# Applied Deep Learning

Dr. Philippe Blaettchen
Bayes Business School (formerly Cass)

www.bayes.city.ac.uk

**Goals:** Introduce recurrent neural networks (RNNs) as a means to work with sequence data
- Understand the importance of sequences in everyday life and machine learning
- Recognize the complexities of working with sequences
- Grasp the theoretical foundations of RNNs

**How will we do this?**
- We introduce sequence data and discuss its prevalence
- We then discuss the concept of recurrence of neurons, neural network layers, and entire networks
- We briefly discuss how gradient descent works for neural networks
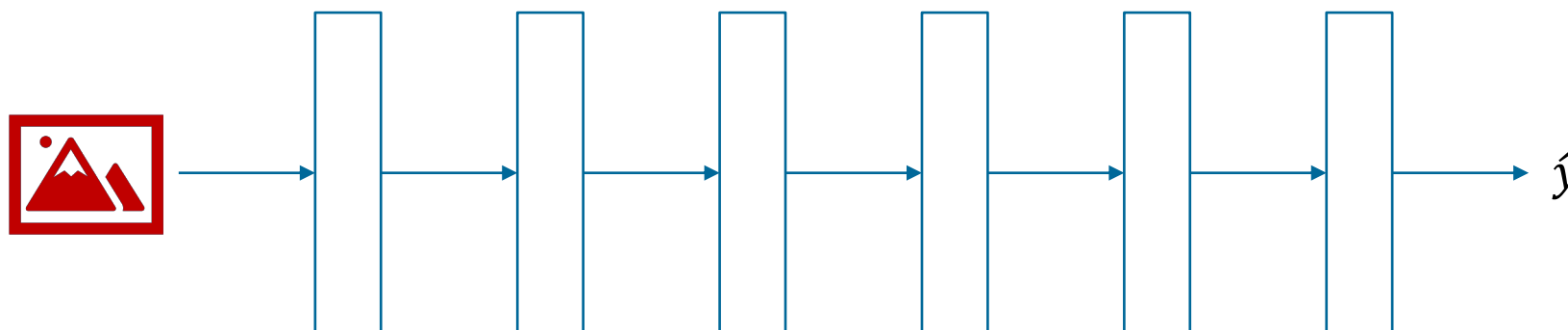- We then consider different RNN architectures as well as the challenges in training them

# Working with sequences

But what if the order between  and  matters?

Sequences are collections of multiple elements (i.e., data points), where:

- The order matters

- Elements may be repeated

- The length is variable

"Why do we care about sequences?"

"Why" "do" "we" "care" "about" "sequences" "?"

≠

"care" "about" "sequences" "Why" "do" "we" "?"

(unless you are  )

"a" "b" "o" "u" "t"

↑

"Why" "do" "we" "care" "about" "sequences" "?"

↓                    ↓

"w" "e"          "s" "e" "q" "u" "e" "n" "c" "e" "s"

Keep in mind:
- The order matters
- We may repeat individual elements
- Sequences vary in length

Warum kümmern wir uns um Sequenzen?

"Why" "do" "we" "care" "about" "sequences" "?"

¿Por qué nos preocupamos por las secuencias?

"And here I am, for all my lore, The wretched fool I was before"

# Modeling sequences with neural networks

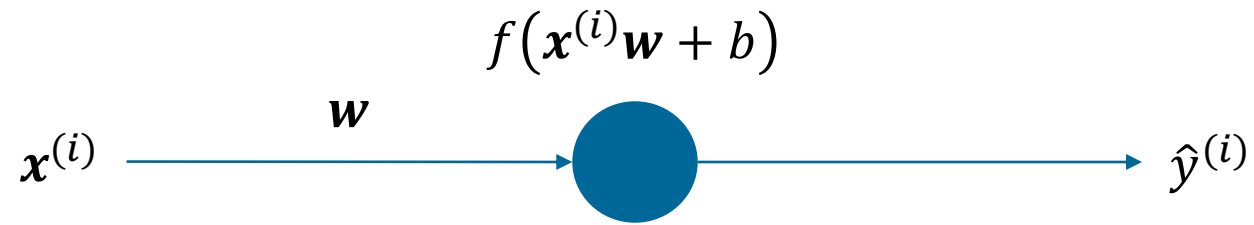$x^{<1>}$

$x^{<2>}$

$\vdots$

$x^{<T_x>}$

$y^{<1>}$

$y^{<2>}$

$\vdots$

$y^{<T_y>}$

Issues:
- Sequence lengths vary
- No definition of order
- Lack of parameter sharing: imagine a minute-long ECG

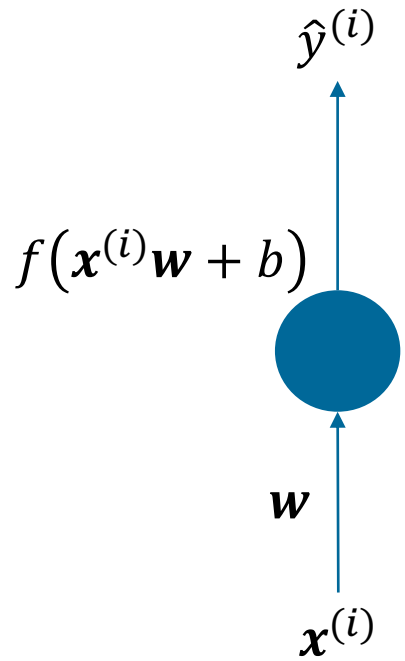$$f\left(\boldsymbol{x}^{(i)}\boldsymbol{w} + b\right)$$

$$\boldsymbol{w}$$

$$\boldsymbol{x}^{(i)} \longrightarrow \bigcirc \longrightarrow \hat{y}^{(i)}$$

$$\hat{y}^{(i)}$$

$$f\left(\boldsymbol{x}^{(i)}\boldsymbol{w} + b\right)$$

$$\boldsymbol{w}$$

$$\boldsymbol{x}^{(i)}$$

# A recurrent neuron



$$f\left(x^{(i)<4>}w + \hat{y}^{(i)<3>}\widetilde{w} + b\right)$$

$\hat{y}^{(i)}$

$\hat{y}^{(i)<1>}$  $\hat{y}^{(i)<2>}$  $\hat{y}^{(i)<3>}$  $\hat{y}^{(i)<4>}$

$w$

$w$  $w$  $w$  $w$

$x^{(i)}$

$x^{(i)<1>}$  $x^{(i)<2>}$  $x^{(i)<3>}$  $x^{(i)<4>}$

Time

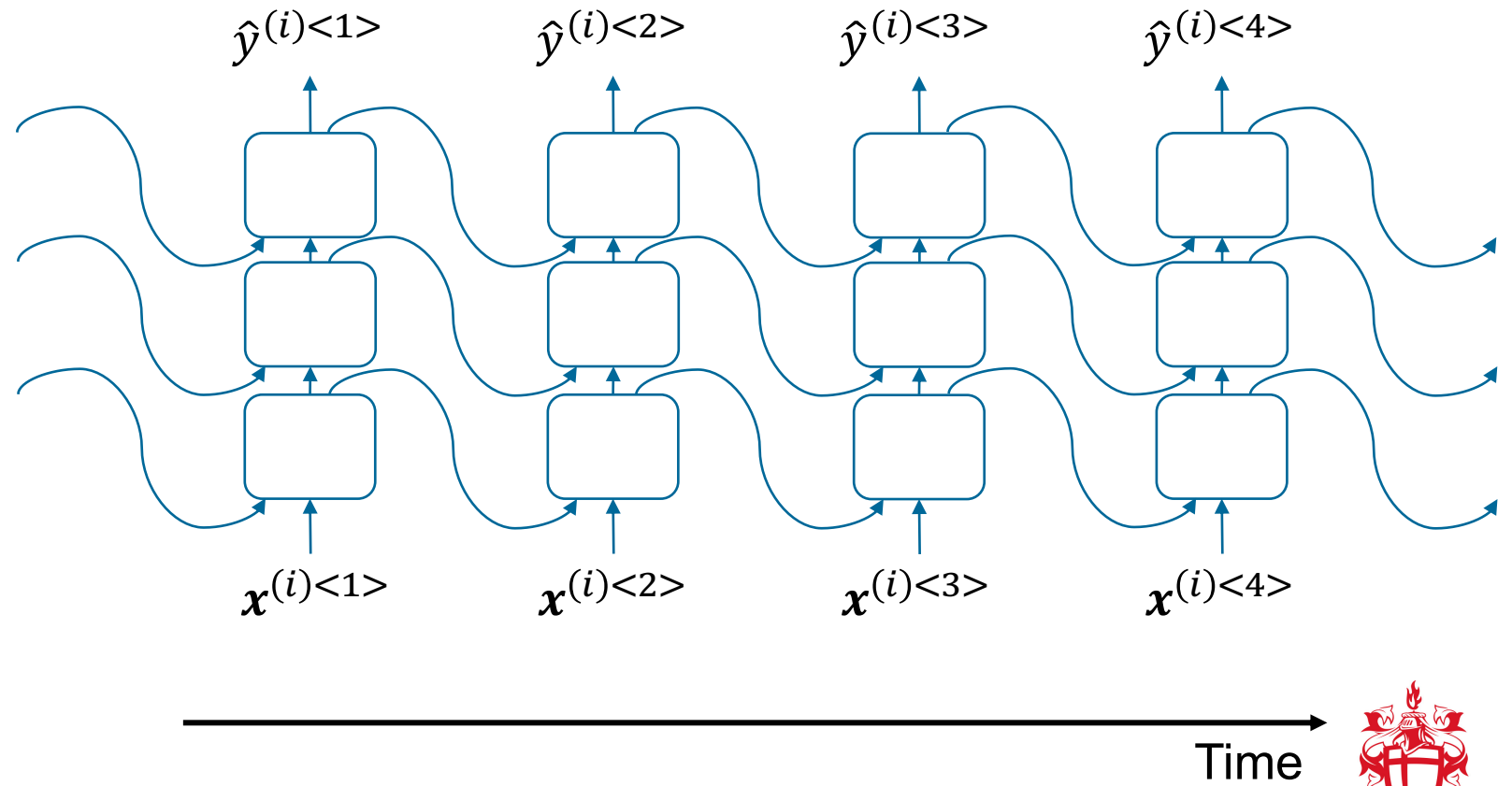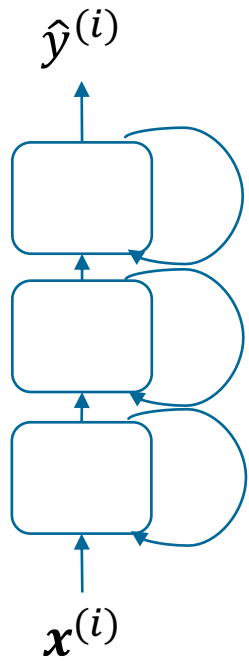# Layers of recurrent neurons – a recurrent neural network (RNN)

# Layers of recurrent neurons – a recurrent neural network (RNN)

$\hat{y}^{(i)}$

$\hat{y}^{(i)<1>}$     $\hat{y}^{(i)<2>}$     $\hat{y}^{(i)<3>}$     $\hat{y}^{(i)<4>}$

$x^{(i)}$

$x^{(i)<1>}$     $x^{(i)<2>}$     $x^{(i)<3>}$     $x^{(i)<4>}$

Time

# Deep RNNs

- At each time step, take the input and the "memory" (or *state*) from the previous time step to compute the output

- Use the same parameters (and, also, activation functions) across different time steps
  - Similar idea to parameter sharing in convolutional layers → we want to detect (recurring) patterns

- Usually, the loss is computed by summing up the losses on all time steps (but there are many variations)
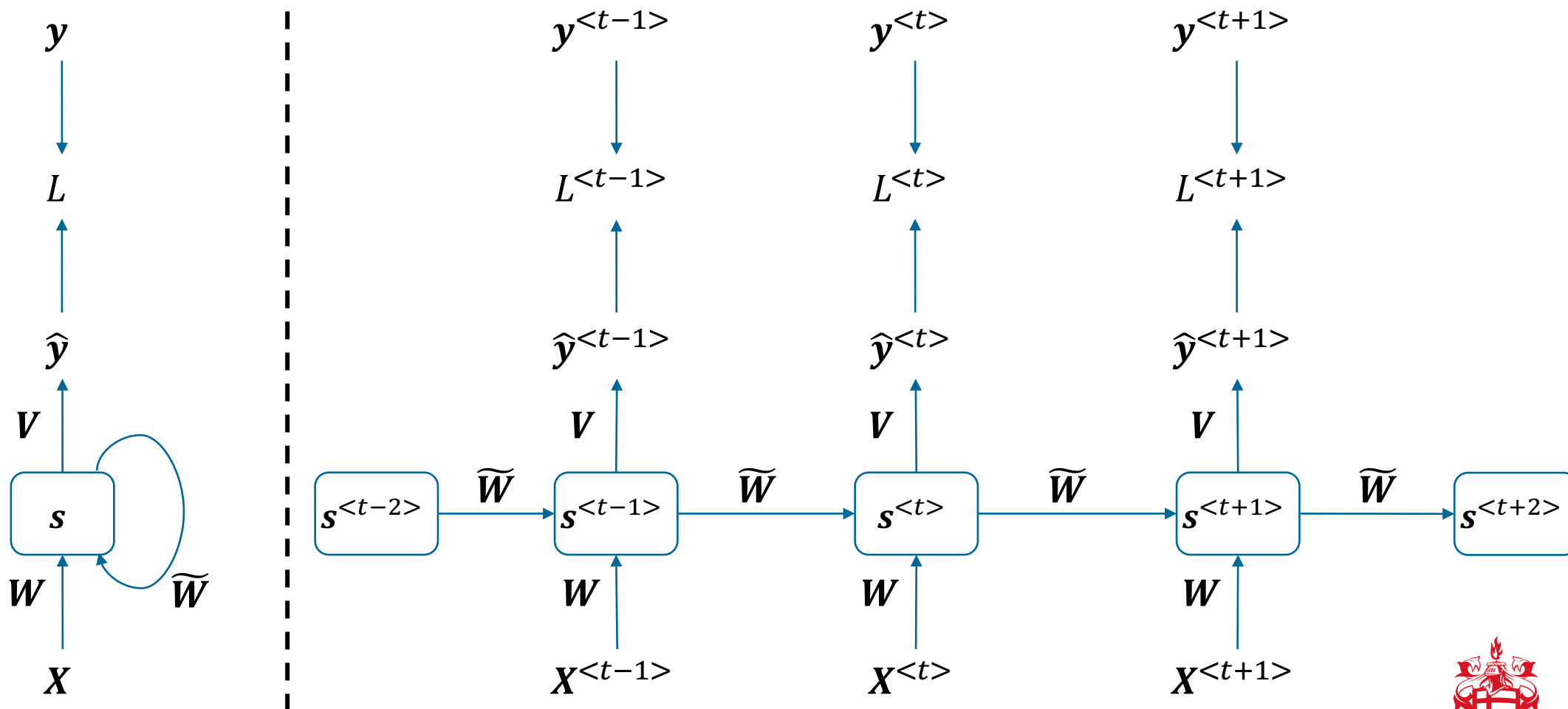
# Back-propagation through time

In principle, unfold the sequence to get to a computational graph, and use back-propagation

- "Back-propagation through time" algorithm (BPTT)

- Once computational graph has been established, can apply any of the known gradient-based optimization algorithms

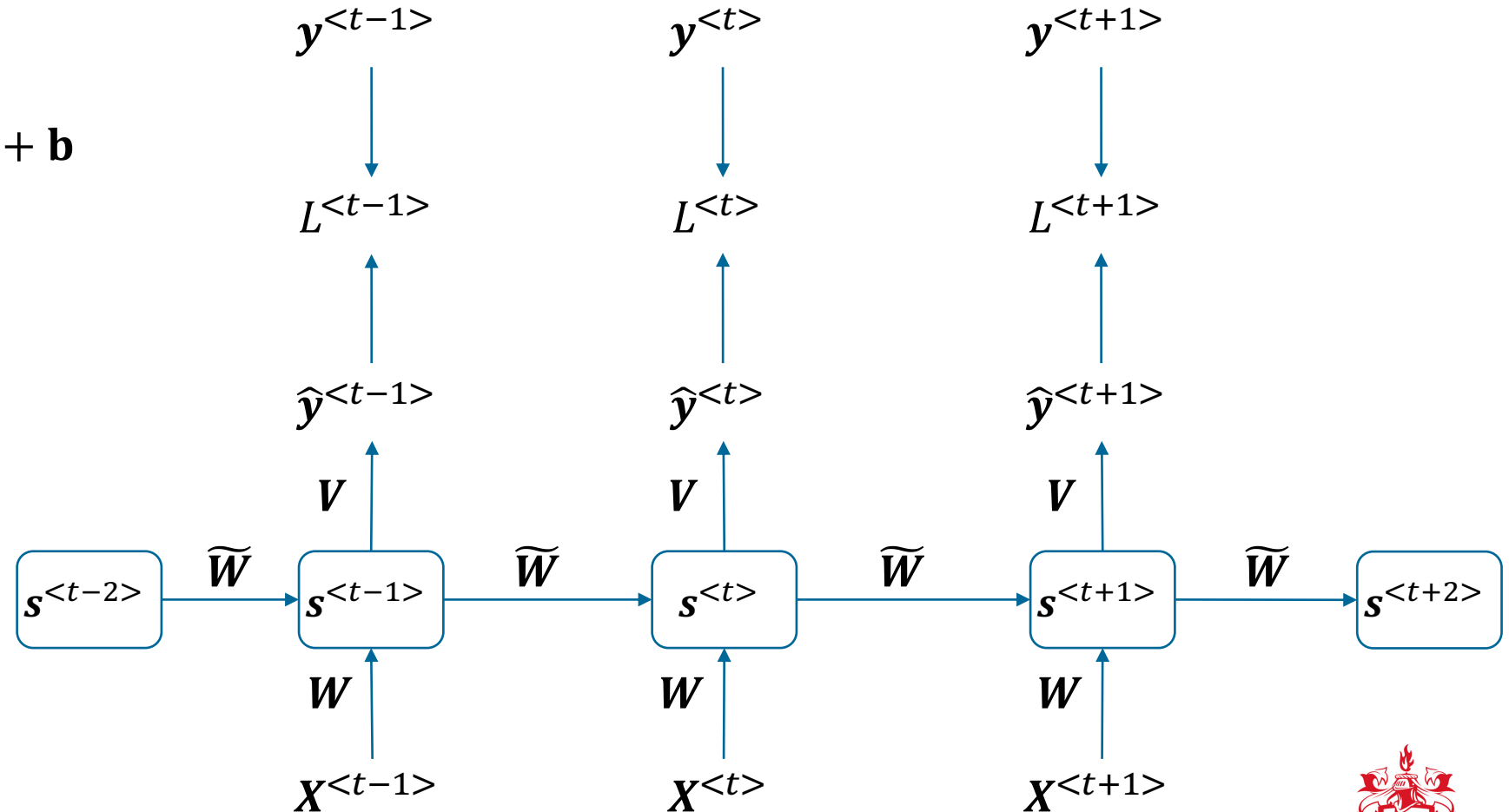$$a^{<t>} = x^{<t>}W + s^{<t-1>}\widetilde{W} + \mathbf{b}$$

$$s^{<t>} = f(a^{<t>})$$

$$\hat{y}^{<t>} = g(s^{<t>}V + c)$$

Gradient of loss to $\hat{y}^{<t>}$:

$$\frac{\partial L}{\partial \hat{y}^{<t>}} = \frac{\partial L}{\partial L^{<t>}} \frac{\partial L^{<t>}}{\partial \hat{y}^{<t>}}$$
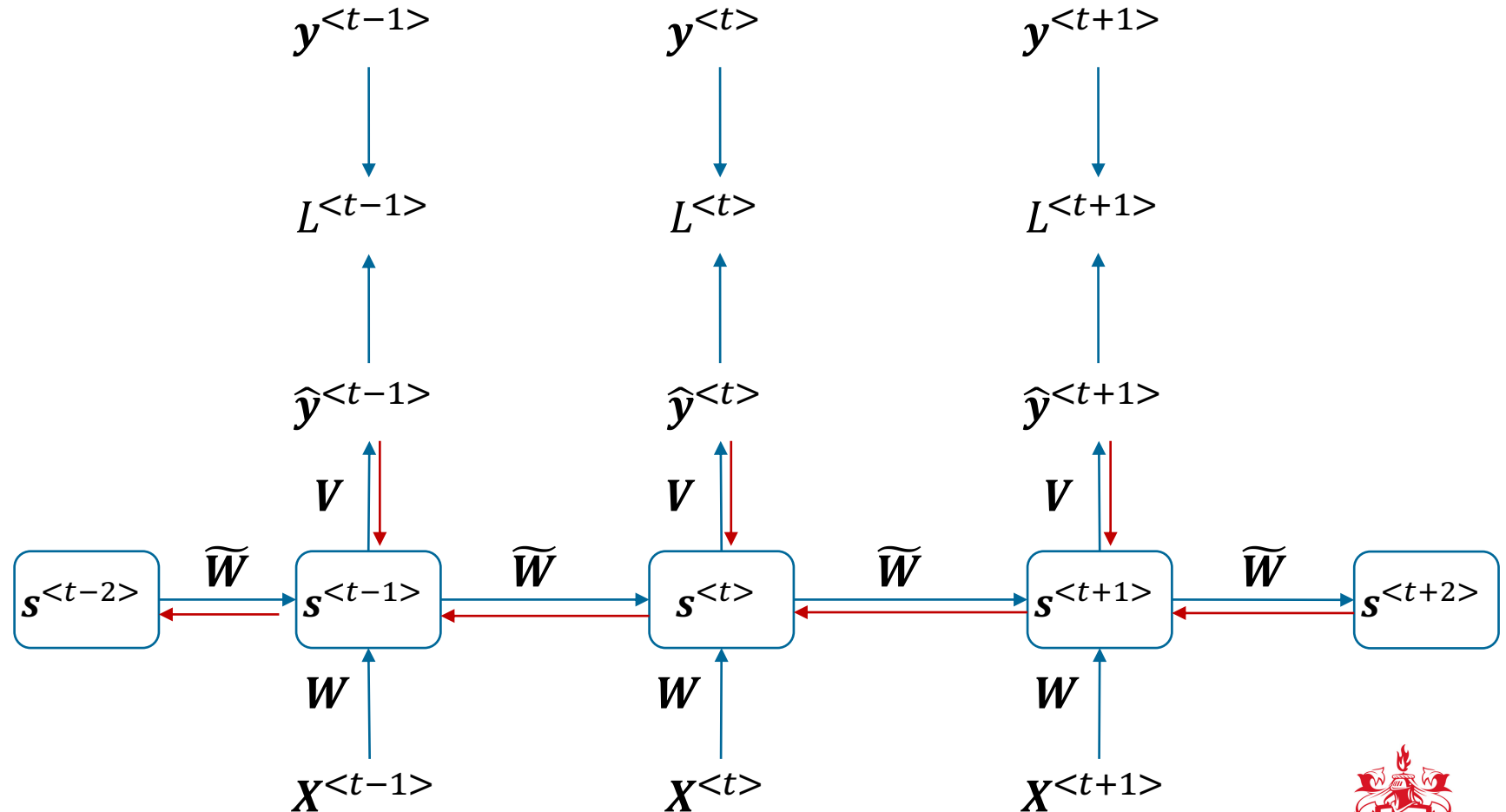
We have $L = \sum_t L^{<t>}$, so the gradient is 1

Gradient of loss to $s^{<t>}$:

$$\frac{\partial L}{\partial s^{<t>}} = \frac{\partial L}{\partial \hat{y}^{<t>}} \frac{\partial \hat{y}^{<t>}}{\partial s^{<t>}} \ ?$$
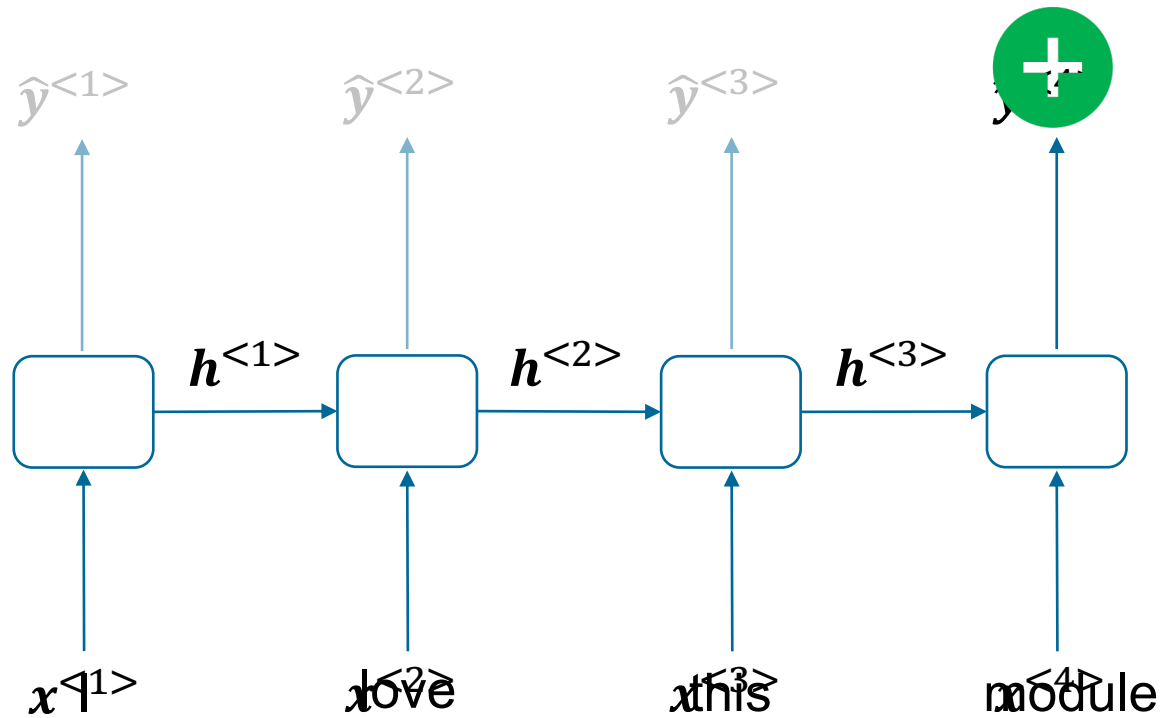
No, because future periods (and their losses) also depend on $s^{<t>}$ !

# RNN variants and their applications

For example:

- Video activity recognition

- DNA sequence probing

- Sentiment classification

# Vector-to-sequence networks



For example:

- Text generation

- Music generation

- Image captions

$\hat{y}^{<1>}$ Ich

$\hat{y}^{<2>}$ liebe

$\hat{y}^{<3>}$ dieses

$\hat{y}^{<4>}$ Modul

$h^{<1>}$ $h^{<2>}$ $h^{<3>}$

$x^{<1>}$ I

$x^{<2>}$ love

$x^{<3>}$ this

$x^{<4>}$ module

For example:

- Speech recognition

- Price predictions

- Translations

# Encoder-decoder networks



For example:

- Translations

- Dialogue

# Bidirectional RNNs – looking into the future

$\hat{y}^{<1>}$ Le    $\hat{y}^{<2>}$ Royaume    $\hat{y}^{<3>}$ Uni

$x^{<1>}$ The    $x^{<2>}$ United    Kingdom $x^{<3>}$

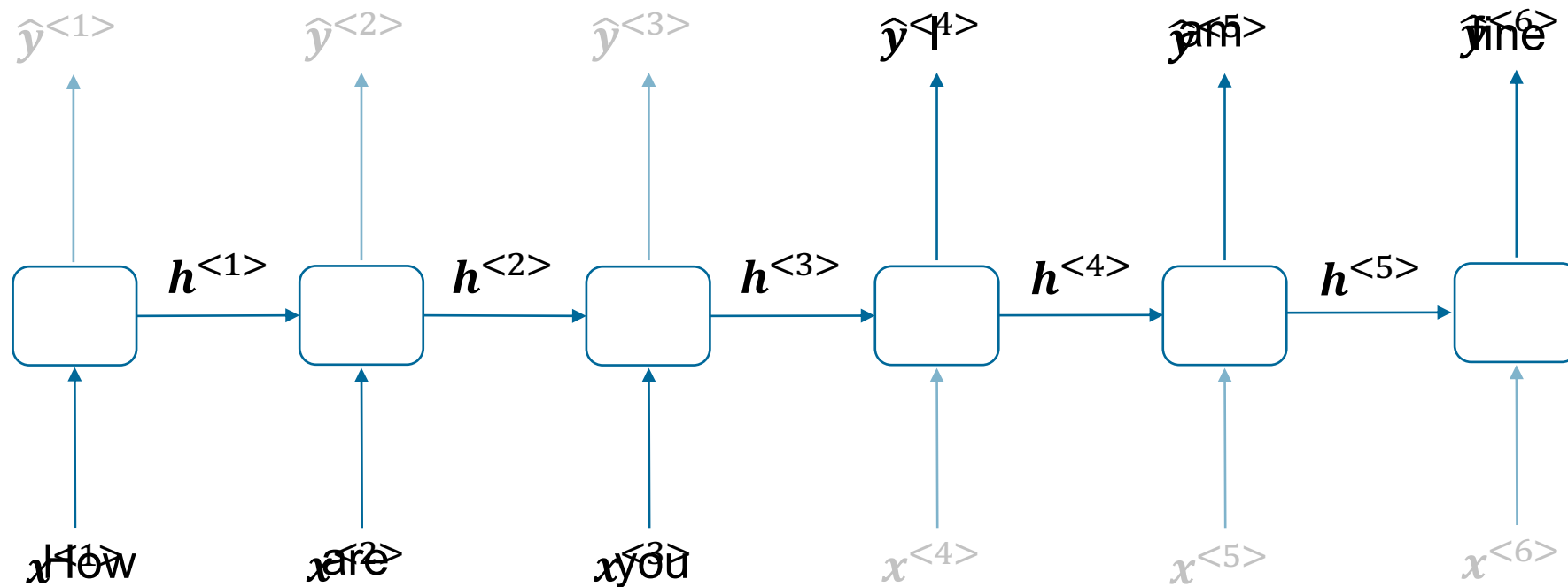For example:

- All sorts of NLP

- Also, in combination with the previous

# The issues with training RNNs

setup

# Problem 1 – vanishing and exploding gradients

- In principle, same as with other networks

- Before, we mostly focused on vanishing gradients
  → use of non-saturating activation functions such as ReLU

- With RNNs, exploding gradients become more of a problem
  - Same weights used for different time steps can lead to self-reinforcing increases of gradients
  → We frequently use saturating activation functions, such as tanh, or other methods such as gradient clipping

- Vanishing gradients are still a problem (sometimes even more so than in other networks):



Now        fair        Hippolyta            revenue

- This is essentially a very very deep neural network!
  - Some information is lost at each time step

- After just a few time steps, there is virtually no more information about the first input

The BA students, which had been working for days on end, was finally done with their project.

The BA students, which had been working for days on end, <span style="color:red">was</span> finally done with their project.

See you in class!

# Sources

- DeepLearning.AI, n.d.: deeplearning.ai
- Garnelo, 2020, Lecture 6: Sequences and Recurrent Networks: https://storage.googleapis.com/deepmind-media/UCLxDeepMind_2020/L6%20-%20UCLxDeepMind%20DL2020.pdf
- Géron, 2019, Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow
- Goodfellow, Bengio, Courville, 2016, The Deep Learning Book: http://www.deeplearningbook.org
- Liang, 2016, Introduction to Deep Learning: https://www.cs.princeton.edu/courses/archive/spring16/cos495/
- Soleimany, 2022, Deep Sequence Modeling: http://introtodeeplearning.com/slides/6S191_MIT_DeepLearning_L2.pdf

BAYES
BUSINESS SCHOOL
CITY, UNIVERSITY OF LONDON