

# Machine Learning (part II)

## Recurrent Neural Networks

Angelo Ciaramella

# Recurrent Neural Networks

- RNNs
  - family of neural networks for processing **sequential data**
  - specialized for processing a sequence of values

$$x^{(1)}, \dots, x^{(\tau)}$$

- **early ideas** found in machine learning and statistical models of the 1980s
  - **sharing parameters** across different parts of a model



# Recurrent Neural Networks

---

- Related idea
  - use of convolution across a 1-D temporal sequence
    - time-delay neural networks
- RNNs
  - minibatches of sequences
  - may also be applied in two dimensions across spatial data such as images



# Adaptive filters

---

- Adaptive filter
  - The parameters are estimated
    - learning algorithm
  - An error function is used
  - e.g., Linear Artificial Neural Network ([Adaline](#))





# Adaptive filters

---

## ■ Hospital

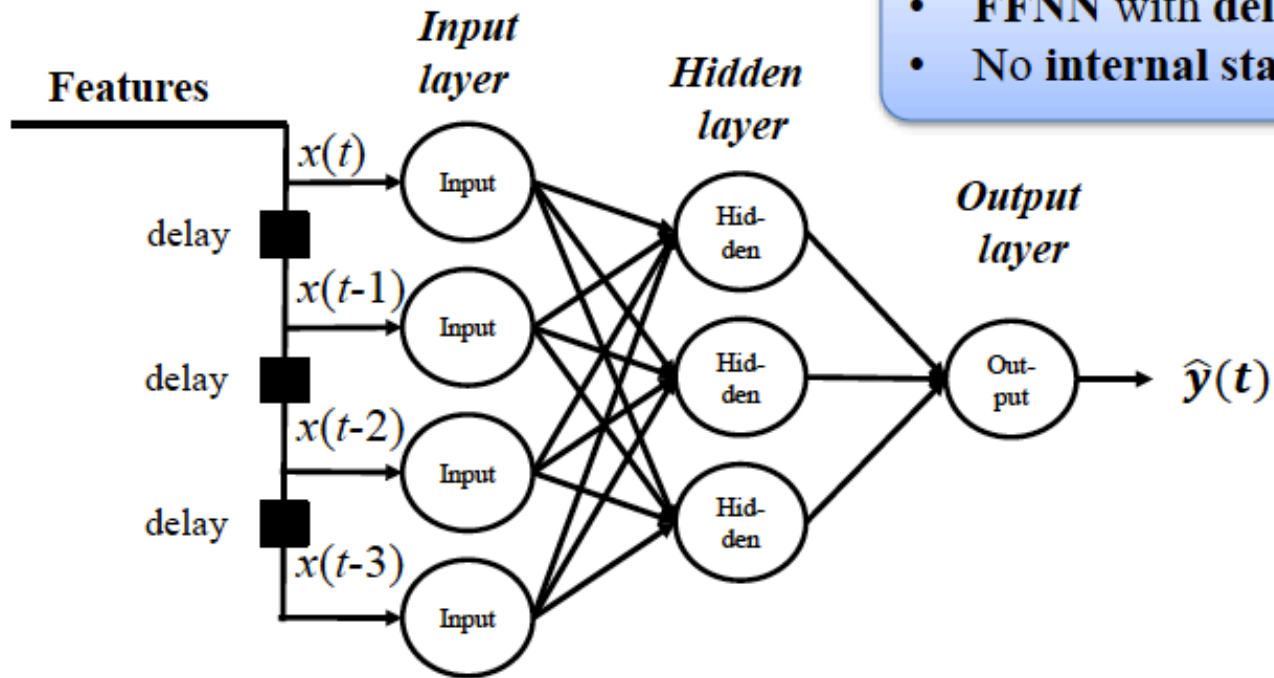
- ECG (electrocardiogram) corrupted by noise at 50 Hz (electricity)
- The current can vary between 47 Hz and 53 Hz
- A filter for the elimination of static noise at 50 Hz could give errors
- An adaptive filter can learn from the current shape of noise

## ■ Helicopter

- Pilot speaking with noise from rotating propeller
- The noise has not a spectrum well defined
- An adaptive filter learns the shape of the noise
- The noise can be subtracted from the signal for only the pilot's voice



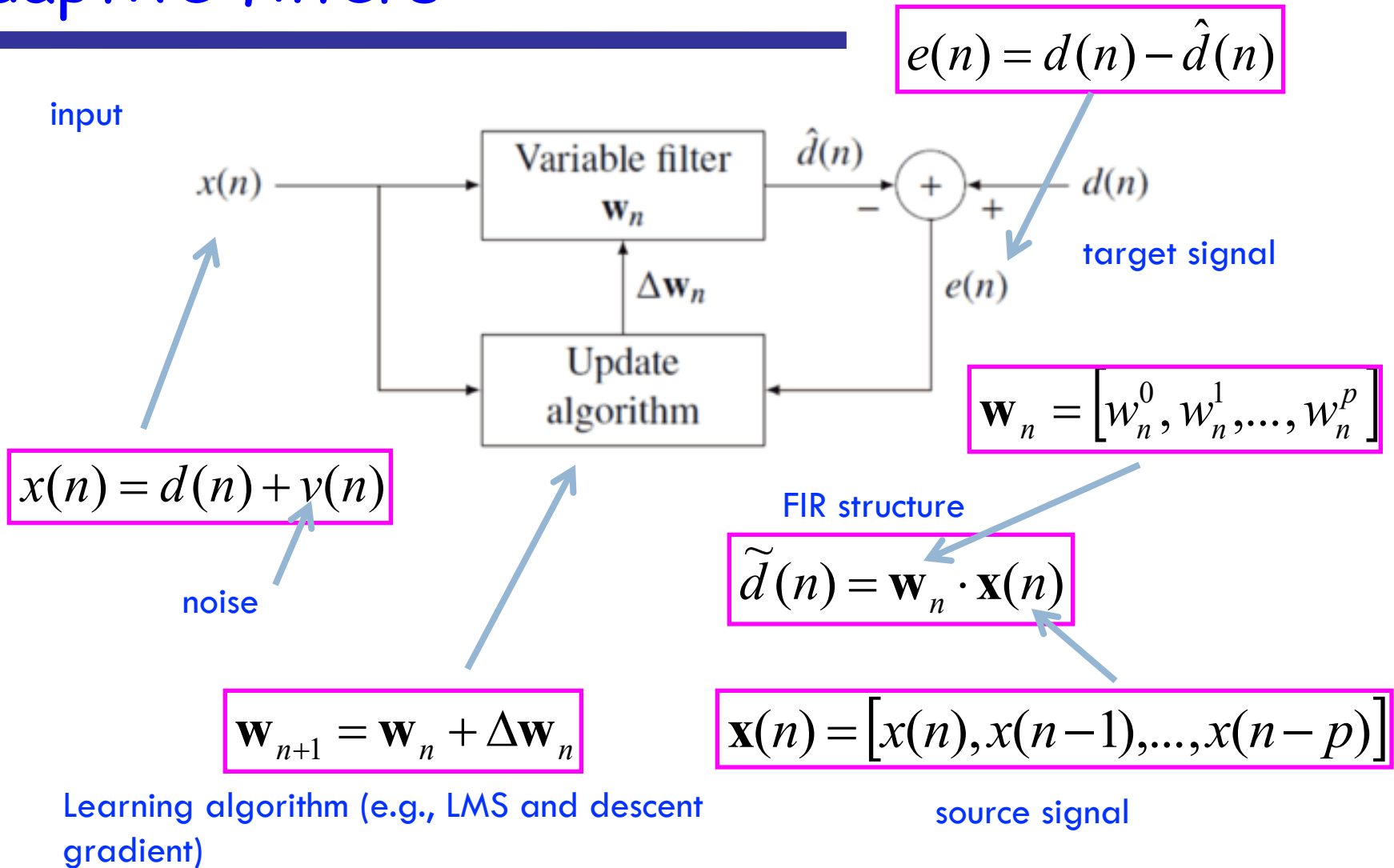
# Adaline



- **FFNN with delayed inputs**
- **No internal state**



# Adaptive filters



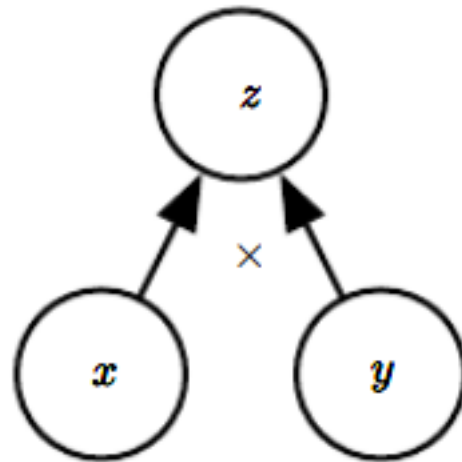
# Computational graphs

---

- Related idea
  - formalize the structure of a set of **computations**
  - introduce the idea of an **operation**
    - an operation is a **simple function** of **one or more variables**



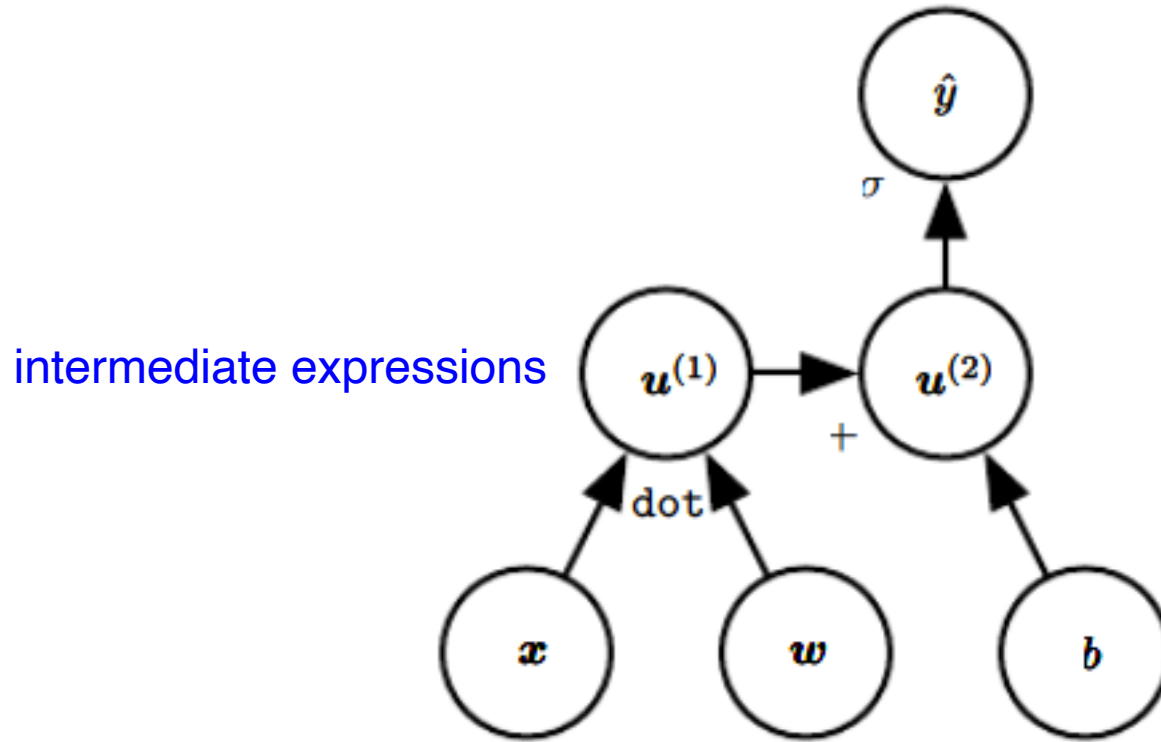
# Computational graphs



Graph using the  $\times$  operation to compute  $z = xy$



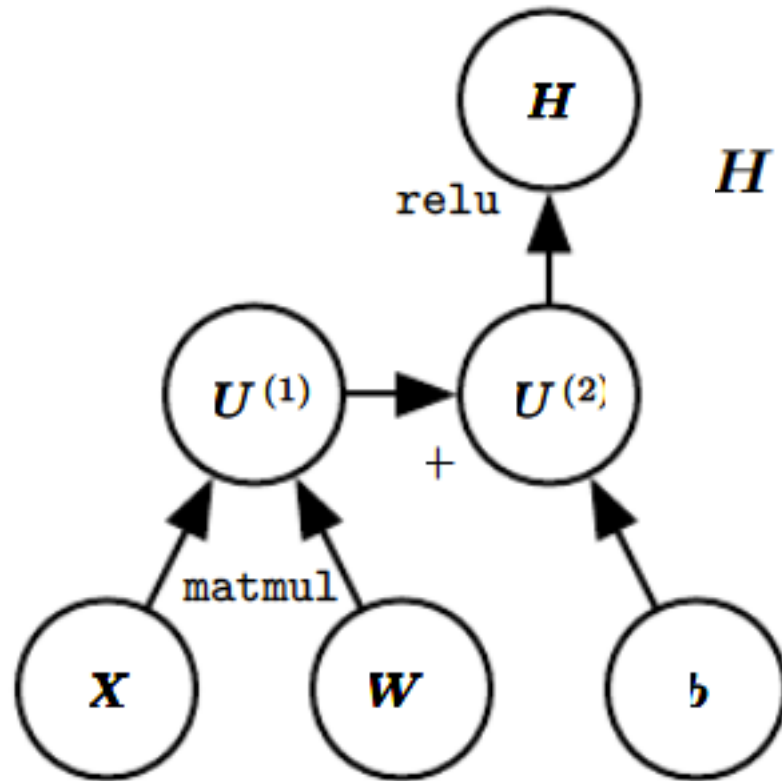
# Computational graphs



logistic regression prediction  $\hat{y} = \sigma(x^\top w + b)$



# Computational graphs



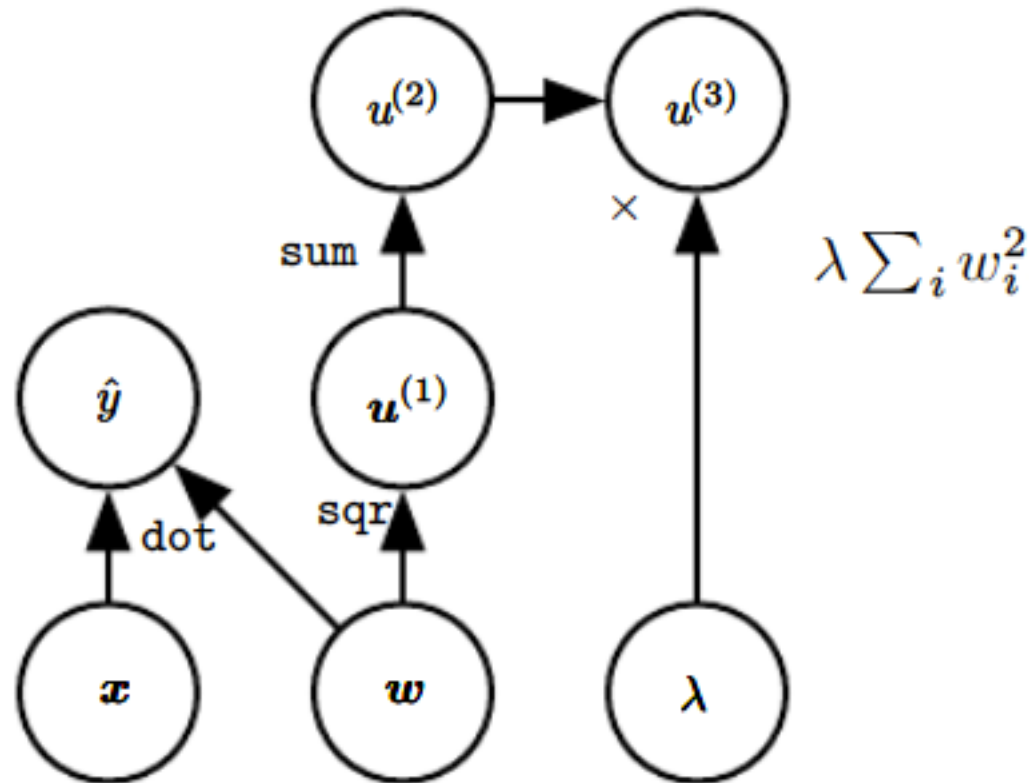
$$H = \max\{0, XW + b\}$$

Minibatch of inputs  $X$





# Computational graphs



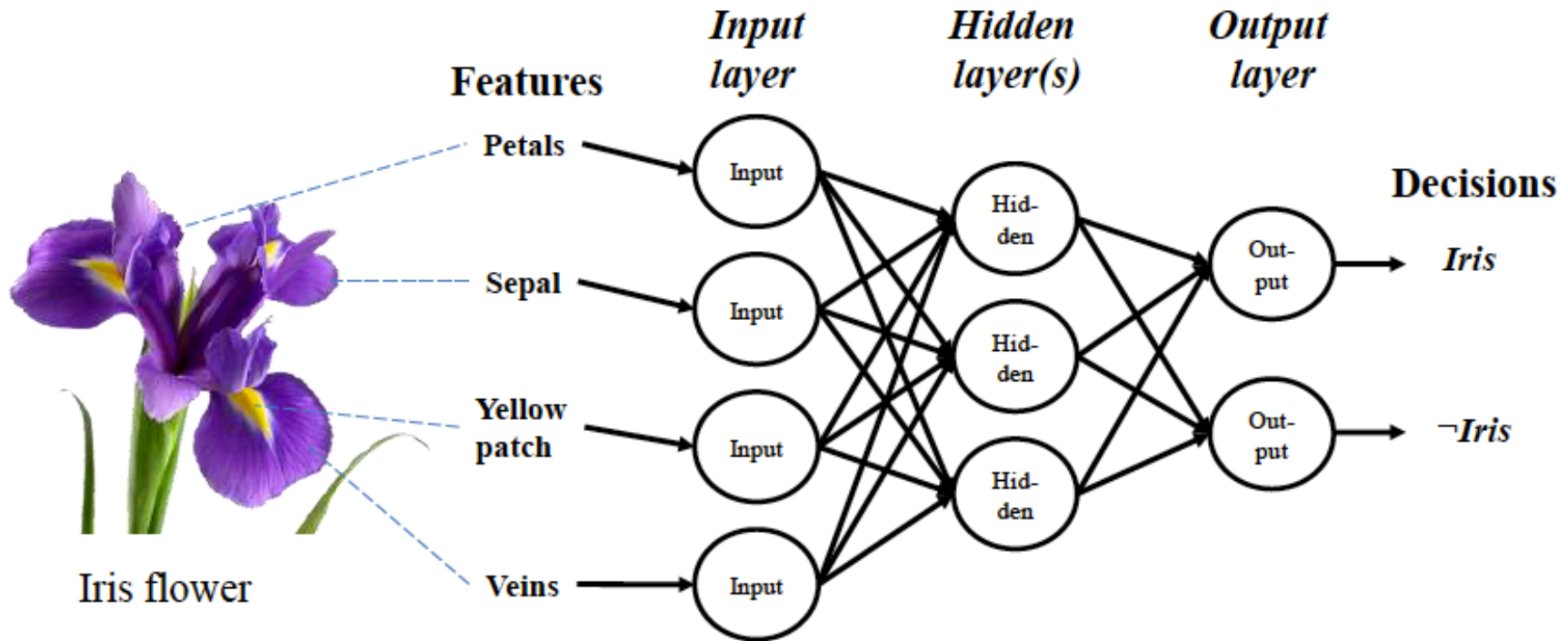
More than one operation of a linear regression model



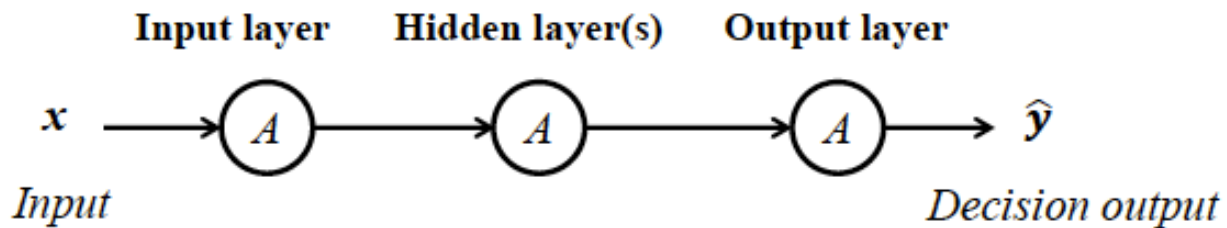
- Artificial Neural Networks
  - exhibit temporal dynamic behavior
  - can use their internal state (memory) to process sequences of inputs
  - models sequences
    - Time series
    - Natural Language
    - Speech
    - Convert non-sequences to sequences, eg: feed an image as a sequence of pixels!



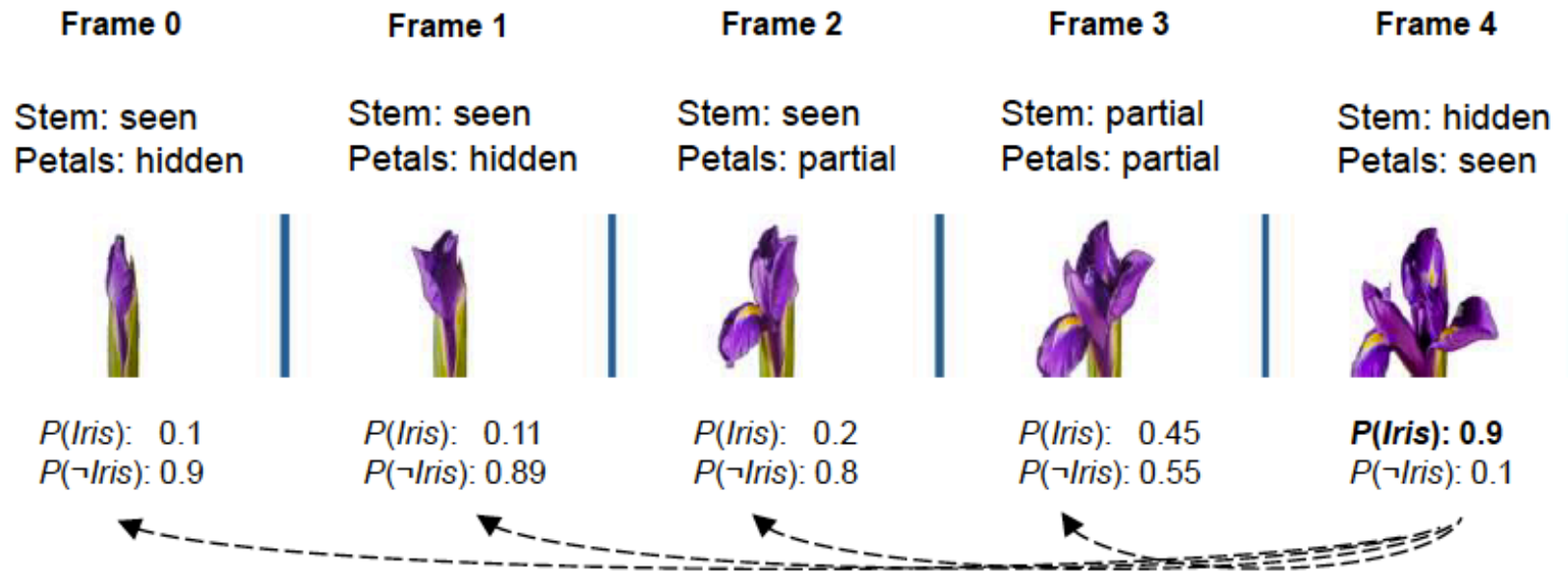
# Feed-forward NN



## Simplified representation:



# Temporal dependencies

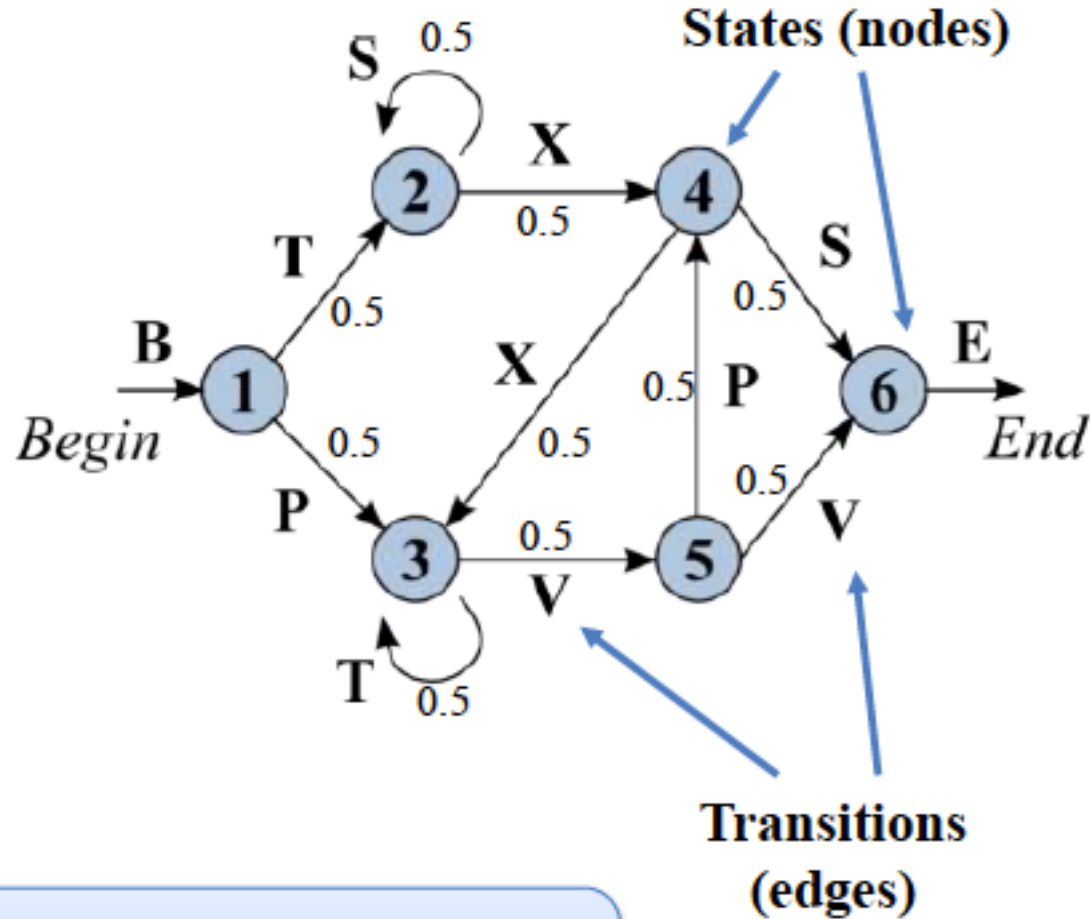


Decision on  
sequence of  
observations



# Reber Grammar

Problem that can not be solved without memory

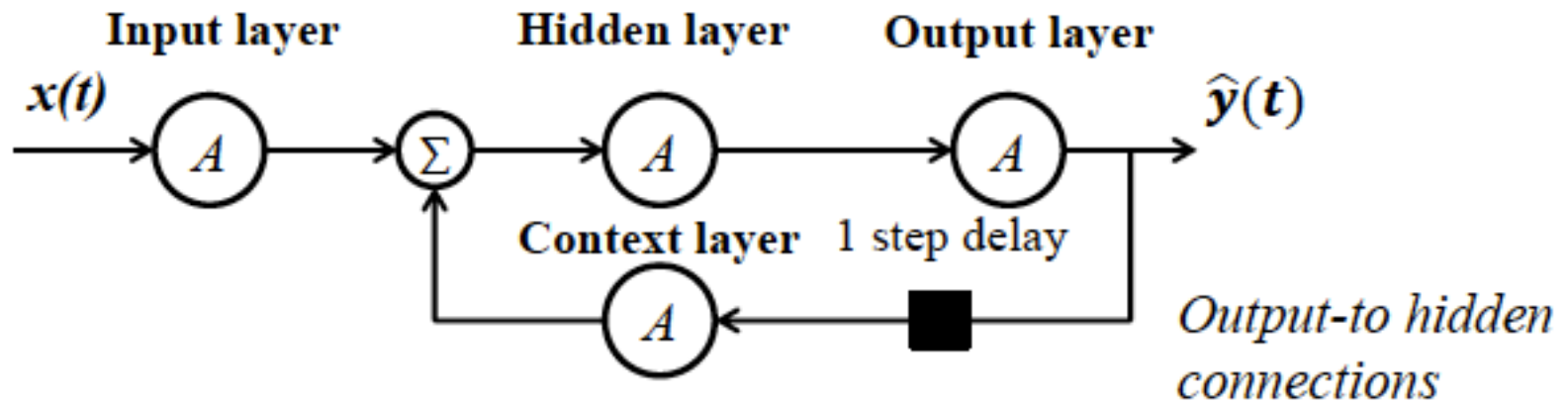


**Transitions** have equal probabilities:  
 $P(1 \rightarrow 2) = P(1 \rightarrow 3) = 0.5$



# Jordan's sequential network

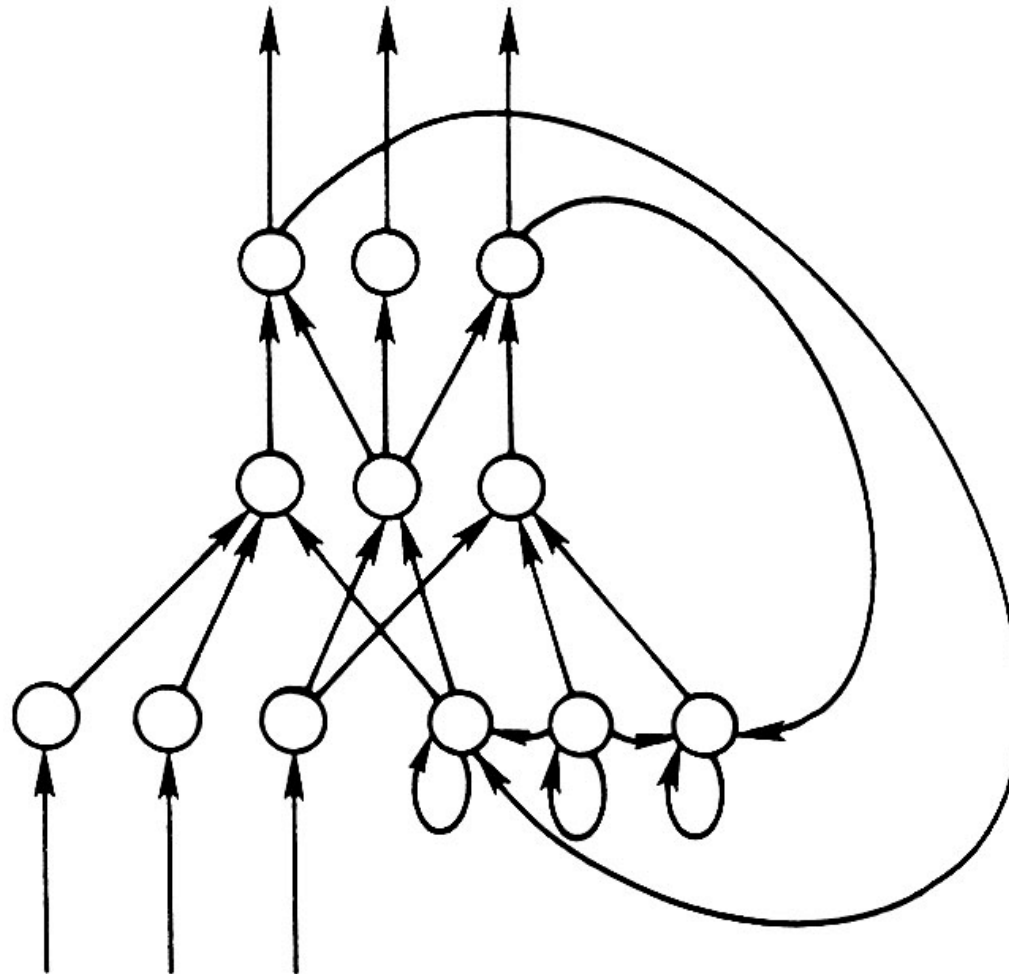
Limited short-term memory



M.I. Jordan NN (1986).



# Jordan's sequential network

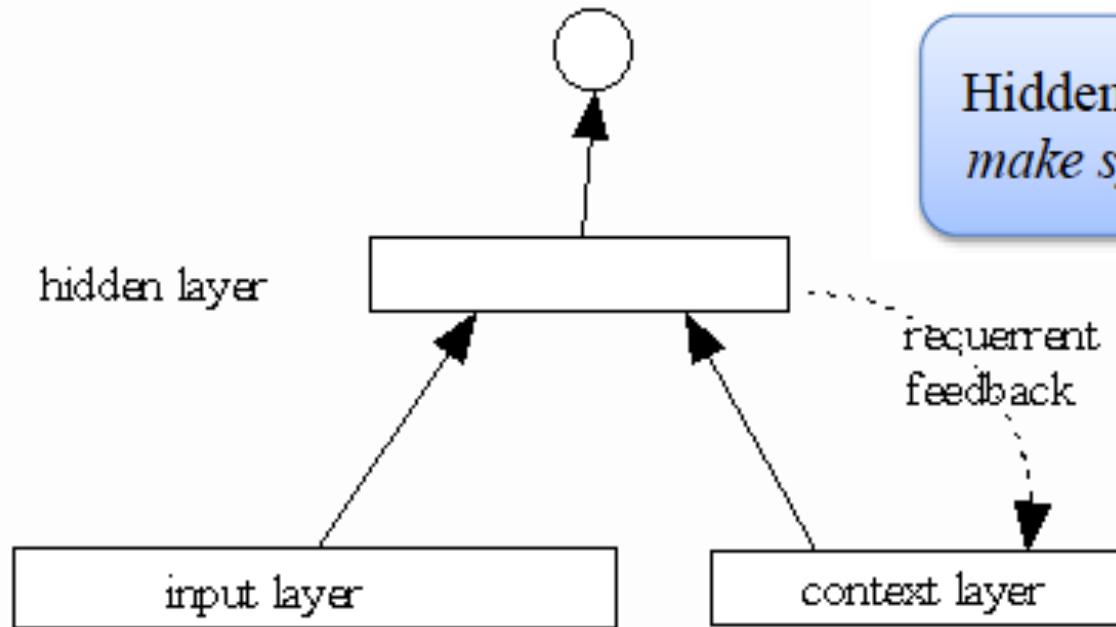


Jordan NN has been applied to categorize a class of English syllables.

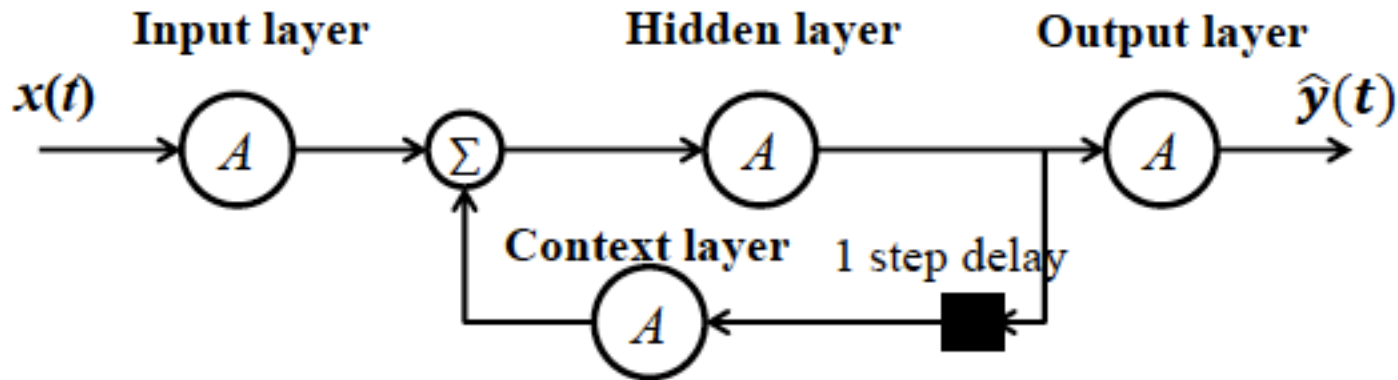




# Simple recurrent network



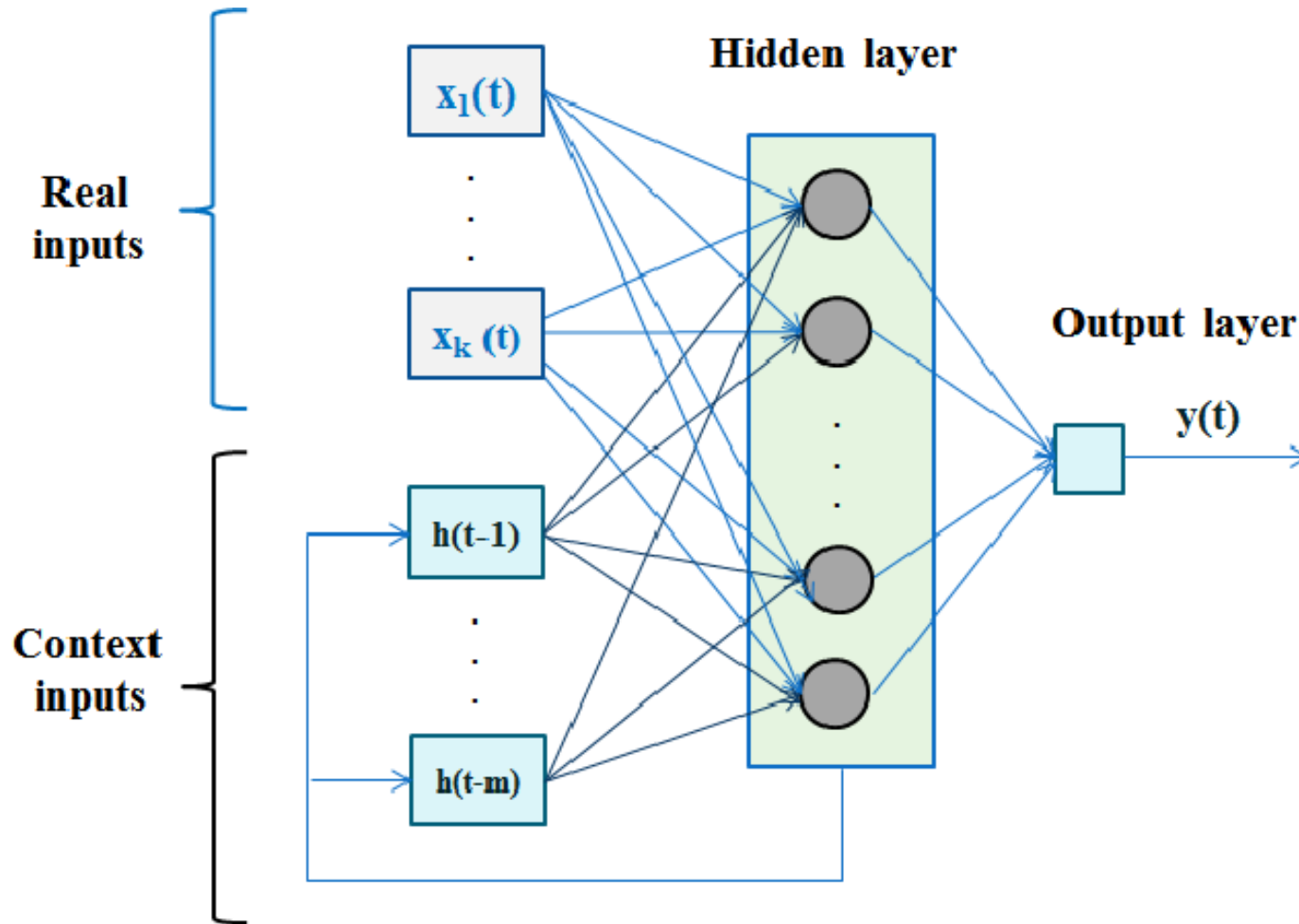
Hidden-to hidden connections  
*make system Turing-complete*



Elman RNN (1990).



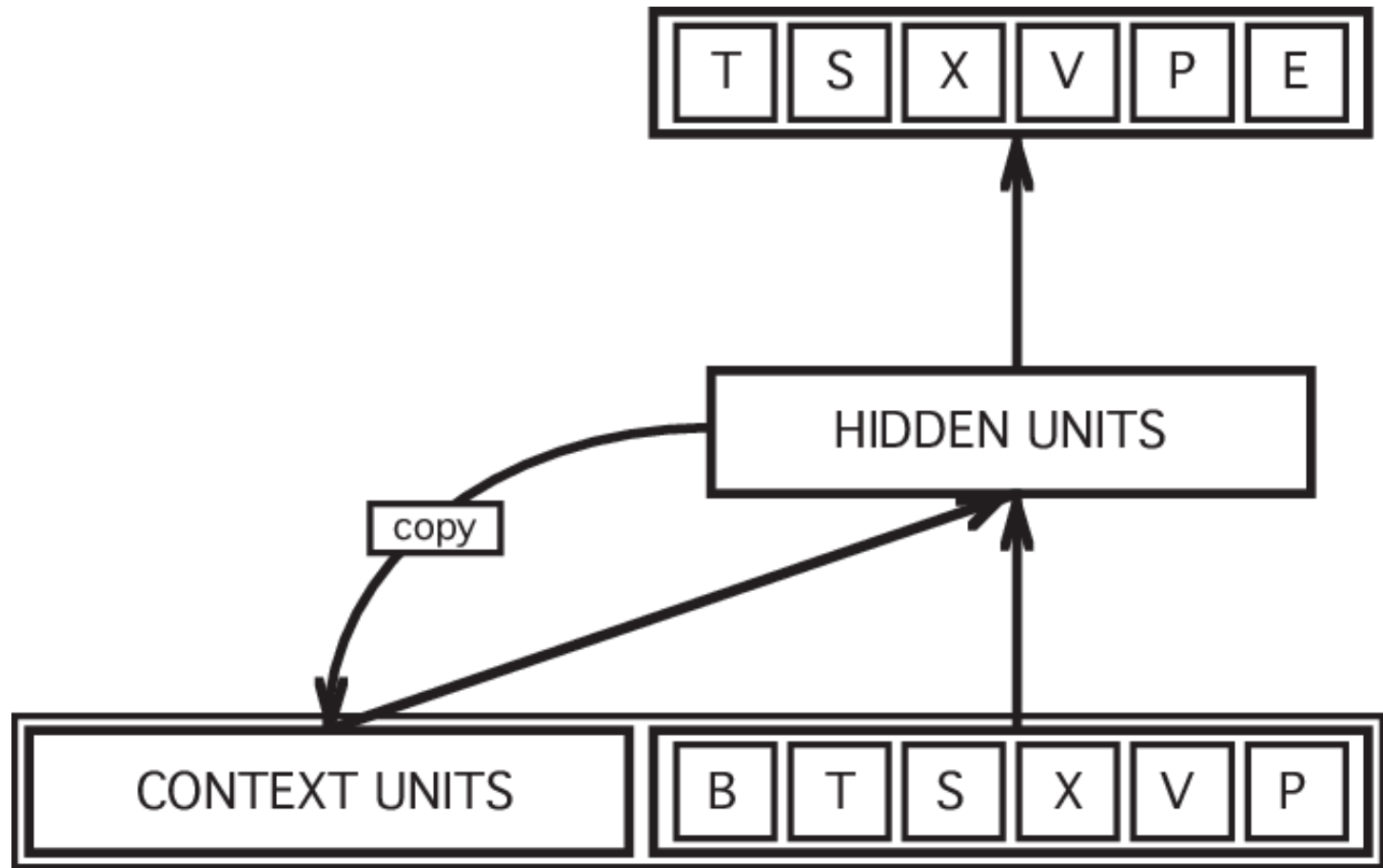
# Simple recurrent network



Elman RNN. Basic RNN structure called “Vanilla” RNN



# Simple recurrent network



Elman SRN. A total of 60,000 randomly generated strings are used for training.



# Applications of RNNs

A person riding a motorcycle on a dirt road.



Image Captioning

RNN Generated Music  
RNN Generated Eminem  
rapper

... and more!

VIOLA:

Why, Salisbury must find his flesh and thought  
That which I am not aps, not a man and in fire,  
To show the reining of the raven and the wars  
To grace my hand reproach within, and not a fair are hand,  
That Caesar and my goodly father's world;

Write like Shakespeare

In reply to Thomas Paine



**DeepDrumpf** @DeepDrumpf · Mar 20

There will be no amnesty. It is going to pass because the people are going to be gone. I'm giving a mandate. #ComeyHearing @Thomas1774Paine



1



12



17

.. and Trump

## Twitterbot

I'm a Neural Network trained on Trump's transcripts. Priming text in [ ]s. Donate (<http://www.gofundme.com/deepdrumpf>) to interact! Created by [@hayesbh](#).



# Recurrent Neural Networks

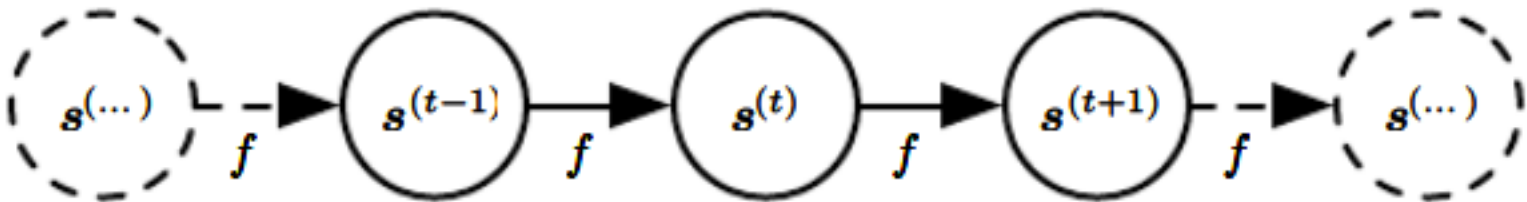
- Dynamical system (recurrent expression)

$$\mathbf{s}^{(t)} = f(\mathbf{s}^{(t-1)}; \boldsymbol{\theta})$$

State of the system

- Example

$$\begin{aligned}\mathbf{s}^{(3)} &= f(\mathbf{s}^{(2)}; \boldsymbol{\theta}) \\ &= f(f(\mathbf{s}^{(1)}; \boldsymbol{\theta}); \boldsymbol{\theta})\end{aligned}$$



Unfolded computational graph



# Recurrent Neural Networks

- Dynamical system driven by an external signal

$$\mathbf{s}^{(t)} = f(\mathbf{s}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta})$$

- RNNs

$$\mathbf{h}^{(t)} = f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta})$$

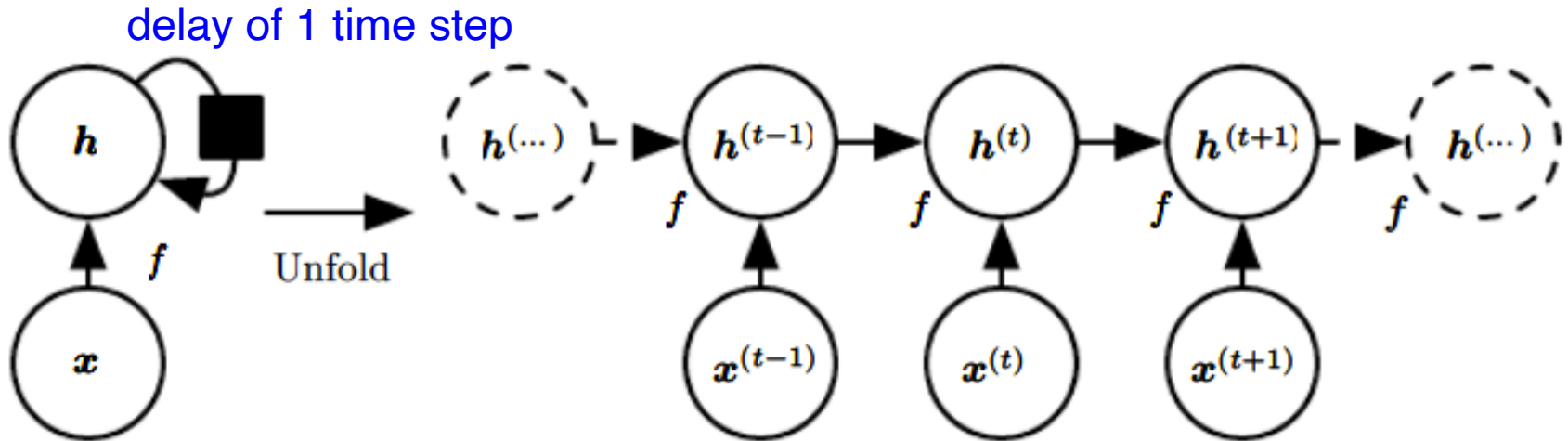
state of the hidden units of the network



# Recurrent Neural Networks

A RNN with no outputs

Circuit diagram



This recurrent network just processes information from the input  $x$  by incorporating it into the state  $h$  that is passed forward through time





# Recurrent Neural Networks

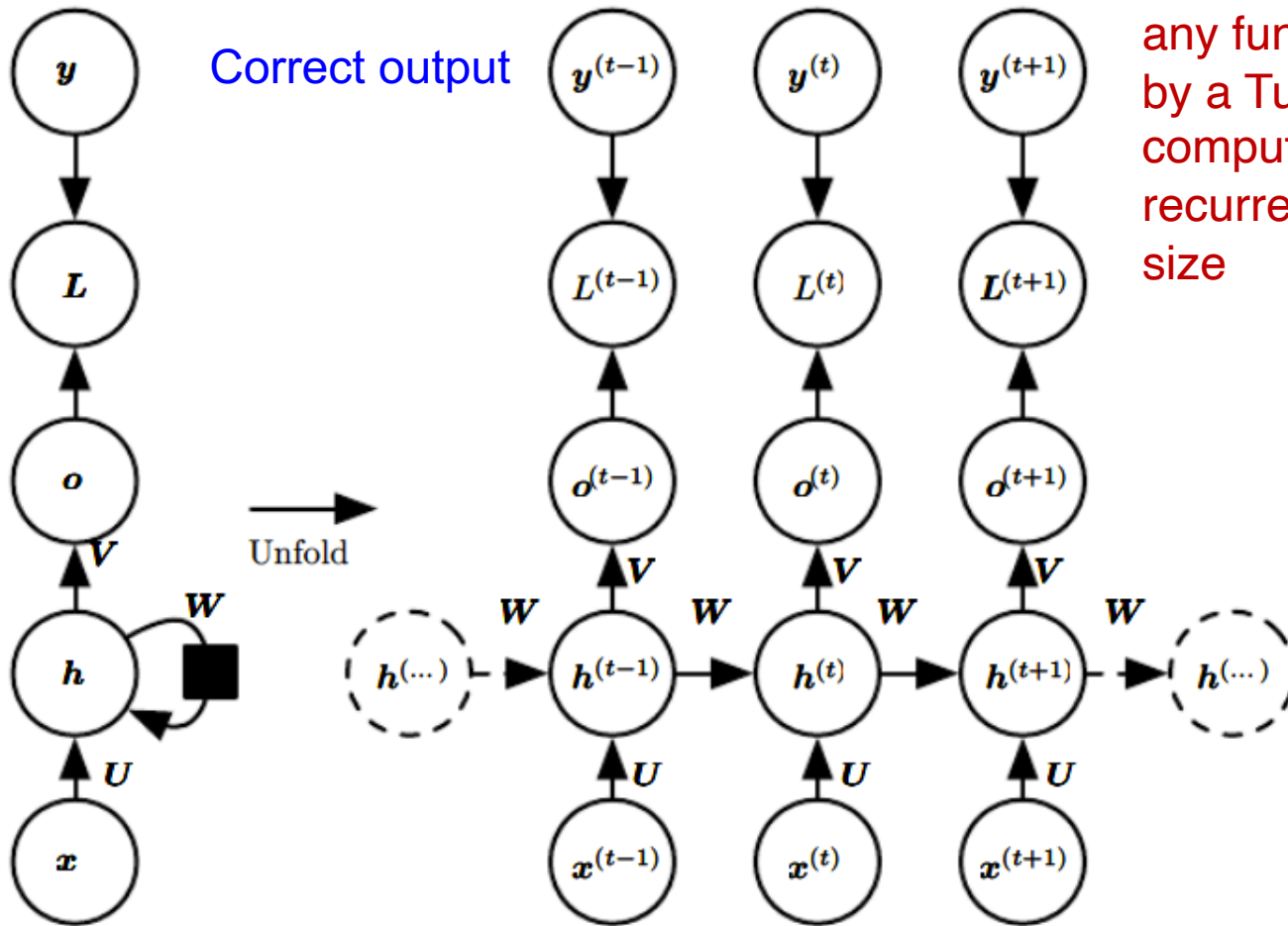
---

- Representation of the unfolded recurrence

$$\begin{aligned} \mathbf{h}^{(t)} &= g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \mathbf{x}^{(t-2)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \\ &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}) \end{aligned}$$



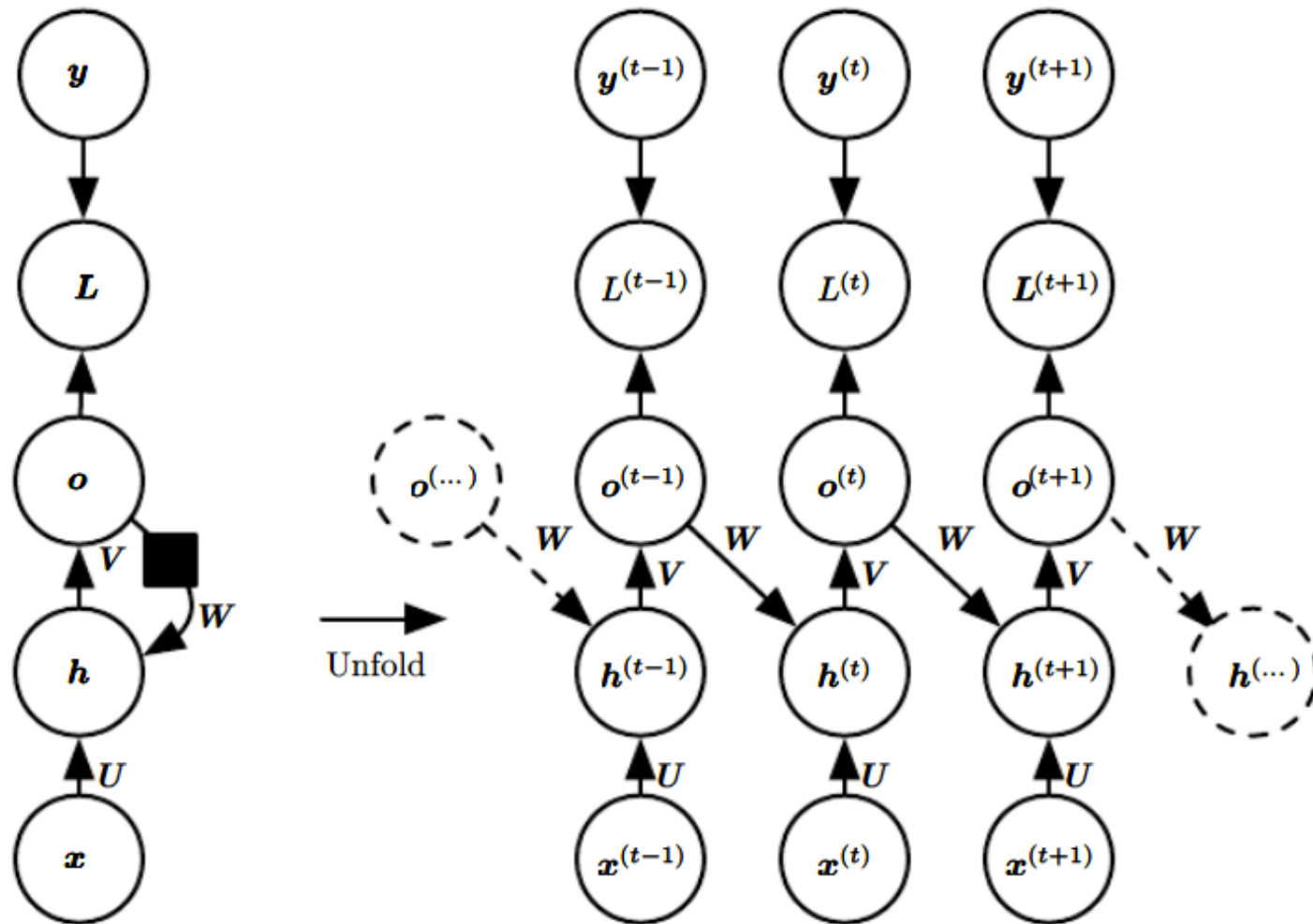
# Recurrent Neural Networks



universal –  
any function computable  
by a Turing machine can be  
computed by such a  
recurrent network of a finite  
size

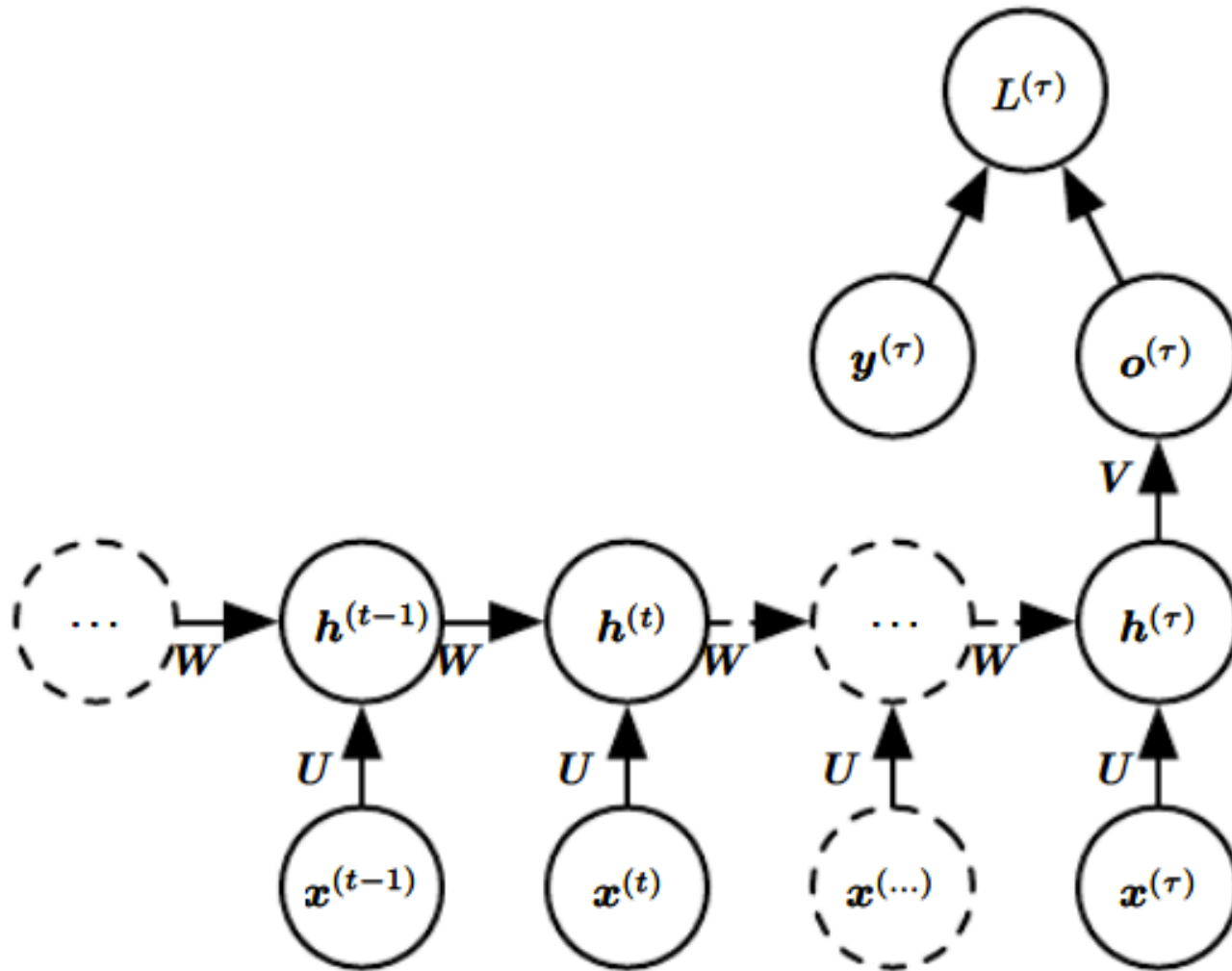
Recurrent networks that produce an output at each time step and have recurrent connections between hidden units

# Recurrent Neural Networks



Recurrent networks that produce an output at each time step and have recurrent connections only from the output at one time step to the hidden units at the next time step

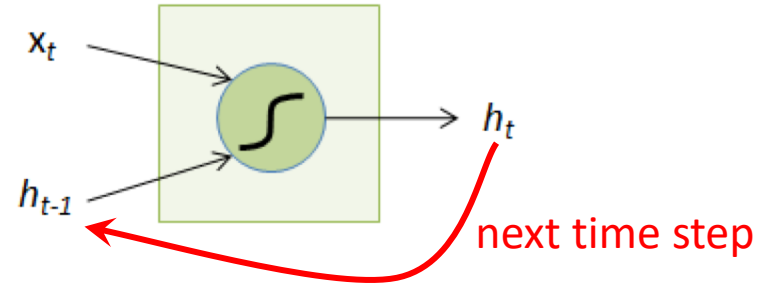
# Recurrent Neural Networks



Recurrent networks with recurrent connections between hidden units, that read an entire sequence and then produce a single output

# Vanilla RNN cell

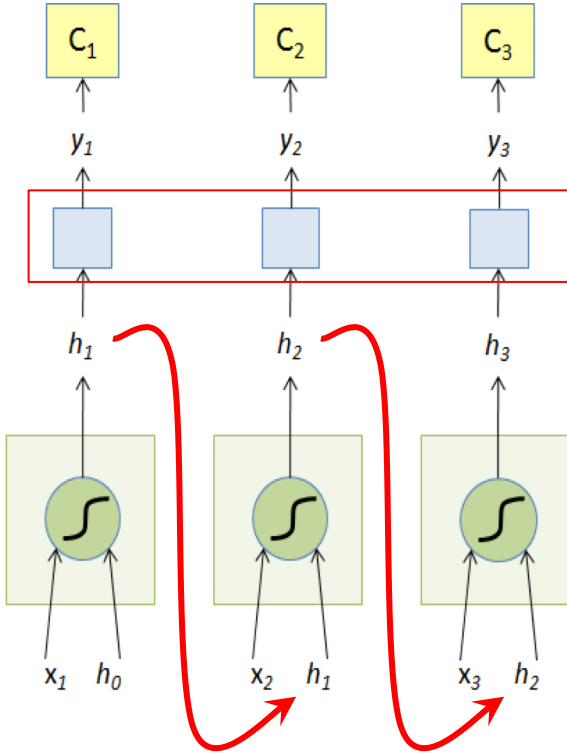
$x_t$ : Input at time t  
 $h_{t-1}$ : State at time t-1



$$h_t = f(W_h h_{t-1} + W_x x_t)$$



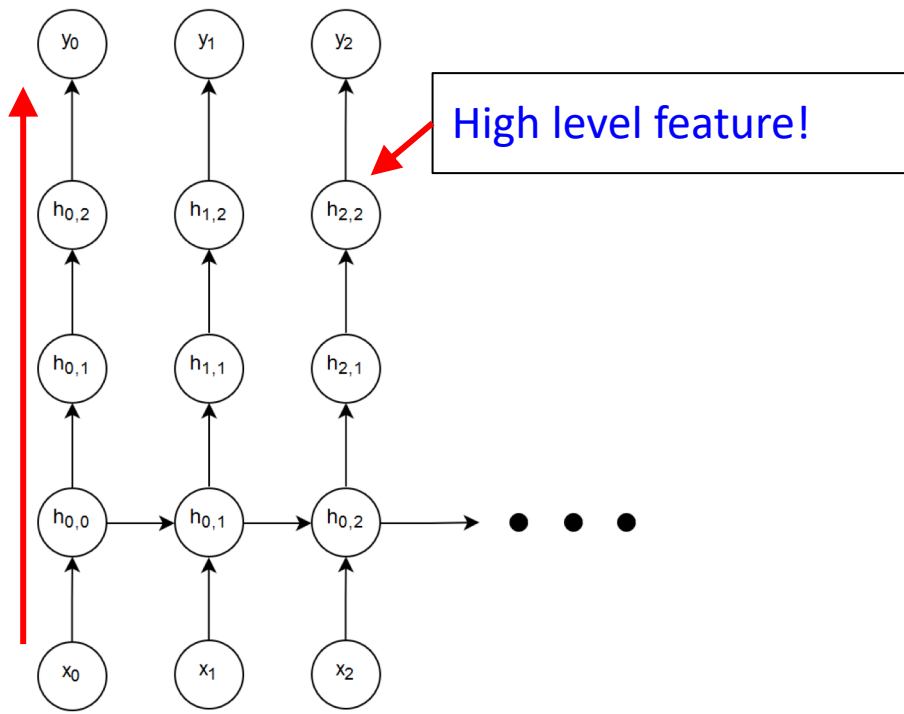
# Unfolding



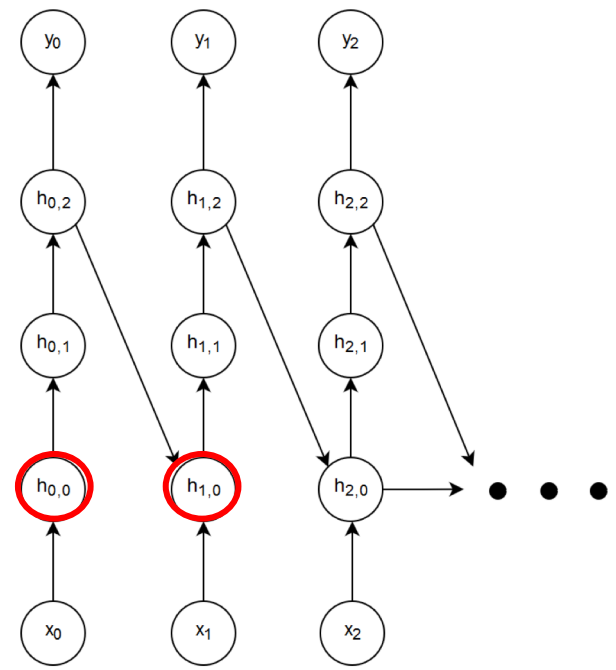
$$h_t = f(W_h h_{t-1} + W_x x_t)$$

Weights shared over time!

# Deep RNN



Feedforward depth = 4

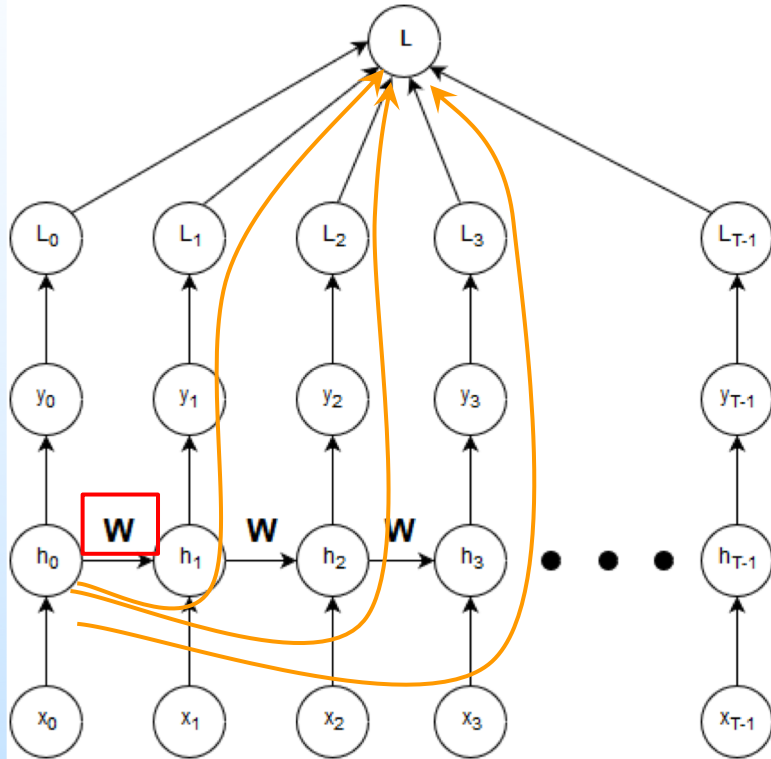


Recurrent depth = 3





# Backpropagation Through Time (BPTT)



Objective is to update the weight matrix:

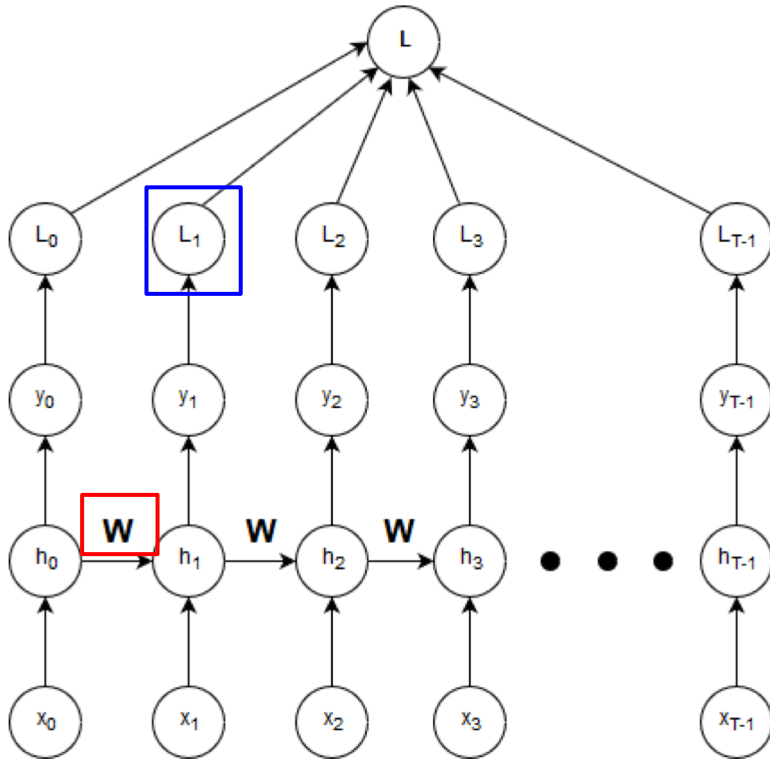
$$\mathbf{W} \rightarrow \mathbf{W} - \alpha \frac{\partial L}{\partial \mathbf{W}}$$

Issue:  $\mathbf{W}$  occurs each timestep  
**Every** path from  $\mathbf{W}$  to  $L$  is one dependency

Find all paths from  $\mathbf{W}$  to  $L$ !

(note: dropping subscript  $h$  from  $\mathbf{W}_h$  for brevity)

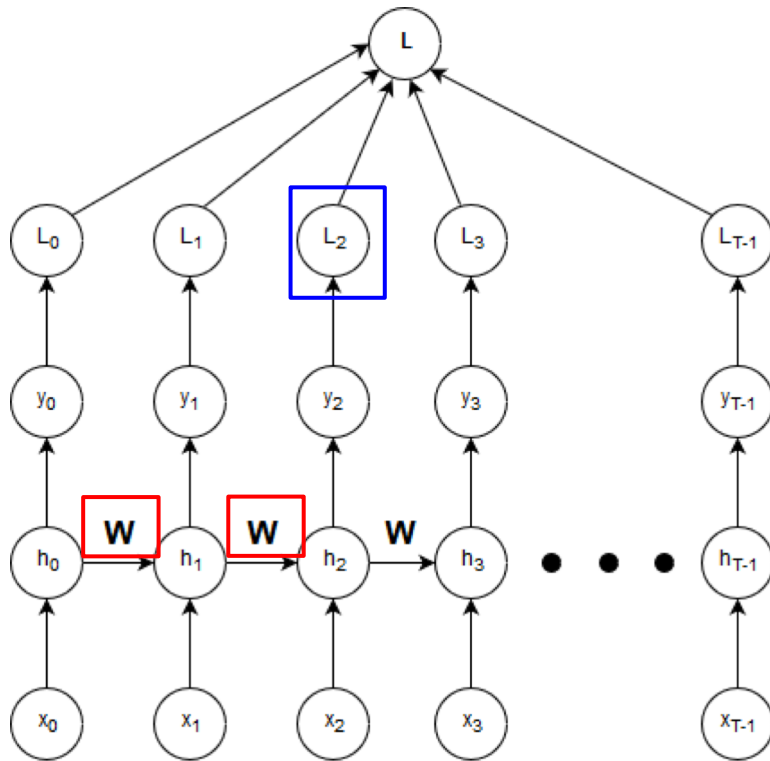
# Backpropagation Through Time (BPTT)



How many paths exist from  $W$  to  $L$  through  $L_1$ ?

Just 1. Originating at  $h_0$ .

# Backpropagation Through Time (BPTT)

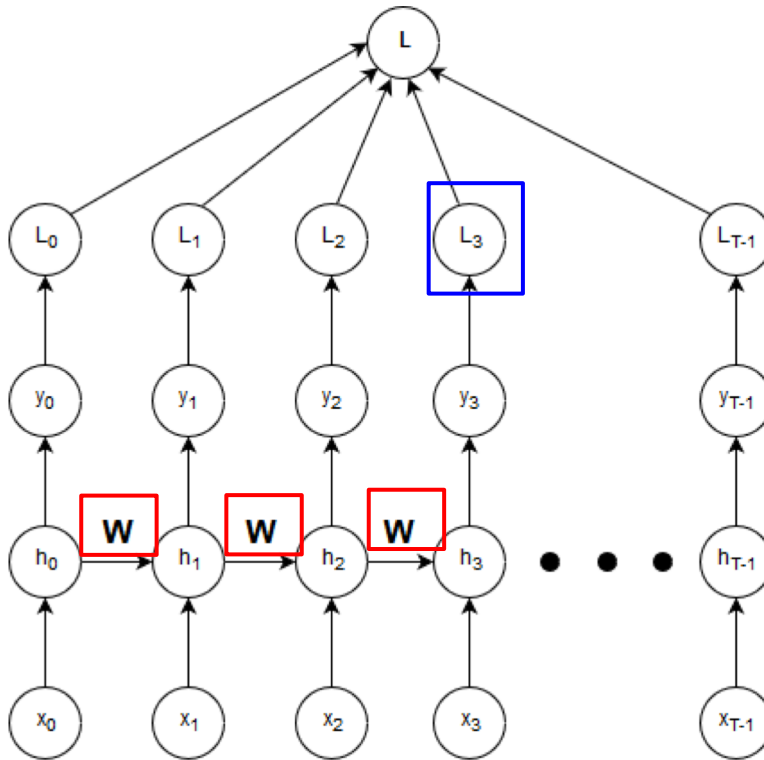


How many paths from  $W$  to  $L$   
through  $L_2$ ?

2. Originating at  $h_0$  and  $h_1$ .



# Backpropagation Through Time (BPTT)



And 3 in this case.

Origin of path = basis for  $\Sigma$

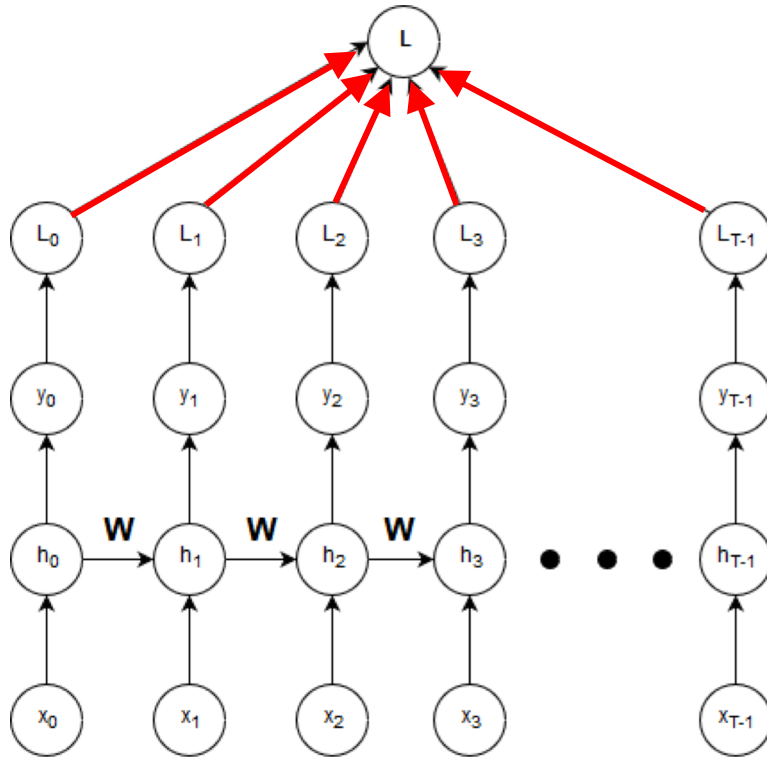
$$\frac{\partial L}{\partial \mathbf{W}}$$

The gradient has two summations:

- 1: Over  $L_j$
- 2: Over  $h_k$



# Backpropagation Through Time (BPTT)

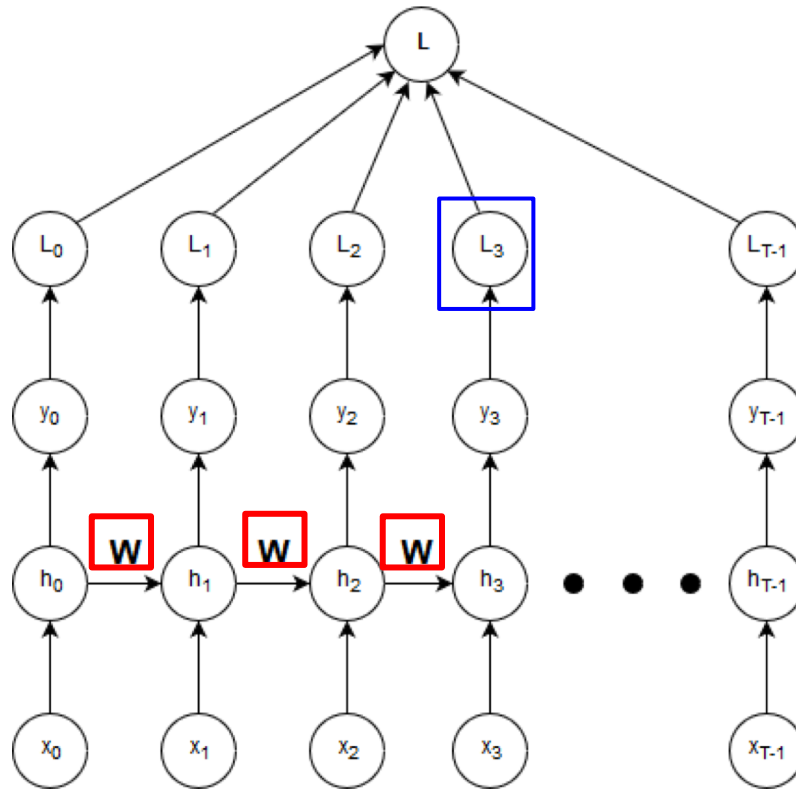


First summation over  $L$

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_{j=0}^{T-1} \frac{\partial L_j}{\partial \mathbf{W}}$$



# Backpropagation Through Time (BPTT)

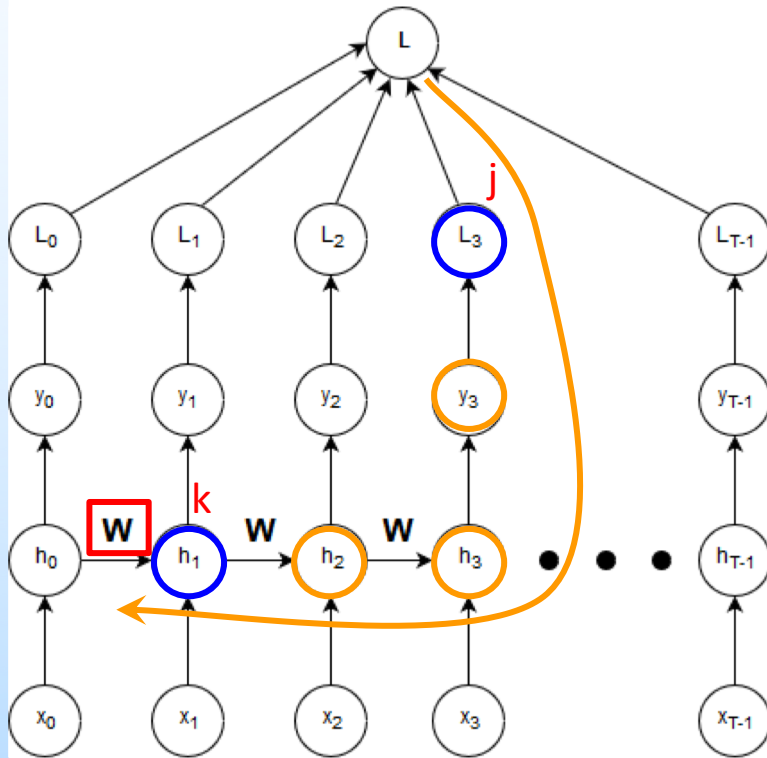


- **Second summation over h:** Each  $L_j$  depends on the weight matrices *before it*

$$\frac{\partial L_j}{\partial \mathbf{W}} = \sum_{k=1}^j \frac{\partial L_j}{\partial h_k} \frac{\partial h_k}{\partial \mathbf{W}}$$

$L_j$  depends on all  $h_k$  before it.

# Backpropagation Through Time (BPTT)



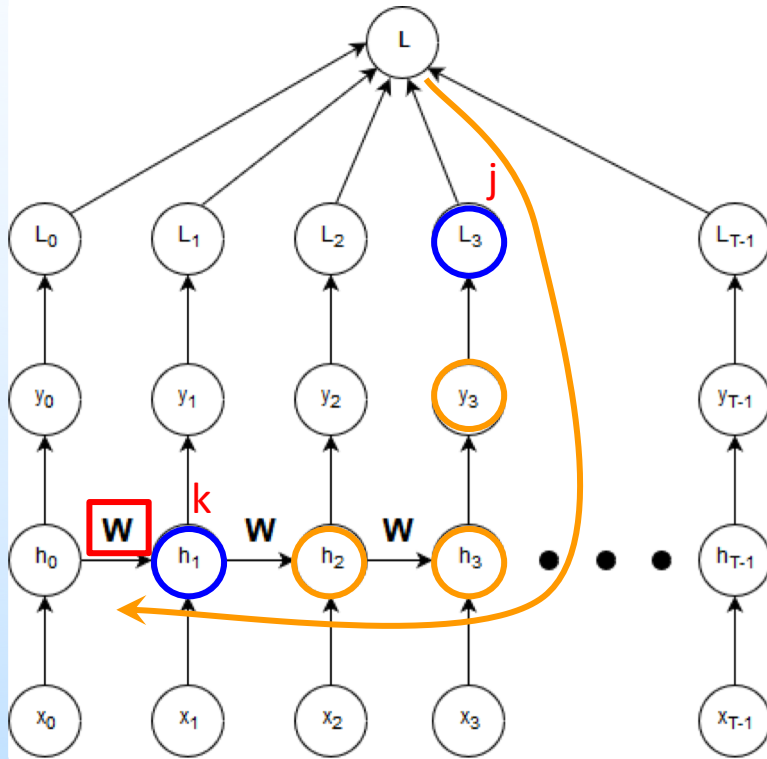
$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \frac{\partial L_j}{\partial h_k} \frac{\partial h_k}{\partial W}$$

No explicit of  $L_j$  on  $h_k$

Use chain rule to fill missing steps

$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \frac{\partial h_j}{\partial h_k} \frac{\partial h_k}{\partial W}$$

# Backpropagation Through Time (BPTT)



$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \frac{\partial L_j}{\partial h_k} \frac{\partial h_k}{\partial W}$$

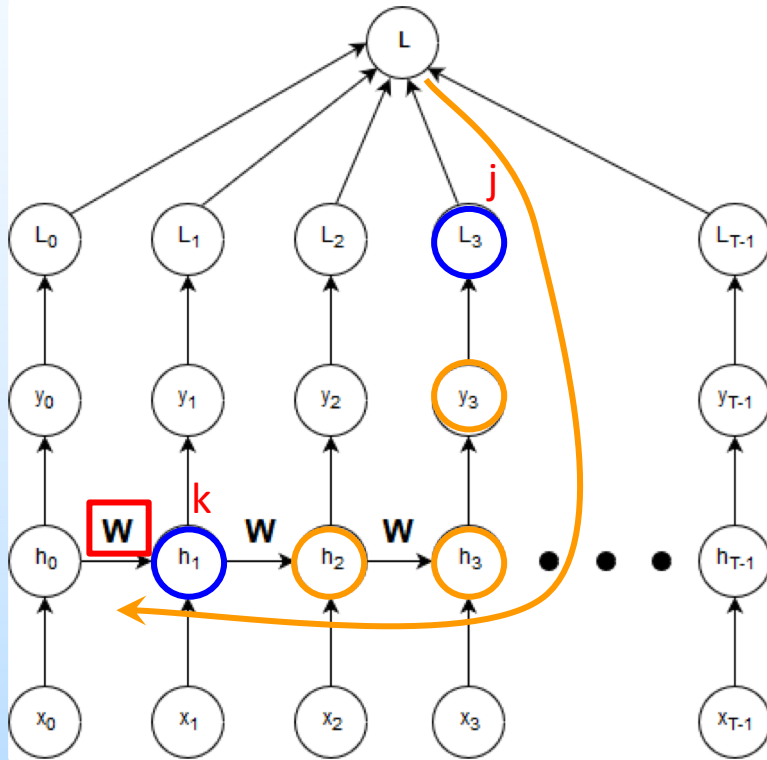
No explicit of  $L_j$  on  $h_k$

Use chain rule to fill missing steps

$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \frac{\partial h_j}{\partial h_k} \frac{\partial h_k}{\partial W}$$



# Backpropagation Through Time (BPTT)



$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \frac{\partial L_j}{\partial h_k} \frac{\partial h_k}{\partial W}$$

No explicit of  $L_j$  on  $h_k$

Use chain rule to fill missing steps

$$\frac{\partial L_j}{\partial W} = \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \frac{\partial h_j}{\partial h_k} \frac{\partial h_k}{\partial W}$$

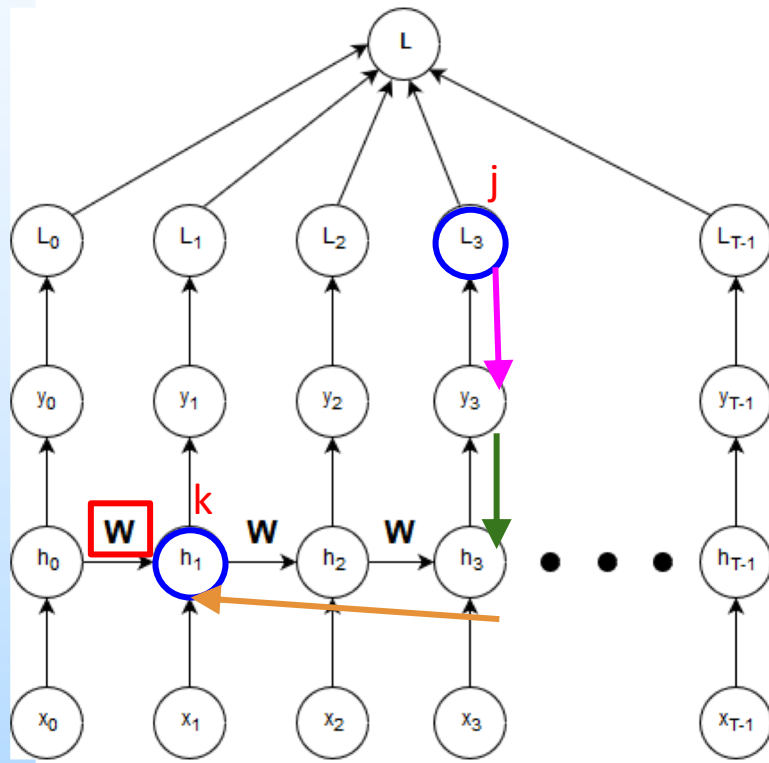
# Backpropagation Through Time (BPTT)

$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{j=0}^{T-1} \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}_h}$$

The final Backpropagation equation



# Backpropagation Through Time (BPTT)



$$\frac{\partial L}{\partial \mathbf{W}_h} = \sum_{j=0}^{T-1} \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}_h}$$

- Often, to reduce memory requirement, we truncate the network
- Inner summation runs from  $j-p$  to  $j$  for some  $p \implies$  truncated BPTT

# Backpropagation Trough Time (BPTT)

$$\frac{\partial L}{\partial \mathbf{W}} = \sum_{j=0}^{T-1} \sum_{k=1}^j \frac{\partial L_j}{\partial y_j} \frac{\partial y_j}{\partial h_j} \left( \prod_{m=k+1}^j \frac{\partial h_m}{\partial h_{m-1}} \right) \frac{\partial h_k}{\partial \mathbf{W}}$$

$$h_m = f(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m)$$

$$\frac{\partial h_m}{\partial h_{m-1}} = \mathbf{W}_h^T \text{diag}(f'(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m))$$

Expanding the Jacobian



# Backpropagation Through Time (BPTT)

$$\frac{\partial h_j}{\partial h_k} = \prod_{m=k+1}^j \mathbf{W}_h^T \text{diag}(f'(\mathbf{W}_h h_{m-1} + \mathbf{W}_x x_m))$$

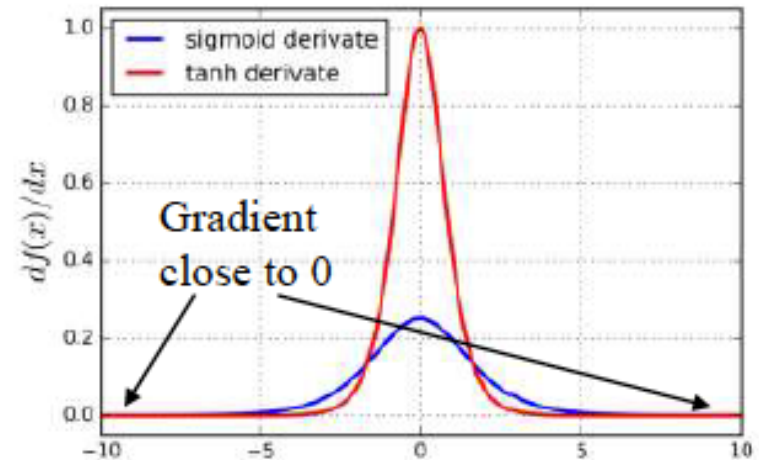
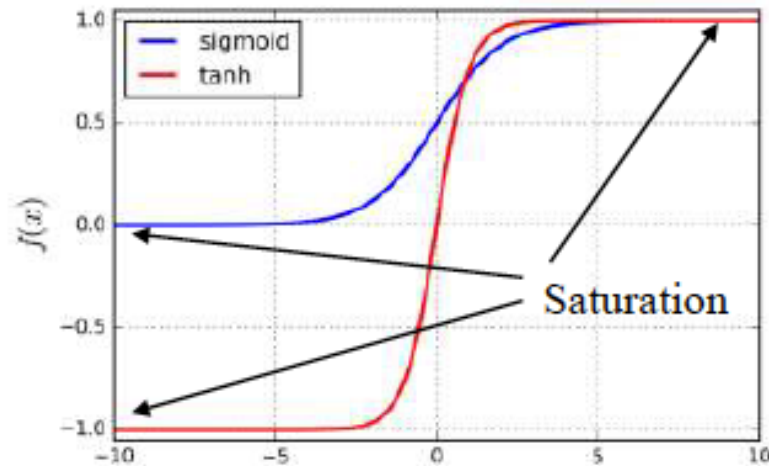
Weight Matrix

Derivative of activation function

Repeated matrix multiplications leads to **vanishing and exploding gradients**.



# Vanishing gradients



$$\frac{\partial h_t}{\partial h_0} = \frac{\partial h_t}{\partial h_{t-1}} \cdot \dots \cdot \frac{\partial h_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial h_1} \cdot \frac{\partial h_1}{\partial h_0}$$

Known problem for deep feed-forward networks.

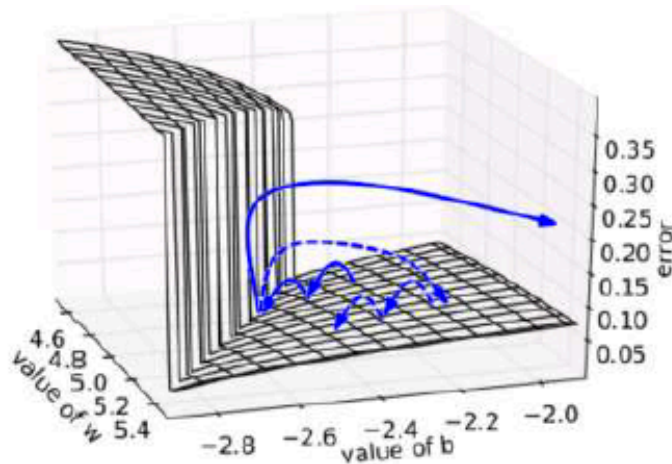
For recurrent networks (even shallow) makes **impossible to learn long-term dependencies!**

## Considerations

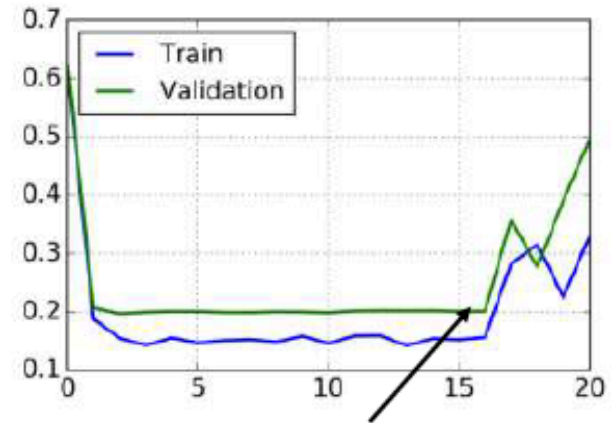
- Smaller weight parameters lead to faster gradients vanishing
- Very big initial parameters make the gradient descent to diverge fast (explode)



# Exploding gradients



Pascanou R. et al. On the difficulty of training recurrent neural networks. arXiv (2012)



Network can not converge and weigh parameters do not stabilize

Large increase in the norm of the gradient during training

**Diagnostics:** NaNs; Cost function large fluctuations

**Solutions:**

- Use gradient clipping
- Try reduce learning rate
- Change loss function by setting constrains on weights ( $L1/L2$  norms)



# Eigenvalues and Stability

Consider identity activation function

If Recurrent Matrix  $\mathbf{W}_h$  is a diagonalizable:

$$\mathbf{W}_h = \mathbf{Q}^{-1} * \mathbf{\Lambda} * \mathbf{Q}$$

$\mathbf{Q}$  matrix composed of  
eigenvectors of  $\mathbf{W}_h$

$\mathbf{\Lambda}$  is a diagonal matrix with  
eigenvalues placed on the  
diagonals

Computing powers of  $\mathbf{W}_h$  is  
simple:

$$\mathbf{W}_h^n = \mathbf{Q}^{-1} * \mathbf{\Lambda}^n * \mathbf{Q}$$





# Eigenvalues and Stability

$$\Lambda = \begin{bmatrix} -0.6180 & 0 \\ 0 & 1.6180 \end{bmatrix} \longrightarrow \Lambda^{10} = \begin{bmatrix} 0.0081 & 0 \\ 0 & 122.9919 \end{bmatrix}$$

Vanishing gradients

Exploding gradients

$$W_h^n = Q^{-1} * \Lambda^n * Q$$



# Fundamental DL problem

---

- DNNs train difficulties
  - Vanishing gradient
  - Exploding gradient
- Solutions
  - Previously proposed
  - Unsupervised pre-training
  - Improve network architecture



# RNNs - forward propagation

- Assume the hyperbolic tangent activation function
- Initial state  $\mathbf{h}^{(0)}$
- Update equation

$$\mathbf{a}^{(t)} = \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)}$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{a}^{(t)})$$

$$\mathbf{o}^{(t)} = \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{o}^{(t)})$$



# RNNs - forward propagation

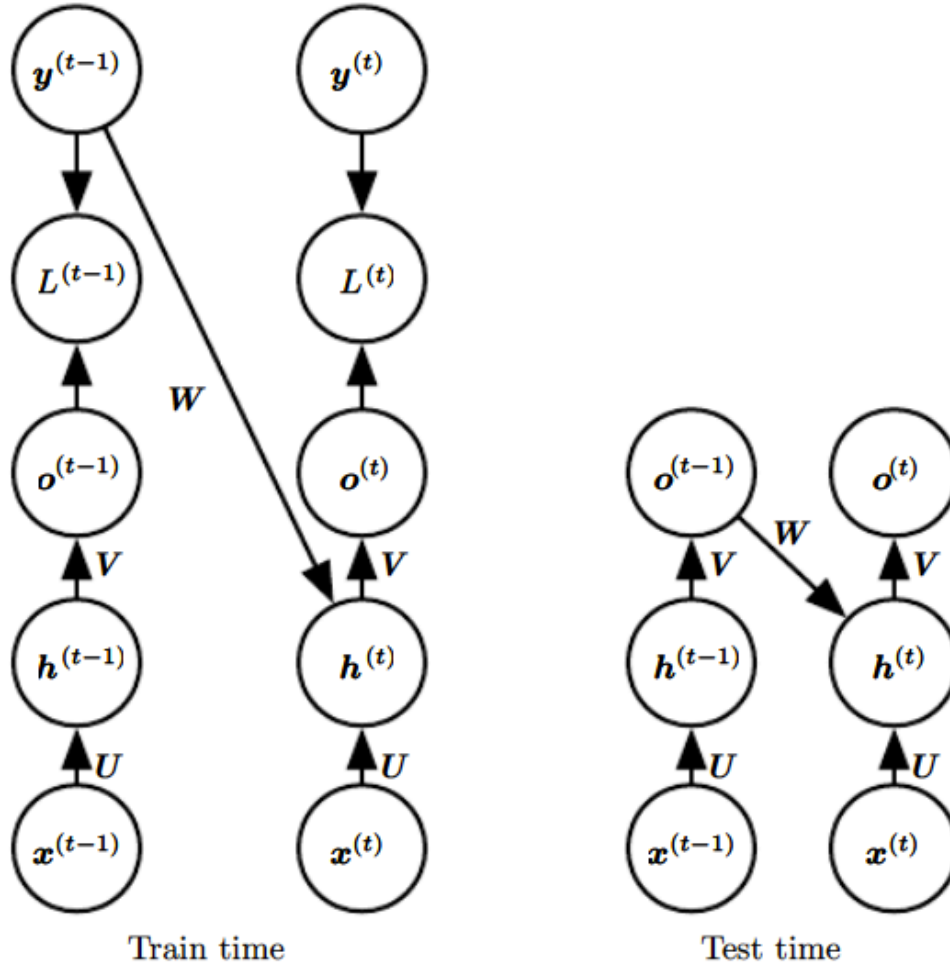
- Total loss

$$\begin{aligned} & L\left(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}\}, \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}\right) \\ &= \sum_t L^{(t)} \\ &= - \sum_t \log p_{\text{model}}\left(y^{(t)} \mid \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\}\right) \end{aligned}$$

Negative log-likelihood



# RNNs - Teacher forcing



# RNNs - learning

- back-propagation through time (BPTT) algorithm
- For each node  $N$  we need to compute the gradient recursively
  - based on the gradient computed at nodes that follow it in the graph

$$\nabla_{\mathbf{N}} L$$

- Start the recursion

$$\frac{\partial L}{\partial L(t)} = 1$$



# RNNs - learning

- Gradient on the outputs at time step  $t$ , for all  $i$ ,  $t$ ,

$$(\nabla_{\boldsymbol{\alpha}^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}}$$

$$\log \text{softmax}(\mathbf{z})_i = z_i - \log \sum_j \exp(z_j)$$

$$\log \sum_j \exp(z_j) \approx \max_j z_j = z_i$$



# RNNs - learning

- Backwards starting from the end of the sequence

$$\nabla_{\mathbf{h}(\tau)} L = \mathbf{V}^\top \nabla_{\mathbf{o}(\tau)} L$$

- Back-propagate gradients through time

$$\begin{aligned} \nabla_{\mathbf{h}^{(t)}} L &= \left( \frac{\partial \mathbf{h}^{(t+1)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{h}^{(t+1)}} L) + \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{h}^{(t)}} \right)^\top (\nabla_{\mathbf{o}^{(t)}} L) \\ &= \mathbf{W}^\top (\nabla_{\mathbf{h}^{(t+1)}} L) \text{diag} \left( 1 - \left( \mathbf{h}^{(t+1)} \right)^2 \right) + \mathbf{V}^\top (\nabla_{\mathbf{o}^{(t)}} L) \end{aligned}$$

Once the gradients on the internal nodes of the computational graph are obtained, we can obtain the gradients on the parameter nodes





# RNNs - learning

- For all the parameters

$$\nabla_{\mathbf{c}}L = \sum_t \left( \frac{\partial \mathbf{o}^{(t)}}{\partial \mathbf{c}} \right)^\top \nabla_{\mathbf{o}^{(t)}}L = \sum_t \nabla_{\mathbf{o}^{(t)}}L$$

$$\nabla_{\mathbf{b}}L = \sum_t \left( \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{b}^{(t)}} \right)^\top \nabla_{\mathbf{h}^{(t)}}L = \sum_t \text{diag} \left( 1 - \left( \mathbf{h}^{(t)} \right)^2 \right) \nabla_{\mathbf{h}^{(t)}}L$$

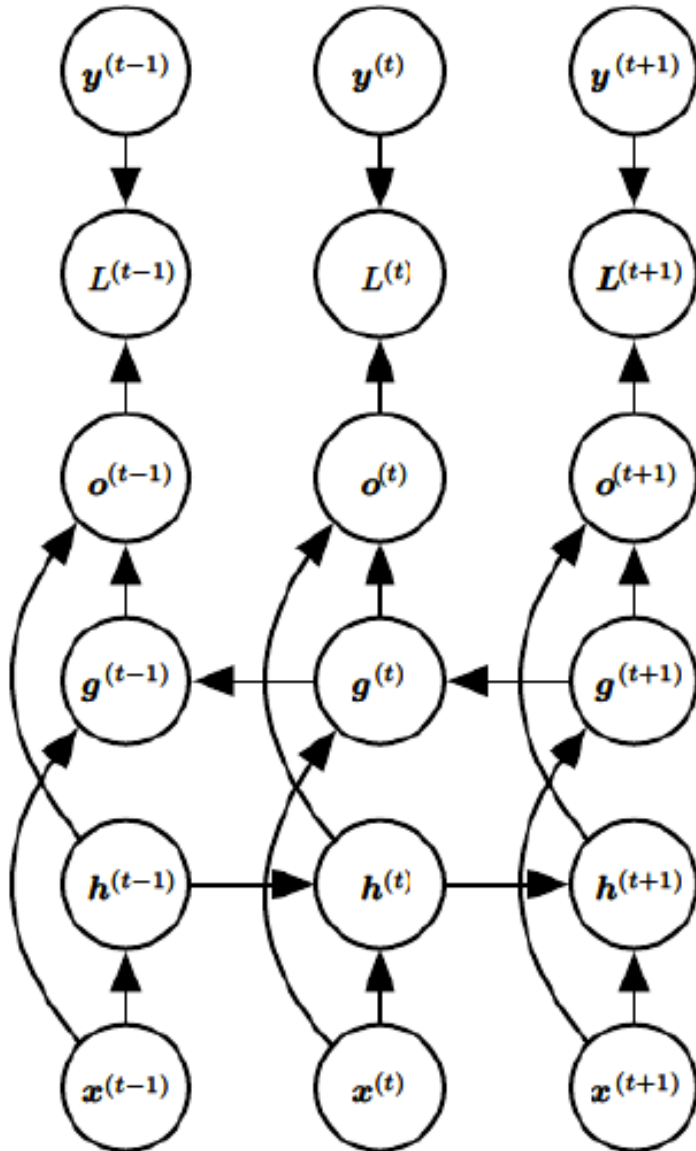
$$\nabla_{\mathbf{v}}L = \sum_t \sum_i \left( \frac{\partial L}{\partial o_i^{(t)}} \right) \nabla_{\mathbf{v}o_i^{(t)}} = \sum_t (\nabla_{\mathbf{o}^{(t)}}L) \mathbf{h}^{(t)\top}$$

$$\begin{aligned} \nabla_{\mathbf{w}}L &= \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{w}^{(t)}h_i^{(t)}} \\ &= \sum_t \text{diag} \left( 1 - \left( \mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}}L) \mathbf{h}^{(t-1)\top} \end{aligned}$$

$$\begin{aligned} \nabla_{\mathbf{U}}L &= \sum_t \sum_i \left( \frac{\partial L}{\partial h_i^{(t)}} \right) \nabla_{\mathbf{U}^{(t)}h_i^{(t)}} \\ &= \sum_t \text{diag} \left( 1 - \left( \mathbf{h}^{(t)} \right)^2 \right) (\nabla_{\mathbf{h}^{(t)}}L) \mathbf{x}^{(t)\top} \end{aligned}$$



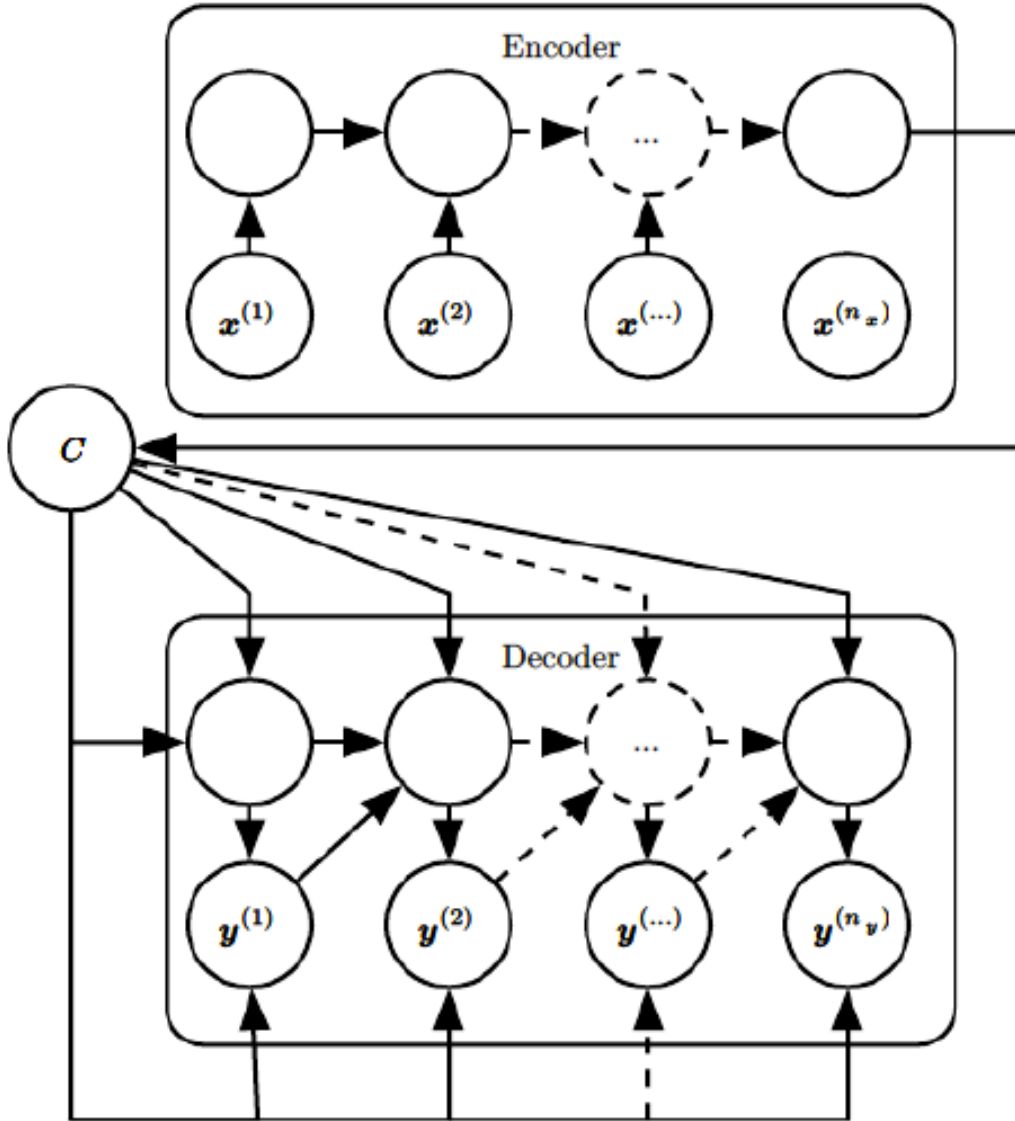
# RNNs - Bidirectional



prediction of  $y(t)$  which may depend on the whole input sequence e.g., speech recognition



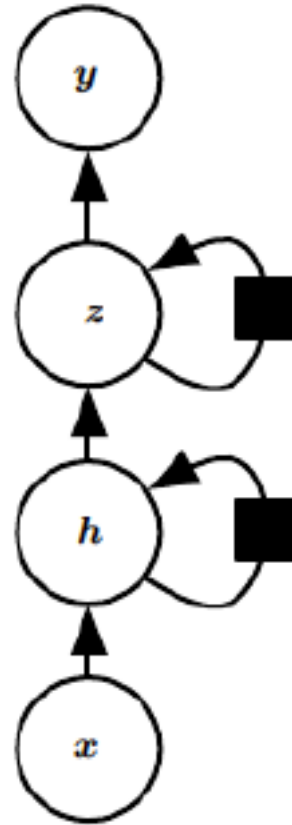
# RNNs - Bidirectional



encoder-decoder or  
sequence-to-sequence  
RNN architecture



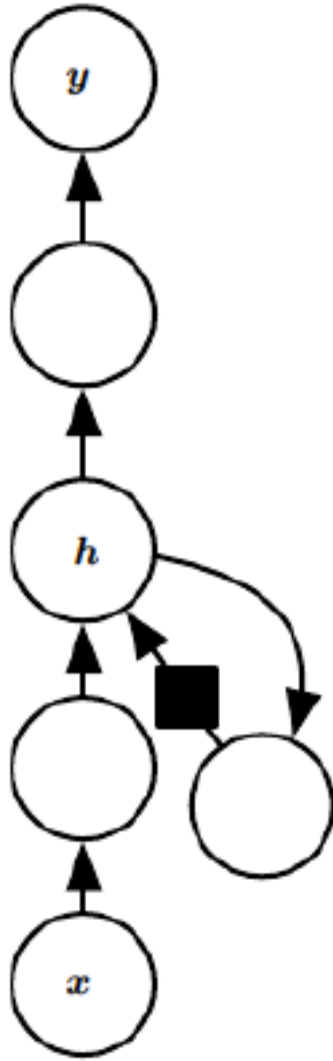
# Deep Recurrent Networks



The hidden recurrent state can be broken down into groups organized hierarchically



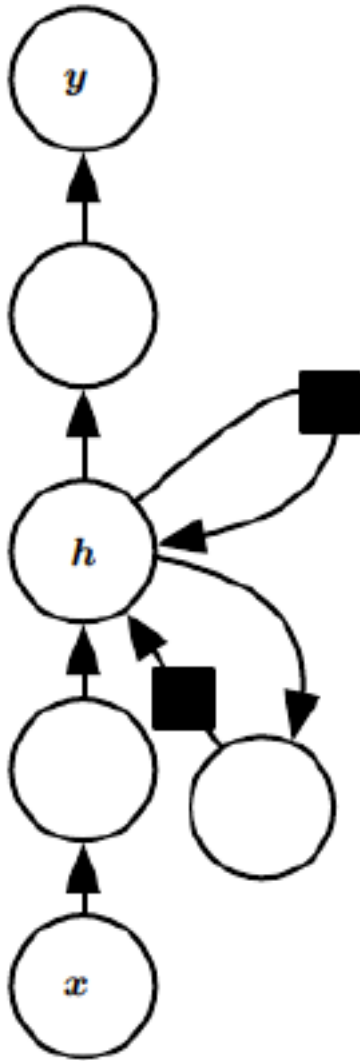
# Deep Recurrent Networks



Deeper computation (e.g., an MLP) can be introduced in the input-to-hidden, hidden-to-hidden and hidden-to-output parts. This may lengthen the shortest path linking different time steps.



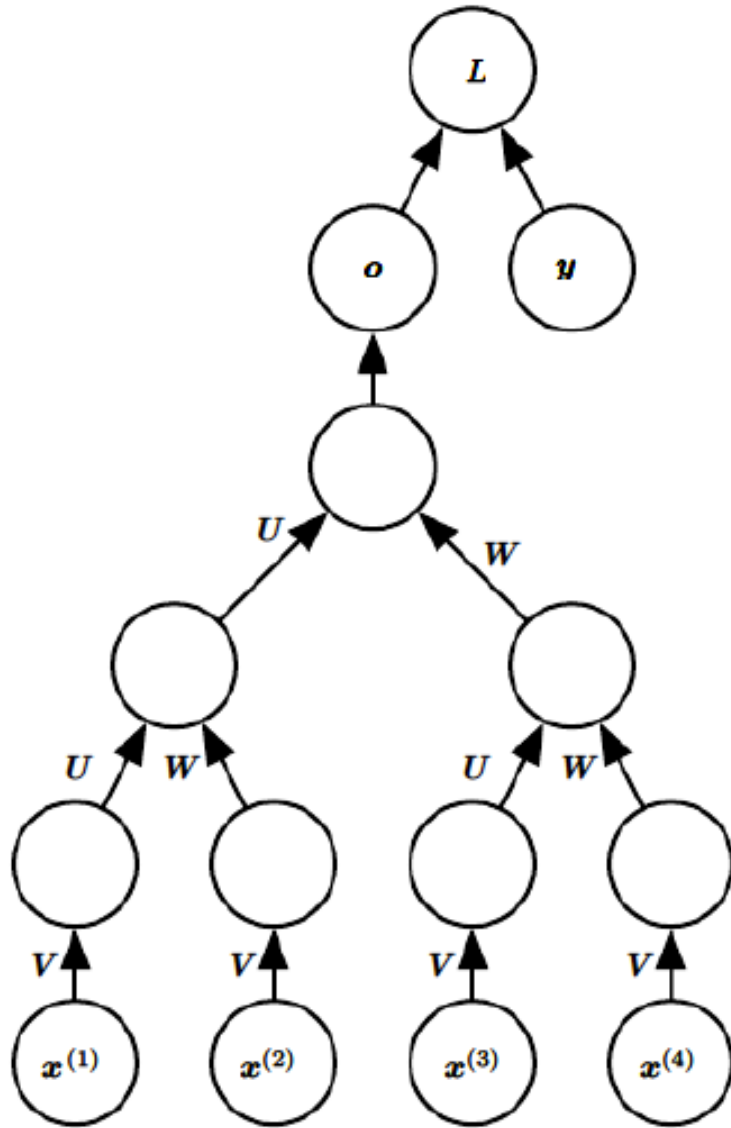
# Deep Recurrent Networks



The path-lengthening effect can be mitigated by introducing skip connections



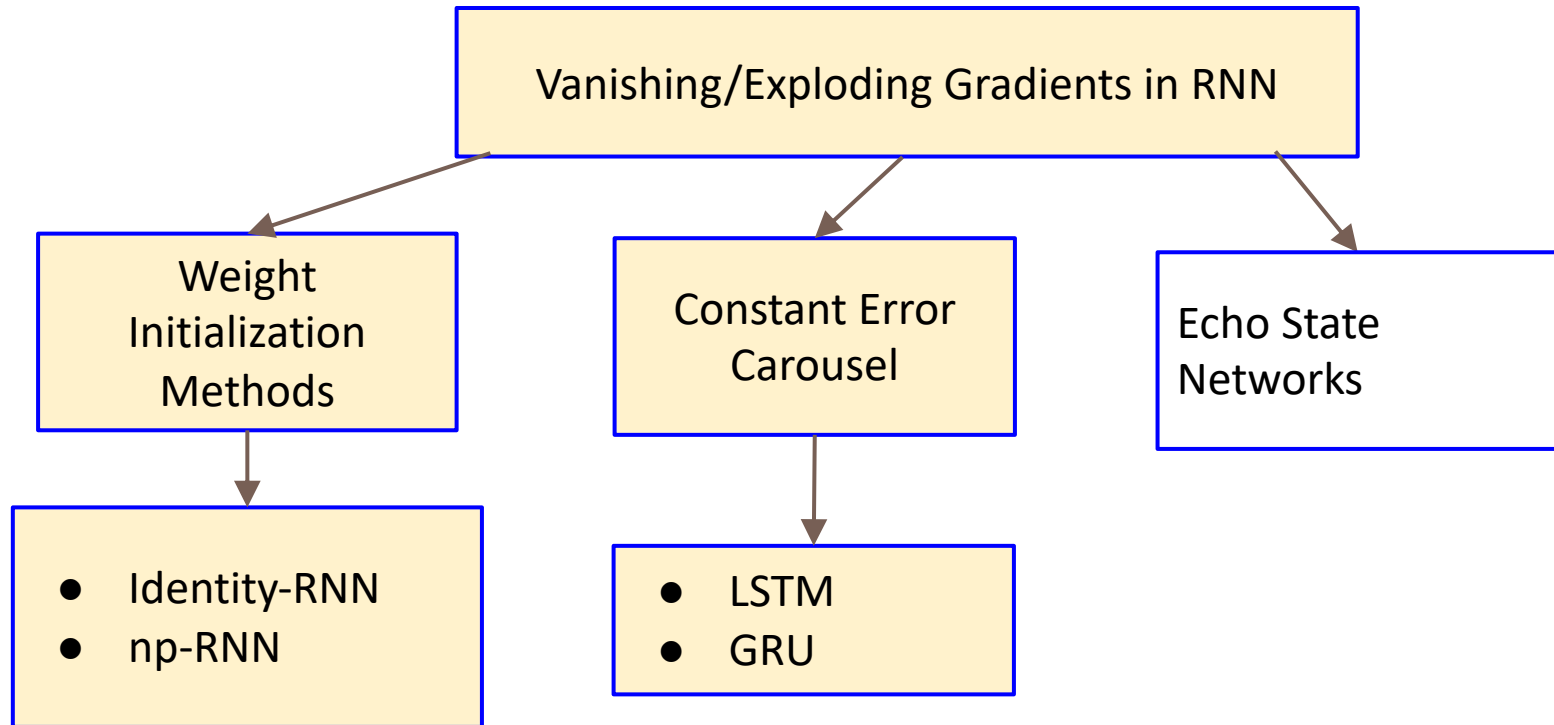
# Recursive NNs



Generalization of recurrent networks  
Applied for structured data



# Long-Term dependencies





# Long-Term dependencies

---

- **Random**  $W_h$  initialization of RNN has no constraint on eigenvalues
  - vanishing or exploding gradients in the initial epoch
- **Careful initialization** of  $W_h$  with suitable eigenvalues
  - allows the RNN to learn in the initial epochs
  - hence can generalize well for further iterations



# Long-Term dependencies

---

## ■ Trick #1 (IRNN)

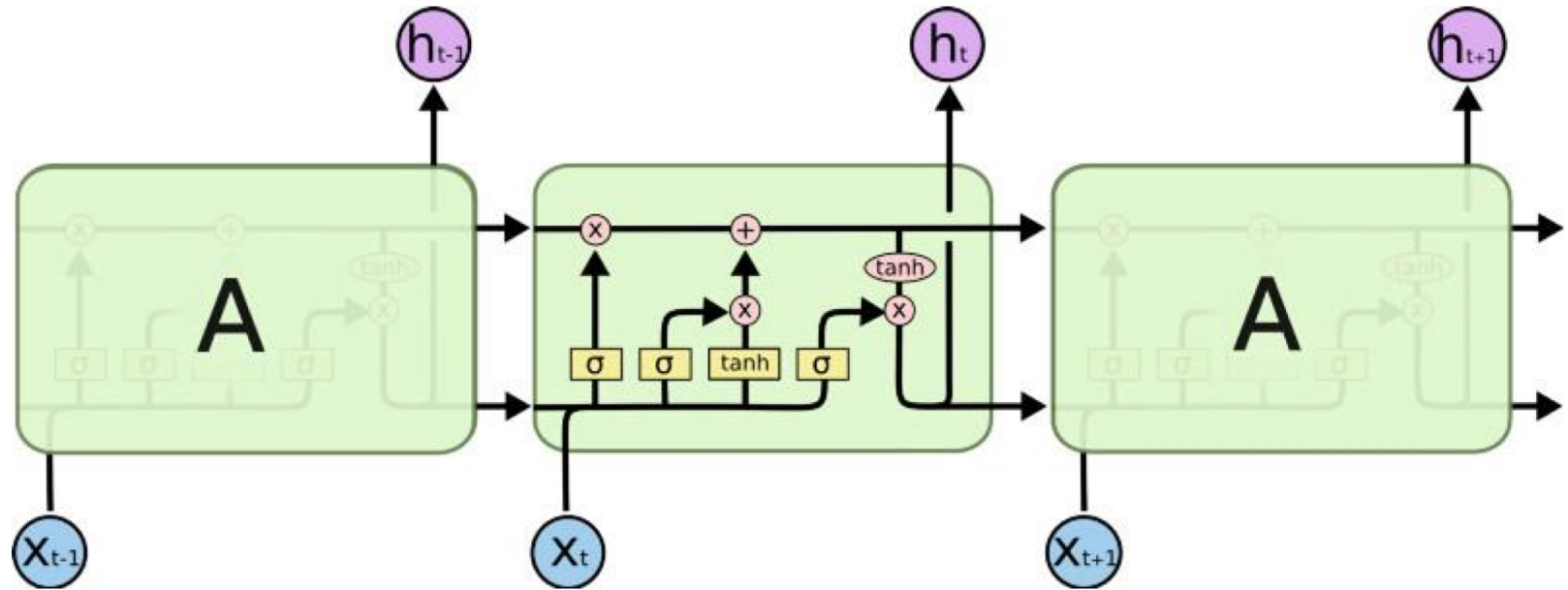
- $W_h$  initialized to Identity
- Activation function: ReLU

## ■ Trick# 2 (np-RNN)

- $W_h$  positive definite (+ve real eigenvalues)
- At least one eigenvalue is 1, others all less than equal to one
- Activation function: ReLU



# Long Short-Term Memory



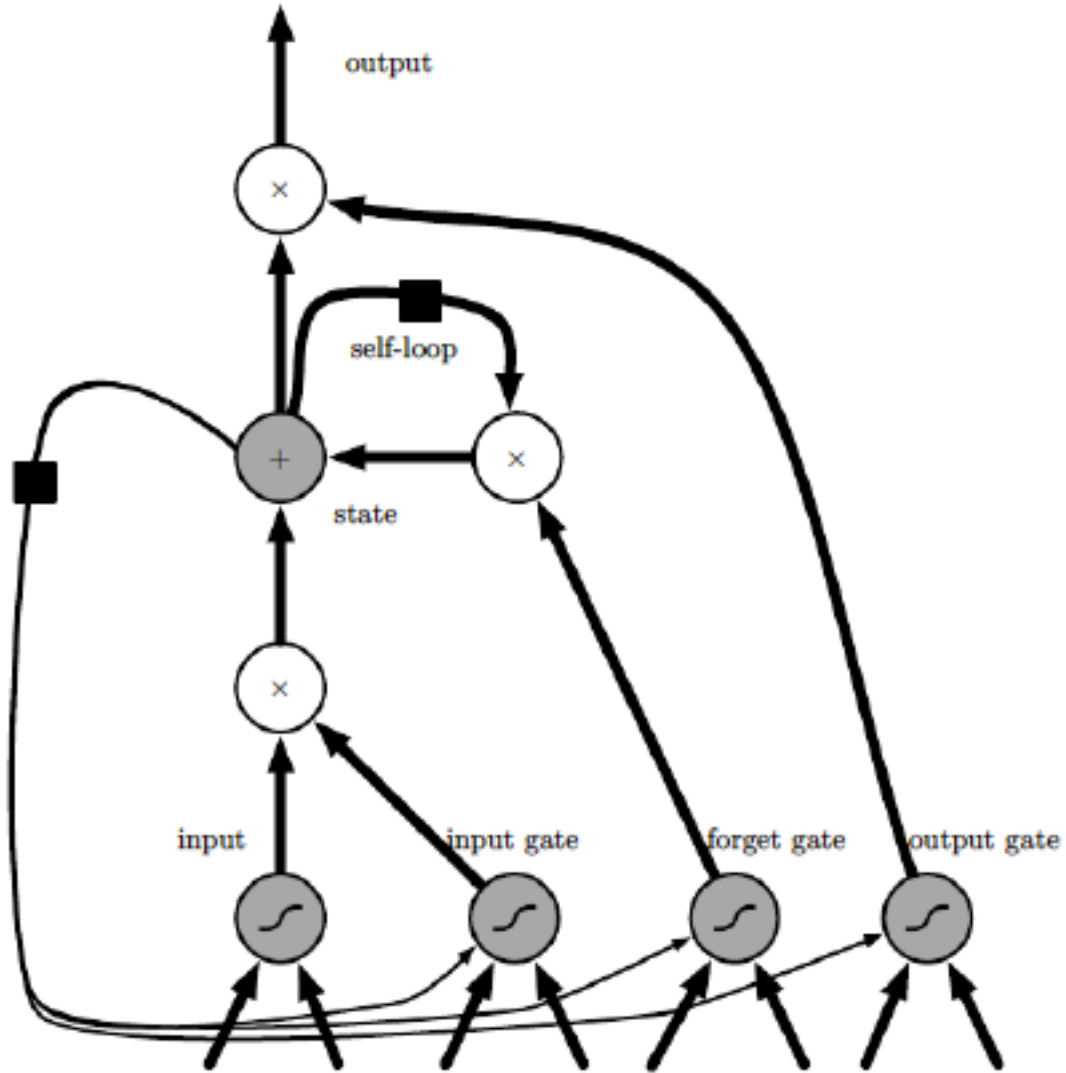
# Gated RNNs

---

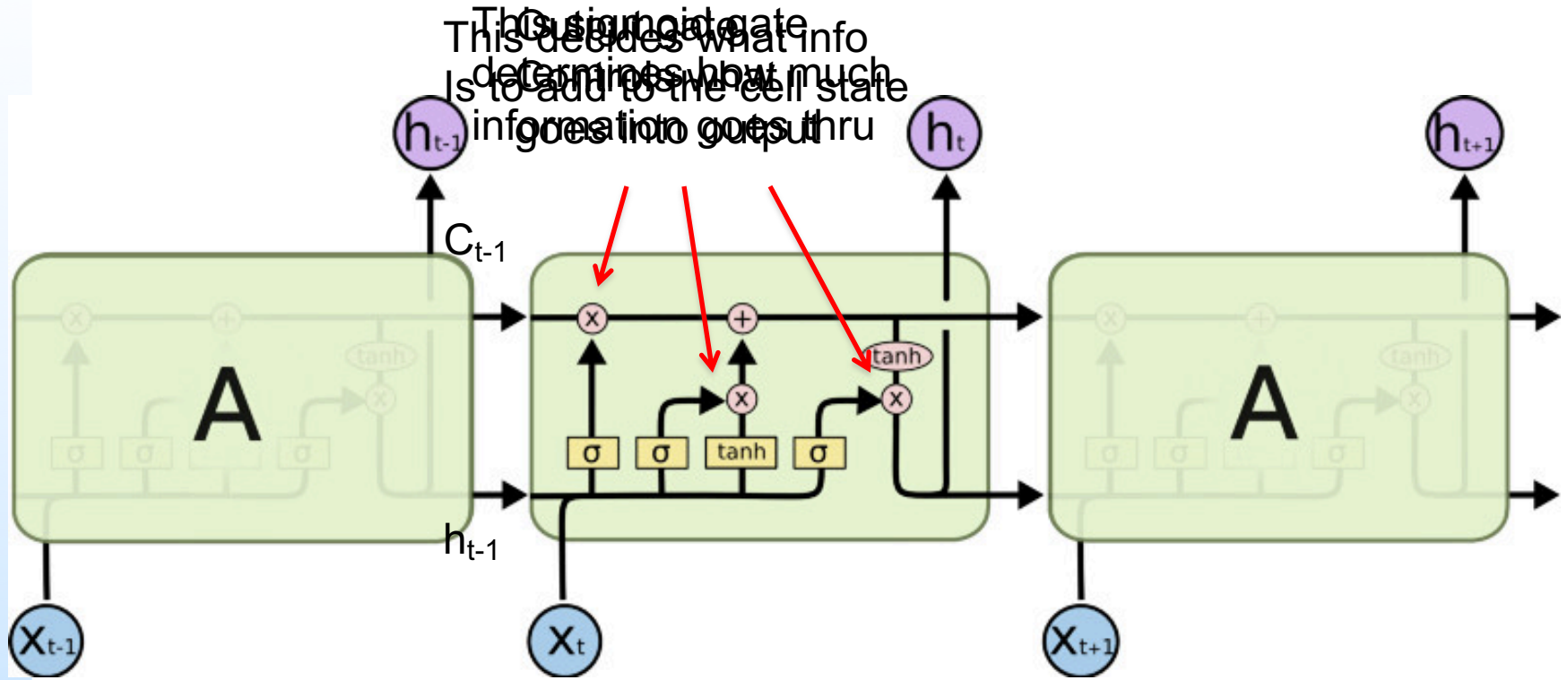
- Gated RNNs
  - Long Short-Term memory
  - Gated Recurrent Unit
- Idea
  - creating paths through time that have derivatives that neither vanish nor explode



# Gated RNNs

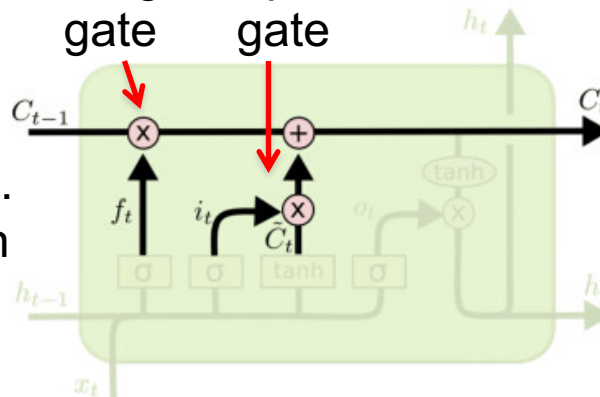


# LSTM cell



The sigmoid gate determines what info is to add to the cell state information goes thru

Forget gate input gate



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

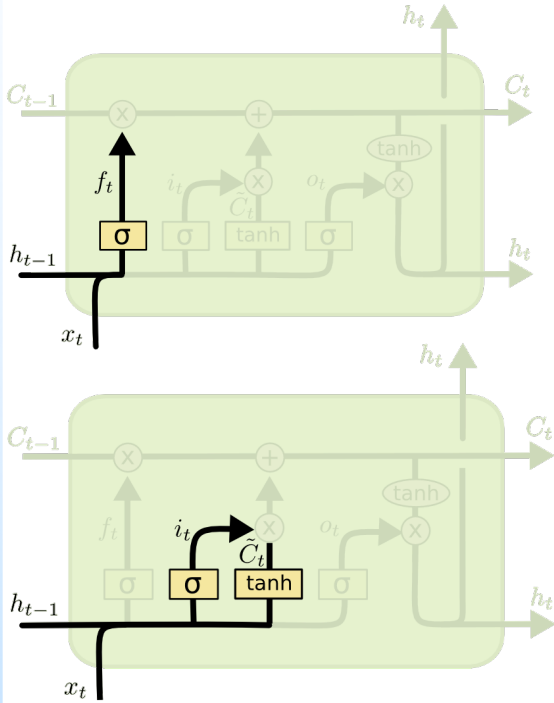
The core idea is this cell state  $C_t$ , it is changed slowly, with only minor vanishing gradient problem in linear interactions. It is very easy for information to flow along it unchanged.

Why sigmoid or tanh:

- Sigmoid: 0-1 gating as switch.



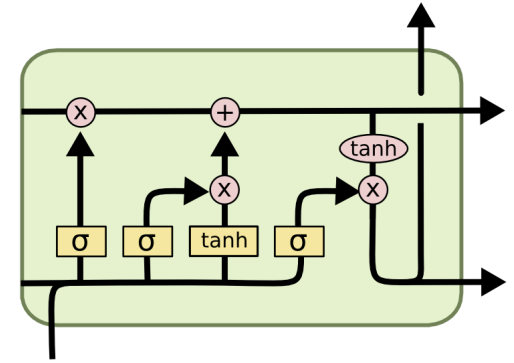
# LSTM cell



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

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

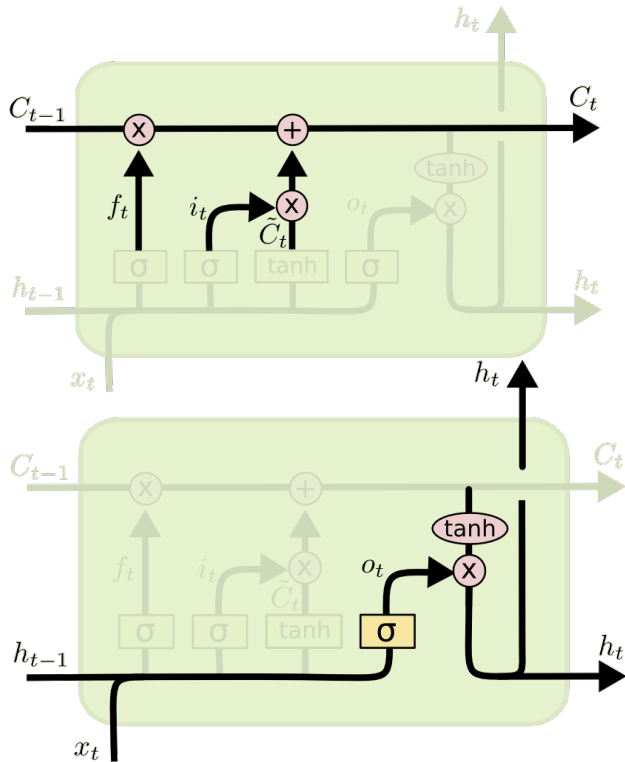
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



$i_t$  decides what component is to be updated.  
 $C'_t$  provides change contents



# LSTM cell



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

Updating the cell state

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

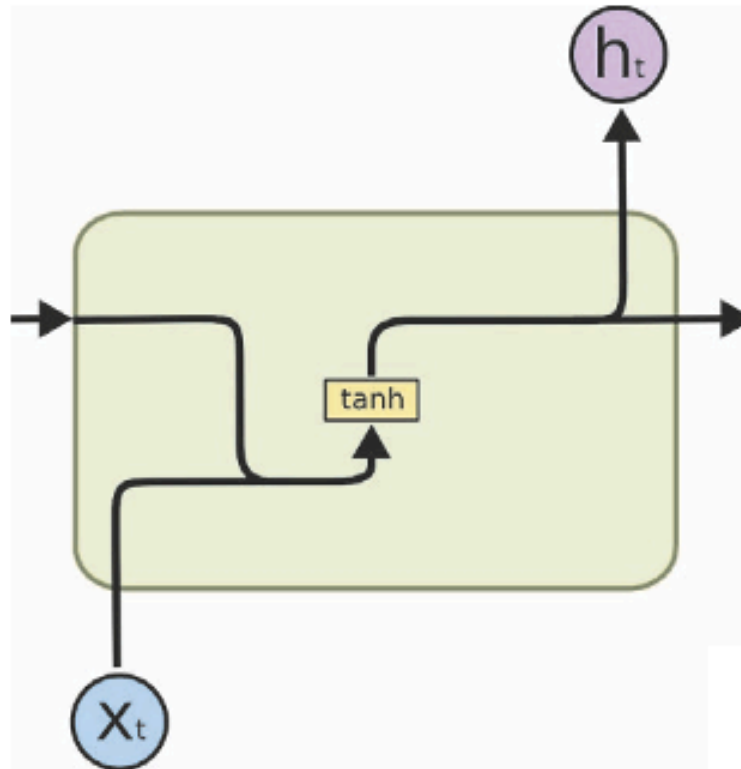
Decide what part of the cell state to output

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

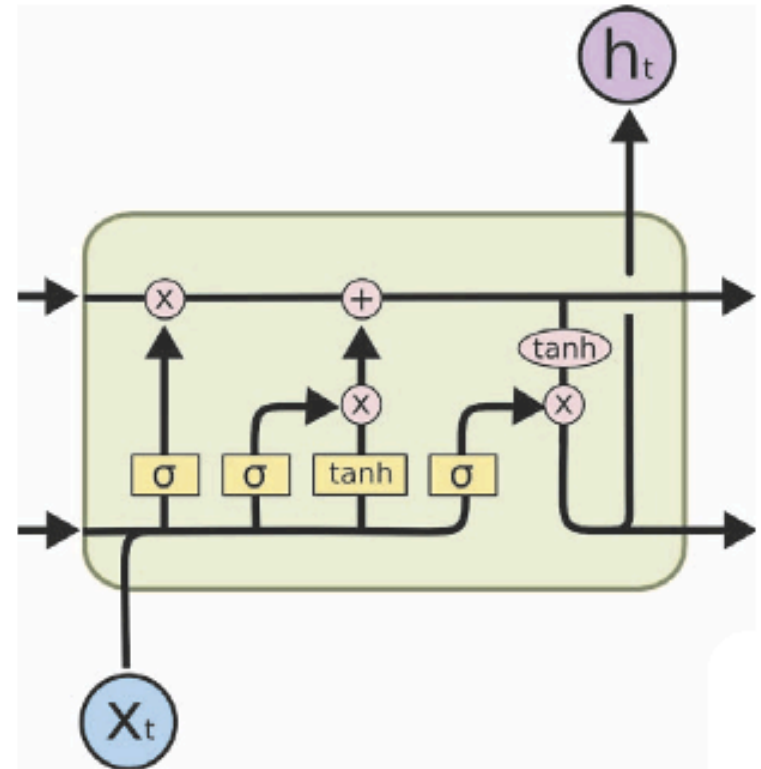




# RNN vs LSTM



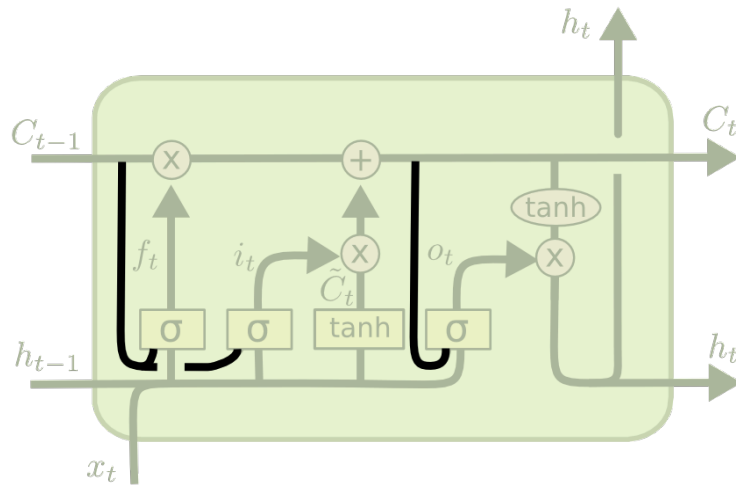
(a) RNN



(b) LSTM



# Peephole LSTM



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

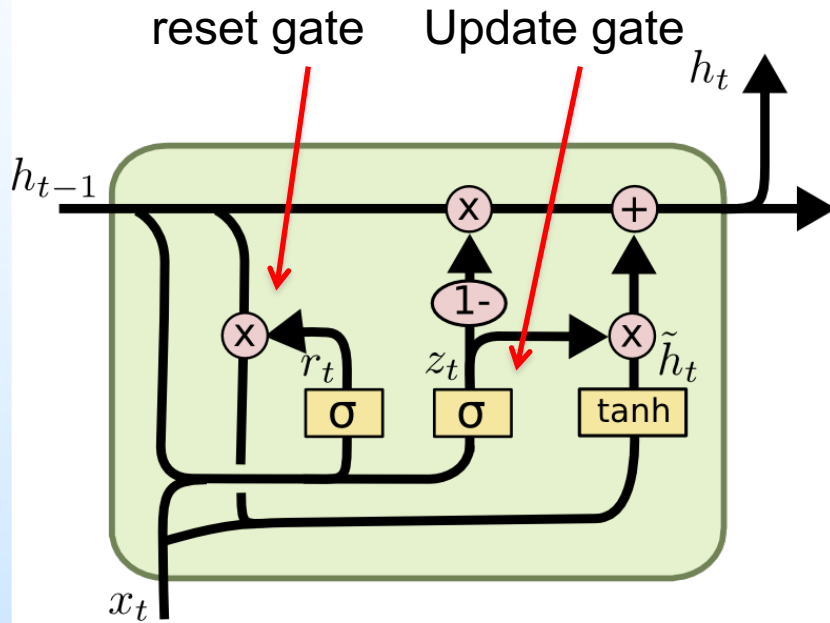
$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

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

Allows “**peeping** into the memory”. Can learn the fine distinction between sequences of spikes separated by either 50 or 49 discrete time steps



# Gated Recurrent Unit (GRU)



$$z_t = \sigma (W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma (W_r \cdot [h_{t-1}, x_t])$$

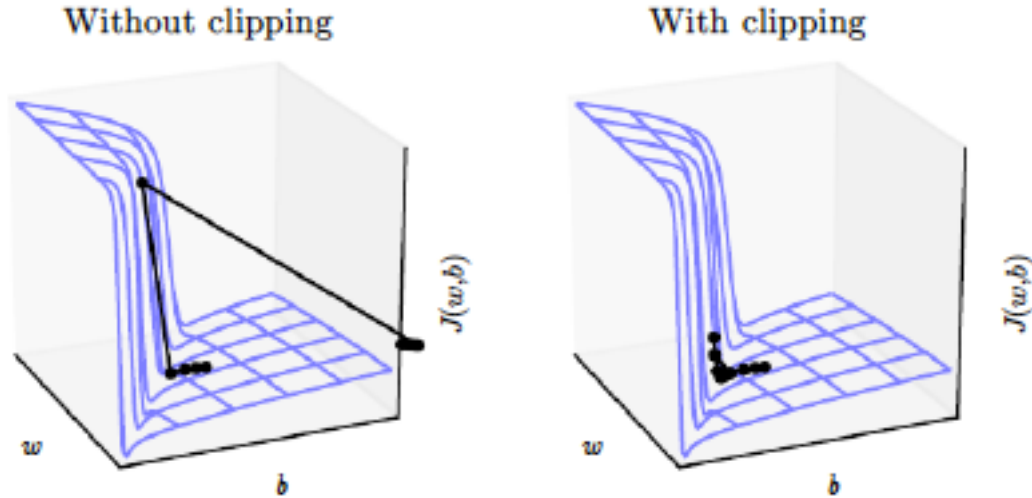
$$\tilde{h}_t = \tanh (W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

It combines the **forget** and **input** into a single **update gate**. It also merges the cell state and hidden state. This is simpler than LSTM. There are many other variants too.



# Clipping gradients



parameter gradient is very large

“landscape” in which one finds “cliffs”

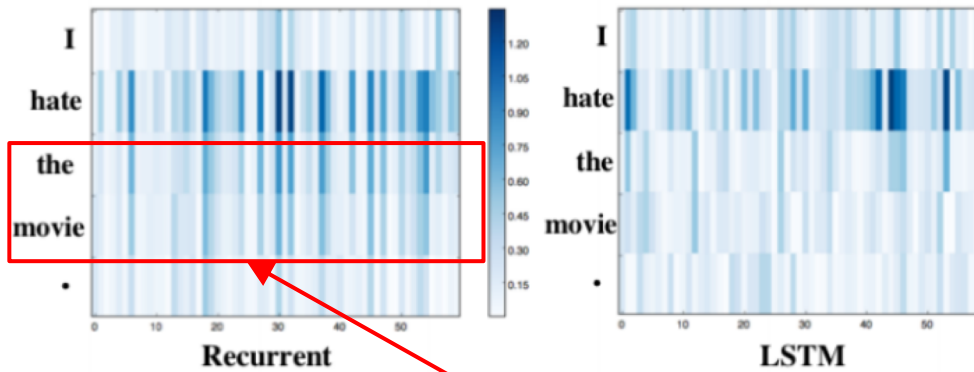
$$\text{if } \|g\| > v \\ g \leftarrow \frac{gv}{\|g\|}$$

Clipping the gradient



# RNN vs LSTM

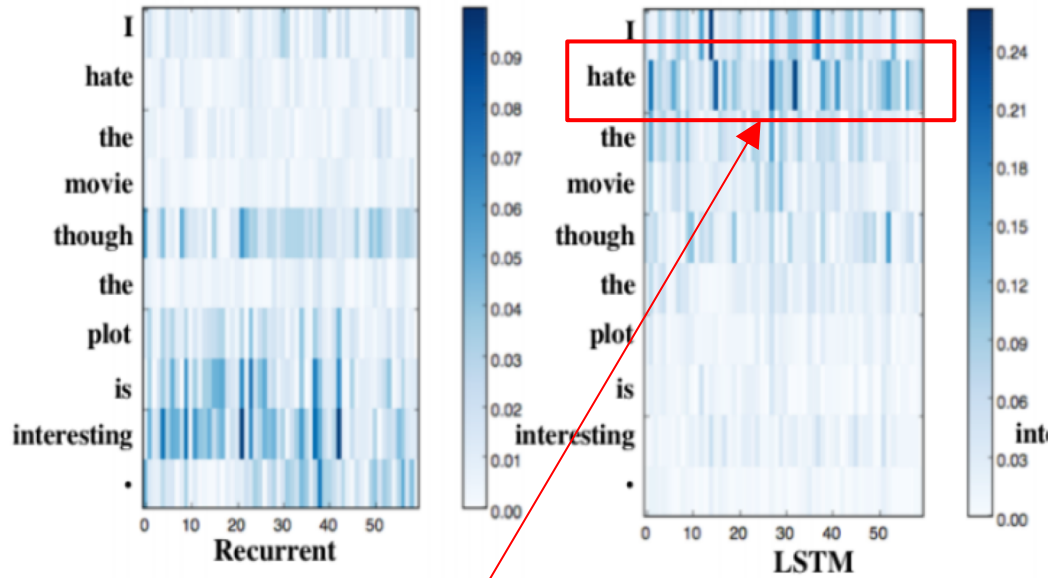
Saliency Heatmap



Recent words more salient

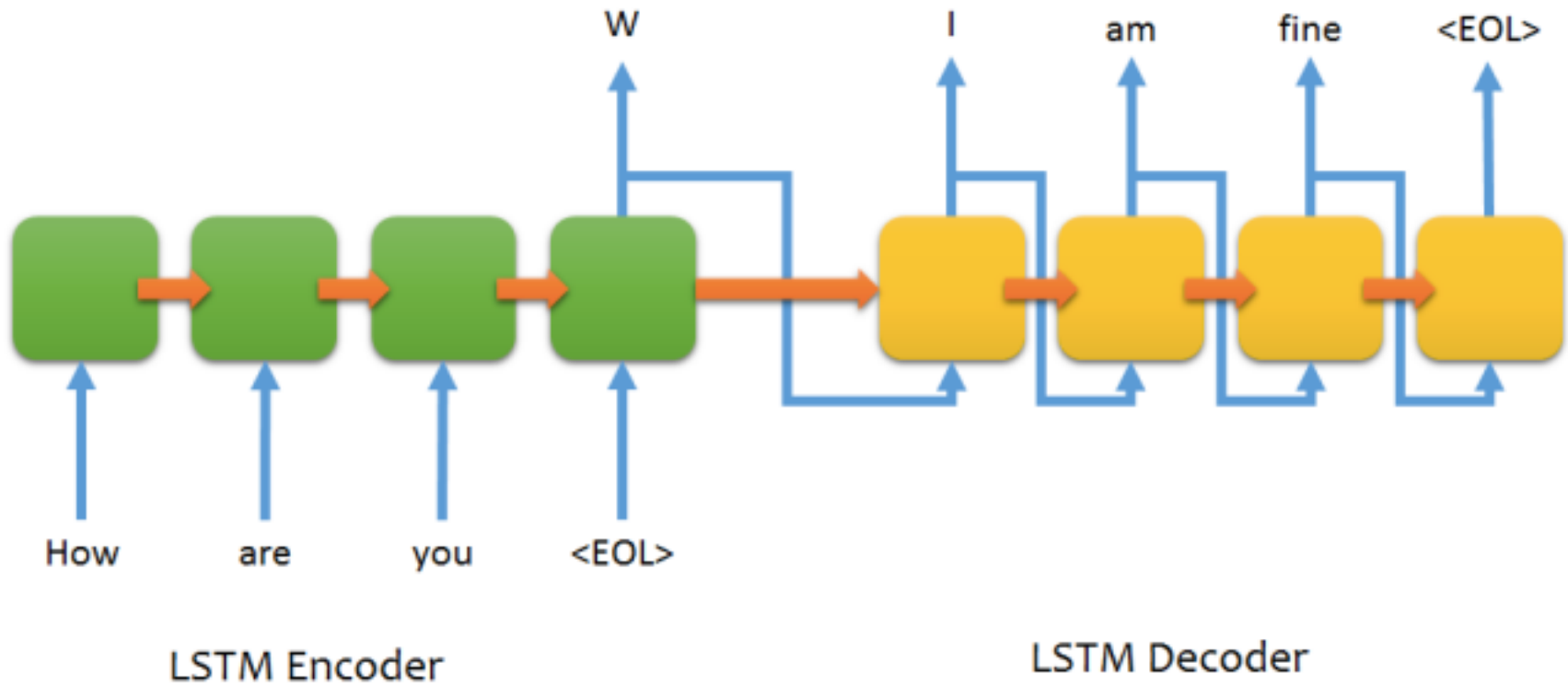
# RNN vs LSTM

Saliency Heatmap



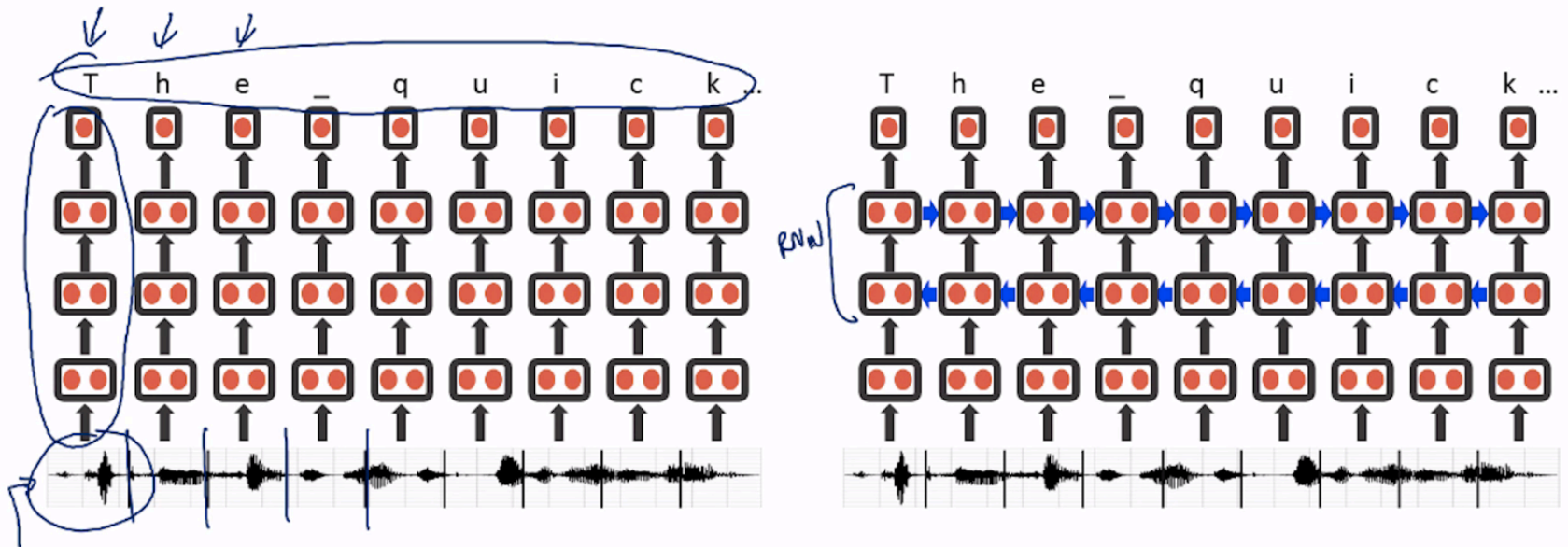
LSTM captures long term dependencies

# Sequence to sequence chat model



# Speech recognition RNN

Speech recognition example (Deep Speech)





# Reservoir computing

---

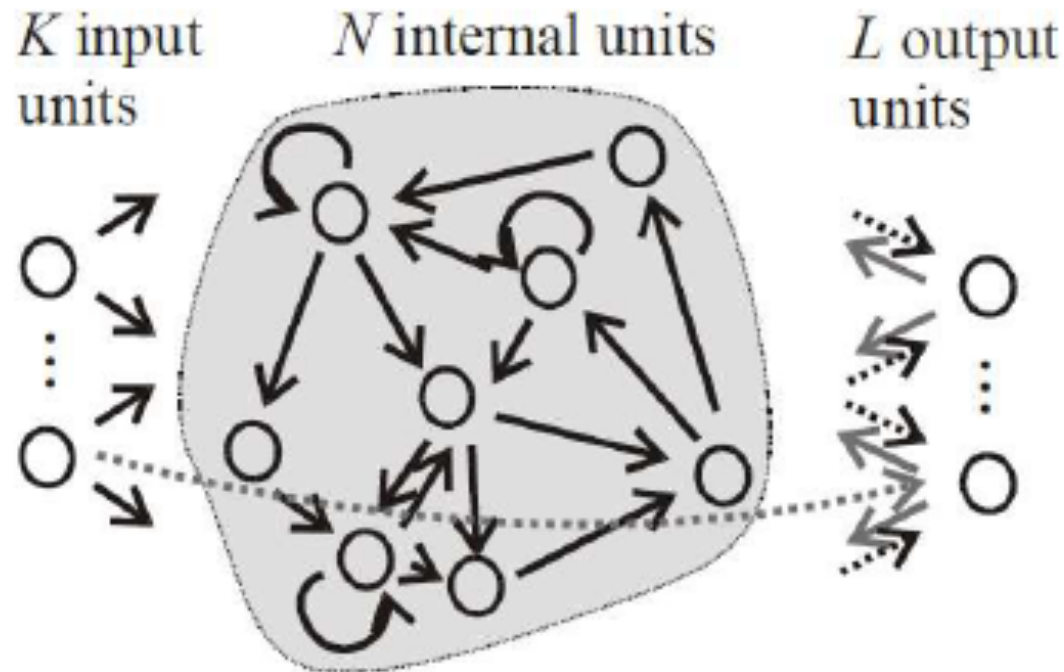
- The equivalent idea for RNNs
  - fix the input- hidden connections and the hidden-hidden connections at random values
  - only learn the hidden-output connections
- The learning is then very simple (assuming linear output units)
- Its important to set the random connections very carefully so the RNN does not explode or die
- See also Liquid State Machine



# Reservoir computing

Herbert Jaeger, 2001

## Echo State Network



# Reservoir computing

