

# Recurrent Neural Networks

Deep Learning  
CS 435/635

Course Instructor: Chandresh  
AI Lab Coordinator @IIT Indore

# Course Content

Module I: History of Deep Learning, Sigmoid Neurons, Perceptrons, and learning algorithms. Multilayer Perceptrons (MLPs), Representation Power of MLPs,.

Module II: Feedforward Neural Networks. Backpropagation. first and second-order training methods. NN Training tricks, Regularization

Module III: Introduction to Autoencoders and their characteristics, relation to PCA, Regularization in autoencoders, and Types of autoencoders, Variational Autoencoder

Module IV: Architecture of Convolutional Neural Networks (CNN), types of CNNs. Image classification, Pre-training vs fine-tuning.- representation learning, Object Detection and Semantic Segmentation

Module V: **Architecture of Recurrent Neural Networks (RNN)**, Backpropagation through time. Encoder-Decoder Models, Attention Mechanism. Advanced Topics: Transformers and BERT.

Module VI: Gen AI- Deep generative models: VAE, GAN,

# Acknowledgement

- Russ Salakhutdinov and Hugo Larochelle's class on Neural Networks
- Neural Networks and Deep Learning course by Danna Gurari University of Colorado Boulder

# Today's Topics

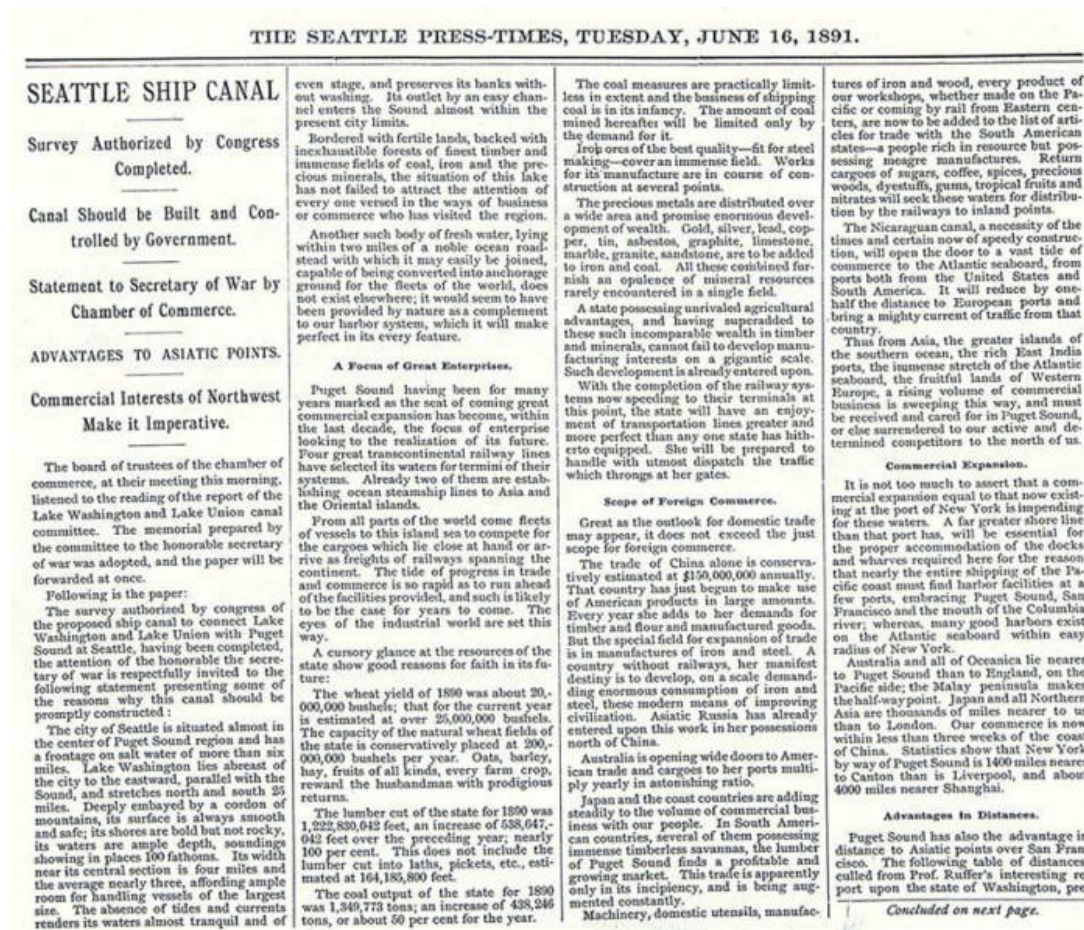
- Deep learning for sequential data
- Recurrent neural networks (RNNs)
- Gated RNNs
- Programming tutorial

# Today's Topics

- Deep learning for sequential data
- Recurrent neural networks (RNNs)
- Gated RNNs
- Programming tutorial

# Sequence Definition: Data of Arbitrary Length

e.g., Document

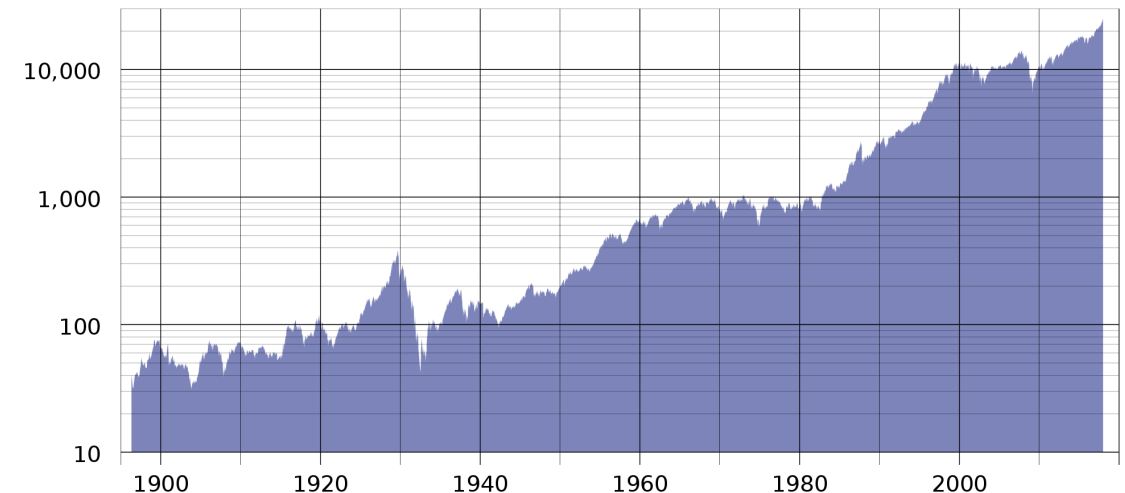


e.g., Images



e.g., Time-Series Data

Dow Jones Industrial Average



e.g., sentences, audio samples, brain waves, radio waves, air temperature



# Properties of Sequences?

e.g., Document

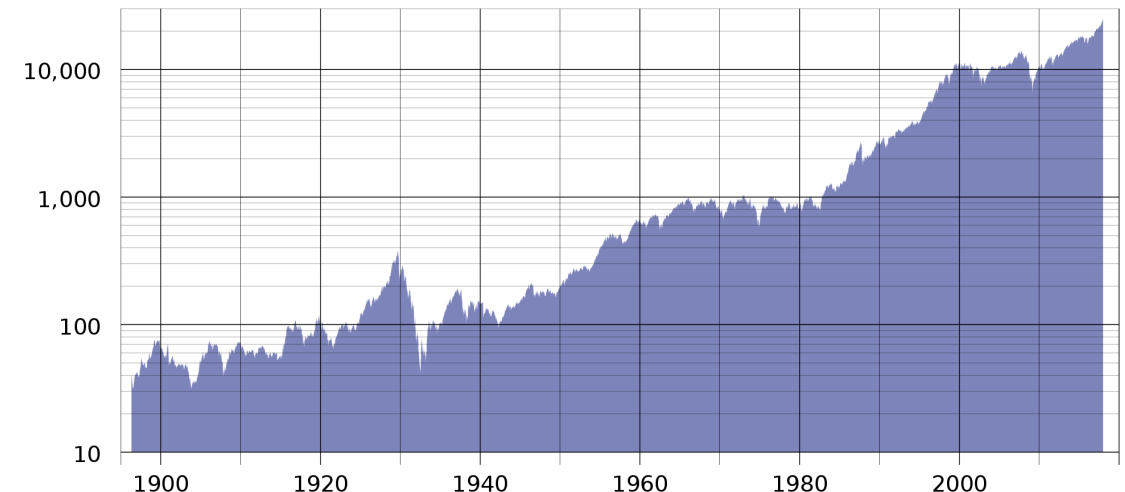
- \* Elements of a sequence occur in a certain order
- \* Elements depend on each other

e.g., Images



e.g., Time-Series Data

Dow Jones Industrial Average



e.g., sentences, audio samples, brain waves, radio waves, air temperature

THE SEATTLE PRESS-TIMES, TUESDAY, JUNE 16, 1891.

## SEATTLE SHIP CANAL.

### Survey Authorized by Congress Completed.

### Canal Should be Built and Controlled by Government.

### Statement to Secretary of War by Chamber of Commerce.

### ADVANTAGES TO ASIATIC POINTS.

### Commercial Interests of Northwest Make it Imperative.

The board of trustees of the chamber of commerce, at their meeting this morning, listened to the reading of the report of the Lake Washington and Lake Union canal committee. The memorial prepared by the committee to the honorable secretary of war was adopted, and the paper will be forwarded at once.

Following is the paper:

The survey authorized by congress of the proposed ship canal to connect Lake Washington and Lake Union with Puget Sound at Seattle, having been completed, the attention of the honorable the secretary of war is respectfully invited to the following statement presenting some of the reasons why this canal should be promptly constructed:

The city of Seattle is situated almost in the center of Puget Sound region and has a frontage on salt water of more than six miles. Lake Washington lies abreast of the city to the eastward, parallel with the Sound, and stretches north and south 25 miles. Deeply embayed by a cordon of mountains, its surface is always smooth and safe; its shores are bold but not rocky, its waters are ample depth, soundings showing in places 100 fathoms. Its width near its central section is four miles and the average nearly three, affording ample room for handling vessels of the largest size. The absence of tides and currents renders its waters almost tranquil and of even stage, and preserves its banks without washing. Its outlet by an easy channel enters the Sound almost within the present city limits.

Bordered with fertile lands, backed with inexhaustible forests of finest timber and immense fields of coal, iron and the precious minerals, the situation of this lake has not failed to attract the attention of every one versed in the ways of business or commerce who has visited the region.

Another such body of fresh water, lying within two miles of a noble ocean roadstead with which it may easily be joined, capable of being converted into anchorage ground for the fleets of the world, does not exist elsewhere; it would seem to have been provided by nature as a complement to our harbor system, which it will make perfect in its every feature.

### A Focus of Great Enterprises.

Puget Sound having been for many years marked as the seat of coming great commercial expansion has become, within the last decade, the focus of enterprise looking to the realization of its future. Four great transcontinental railway lines have selected its waters for termini of their systems. Already two of them are establishing ocean steamship lines to Asia and the Oriental islands.

From all parts of the world come fleets of vessels to this island sea to compete for the cargoes which lie close at hand or arrive as freights of railways spanning the continent. The tide of progress in trade and commerce is so rapid as to run ahead of the facilities provided, and such is likely to be the case for years to come. The eyes of the industrial world are set this way.

A cursory glance at the resources of the state show good reasons for faith in its future:

The wheat yield of 1890 was about 20,000,000 bushels; that for the current year is estimated at over 25,000,000 bushels. The capacity of the natural wheat fields of the state is conservatively placed at 200,000,000 bushels per year. Oats, barley, hay, fruits of all kinds, every farm crop reward the husbandman with prodigious returns.

The lumber cut of the state for 1890 was 1,222,830,042 feet, an increase of 538,047,042 feet over the preceding year; nearly 100 per cent. This does not include the lumber cut into laths, pickets, etc., estimated at 164,183,800 feet.

The coal output of the state for 1890 was 1,340,773 tons; an increase of 438,246 tons, or about 50 per cent for the year.

The coal measures are practically limitless in extent and the business of shipping coal is in its infancy. The amount of coal mined hereafter will be limited only by the demand for it.

Trop ores of the best quality—fit for steel making—cover an immense field. Works for its manufacture are in course of construction at several points.

The precious metals are distributed over a wide area and promise enormous development of wealth. Gold, silver, lead, copper, tin, asbestos, graphite, limestone, marble, granite, sandstone, are to be added to iron and coal. All these combined furnish an opulence of mineral resources rarely encountered in a single field.

A state possessing unrivaled agricultural advantages, and having superadded to these such incomparable wealth in timber and minerals, cannot fail to develop manufacturing interests on a gigantic scale. Such development is already entered upon.

With the completion of the railway systems now speeding to their terminals at this point, the state will have an enjoyment of transportation lines greater and more perfect than any one state has hitherto equipped. She will be prepared to handle with utmost dispatch the traffic which throngs at her gates.

### Scope of Foreign Commerce.

Great as the outlook for domestic trade may appear, it does not exceed the just scope for foreign commerce.

The trade of China alone is conservatively estimated at \$150,000,000 annually. That country has just begun to make use of American products in large amounts. Every year she adds to her demands for timber and flour and manufactured goods. But the special field for expansion of trade is in manufactures of iron and steel. A country without railways, her manifest destiny is to develop, on a scale demanding enormous consumption of iron and steel, these modern means of improving civilization. Asiatic Russia has already entered upon this work in her possessions north of China.

Australia is opening wide doors to American trade and cargoes to her ports multiply yearly in astonishing ratio.

Japan and the coast countries are adding steadily to the volume of commercial business with our people. In South American countries, several of them possessing immense timberless savannas, the lumber of Puget Sound finds a profitable and growing market. This trade is apparently only in its incipency, and is being augmented constantly.

Machinery, domestic utensils, manufactures of iron and wood, every product of our workshops, whether made on the Pacific or coming by rail from Eastern centers, are now to be added to the list of articles for trade with the South American states—a people rich in resource but possessing meagre manufactures. Return cargoes of sugars, coffee, spices, precious woods, dyestuffs, gums, tropical fruits and nitrates will seek these waters for distribution by the railways to inland points.

The Nicaraguan canal, a necessity of the times and certain now of speedy construction, will open the door to a vast tide of commerce to the Atlantic seaboard, from ports both from the United States and South America. It will reduce by one-half the distance to European ports and bring a mighty current of traffic from that country.

Thus from Asia, the greater islands of the southern ocean, the rich East India ports, the immense stretch of the Atlantic seaboard, the fruitful lands of Western Europe, a rising volume of commercial business is sweeping this way, and must be received and cared for in Puget Sound, or else surrendered to our active and determined competitors to the north of us.

### Commercial Expansion.

It is not too much to assert that a commercial expansion equal to that now existing at the port of New York is impending for these waters. A far greater shore line than that port has, will be essential for the proper accommodation of the docks and wharves required here for the reason that nearly the entire shipping of the Pacific coast must find harbor facilities at a few ports, embracing Puget Sound, San Francisco and the mouth of the Columbia river; whereas, many good harbors exist on the Atlantic seaboard within easy radius of New York.

Australia and all of Oceania lie nearer to Puget Sound than to England, on the Pacific side; the Malay peninsula makes the half-way point. Japan and all Northern Asia are thousands of miles nearer to us than to London. Our commerce is now within less than three weeks of the coast of China. Statistics show that New York by way of Puget Sound is 1400 miles nearer to Canton than is Liverpool, and about 4000 miles nearer Shanghai.

### Advantages in Distances.

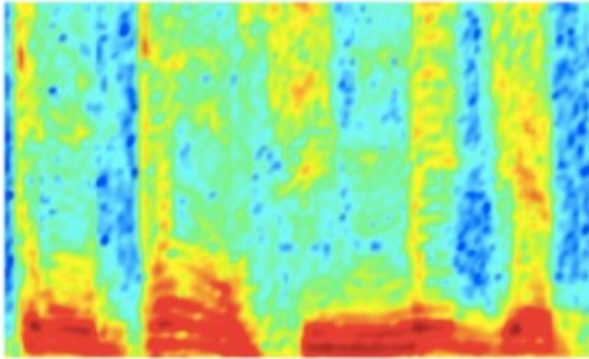
Puget Sound has also the advantage in distance to Asiatic points over San Francisco. The following table of distances, culled from Prof. Ruffer's interesting report upon the state of Washington, presented constantly.

Concluded on next page.

# Sequence Sources

- \* Elements of a sequence occur in a certain order
- \* Elements depend on each other

## AUDIO



Audio Spectrogram

## IMAGES

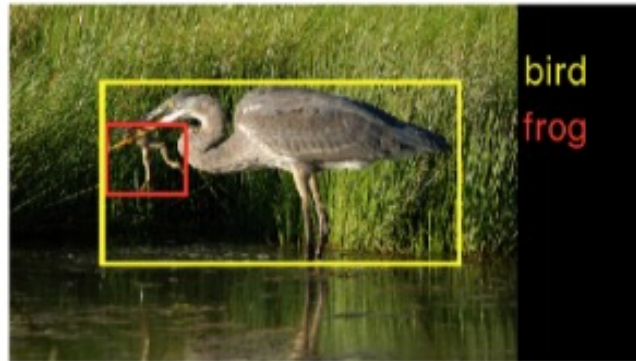
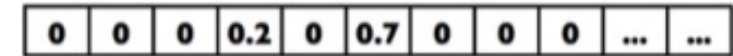


Image pixels

## TEXT

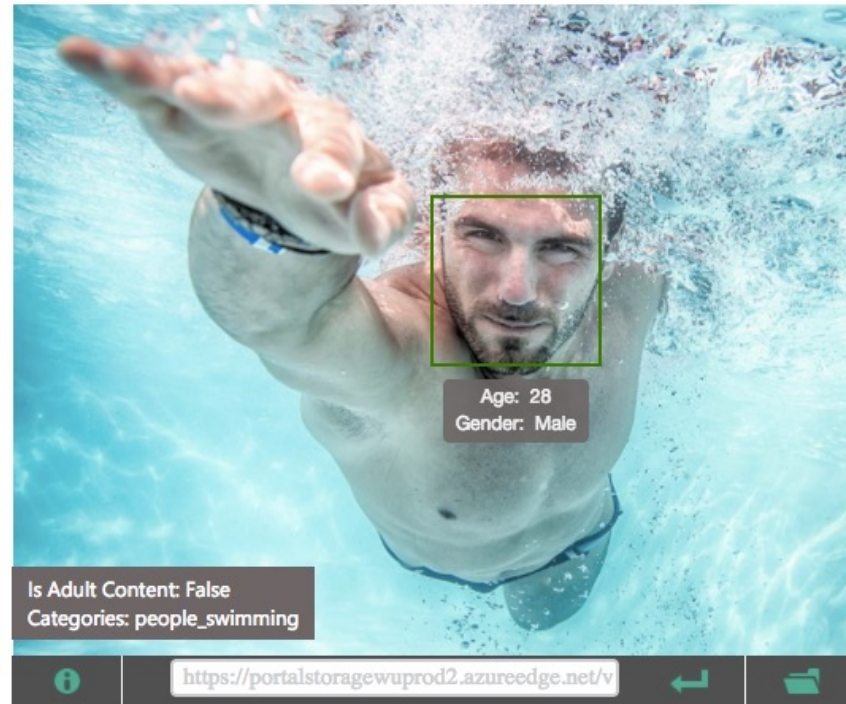


Word, context, or  
document vectors



# Sequence Applications: One-to-Many

- **Input:** fixed-size
- **Output:** sequence
- e.g., image captioning



Features:	
Feature Name	Value
Description	{ "type": 0, "captions": [ { "text": "a man swimming in a pool of water", "confidence": 0.7850108693093019 } ] }
Tags	[ { "name": "water", "confidence": 0.9996442794799805 }, { "name": "sport", "confidence": 0.9504992365837097 }, { "name": "swimming", "confidence": 0.9062818288803101, "hint": "sport" }, { "name": "pool", "confidence": 0.8787588477134705 }, { "name": "water sport", "confidence": 0.631849467754364, "hint": "sport" } ]
Image Format	jpeg
Image Dimensions	1500 x 1155
Clip Art Type	0 Non-clipart
Line Drawing Type	0 Non-LineDrawing
Black & White Image	False


Captions: <https://www.microsoft.com/cognitive-services/en-us/computer-vision-api>

# Sequence Applications: Many-to-One


- **Input:** sequence
- **Output:** fixed-size
- e.g., sentiment analysis (hate? love?, etc)


**CRITIC REVIEWS FOR *STAR WARS: THE LAST JEDI***

All Critics (371) | Top Critics (51) | Fresh (336) | Rotten (35)


 What's most interesting to me about The Last Jedi is Luke's return as the mentor rather than the student, grappling with his failure in this new role, and later aspiring to be the wise and patient teacher.

December 26, 2017 | Rating: 3/4 | [Full Review...](#)

 **Leah Pickett**  
Chicago Reader  
★ Top Critic

 Fanatics will love it; for the rest of us, it's a tolerably good time.

December 15, 2017 | Rating: B | [Full Review...](#)

 **Peter Rainer**  
Christian Science Monitor  
★ Top Critic

[https://www.rottentomatoes.com/m/star\\_wars\\_the\\_last\\_jedi](https://www.rottentomatoes.com/m/star_wars_the_last_jedi)

# Sequence Applications: Many-to-Many

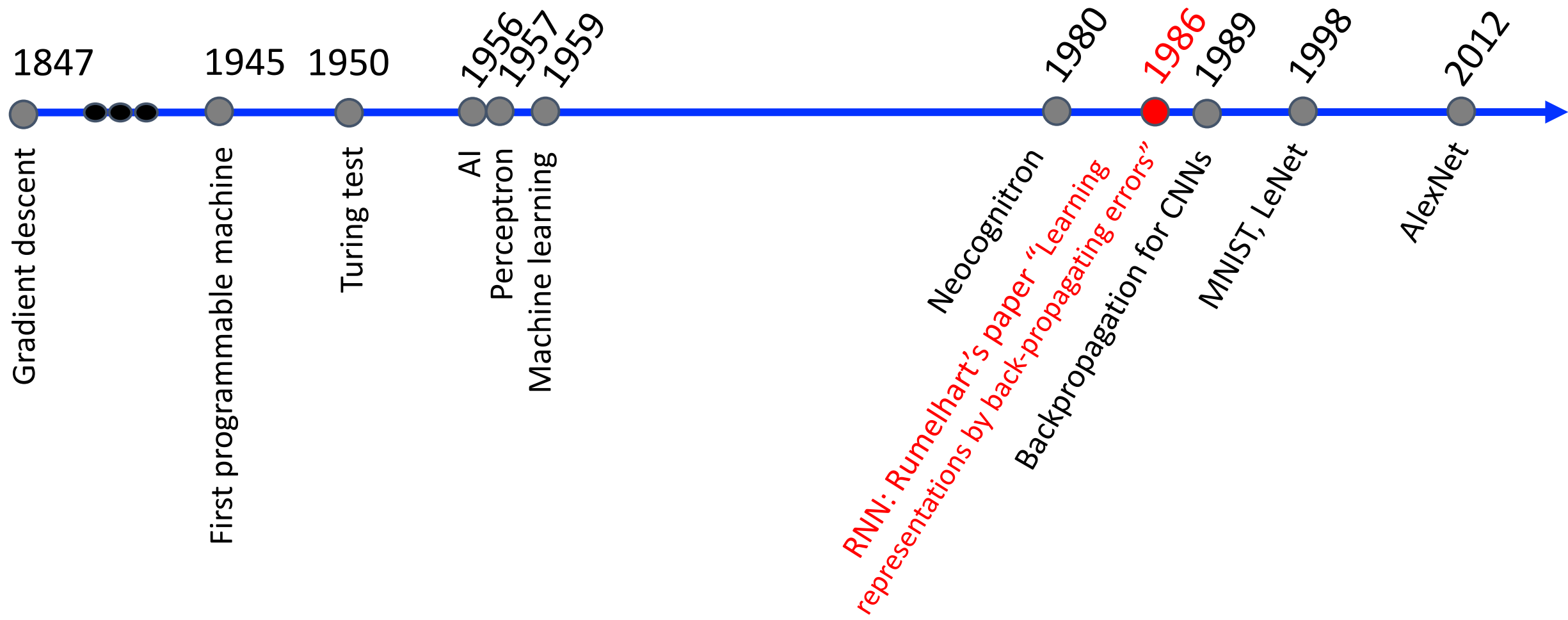
- **Input:** sequence
- **Output:** sequence
- e.g., language translation



# Today's Topics

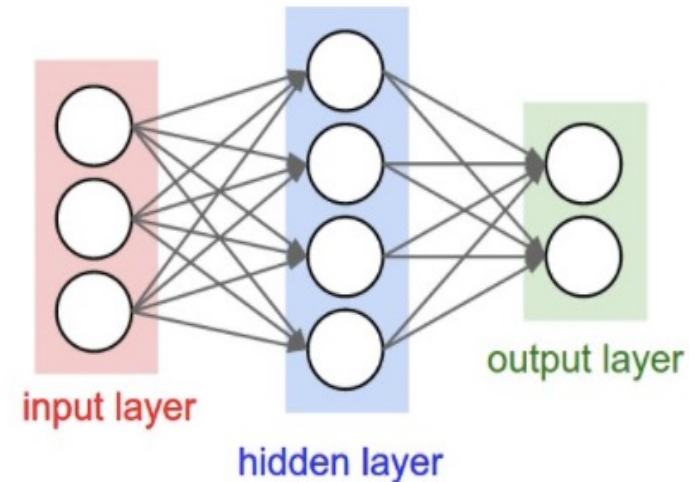
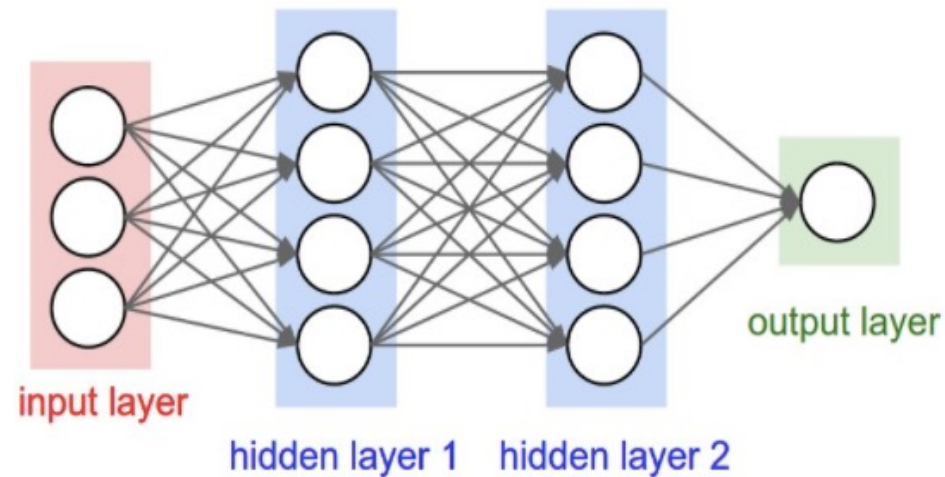
- Deep learning for sequential data
- Recurrent neural networks (RNNs)
- Gated RNNs
- Programming tutorial

# Historical Context





# Recall: Feedforward Neural Networks



**Problem:** many model parameters

**Problem:** inputs/output sizes are fixed

**Problem:** no memory of past since weights learned independently

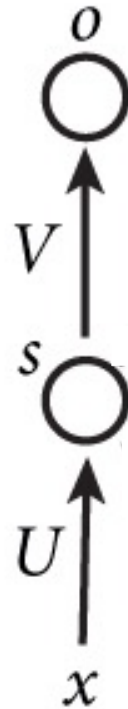
**Each layer serves as input to the next layer with no loops**

# Recurrent Neural Networks (RNNs)

- Main idea: use hidden state to **capture information about the past**

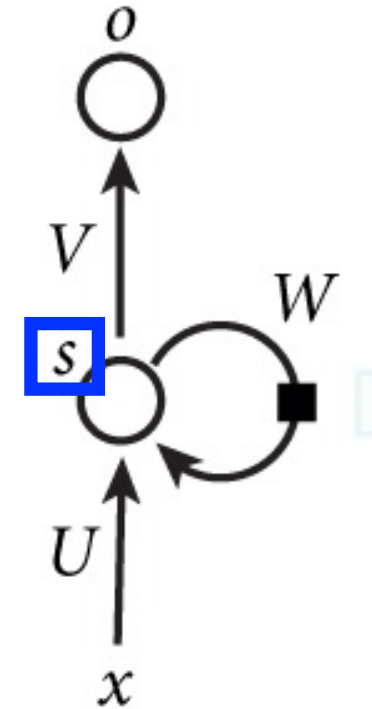
## Feedforward Network

Each layer receives input from the previous layer with no loops



## Recurrent Network

Each layer receives input from the previous layer **and the output from the previous time step (i.e., time-delayed connections)**



# Recurrent Neural Networks (RNNs)

- Main idea: use hidden state to **capture information about the past**

Recurrent: same function is applied to previous result

Model parameters

$$s_t = f_m(s_{t-1}, x_t)$$

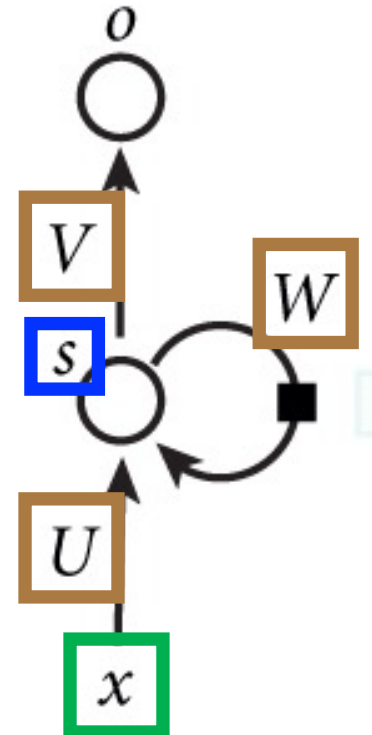
New  
state

Old  
state

Input at  
time step

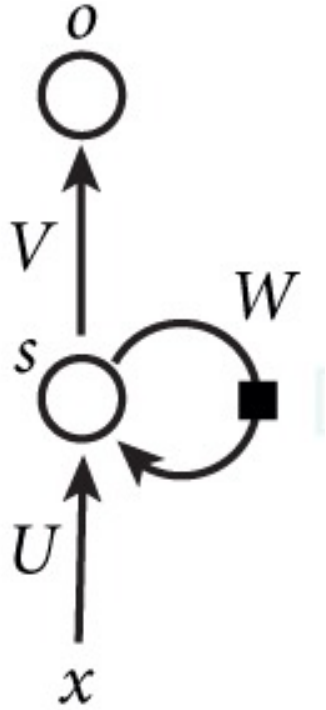
## Recurrent Network

Each layer receives input from the previous layer and the output from the previous time step (i.e., time-delayed connections)



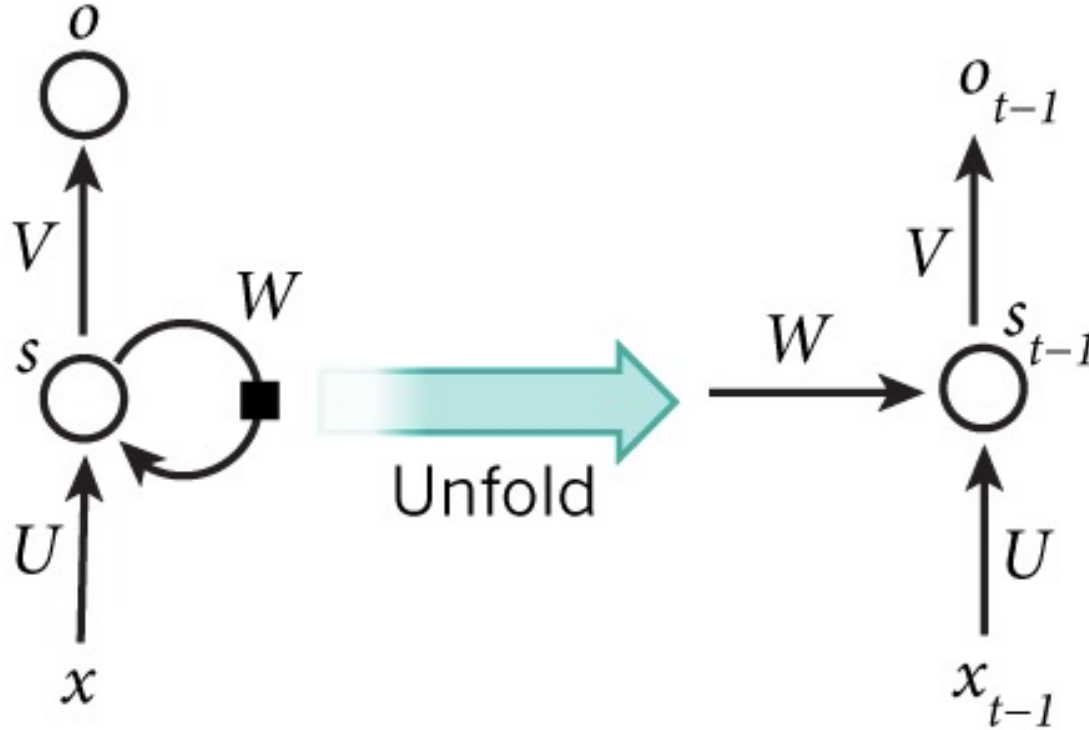
# RNN: Time Step 1

- Main idea: use hidden state to capture information about the past



# RNN: Time Step 1

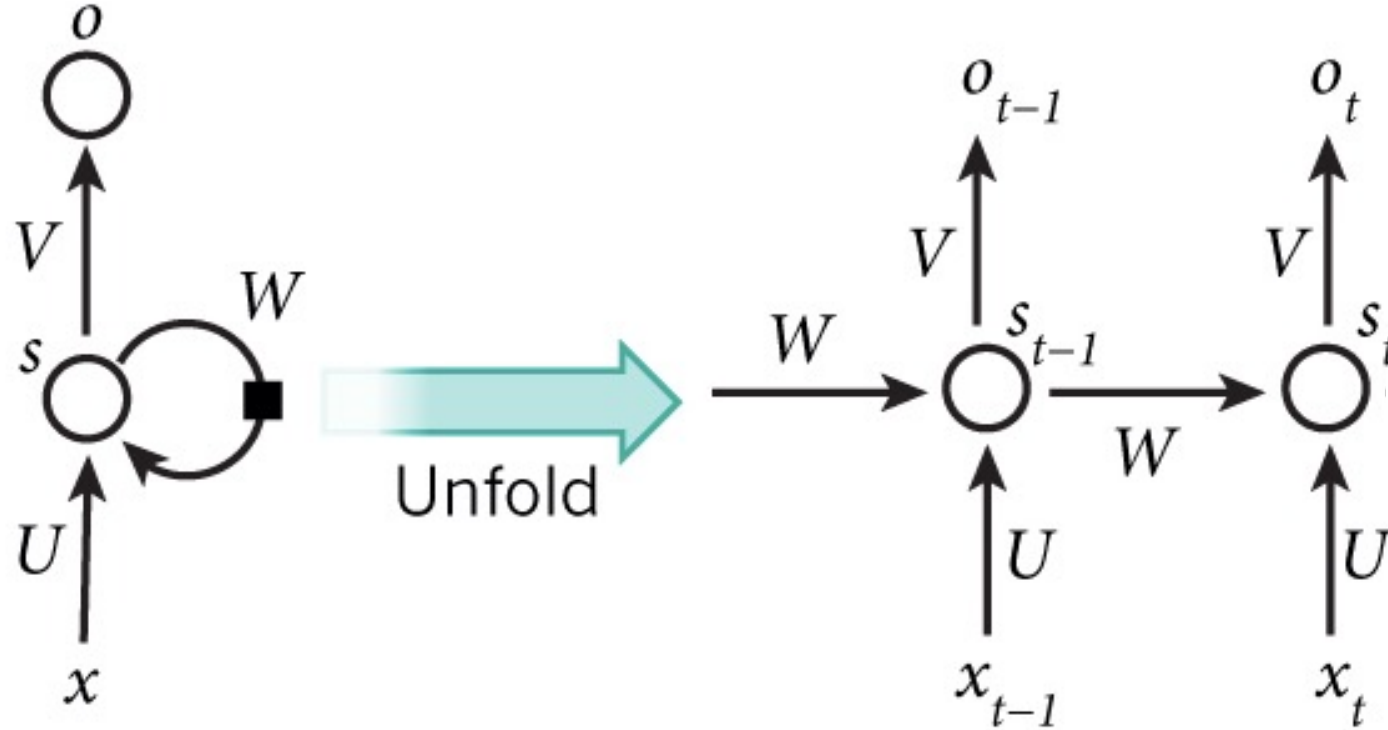
- Main idea: use hidden state to capture information about the past





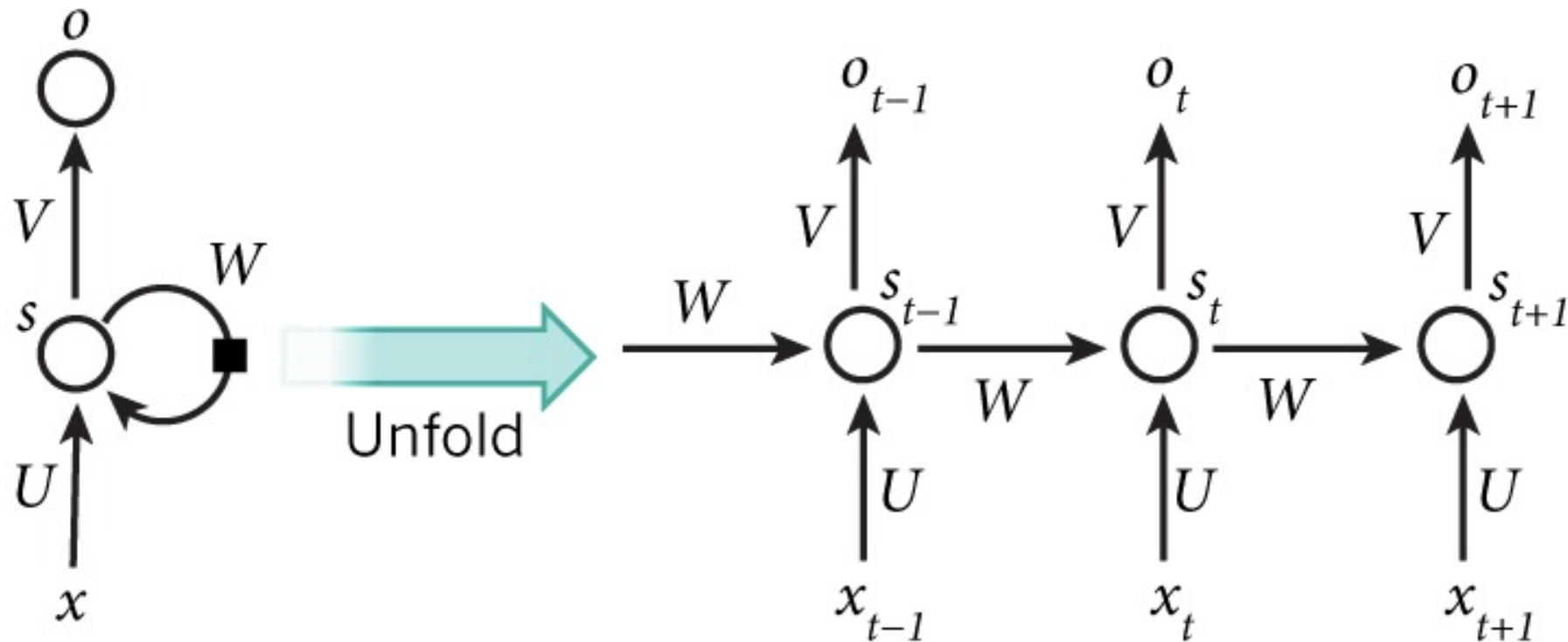
# RNN: Time Step 2

- Main idea: use hidden state to capture information about the past



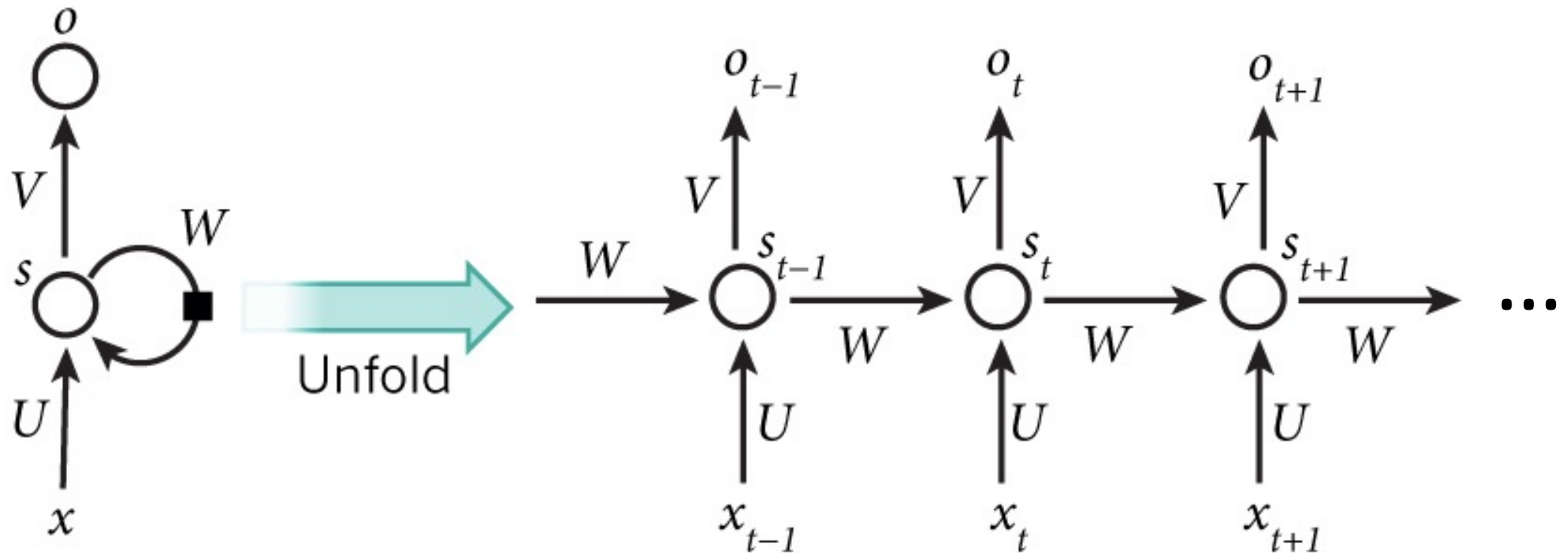
# RNN: Time Step 3

- Main idea: use hidden state to capture information about the past



# RNN: And So On...

- Main idea: use hidden state to capture information about the past



# RNN: Model Parameters and Inputs

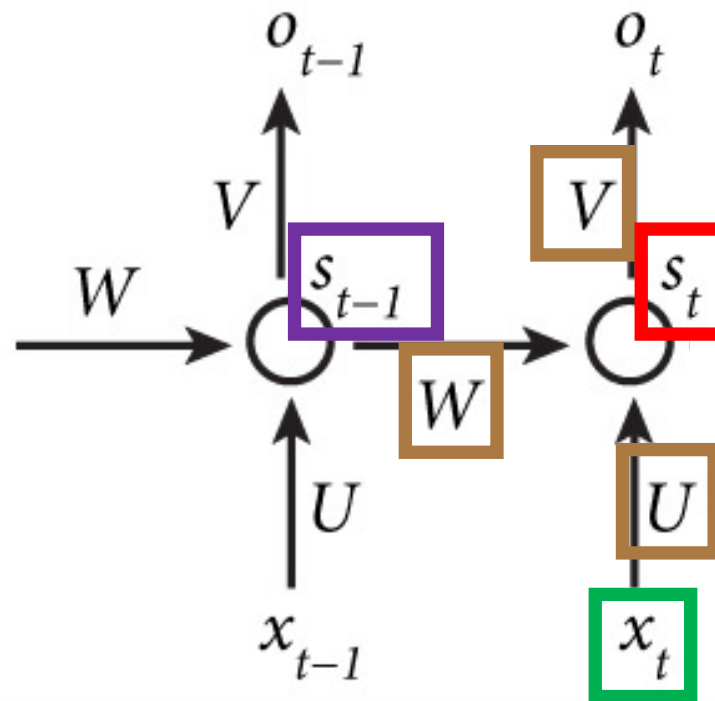
- Main idea: use hidden state to capture information about the past

Recurrence formula applied at every time step:

Model parameters

$$s_t = f_m(s_{t-1}, x_t)$$

New state      Old state      Input at time step



# RNN: Model Parameters and Inputs

- Main idea: use hidden state to capture information about the past

Recurrence formula applied at every time step:

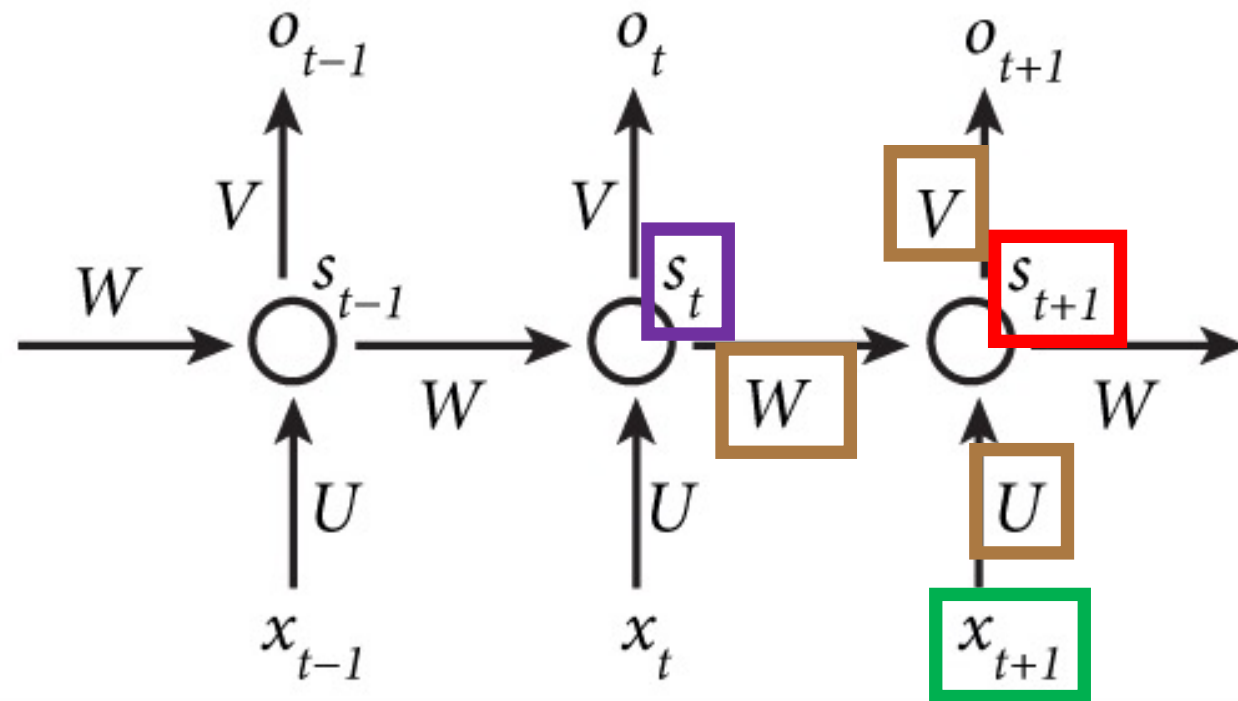
Model parameters

$$s_t = f_m(s_{t-1}, x_t)$$

New  
state

Old  
state

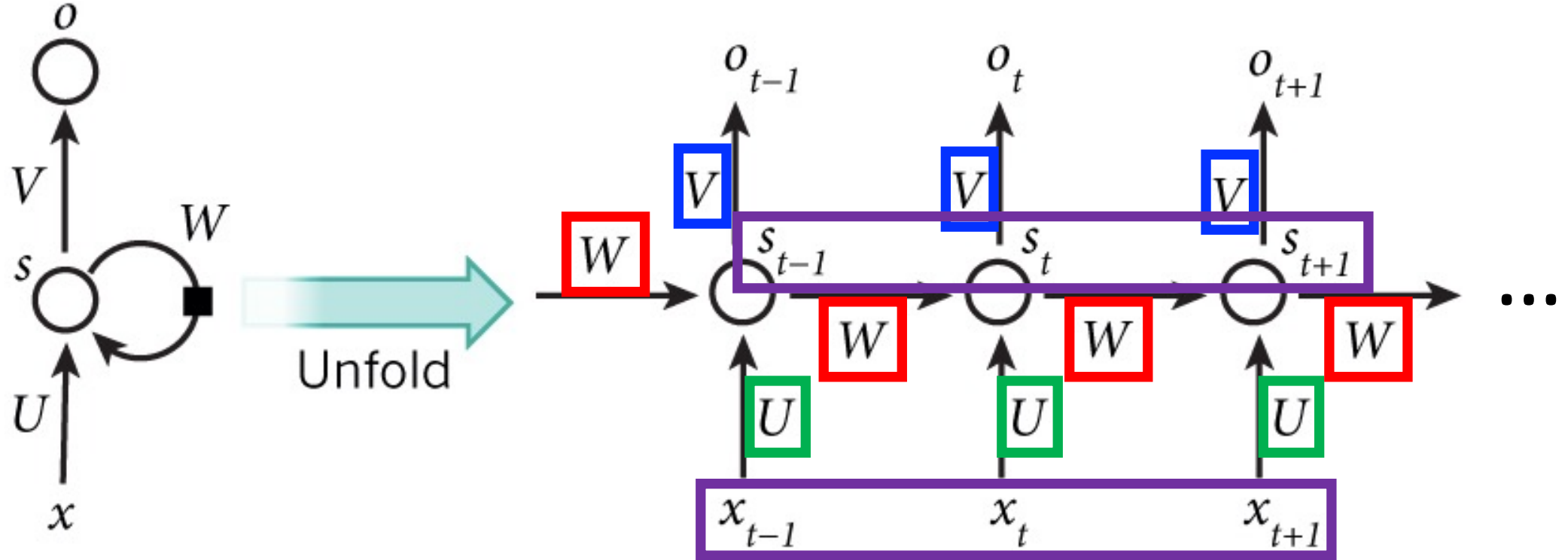
Input at  
time step





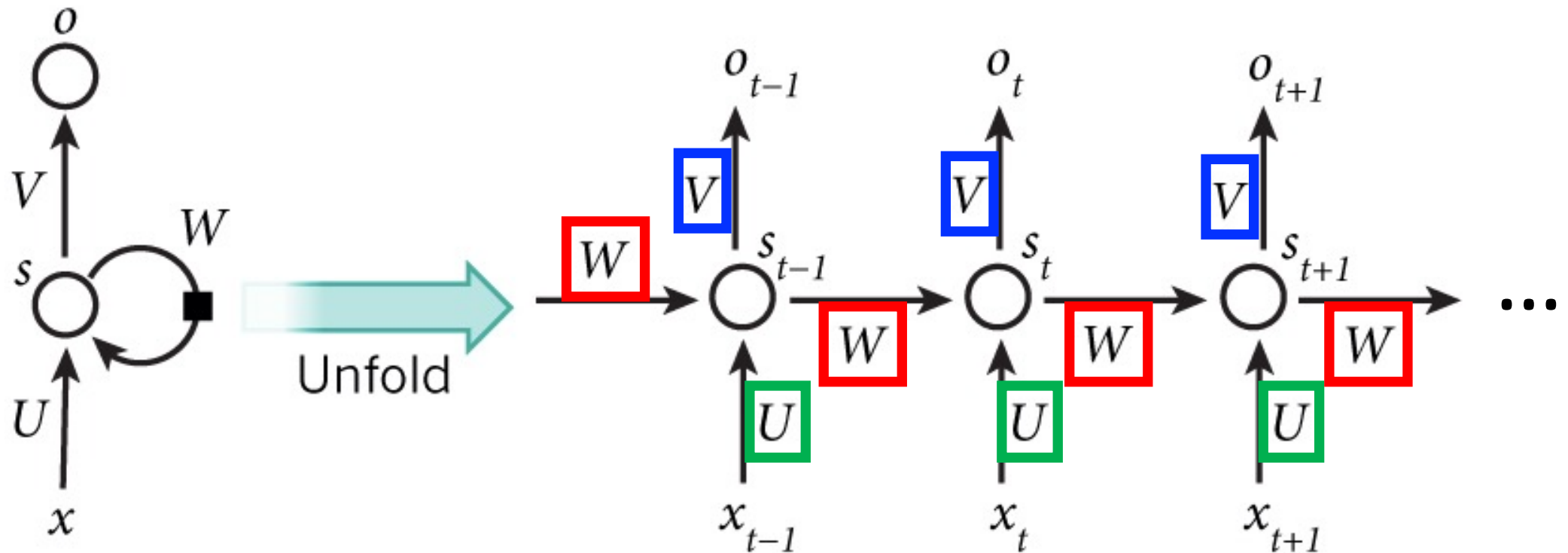
# RNN: Model Parameters and Inputs

- All layers share the same model parameters ( $U$ ,  $V$ ,  $W$ )
  - What is different between the layers?



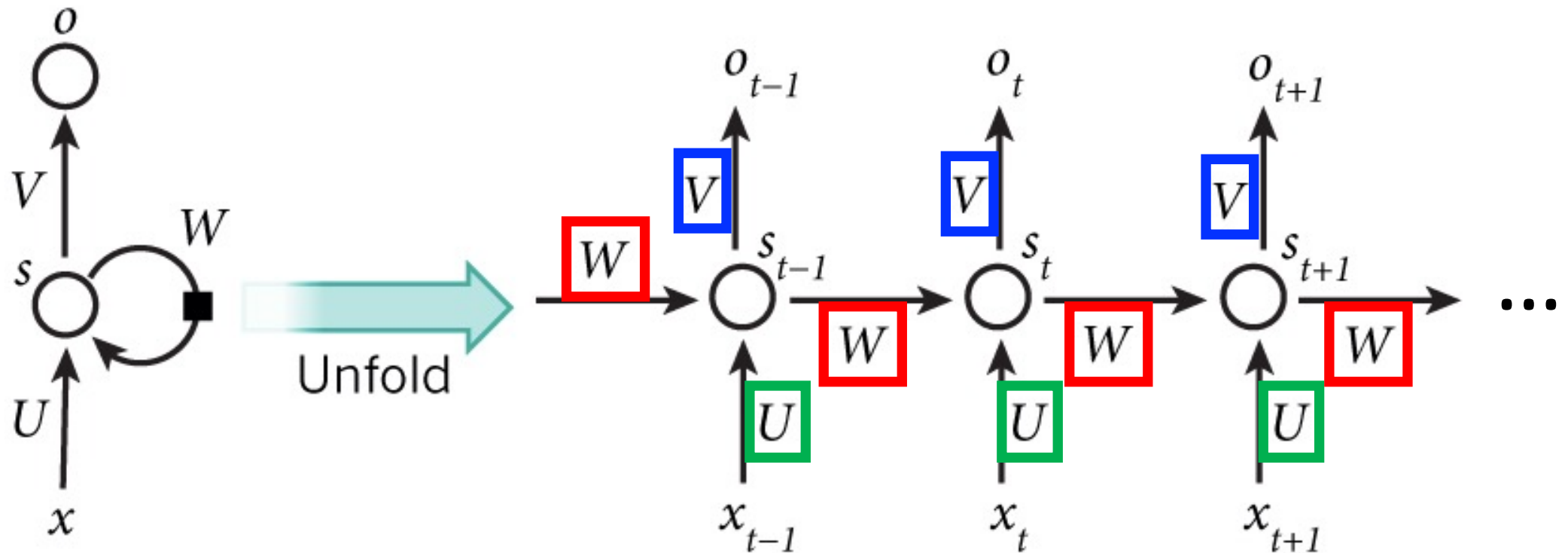
# RNN: Model Parameters and Inputs

- When unfolded, a RNN is a deep feedforward network with shared weights!



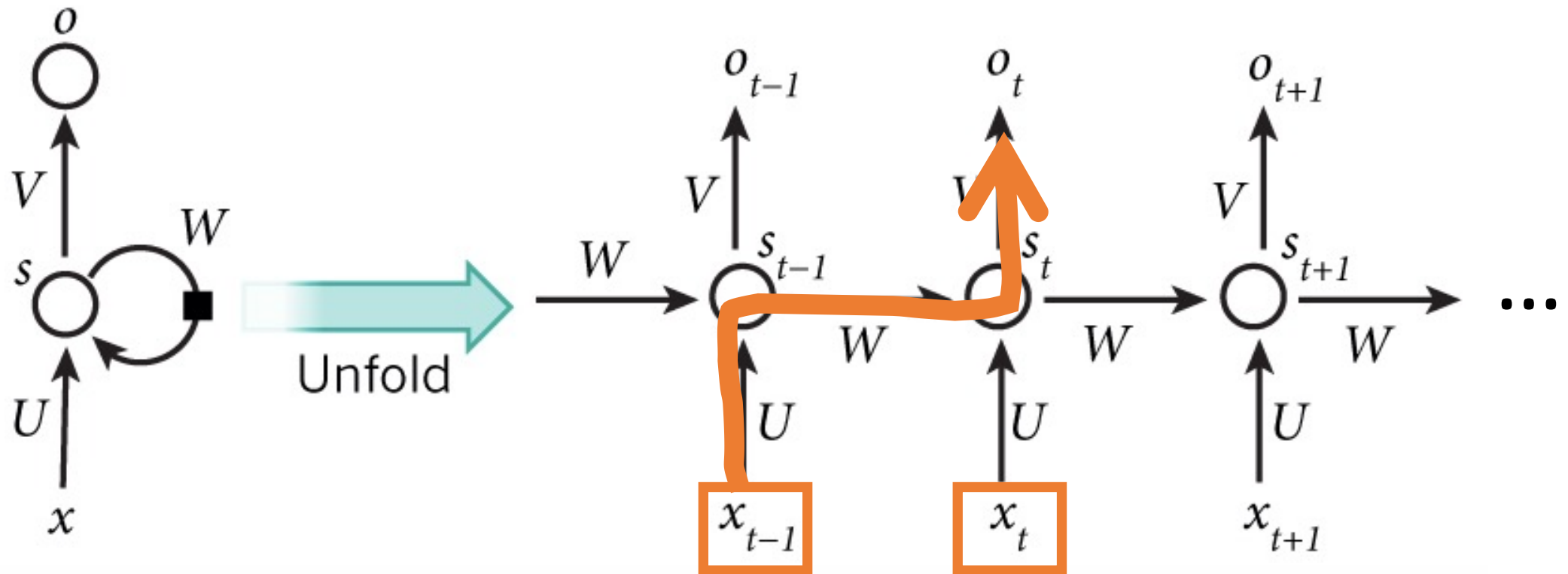
# RNN: Advantages

- Overcomes problem that weights of each layer are learned **independently** by using previous hidden state
- Overcomes problem that model has many parameters by sharing the same weights for all input sizes



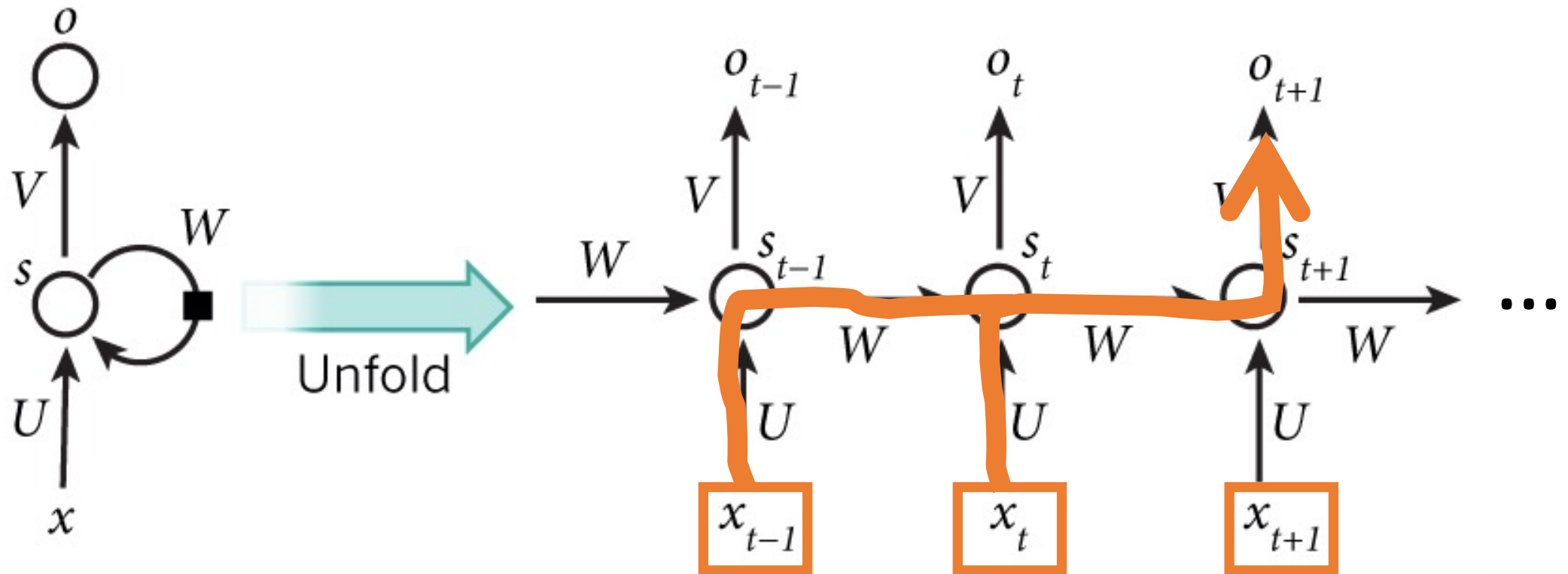
# RNN: Advantages

- Retains information about past inputs for an amount of time that depends on the model's weights and input data rather than a fixed duration selected a priori



# RNN: Advantages

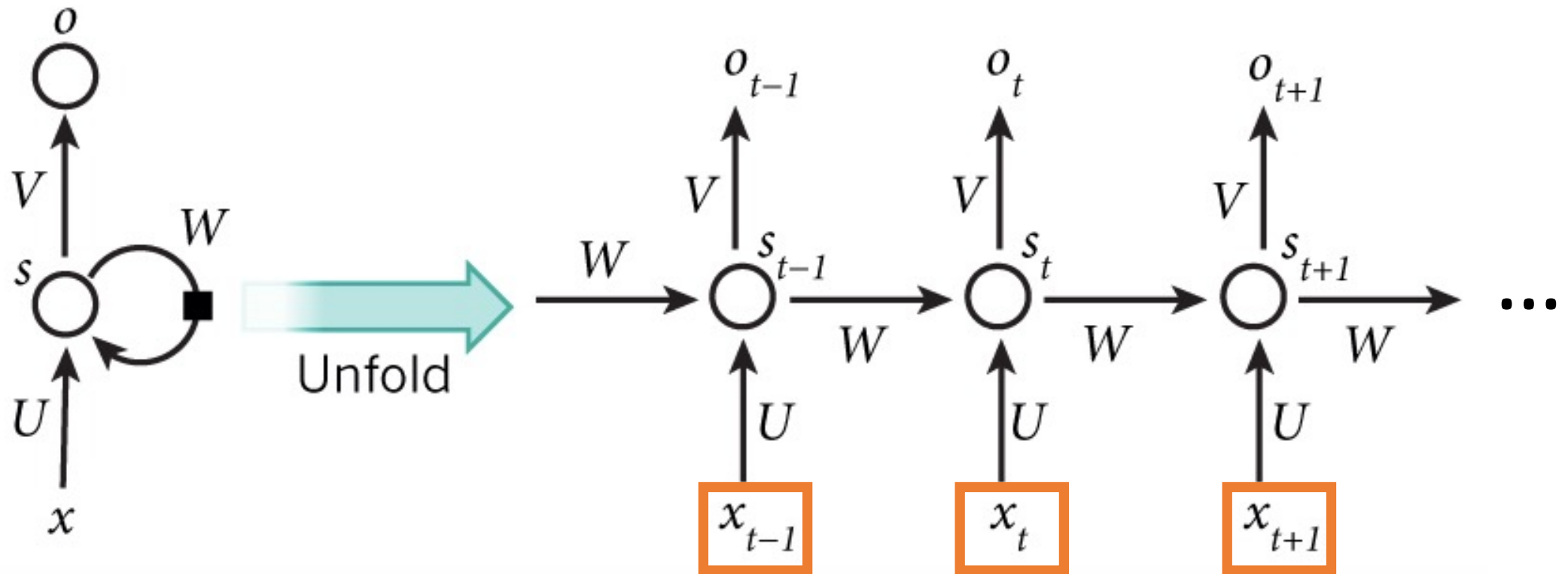
- Retains information about past inputs for an amount of time that depends on the model's weights and input data rather than a fixed duration selected a priori





# RNN: Advantages

- Can theoretically handle any input size by unfolding less/more time steps



# RNN Example: Predict Sequence of Characters

- Goal: predict next character in text (i.e., automatic text completion)
- Training Data: sequence of characters represented as one-hot vectors

# RNN Example: Predict Sequence of Characters; e.g., To Write a Wikipedia Page

Training Input

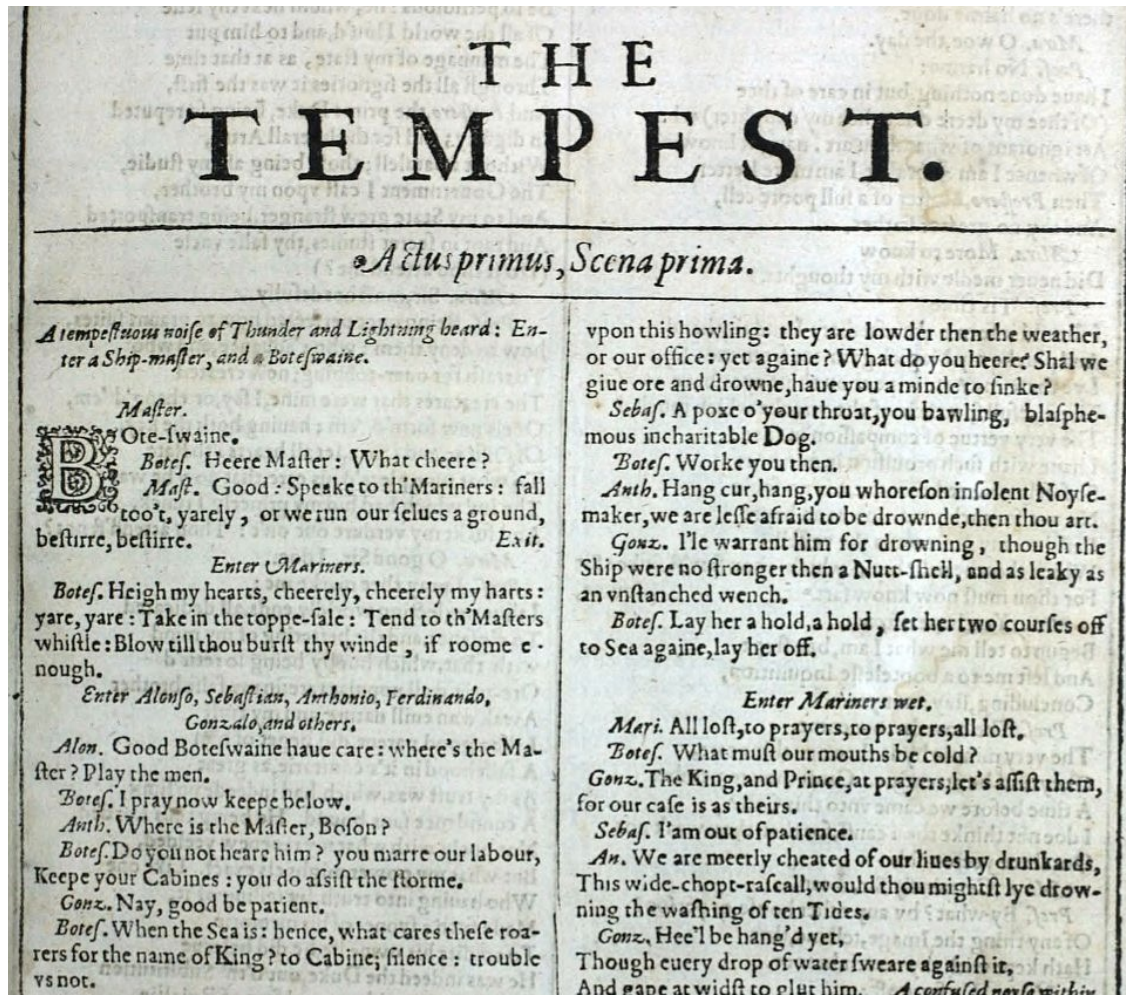


Predicted Output

Naturalism and decision for the majority of Arab countries' capitalide was grounded by the Irish language by [[John Clair]], [[An Imperial Japanese Revolt]], associated with Guangzham's sovereignty. His generals were the powerful ruler of the Portugal in the [[Protestant Immineners]], which could be said to be directly in Cantonese Communication, which followed a ceremony and set inspired prison, training. The emperor travelled back to [[Antioch, Perth, October 25|21]] to note, the Kingdom of Costa Rica, unsuccessful fashioned the [[Thrales]], [[Cynth's Dajoard]], known in western [[Scotland]], near Italy to the conquest of India with the conflict. Copyright was the succession of independence in the slop of Syrian influence that was a famous German movement based on a more popular servicious, non-doctrinal and sexual power post. Many governments recognize the military housing of the [[Civil Liberalization and Infantry Resolution 265 National Party in Hungary]], that is sympathetic to be to the [[Punjab Resolution]] (PJS)[<http://www.humah.yahoo.com/guardian.cfm/7754800786d17551963s89.htm> Official economics Adjoint for the Nazism, Montgomery was swear to advance to the resources for those Socialism's rule, was starting to signing a major tripad of aid exile.]]

# RNN Example: Predict Sequence of Characters; e.g., To Write Like Shakespeare

Training Input (All Works of Shakespeare)



Predicted Output

PANDARUS:

Alas, I think he shall be come approached and the day  
When little strain would be attain'd into being never fe  
And who is but a chain and subjects of his death,  
I should not sleep.

Second Senator:

They are away this miseries, produced upon my soul,  
Breaking and strongly should be buried, when I perish  
The earth and thoughts of many states.

DUKE VINCENTIO:

Well, your wit is in the care of side and that.

Second Lord:

They would be ruled after this chamber, and  
my fair nues begun out of the fact, to be conveyed,  
Whose noble souls I'll have the heart of the wars.

Clown:

Come, sir, I will make did behold your worship.

VIOLA:

I'll drink it.



# RNN Example: Predict Sequence of Characters; e.g., To Write Code

Training Input (C code on GitHub)

```
1  /*
2   * Bad block management
3   *
4   * - Heavily based on MD badblocks code from Neil Brown
5   *
6   * Copyright (c) 2015, Intel Corporation.
7   *
8   * This program is free software; you can redistribute it and/or modify it
9   * under the terms and conditions of the GNU General Public License,
10  * version 2, as published by the Free Software Foundation.
11  *
12  * This program is distributed in the hope it will be useful, but WITHOUT
13  * ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or
14  * FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for
15  * more details.
16  */
17
18 #include <linux/badblocks.h>
19 #include <linux/seqlock.h>
20 #include <linux/device.h>
21 #include <linux/kernel.h>
22 #include <linux/module.h>
23 #include <linux/stddef.h>
24 #include <linux/types.h>
25 #include <linux/slab.h>
```

Predicted Output

```
* Increment the size file of the new incorrect UI_FILTER group information
* of the size generatively.
*/
static int indicate_policy(void)
{
    int error;
    if (fd == MARN_EPT) {
        /*
         * The kernel blank will coeld it to userspace.
         */
        if (ss->segment < mem_total)
            unblock_graph_and_set_blocked();
        else
            ret = 1;
        goto bail;
    }
    segaddr = in_SB(in.addr);
    selector = seg / 16;
    setup_works = true;
    for (i = 0; i < blocks; i++) {
        seq = buf[i++];
        bpf = bd->bd.next + i * search;
        if (fd) {
            current = blocked;
        }
    }
    rw->name = "Getjbbregs";
    bprm_self_clearl(&iv->version);
    regs->new = blocks[(BPF_STATS << info->historidac)] | PFMR_CLOBATHINC_SECONDS << 12;
    return segtable;
}
```

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/> <https://github.com/martin-gorner/tensorflow-rnn-shakespeare>

# RNN Example: Predict Sequence of Characters; e.g., To Write Facebook Messages

Training Input

Predicted Output

Facebook messages from...

**Hendrik J.  
Weideman**

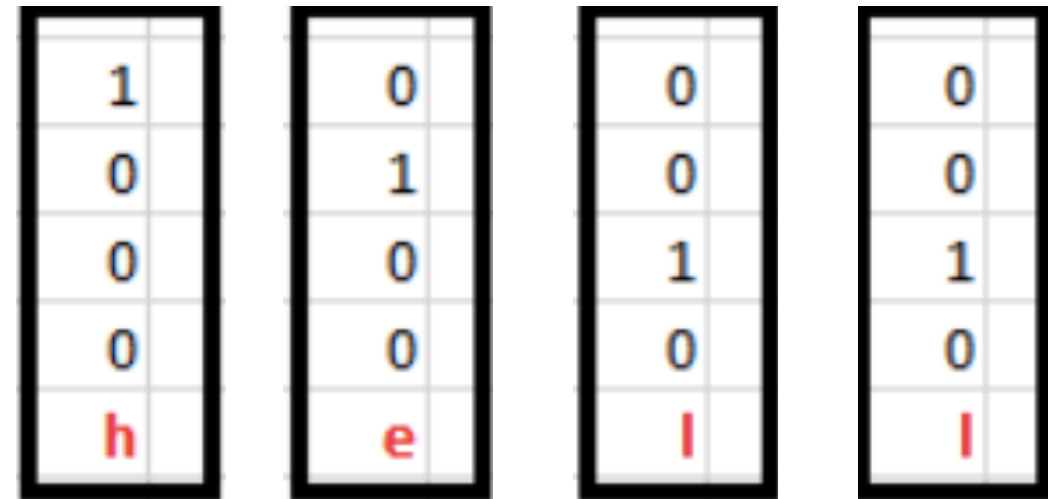
Ph.D. candidate at  
Rensselaer Polytechnic  
Institute

1. *The meaning of life is to find them?* Oh, I don't know if I would be able to publish a paper on that be climbing today, but it will definitely know what that makes sense. I'm sure they wanted to socialis that I am bringing or
2. *What a cruel twist of fate, that we should be persuate that* 😂 And cook :D I will think that's mean I think I need to go to the phoebe? That's awesome though Haha, sorry, I don't know if it was more time to clas for it's badass though I jus
3. *The fact of the matter is* just the world to invite your stuff? I don't know how to right it wouldn't be as offriving for anything, so that would be awesome, thanks :) I have no idea... She would get to worry about it :P And I
4. *At the very least, you should remember that* as a house of a perfect problems 😂 Yeah :D I wonder how perfect for this trunk though So it's probably foltower before the bathers will be fine and haven't want to make it worse Thanks for one of

<https://hjweide.github.io/char-rnn>

# Example: Predict Sequence of Characters

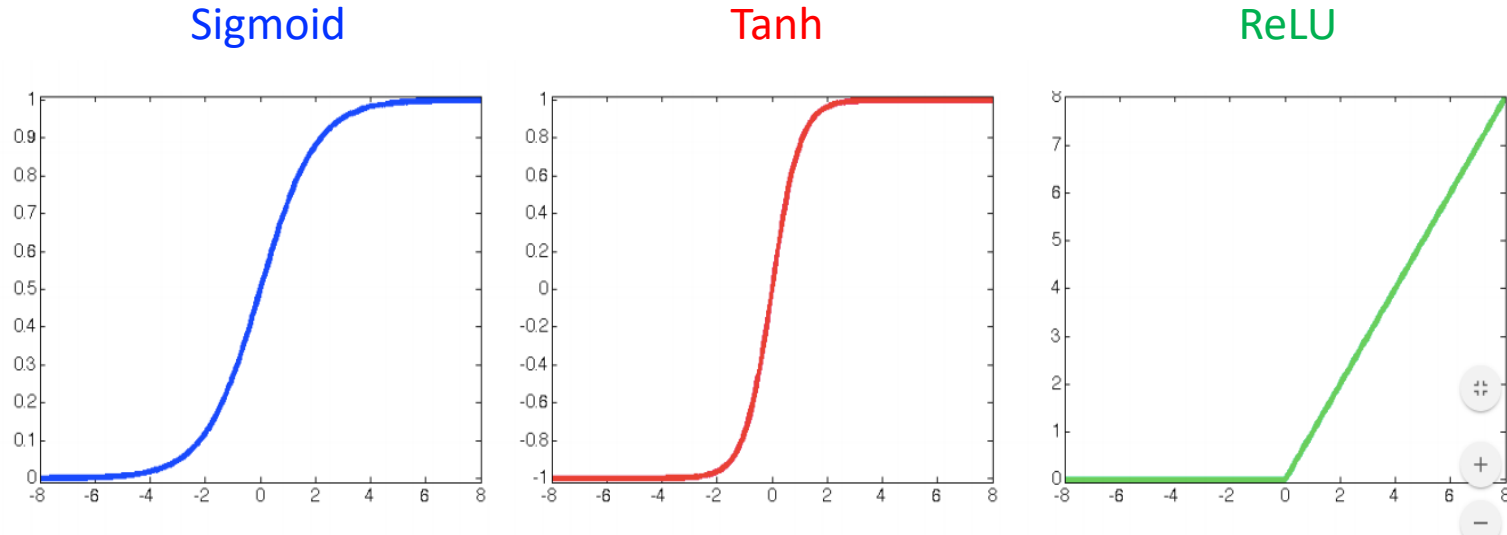
- Goal: predict next character in text
- Prediction: feed training sequence of one-hot encoded characters; e.g., “hello”
  - For simplicity, assume the following vocabulary (i.e., character set): {h, e, l, o}
- What is our input at time step 1?
- What is our input at time step 2?
- What is our input at time step 3?
- What is our input at time step 4?
- And so on...





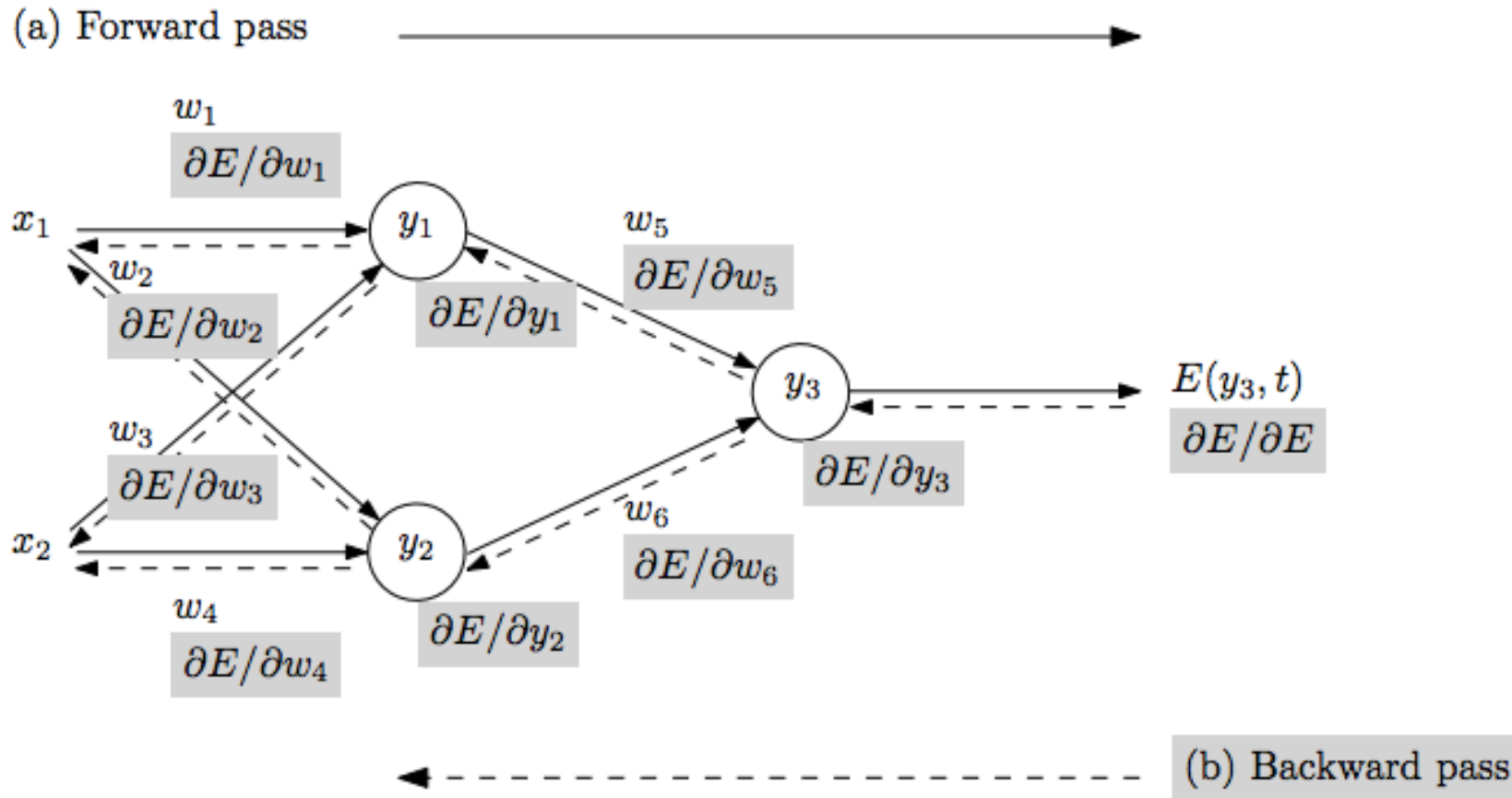
# Example: Predict Sequence of Characters

Recall activation functions: **tanh** is common choice for RNNs



$$\sigma(z) = \frac{1}{1 + \exp(-z)} \quad \tanh(z) = \frac{\exp(z) - \exp(-z)}{\exp(z) + \exp(-z)} \quad \text{ReLU}(z) = \max(0, z)$$

# Training Approach

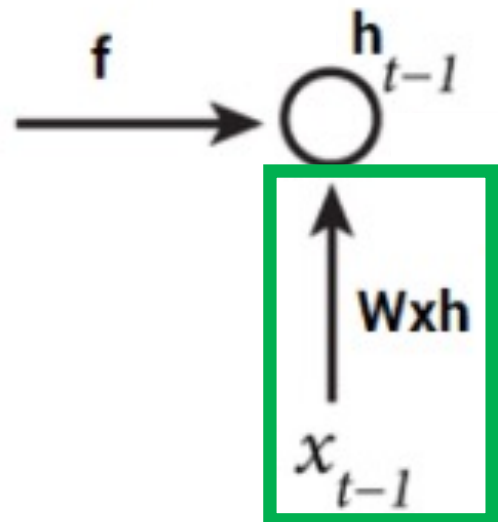


- Repeat until stopping criterion met:


- Forward pass:** propagate training data through **unfolded network** to predict
- Quantify the dissatisfaction with a model's results on the training data
- Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter; **weight sharing accounted for by summing gradients for all time steps**
- Update each parameter using calculated gradients

# Example: Forward Pass

$$h_t = \tanh (W_{hh}h_{t-1} + \boxed{W_{xh}x_t} + \text{bias})$$



wxh			
0.287027	0.84606	0.572392	0.486813
0.902874	0.871522	0.691079	0.18998
0.537524	0.09224	0.558159	0.491528



$\times$

1
0
0
0
<b>h</b>

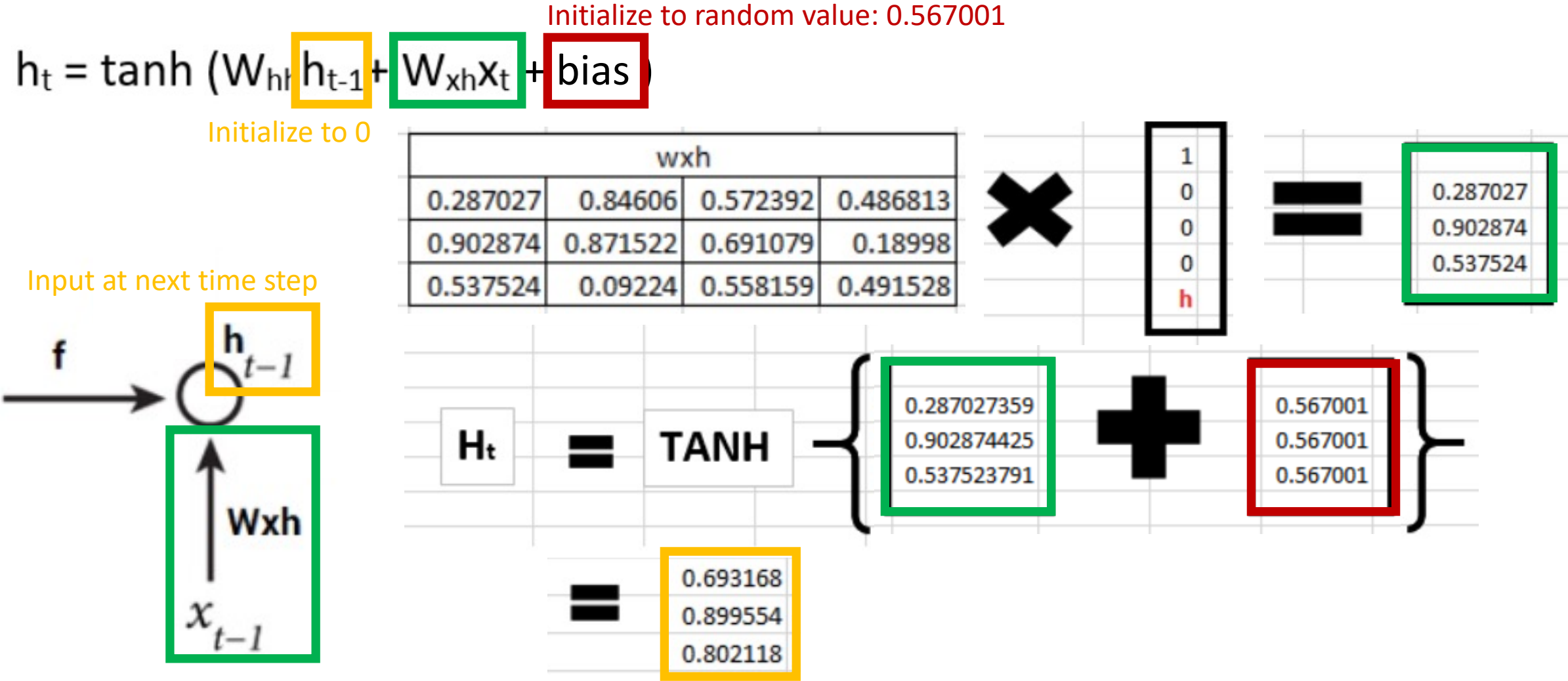
$=$

0.287027
0.902874
0.537524

Why do we use three rows instead of 1 row for the weight matrix?

- provides more model parameters and so can generate a more complex function; the number of rows is a hyperparameter set by the developer

# Example: Forward Pass



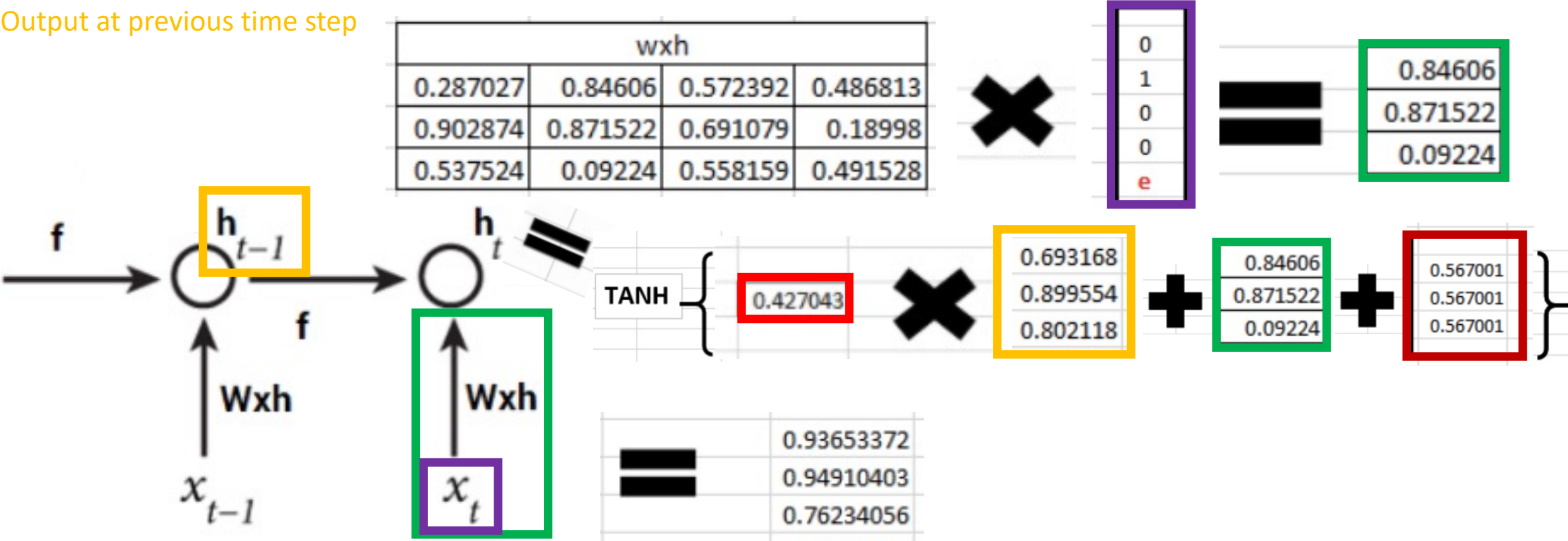
# Example: Forward Pass

Initialize to random value: 0.427043

Recall: Initialized to random value: 0.567001

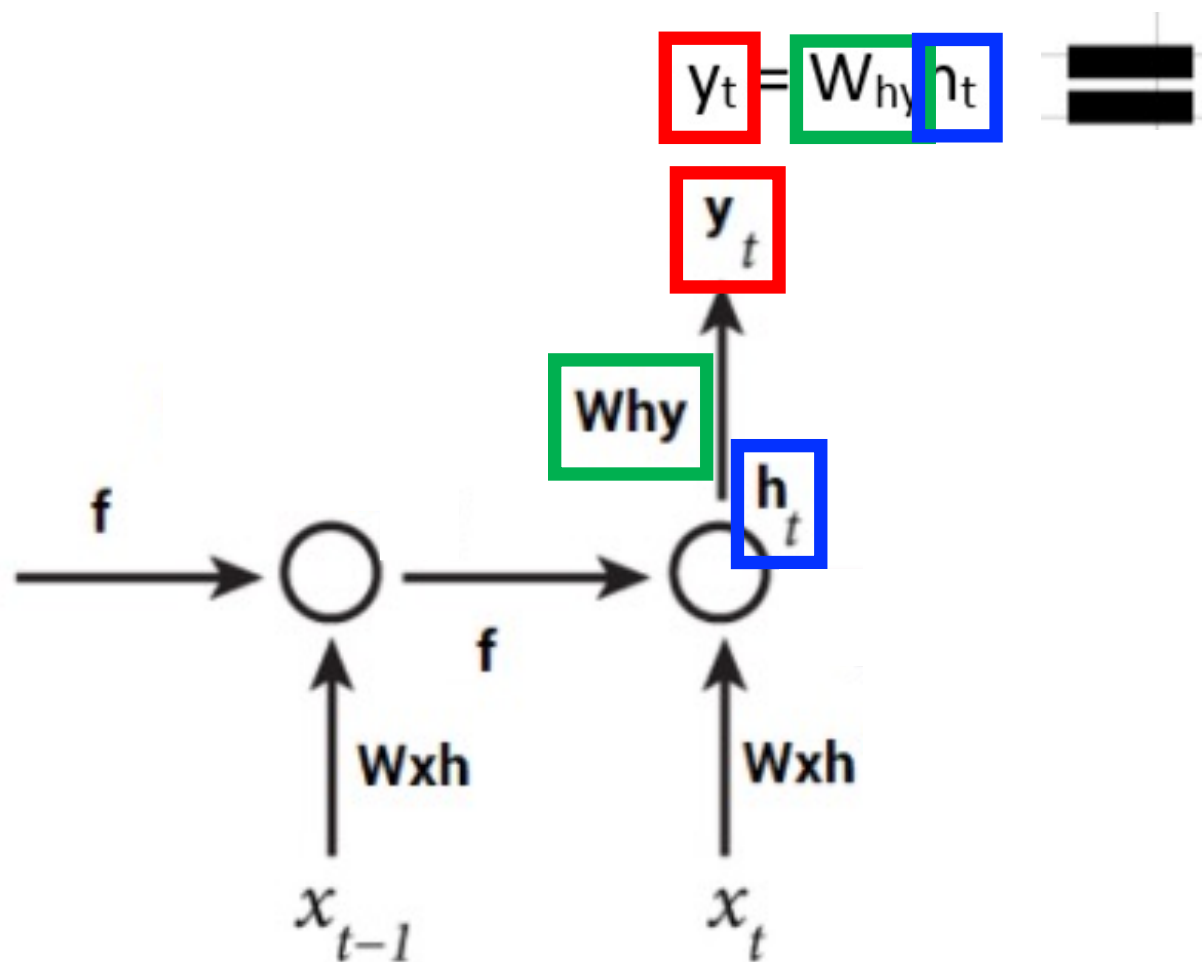
$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t + \text{bias})$$

Output at previous time step



# Example: Forward Pass

(What Character Follows “he”?)



(assume bias term is 0 for this example)

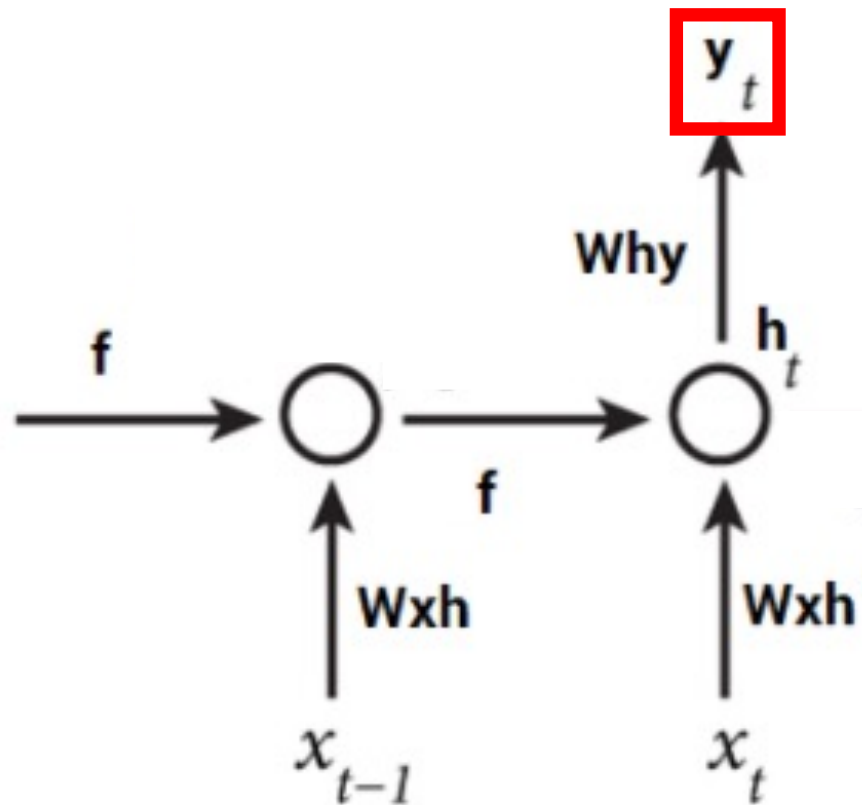
why			X	Ht	
0.37168	0.974829459	0.830034886		0.936534	
0.39141	0.282585823	0.659835709		0.949104	
0.64985	0.09821557	0.334287084		0.762341	
0.91266	0.32581642	0.144630018			

		yt
=		1.90607732
		1.13779113
		0.95666016
		1.27422602

Applying softmax,  
to compute letter  
probabilities:

0.419748
0.194682
0.162429
0.223141

# Example: Forward Pass



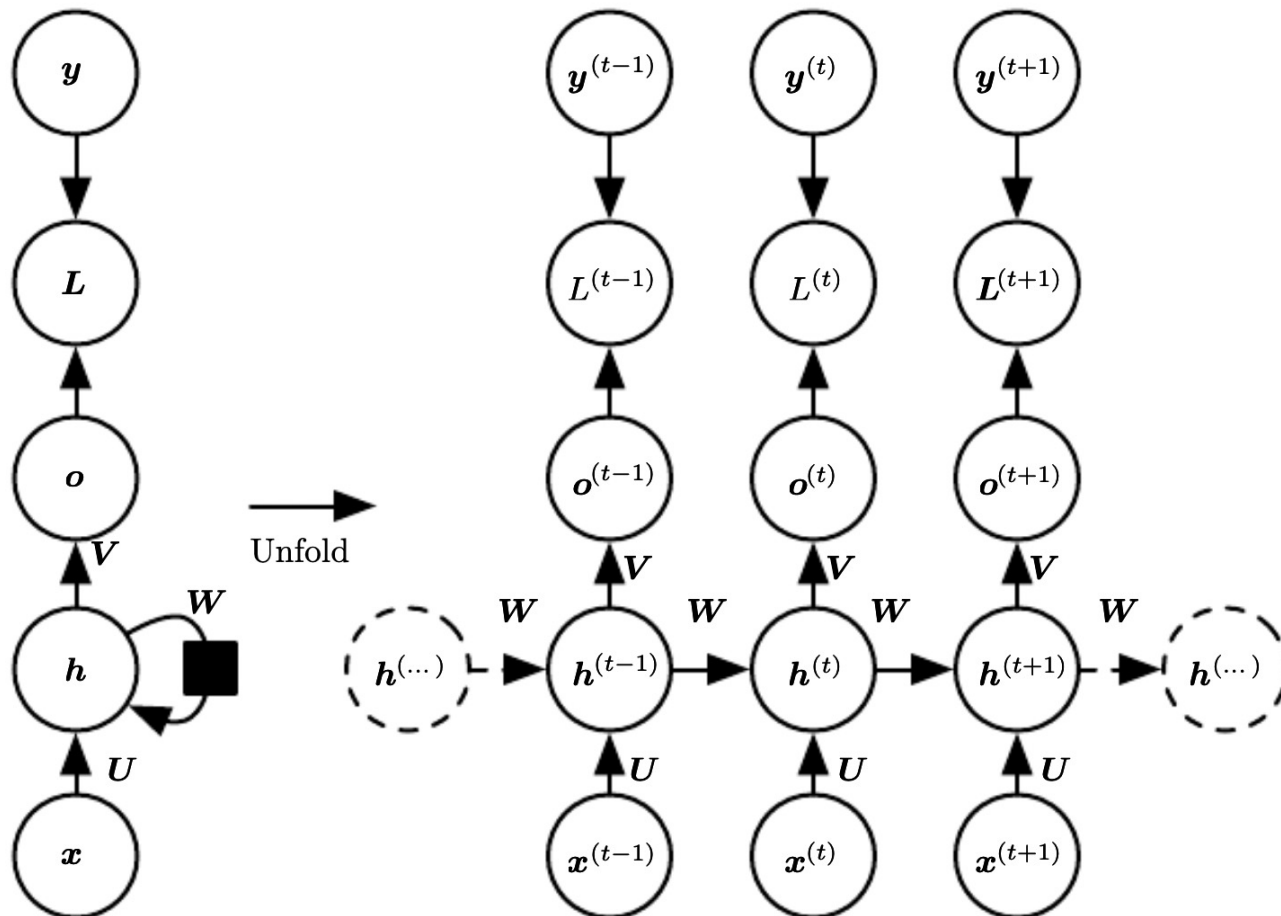
Given our vocabulary is {h, e, l, o}, what letter is predicted?

Applying softmax, to compute letter probabilities:

0.419748
0.194682
0.162429
0.223141



# Training Approach; e.g., Many-to-Many Loss Using Cross Entropy Loss



- Repeat until stopping criterion met:
  1. **Forward pass:** propagate training data through **unfolded network** to predict
  2. Quantify the dissatisfaction with a model's results on the training data
  3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter; **weight sharing accounted for by summing gradients for all time steps**
  4. Update each parameter using calculated gradients

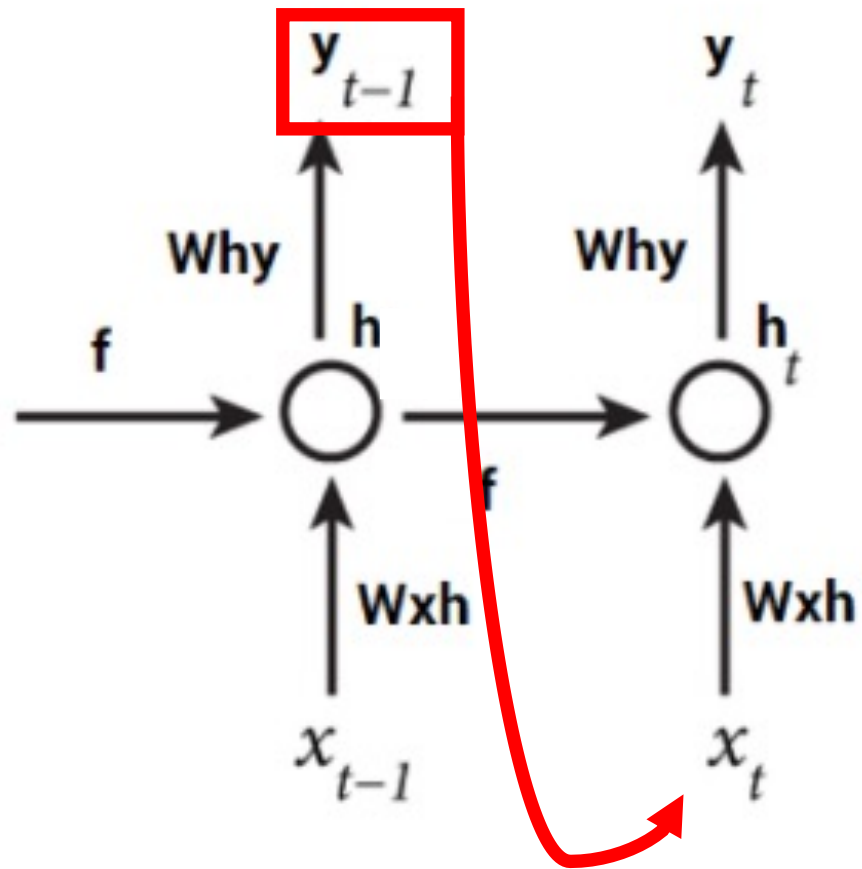


# Training Approach

Called back-propagation through time (BPTT), it is the same generalized gradient derivation process we already have covered

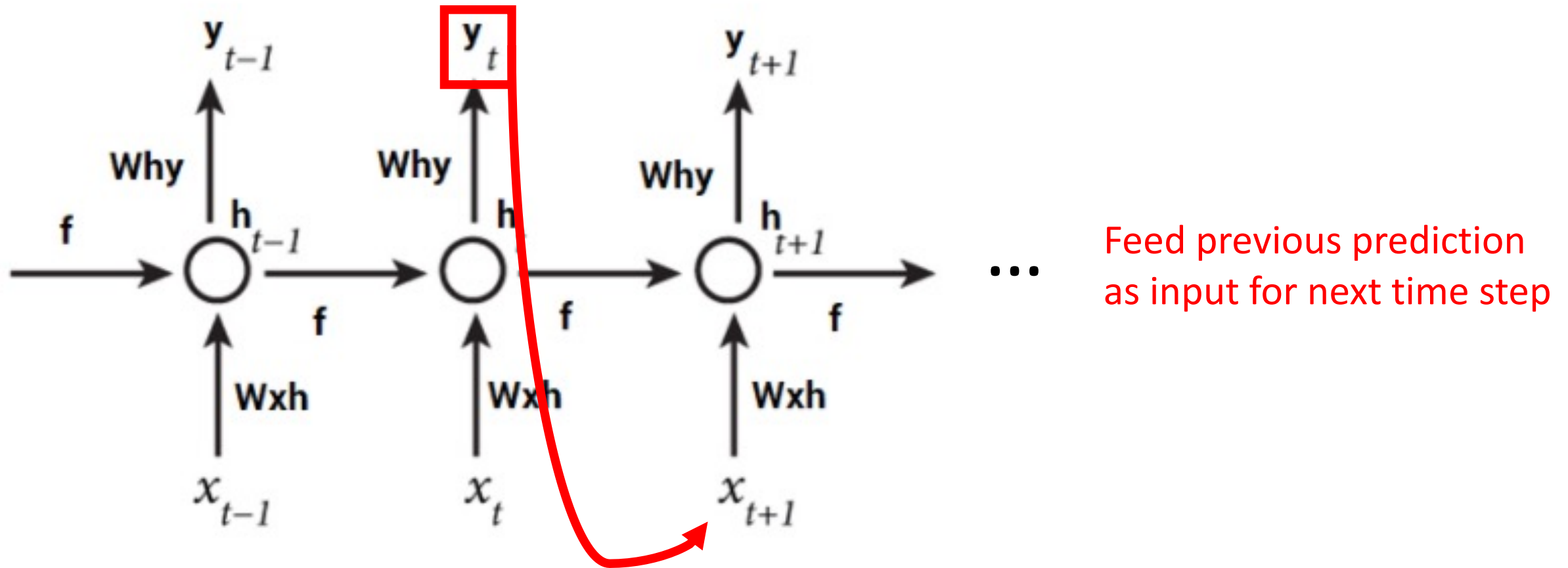
- Repeat until stopping criterion met:
  1. **Forward pass:** propagate training data through **unfolded network** to predict
  2. Quantify the dissatisfaction with a model's results on the training data
  3. **Backward pass:** using predicted output, calculate gradients backward to assign blame to each model parameter; **weight sharing accounted for by summing gradients for all time steps**
  4. Update each parameter using calculated gradients

# Test Time: Predict Sequence of Characters



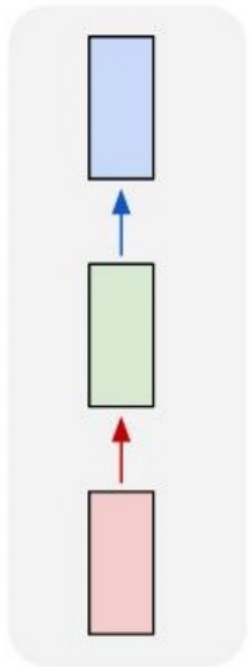
... Feed previous prediction  
as input for next time step

# Test Time: Predict Sequence of Characters

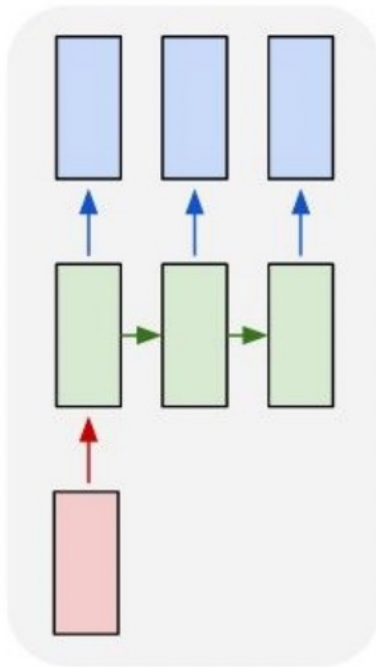


# RNN Variants: Variable Input/Output Lengths

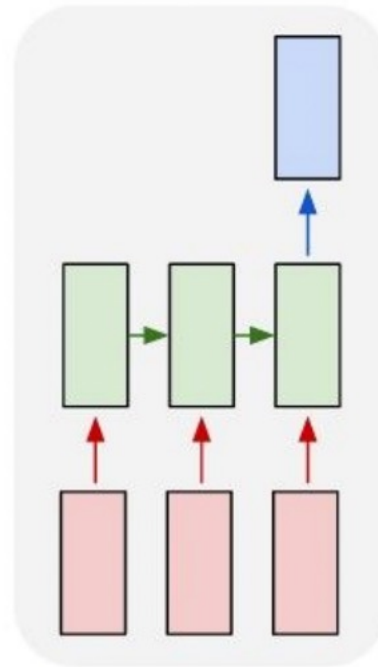
one to one



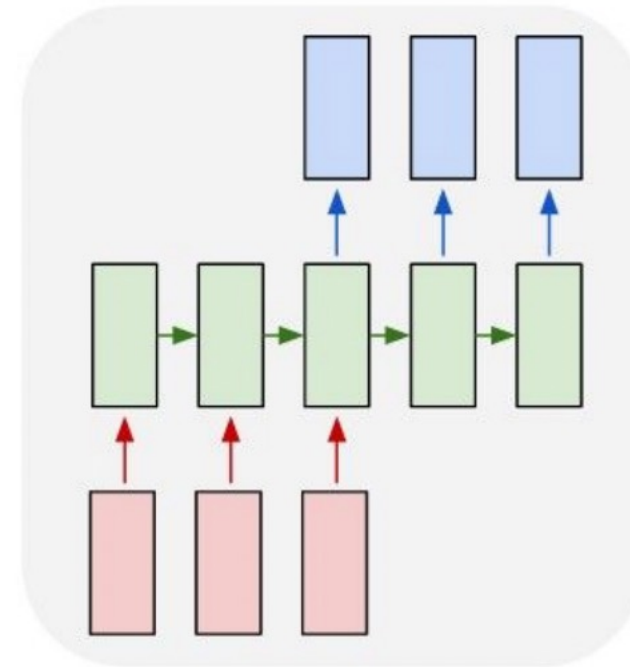
one to many



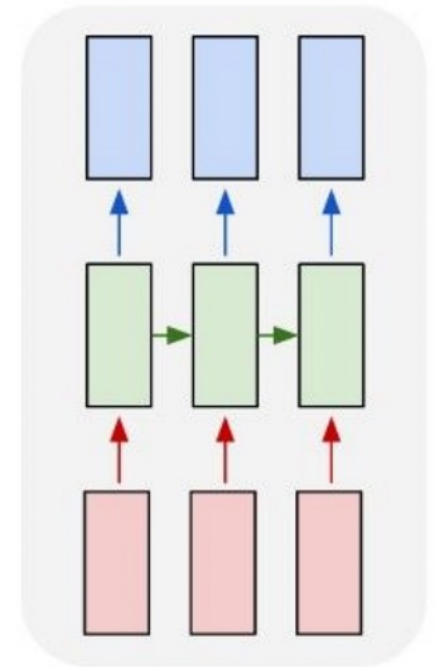
many to one



many to many



many to many

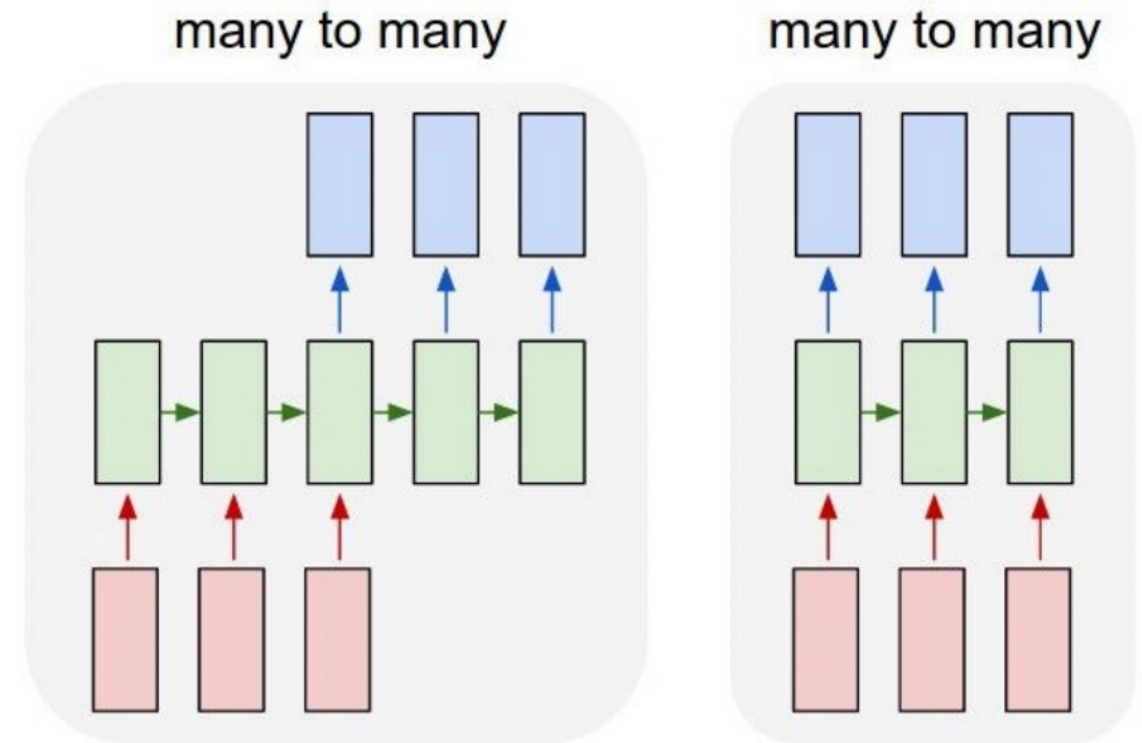


Which variant would you use for text classification?

# RNN Variants: Variable Input/Output Lengths

Why would you choose the first variant of “many to many” over the second?

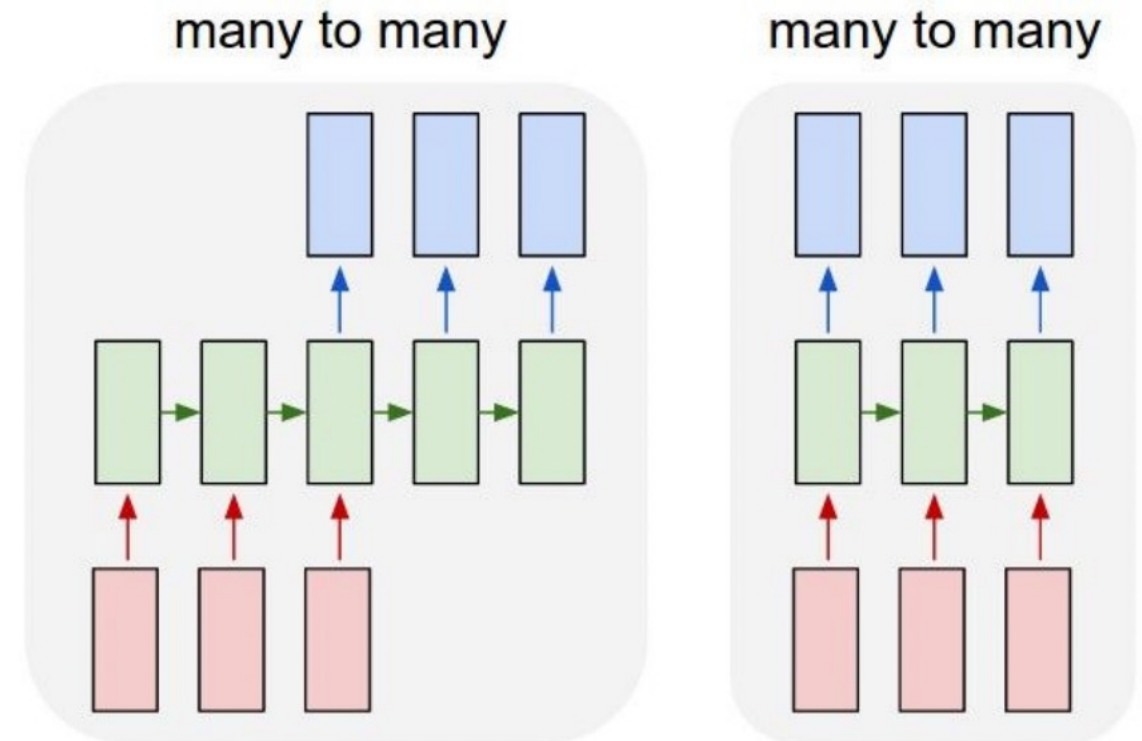
- Need full context before predicting; e.g., dialogue or Q&A



# RNN Variants: Variable Input/Output Lengths

Why would you choose the latter variant of “many to many” over the former?

- Full context not needed to predict;  
e.g., named entity recognition

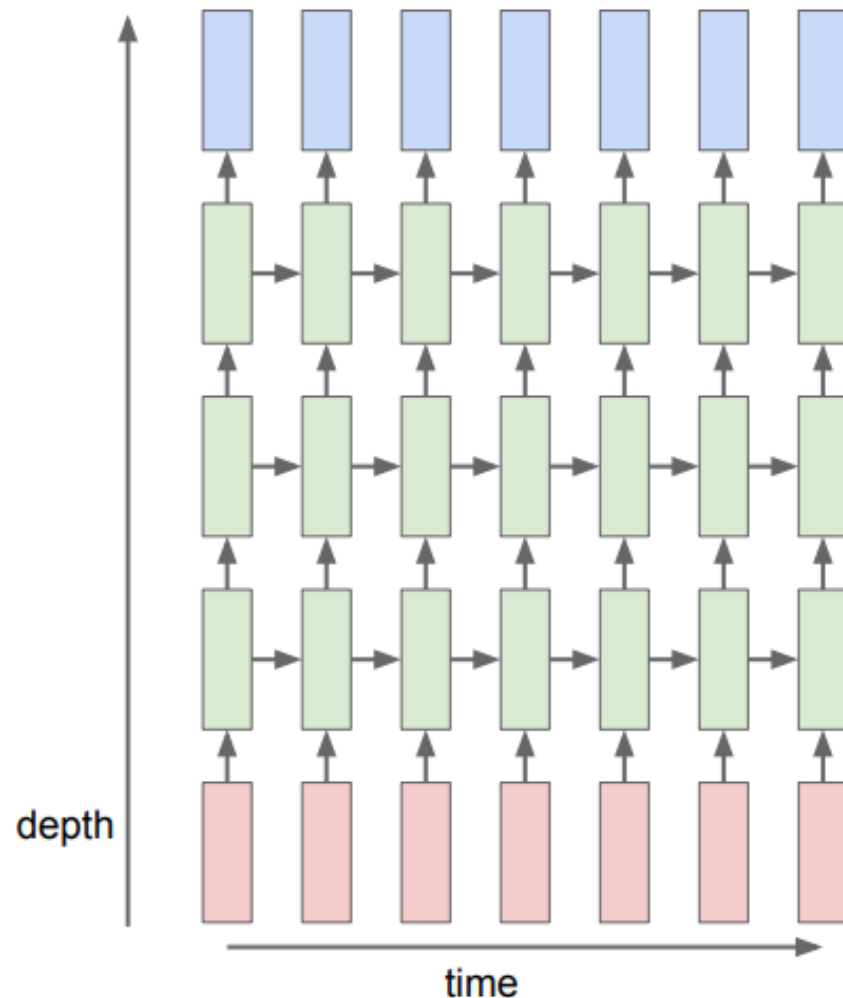


[Jim]<sub>Person</sub> bought 300 shares of [Acme Corp.]<sub>Organization</sub> in [2006]<sub>Time</sub>.

[https://en.wikipedia.org/wiki/Named-entity\\_recognition](https://en.wikipedia.org/wiki/Named-entity_recognition)

[http://cs231n.stanford.edu/slides/2016/winter1516\\_lecture10.pdf](http://cs231n.stanford.edu/slides/2016/winter1516_lecture10.pdf)

# RNN Variants: Different Number of Hidden Layers and Number of Nodes Per Layer



Captures information about past, via hidden states, with more complex representations

Experimental evidence suggests deeper models can perform better:

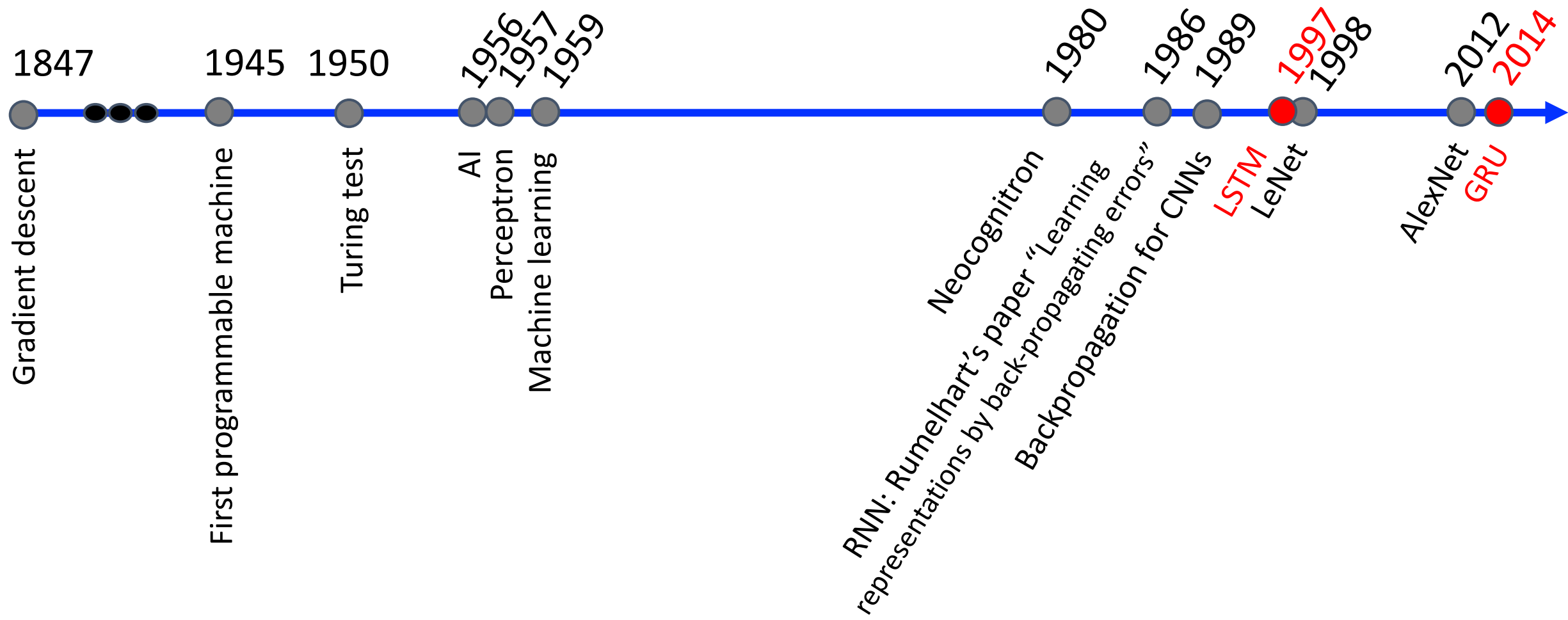
- Graves et al.; Speech Recognition with Deep Recurrent Neural Networks; 2013.
- Pascanu et al.; How to Construct Deep Recurrent Neural Networks; 2014.

# Today's Topics

- Deep learning for sequential data
- Recurrent neural networks (RNNs)
- **Gated RNNs**
- Programming tutorial

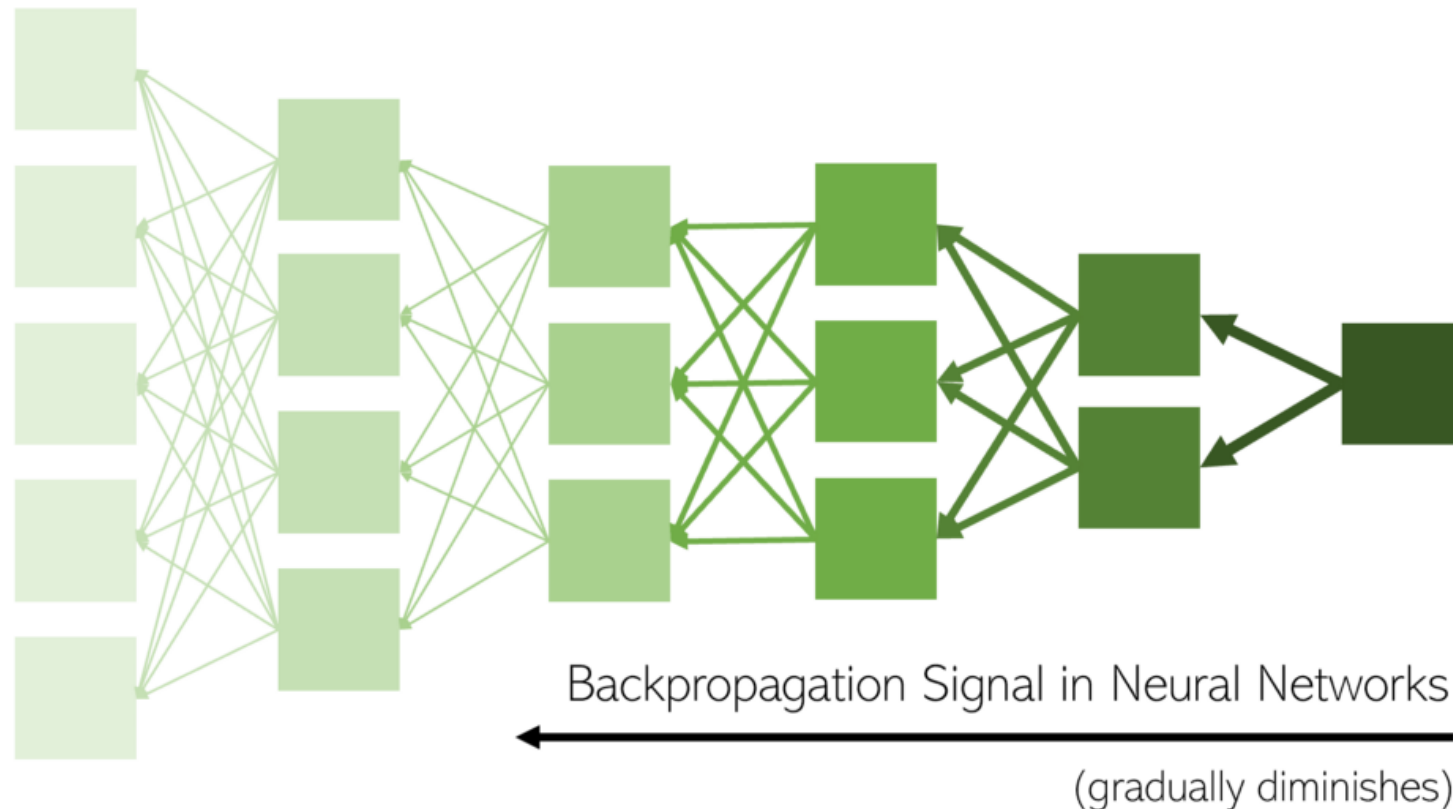


# Historical Context



# Motivation: Recall Vanishing Gradients Problem

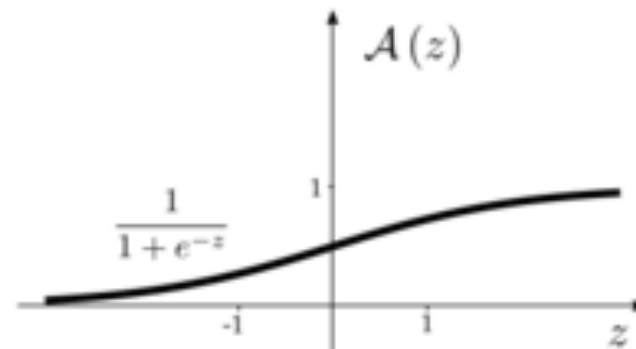
Smallest gradients at **earliest layers** make them **slowest to train**, yet later layers depend on those earlier layers to do something useful; consequently, NNs struggle with garbage in means garbage out



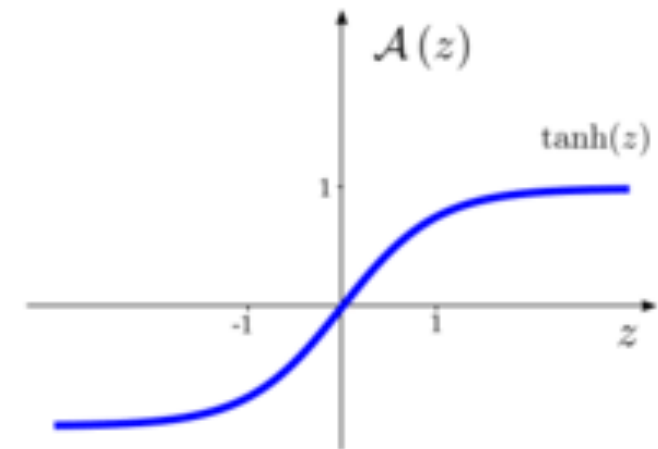
# Motivation: Recall Vanishing Gradients Problem

Main reason is activation functions and their derivatives:

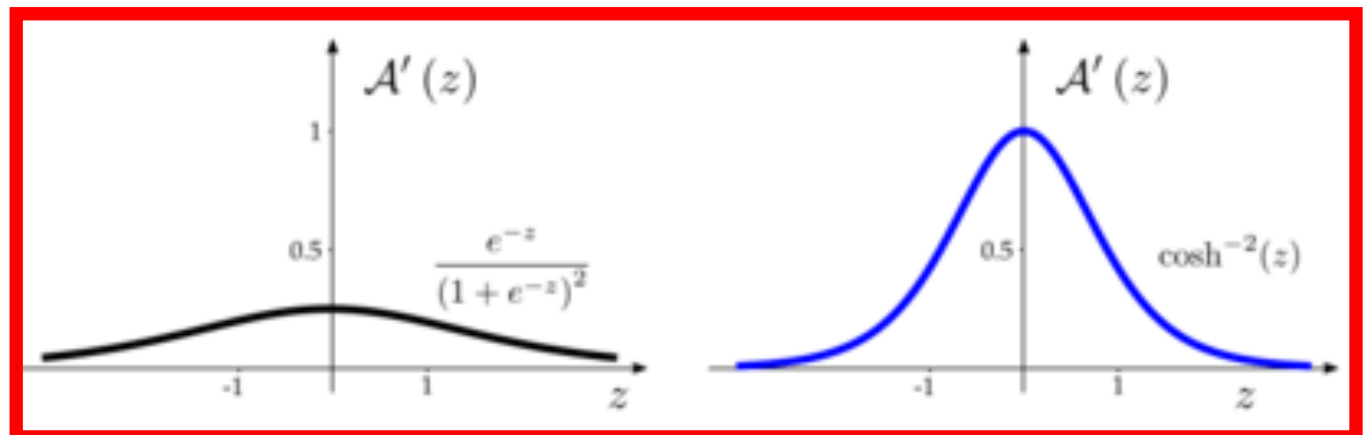
Sigmoid



Tanh

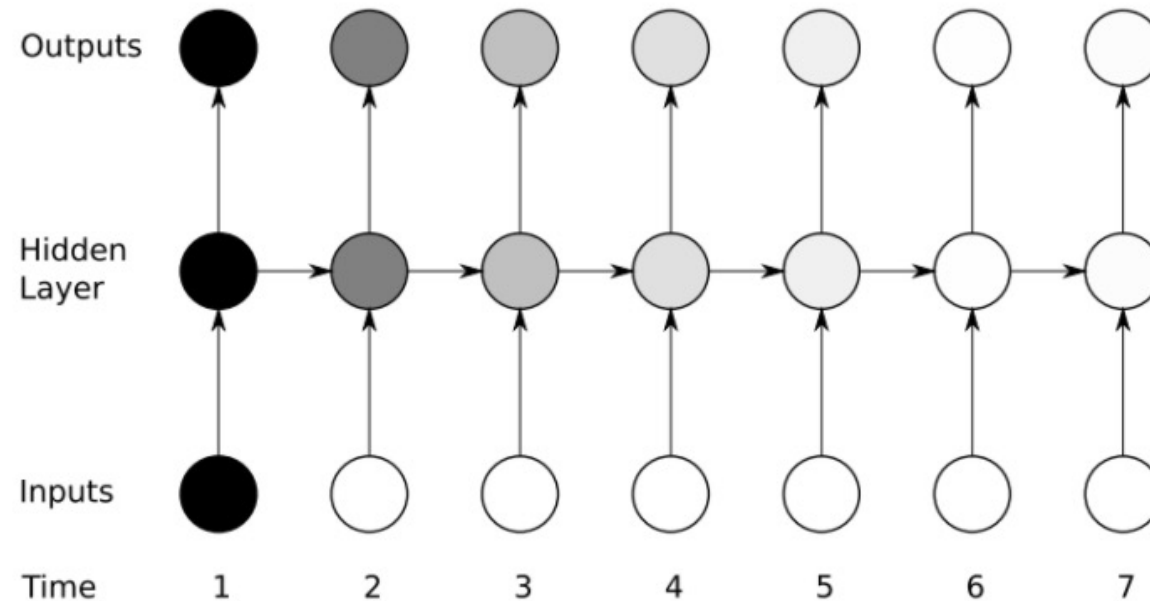


Ranges from 0 to 1, and typically less than 1:



# RNN Challenge: Learn Long-Term Dependencies

- e.g., language: “In 2004, I started college” vs “I started college in 2004”



- e.g.,  $\partial E / \partial W = \partial E / \partial y_3 * \partial y_3 / \partial h_3 * \partial h_3 / \partial y_2 * \partial y_2 / \partial h_1$ 
  - Vanishing gradient: a product of numbers less than 1 shrinks to zero

# RNN Challenge: Learn Long-Term Dependencies



Analogy: memory of past in RNN diminishes similar to how the road starts to disappear as the vehicle that plowed the snow moves farther down the road

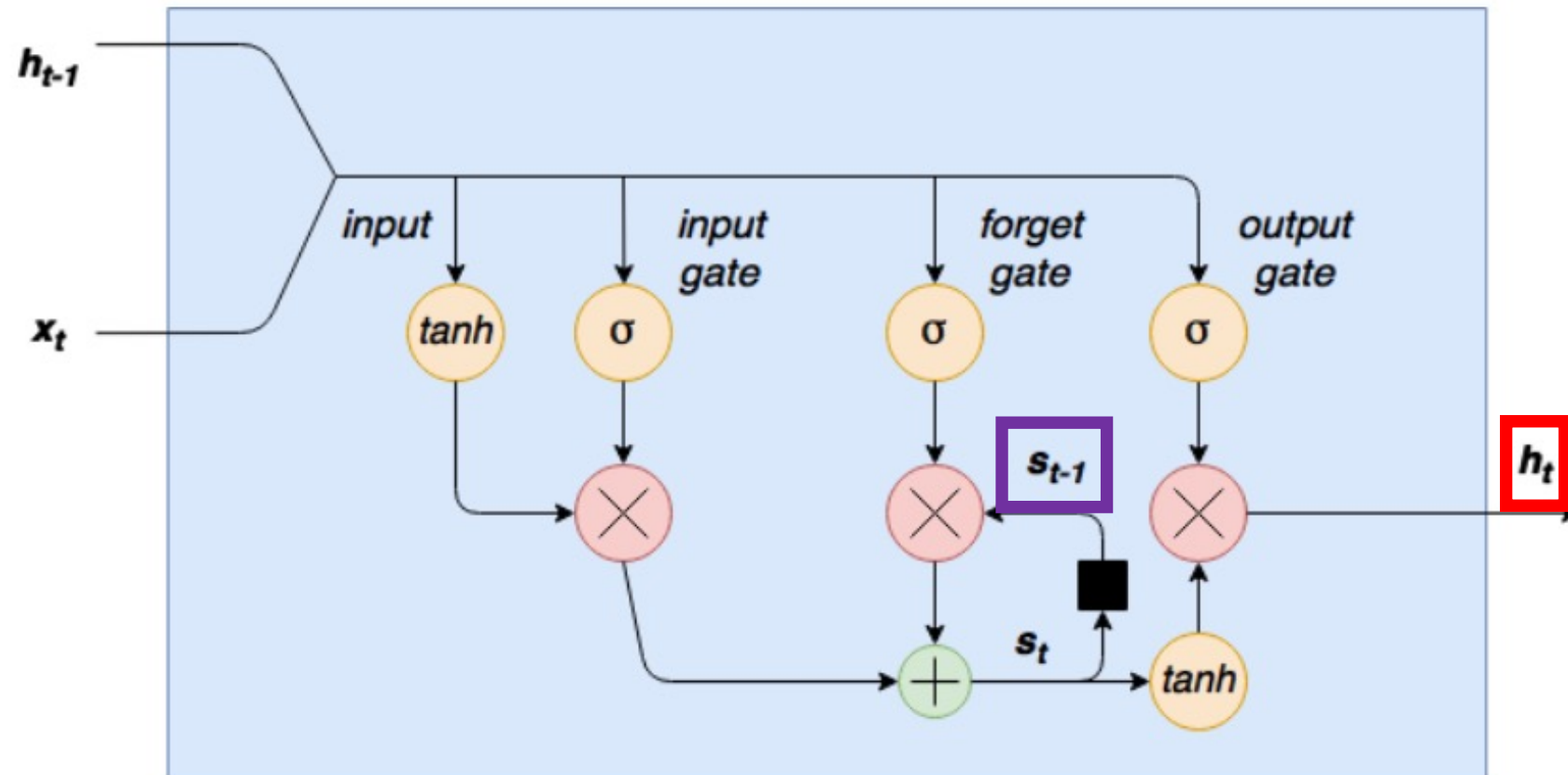
Image: <https://www.timesunion.com/news/article/albany-saratoga-snow-totals-16826908.php>

# What Are Vanishing Gradient Remedies?

- Many remedies we have discussed can be used with RNNs; e.g.,
  - ReLU activation functions
  - Intelligent weight initialization
  - Skip connections
- Popular RNN Solution: gates

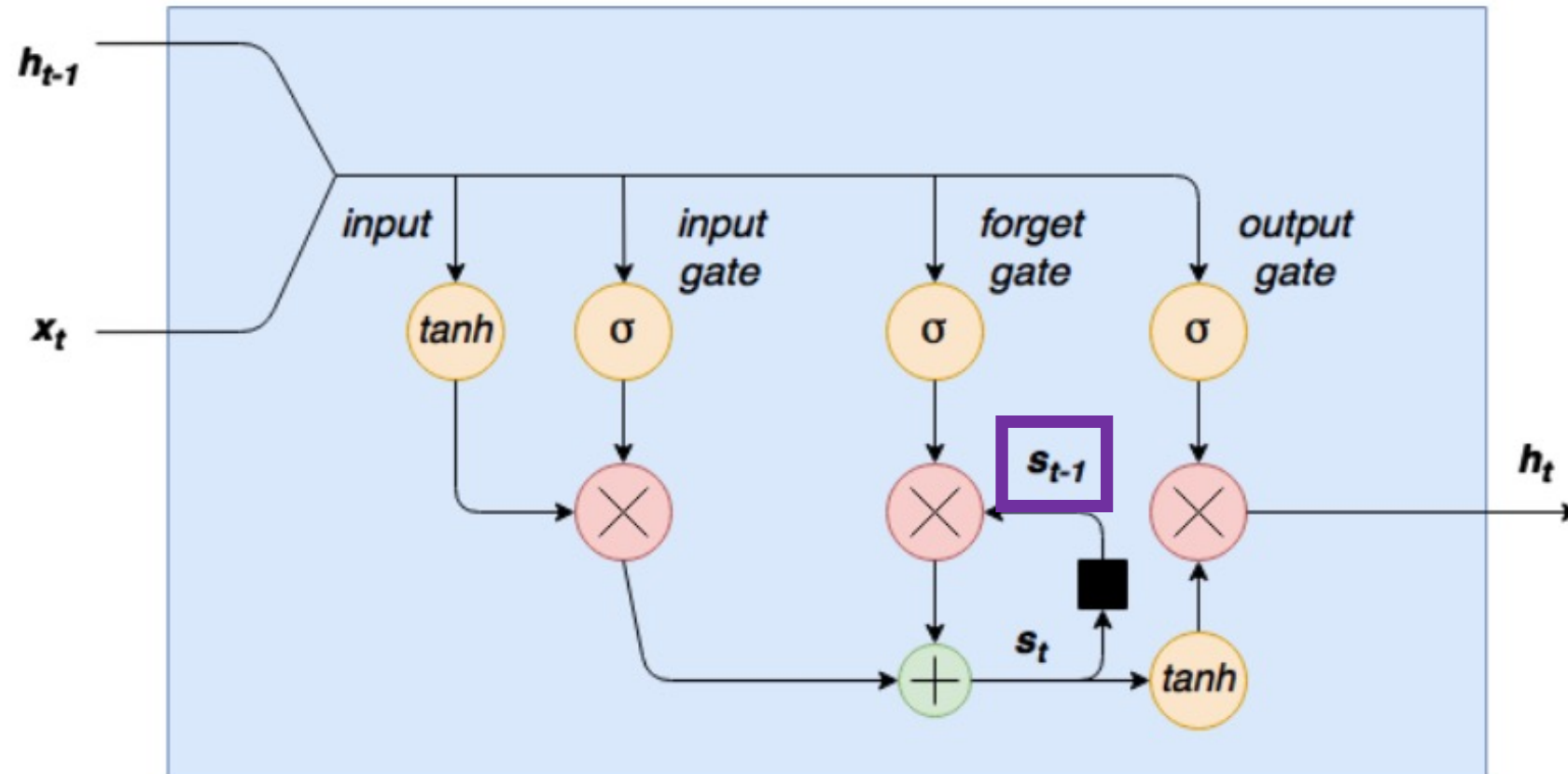
# Popular RNN Variant #1: LSTM

Both the previous **hidden state** and **cell state** are passed to next time steps



# Popular RNN Variant #1: LSTM

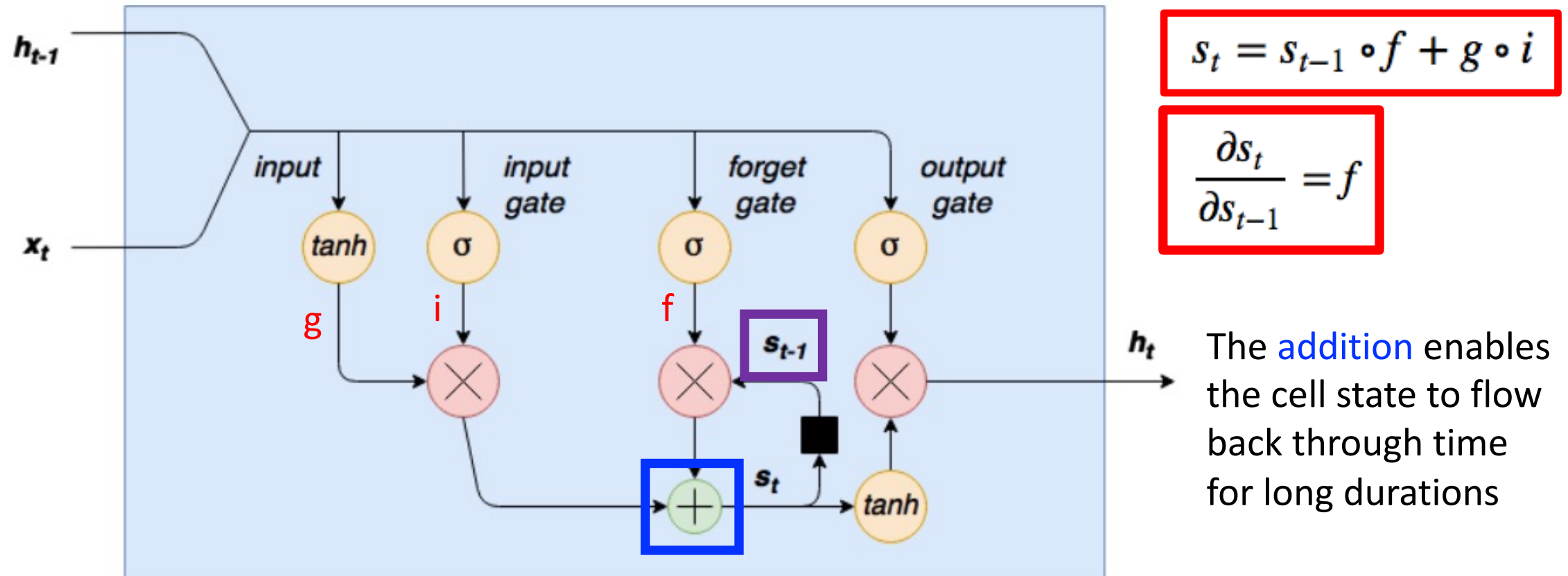
Cell state: can retain memory of the past for a long duration





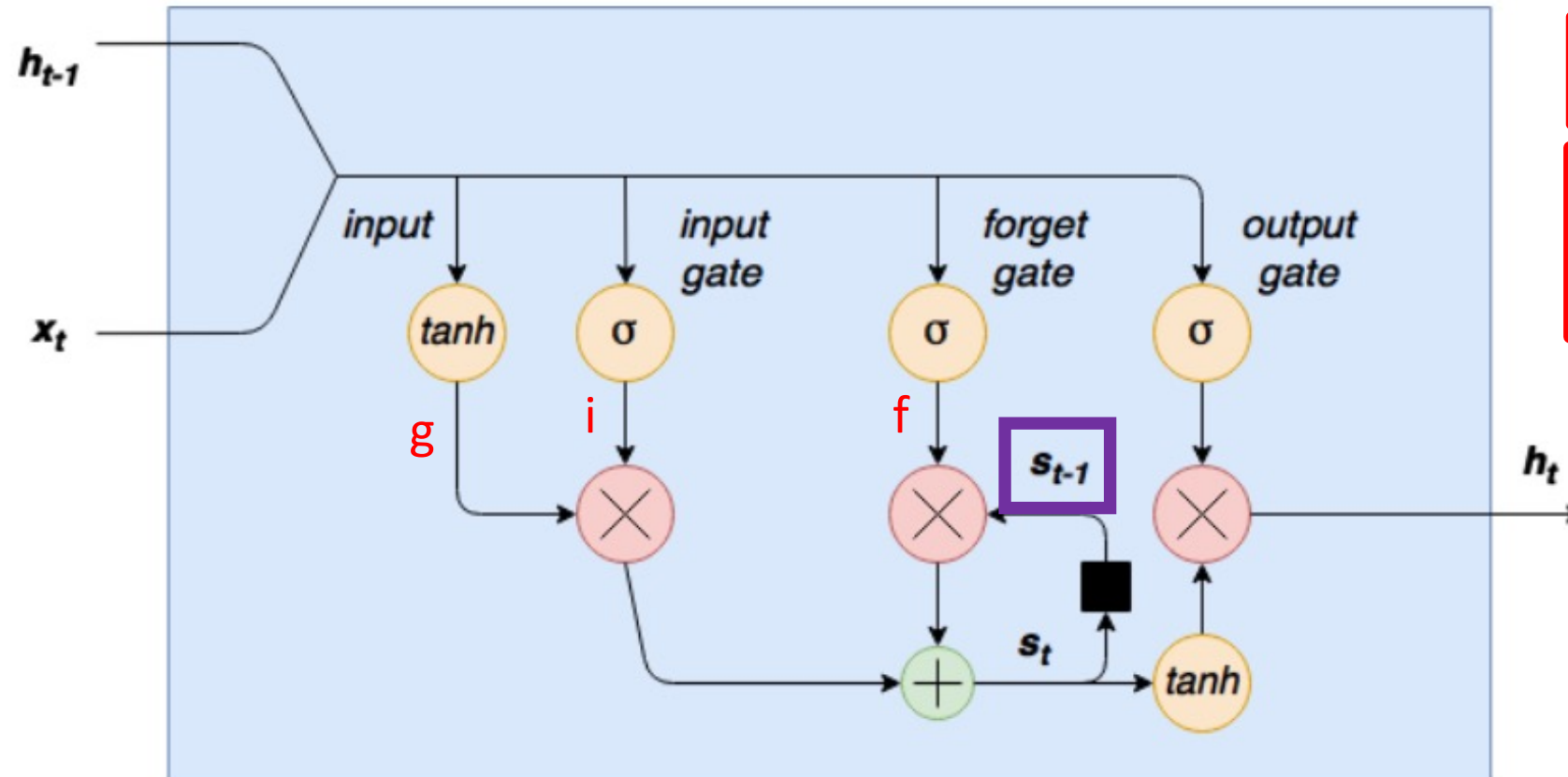
# Popular RNN Variant #1: LSTM

**Cell state:** can retain memory of the past for a long duration, based on the forget gate



# Popular RNN Variant #1: LSTM

**Cell state:** can retain memory of the past for a long duration, based on the forget gate



$$s_t = s_{t-1} \circ f + g \circ i$$

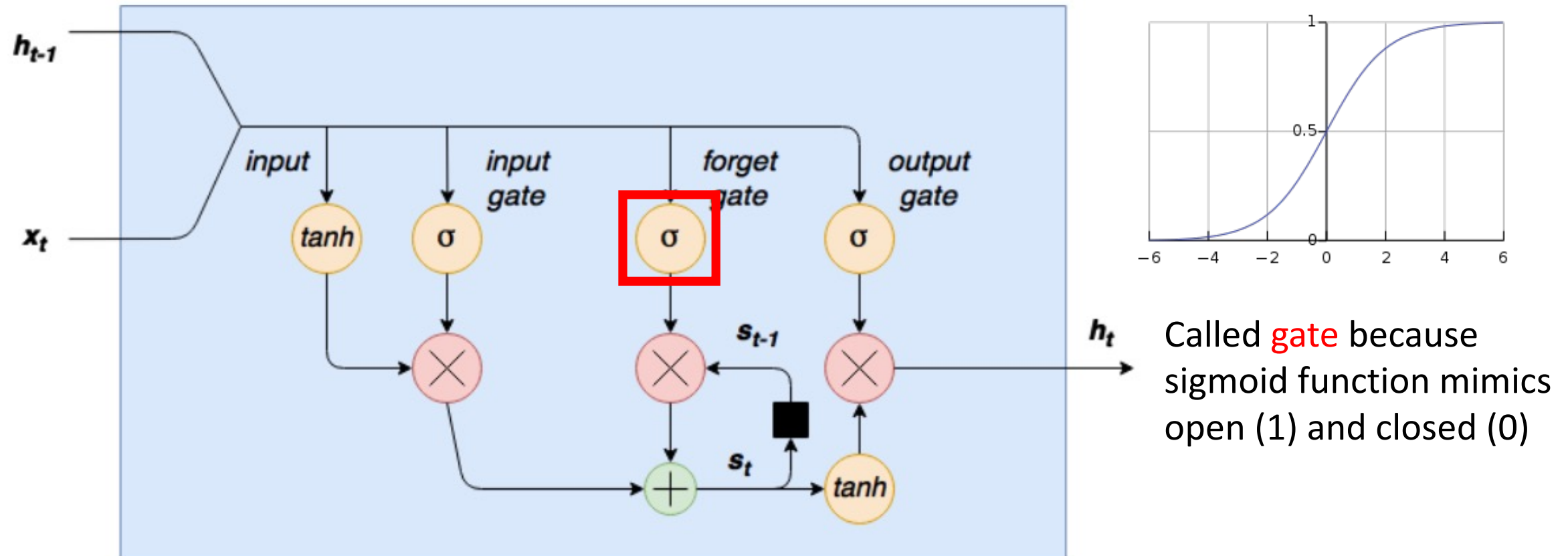
$$\frac{\partial s_t}{\partial s_{t-1}} = f$$

Forget gate range: 0 to 1

- What happens if 0?
  - Past discarded
- What happens if 1?
  - Past retained

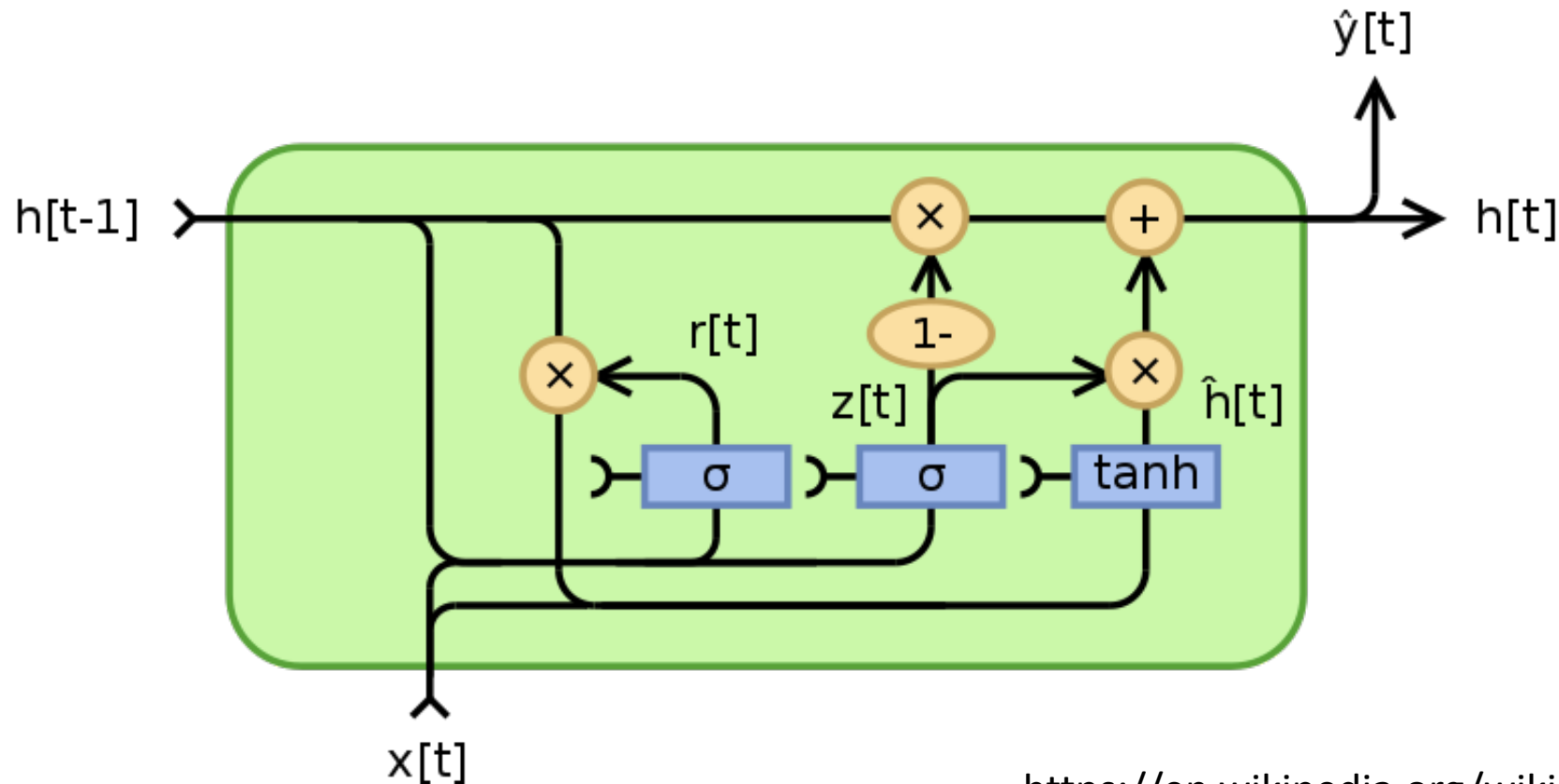
# Popular RNN Variant #1: LSTM

**Cell state:** can retain memory of the past for a long duration, based on the forget gate



# Popular RNN Variant #1: Gated Recurrent Unit

- Simplifies LSTM by merging: (1) cell and hidden states and (2) forget and input gates



[https://en.wikipedia.org/wiki/Gated\\_recurrent\\_unit](https://en.wikipedia.org/wiki/Gated_recurrent_unit)

# Today's Topics

- Deep learning for sequential data
- Recurrent neural networks (RNNs)
- Gated RNNs
- Programming tutorial

# Today's Topics

- Deep learning for sequential data
- Recurrent neural networks (RNNs)
- Gated RNNs
- Programming tutorial



*The End*