



THE UNIVERSITY OF
MELBOURNE

Temporal Convolutional Network & Recurrent Neural Network

COMP90051 Statistical Machine Learning

QiuHong Ke

Copyright: University of Melbourne

Sit down or stand up?

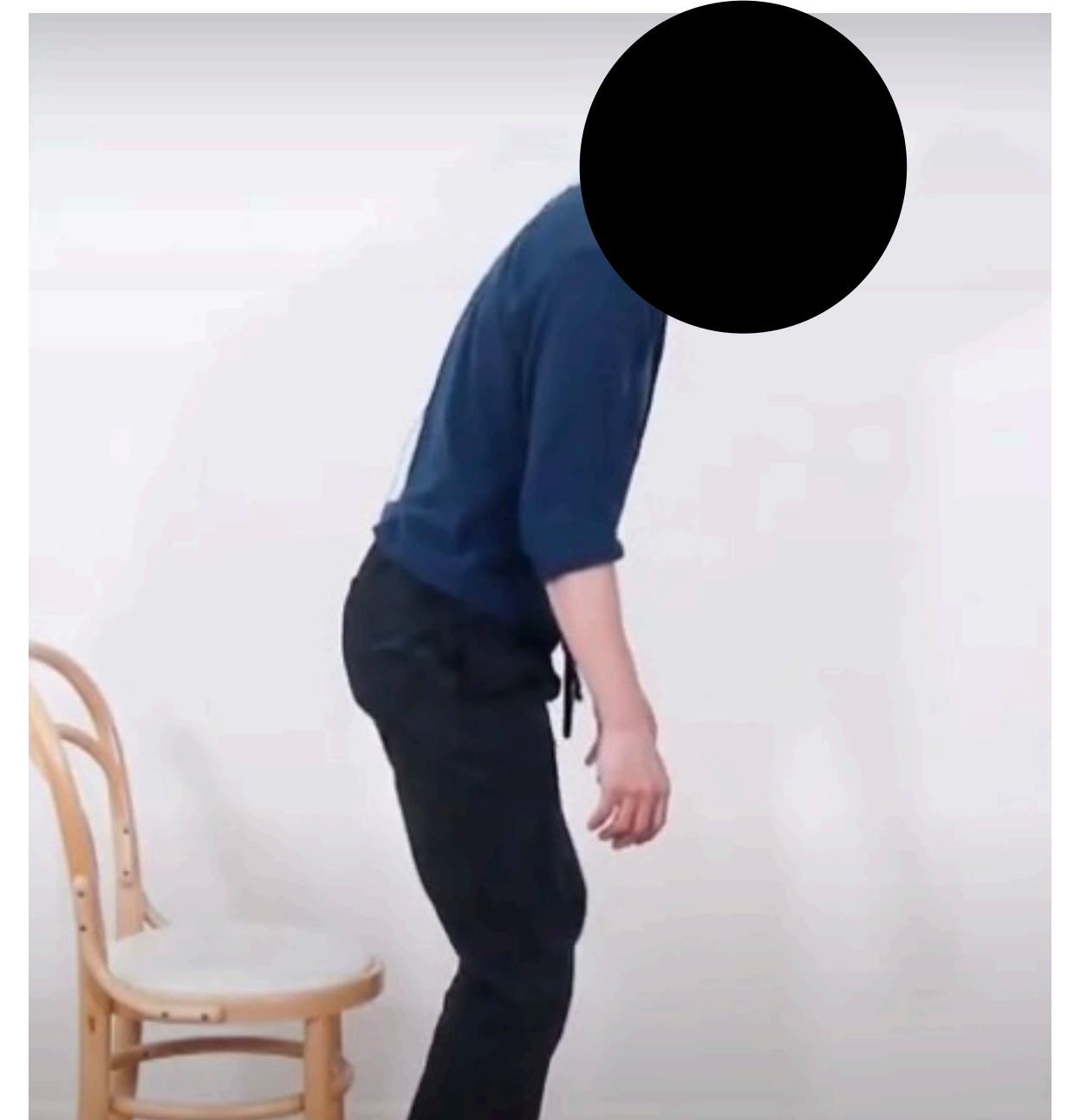
We need to learn the temporal information to understand the action



...



...



temporal information: temporal evolution (changes) of the human pose

Temporal information does matter

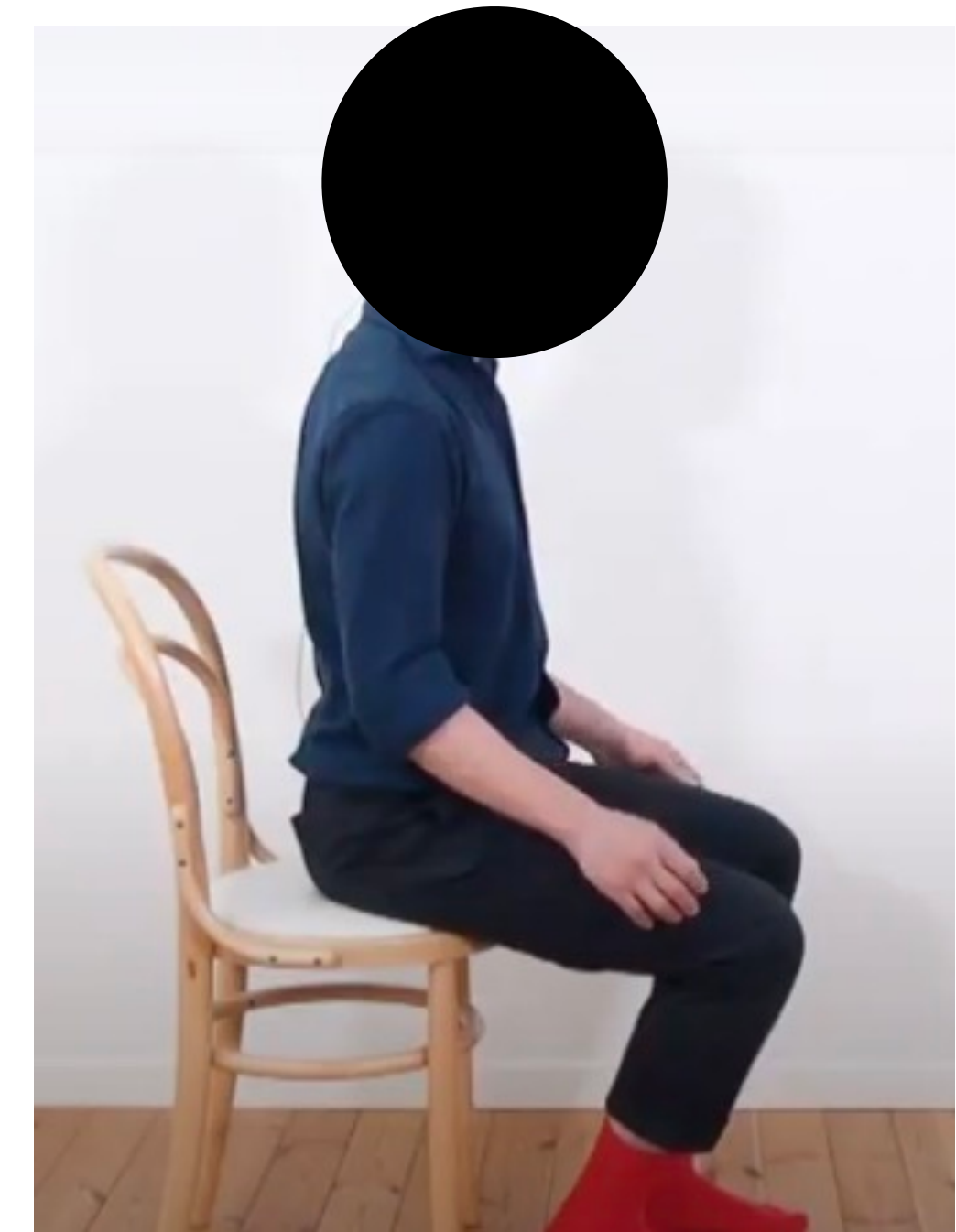
Change the order: different actions



...



...



What is the next word ?

This game is really fun

This cookie is really tasty

Like?Dislike? Sentiment analysis

I bought this charger in Jul 2003 and it worked OK for a while. The design is nice and convenient. 😊

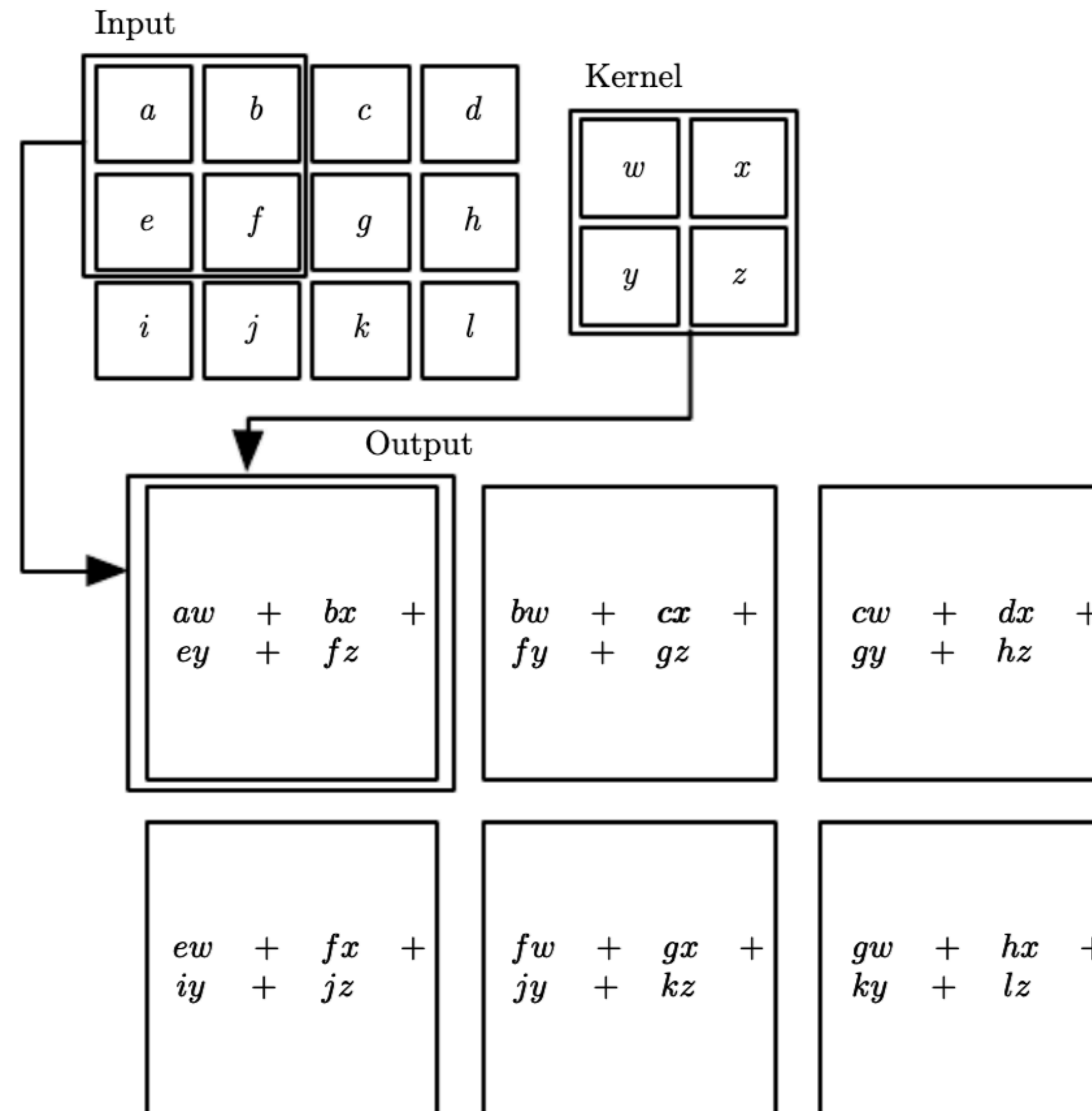
However, after about a year, the batteries would not hold a charge. Might as well just get alkaline disposables, or look elsewhere for a charger that comes with batteries that have better staying power. 😞

It is important to understand the context the sentence

Outline

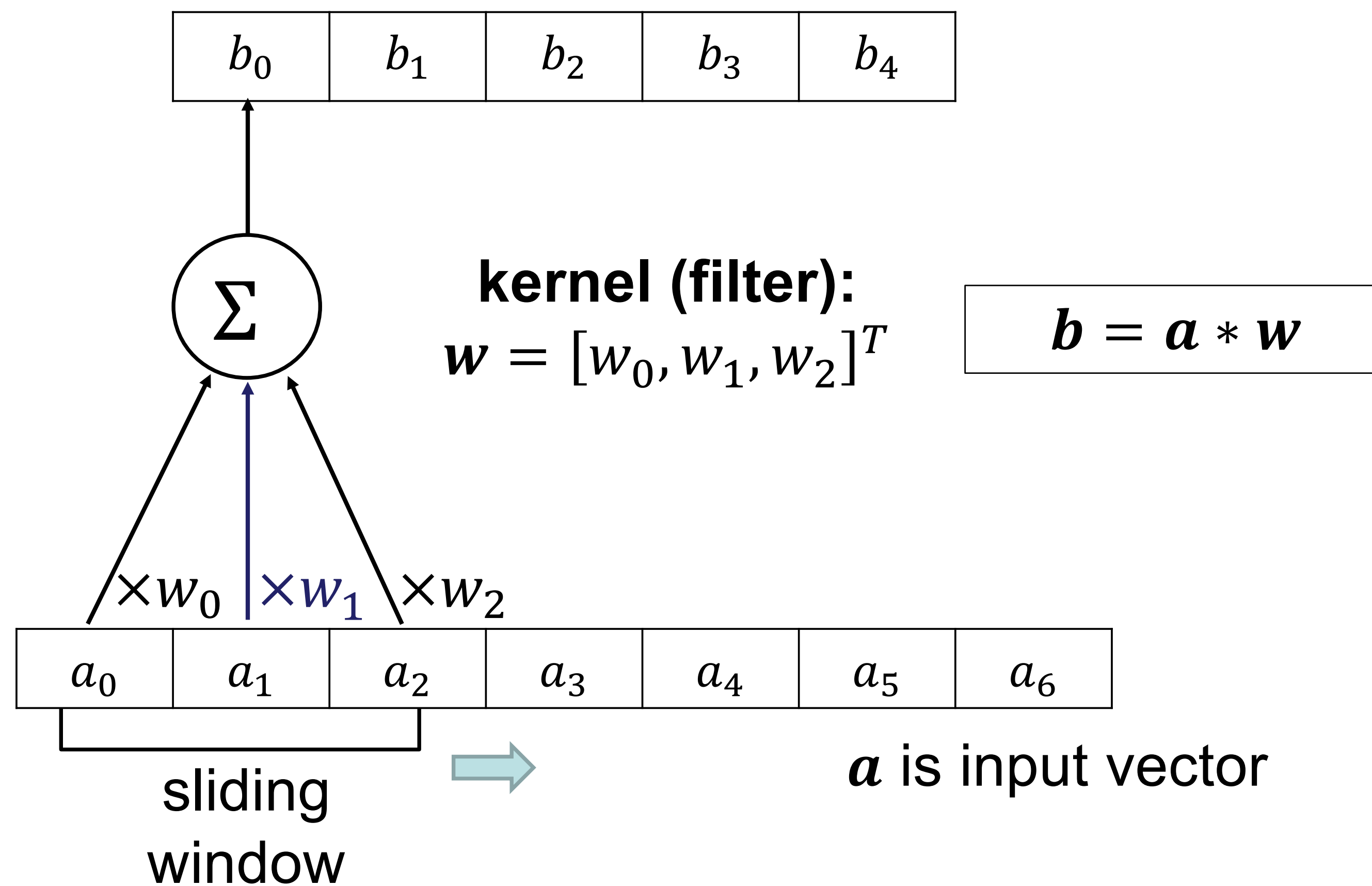
- Temporal convolutional network (TCN)
- Recurrent neural network (RNN)

Recap: 2D Convolution

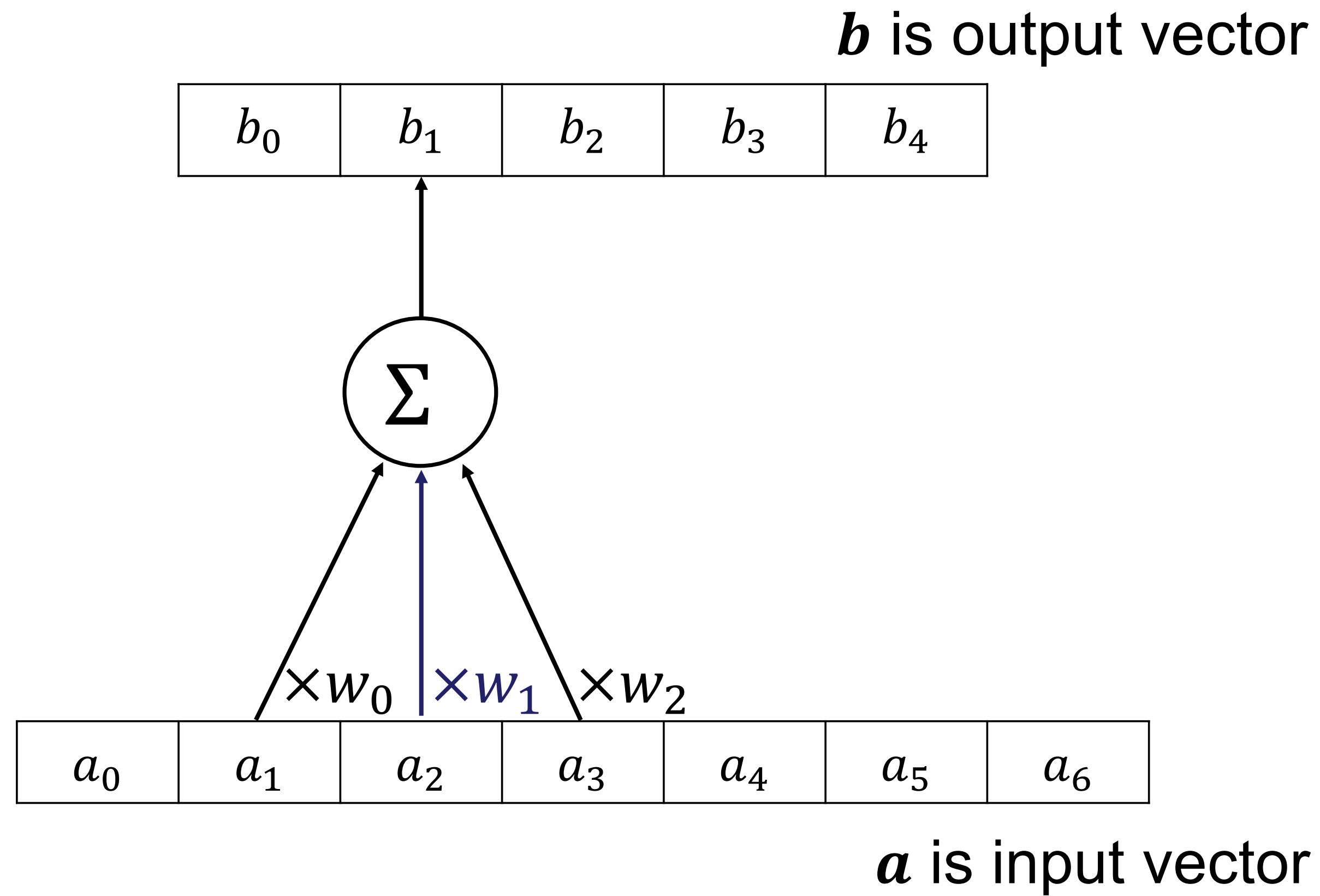


Recap: 1D Convolution

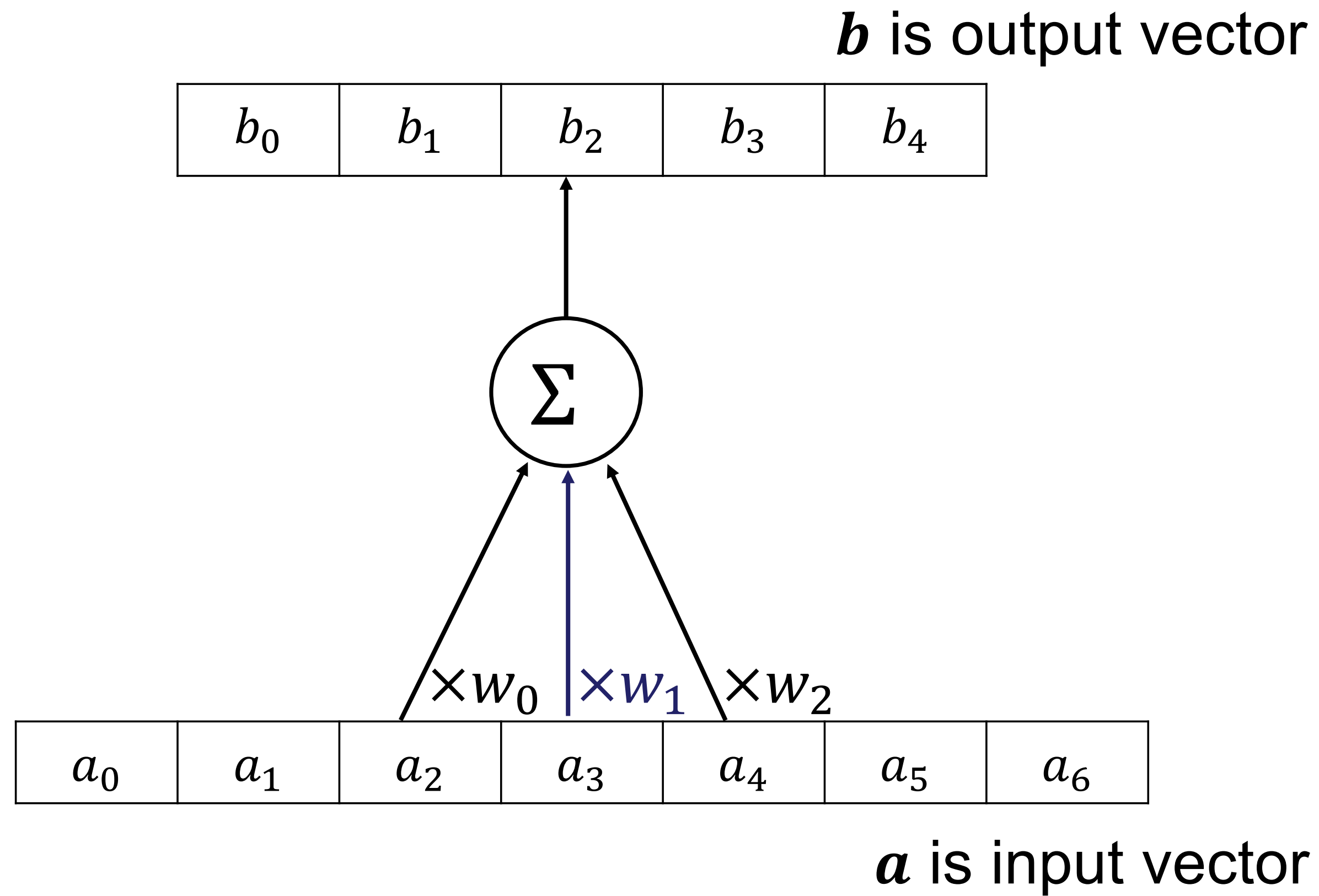
b is output vector (feature map)



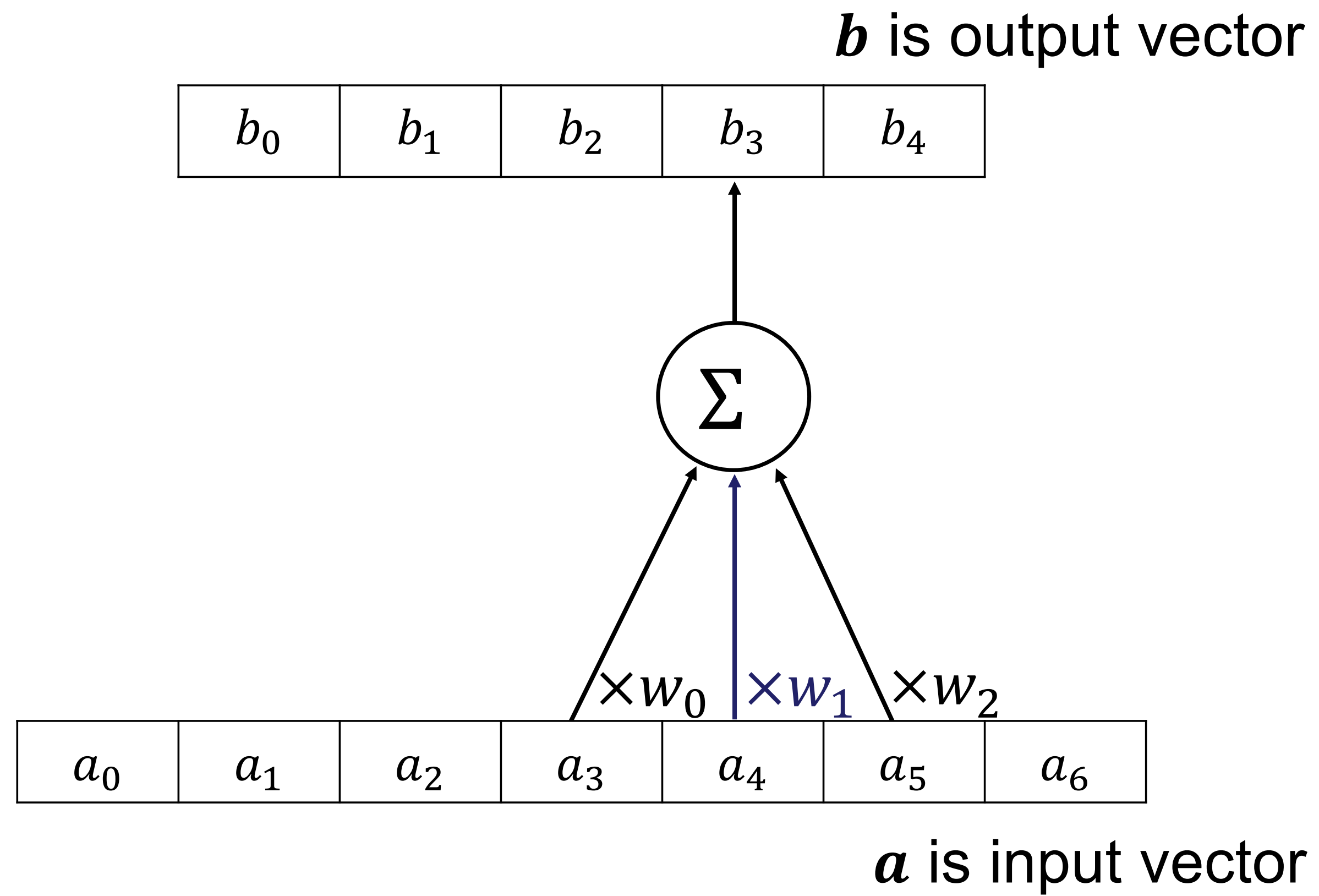
Recap: 1D Convolution



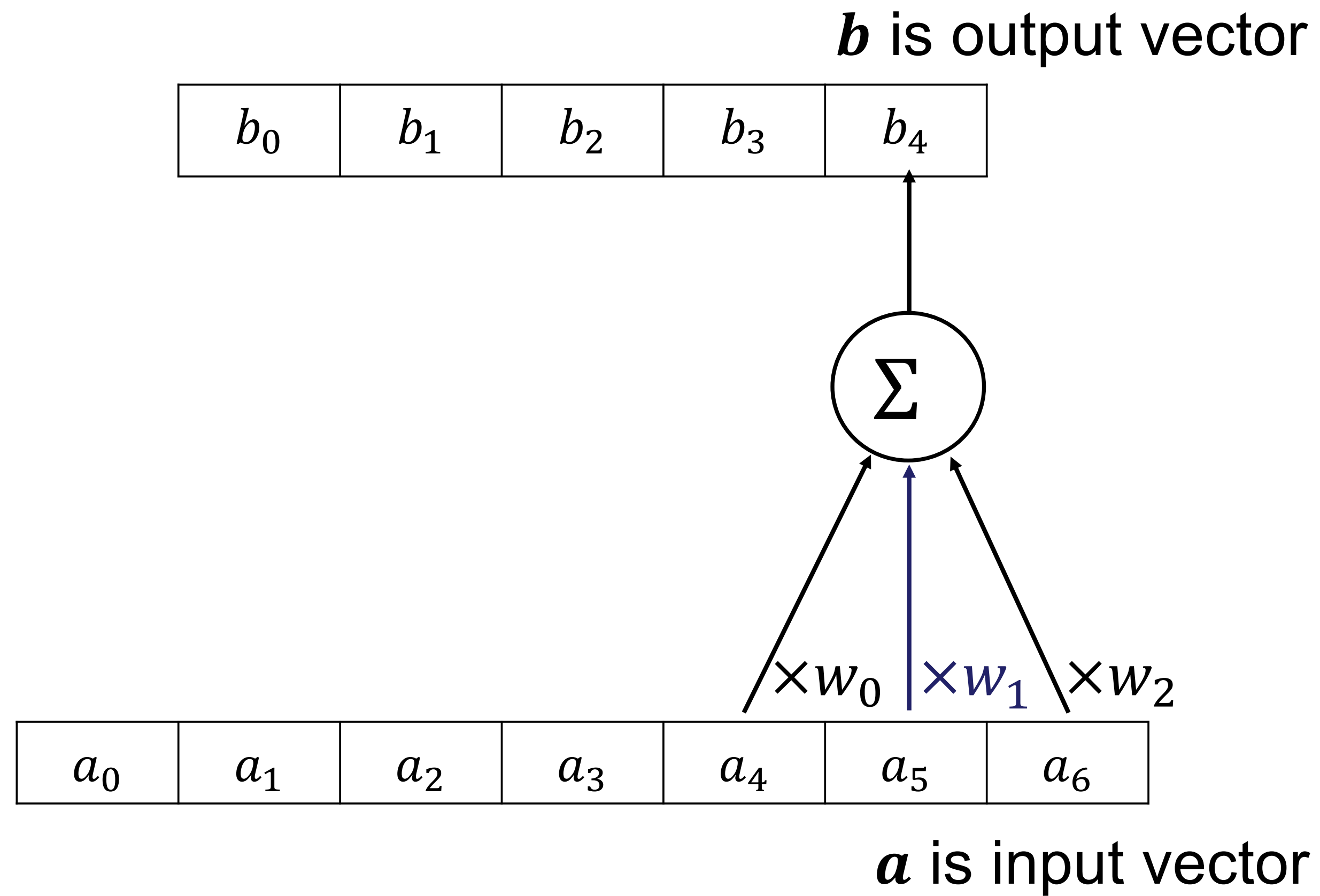
Recap: 1D Convolution



Recap: 1D Convolution

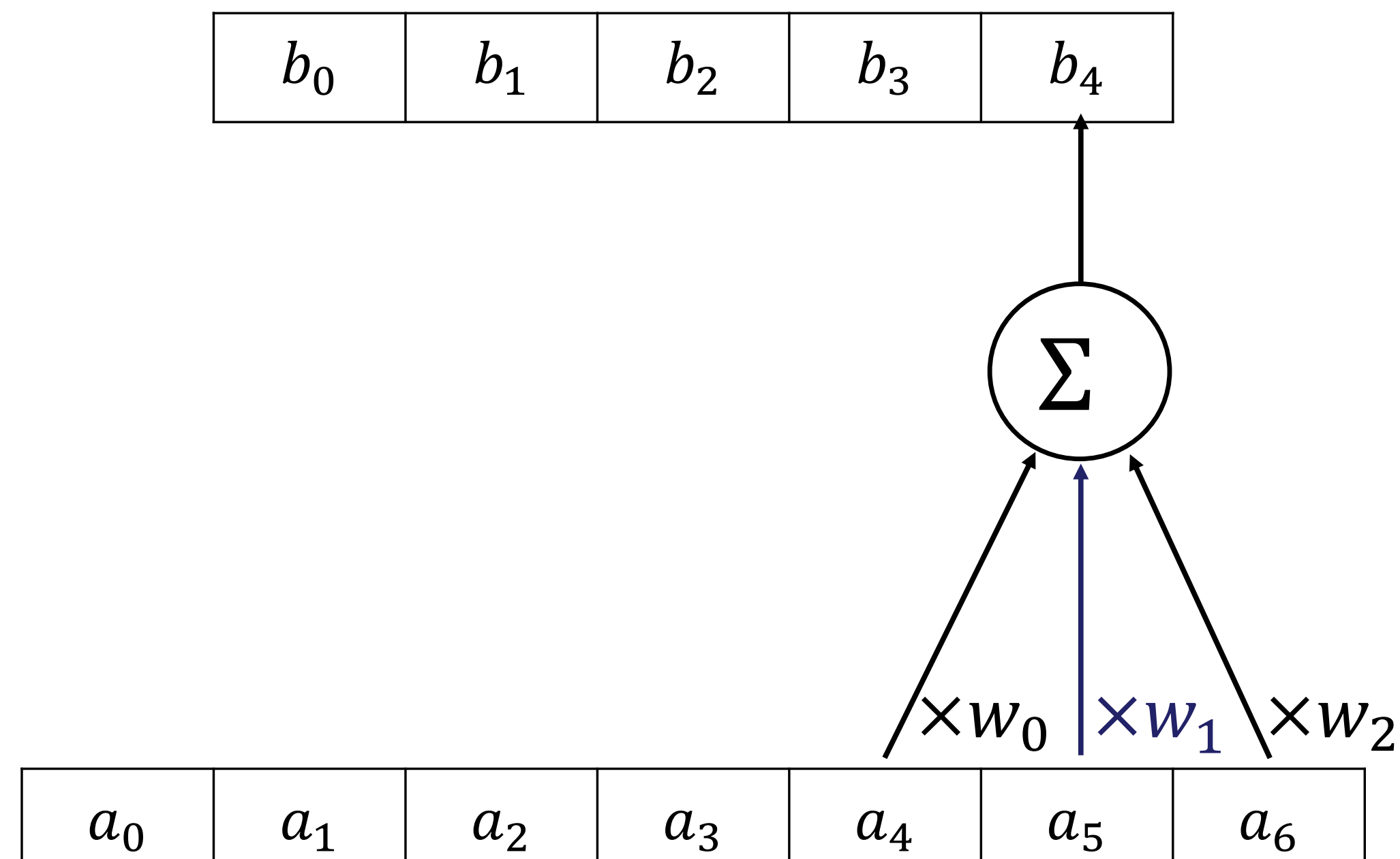


Recap: 1D Convolution



Temporal convolution: 1D convolution (on the temporal dimension)

1D input feature vector:
special case of sequence data



A sequence that consists of 7 time-steps:
One feature in each step

Temporal convolution on a sequence of feature vectors

Input:

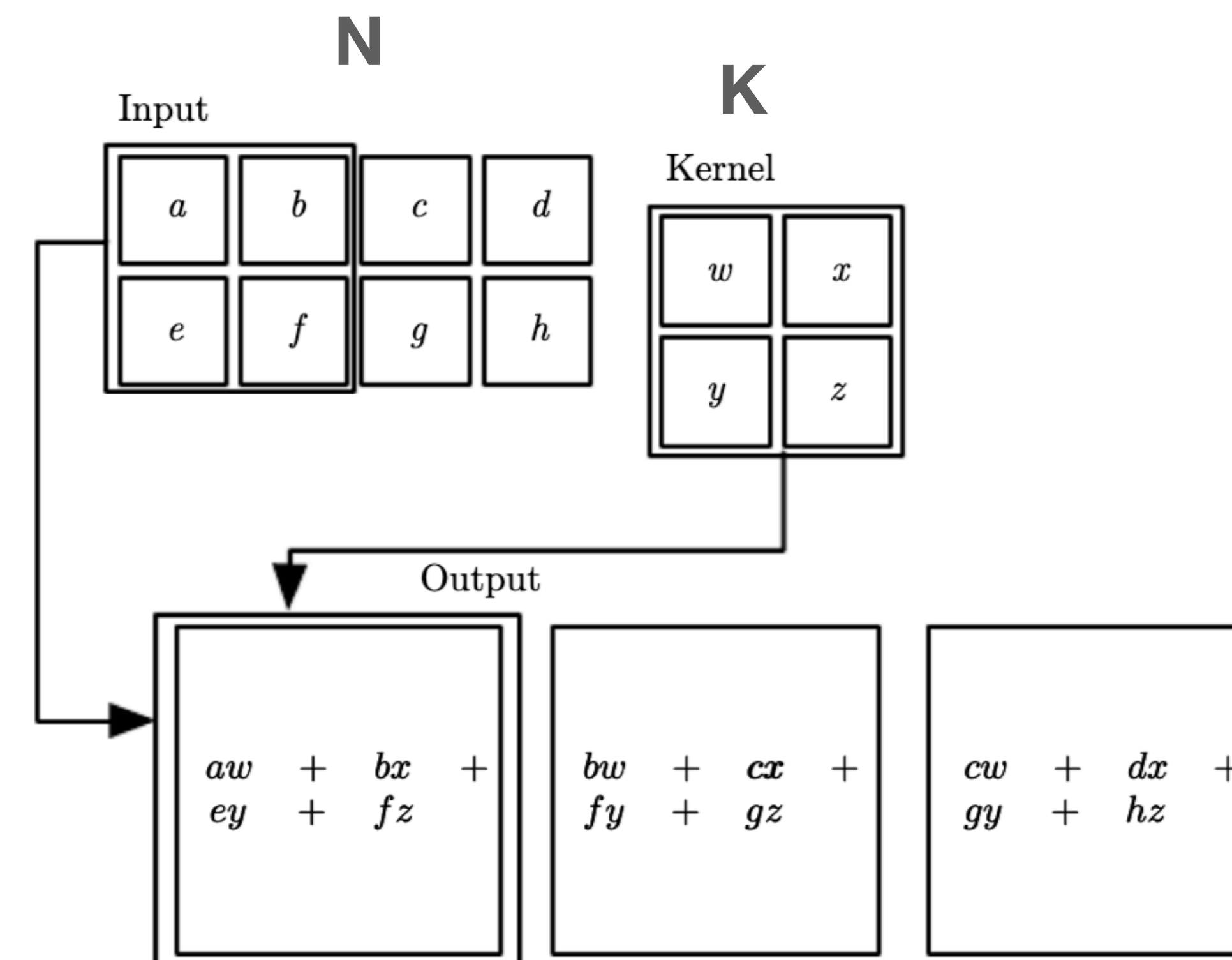
- a sequence of $N=4$ time-steps
- Each time-step is a feature vector (dimension: $d=2$ (2 features))

Kernel: a weight matrix ($d \times K$)

- One dimension is the same as d
- Another dimension: size of temporal convolutional window K

Output: A sequence

- Each time-step is a **scalar (1 kernel)**
- The number of time-steps (same as 2D):
 - No padding: ceiling $(N-K+1)/\text{stride}$
 - With padding: ceiling (N/stride)



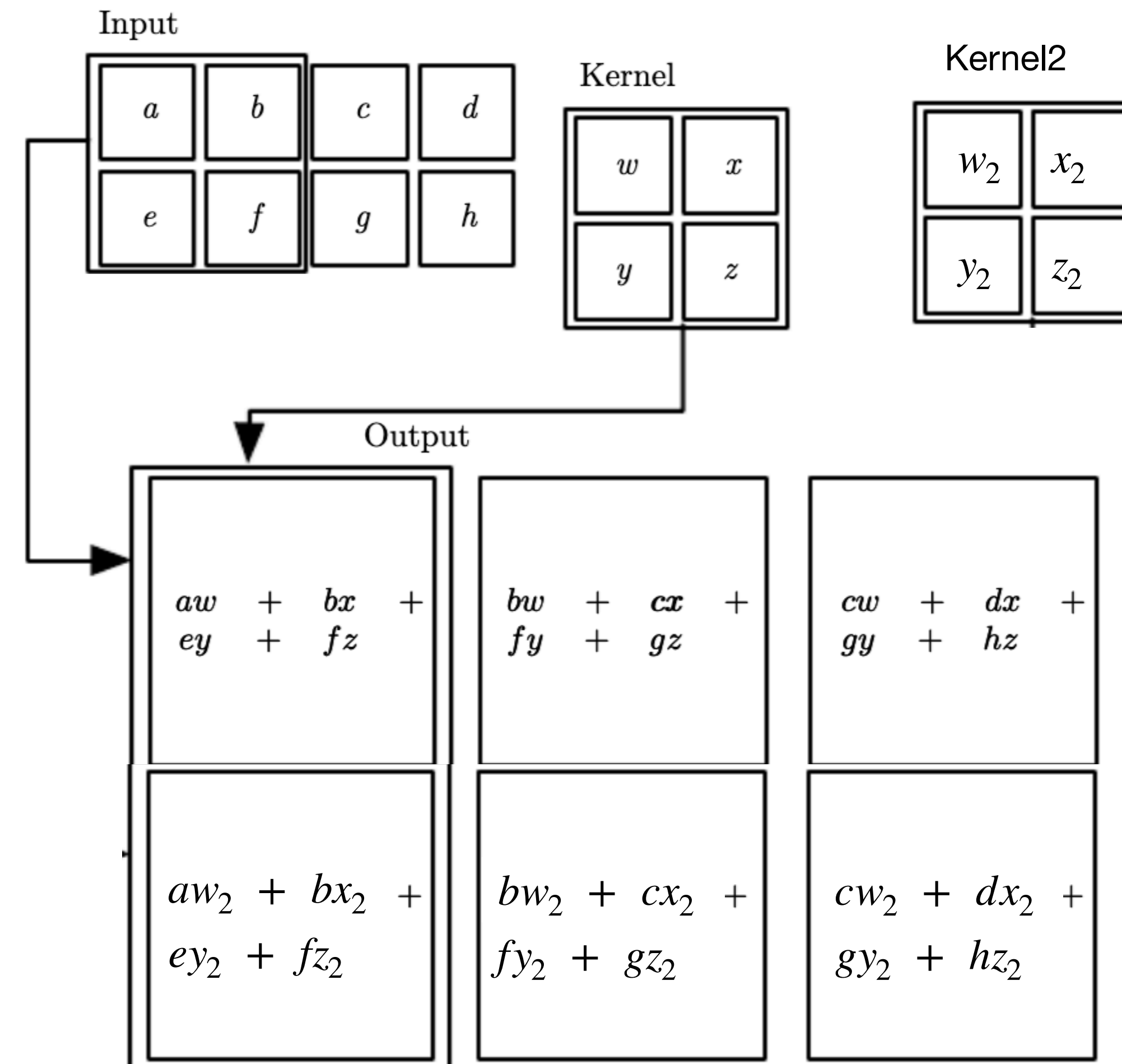
Temporal convolution on a sequence of feature vectors

If there are m kernels in the layer

Output: m sequences ==

a sequence of feature vectors:

- The number of time-steps:
 - No padding: ceiling $(N-K+1)/\text{stride}$
 - With padding: ceiling (N/stride)
- Each time-step is a vector
 - Dimension of each vector == kernel number m



Temporal convolution for text

I bought this charger in Jul 2003 and it worked OK for a while. The design is nice and convenient. However, after about a year, the batteries would not hold a charge. Might as well just get alkaline disposables, or look elsewhere for a charger that comes with batteries that have better staying power.

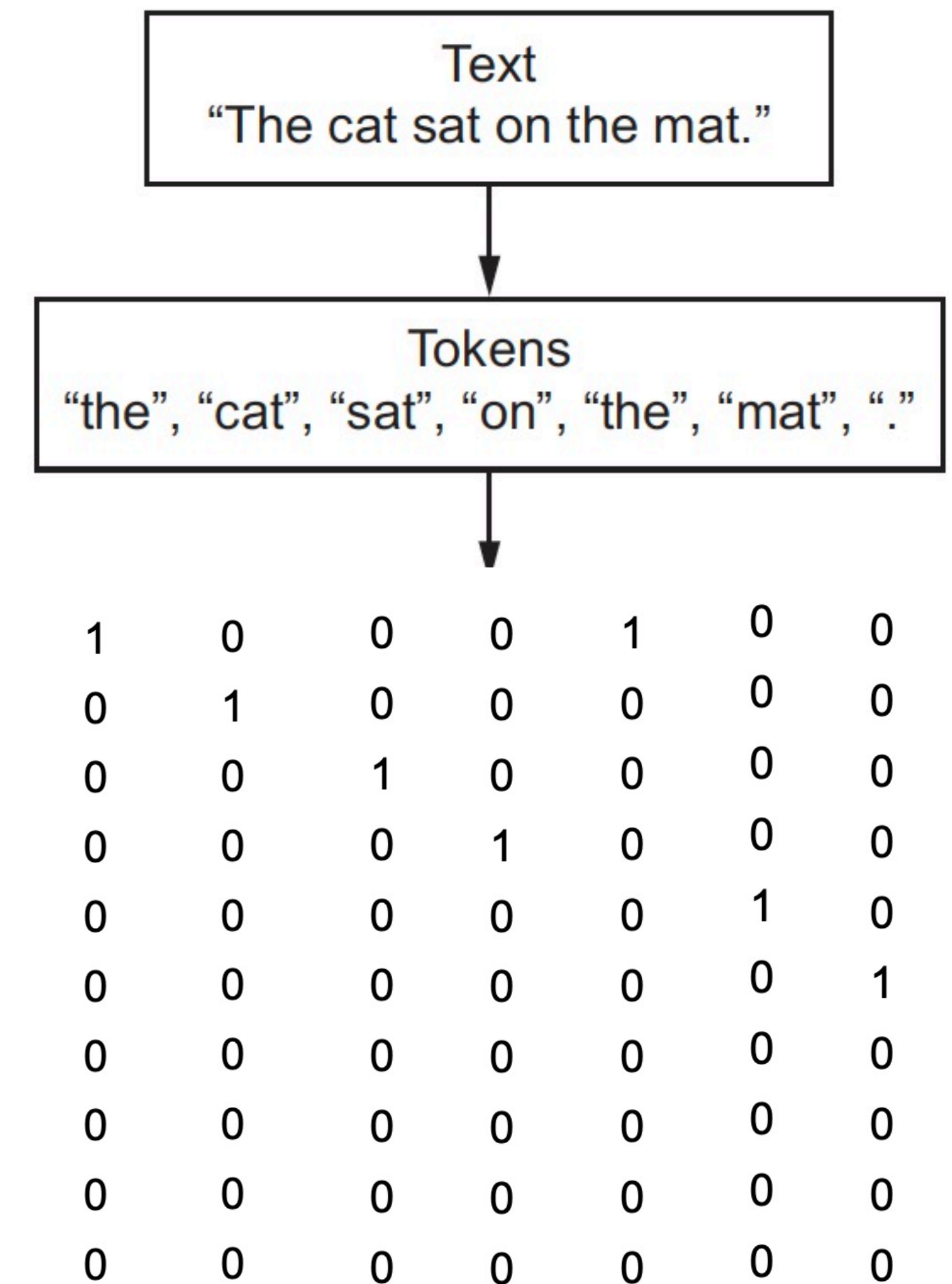
How to represent words at each time-step before temporal convolution

- One-hot encoding
- Word embedding

One-hot encoding

- Create a feature vector
- Dimension of the vector == size of vocabulary
- If the word is the i^{th} word in the vocabulary, then the i^{th} element of the vector is 1, all the others are 0
- Sparse and high-dimensional

dic={"the", "cat", "sat", "on", "mat", ".",
"these", "are", "other", "words"}



Word embedding: learn low-dimensional embedding from data

One-hot word vectors:

- Sparse
- High-dimensional
- Hardcoded



Word embeddings:

- Dense
- Lower-dimensional
- Learned from data

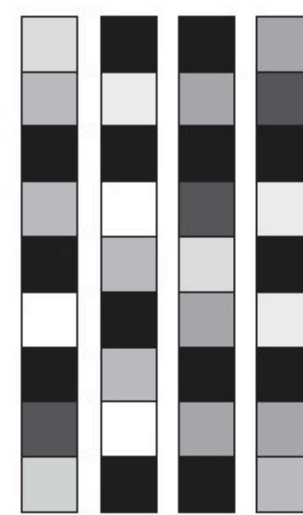


Figure 6.3 in Deep learning with python by Francois Chollet

Text
"The cat sat on the mat."

Tokens
"the", "cat", "sat", "on", "the", "mat", "."

Vector encoding of the tokens						
0.0	0.0	0.4	0.0	0.0	1.0	0.0
0.5	1.0	0.5	0.2	0.5	0.5	0.0
1.0	0.2	1.0	1.0	1.0	0.0	0.0
the	cat	sat	on	the	mat	.

Figure 6.1 in Deep learning with python by Francois Chollet

Word embedding: learn low-dimensional embedding from data

```
Emb_layer = tf.keras.layers.Embedding(vocabulary_size, emb_dim)
```

Word index → Embedding layer → Corresponding word vector

Embedding layer:

- Create a weight matrix of d (emb_dim) rows, N (vocabulary_size) columns
- Return the i^{th} column as the feature vector for the i^{th} word
- Learn word embeddings jointly with the main task or load pretrained word embeddings

Example

Input:

Word
embedding

I	love	deep	learning
0.1	0.4	0.2	-0.5
0.3	0.3	0.1	0.3

Kernel:

1	0	0
0	0	1

Output:

element-wise
multiplication
and sum

0.2

Example

Input:

Word
embedding

I	love	deep	learning
0.1	0.4	0.2	-0.5
0.3	0.3	0.1	0.3

Kernel:

1	0	0
0	0	1

Output:

element-wise
multiplication
and sum

0.2 0.7

Example

Input:

Word
embedding

I	love	deep	learning
0.1	0.4	0.2	-0.5
0.3	0.3	0.1	0.3

2 Kernels:

0	0	1
0	1	0

Output:

element-wise
multiplication
and sum

0.0	0.4
0.5	

Example

Input:

Word
embedding

I	love	deep	learning
0.1	0.4	0.2	-0.5
0.3	0.3	0.1	0.3

2 Kernels:

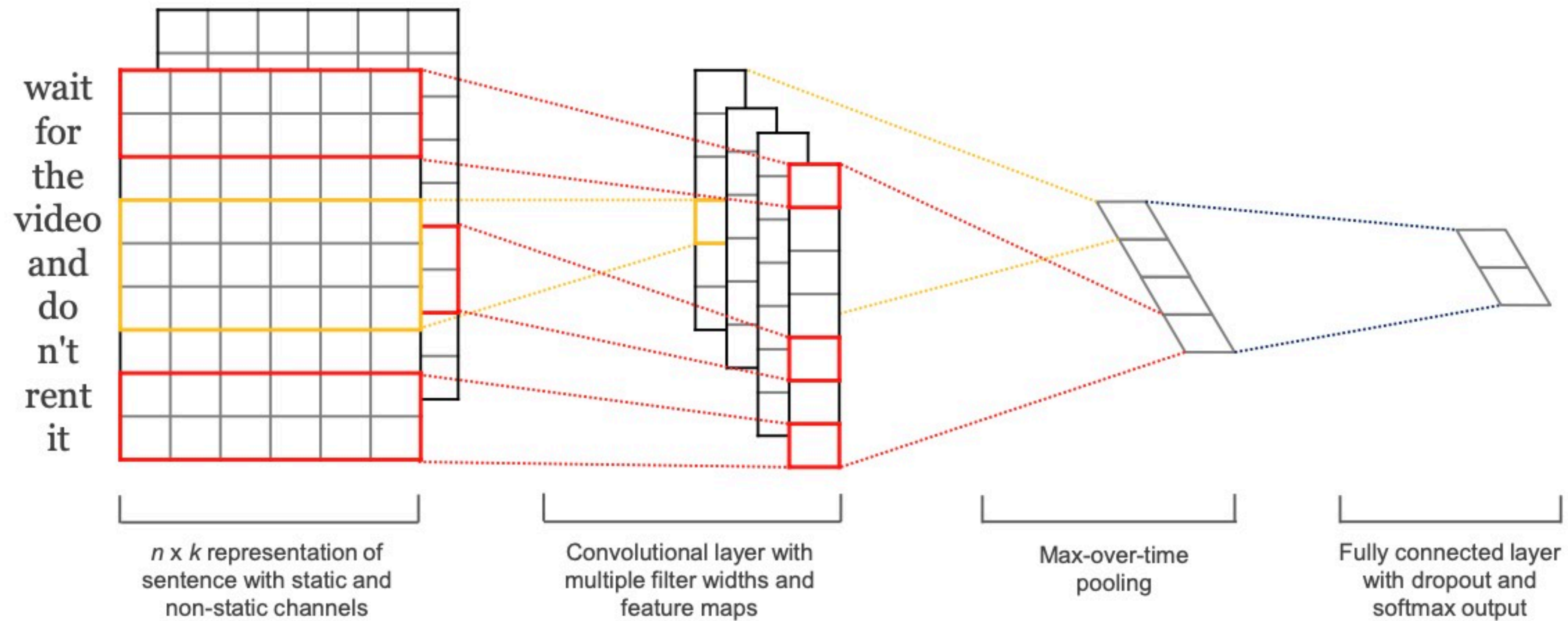
0	0	1
0	1	0

Output:

element-wise
multiplication
and sum

0.0	0.4
0.5	-0.4

TCN for Sentence Classification



Recap of Convolution on Multiple-channel input



R



Kernel: same channel (depth)

$* K(\text{Channel 1})$

G



$* K(\text{Channel 2})$ \longrightarrow Element-wise sum \longrightarrow One channel

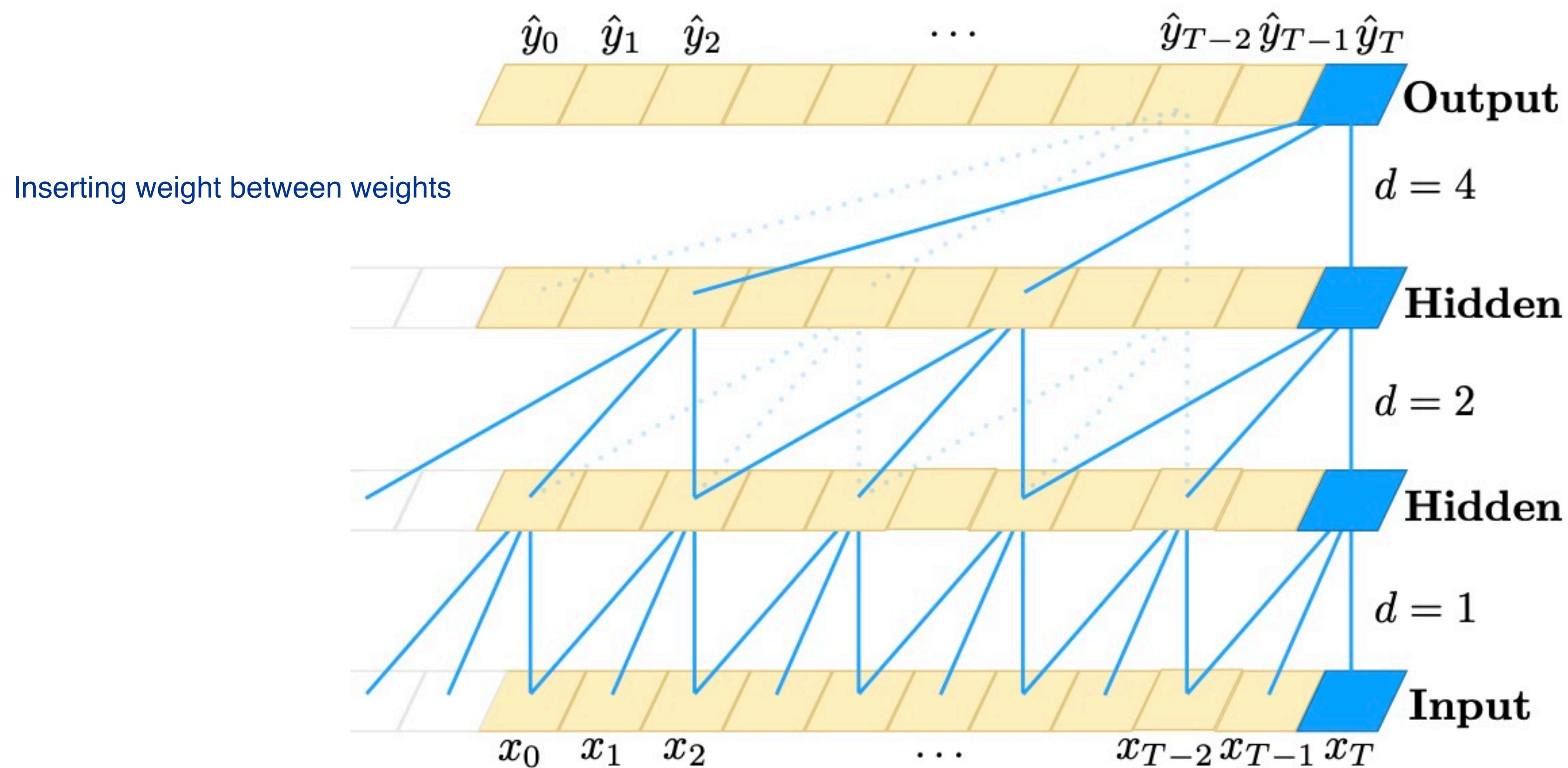
B



$* K(\text{Channel 3})$

A special version: Dilated casual convolution

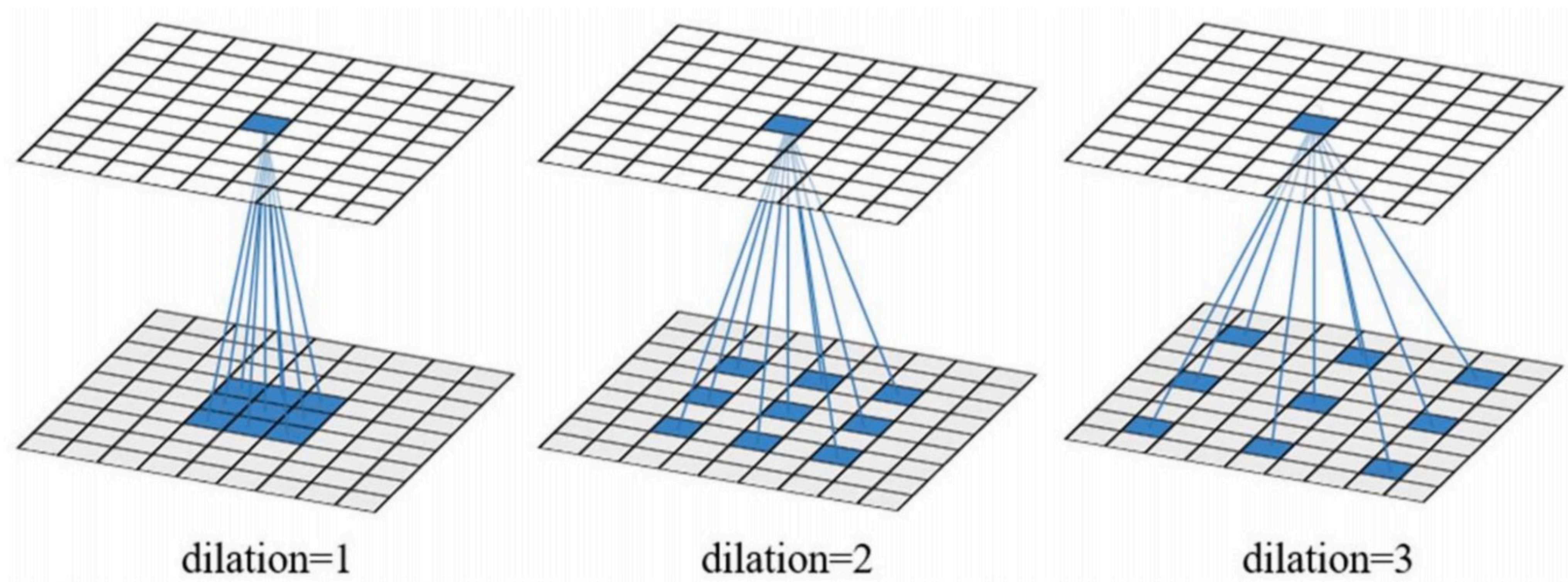
Dilated: expand the alignment of the kernel weights by dilation factor



Dilation factor (d): step between every two adjacent filter taps.

Larger dilation factor: the weights are placed more far away (i.e., more sparse),

Dilated convolution in 2D

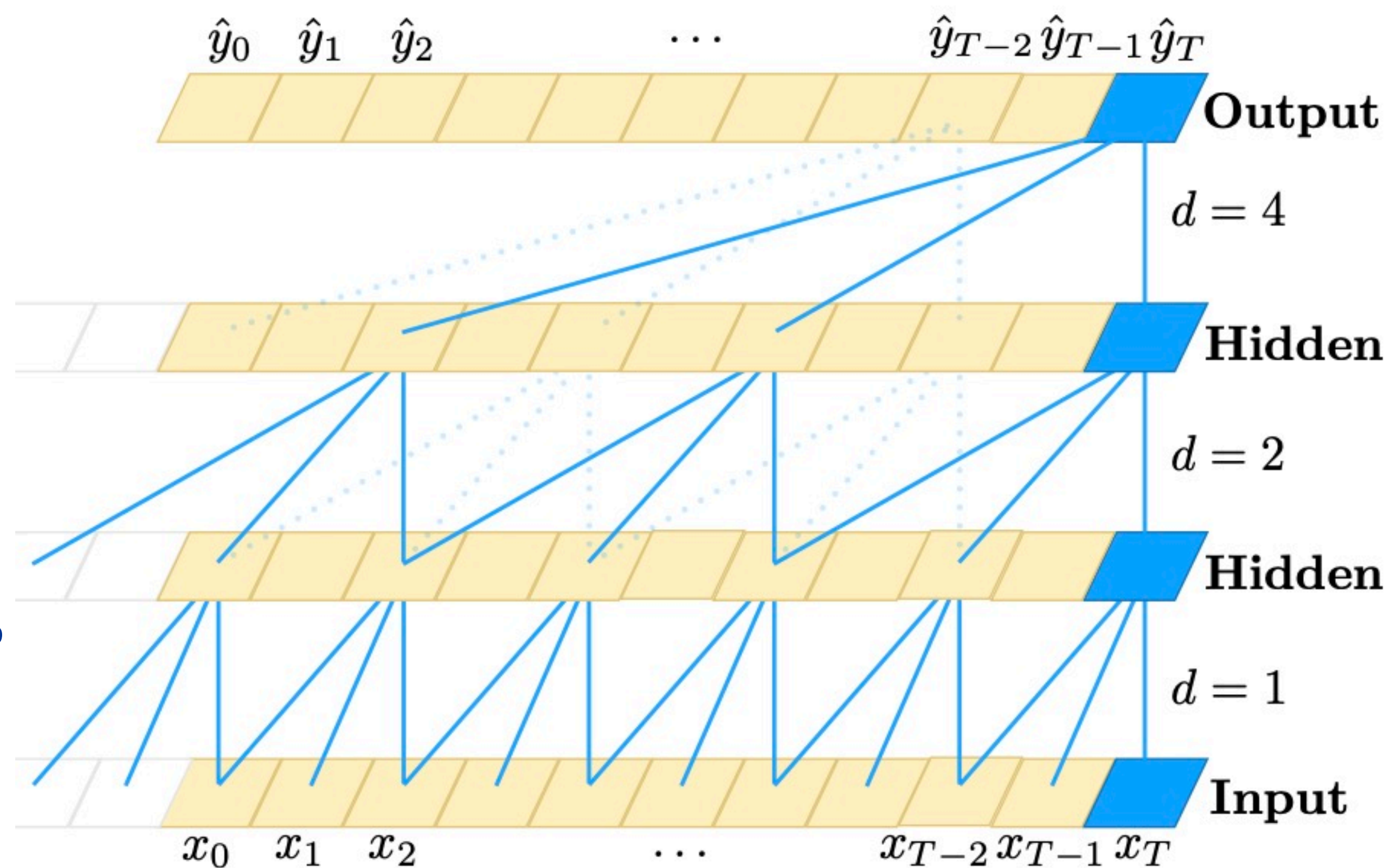


Increasing receptive fields

- Stacking multiple convolutional layers
- max-Pooling layer
- Dilated convolution

Dilated casual convolution

Causal: The output at time-step t does not depend on the input information after time-step t (e.g., real-time recognition, prediction)

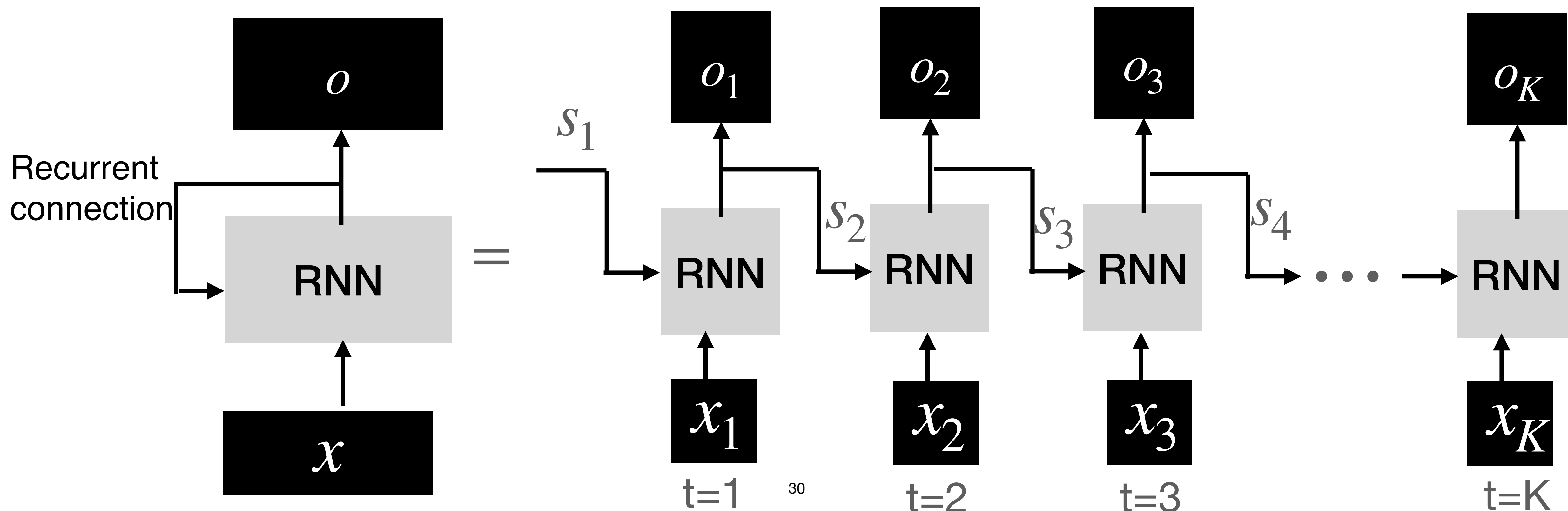


We perform padding between time step features where the current features are only depending on historical data and current data

Networks with loops

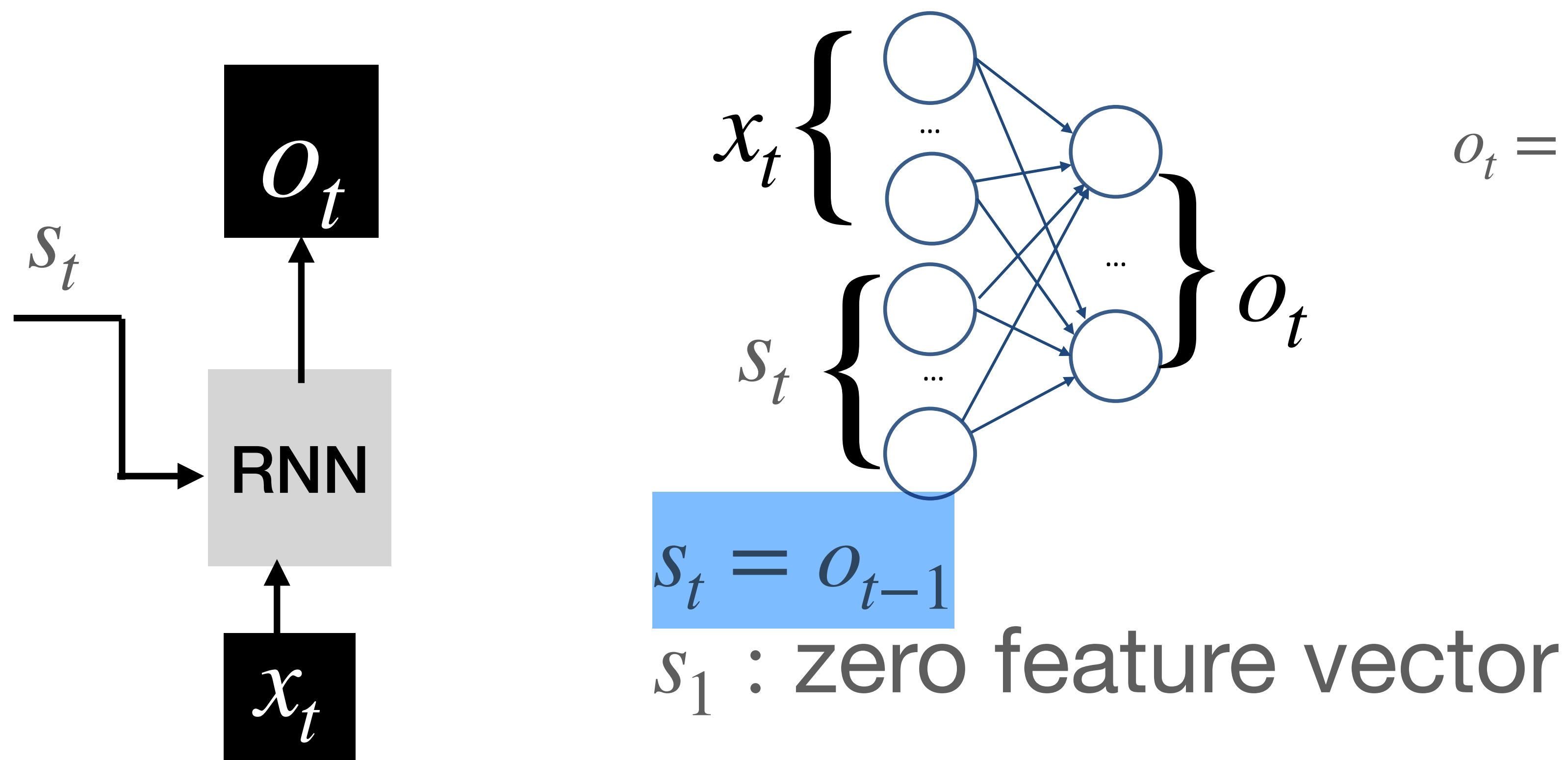
Process sequence step by step

Recurrent connection: the output hidden feature vector (state) of each step is connected to the next step



RNN: Process sequence step by step

At each step t : combine current timestep x_t (feature vector of input sequence at timestep t) and historical information, i.e., state s_t (feature vector) to generate o_t



$$o_t = \text{activation}(W \cdot x_t + U \cdot s_t + b)$$

Shared over time

A simple RNN

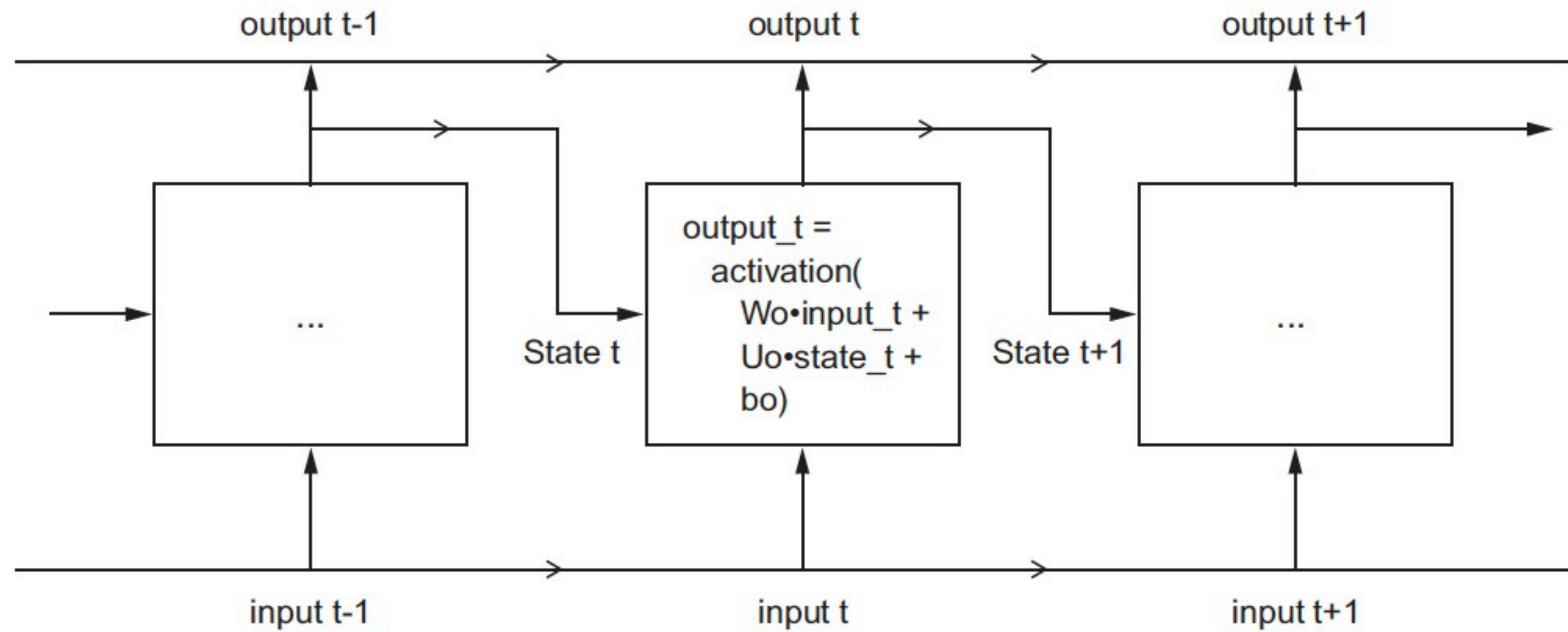
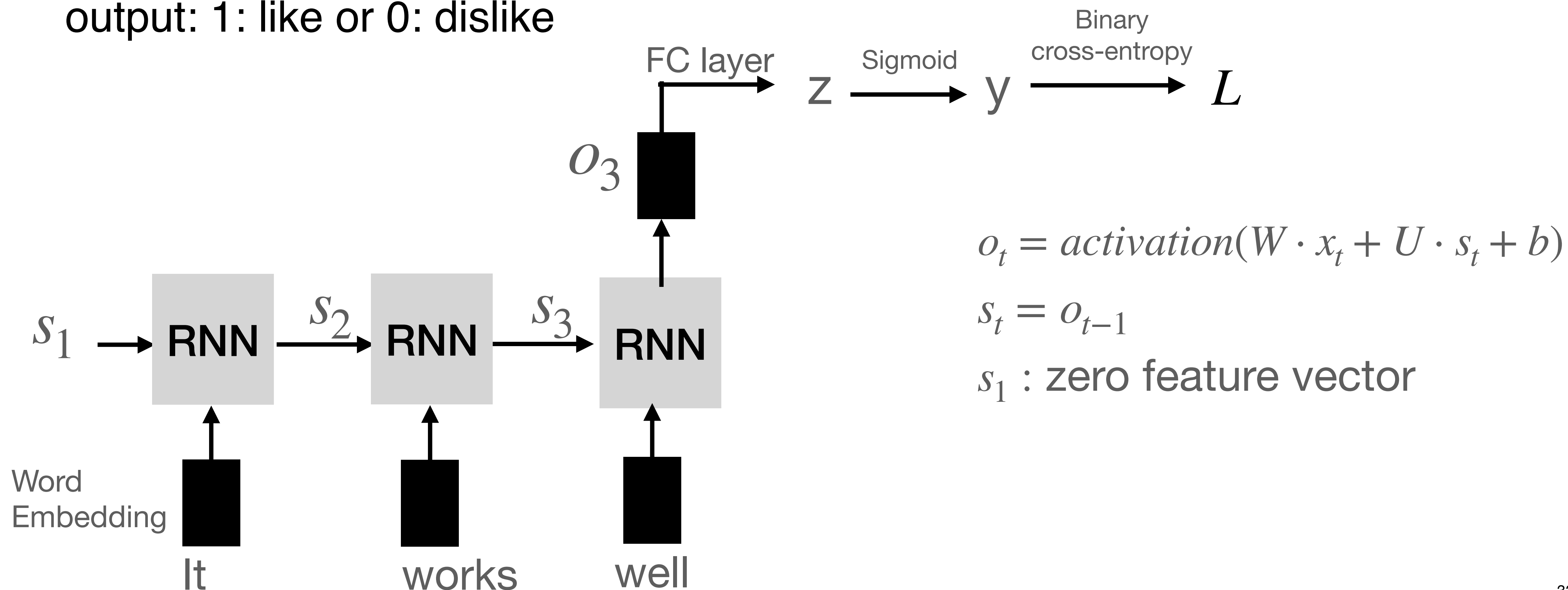


Figure 6.13 in Deep learning with python by Francois Chollet

Example: RNN for sentiment analysis

Input: review (a piece of text) e.g., “It works well”

output: 1: like or 0: dislike



Recap of Chain rule

$$\text{Given } z = g(u) \quad u = f(x) \quad \Rightarrow \quad \frac{dz}{dx} = \frac{dz}{du} \frac{du}{dx}$$

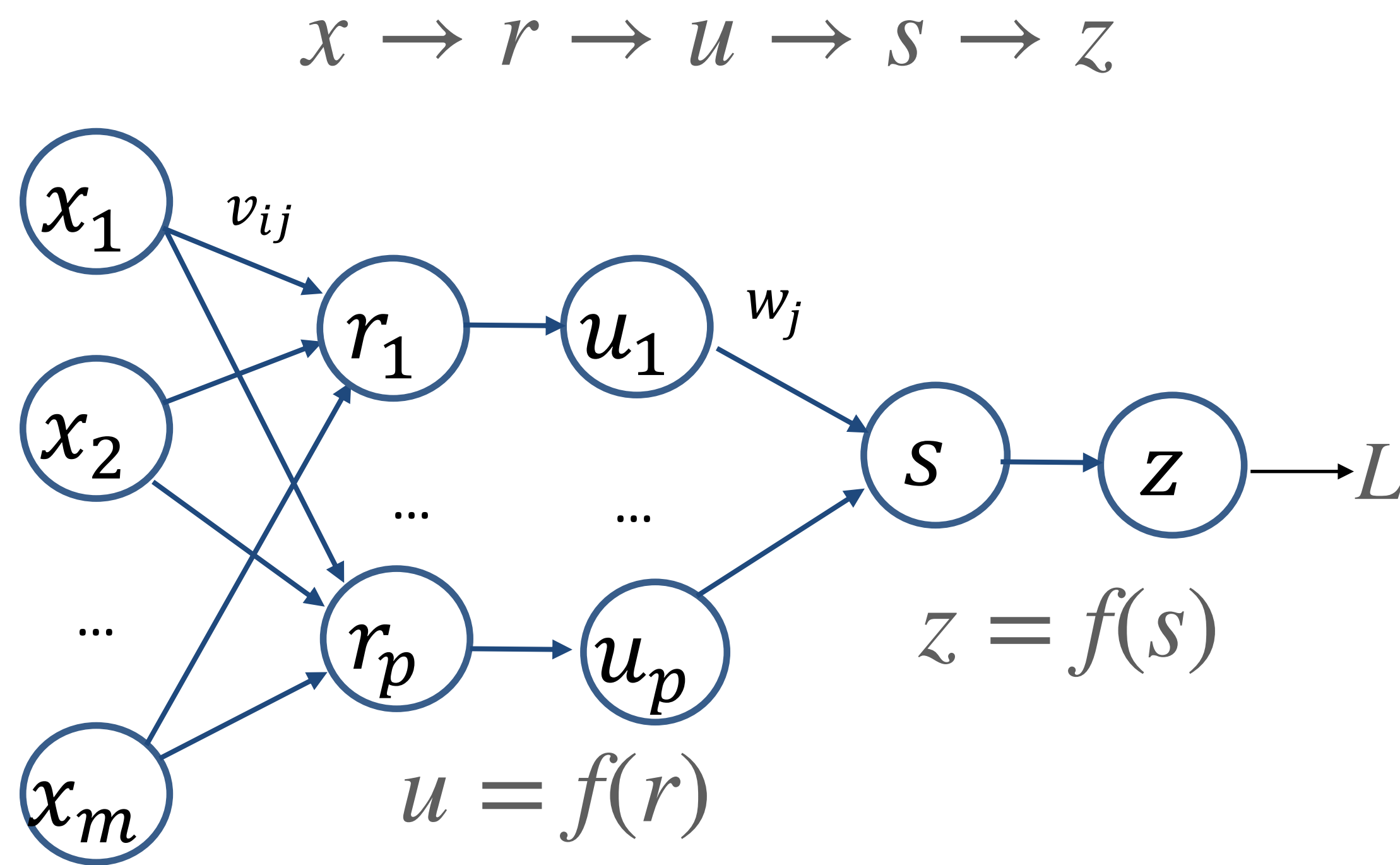
$$\text{Example : } z = \sin(x^2)$$

$$\begin{array}{ccc} z = \sin(u) & \frac{dz}{du} = \cos(u) & \\ \uparrow & \downarrow & \\ u = x^2 & \frac{dz}{dx} = \frac{dz}{du} \frac{du}{dx} = 2x \cos(u) & \end{array}$$

Recap of Multi-layer perceptron: Function composition

Forward prediction

$$\begin{aligned}
 z &= f(s) = \text{sigmoid}(s) = \frac{1}{1 + e^{-s}} \\
 &\quad \uparrow \\
 s &= \sum_{j=0}^p w_j u_j \\
 &\quad \uparrow \\
 u_j &= f(r_j) = \text{sigmoid}(r_j) = \frac{1}{1 + e^{-r_j}} \\
 &\quad \uparrow \\
 r_j &= \sum_{i=0}^m x_i v_{ij}
 \end{aligned}$$



f : activation functions

Recap of Multi-layer perceptron: Chain rule

Backward propagation:

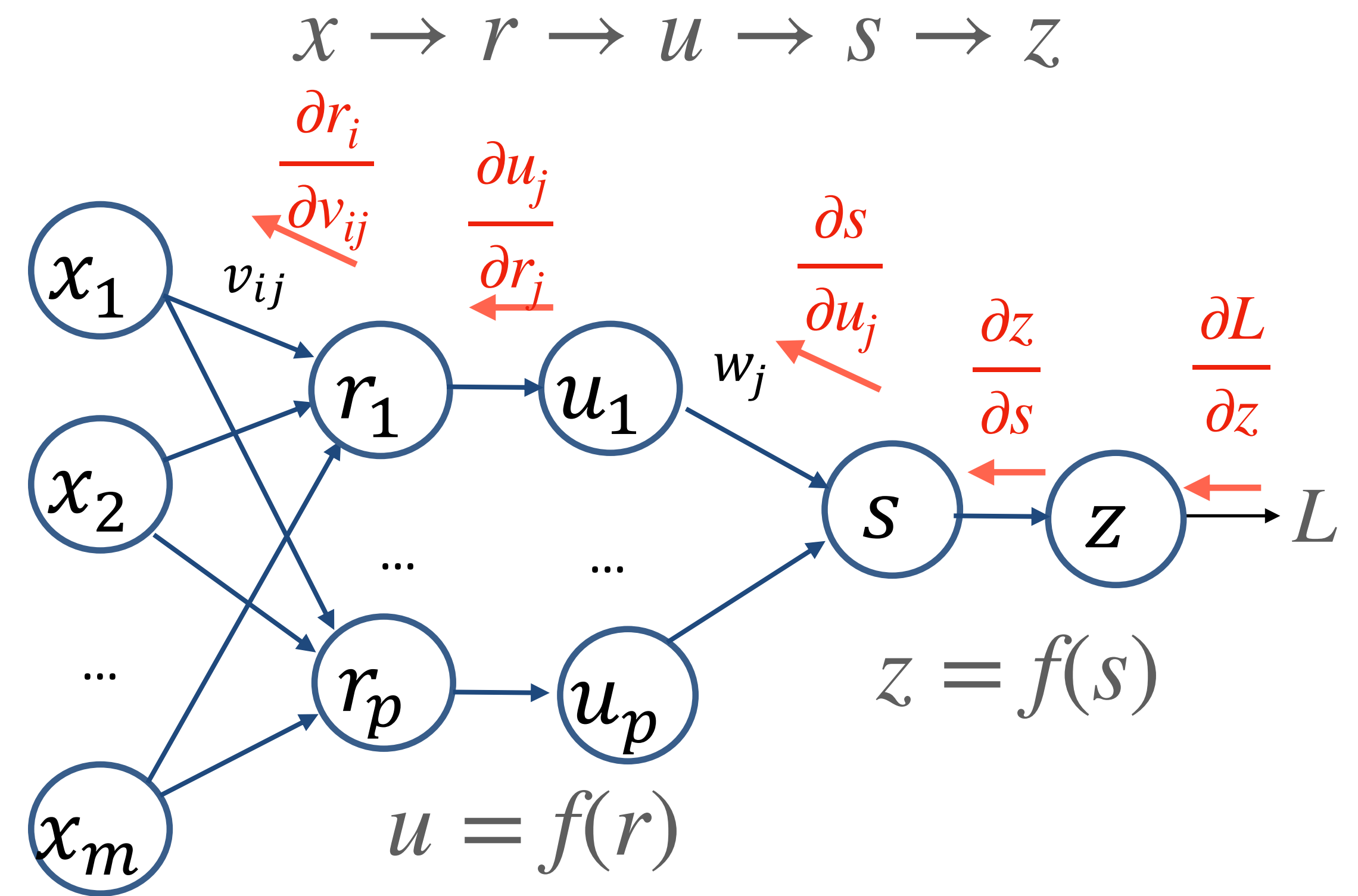
$$\frac{\partial L}{\partial v_{ij}} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial s} \frac{\partial s}{\partial u_j} \frac{\partial u_j}{\partial r_j} \frac{\partial r_j}{\partial v_{ij}}$$

$$s = \sum_{j=0}^p u_j w_j$$

$$u_j = \text{sigmoid}(r_j) = \frac{1}{1 + e^{-r_j}}$$

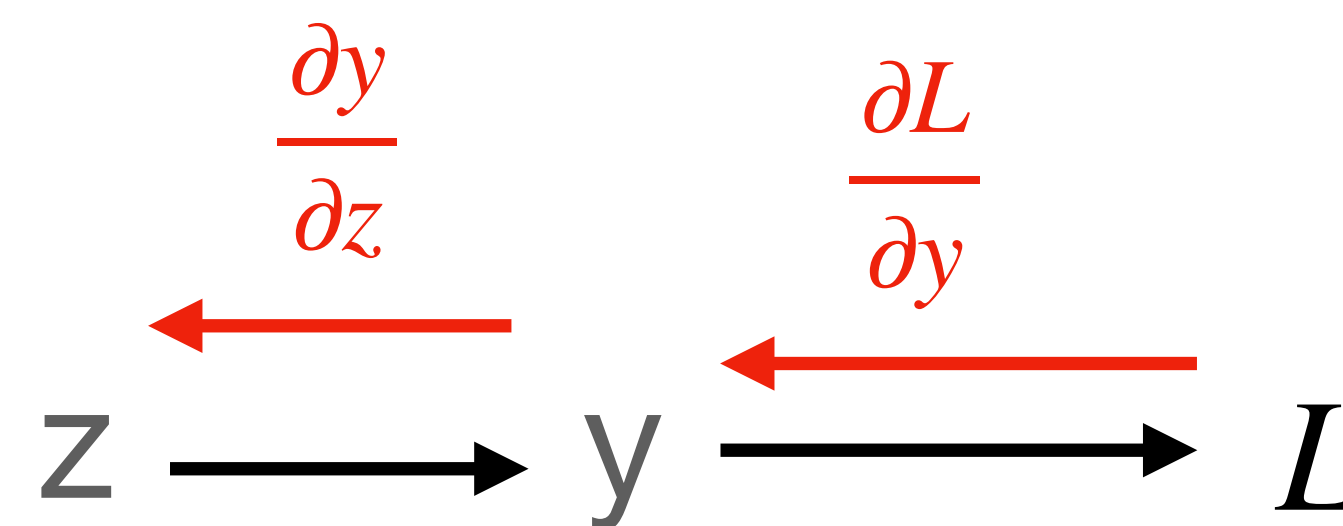
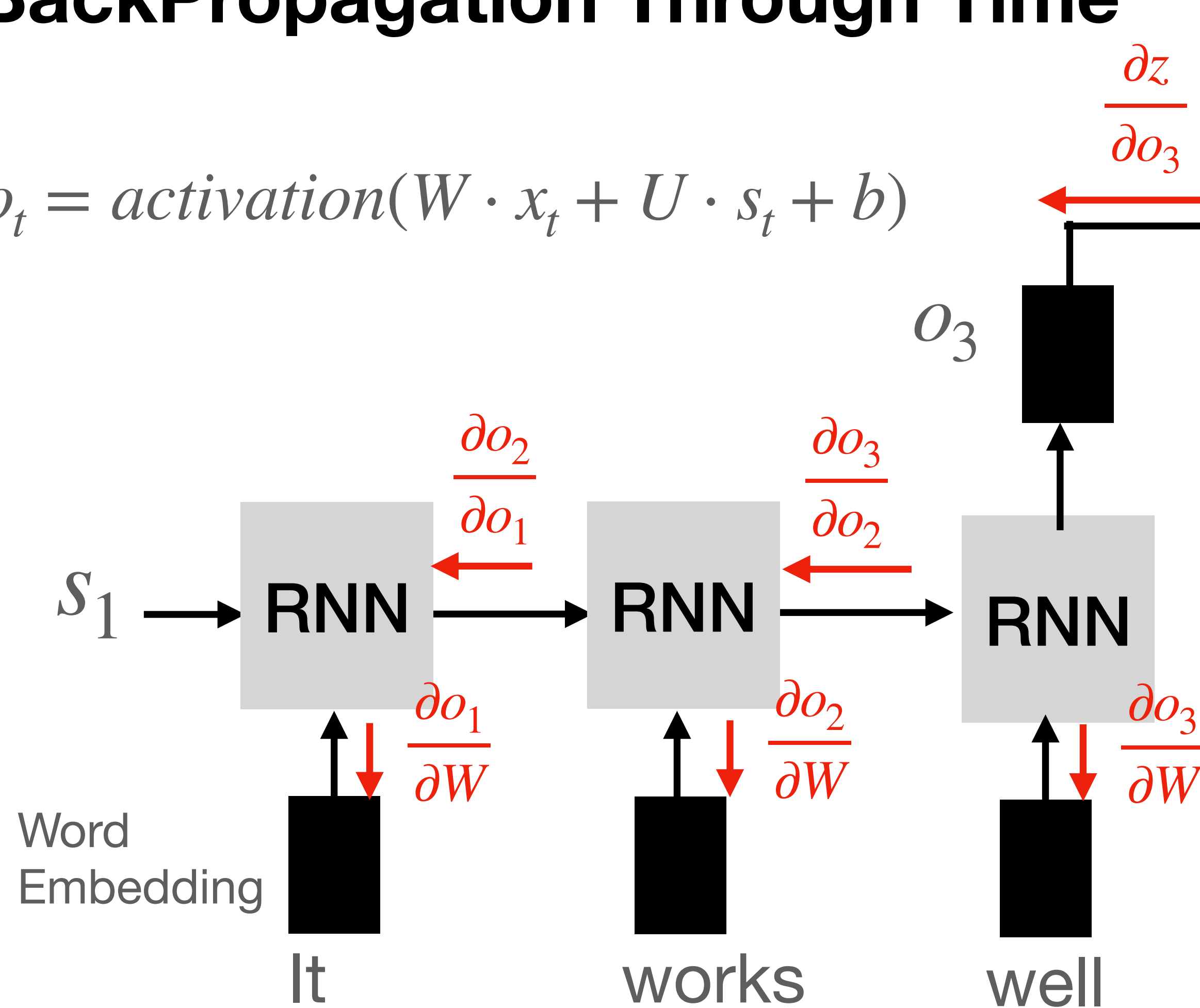
$$r_j = \sum_{i=0}^m x_i v_{ij}$$

Forward



BackPropagation Through Time

$$o_t = \text{activation}(W \cdot x_t + U \cdot s_t + b)$$



$$\begin{aligned} \frac{\partial L}{\partial W} = & \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial o_3} \frac{\partial o_3}{\partial W} + \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial W} \\ & + \frac{\partial L}{\partial y} \frac{\partial y}{\partial z} \frac{\partial z}{\partial o_3} \frac{\partial o_3}{\partial o_2} \frac{\partial o_2}{\partial o_1} \frac{\partial o_1}{\partial W} \end{aligned}$$

A special RNN: Long short-term memory (LSTM)

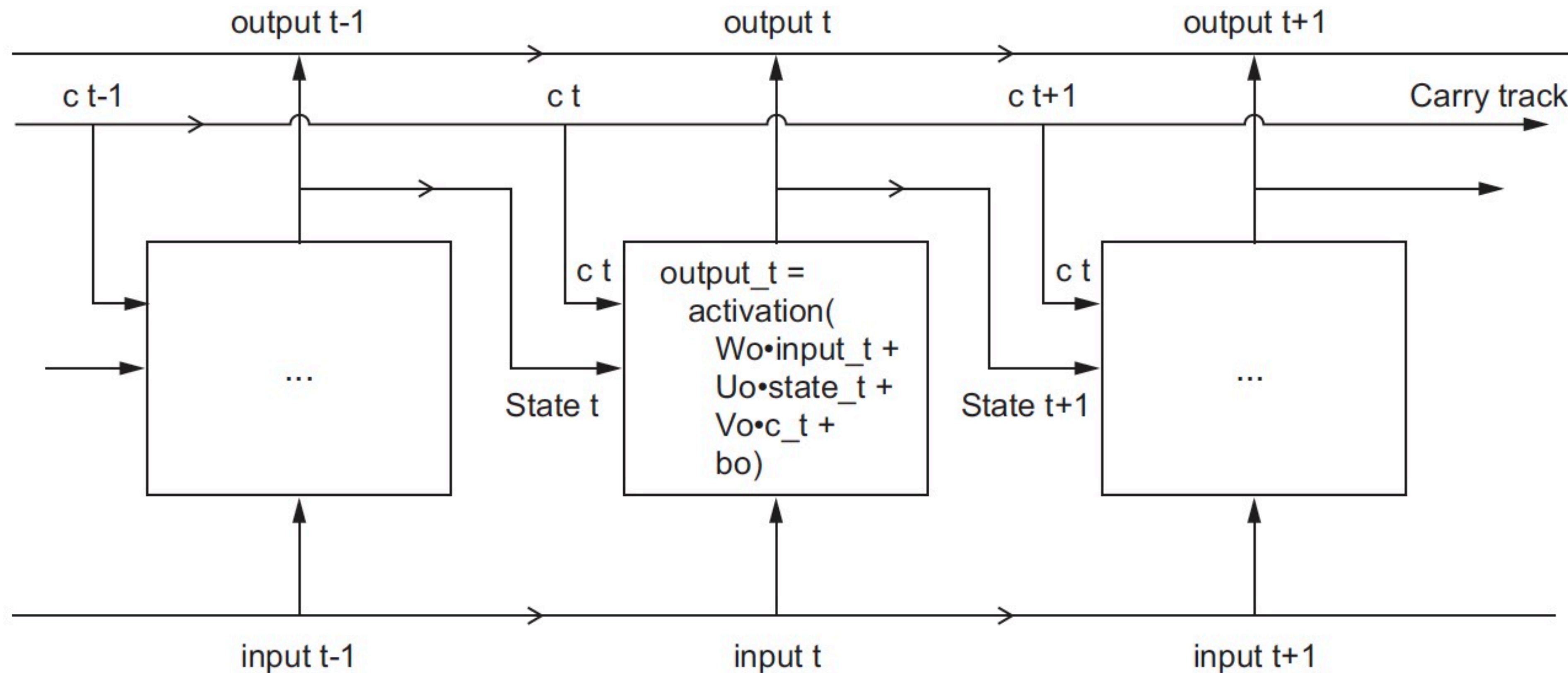


Figure 6.14 in Deep learning with python by Francois Chollet

A special RNN: Long short-term memory (LSTM)

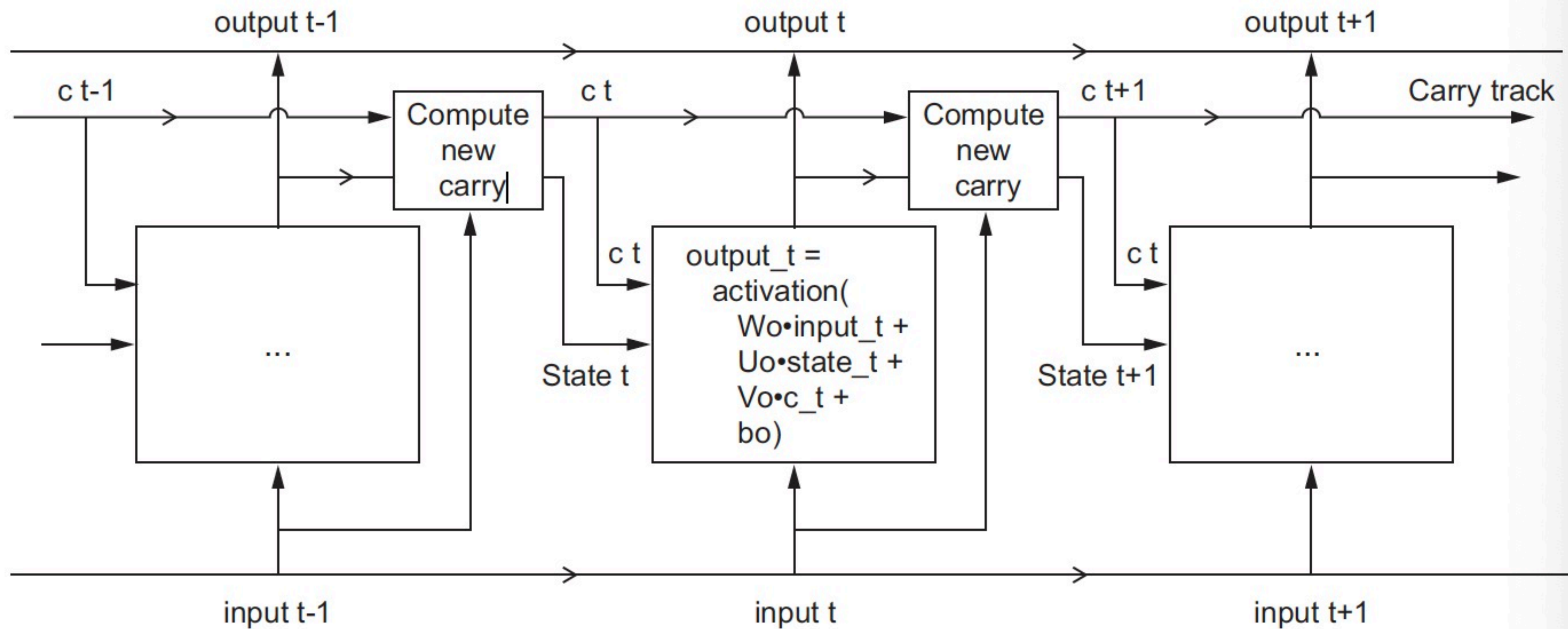
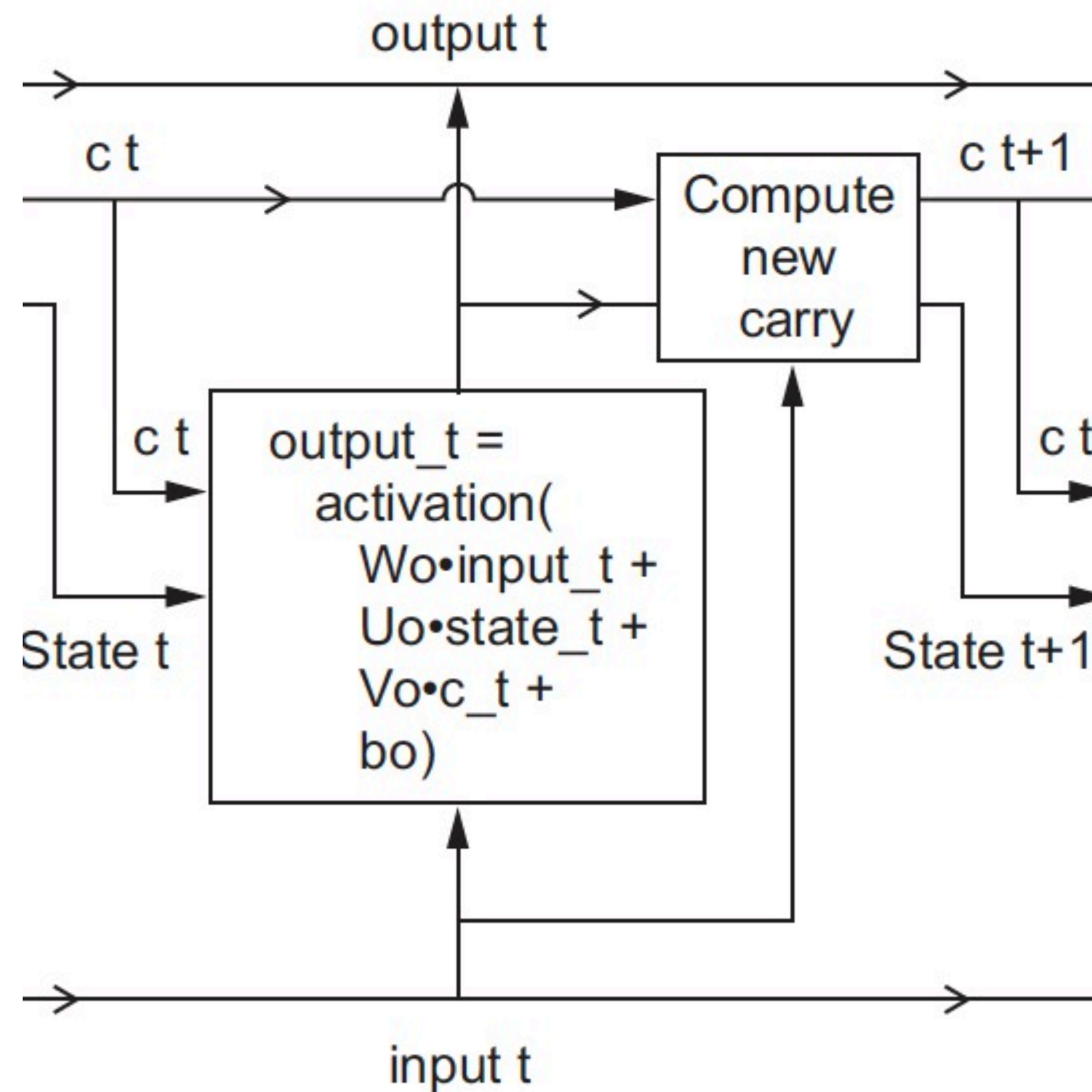


Figure 6.15 in Deep learning with python by Francois Chollet

A special RNN: Long short-term memory (LSTM)



Forget gate: $f_t = \text{sigmoid}(W_f \cdot x_t + U_f \cdot s_t + b_f)$

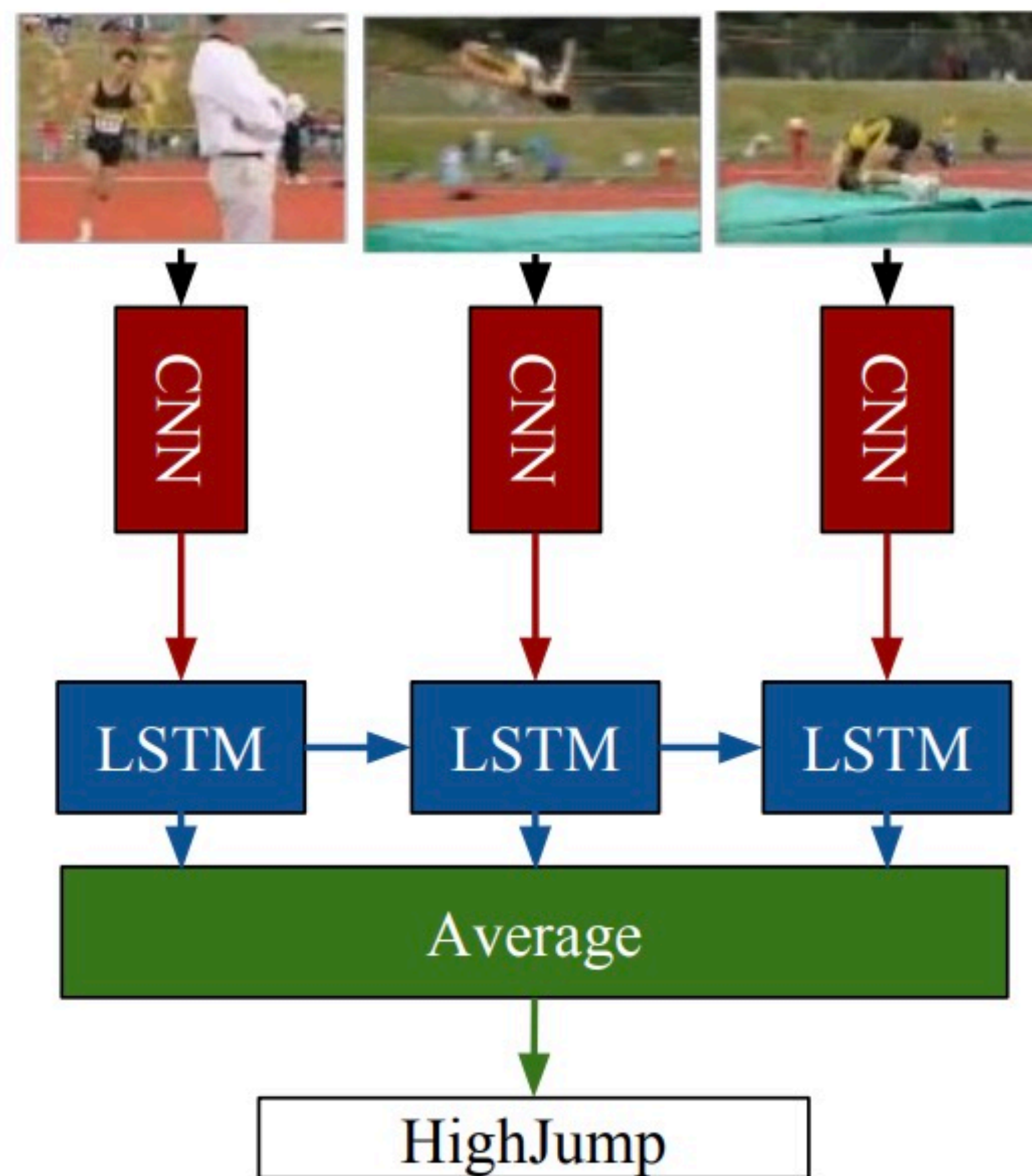
Input gate: $i_t = \text{sigmoid}(W_i \cdot x_t + U_i \cdot s_t + b_i)$

$$\tilde{c}_{t+1} = \tanh(W_c \cdot x_t + U_c \cdot s_t + b_c)$$

$$c_{t+1} = c_t * f_t + \tilde{c}_{t+1} * i_t$$

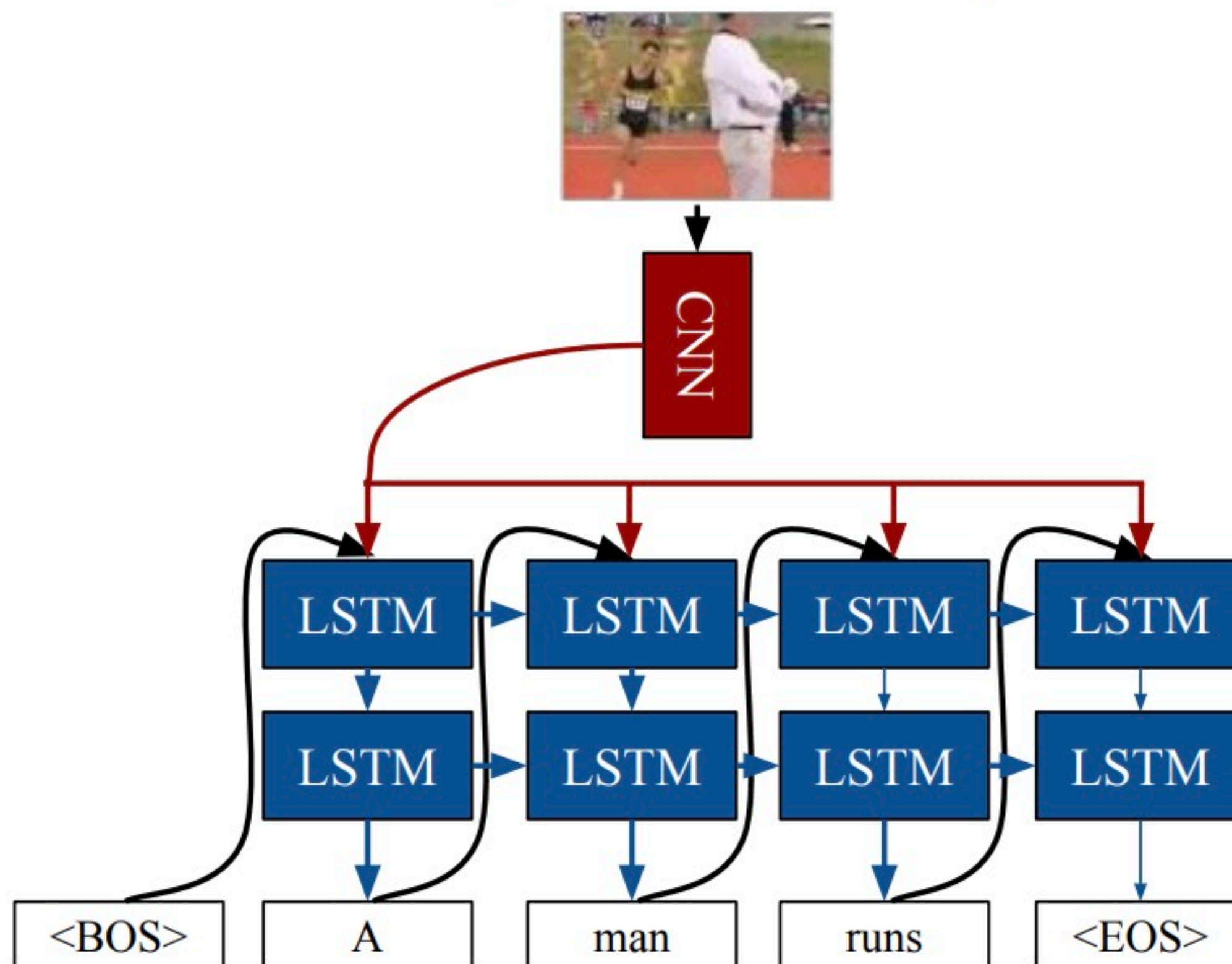
Applications

Activity Recognition Sequences in the Input



Applications

Image Captioning Sequences in the Output



Summary

- How does temporal convolution work?
- How does dilated casual convolution work?
- How does recurrent neural network work?
- How does LSTM work?

Next: Graph Convolution Network (Guest lecture)