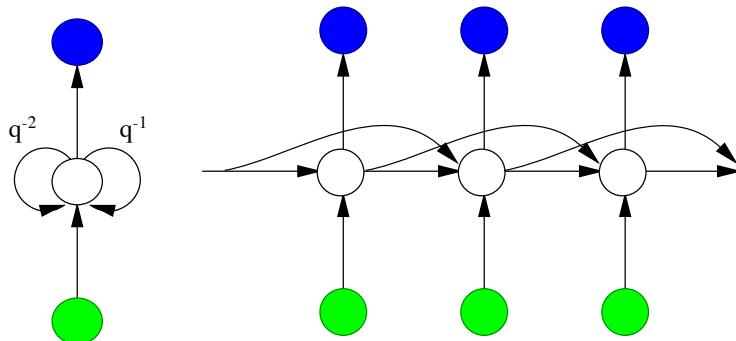


$$\mathbf{o}^{(t)} = g(\mathbf{V}^h \mathbf{h}^{(t)} + \mathbf{V}^x \mathbf{x}^{(t)} + \mathbf{b})$$

A.A. 2022/23: Deep Learning

Some Additional Architectural Features for RNN

higher-order states (e.g., 2nd order)

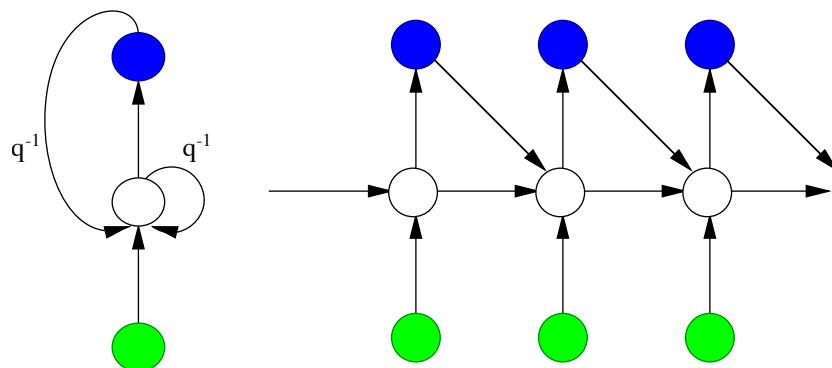


$$\mathbf{h}^{(t)} = f(\mathbf{U} \mathbf{x}^{(t)} + \mathbf{W}^{(1)} \mathbf{h}^{(t-1)} + \mathbf{W}^{(2)} \mathbf{h}^{(t-2)} + \mathbf{c})$$

A.A. 2022/23: Deep Learning

Some Additional Architectural Features for RNN

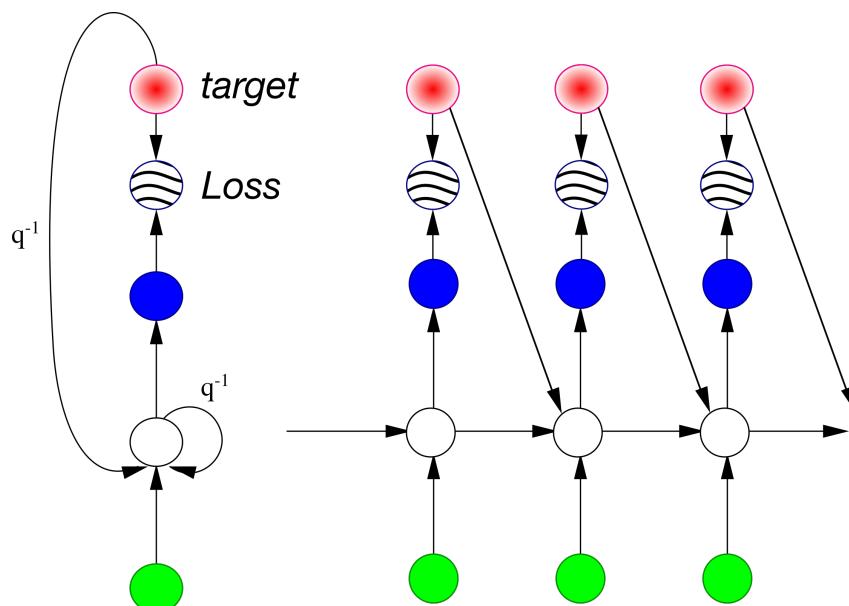
feedback from output (context)



$$\mathbf{h}^{(t)} = f(\mathbf{Ux}^{(t)} + \mathbf{Wh}^{(t-1)} + \mathbf{Zo}^{(t-1)})$$

Some Additional Architectural Features for RNN

feedback connections allow to *force* the target signal: **teacher forcing**



Some Additional Architectural Features for RNN

Example of Teacher Forcing

$$\mathbf{h}^{(t)} = \tanh(\mathbf{U}\mathbf{x}^{(t)} + \mathbf{W}\mathbf{o}^{(t-1)} + \mathbf{b}),$$

$$\mathbf{o}^{(t)} = \mathbf{V}\mathbf{h}^{(t)} + \mathbf{c},$$

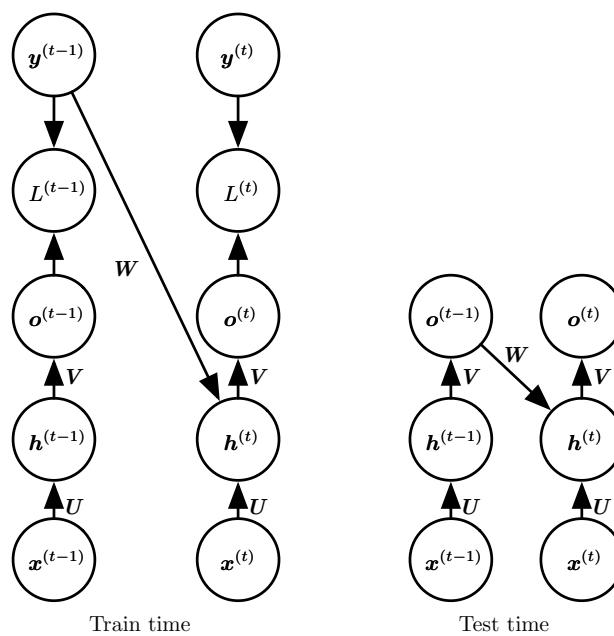
$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)}),$$

$$(\text{Loss}) \quad L = - \sum_t \log p_{\text{model}} \left(\mathbf{y}^{(t)} | \mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t-1)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)} \right)$$

e.g., for two steps the conditional maximum likelihood criterion is:

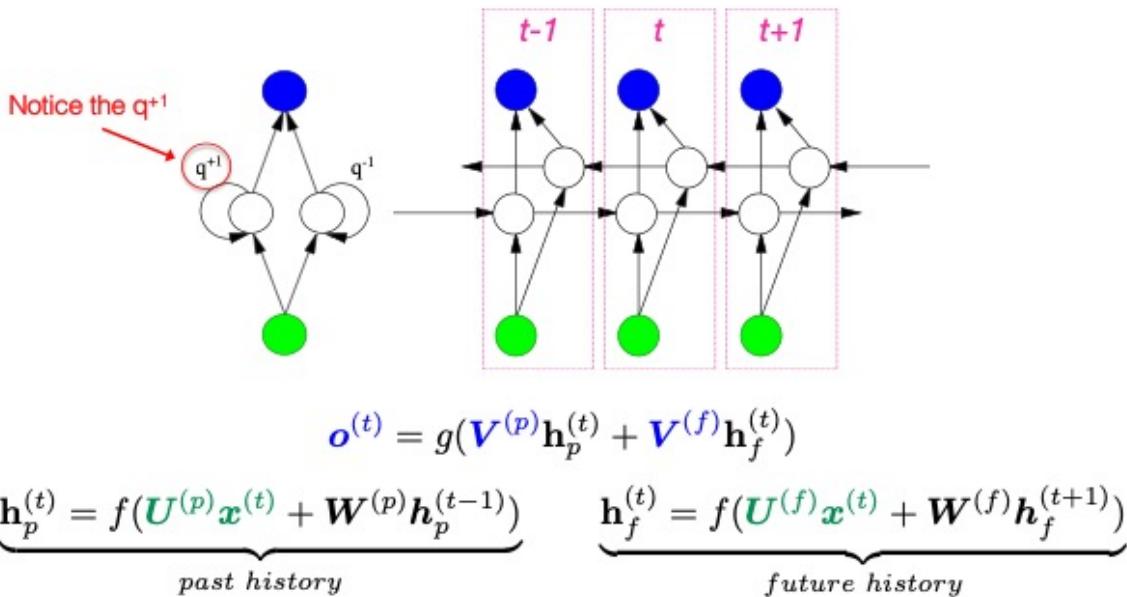
$$\begin{aligned} & \log p_{\text{model}} \left(\mathbf{y}^{(1)}, \mathbf{y}^{(2)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)} \right) \\ &= \log p_{\text{model}} \left(\mathbf{y}^{(2)} | \mathbf{y}^{(1)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)} \right) + \log p_{\text{model}} \left(\mathbf{y}^{(1)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)} \right) \end{aligned}$$

Example of Teacher Forcing



Some Additional Architectural Features for RNN

Bidirectional RNN: when the sequence is not temporal (or off-line processing)



e.g. DNA sequence in bioinformatics

A.A. 2022/23: Deep Learning

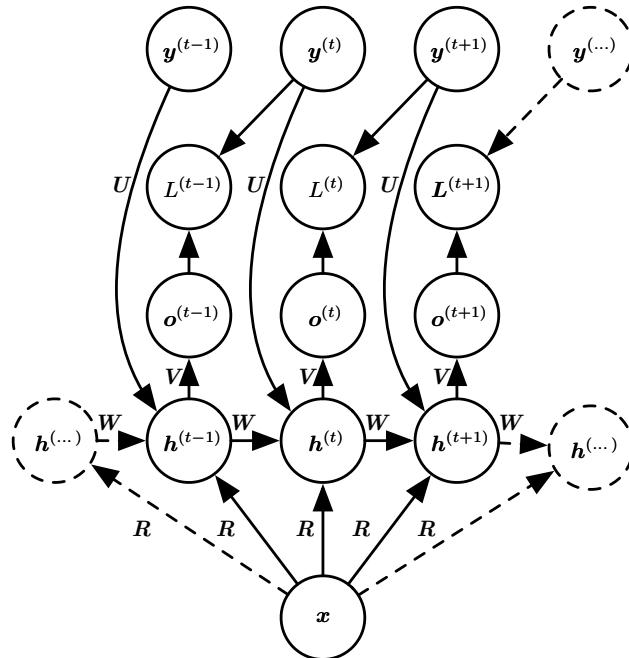
Some Additional Architectural Features for RNN

- All these architectural features (and others...) are *orthogonal*, i.e. they can be combined together
- We will see examples of such combinations when looking at some recent architectures for machine translation and image captioning (can you figure out which kind of *transductions* they call for ?)
- For these applications there are typical architectural *building blocks* that are worth discussing right now:
 - generation of a sequence from a vector (*1 to n transduction*)
 - sequence to sequence transduction via Encoder-Decoder (*n to m transduction*)

A.A. 2022/23: Deep Learning

Some Building Blocks for RNN

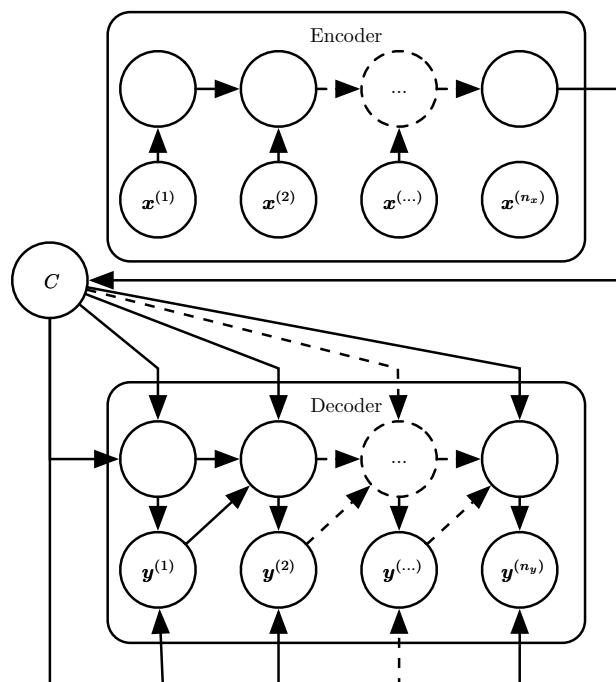
Generation of a sequence from a vector



A.A. 2022/23: Deep Learning

Some Building Blocks for RNN

Encoder-Decoder



A.A. 2022/23: Deep Learning