

## Lec 10: Introduction to Long-Short Term Memory



hsyi@kisti.re.kr

Hongsuk Yi (이홍석)

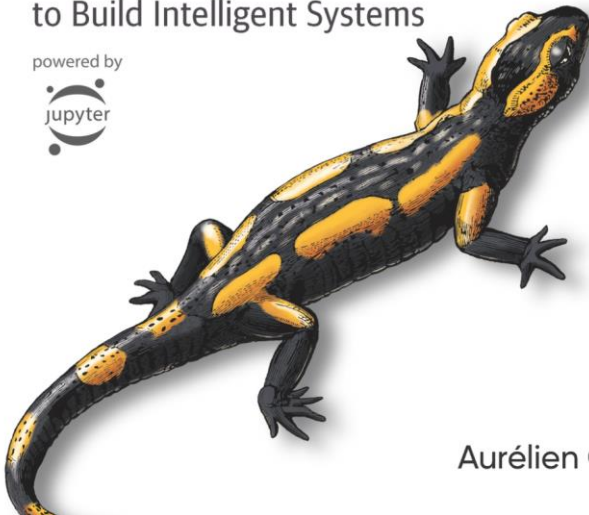


- ❖ Introduction to Recurrent Neural Network
  - ✓ Simple RNN, BPTT, Memory Cell
  - ✓ Code: Implementing an RNN with Keras
- ❖ **Introduction to Long-Short Term Memory**
  - ✓ Cell state, LSTM, and GRU, and Applications
  - ✓ A Visual Guide to Recurrent Layers in Keras
  - ✓ Code: A simple LSTM layers
- ❖ Text generation with RNN
  - ✓ Tokenizer, Character-Level Language model
  - ✓ Code: Alice's Adventures in Wonderland
- ❖ Sequence to Sequence Learning model with RNN
  - ✓ Introduction to Seq2Seq and Attention model
  - ✓ Code: Character-Level Neural Machine Translation

O'REILLY®

## Hands-On Machine Learning with Scikit-Learn, Keras & TensorFlow

Concepts, Tools, and Techniques  
to Build Intelligent Systems



Aurélien Géron

2nd Edition  
Updated for  
TensorFlow 2

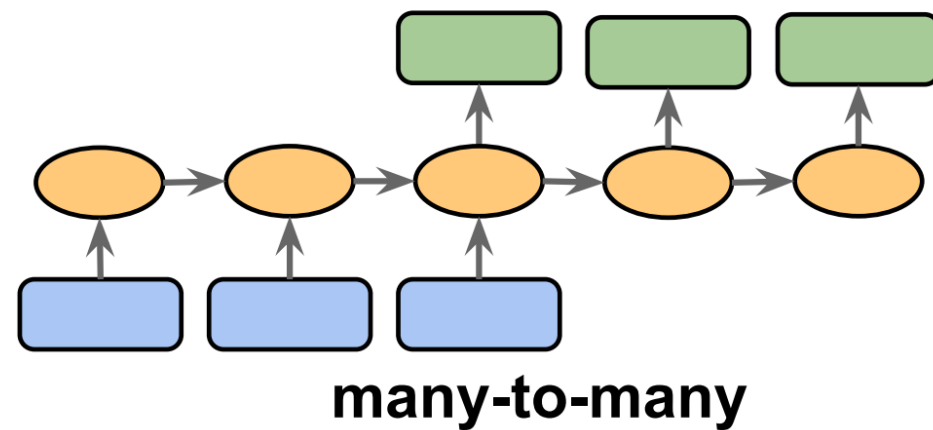
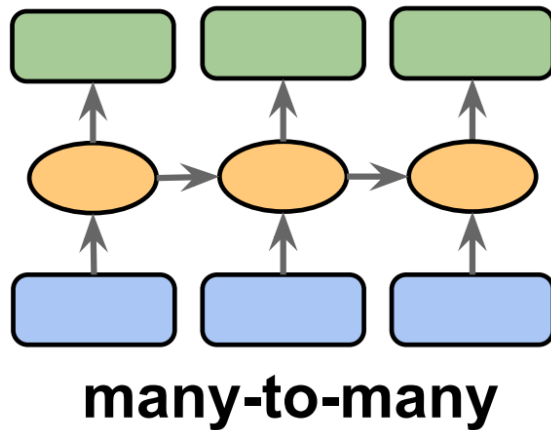
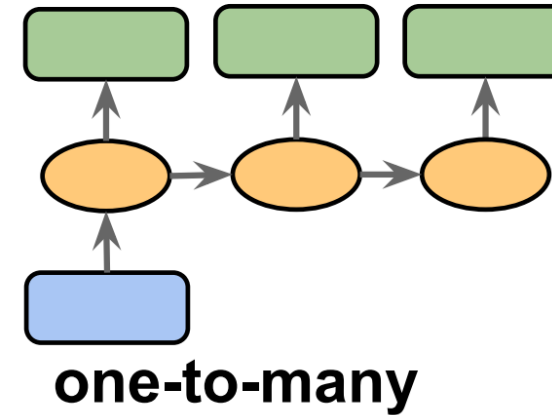
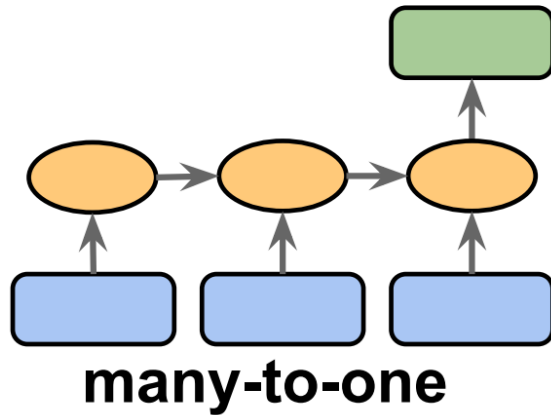
딥 러닝을 이용한 자연어 처리 입문

<https://wikidocs.net/book/2155>



# Reviewing the last class: RNN

- ❖ RNN can handle interactions more flexibly
- ❖ they are applied in a step-by-step fashion to a sequential input
  - ✓ a sequence of words, characters, or something else
- ❖ RNNs use a **state** representing what has happened previously
  - ✓ after each step, a new state is computed



<https://chalmers.instructure.com/courses/16100>

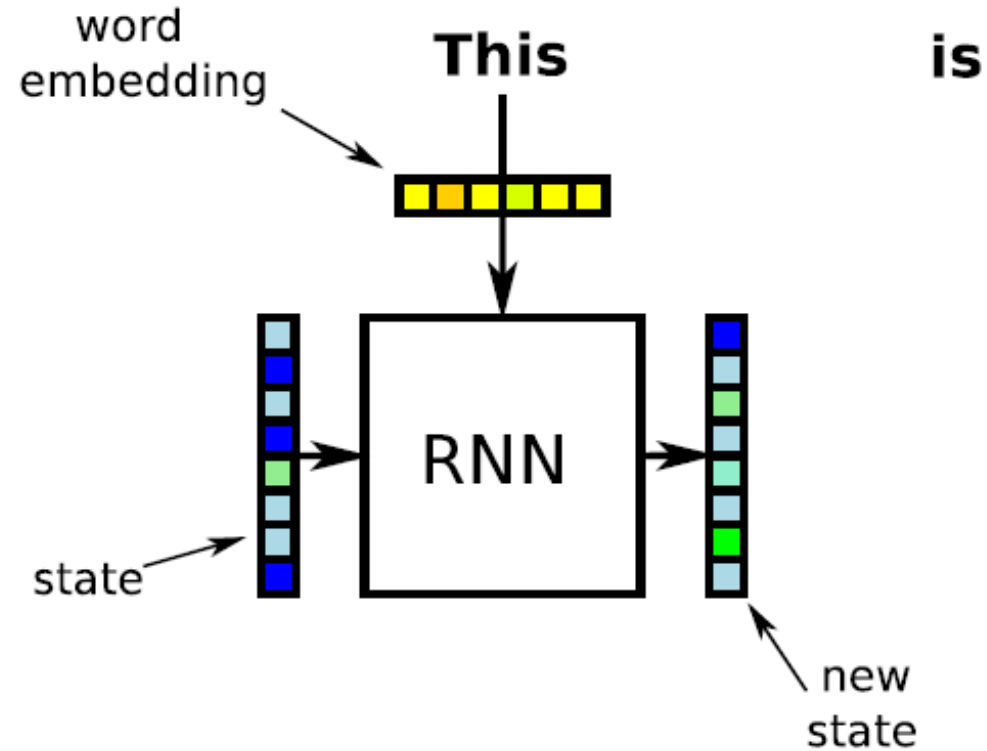


image borrowed from [Richard Johansson](#) (*Chalmers Technical University and University of Gothenburg*)



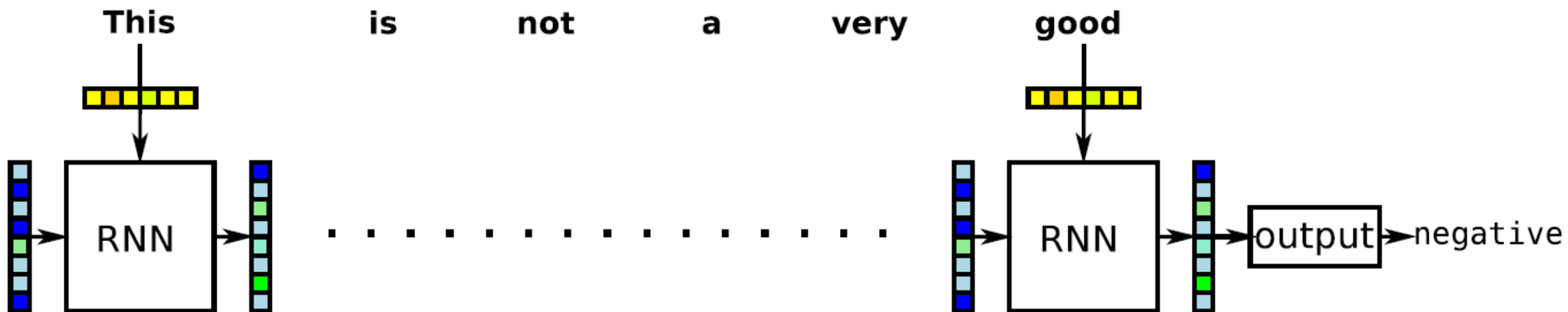


image borrowed from [Richard Johansson](#) (*Chalmers Technical University and University of Gothenburg*)



## ❖ the simple RNN looks similar to a feedforward NN

- ✓ the next state is computed like a hidden layer in a feedforward NN
- ✓ the output is identical to the state representation:

$$y_t = s_t$$
$$s_t = g(\mathbf{W} \cdot (s_{t-1} \oplus x_t) + \mathbf{b})$$

activation is typically tanh

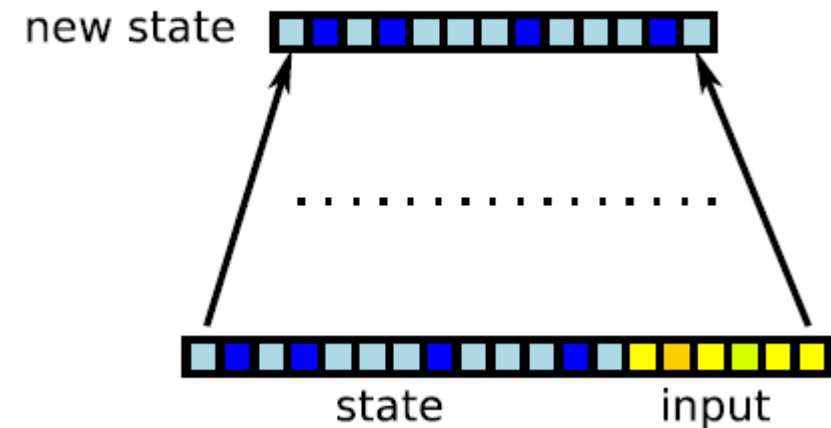
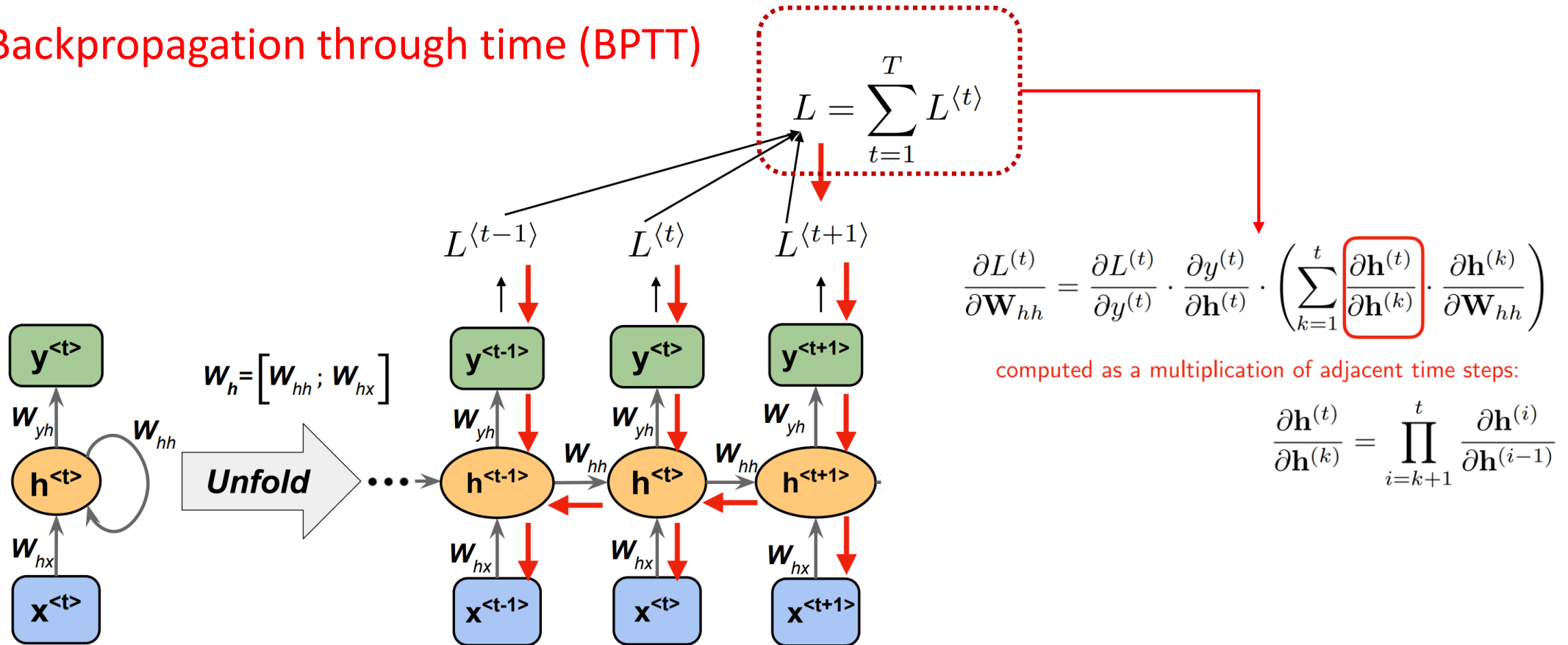


image borrowed from [Richard Johansson](#) (Chalmers Technical University and University of Gothenburg)

- ✓ To train an RNN the trick is to unroll it through time and then simply use regular backpropagation.

## Backpropagation through time (BPTT)



- ❖ simple RNNs suffer from the problem of **vanishing gradients** (Hochreiter, 1998)

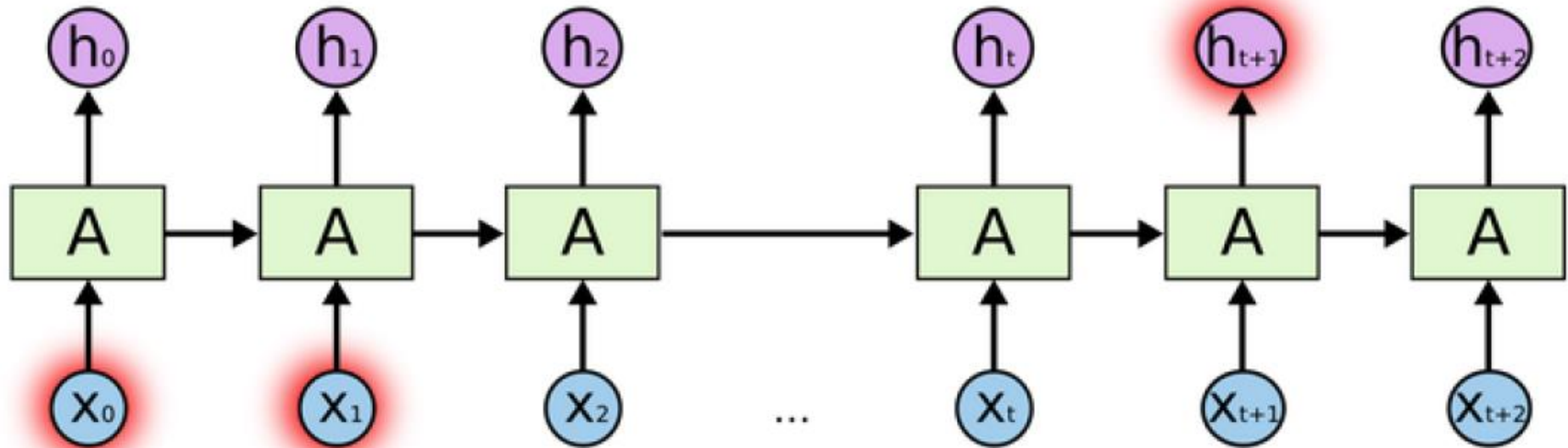


image borrowed from [Richard Johansson](#) (*Chalmers Technical University and University of Gothenburg*)

## ❖ Truncated backpropagation through time (TBPTT)

- ✓ simply limits the number of time steps the signal can backpropagate after each forward pass.
  - E.g., even if the sequence has 100 elements/steps, we may only backpropagate through 20 or so

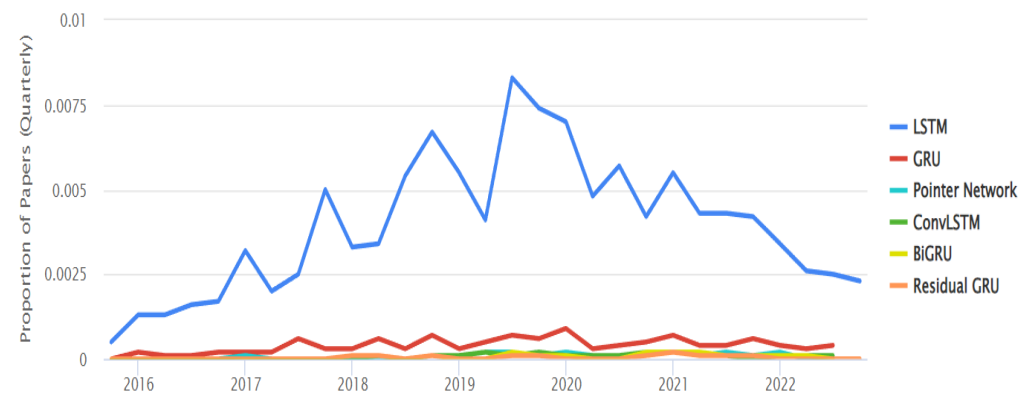
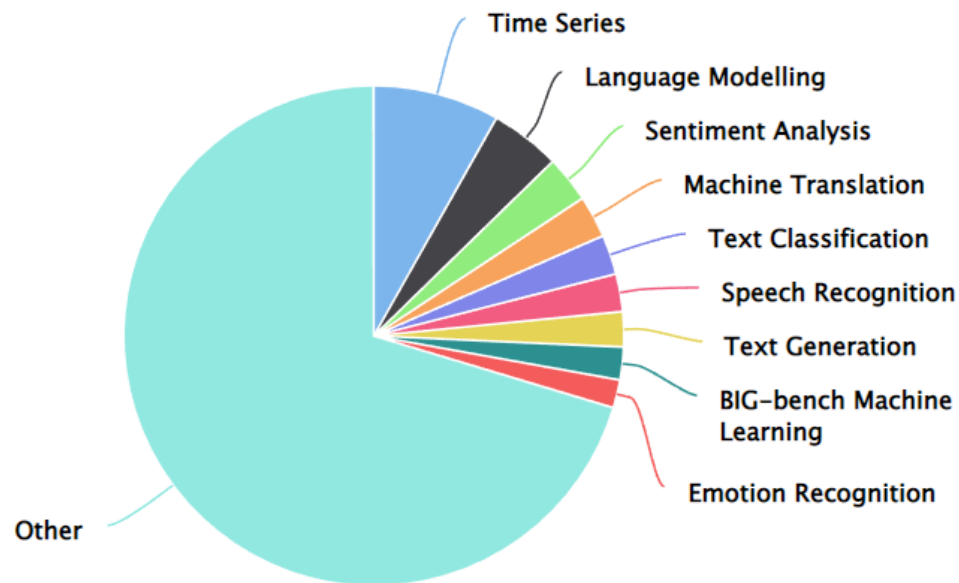
## ❖ Long short-term memory (LSTM)

- ✓ uses a memory cell for modeling long-range dependencies and avoid vanishing gradient problems

## ❖ Gradient Clipping

- ✓ set a max value for gradients if they grow to large
  - solves only exploding gradient problem)

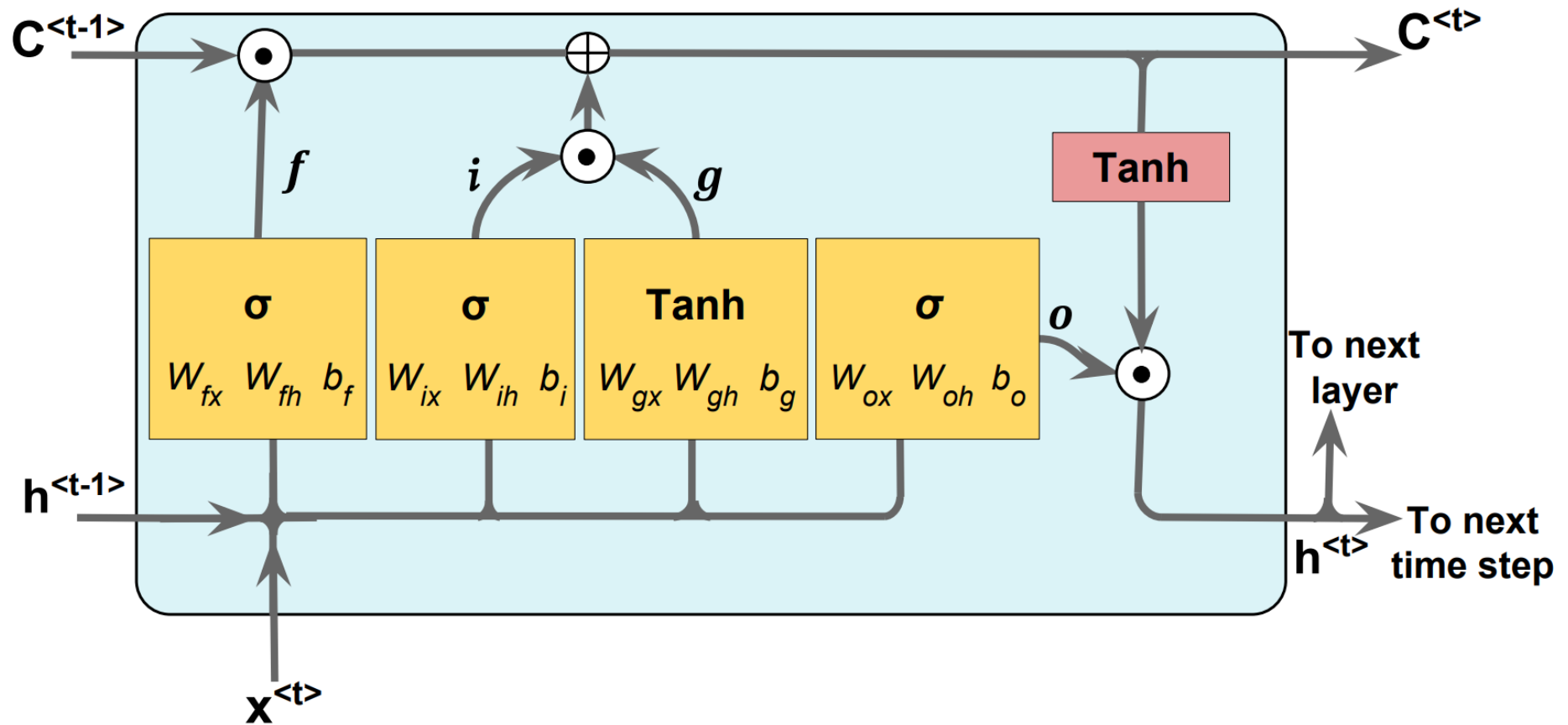
# Long-Short Term Memory (장단기 메모리)



⚠ This feature is experimental; we are continuously improving our matching algorithm.

<source> <https://paperswithcode.com/method/lstm>

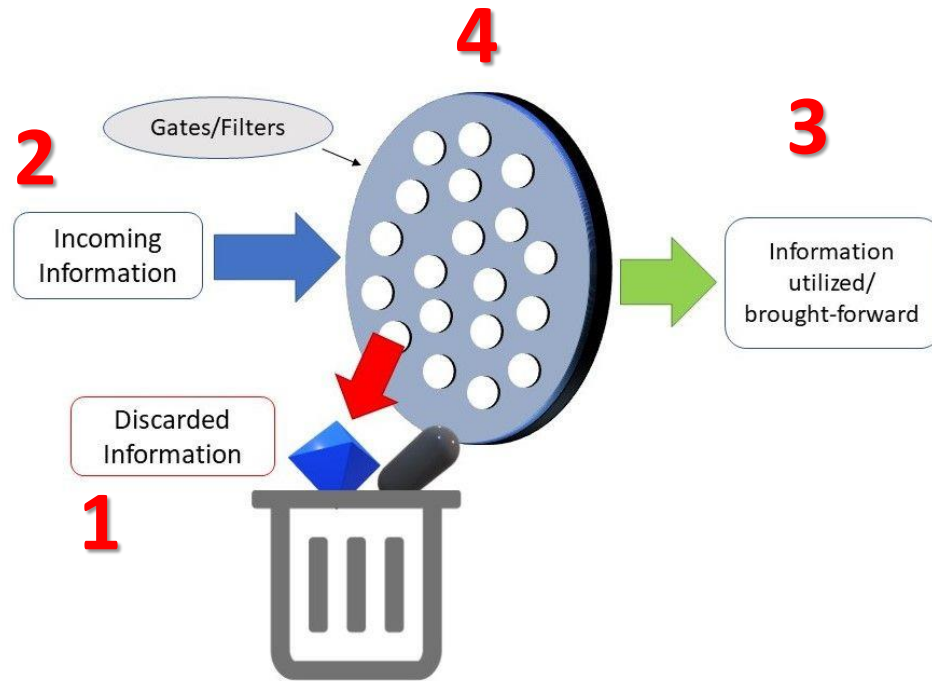
# Long-short term memory (LSTM)





## ❖ Cell state

- ✓ Gates control the flow of information to/from the memory



## ❖ Forget Gate

- ✓ whether we should keep the information from the previous timestamp or forget it

## ❖ Input Gate

- ✓ Decide how much this unit adds to the current state

## ❖ New information: Memory Upgate

- ✓ The cell state vector aggregates the two components (old memory via the forget gate and new memory via the input gate)

## ❖ Output Gate

- ✓ Decide what part of the current cell state makes it to the output

❖ **gating architectures** allow information flow to be controlled more carefully

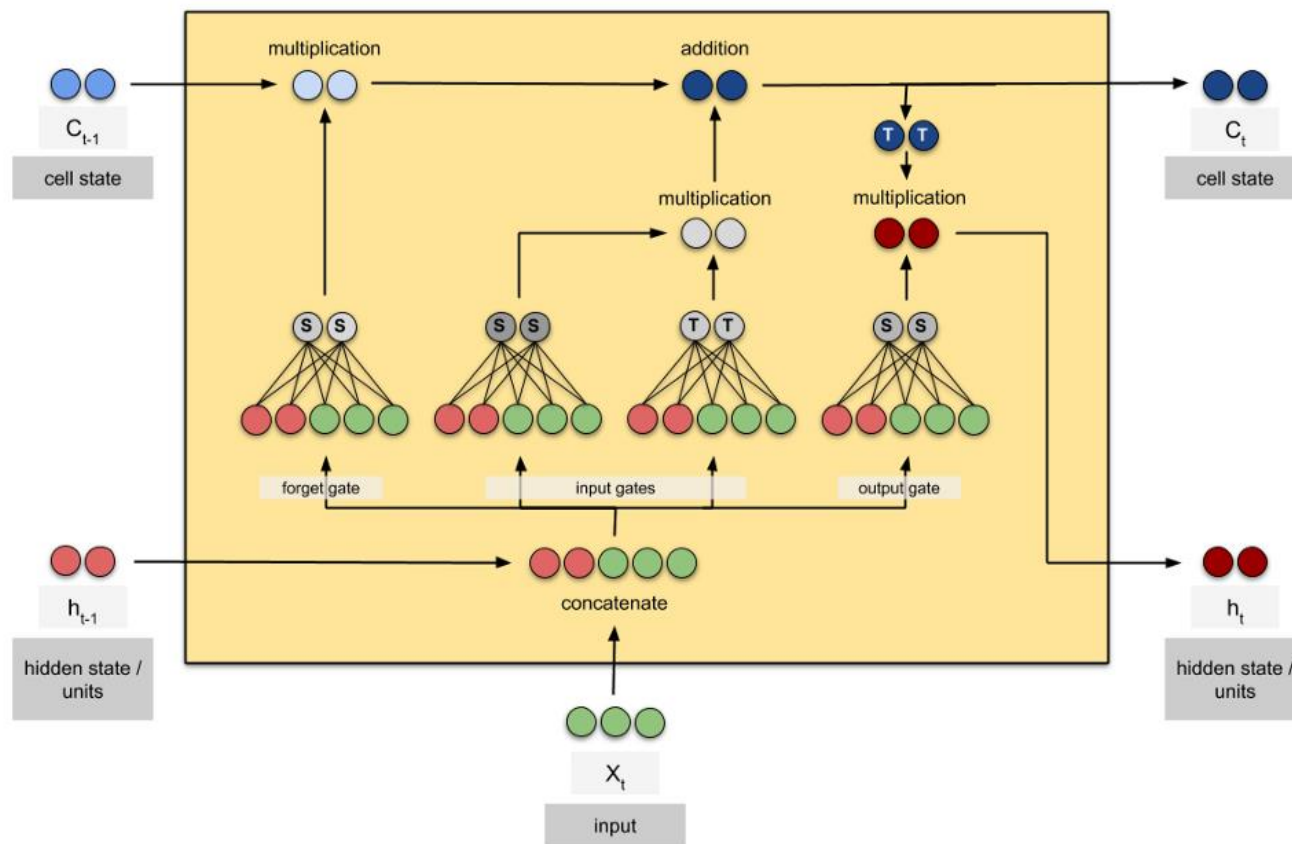
✓ should we copy the previous state, or replace it?

$$y_t = s_t \quad s_t = g(\mathbf{W} \cdot (s_{t-1} \oplus x_t) + \mathbf{b})$$

$$\begin{array}{c}
 \begin{bmatrix} 8 \\ 11 \\ 3 \\ 7 \\ 5 \\ 15 \end{bmatrix} \\
 s'
 \end{array}
 \leftarrow
 \begin{array}{cc}
 \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ 0 \\ 1 \end{bmatrix} & \odot & \begin{bmatrix} 10 \\ 11 \\ 12 \\ 13 \\ 14 \\ 15 \end{bmatrix} \\
 \mathbf{g} & & \mathbf{x}
 \end{array}
 +
 \begin{array}{cc}
 \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} & \odot & \begin{bmatrix} 8 \\ 9 \\ 3 \\ 7 \\ 5 \\ 8 \end{bmatrix} \\
 (1 - \mathbf{g}) & & \mathbf{s}
 \end{array}$$

image borrowed from [Richard Johansson](#) (Chalmers Technical University and University of Gothenburg)

- ❖ Gates are controlled by a concatenation of the output from the previous time step and the current input and **optionally the cell state vector**.

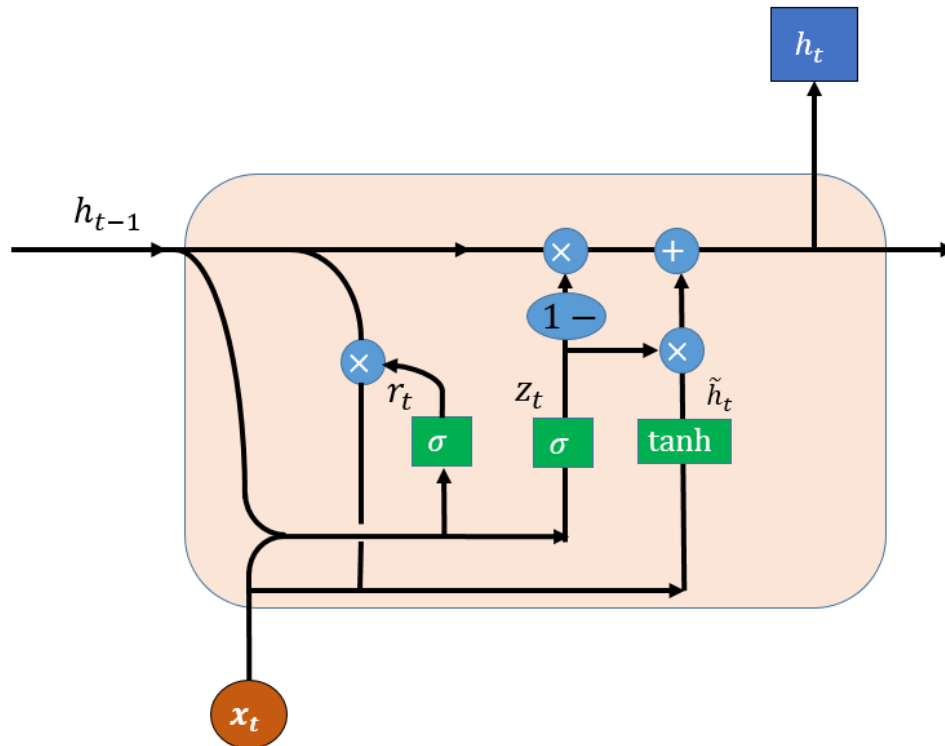


S: Sigmoid  
T: tanh

# Gated Recurrent Units (GRU)

❖ Just like LSTM, GRU uses gates to control the flow of information.

✓ GRU simplifies the architecture of LSTM, which was similar in performance to LSTM and was complex



$$r_t = \sigma(W_{xr}x_t + W_{hr}h_{t-1} + b_r)$$

$$z_t = \sigma(W_{xz}x_t + W_{hz}h_{t-1} + b_z)$$

$$g_t = \tanh(W_{hg}(r_t \circ h_{t-1}) + W_{xg}x_t + b_g)$$

$$h_t = (1 - z_t) \circ g_t + z_t \circ h_{t-1}$$

source: [https://www.researchgate.net/figure/Structure-of-a-GRU-cell\\_fig1\\_334385520](https://www.researchgate.net/figure/Structure-of-a-GRU-cell_fig1_334385520)

- ❖ Even with gated RNNs, it can be hard to cram the useful information into **the last state**

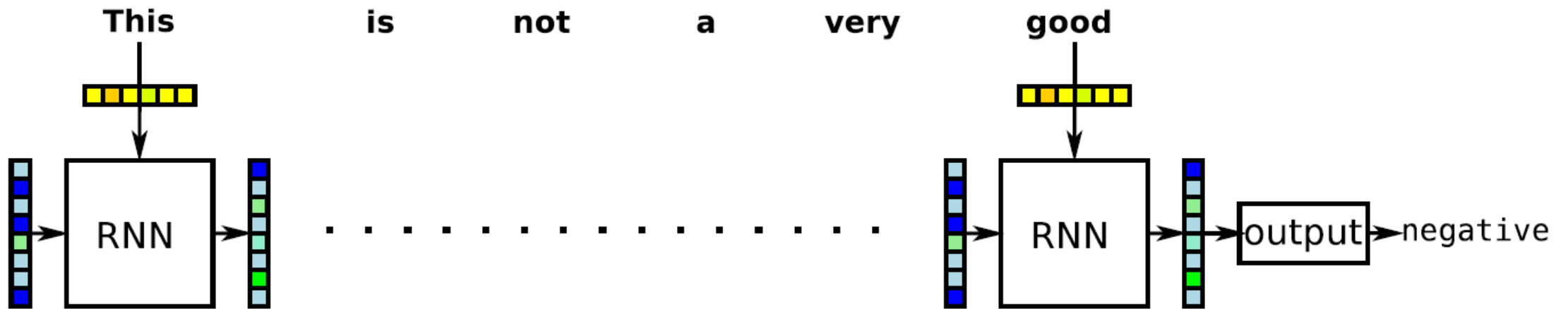
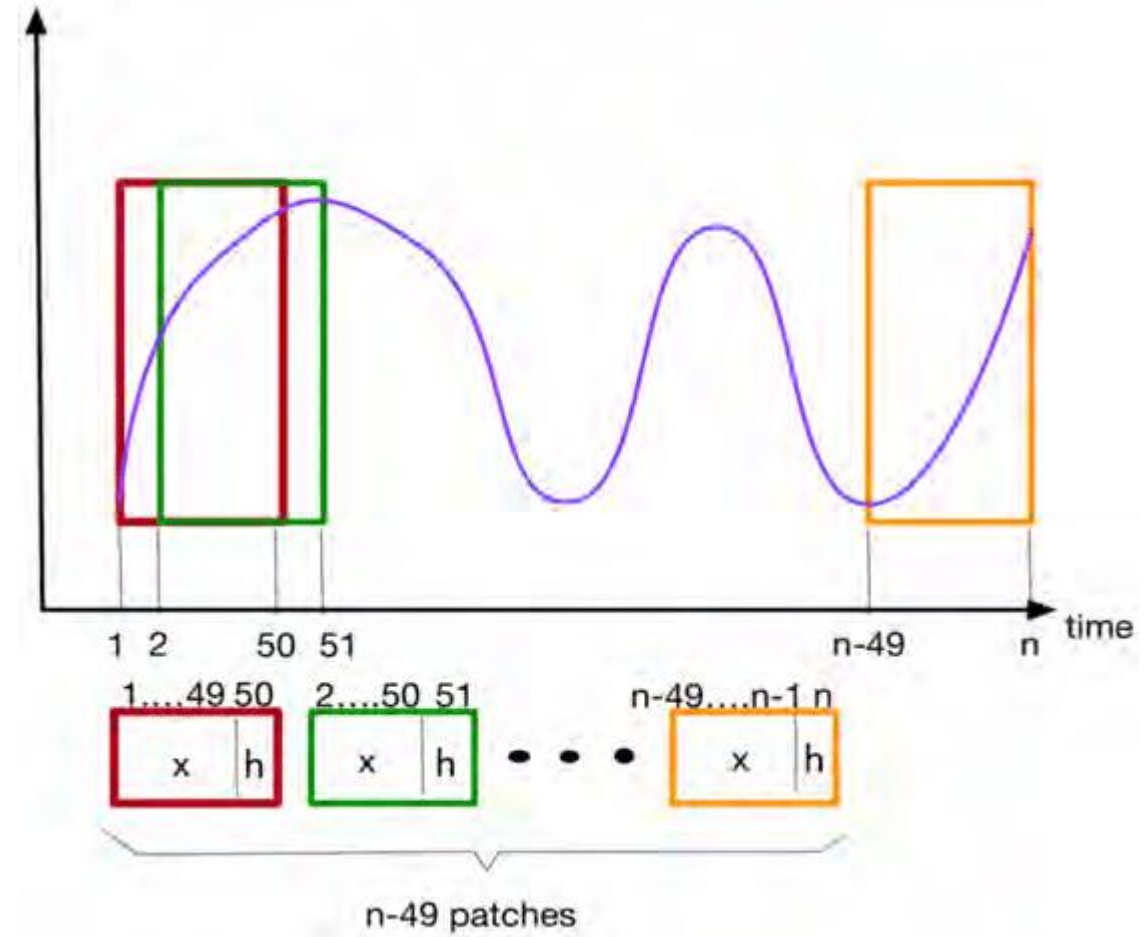
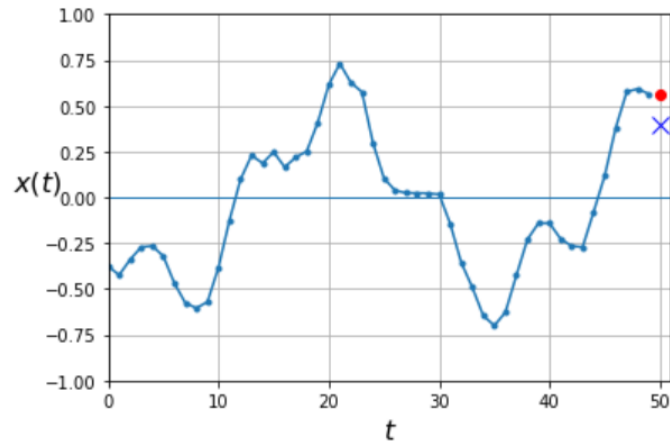


image borrowed from [Richard Johansson](#) (*Chalmers Technical University and University of Gothenburg*)

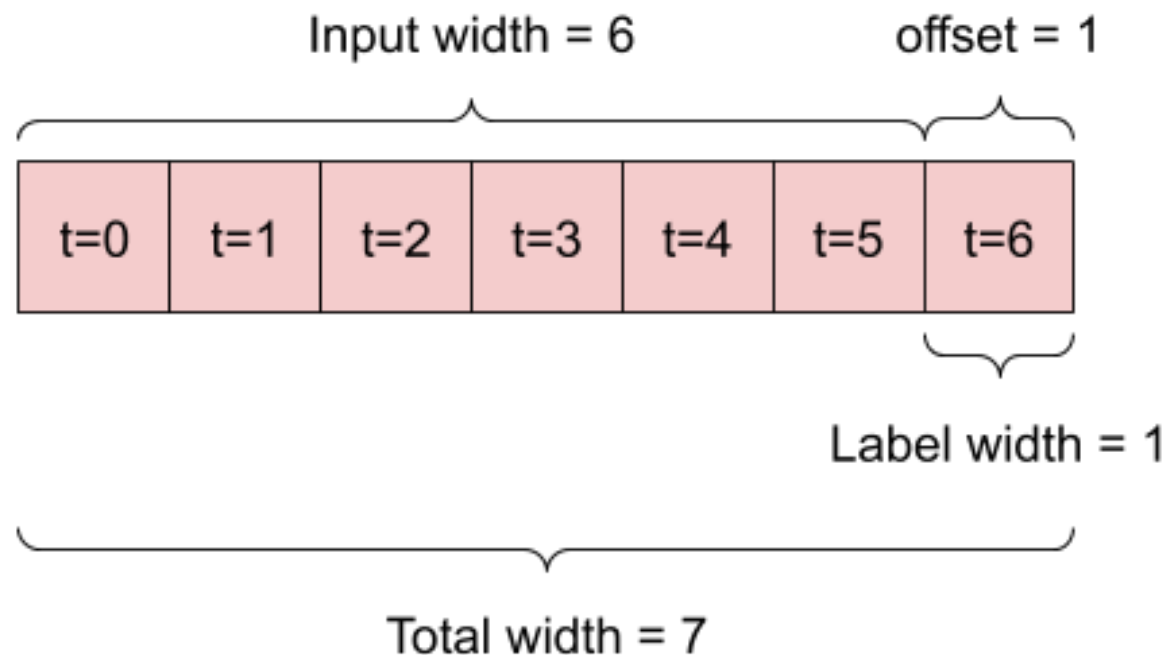


# Simple LSTM layer

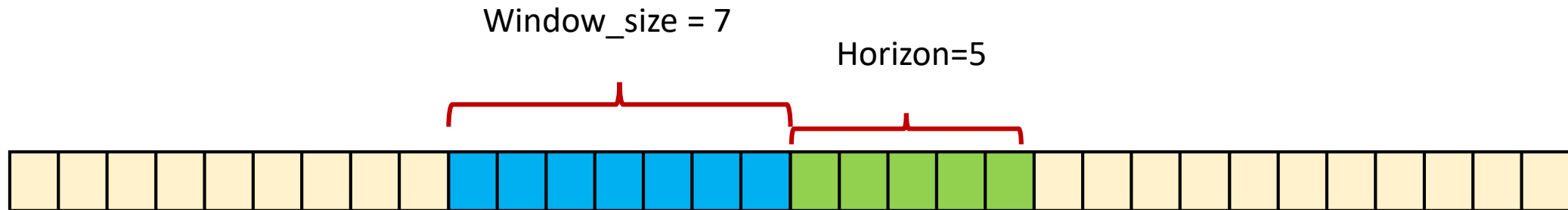
# RNN Sliding window size and horizon



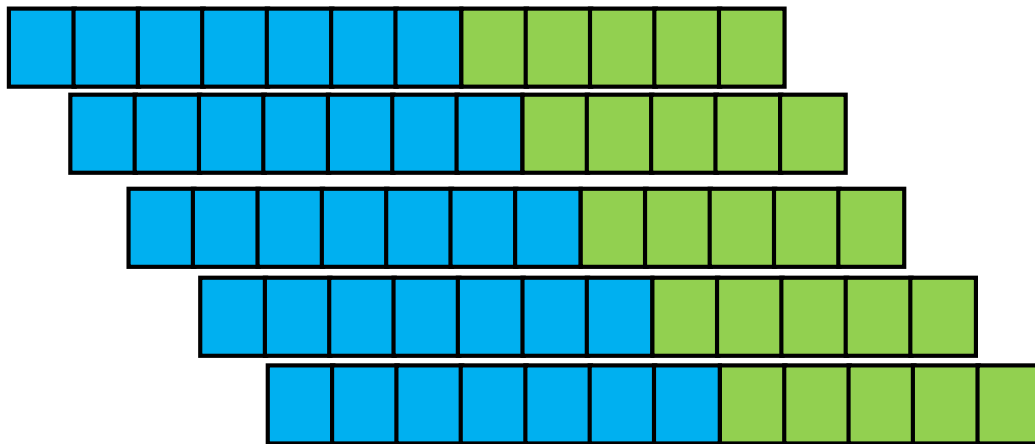
- ❖ A model : prediction one hour into the future, given six hours of history



Step1: set the number of window\_size, horizon



Step2: set the number of window\_size, horizon

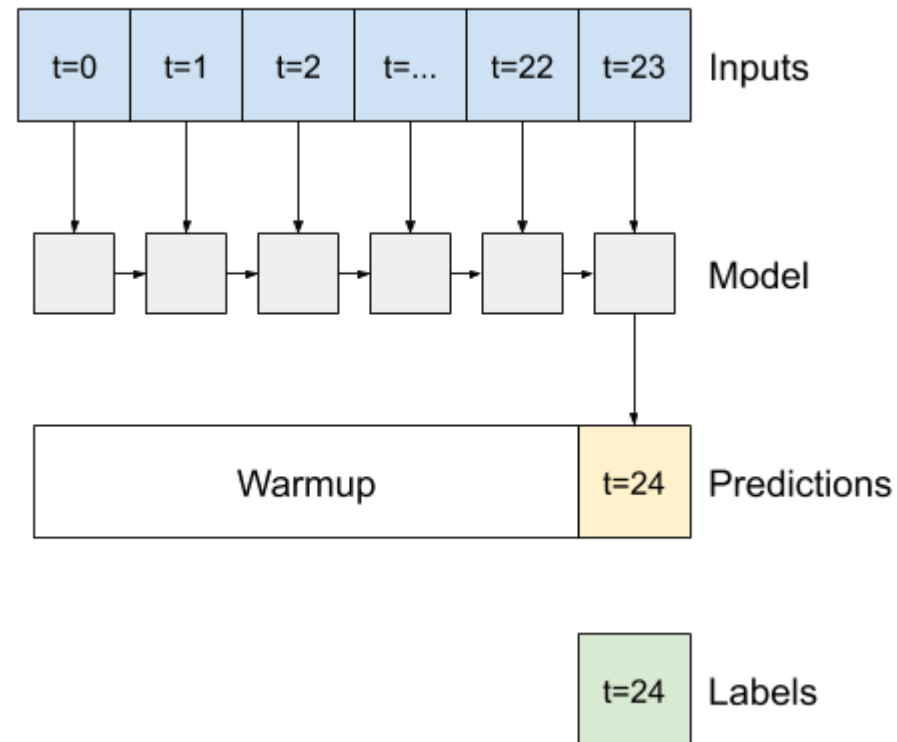


## ❖ LSTM has 3 important parameters

- ✓ **neurons** : dimensionality of the output space
- ✓ **return\_sequences**: whether to return the last output. (hidden state, memory cell, h)
  - in the output sequence, or the full sequence.
  - Default: False.
- ✓ **return\_state**: whether to return the last state in addition to the output. (cell state, c)
  - Default: False.

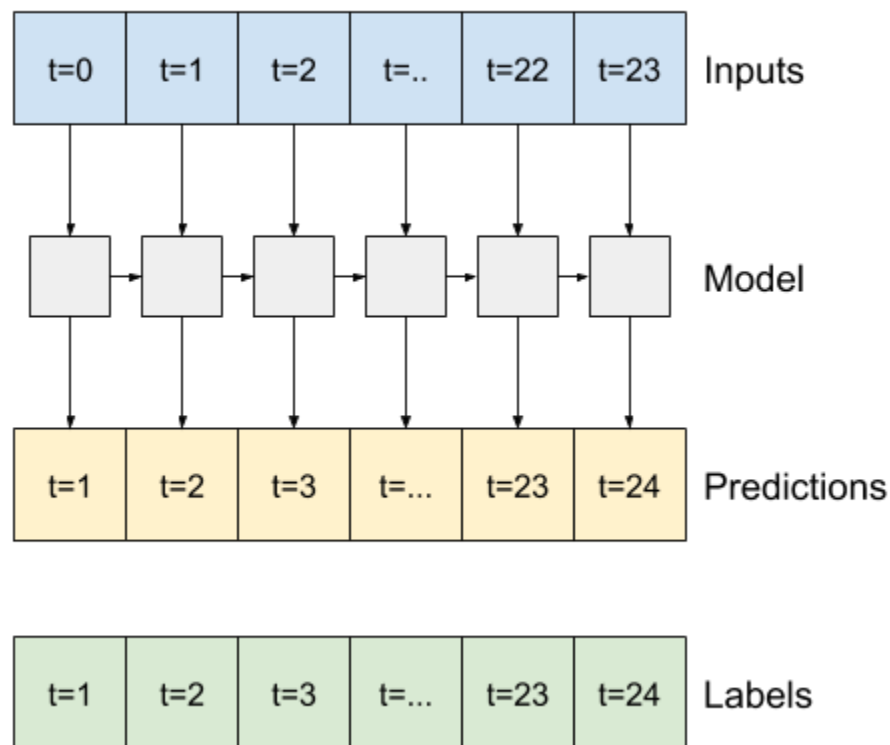
## ❖ Return\_sequences=False

- ✓ the layer only returns the output of the final time step
- ✓ giving the model time to warm up its internal state before making a single prediction:



## ❖ Return\_sequences=True

- ✓ the layer returns an output for each input. This is useful for Stacking RNN layers.
- ✓ Training a model on multiple time steps simultaneously





**a sequence modeling  
problem :**  
predict the next word

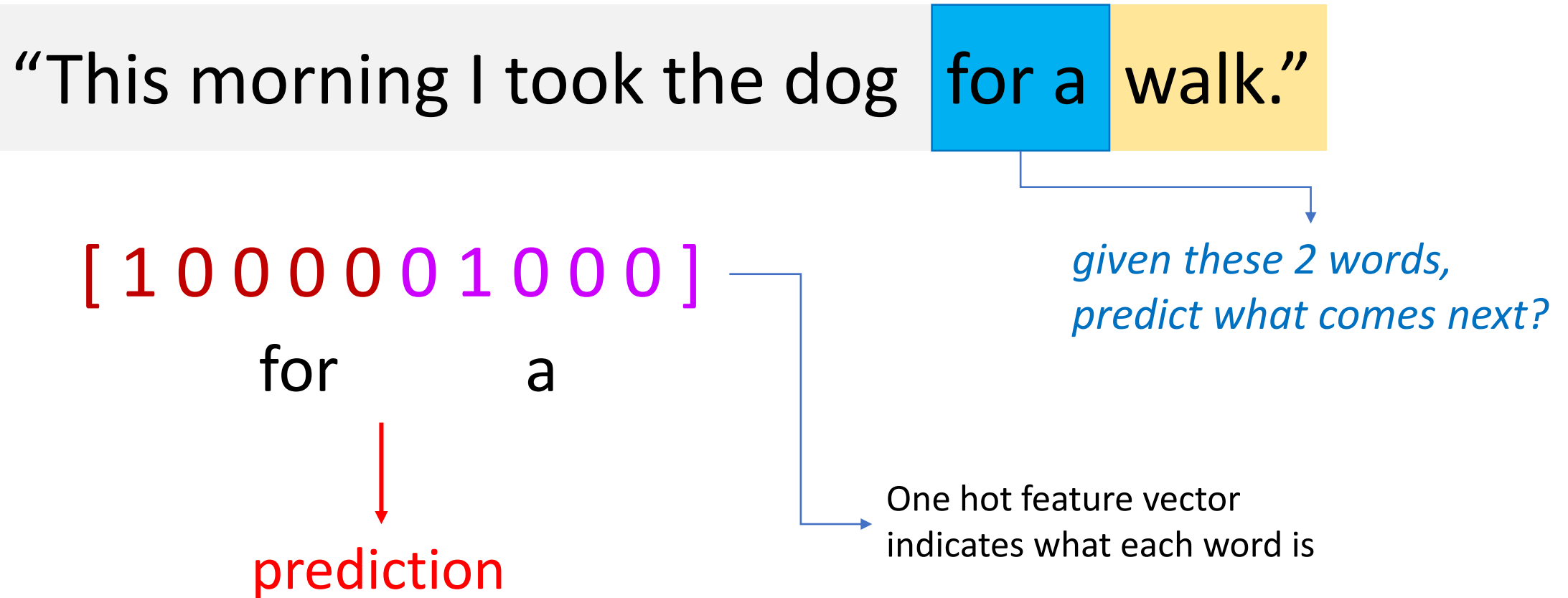
## a sequence modeling problem

“This morning I took the dog for a walk.”

*given these words*

*predict what comes next?*

## ❖ idea: use a fixed window



**problem:** we can't model long-term dependencies

"In **France**, I had a great time and I learnt some of the \_\_\_\_\_  
**language.**"

We need information from the far past and future to accurately guess the correct word.

## idea: use entire sequence, as a set of counts

Bag-of-words model:

In this model, a text (such as a sentence or a document) is represented as the bag (multiset) of its words, disregarding grammar and even word order but keeping multiplicity.

This morning I took the dog for a

[0 1 0 0 1 0 0 ... 0 0 1 1 0 0 0 1]

"bag of words"

prediction



**problem:** counts don't preserve order

“The food was good, not bad at all.”

VS

“The food was bad, not good at all.”

**idea:** use a really big fixed window

"This morning I took the dog for a walk."

*given these 7  
words, predict  
the next word*

[ 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 ... ]

morning I took the dog ...

↓  
prediction



## problem: no parameter sharing

this morning  
[ 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 ... ]

each of these inputs has a *separate parameter*

this morning  
[ 0 0 0 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 1 0 0 ... ]

things we learn about the sequence **won't transfer** if they appear **at different points** in the sequence.

to model sequences, we need:

1. to deal with **variable-length** sequences
2. to maintain **sequence order**
3. to keep track of **long-term dependencies**
4. to **share parameters** across the sequence

# A Visual Guide to Recurrent Layers in Keras

source: <https://amitness.com/2020/04/recurrent-layers-keras/>

## ❖ Let's take a simple example of encoding

For simplicity, let's assume we used some word embedding to convert each word into 2 numbers.



I am groot

Credits: Marvel Studios

Word	E1	E2
I	0.5	0.4
am	0.3	0.1
groot	0.7	0.5

We could either use one-hot encoding, pretrained word vectors, or learn word embeddings from scratch

<https://amitnness.com/2020/04/recurrent-layers-keras/>

## ❖ SimpleRNN with a Dense layer

- ✓ to build an architecture for something like sentiment analysis or text classification.

```
import tensorflow as tf
from tensorflow.keras.layers import SimpleRNN  #Dense, LSTM
# from tensorflow.keras.models import Sequential

x = tf.random.normal((1, 3, 2))

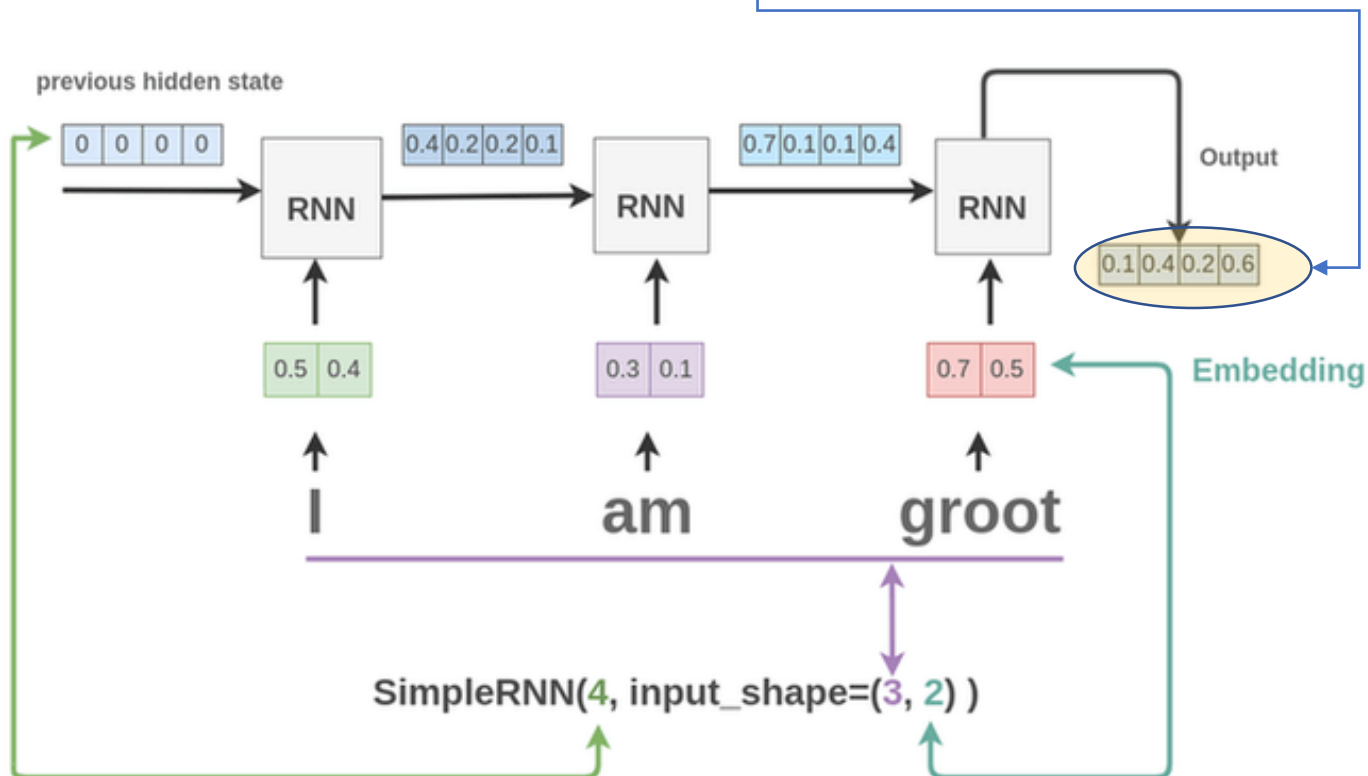
layer = SimpleRNN(4, input_shape=(3, 2))
output = layer(x)

print(output.shape)
print(x)

(1, 4)
tf.Tensor(
[[[ 0.6887584  1.3883604 ]
 [ 0.01564607 -1.4314882 ]
 [-0.05214449 -0.65099174]]], shape=(1, 3, 2), dtype=float32)
```

- ❖ we treat each word as a time-step and the embedding as features.

```
model.add(layers.SimpleRNN(4, input_shape=(3, 2)))
```



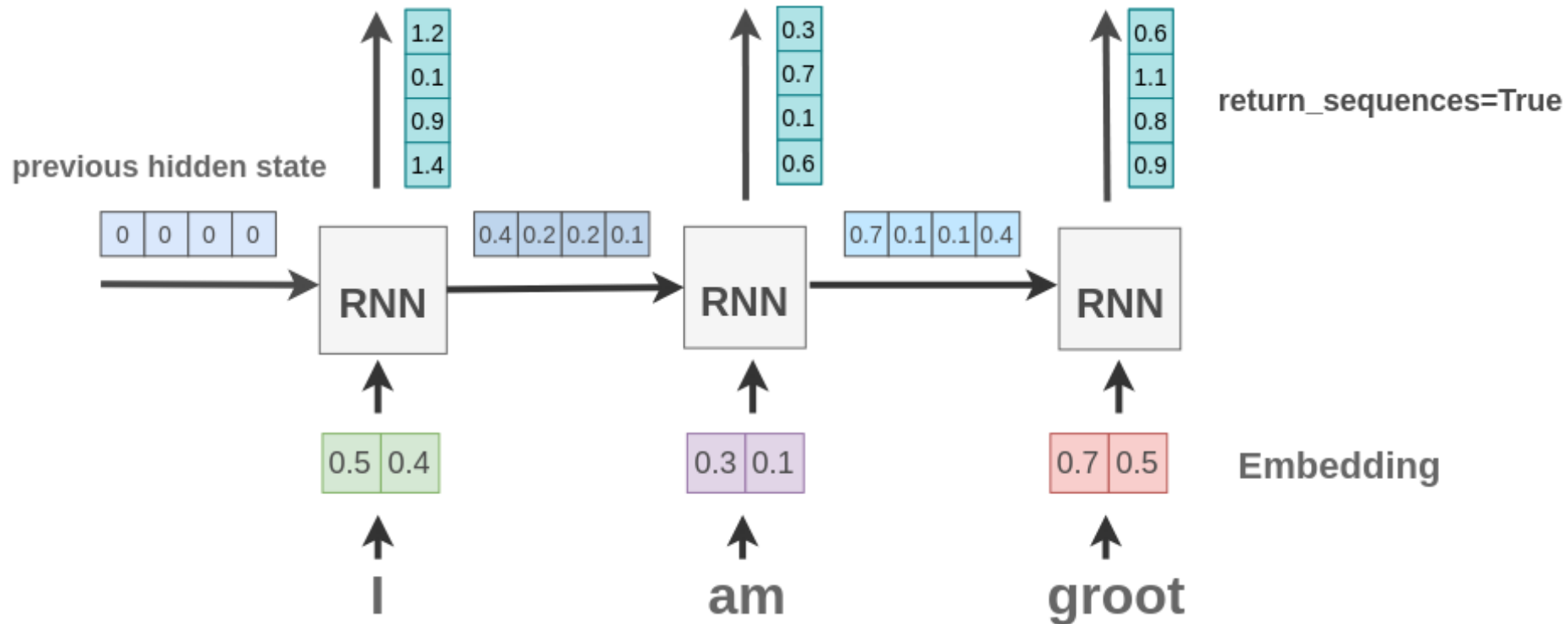
```
# multiple output
layer = SimpleRNN(4, input_shape=(3, 2), return_sequences=True )
output = layer(x)
print(output.shape)
print(output)
```

```
(1, 3, 4)
tf.Tensor(
[[[-0.6854385  0.08265962  0.30888444 -0.30752325]
 [ 0.4584542 -0.1935767 -0.91095936 -0.2416075 ]
 [ 0.7241105 -0.49960855 -0.5059616  0.7261468 ]]], shape=(1, 3, 4), dtype=float32)
```

## ❖ return\_sequences = True

✓ True : the output from each unfolded RNN cell is returned instead of only the last cell.

```
model.add(SimpleRNN(4, input_shape=(3, 2), return_sequences=True))
```





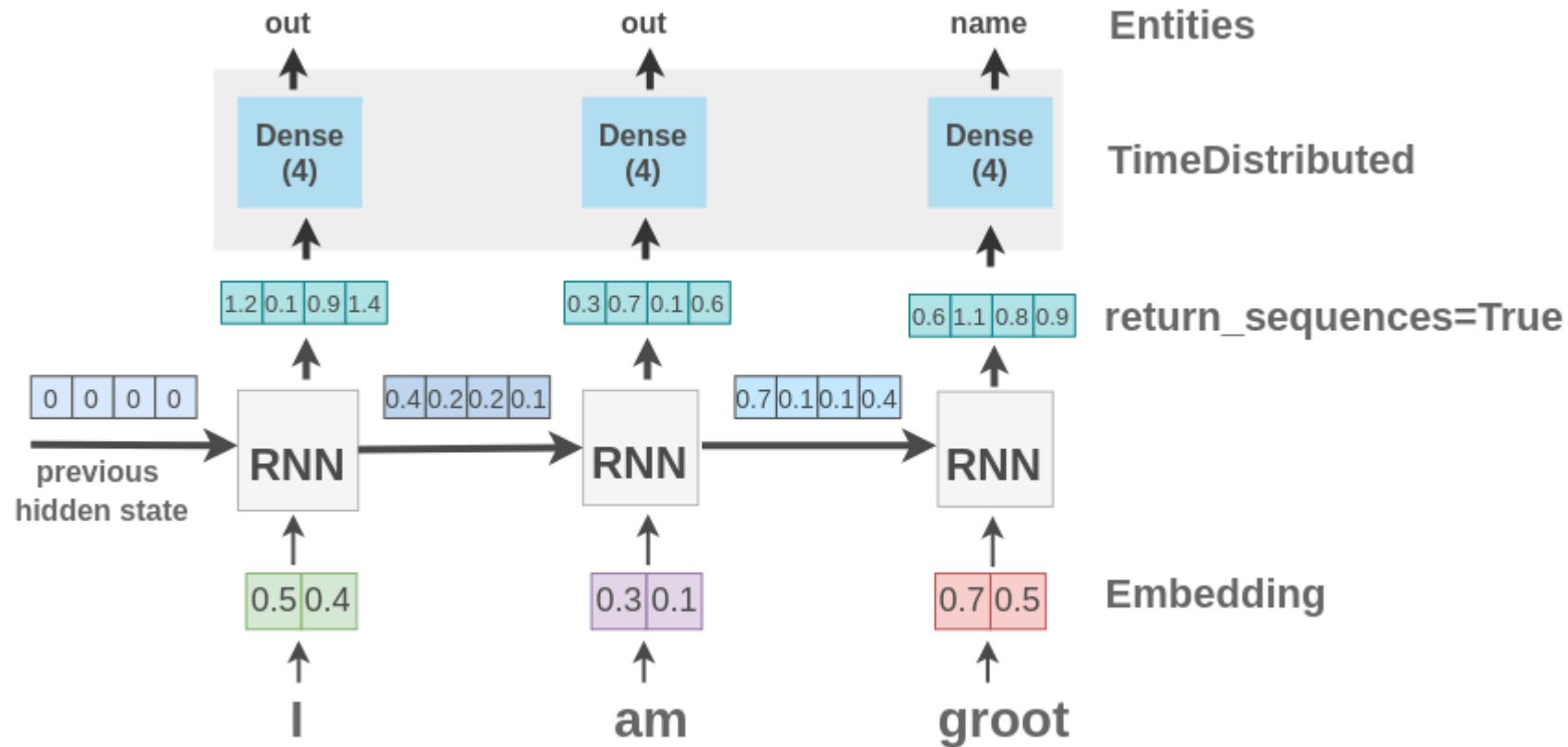
- ❖ Suppose we want to recognize entities in a text.
  - ✓ For example, in our text "I am Groot", we want to identify "Groot" as a name.

Identify entity

I   am   groot  
NAME

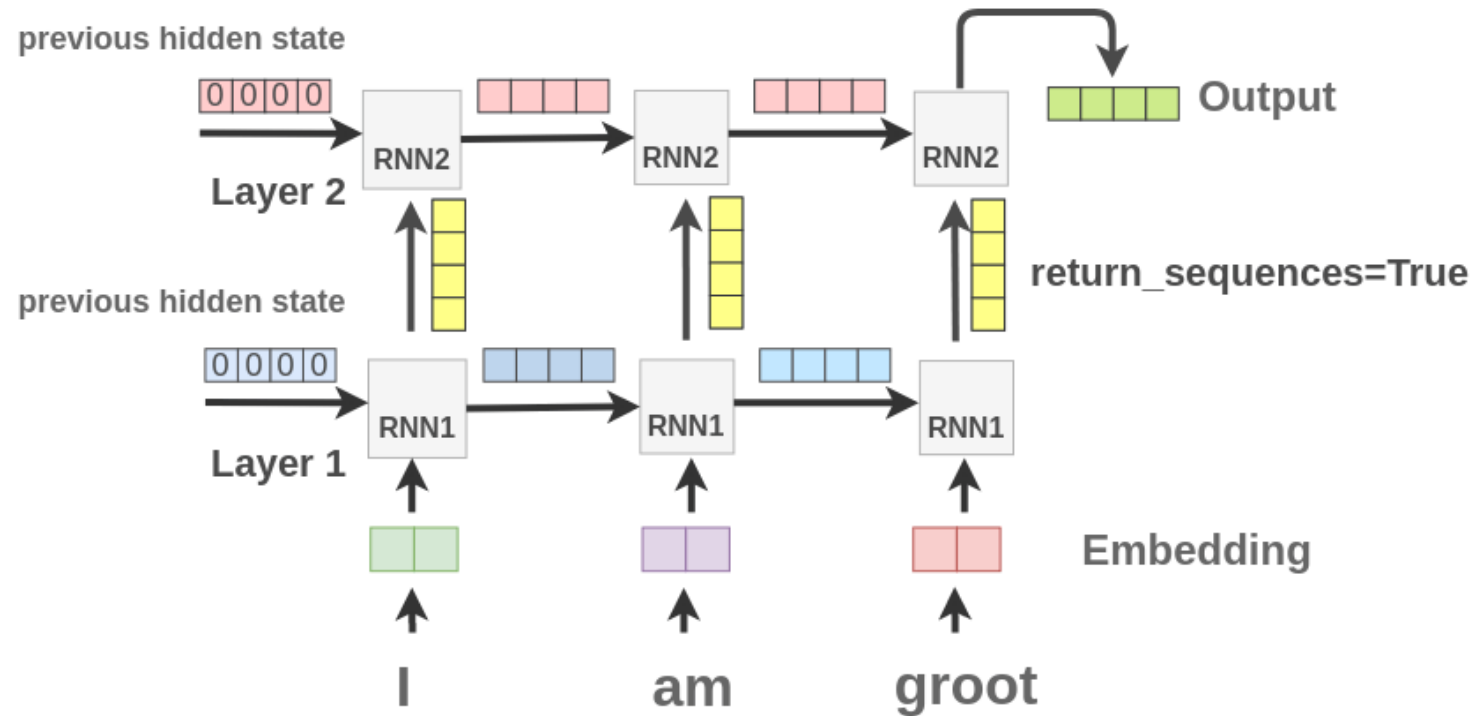
### (3) RNN: TimeDistributed Layer

```
model.add(SimpleRNN(4, input_shape=(3, 2), return_sequences=True))  
model.add(TimeDistributed(Dense(4, activation='softmax')))
```



- ❖ We can also stack multiple recurrent layers one after another in Keras

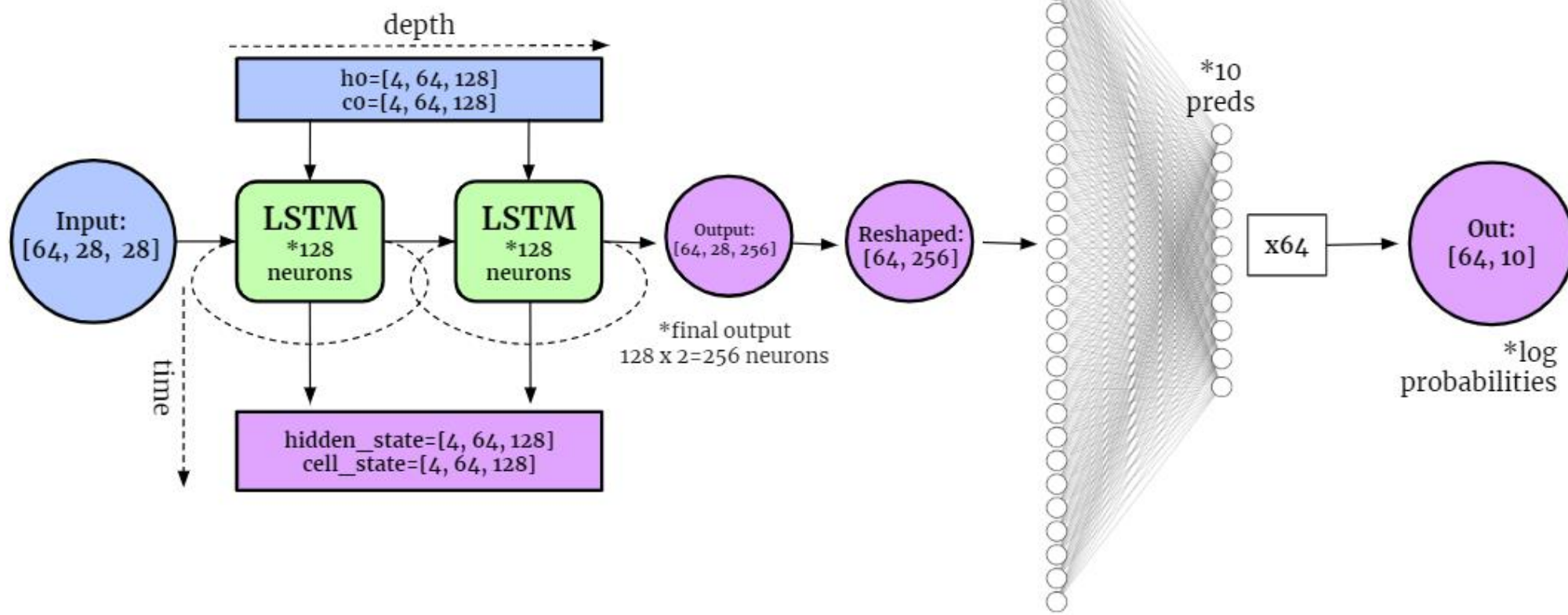
```
model.add(SimpleRNN(4, input_shape=(3, 2), return_sequences=True))  
model.add(SimpleRNN(4))
```



# Let's Code 2: LSTM

## LSTM Example

batch\_size=64  
input\_size=28  
hidden\_size=128  
layer\_size =2  
output\_size=10



```
from keras.models import Model
from keras.layers import Input, Dense, LSTM
import numpy as np

x = np.array([[[1.], [2.], [3.], [4.], [5.]])
y = np.array([[6.]])

xInput = Input(batch_shape=(None, 5, 1))
xLstm = LSTM(3)(xInput)
xOutput = Dense(1)(xLstm)

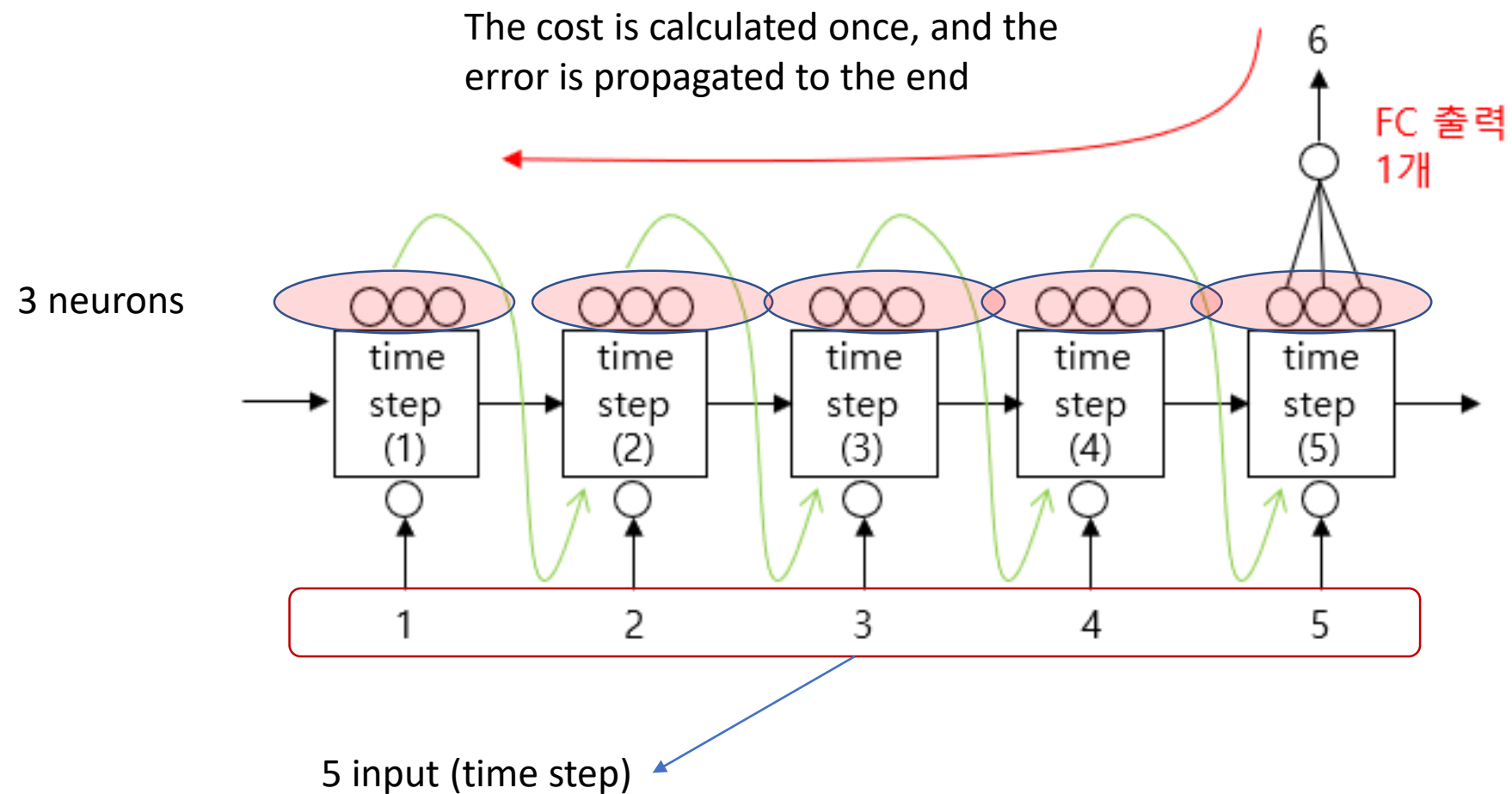
model = Model(xInput, xOutput)
model.compile(loss='mean_squared_error', optimizer='adam')
print(model.summary())

model.fit(x, y, epochs=50, batch_size=1, verbose=0)
model.predict(x, batch_size=1)
```

Model: "model\_4"

Layer (type)	Output Shape	Param #
=====		
input_5 (InputLayer)	[(None, 5, 1)]	0
lstm_6 (LSTM)	(None, 3)	60
dense_6 (Dense)	(None, 1)	4
=====		
Total params: 64		
Trainable params: 64		
Non-trainable params: 0		

return\_sequences=False



# LSTM many-to-many with TimeDistributed Layer

```
import tensorflow as tf
from tensorflow import keras
import numpy as np
x = np.array([[[1.], [2.], [3.], [4.], [5.]])
y = np.array([[[2.], [3.], [4.], [5.], [6.]])
```

```
model2 = keras.models.Sequential([
    keras.layers.LSTM(3, return_sequences=True, input_shape=[5, 1]),
    keras.layers.TimeDistributed(keras.layers.Dense(1) )
])
model2.compile(loss='mean_squared_error', optimizer='adam')
model2.summary()
```

Model: "sequential\_7"

Layer (type)	Output Shape	Param #
=====		
lstm_18 (LSTM)	(None, 5, 3)	60
time_distributed_8 (TimeDistributed)	(None, 5, 1)	4

=====

Total params: 64

Trainable params: 64

Non-trainable params: 0

LSTM Output is 5x3(neurons)

Final Output is 5x1



2022

Korea Institute of Science  
and Technology Information

TRUST  
KISTIL

