

Recurrent Neural Networks

Sudong Lee

 sudonglee@ulsan.ac.kr

 <https://dais.ulsan.ac.kr/>

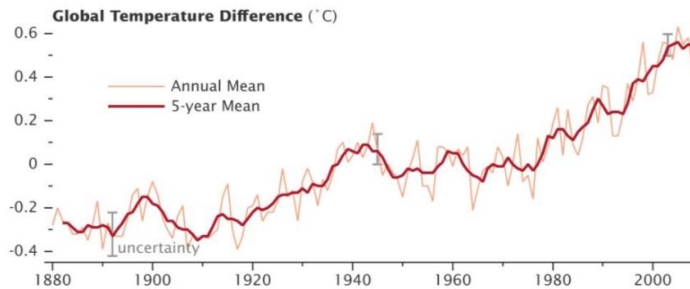
Contents

- Introduction to Sequential Data
- Vanilla Recurrent Neural Network
- LSTM (Long Short-Term Memory)

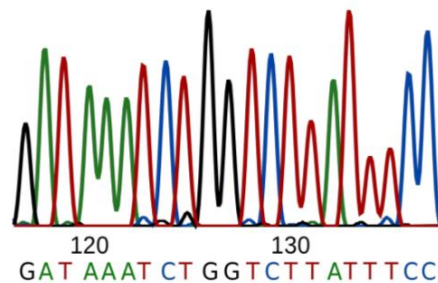
Introduction to Sequential Data

What is Sequential Data?

“Sequential data refer to data in which the order of instances matters.”



Time Series

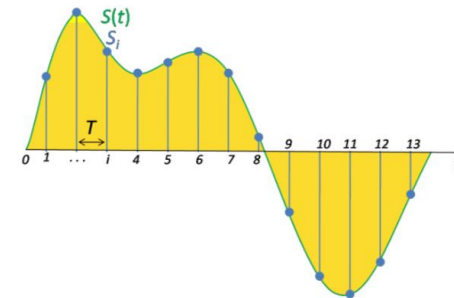


DNA sequencing

“I am a boy.”

“Am I a boy?”

Natural language



Audio sampling

What is Sequential Data?

“Sequential data have the temporal dependency, so the order of instances matters.”

- Notations

- Non-sequential:

$$\mathbf{X}^T = [x_{i,j}] \text{ for } i = 1, \dots, n; j = 1, \dots, p. \quad (n \times p)$$

- Sequential:

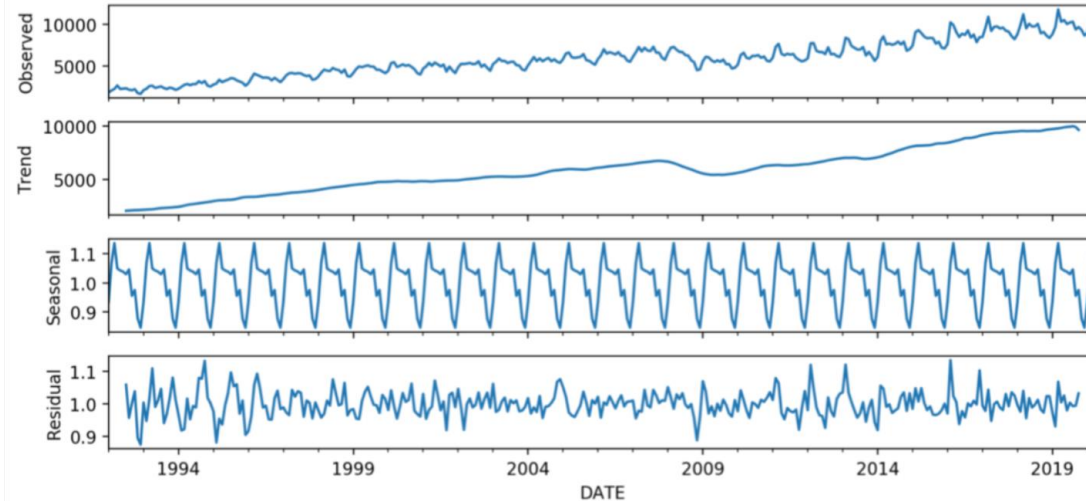
$$\mathbf{X}^T = [x_i(t)] \text{ for } i = 1, \dots, p; t = 1, \dots, T \quad (T \times p)$$

$$\mathbf{x}(t) = [x_1(t), \dots, x_p(t)]^T$$

$$\mathbf{X} = [\mathbf{x}(1), \dots, \mathbf{x}(T)]$$

Types of Sequential Data: Time Series

“Time series data is a sequence of data points collected or recorded at regular intervals.”



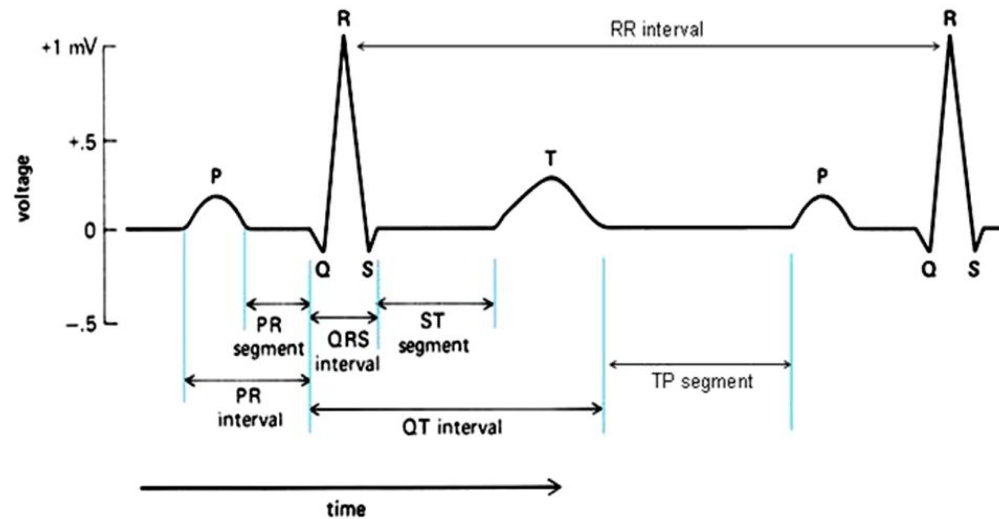
Decomposition of a used car sales

Characteristics of Time Series Data

- **Temporal ordering**
data points are arranged chronologically.
- **Regular intervals**
collected at consistent time intervals (e.g., hourly, daily, ...)
- **Dependency**
each data point is related to or influenced by previous points.

Types of Sequential Data: Signal

“Signal data imply a continuous or discrete series of measurements over time or space.”



ECG (Electrocardiogram) signal data

Characteristics of Signal Data

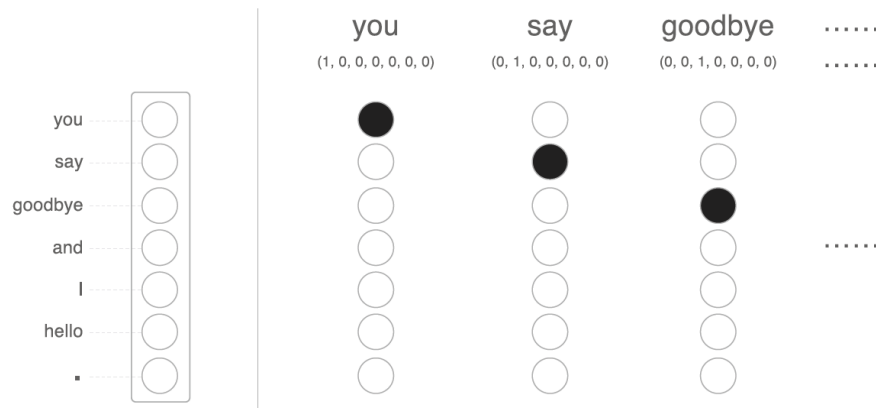
- **Time-varying**
represents how a quantity changes over time.
- **Continuous or discrete**
continuous or discrete in nature.
- **Amplitude and frequency**
characterized by their amplitude (strength) and frequency (rate of change).

Types of Sequential Data: Natural Language

“Natural language data expressed in human language are encoded as vectors.”

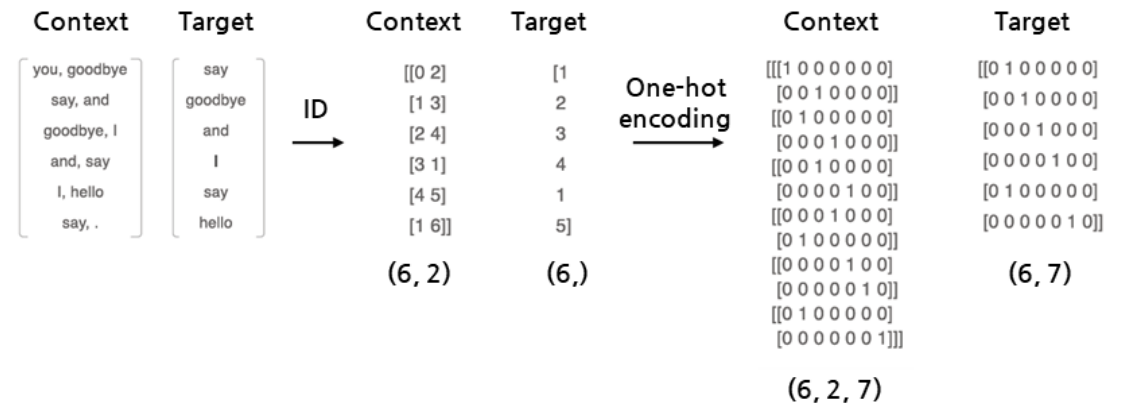
Corpus: “you say goodbye and I say hello.”

Word (text)	ID (label)	One-hot encoding
$\begin{pmatrix} \text{you} \\ \text{goodbye} \end{pmatrix}$	$\begin{pmatrix} 0 \\ 2 \end{pmatrix}$	$\begin{pmatrix} (1, 0, 0, 0, 0, 0, 0) \\ (0, 0, 1, 0, 0, 0, 0) \end{pmatrix}$



One-hot encoding of words using a bag-of-words

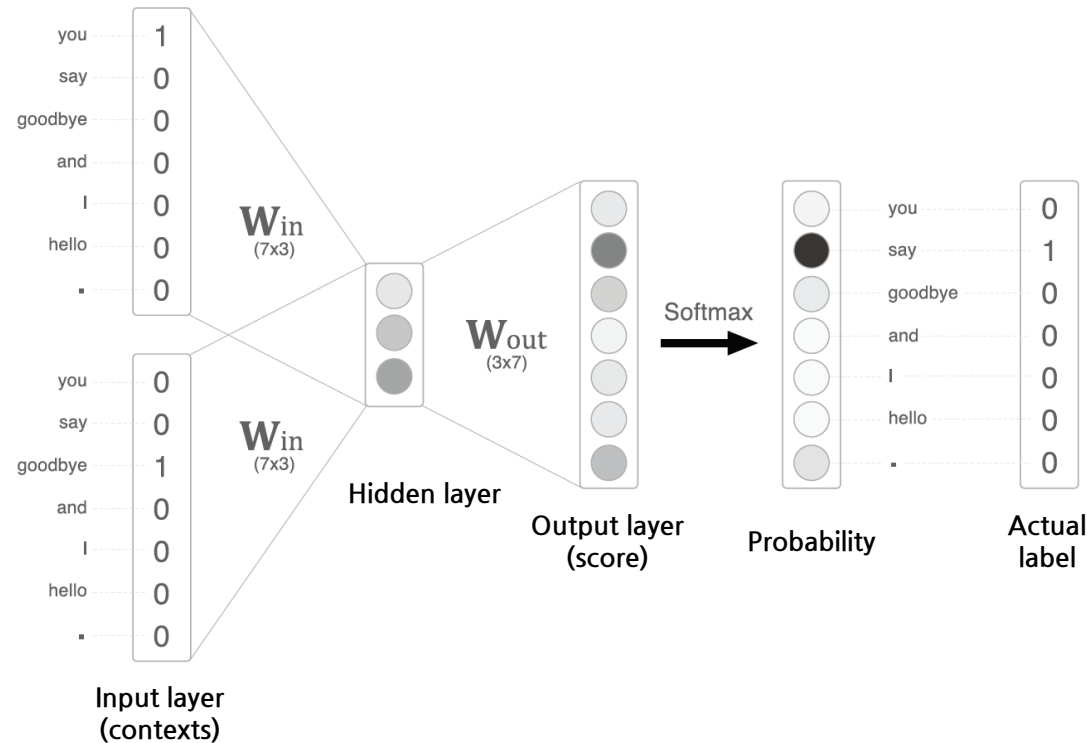
Inference: “you [?] goodbye and I say hello.”



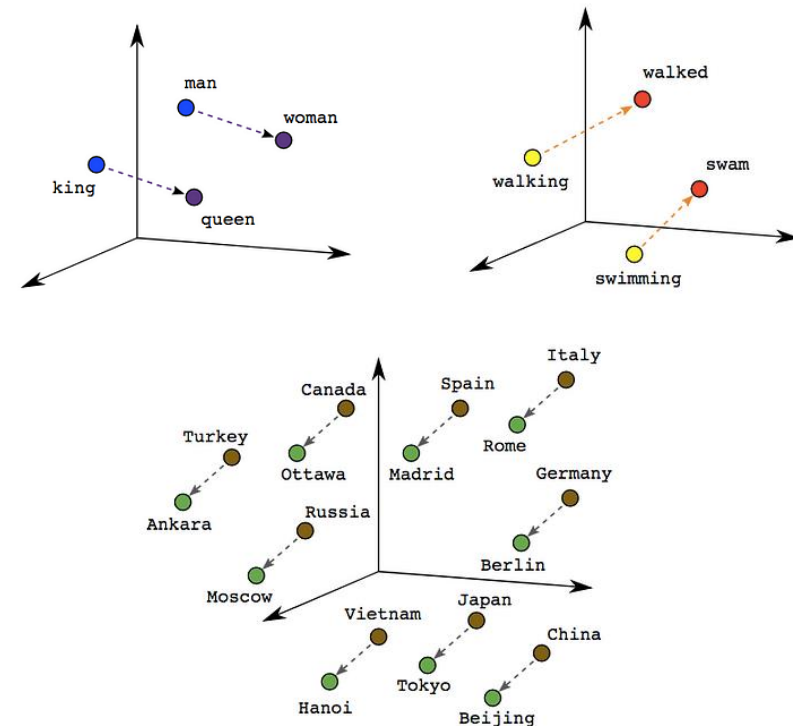
One-hot encoding of words using CBOW (continuous bag of words)

Types of Sequential Data: Natural Language

“Natural language data expressed in human language are encoded as vectors.”



word2vec using CBOW

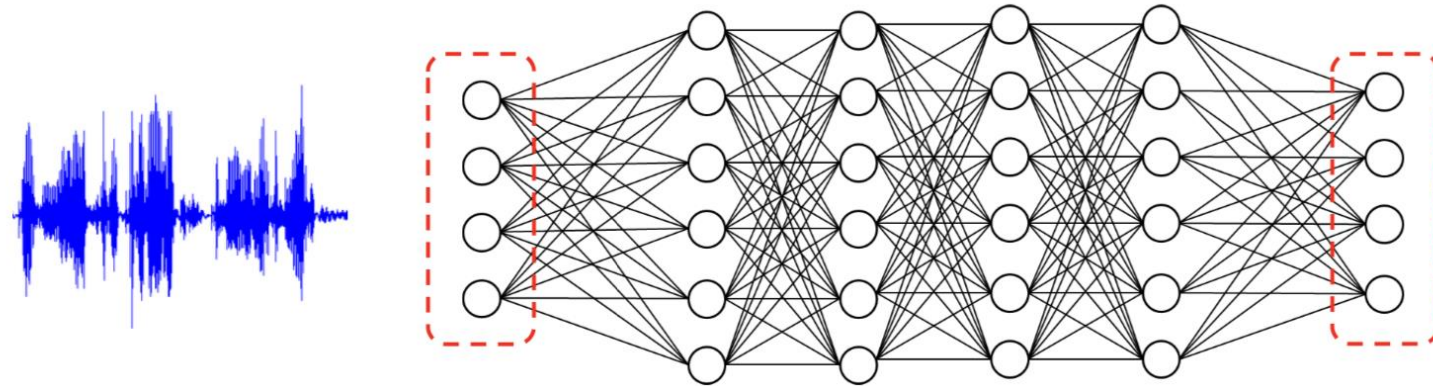


Word embeddings by word2vec

Vanilla Recurrent Neural Network

Representation of Sequential Data

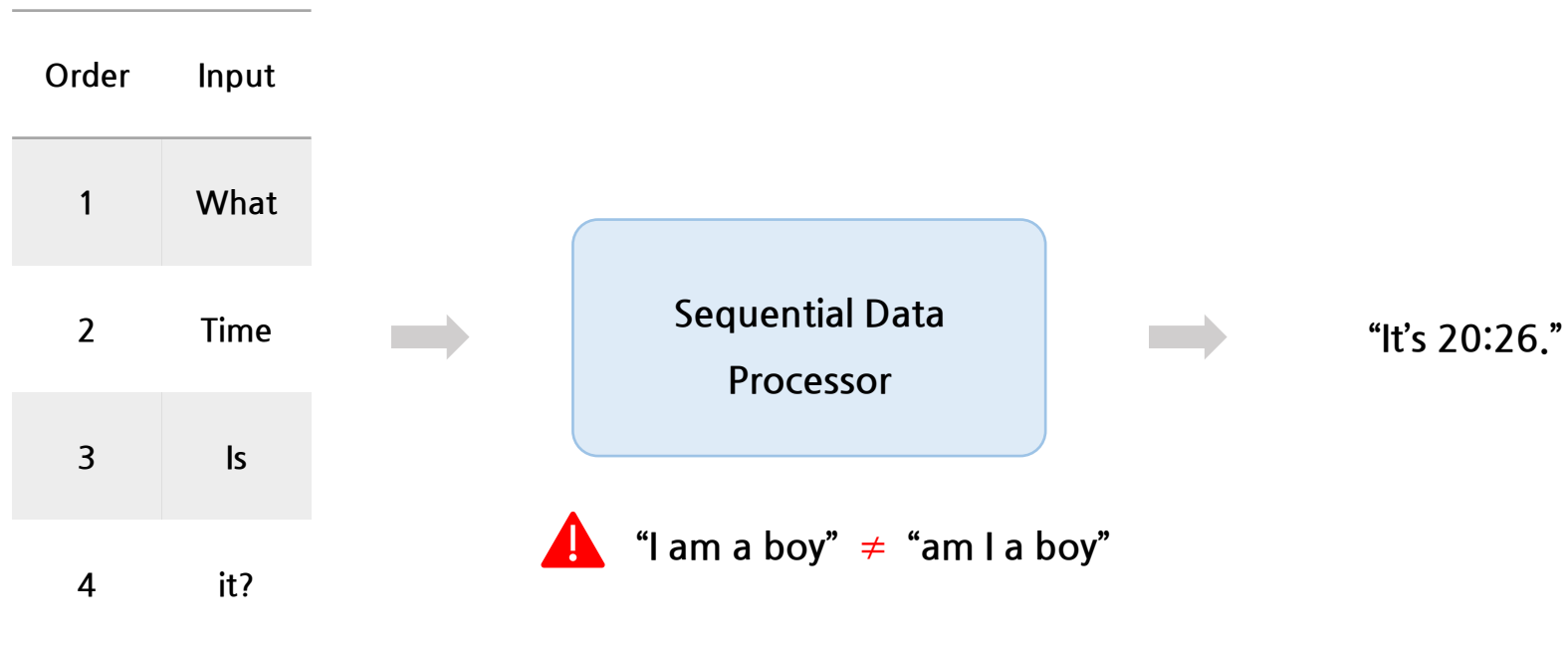
“How can we effectively extract the representation of sequential data using deep learning?”



“I am a boy” \neq “am I a boy”

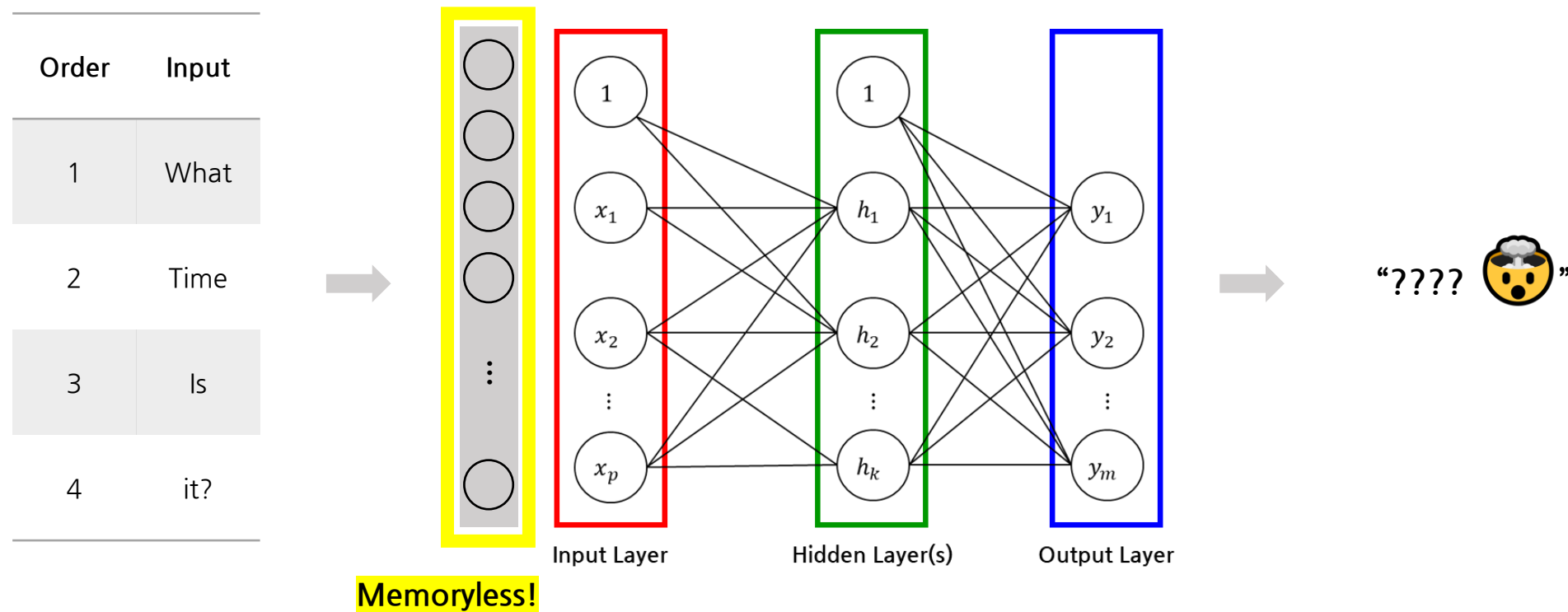
Representation of Sequential Data

“How can we effectively extract the representation of sequential data using deep learning?”



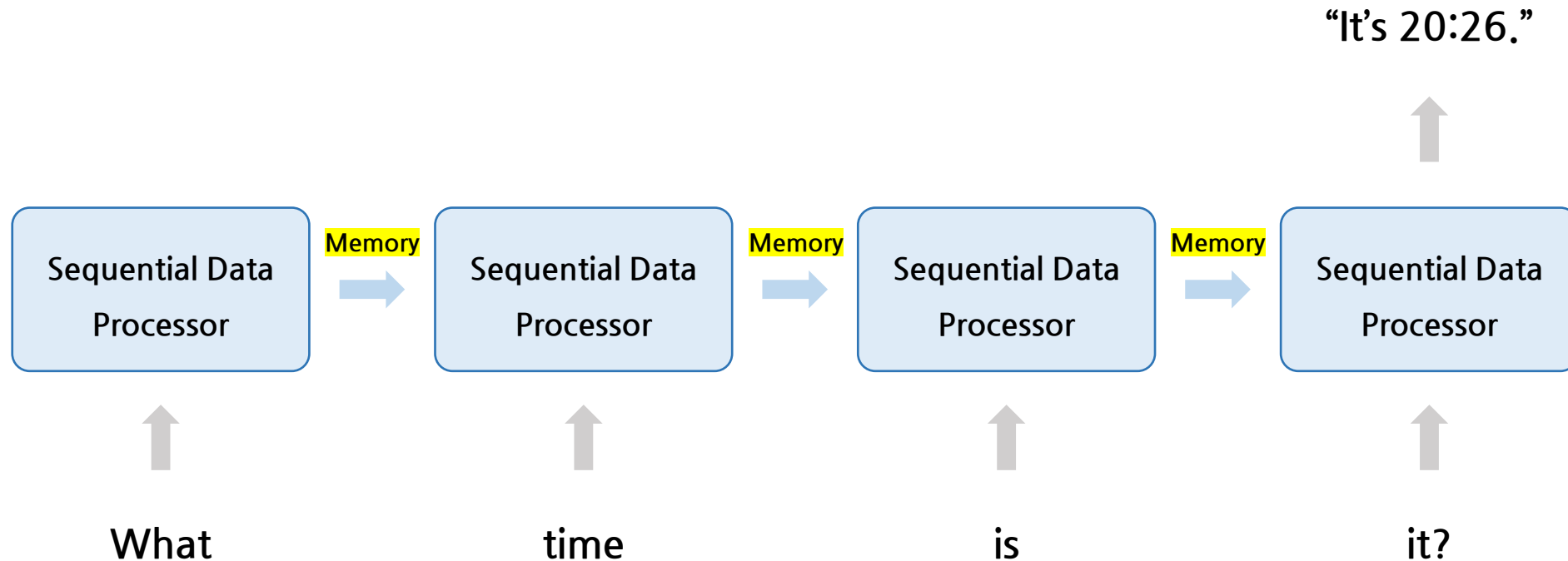
Limitation of Vanilla Neural Network for Sequential Data

“Independent nodes in the input layer are memoryless of the input sequence.”



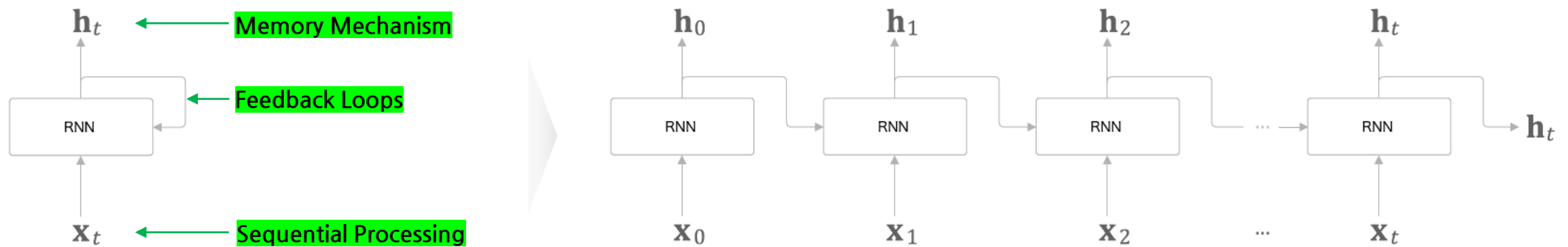
Memory System

“A sequential data processor is a memory system that can memorize the input order.”



Recurrent Neural Network (RNN)

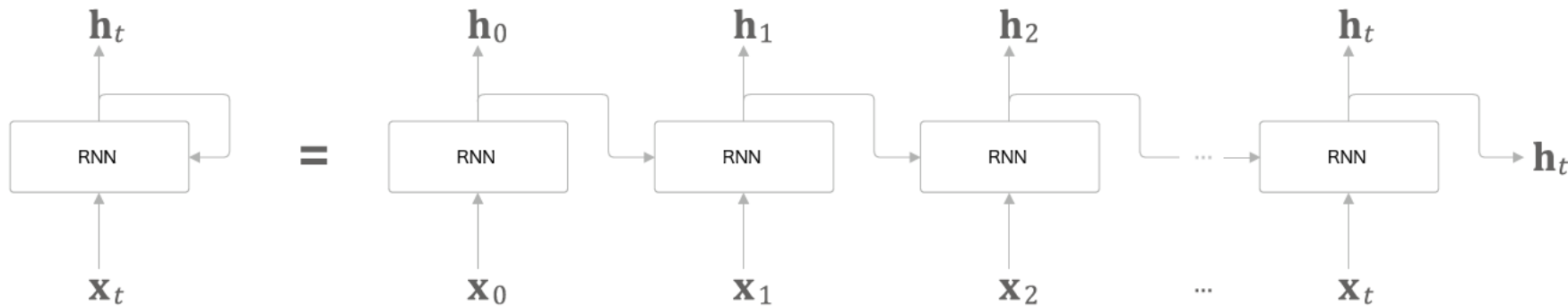
“A recurrent neural network **recurrently** applies a sequence processor using *hidden states*.”



- x_t : the input vector at time t .
- h_t : the hidden state vector at time t .

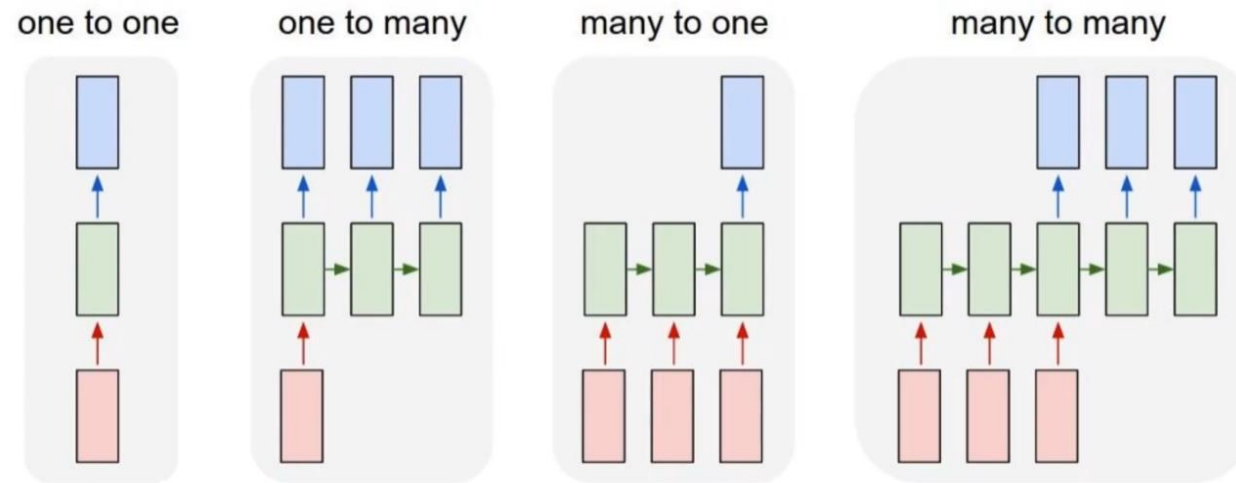
Recurrent Neural Network (RNN)

“A recurrent neural network **recurrently** applies a sequence processor using *hidden states*.”



- $h_t = \tanh(h_{t-1}^T W_{hh} + x_t W_{xh} + b_h)$
- $y_t = \tanh(h_t W_{hy} + b_y)$

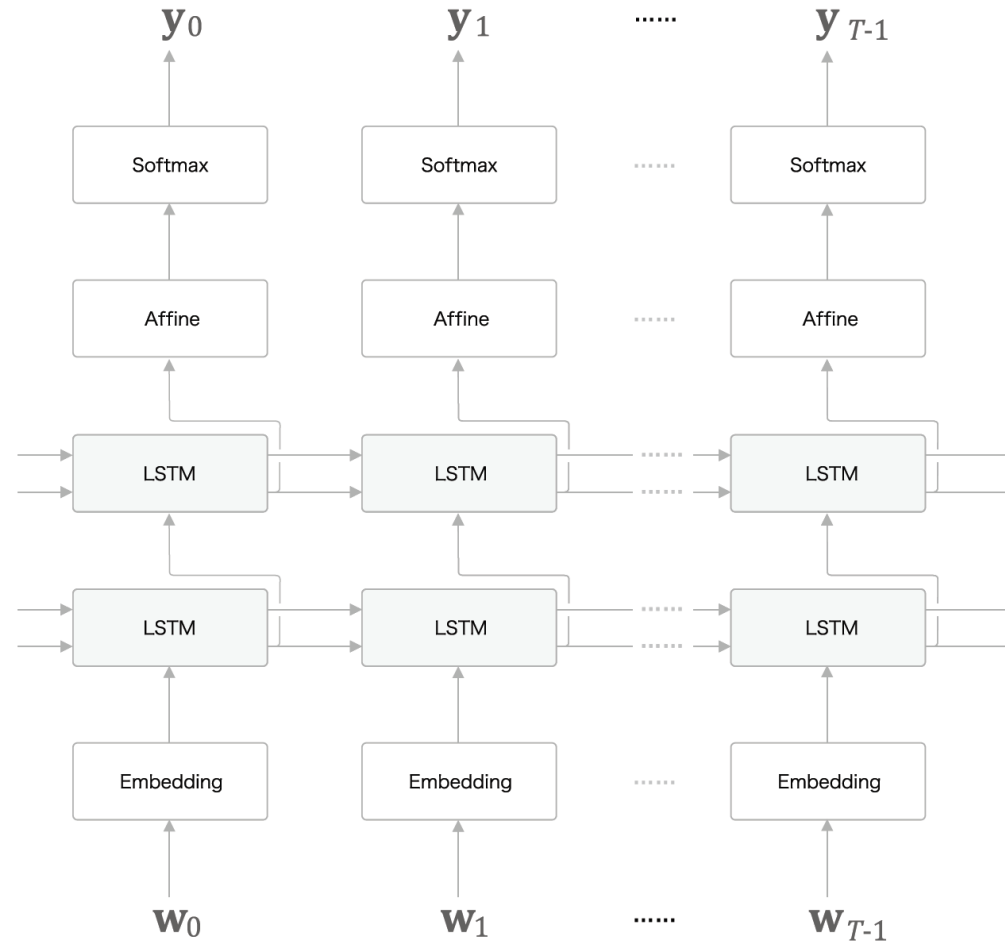
Various Usage of RNN



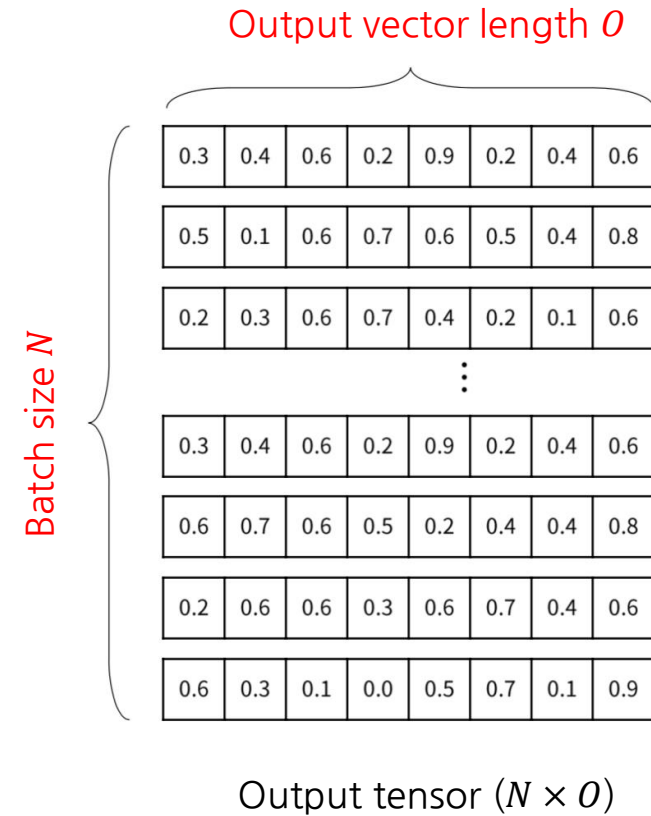
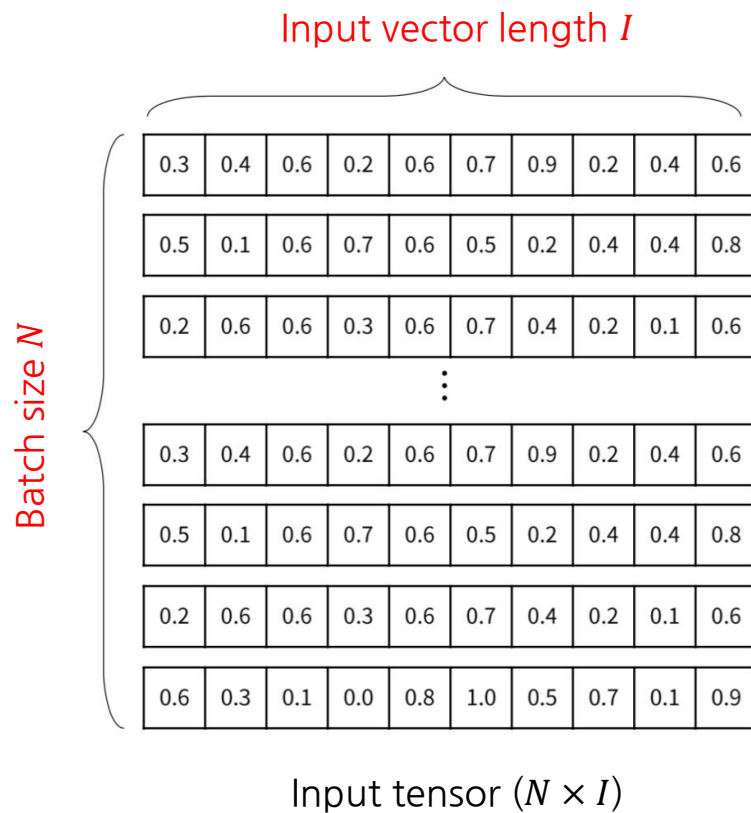
Examples

- ① One-to-one: MLP
- ② One-to-many: description of an input image
- ③ Many-to-one: time series classification
- ④ Many-to-many: machine translation

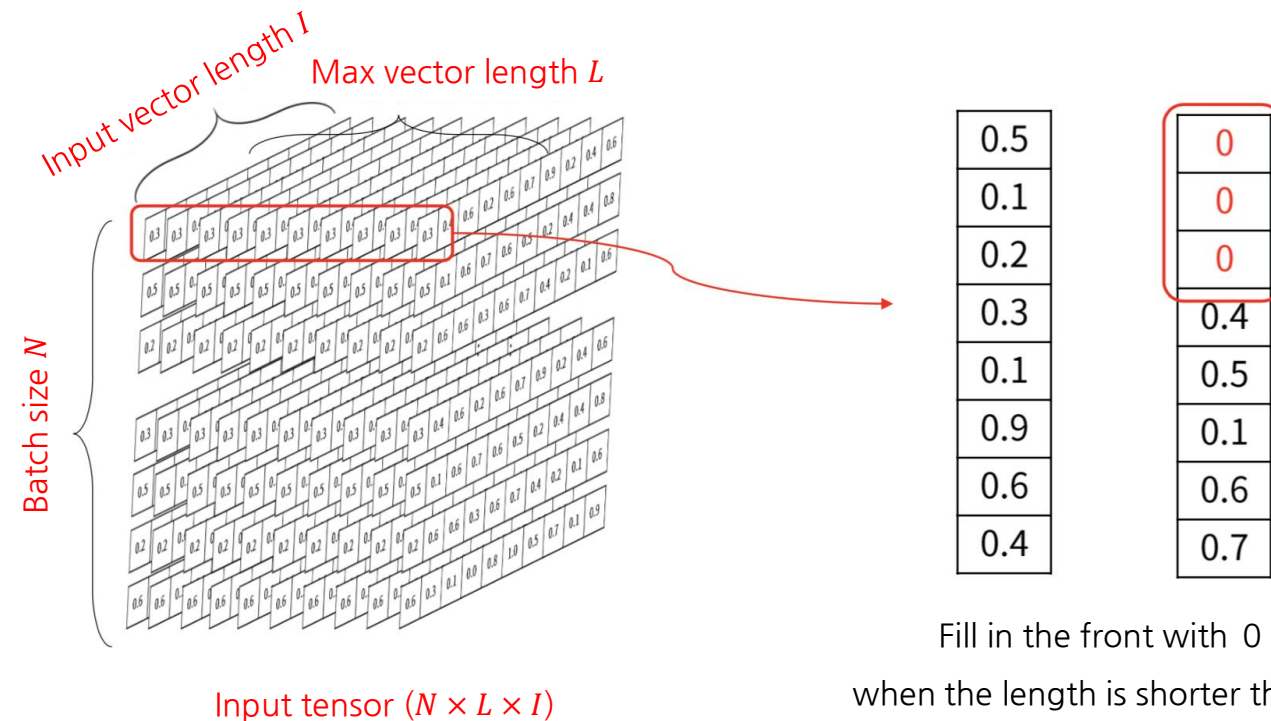
Example: A Multi-layered LSTM for the Many-to-Many Problem



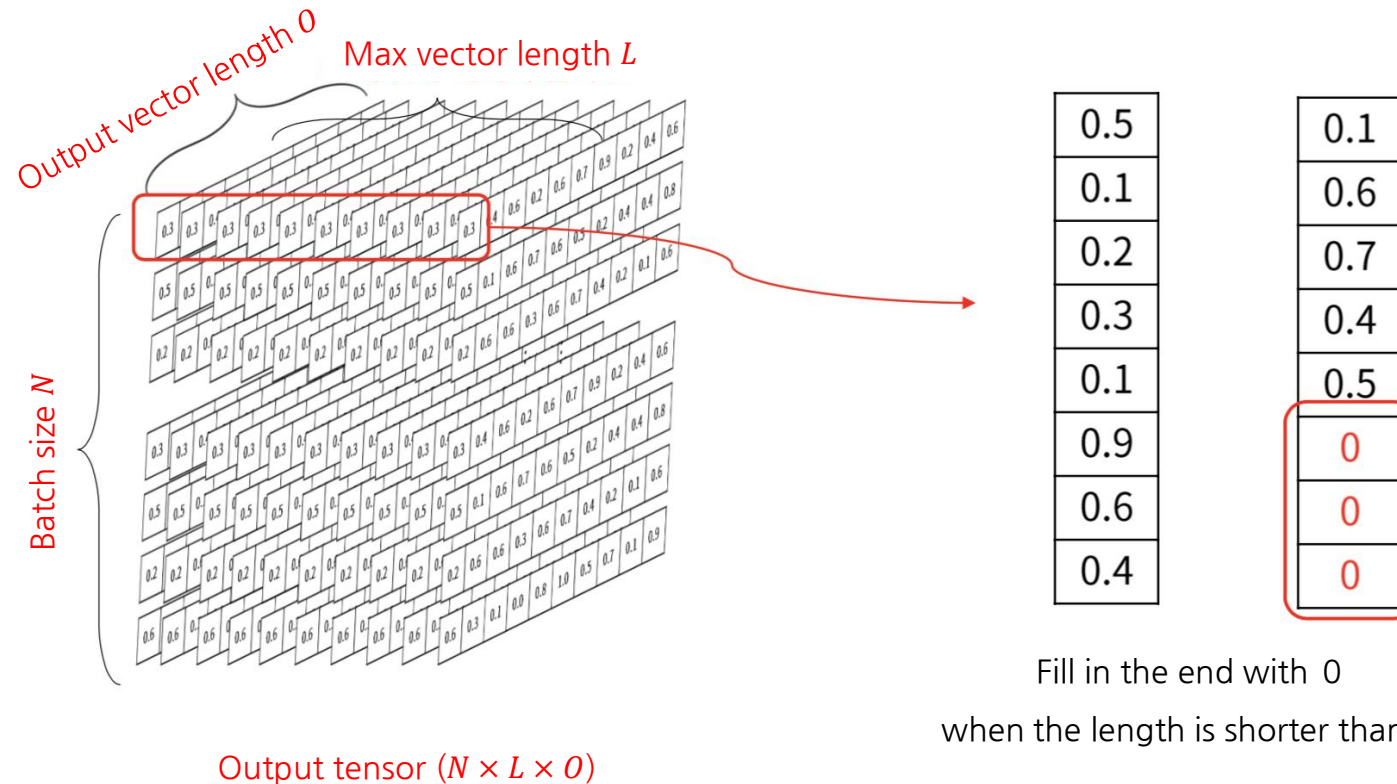
Understanding RNNs from a Tensor Perspective



Understanding RNNs from a Tensor Perspective

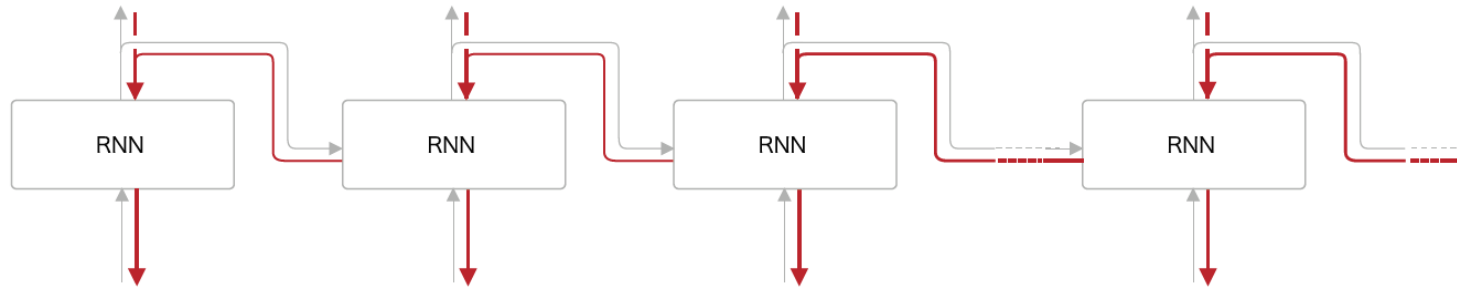


Understanding RNNs from a Tensor Perspective



BPTT: Backpropagation Through Time

“BPTT is an extension of the standard backpropagation algorithm adapted for the recurrent structure of RNN.”



Challenges and Limitations of BPTT

- **Vanishing and exploding gradients:**

As gradients are propagated back through many time steps, they can become very small (vanishing) or very large (exploding).

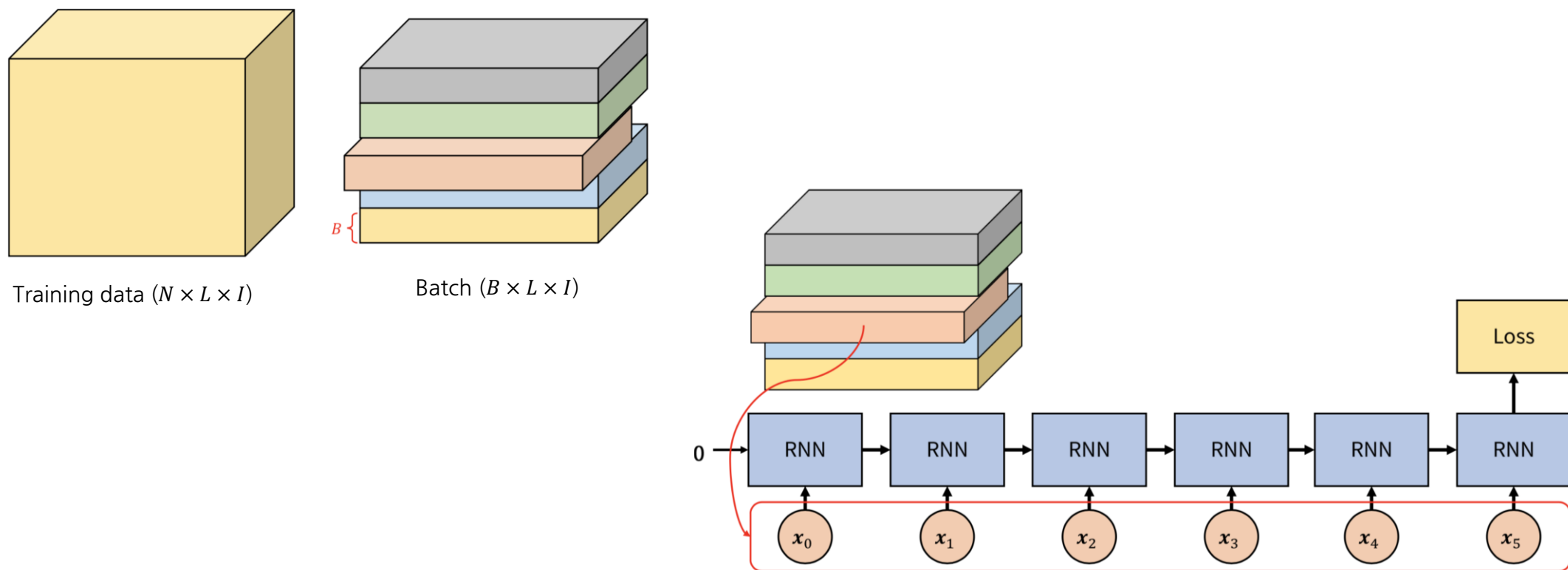
- **Computational complexity:**

Computationally expensive, especially for long sequences, as it requires storing activations for all time steps.

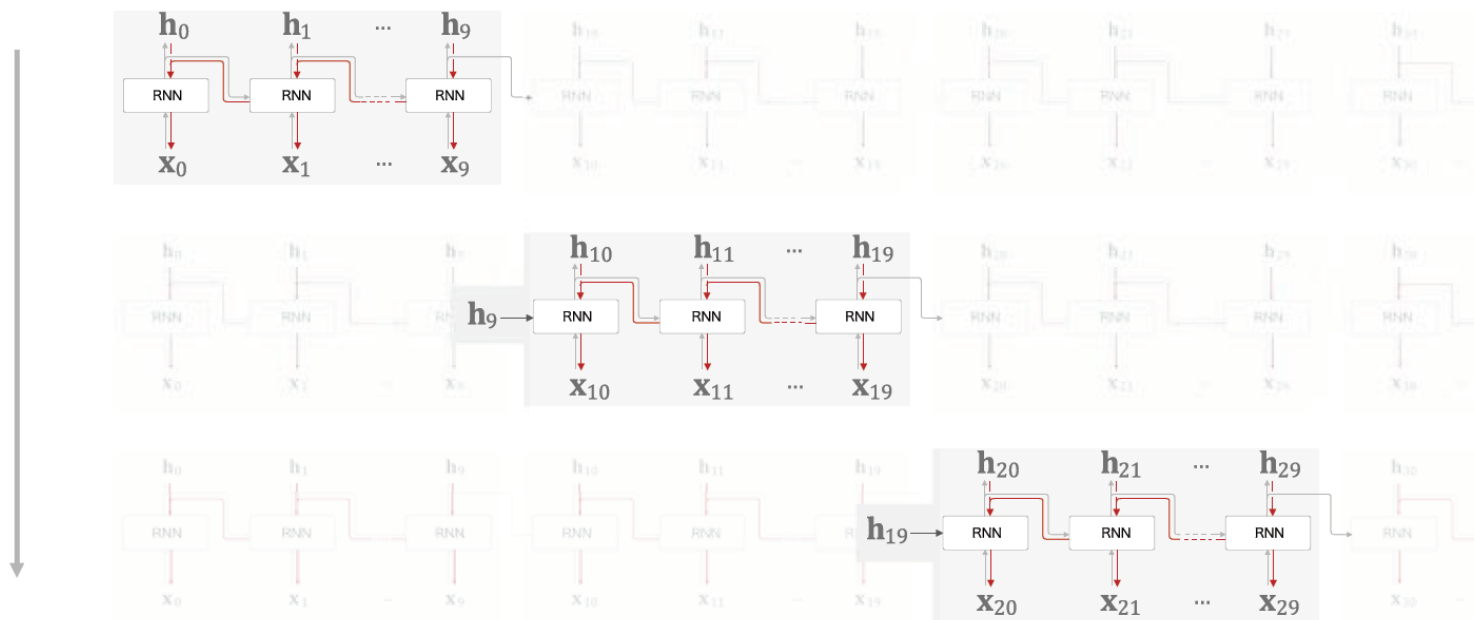
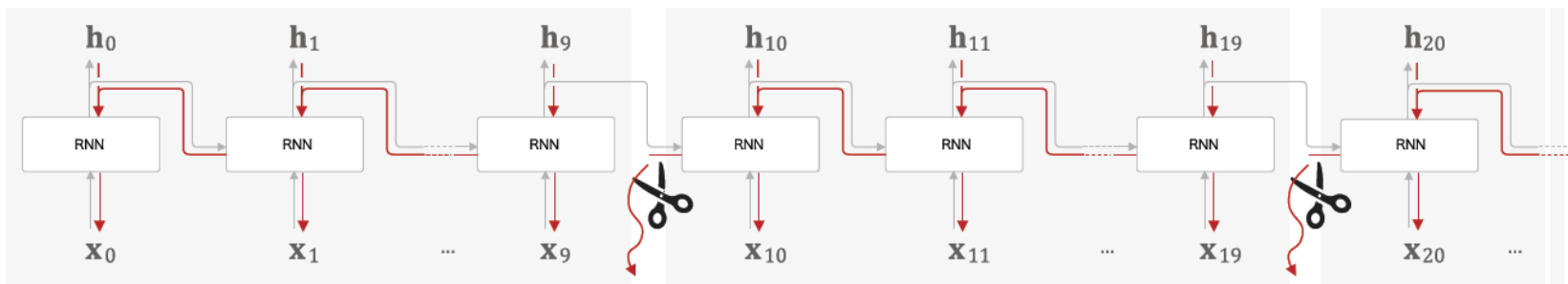


In practice, BPTT is often truncated to a certain number of time steps to manage computational costs, which can limit the network's ability to learn long-term dependencies.

BPTT: Backpropagation Through Time

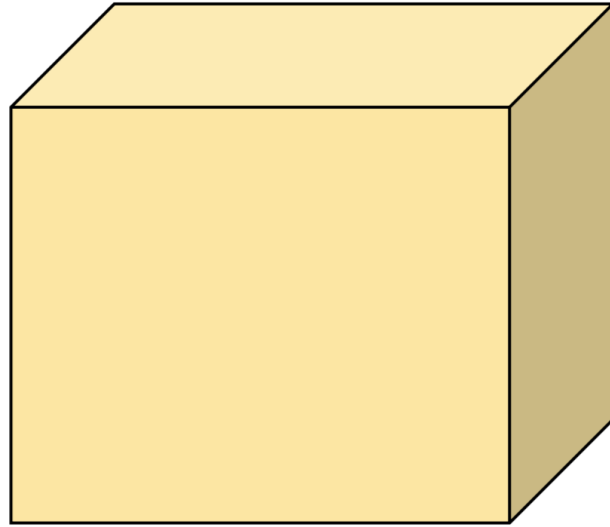


Truncated BPTT

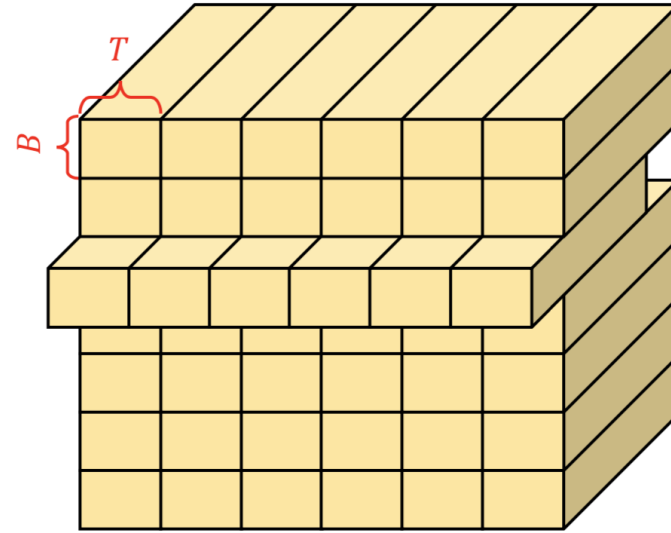


Truncated BPTT

“Truncates the length of sequential data to a constant T length to reduce the size of the computation at once.”



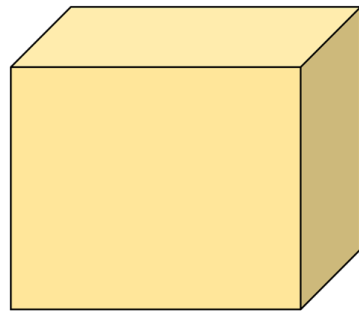
Training data ($N \times L \times I$)



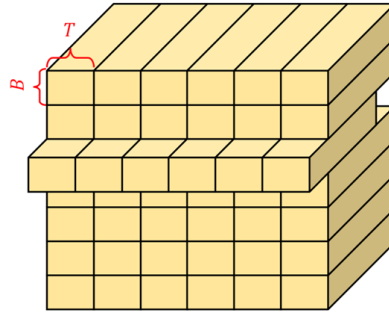
Truncated batch ($B \times T \times I$)

Truncated BPTT

“Truncates the length of sequential data to a constant T length to reduce the size of the computation at once.”

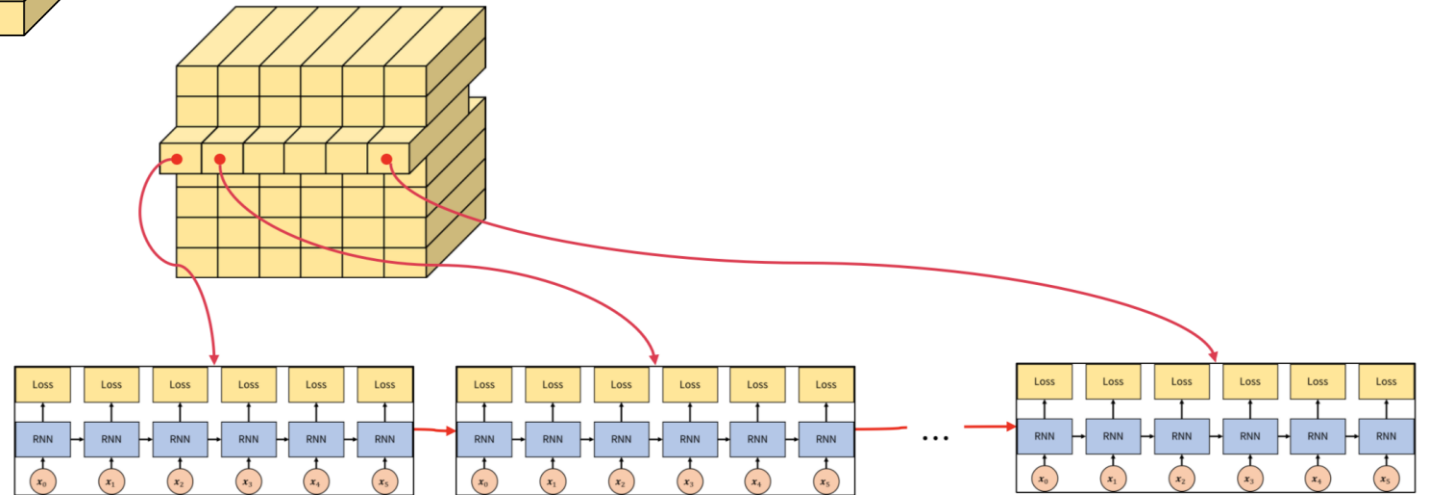
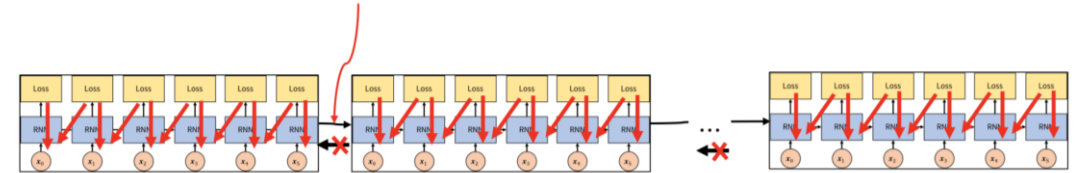


Training data ($N \times L \times I$)



Truncated batch ($B \times T \times I$)

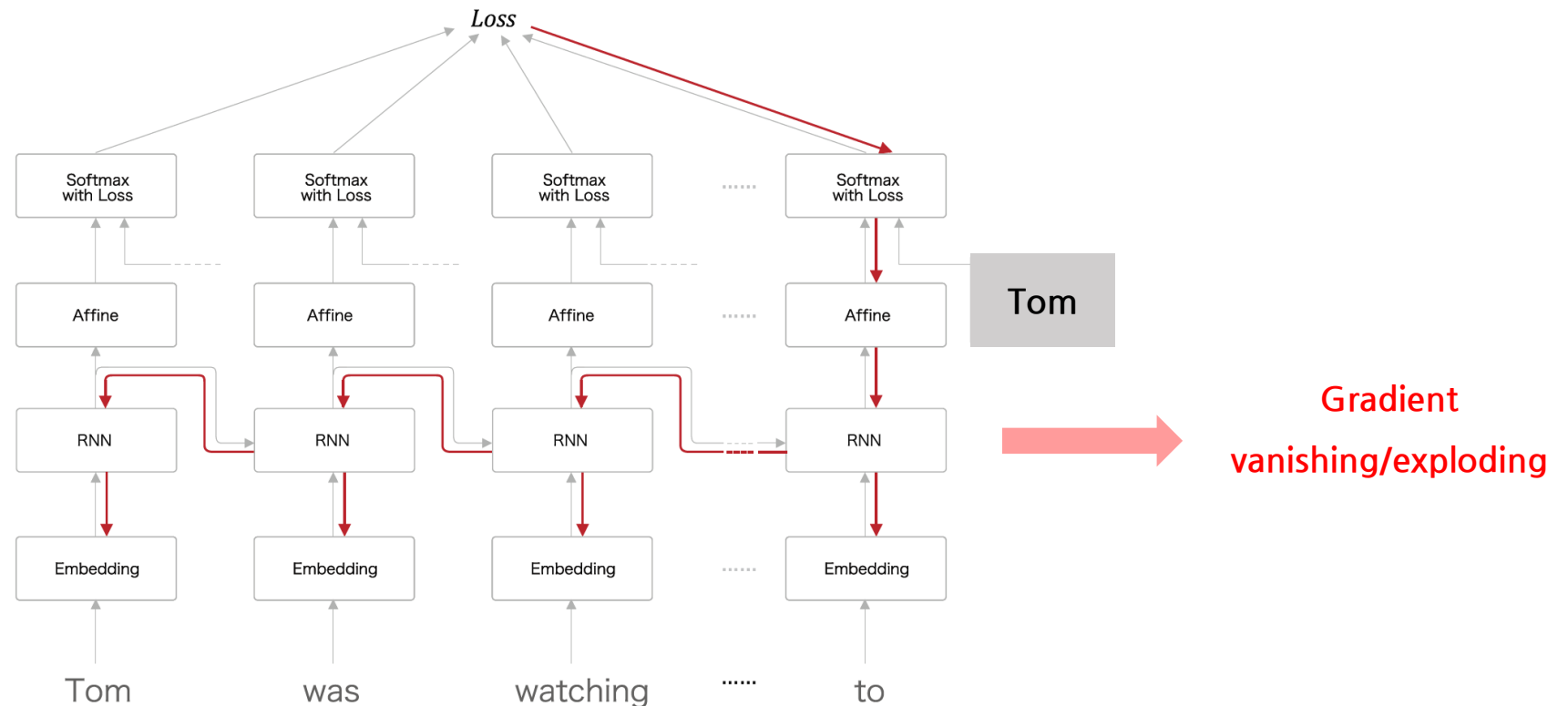
Forward propagation of hidden and cell states



Limitation of Vanilla RNNs

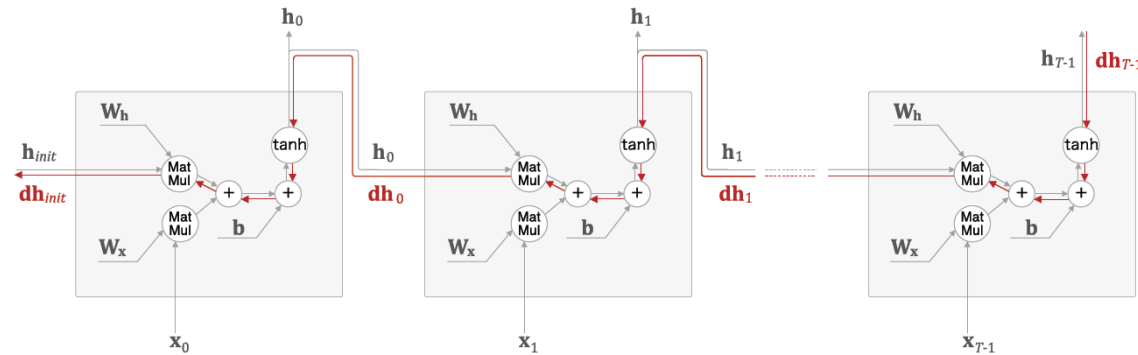
“Vanilla RNNs are limited in their ability to learn long-term memory due to **computational cost** and **gradient issues**.”

“Tom was watching TV in his room. Mary came into the room. Mary said hi to [?].”

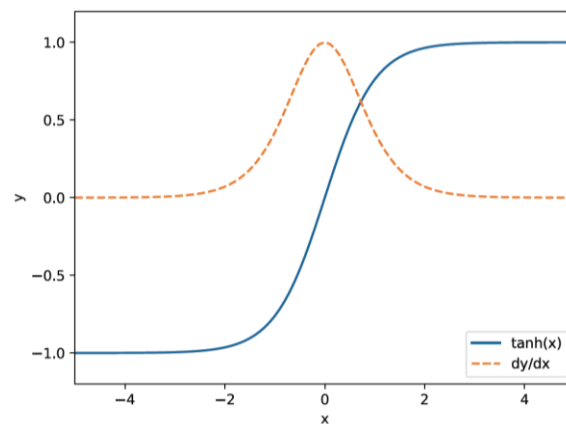


Causes of Gradient Vanishing and Exploding

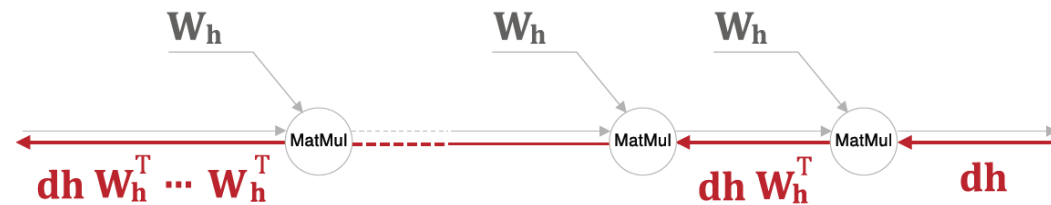
“The structure of RNN inherently causes gradient problems.”



Source ① *tanh*



Source ② *MatMul*

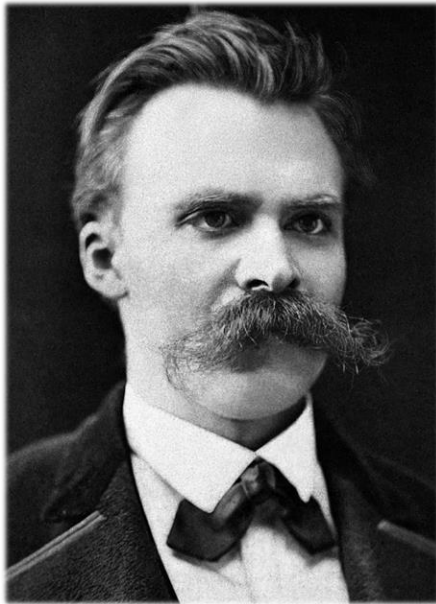


How can we memorize the long-term dependency?

LSTM (Long Short-Term Memory)

Introduction to LSTM

“How can we memorize the long-term dependency?”

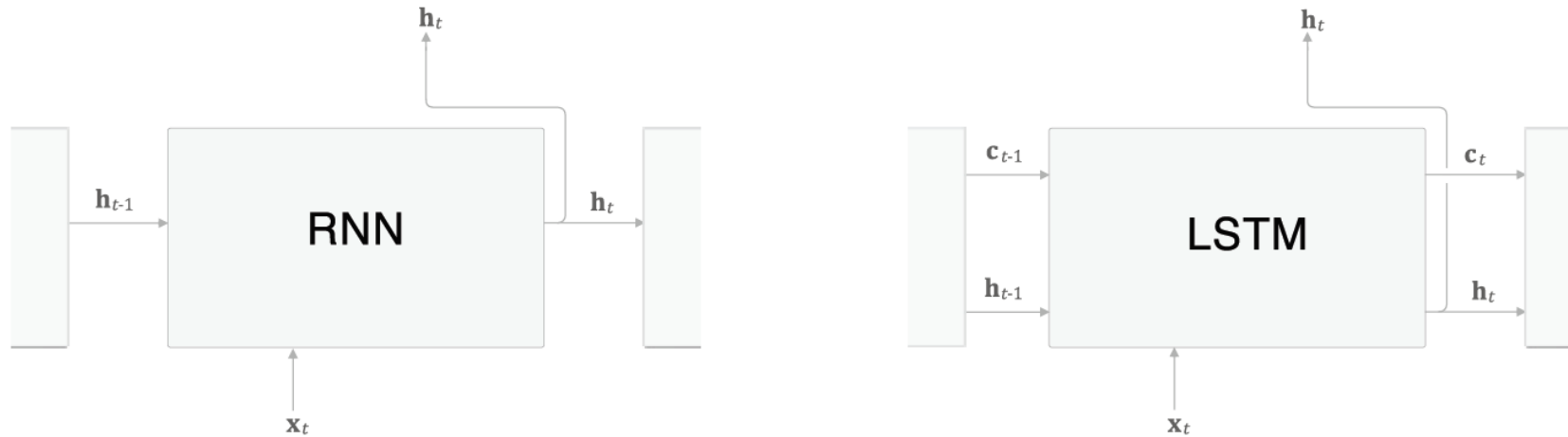


Friedrich Nietzsche

“Blessed are the forgetful
for they get the better even of their blunders.”
망각하는 자는 복이 있나니,
자신의 실수조차 잊기 때문이라.

Introduction to LSTM

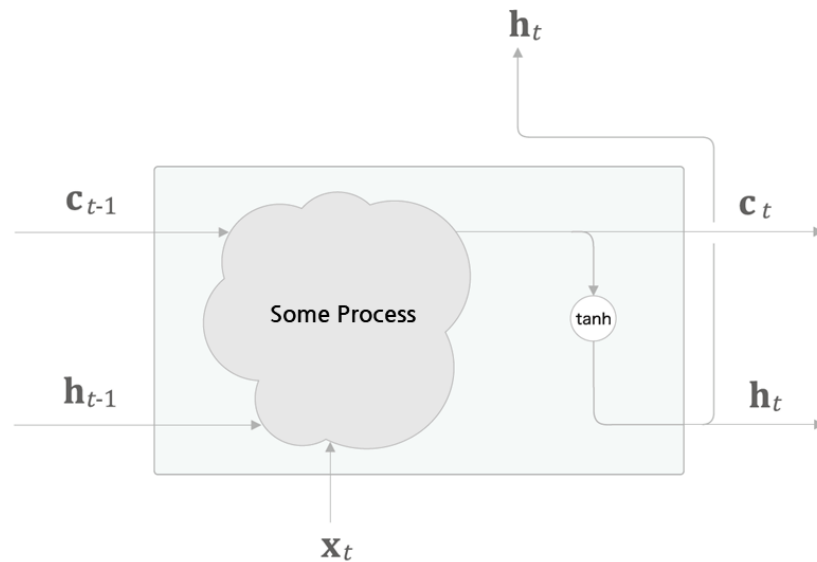
“LSTM (Long Short-Term Memory) handles long-term dependencies in sequential data using the **memory cell** and **gates**.”



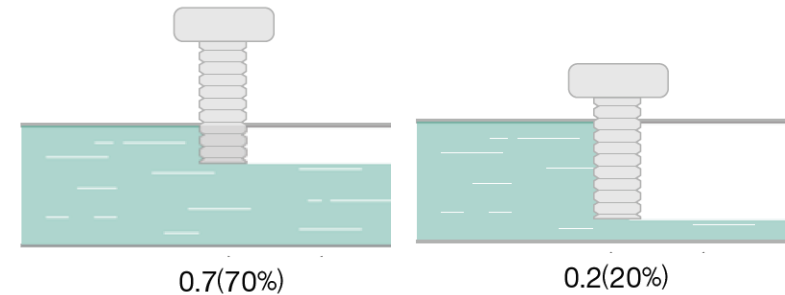
- c_t : the (memory) *cell*/state at time t .

Introduction to LSTM

“LSTM (Long Short-Term Memory) handles long-term dependencies in sequential data using the **memory cell** and **gates**.”



An LSTM unit (layer)

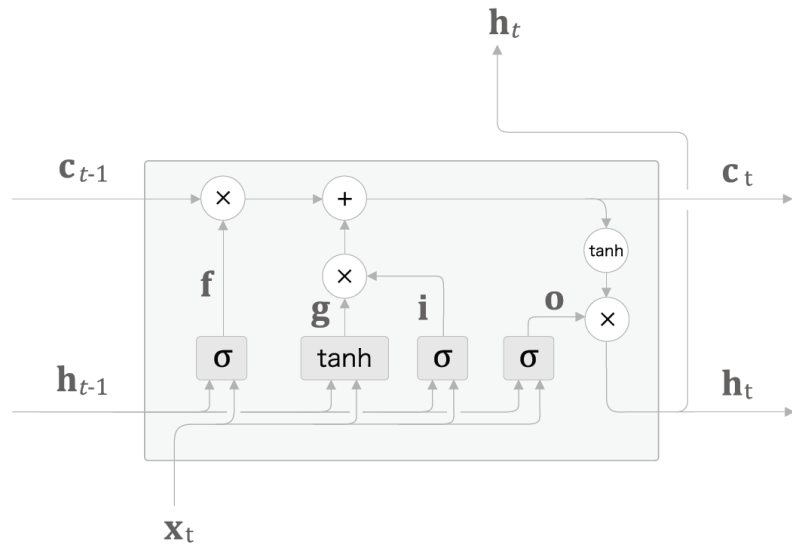


“A gate controls the amount of flow.”

Gate

LSTM at a Glance

“An LSTM unit consists of a memory cell, gates, and hidden state.”



States

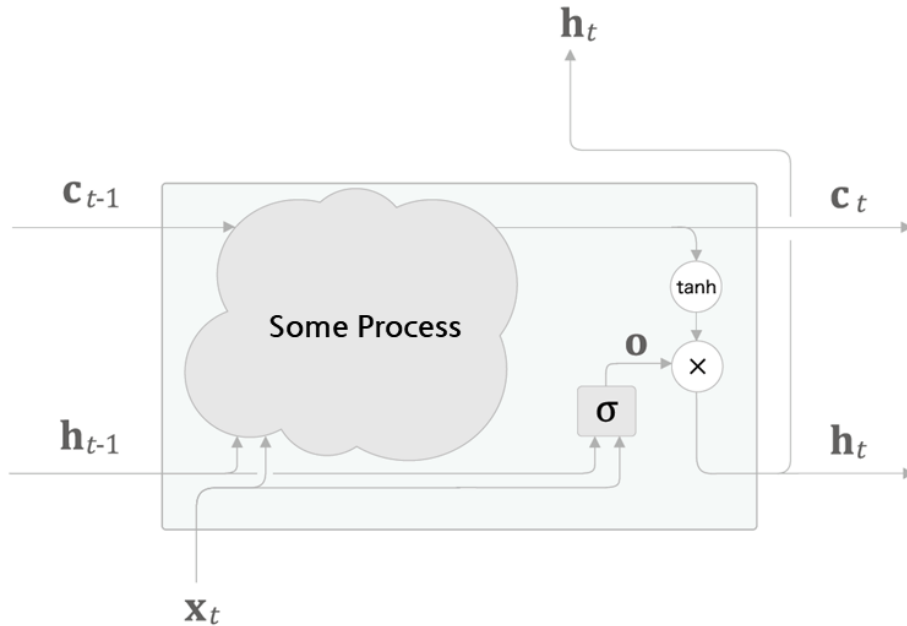
- (Memory) cell state (c): storing information over long periods.
- Hidden state (h): acts as short-term memory, similar to vanilla RNNs.

Gates

- Output gate (o): filters the updated cell state to produce the next hidden state.
- Forget gate (f): determines what information to discard from the previous state.
- Input gate (i): decides which new information (g) should be added to the cell state.

Output Gate

“An output gate filters the updated cell state to produce the next hidden state.”



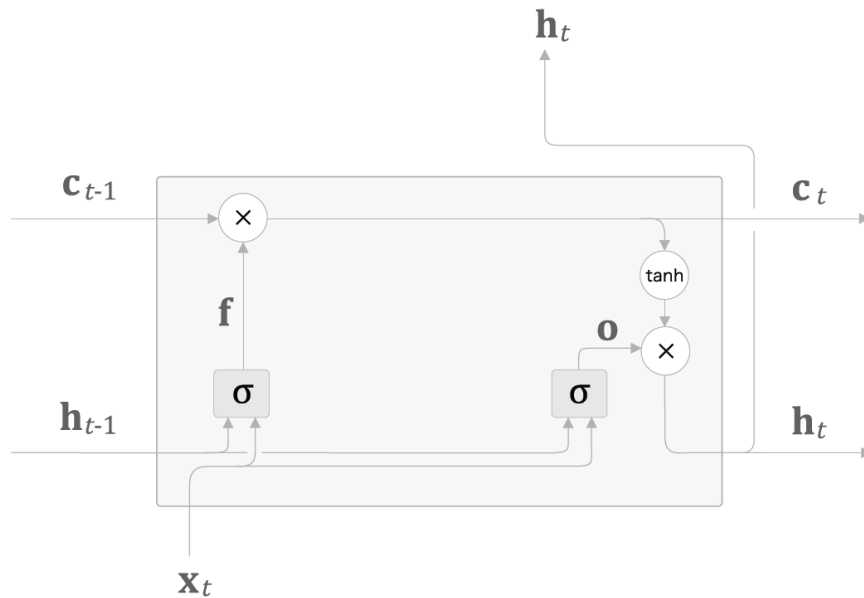
- $\mathbf{o} = \sigma(\mathbf{x}_t \mathbf{W}_x^{(o)} + \mathbf{h}_{t-1} \mathbf{W}_h^{(o)} + \mathbf{b}^{(o)})$

- $\mathbf{h}_t = \mathbf{o} \odot \tanh(\mathbf{c}_t)$

where \odot is elementwise product
(a.k.a., Hadamard product).

Forget Gate

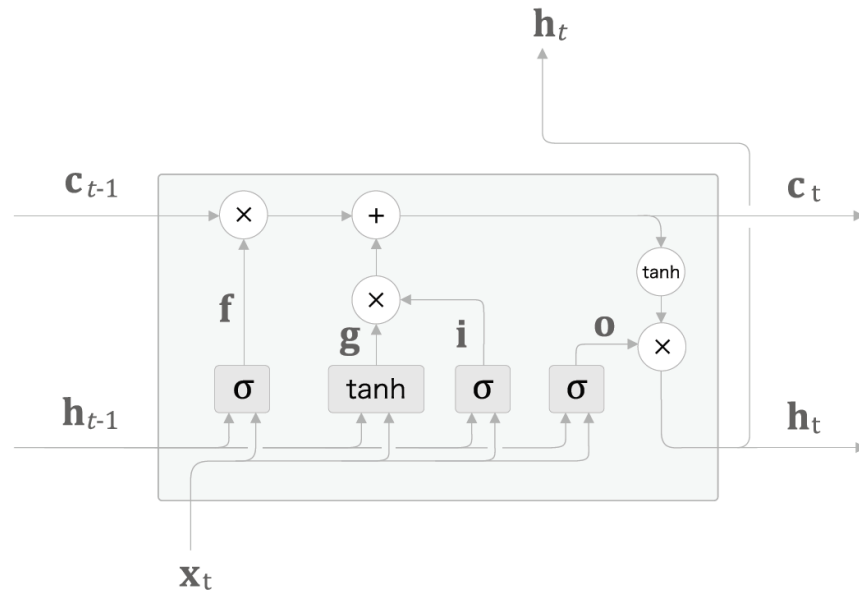
“A forget gate determines what information to discard from the previous state.”



- $f = \sigma(x_t W_x^{(f)} + h_{t-1} W_h^{(f)} + b^{(f)})$
- $c_t = f \odot c_{t-1}$

Input Gate

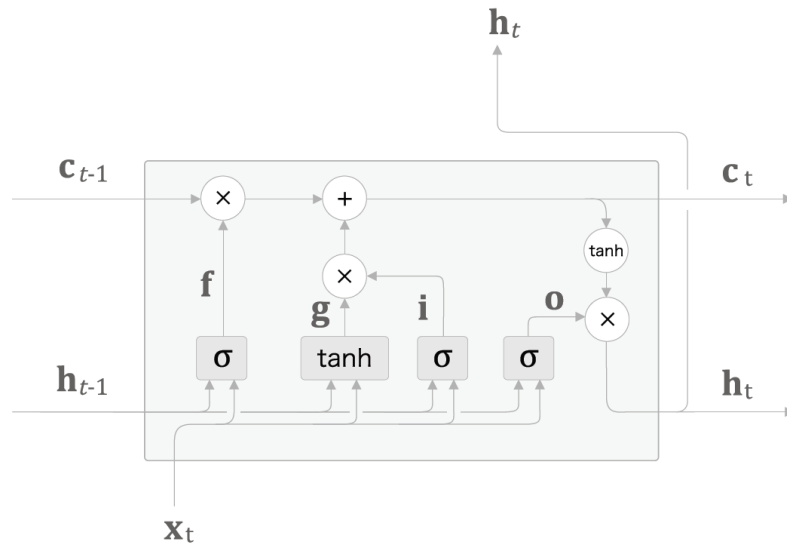
“An input gate decides which new information (g) should be added to the cell state.”



- $g = \tanh(x_t W_x^{(g)} + h_{t-1} W_h^{(g)} + b^{(g)})$
- $i = \sigma(x_t W_x^{(i)} + h_{t-1} W_h^{(i)} + b^{(i)})$
- $c_t = f \odot c_{t-1} + g \odot i$

Recap: LSTM at a Glance

“An LSTM unit consists of a memory cell, gates, and hidden state.”



- $f = \sigma(x_t W_x^{(f)} + h_{t-1} W_h^{(f)} + b^{(f)})$
- $g = \tanh(x_t W_x^{(g)} + h_{t-1} W_h^{(g)} + b^{(g)})$
- $i = \sigma(x_t W_x^{(i)} + h_{t-1} W_h^{(i)} + b^{(i)})$
- $o = \sigma(x_t W_x^{(o)} + h_{t-1} W_h^{(o)} + b^{(o)})$
- $c_t = f \odot c_{t-1} + g \odot i$
- $h_t = o \odot \tanh(c_t)$

tanh and Sigmoid in LSTM

▪ Why do we use ***tanh*** in LSTM?

- Bounded output range
 - *tanh* squishes values between -1 and 1, which helps prevent exploding gradients during training.
 - *tanh* provides the normalization effect during the recursive processing of past values.

- Zero-centered output

Unlike *sigmoid*, which outputs between 0 and 1, *tanh* can output both positive and negative values that allow the hidden state values to increase or decrease, providing more flexibility in modeling sequential data.

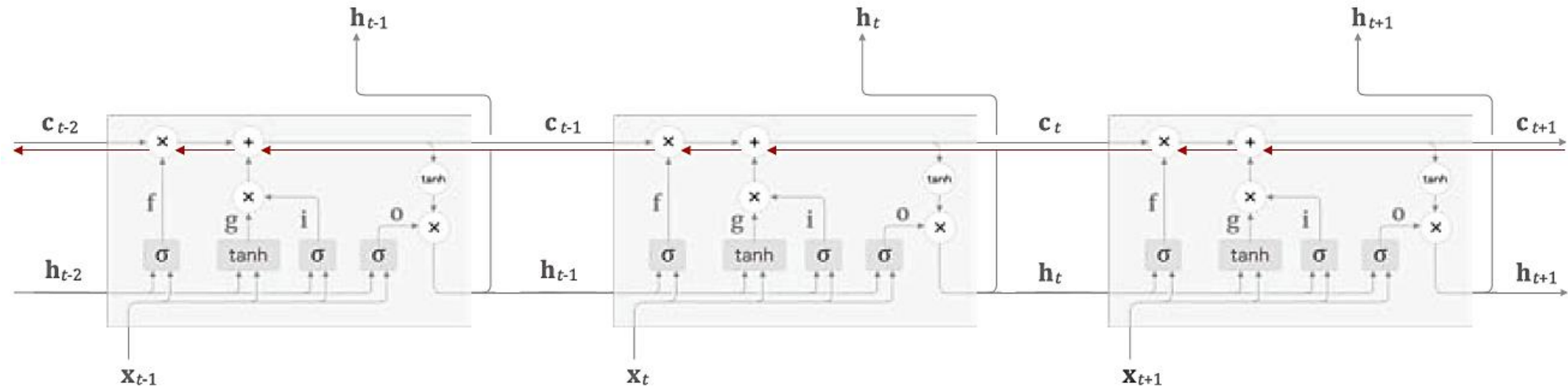
▪ Why do we use ***sigmoid*** in LSTM?

- Gating mechanism

Its output range for 0 to 1 makes it ideal for the gates, as it allows them to control the flow of information.

The Gradient Flow in LSTM

“LSTM addresses the gradient problems by the backpropagation through the memory cells.”



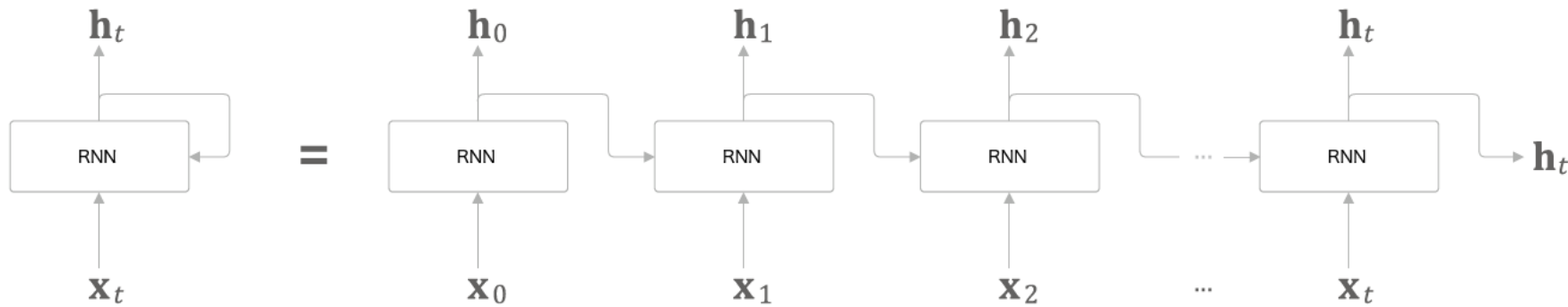
“Only + and \times operations on the path.”:

- + : no change in gradients.
- \times : independent elementwise products instead of recursive matrix multiplications.

Takeaways

Recurrent Neural Network (RNN)

“A recurrent neural network **recurrently** applies a sequence processor using *hidden states*.”

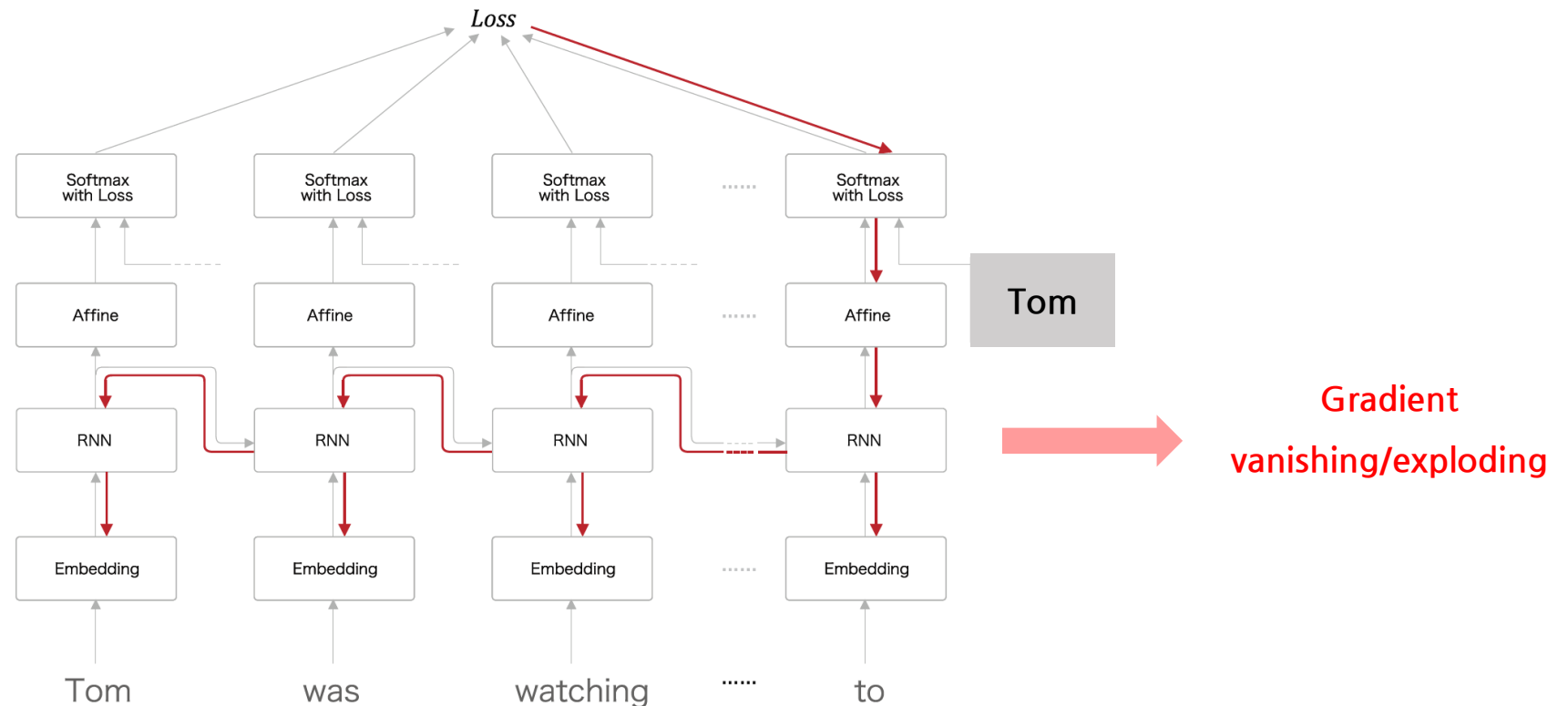


- $h_t = \tanh(h_{t-1}^T W_{hh} + x_t W_{xh} + b_h)$
- $y_t = \tanh(h_t W_{hy} + b_y)$

Limitation of Vanilla RNNs

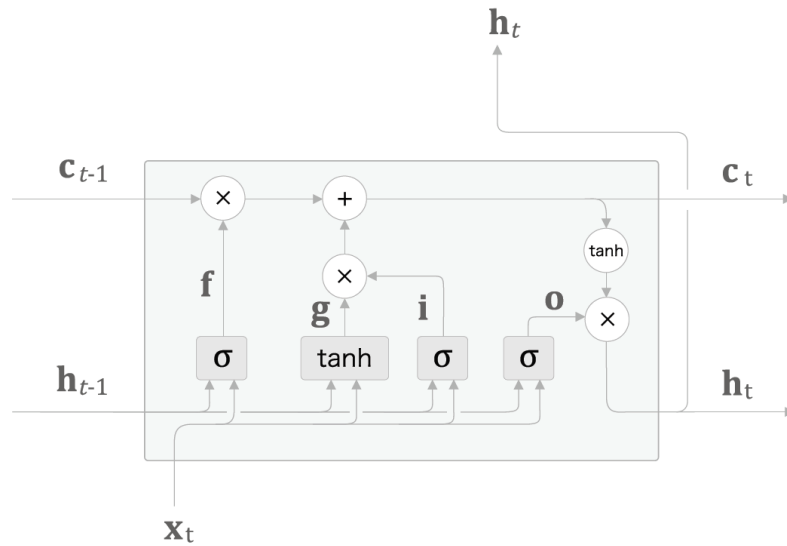
“Vanilla RNNs are limited in their ability to learn long-term memory due to **computational cost** and **gradient issues**.”

“Tom was watching TV in his room. Mary came into the room. Mary said hi to [?].”



Recap: LSTM at a Glance

“An LSTM unit consists of a memory cell, gates, and hidden state.”



- $f = \sigma(x_t W_x^{(f)} + h_{t-1} W_h^{(f)} + b^{(f)})$
- $g = \tanh(x_t W_x^{(g)} + h_{t-1} W_h^{(g)} + b^{(g)})$
- $i = \sigma(x_t W_x^{(i)} + h_{t-1} W_h^{(i)} + b^{(i)})$
- $o = \sigma(x_t W_x^{(o)} + h_{t-1} W_h^{(o)} + b^{(o)})$
- $c_t = f \odot c_{t-1} + g \odot i$
- $h_t = o \odot \tanh(c_t)$

Thank you! 😊