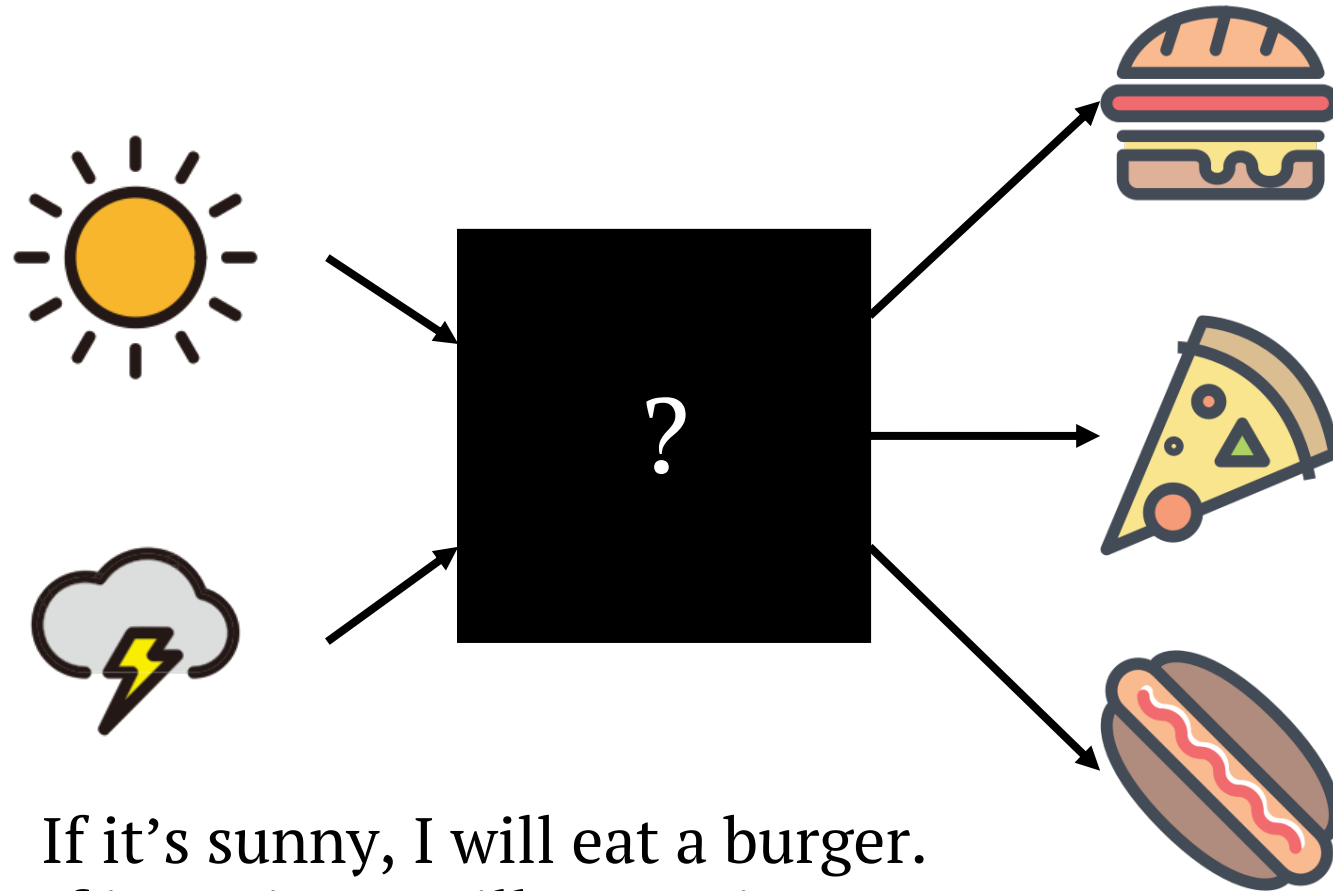


Recurrent Neural Networks

Dr. Unchalisa Taetragool

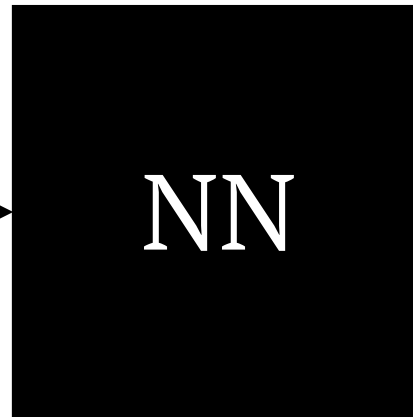
Department of Computer Engineering, Faculty of Engineering
King Mongkut's University of Technology Thonburi

Neural Network: Quiz!



If it's sunny, I will eat a burger.
If it's rainy, I will eat a pizza.

Neural Network: Quiz!



Dinner Schedule

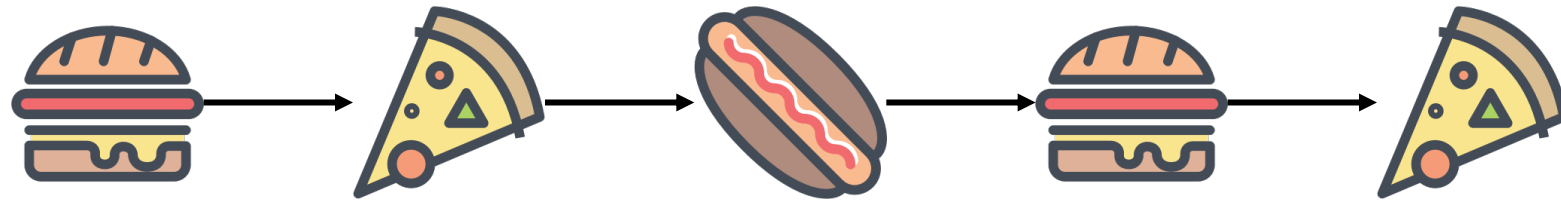
Monday

Tuesday

Wednesday

Thursday

Friday



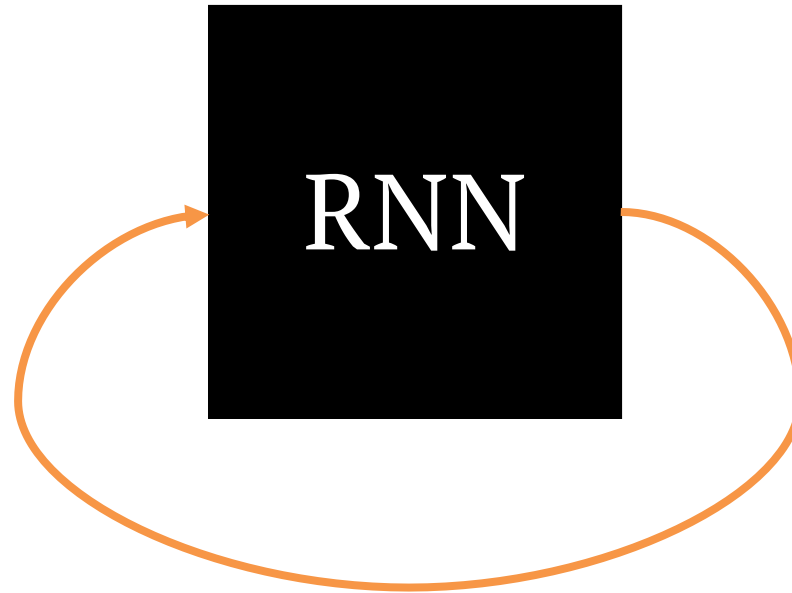
Simple Recurrent Neural Network



RNN



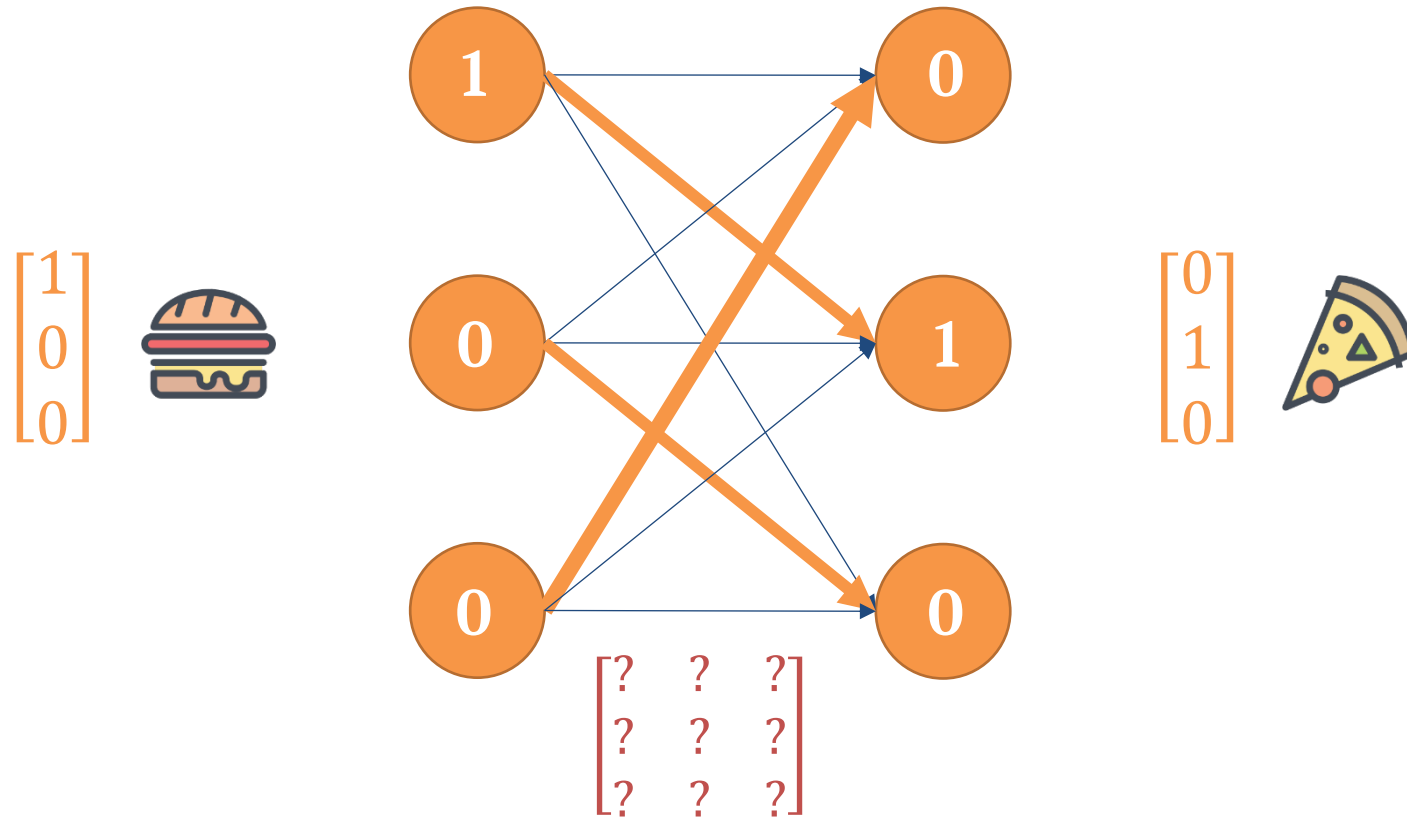
Simple Recurrent Neural Network



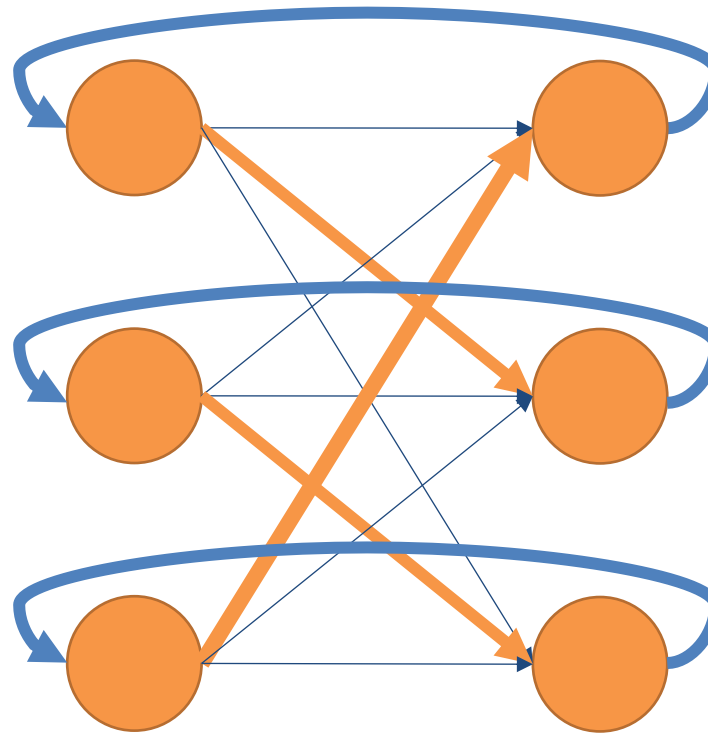
Simple Recurrent Neural Network



Simple Recurrent Neural Network



Simple Recurrent Neural Network



Weather Effect



Sunny
New food



Rain
Leftover

Weather Effect

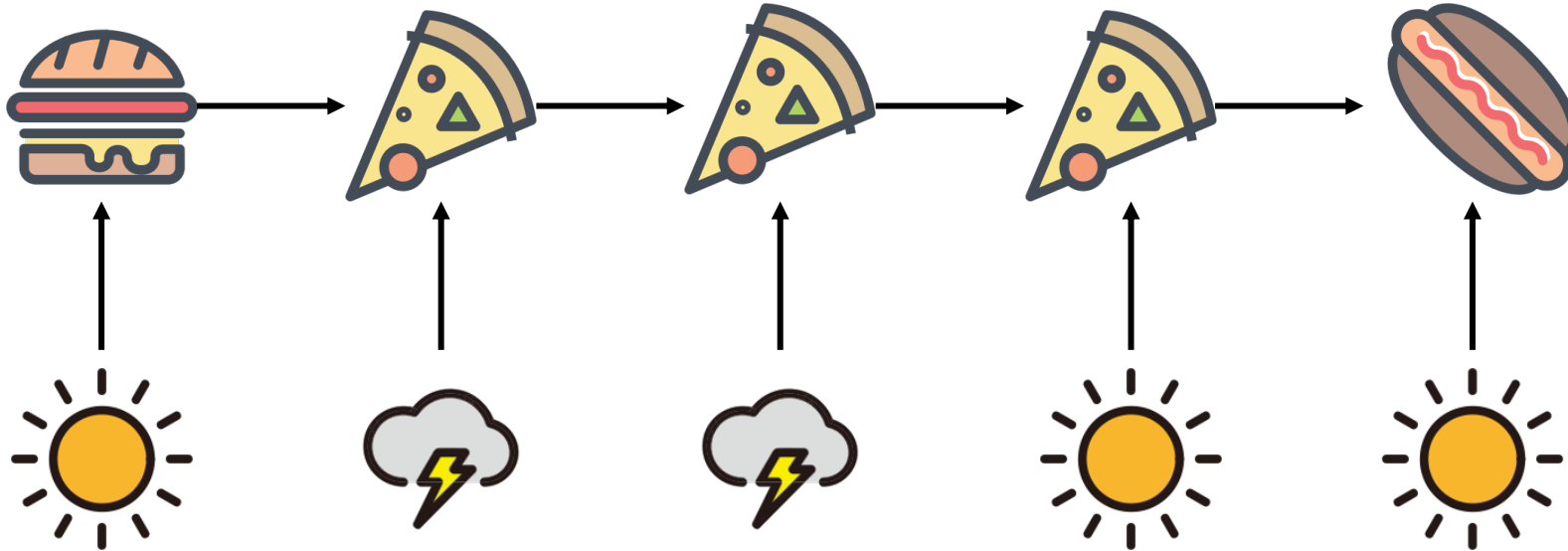
Monday

Tuesday

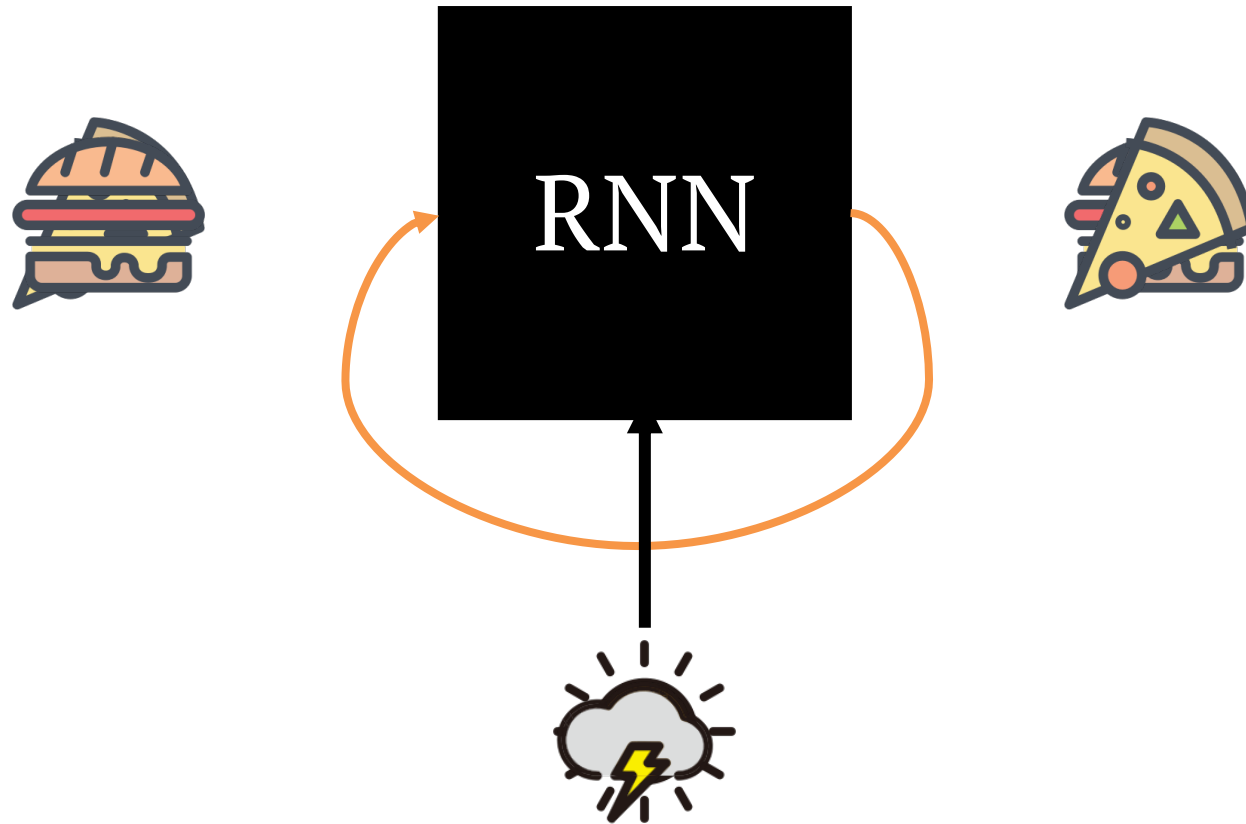
Wednesday

Thursday

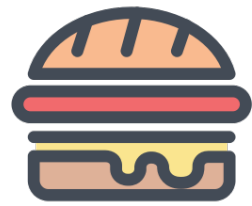
Friday



More Complicated RNN



Vectors



$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$

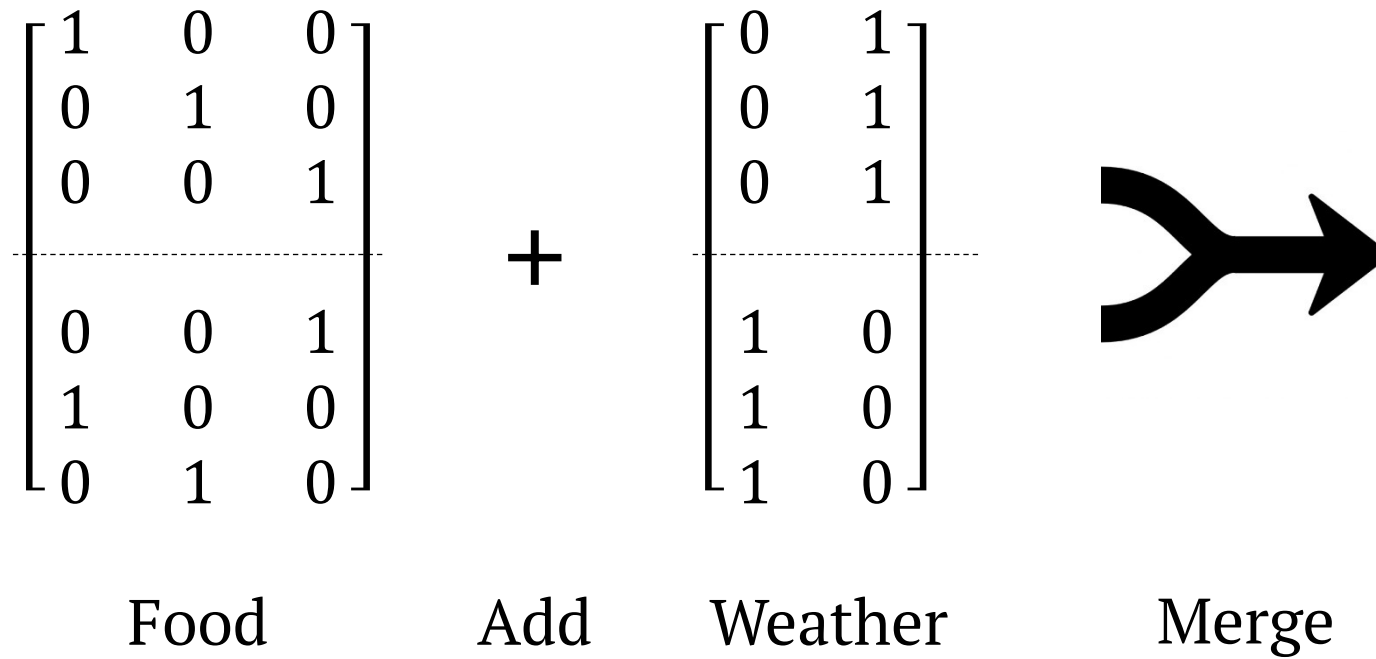
$$\begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$\begin{bmatrix} 0 \\ 1 \end{bmatrix}$$



More Complicated RNN



Food



$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$


Food




=

$$\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

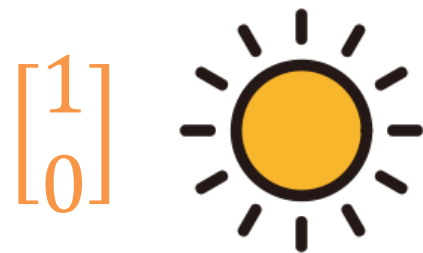
Same



$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$


Next day

Weather



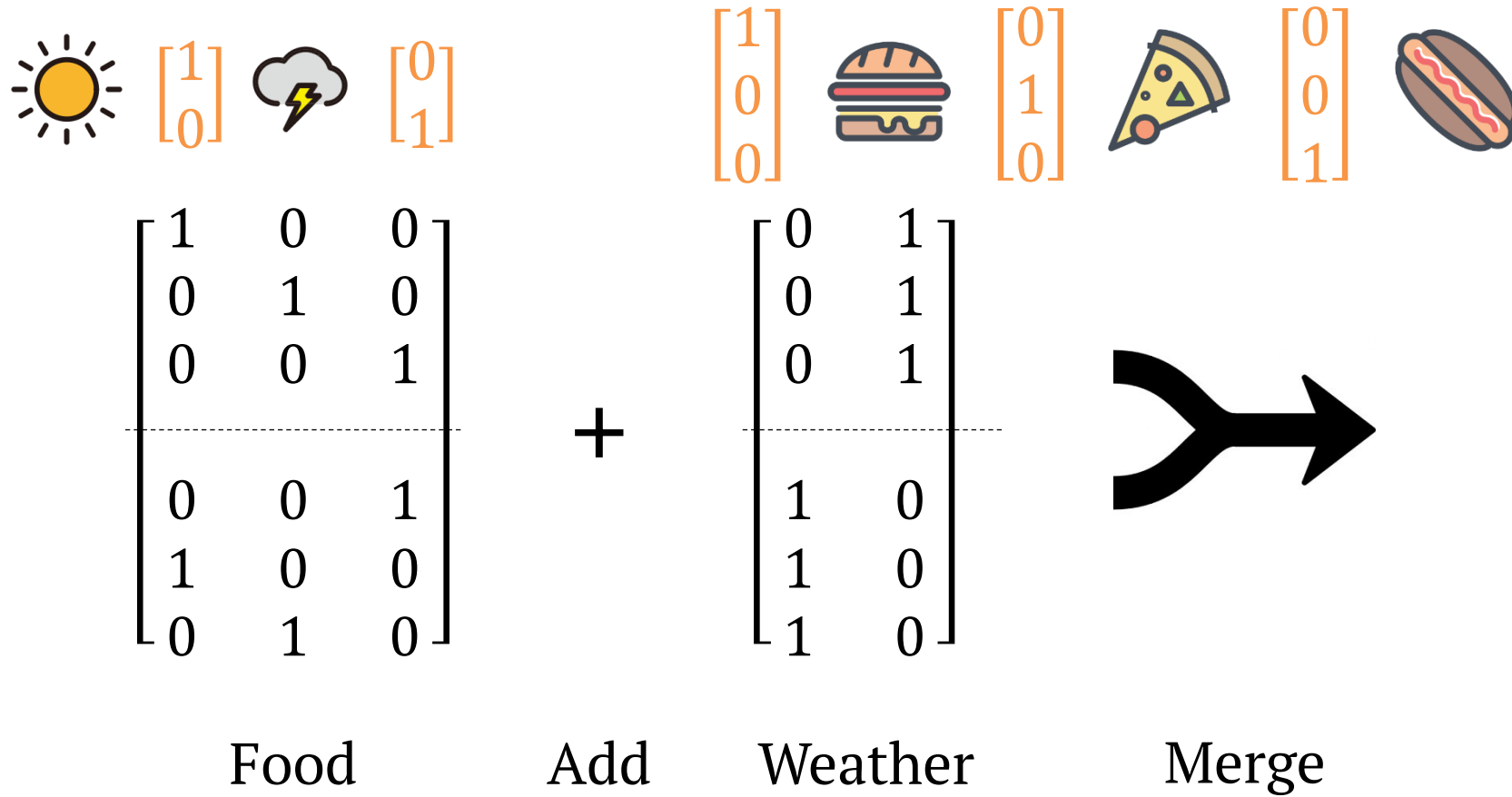
$$\begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ \hline 1 & 0 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}$$

Weather



$$= \begin{bmatrix} 0 \\ 0 \\ 0 \\ \hline 1 \\ 1 \\ 1 \end{bmatrix} \begin{matrix} \text{Leftover} \\ \\ \\ \text{New food} \end{matrix}$$

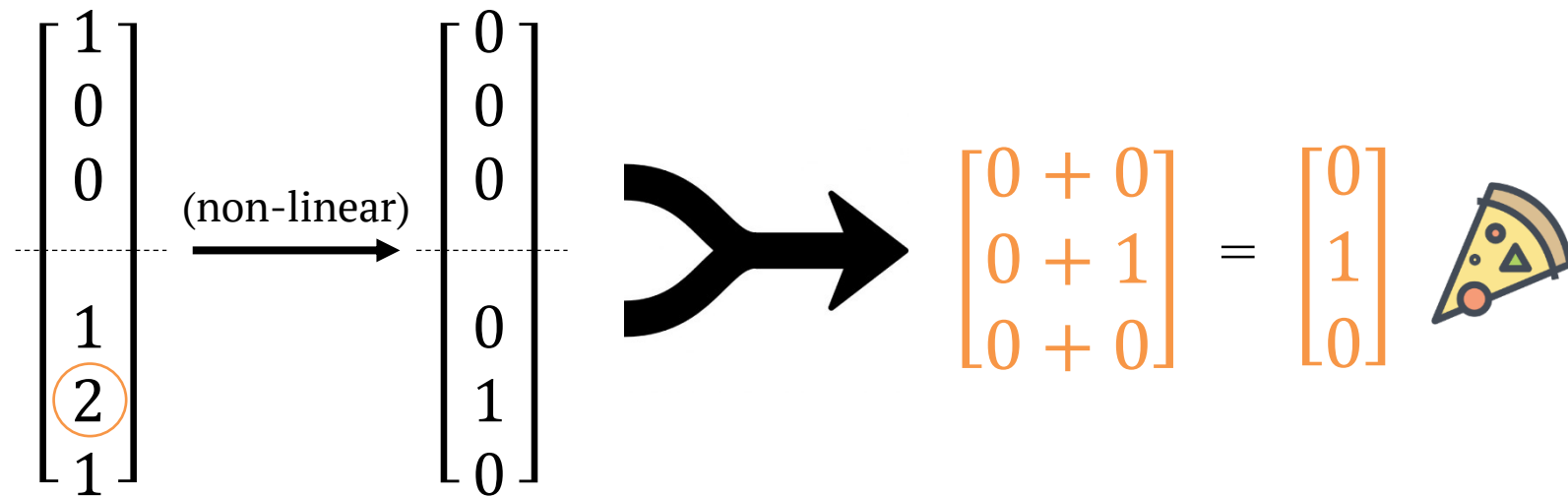
More Complicated RNN



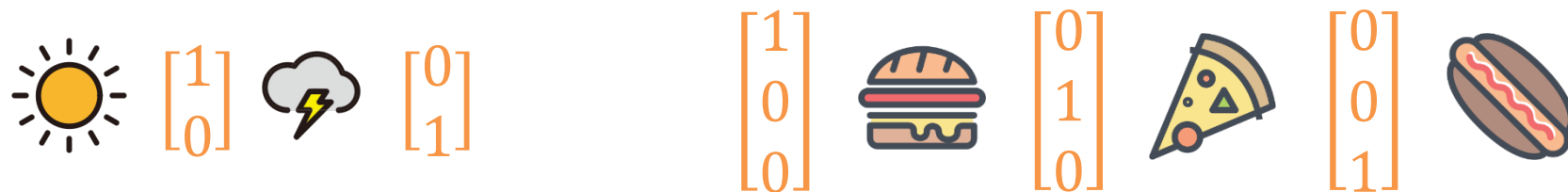
Add  

$$\begin{array}{c}
 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \begin{array}{l} \text{Same} \\ \img alt="burger icon" data-bbox="315 405 375 485" \end{array} \\
 \hline
 \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \begin{array}{l} \img alt="pizza slice icon" data-bbox="315 585 375 685" \end{array} \\
 \text{Next day}
 \end{array}
 +
 \begin{array}{c}
 \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \begin{array}{l} \text{Leftover} \end{array} \\
 \hline
 \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \begin{array}{l} \text{New food} \end{array}
 \end{array}
 =
 \begin{array}{c}
 \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\
 \hline
 \begin{bmatrix} 1 \\ \textcircled{2} \\ 1 \end{bmatrix}
 \end{array}$$

Merge



More Complicated RNN



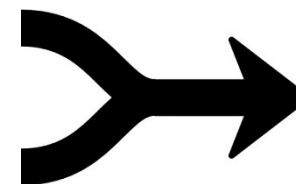
$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ \hline 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}$$

Food

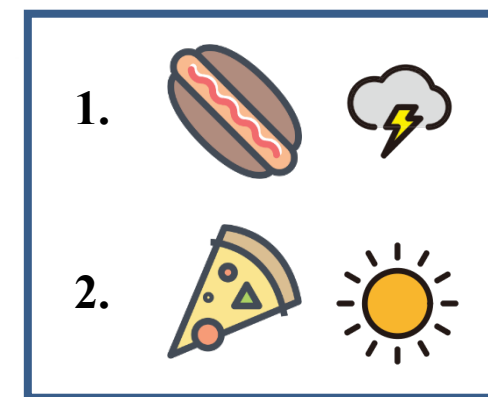
+

$$\begin{bmatrix} 0 & 1 \\ 0 & 1 \\ 0 & 1 \\ \hline 1 & 0 \\ 1 & 0 \\ 1 & 0 \end{bmatrix}$$

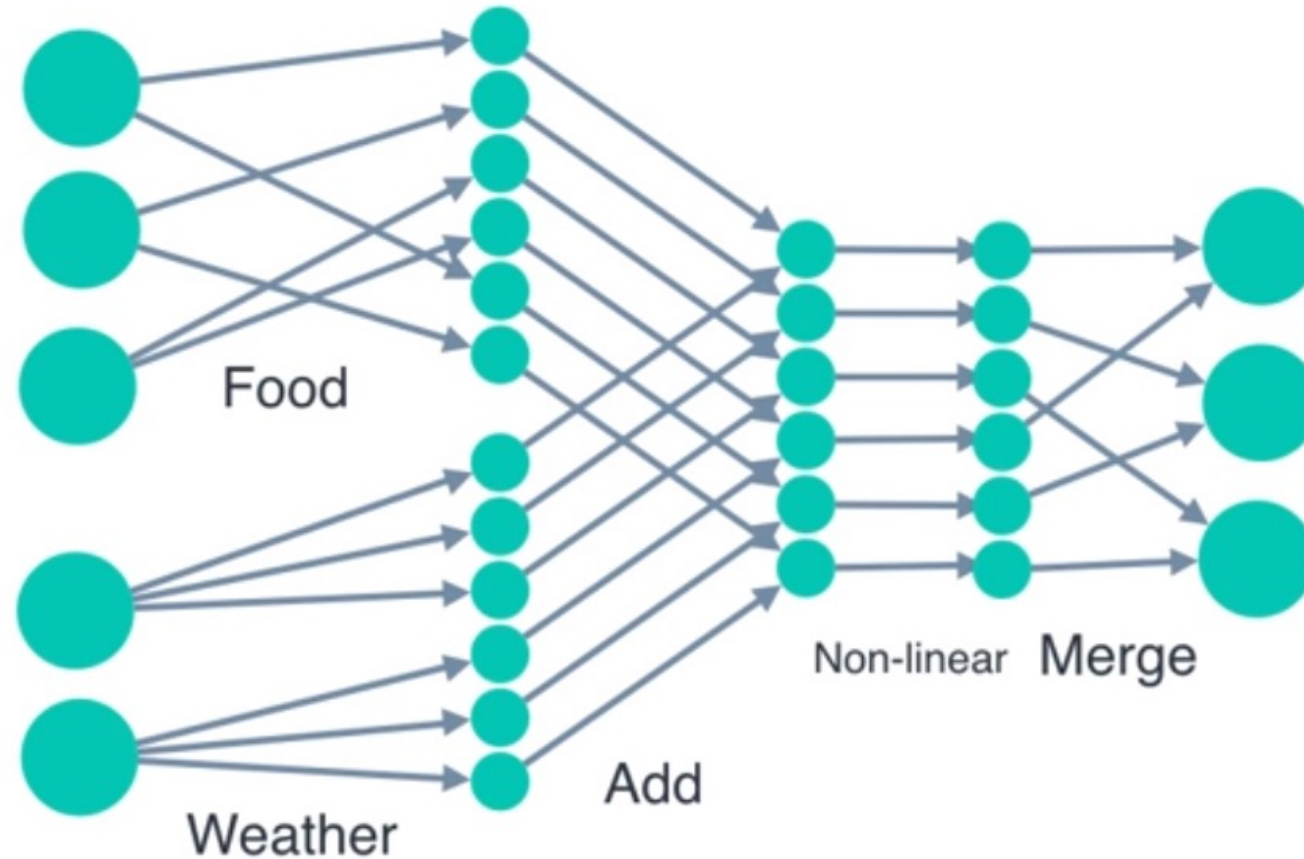
Weather



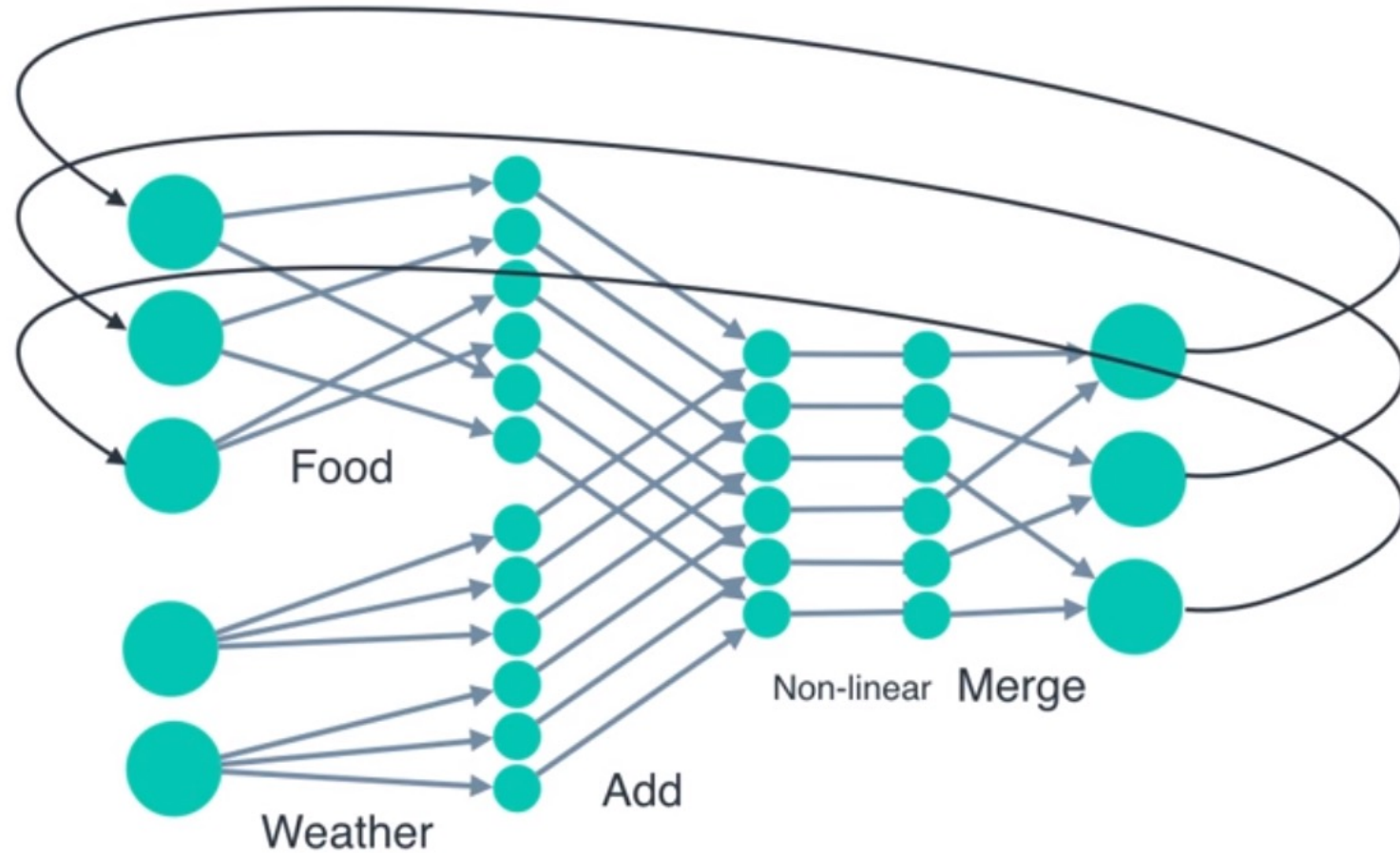
Merge



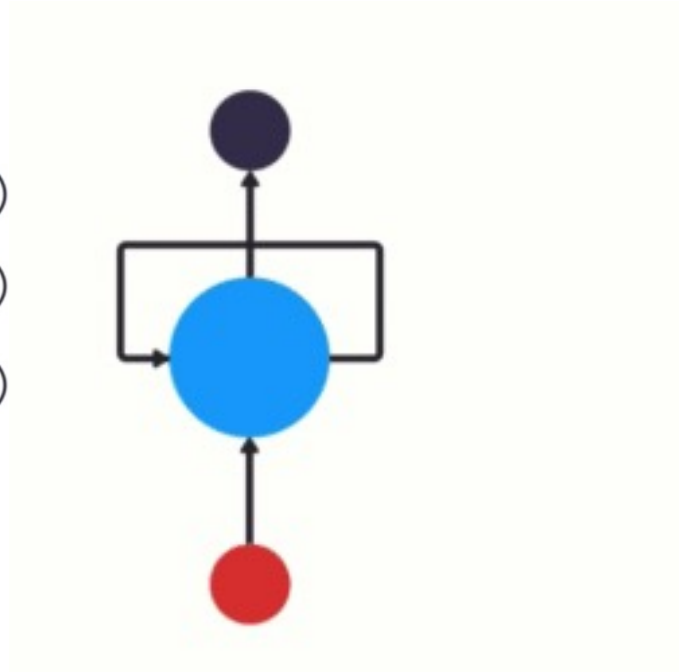
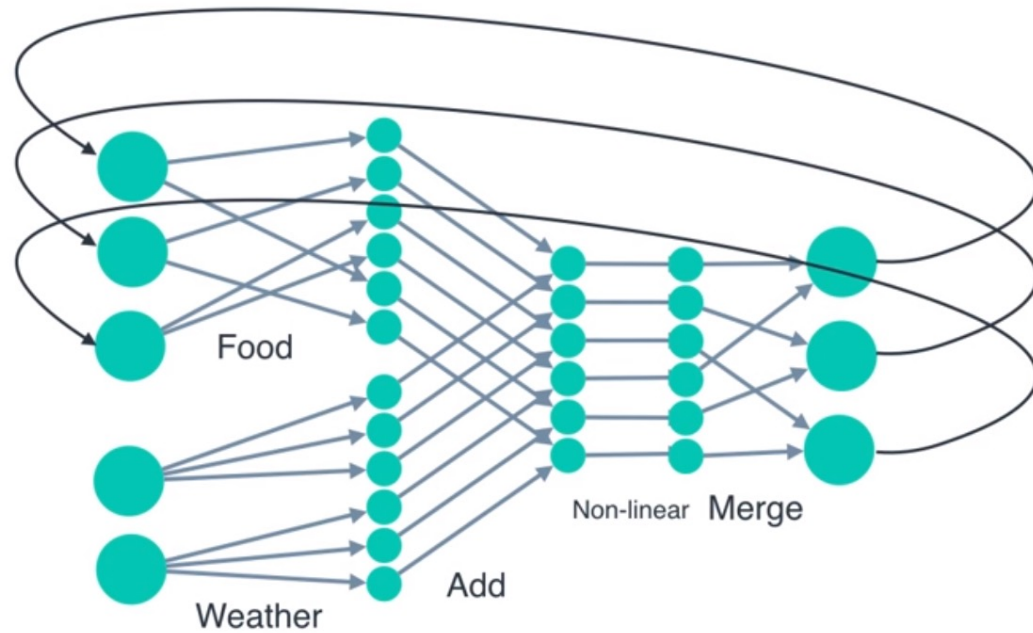
Recurrent Neural Network



Recurrent Neural Network

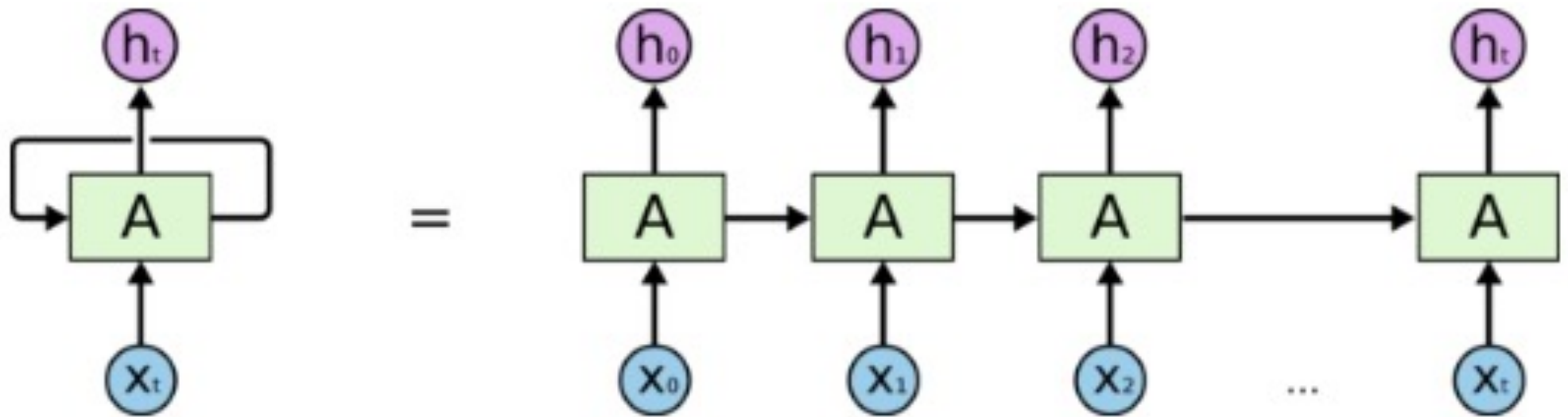


Recurrent Neural Network



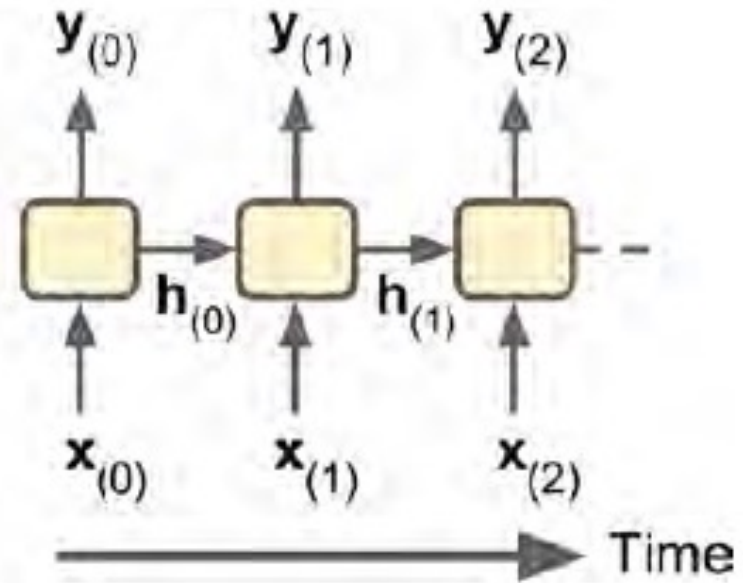
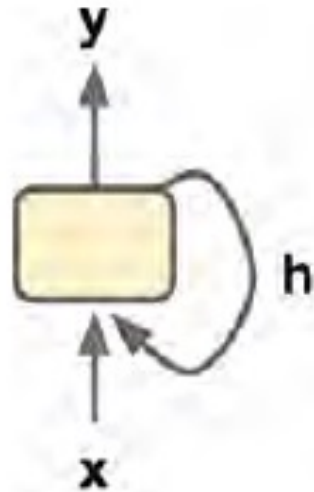
Credit: <https://towardsdatascience.com>

Recurrent Neural Network



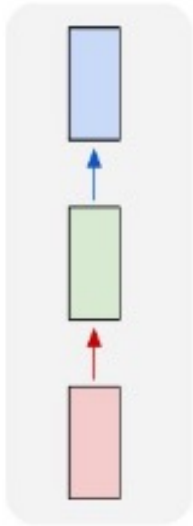
An unrolled recurrent neural network.

Memory Cells

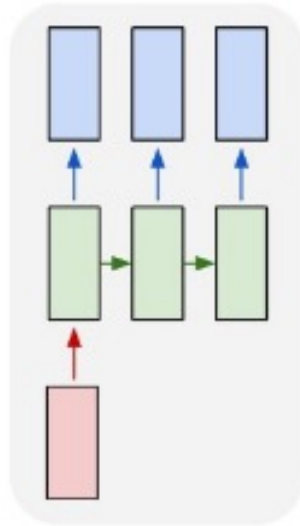


Recurrent Neural Network

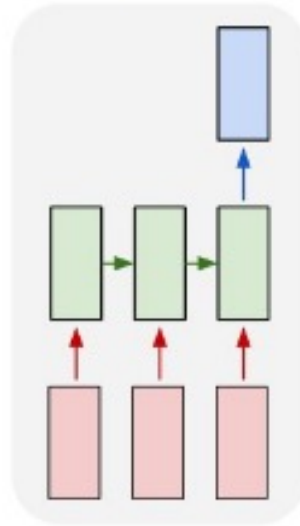
one to one



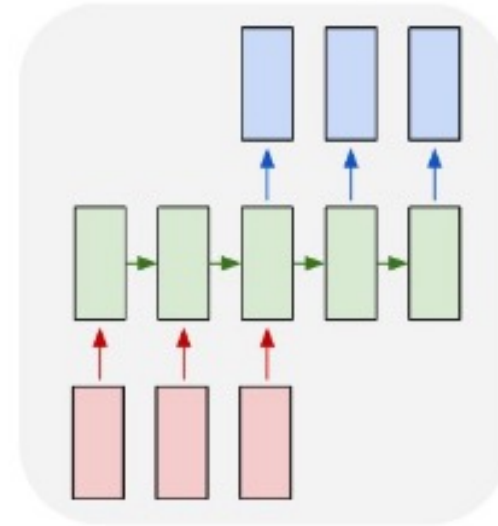
one to many



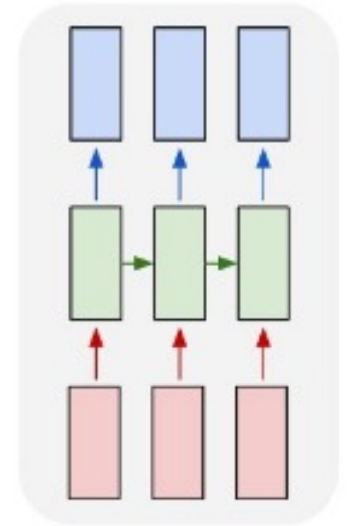
many to one



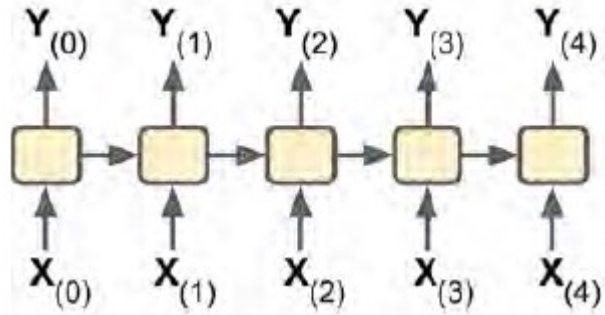
many to many



many to many



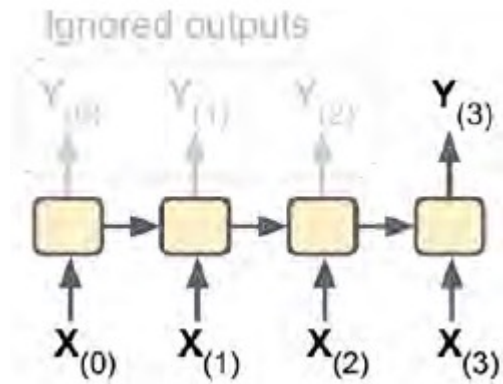
Input and Output Sequences



Sequence to Sequence

I/P – the prices over the last N days

O/P – the prices shifted by one day into the future

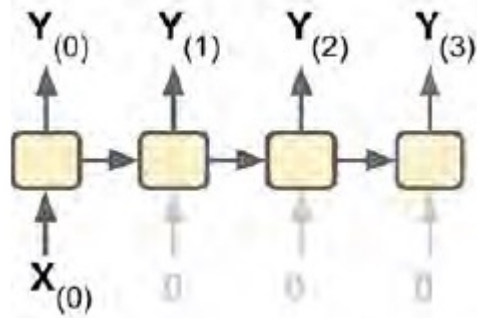


Sequence to Vector

I/P – the sequence of words corresponding to a review

O/P – a sentiment score

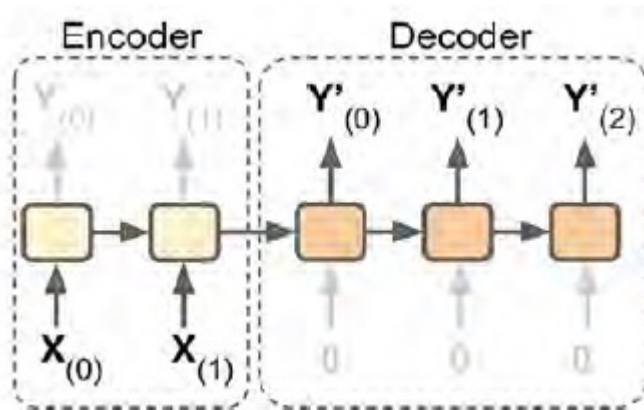
Input and Output Sequences



Vector to Sequence

I/P – an image

O/P – a caption for that image



Delayed Sequence to Sequence

e.g. translating a sentence from one language to another

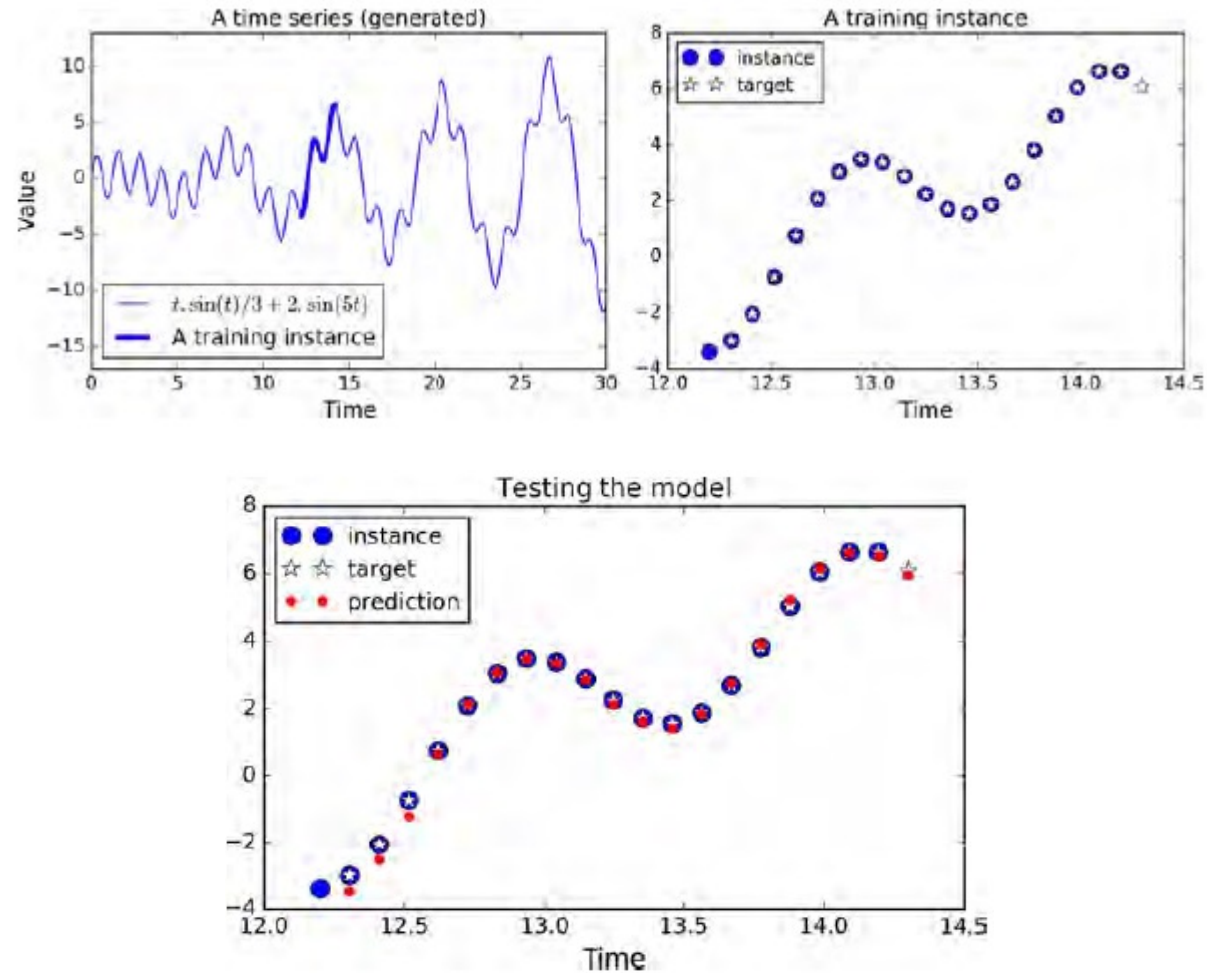
Two-step model called an “Encoder-Decoder”

- Encoder – a sequence-to-vector network
 - Decoder – a vector-to-sequence network
- works much better than trying to translate on the fly with a single sequence-to-sequence RNN

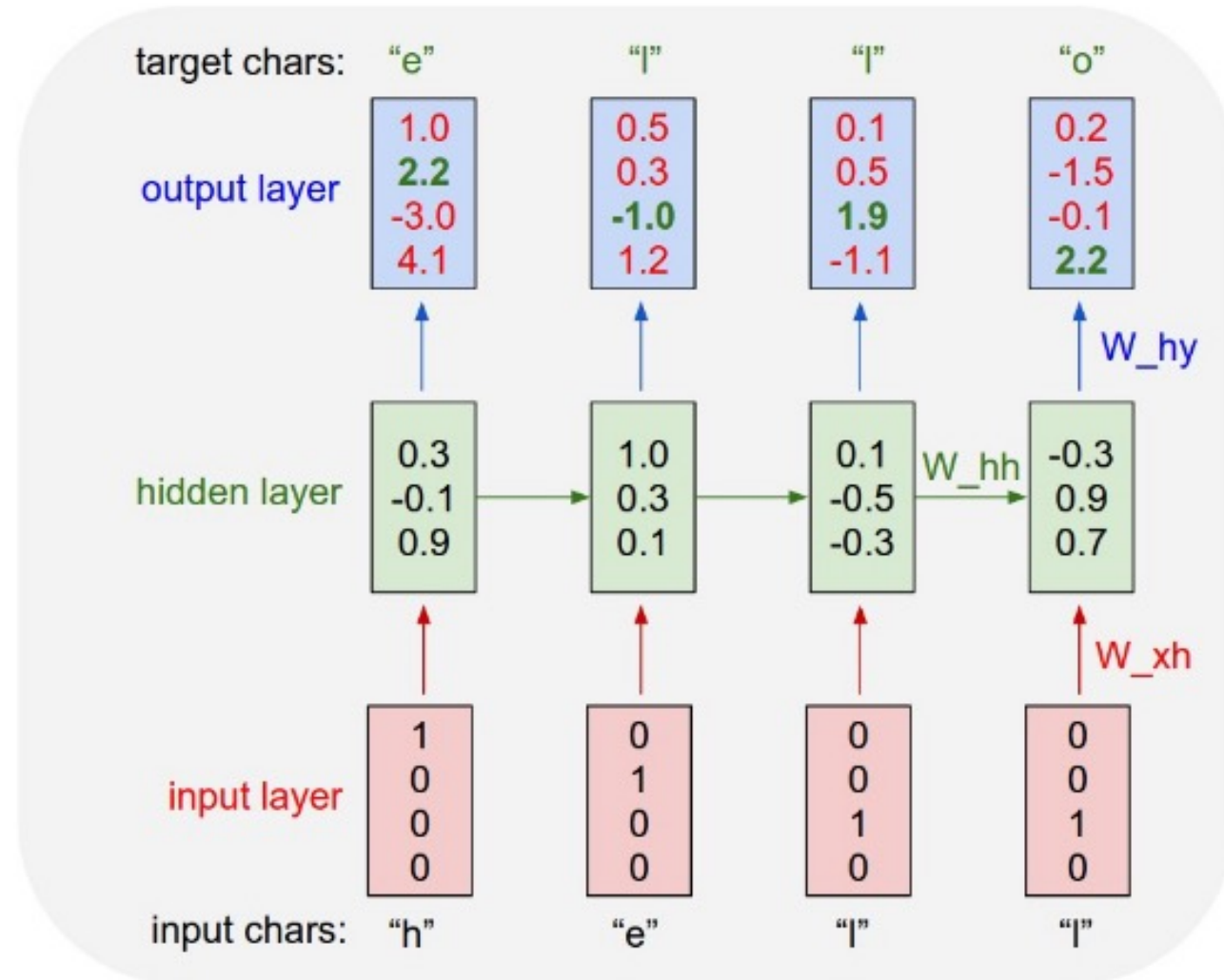
RNN Applications

	x - input		y - output
Speech recognition		⇒	"Fuzzy Wuzzy was a bear. Fuzzy Wuzzy had no hair."
Music generation	\emptyset	⇒	
Sentiment classification	"Decent effort. The plot could have been better."	⇒	
DNA sequence analysis	ACTGTACCCATGTGACTGCCC	⇒	ACTGTACCCATGTGACTGCCC
Machine translation	"El que no arriesga, no gana."	⇒	"If you don't take risks, you cannot win."
Video activity recognition		⇒	Running
Name entity recognition	"Ygritte says Jon Snow knows nothing."	⇒	"Ygritte says Jon Snow knows nothing."

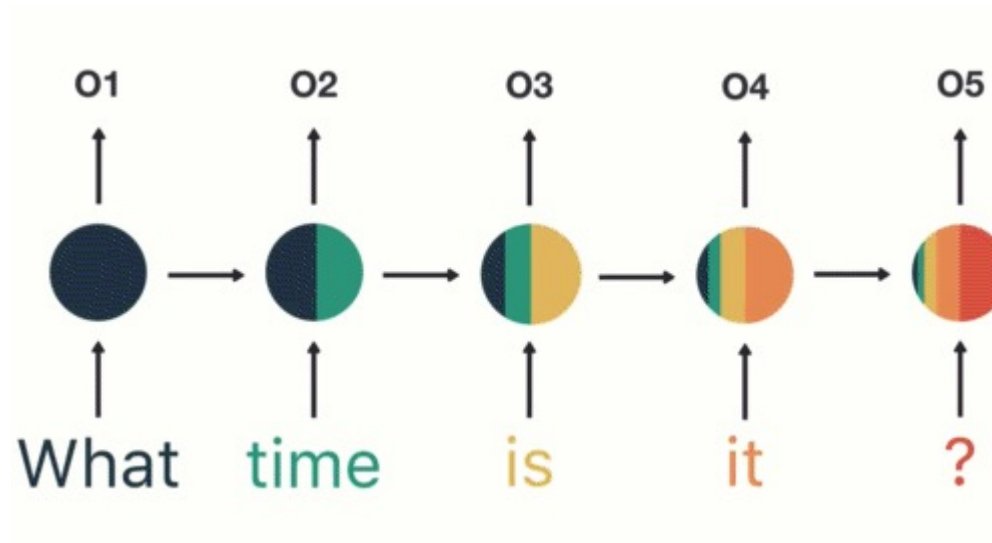
Example: Predict Time Series



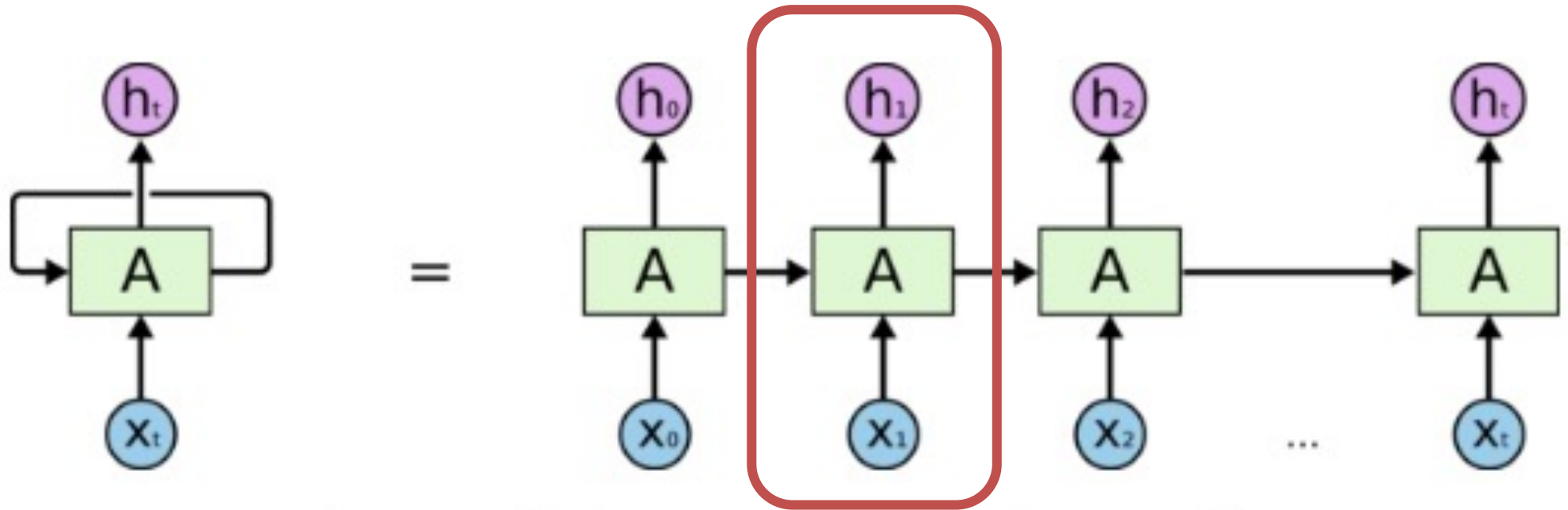
Example: Character-Level Language Models



Example: Encoder

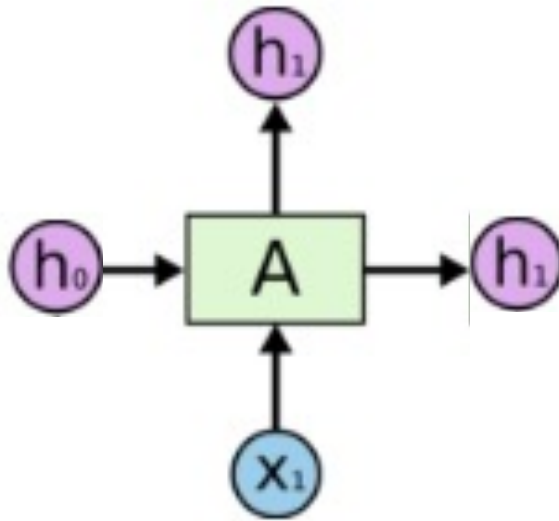


How to Train RNN



An unrolled recurrent neural network.

Training RNN



- Feed forward:
 - A simple one layer of input

$$h_t = u_t X_t + w_t h_{t-1}$$

- Back propagation: **Gradient descent!**

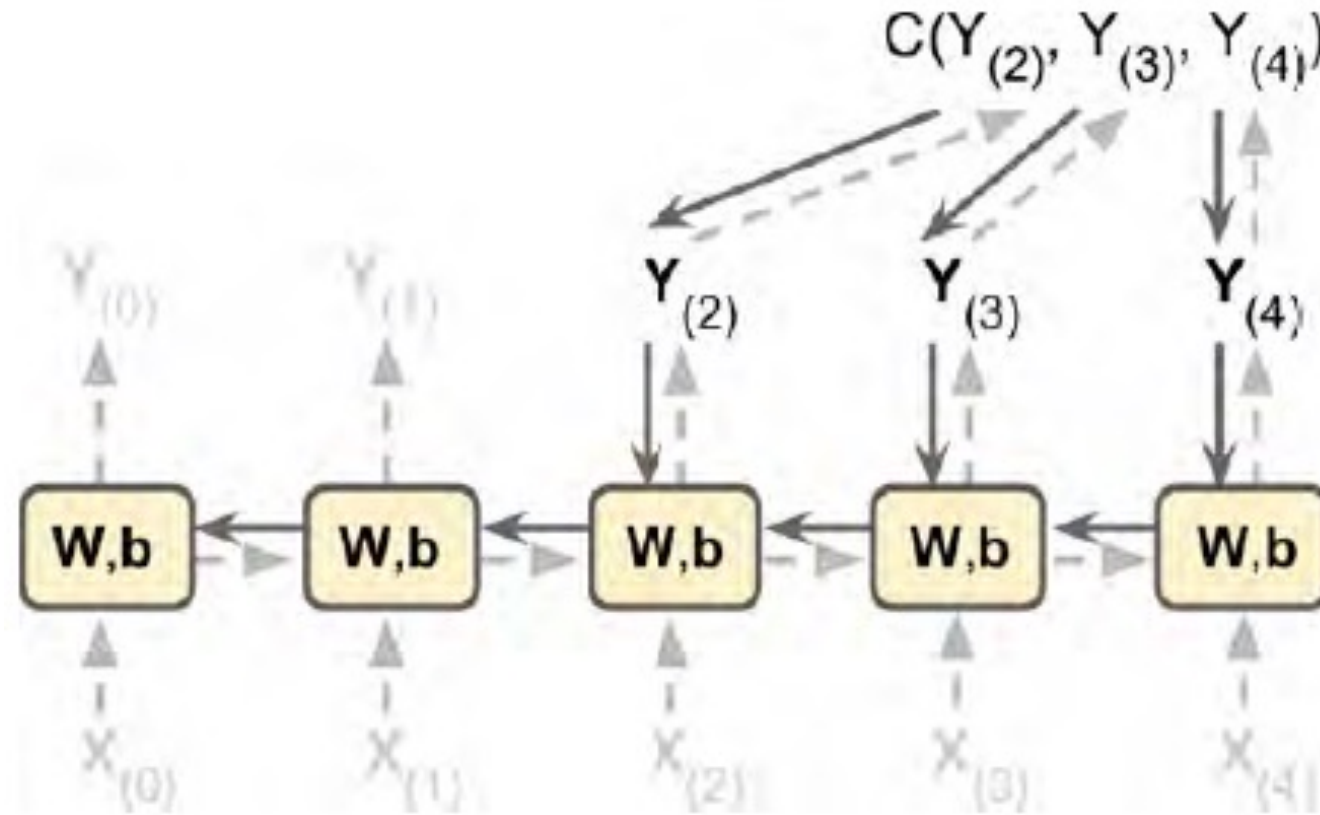
$$E = \frac{1}{N} \sum_{i=1}^n (y_t^i - h_t^i)^2 \text{ then find } \frac{\partial E}{\partial u_t} \text{ and } \frac{\partial E}{\partial w_t}$$

$$u_t^{new} = u_t^{old} - \frac{\partial E}{\partial u_t} \text{ and } w_t^{new} = w_t^{old} - \frac{\partial E}{\partial w_t}$$

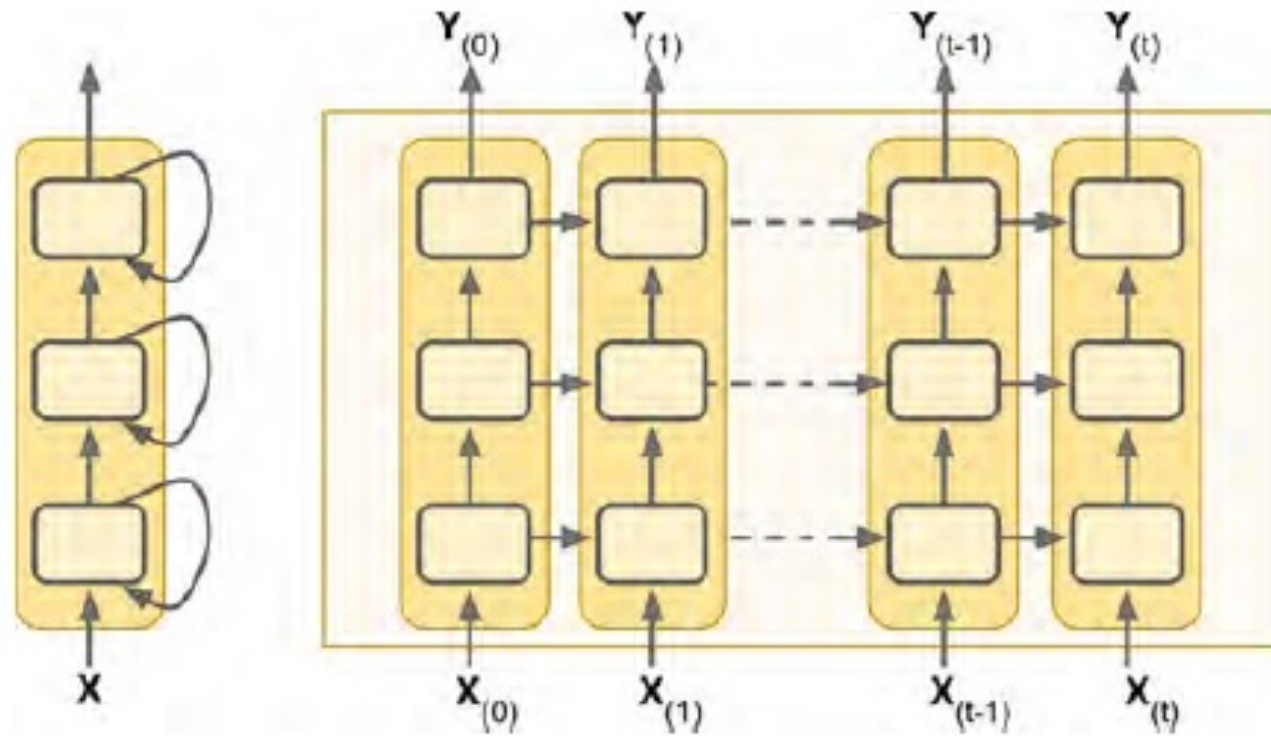
Training RNN: Backpropagation Through Time (BPTT)

- Just like in regular backpropagation,
 - there is a first forward pass through the network
 - then the output sequence is evaluated using a cost function
 - the gradients of that cost function are propagated backward through the network
 - finally, the model parameters are updated using the gradients computed during BPTT

Training RNN: Backpropagation Through Time (BPTT)



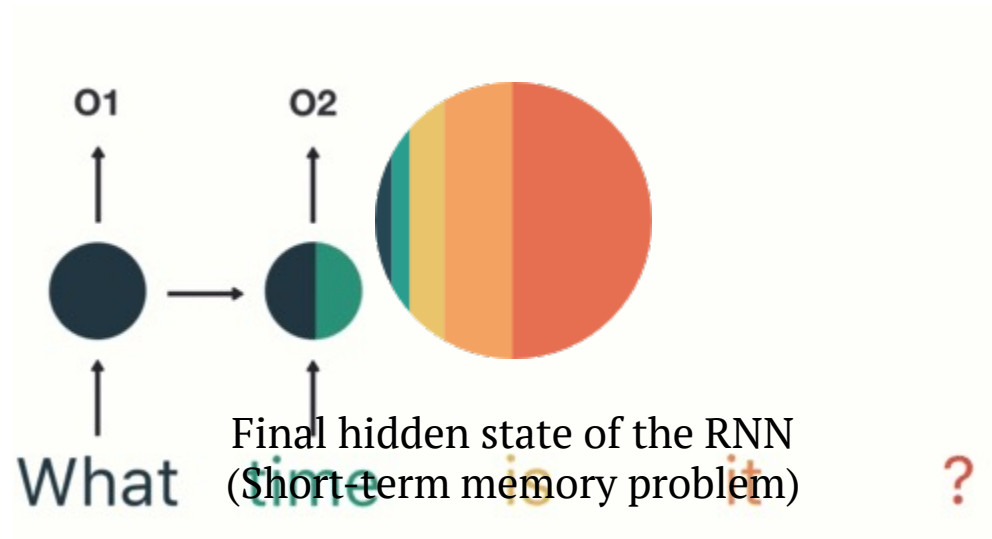
Deep RNN



Difficulties of Training RNN

- To train an RNN on long sequences, you will need to run it over many time steps, making the RNN a very deep network.
- Just like deep neural network, it may suffer from
 - taking forever to train
 - vanishing gradient problem
- To alleviate these problems, you can use
 - good parameter initialization
 - non-saturating activation functions (e.g., ReLU)
 - faster optimizers
 - etc.

Vanishing Gradient



Vanishing Gradient

- The *vanishing gradient* is due to the nature of back-propagation
- Weight adjustment from the gradient descent algorithm

$$w_{ij} := w_{ij} - \alpha * \frac{\partial Cost}{\partial w_{ij}}$$

↑ ↑ ↑ ↑
new weight old weight learning gradient
rate of cost function

- The bigger the gradient, the bigger the adjustments and vice versa.

Vanishing Gradient

- When doing back propagation, each node in a layer calculates its gradient with respect to the effects of the gradients in the layer before it.
- So if the adjustments to the layers before it is small, then adjustments to the current layer will be even smaller.

new weight = weight - learning rate*gradient

$$\boxed{2.0999} = \boxed{2.1} - \boxed{0.001}$$

Not much of a difference update value

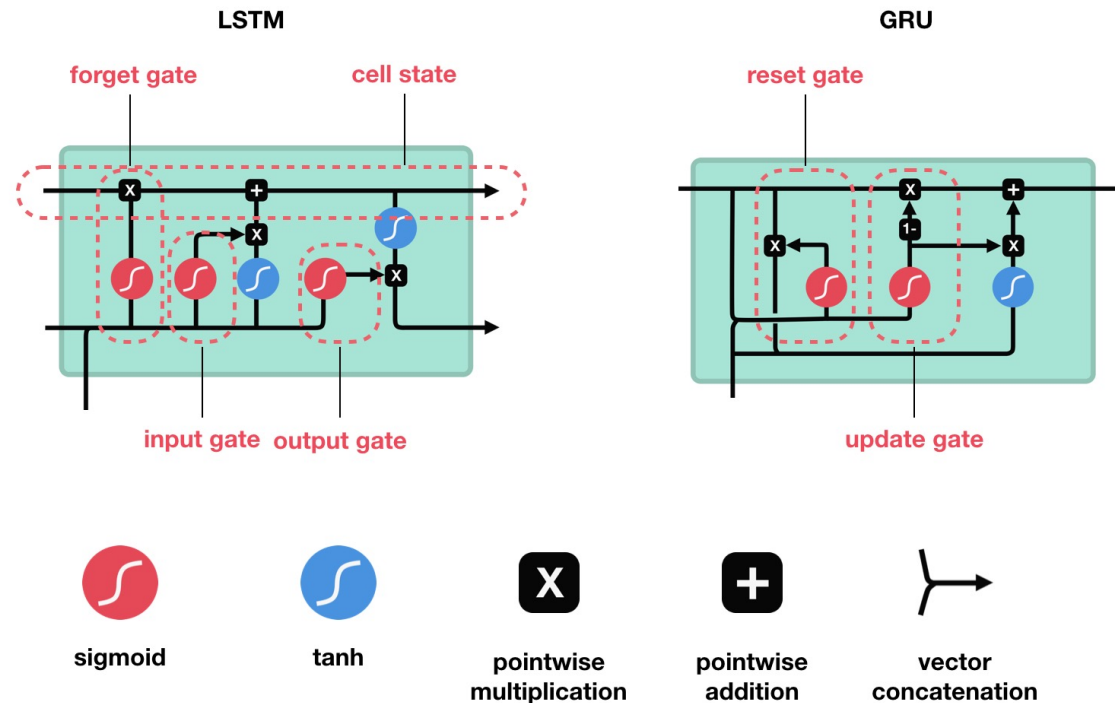
Vanishing Gradient



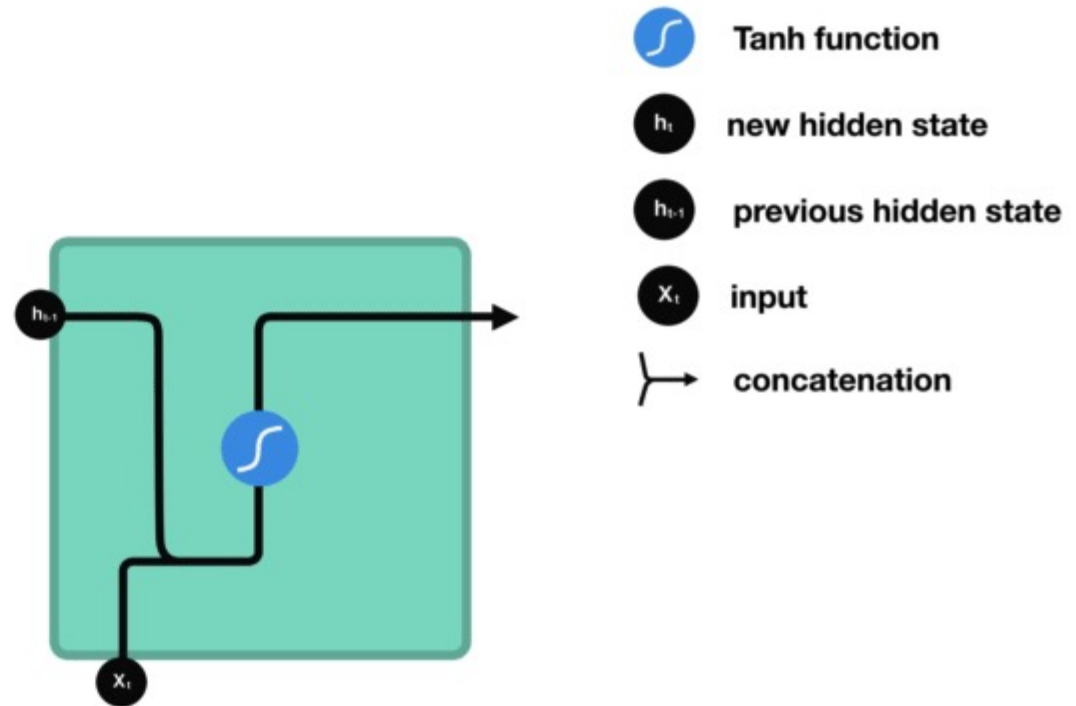
Gradients shrink as it back-propagates through time

Solutions to short-term memory

- Long Short-Term Memory (LSTM)
- Gated Recurrent Units (GRU)



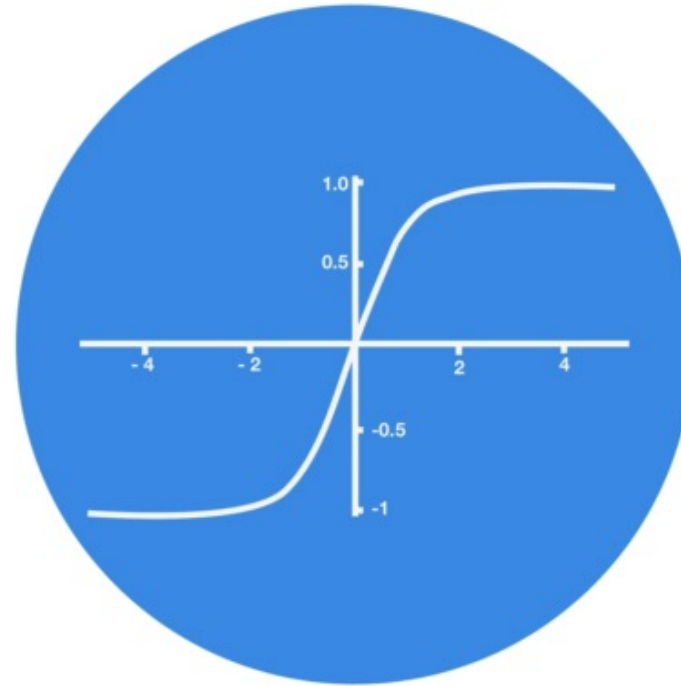
RNN (recap)



Credit: <https://towardsdatascience.com>

Tanh activation

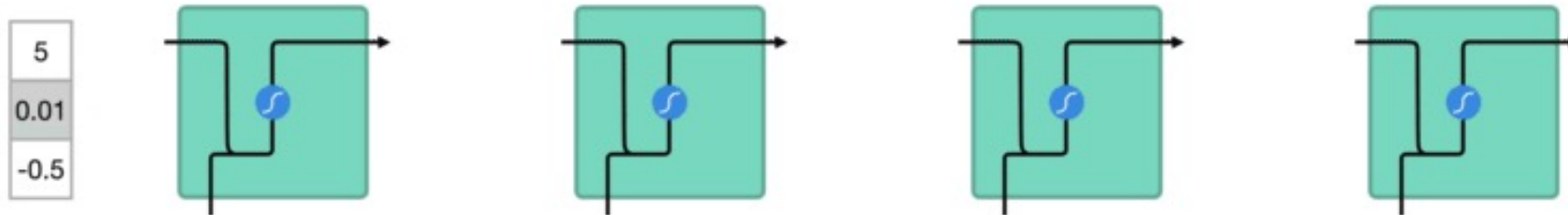
5
0.1
-0.5



Tanh activation



vector transformations without tanh

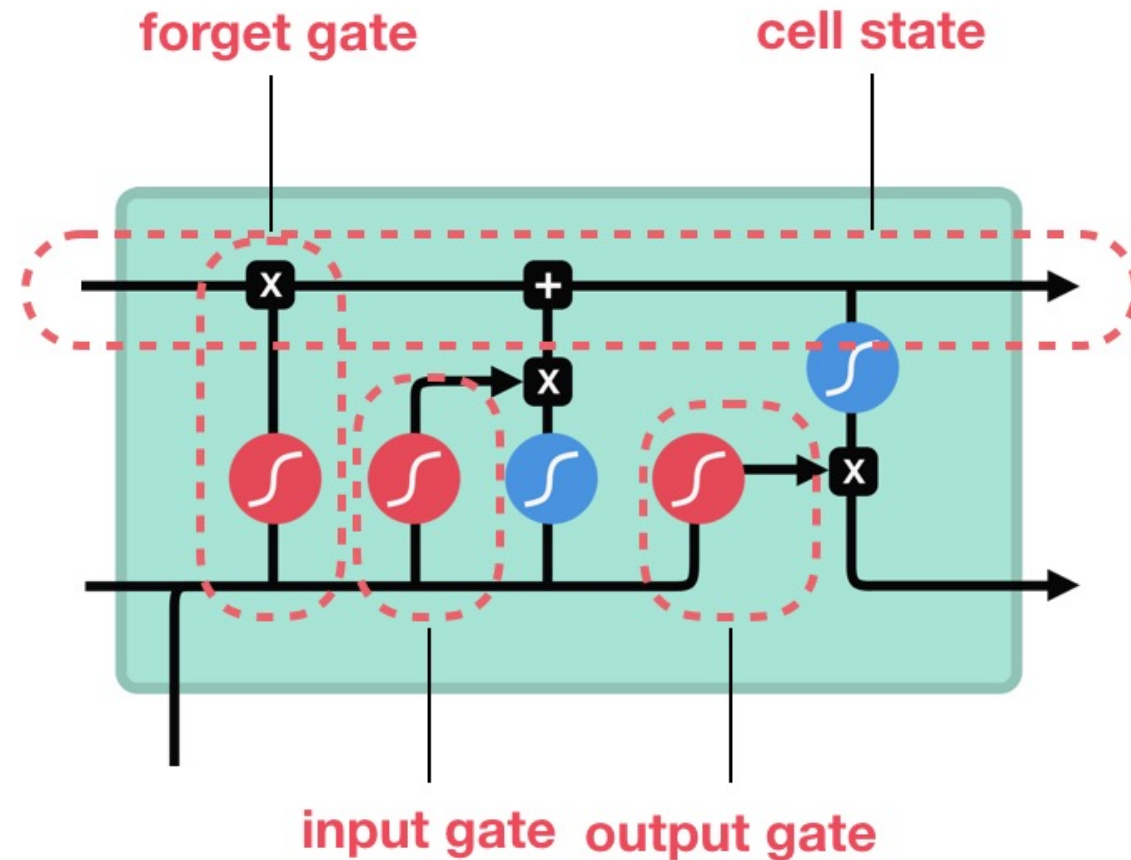


vector transformations with tanh

Long Short-Term Memory (LSTM)

- The core concept of LSTM's are the **cell state**, and various **gates**.
- The **cell state**, in theory, can carry relevant information throughout the processing of the sequence.
- As the cell state goes on its journey, information get's added or removed to the cell state via **gates**.
- The gates are different neural networks that decide which information is allowed on the cell state.
- The gates can learn what information is relevant to keep or forget during training.

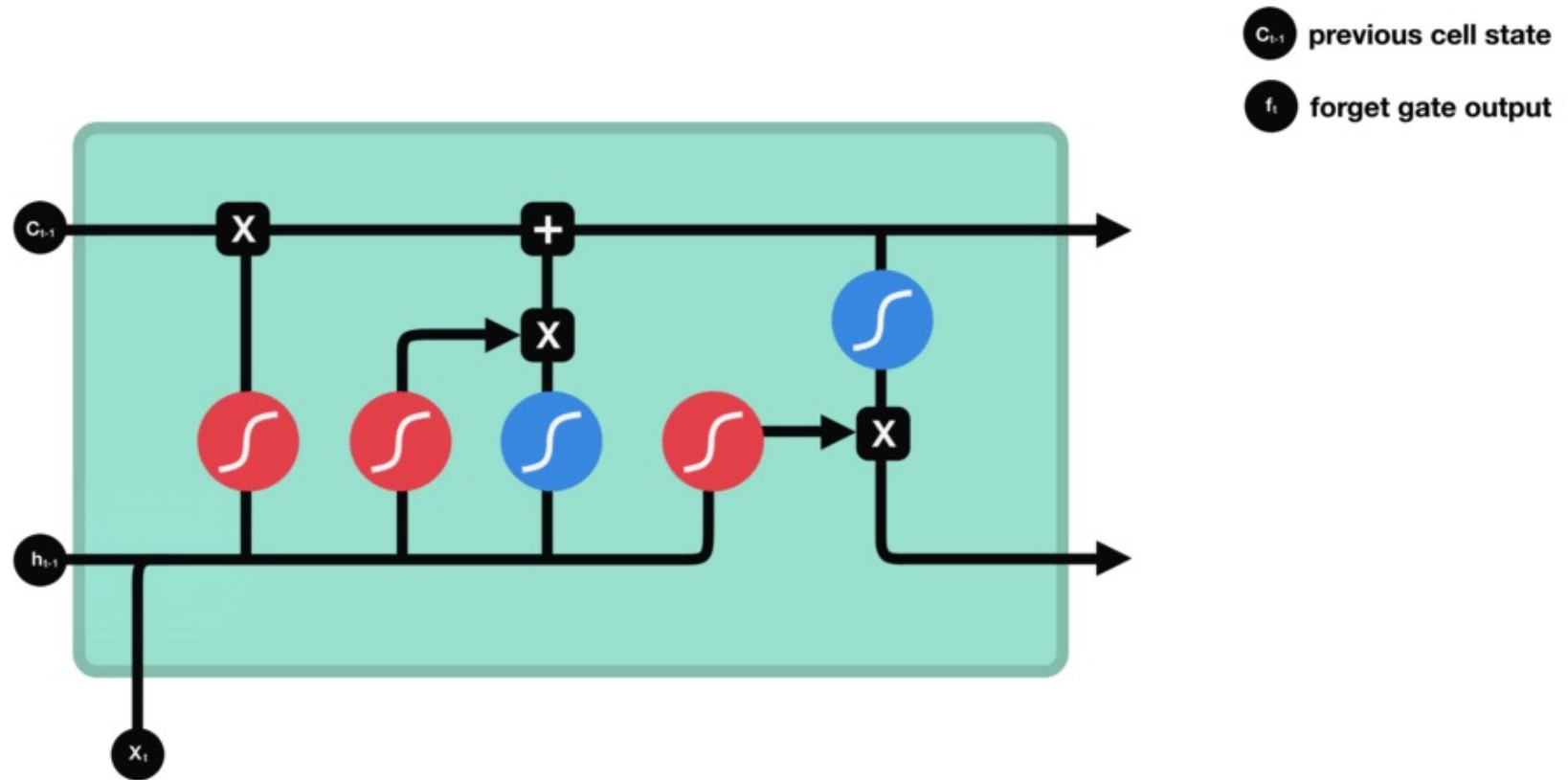
Long Short-Term Memory (LSTM)



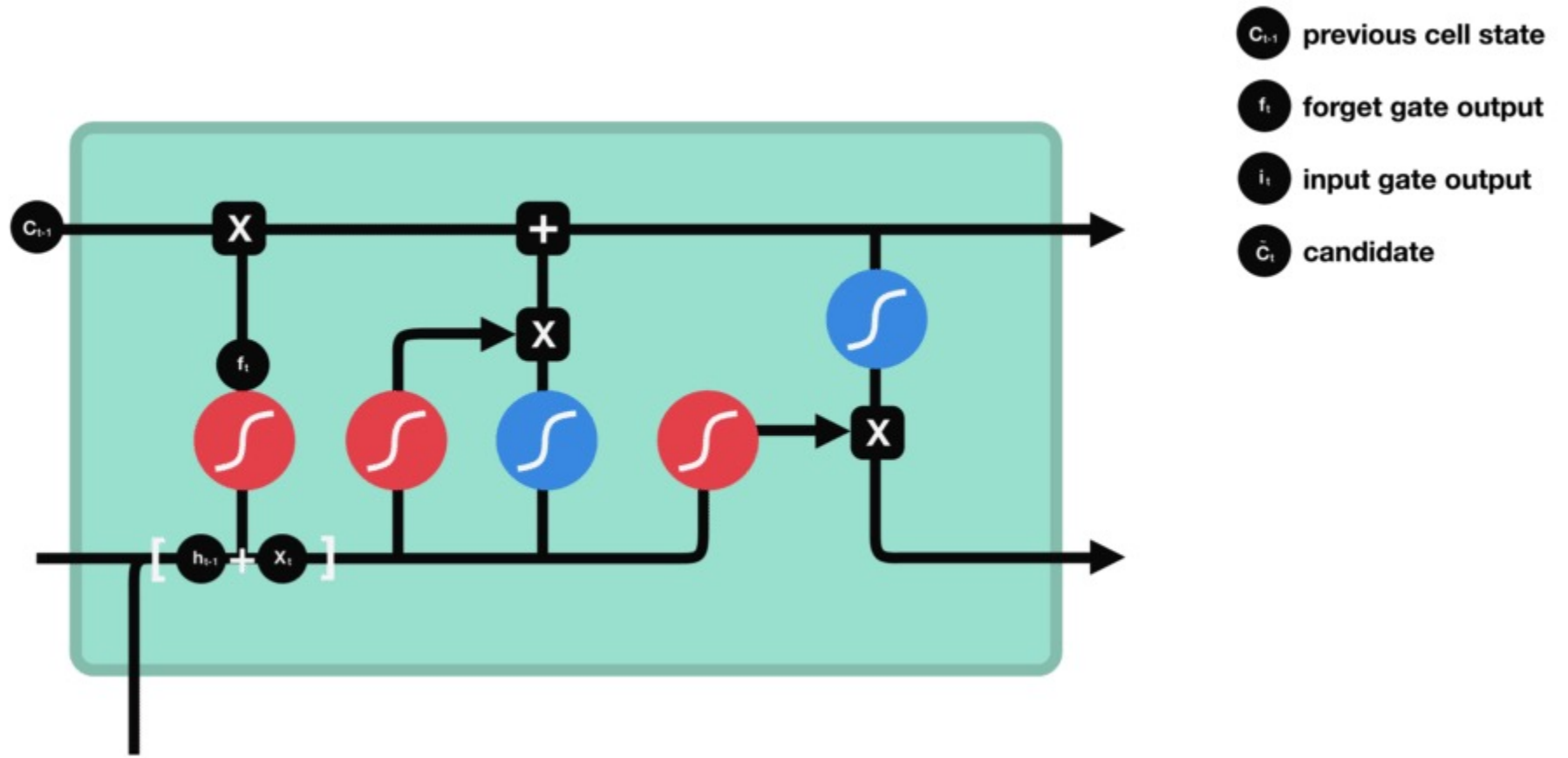
Long Short-Term Memory (LSTM)

- 3 main gates:
 - Forget gate
 - decides what is relevant to keep from prior steps
 - Input gate
 - decides what information is relevant to add from the current step
 - Output gate
 - decides what the next hidden state should be

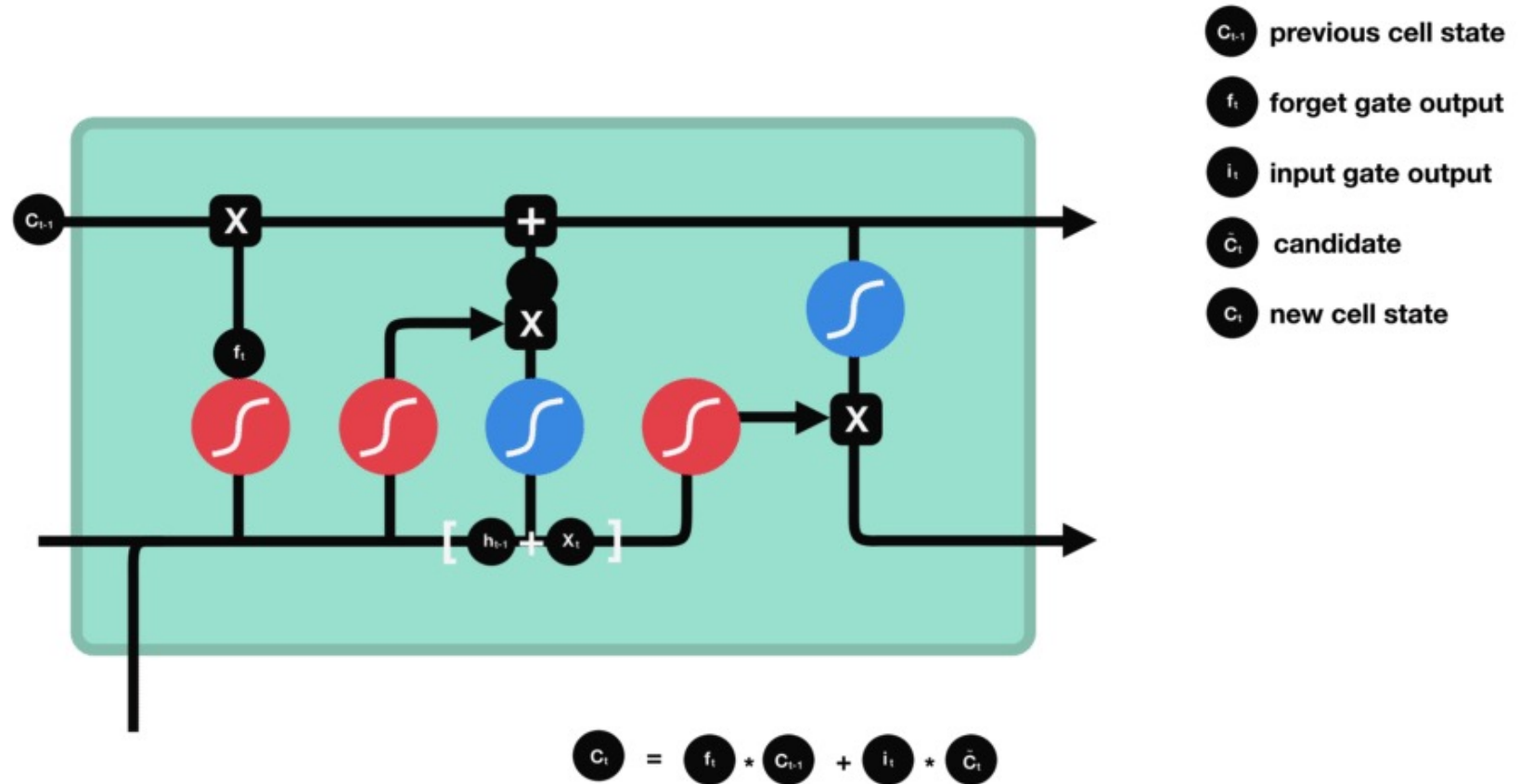
Forget Gate



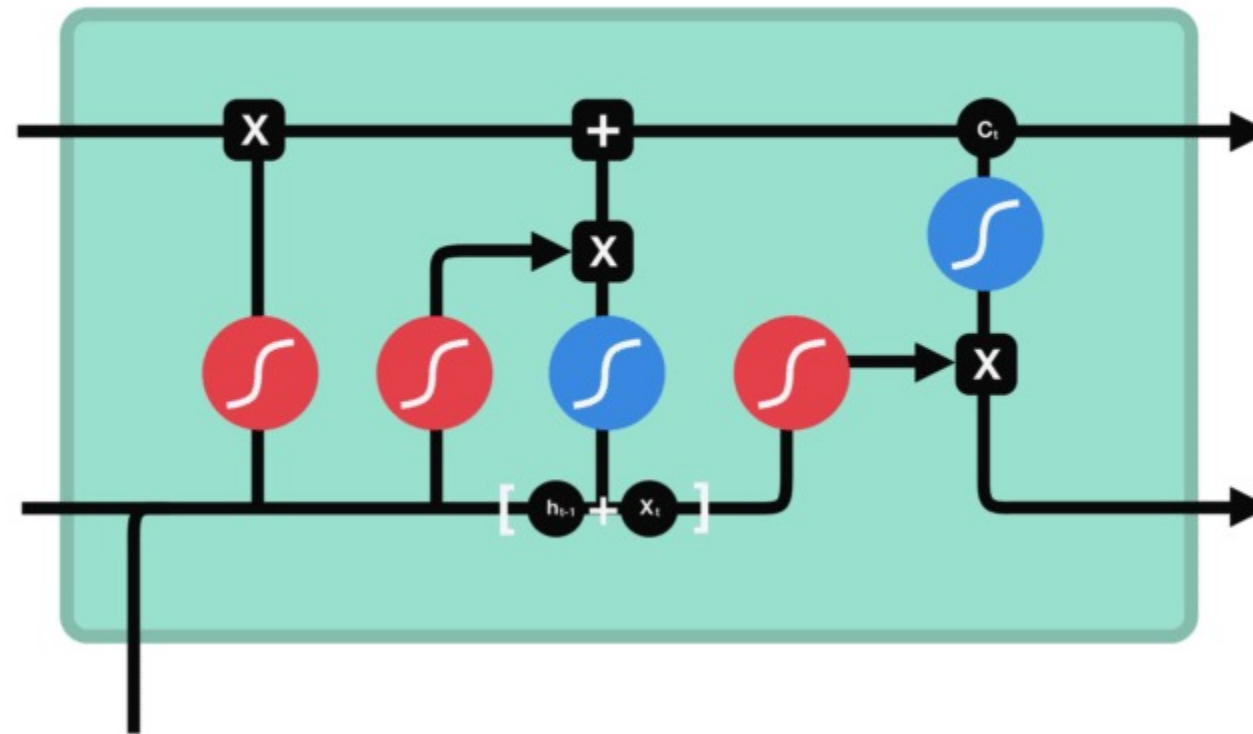
Input Gate



Cell State



Output Gate

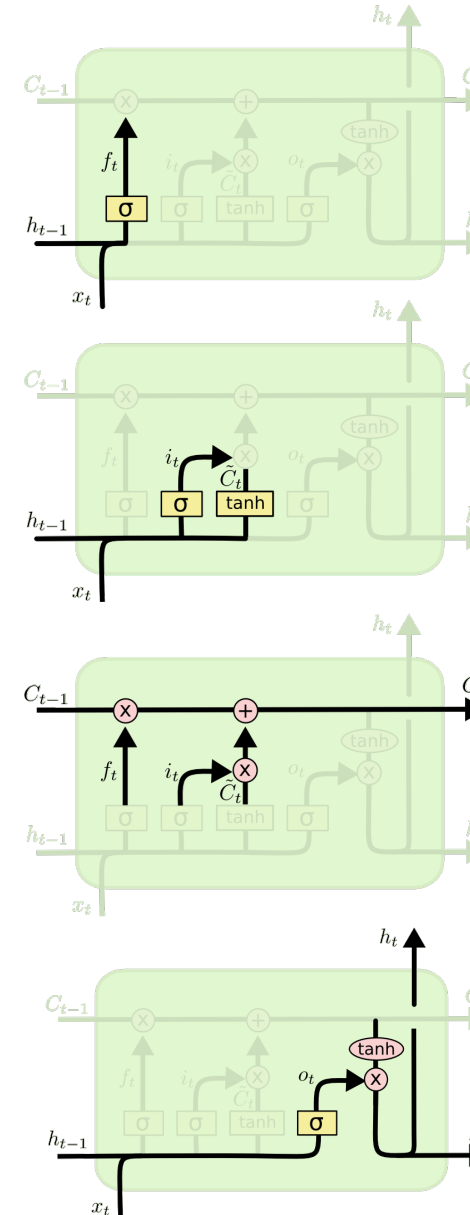


- c_{t-1} previous cell state
- f_t forget gate output
- i_t input gate output
- \tilde{c}_t candidate
- c_t new cell state
- o_t output gate output
- h_t hidden state

LSTM: input and new memory

LSTM cells takes the following input vectors:

- the input x_t
- the past output h_{t-1}
- the past cell state c_{t-1}



Forget Gate

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

Input Gate

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Cell State

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

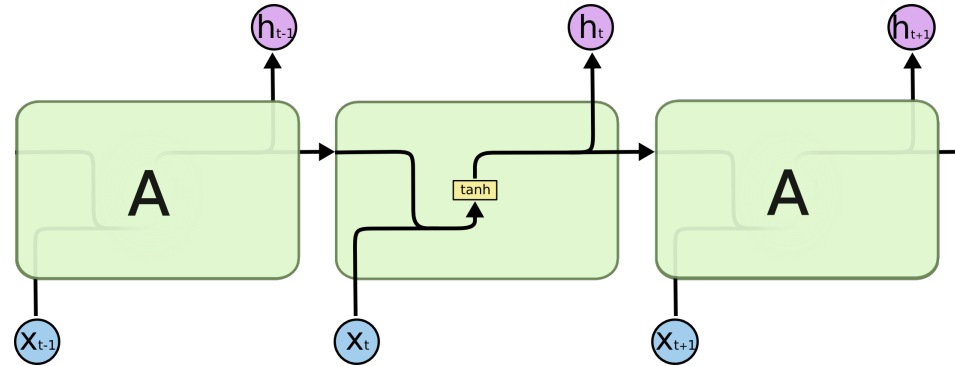
Output Gate

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

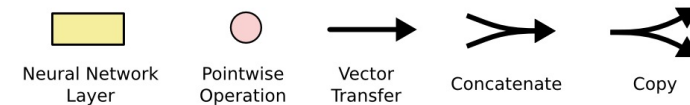
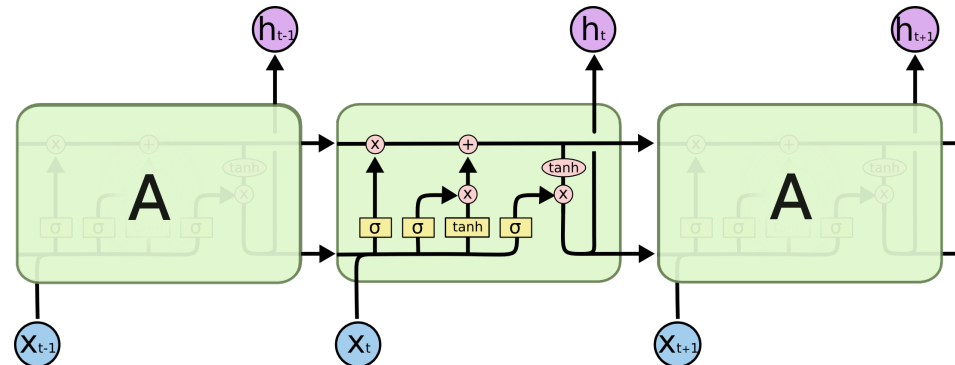
$$h_t = o_t * \tanh(C_t)$$

Standard RNNs to LSTM

Standard



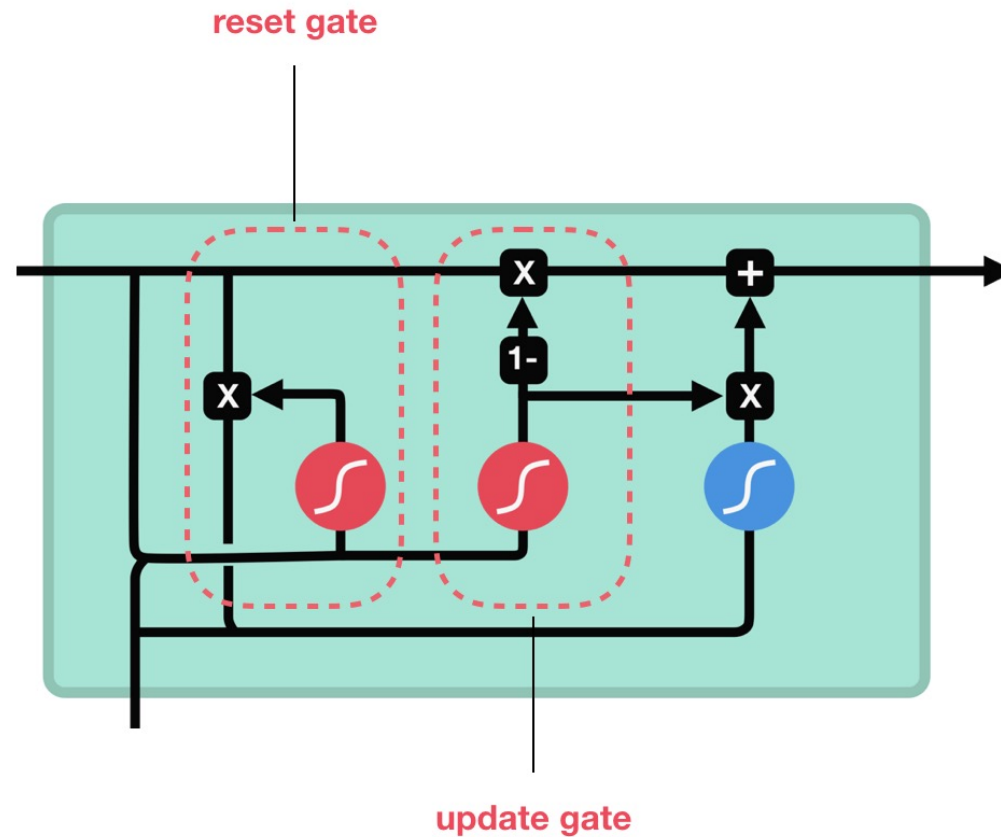
LSTM



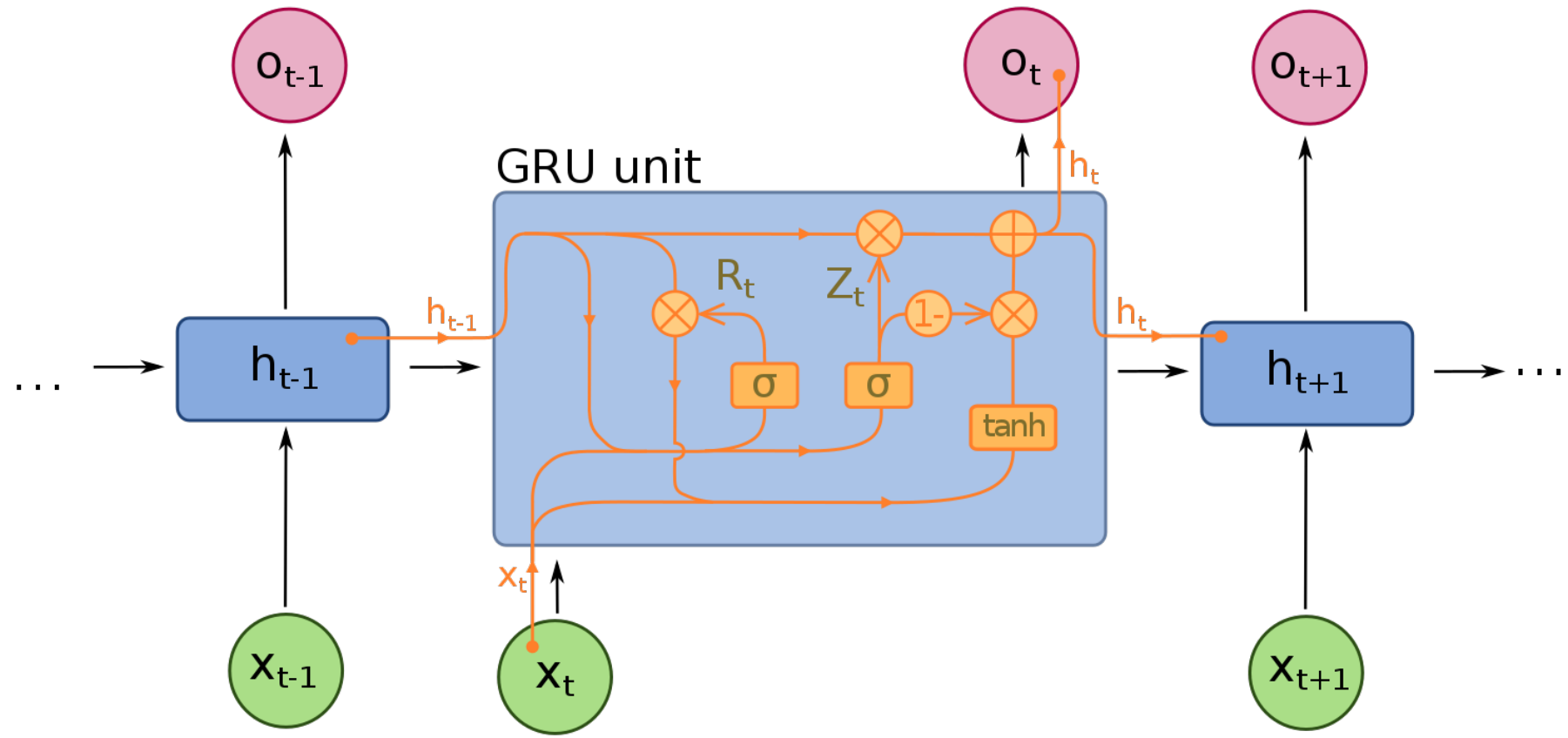
Gated Recurrent Units (GRU)

- GRU is the newer generation of Recurrent Neural networks and is pretty similar to an LSTM
- GRU's got rid of the cell state and used the hidden state to transfer information
- It only has two gates,
 - an update gate
 - acts similar to the forget and input gate of an LSTM
 - decides what information to throw away and what new information to add
 - a reset gate
 - decide how much past information to forget

Gated Recurrent Units (GRU)



GRU Unit



A mostly complete chart of Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

-  Backfed Input Cell
-  Input Cell
-  Noisy Input Cell
-  Hidden Cell
-  Probabilistic Hidden Cell
-  Spiking Hidden Cell
-  Output Cell
-  Match Input Output Cell
-  Recurrent Cell
-  Memory Cell
-  Different Memory Cell
-  Kernel
-  Convolution or Pool

Perceptron (P)



Feed Forward (FF)



Radial Basis Network (RBF)



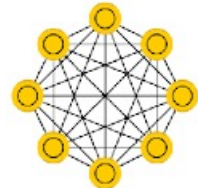
Deep Feed Forward (DFF)



Markov Chain (MC)



Hopfield Network (HN)



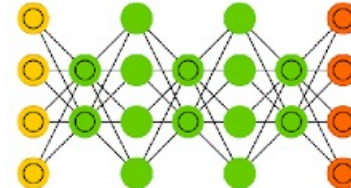
Boltzmann Machine (BM)



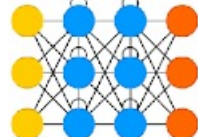
Restricted BM (RBM)



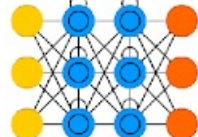
Deep Belief Network (DBN)



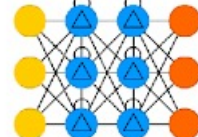
Recurrent Neural Network (RNN)



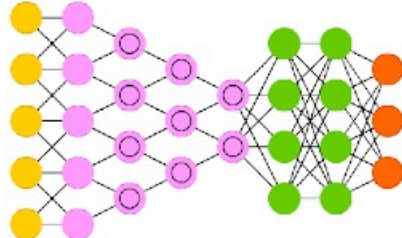
Long / Short Term Memory (LSTM)



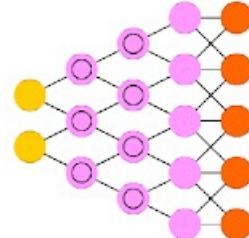
Gated Recurrent Unit (GRU)



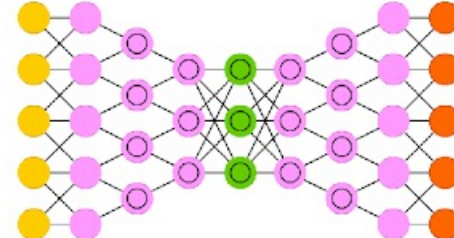
Deep Convolutional Network (DCN)



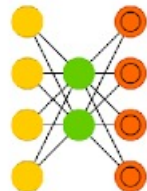
Deconvolutional Network (DN)



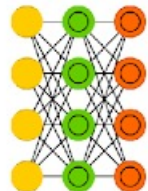
Deep Convolutional Inverse Graphics Network (DCIGN)



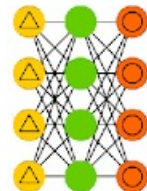
Auto Encoder (AE)



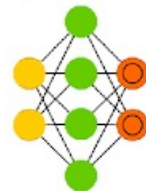
Variational AE (VAE)



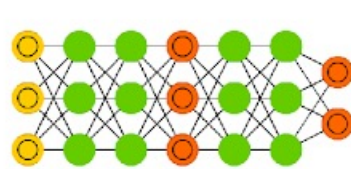
Denosing AE (DAE)



Sparse AE (SAE)



Generative Adversarial Network (GAN)



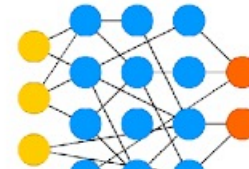
Liquid State Machine (LSM)



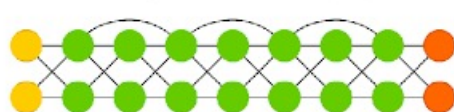
Extreme Learning Machine (ELM)



Echo State Network (ESN)



Deep Residual Network (DRN)



Kohonen Network (KN)



Support Vector Machine (SVM)



Neural Turing Machine (NTM)

