

# Sequence Modeling: Recurrent and Recursive Nets

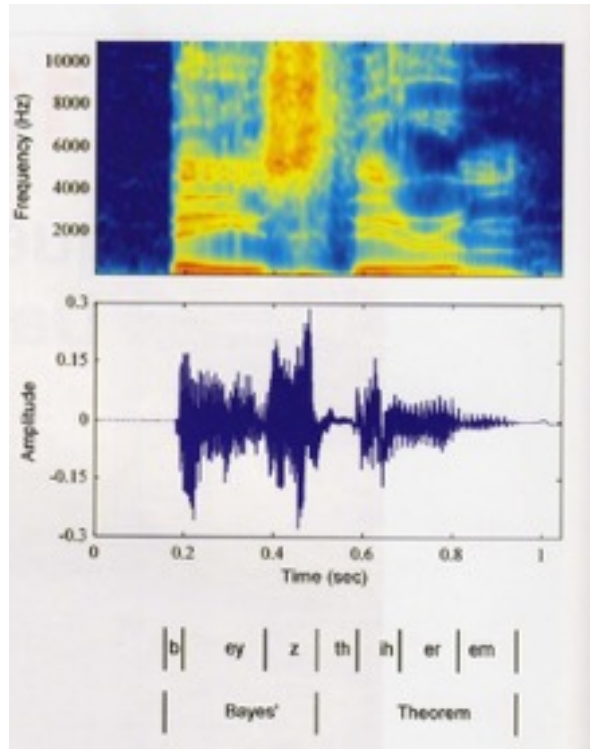
1. Sequential Data and RNN Overview
2. Unfolding Computational Graphs
3. Recurrent Neural Networks
4. Bidirectional RNNs
5. Encoder-Decoder Sequence-to-Sequence Architectures
6. Deep Recurrent Networks
7. Recursive Neural Networks
8. The Challenge of Long-Term Dependencies
9. Echo-State Networks
10. Leaky Units and Other Strategies for Multiple Time Scales
  
10. LSTM and Other Gated RNNs
11. Optimization for Long-Term Dependencies
12. Explicit Memory

## Sequential Data Examples

- Often arise through measurement of time series
  - Acoustic features at successive time frames in speech recognition
  - Sequence of characters in an English sentence
  - Parts of speech of successive words
  - Snowfall measurements on successive days
  - Rainfall measurements on successive days
  - Daily values of currency exchange rate
  - Nucleotide base pairs in a strand of DNA

# Sound Spectrogram to Word Sequence

Bayes Theorem



- Decompose sound waves into frequency, amplitude using Fourier transforms
- Plot of the intensity of the spectral coefficients versus time index
- Successive observations of speech spectrum highly correlated (Markov dependency)

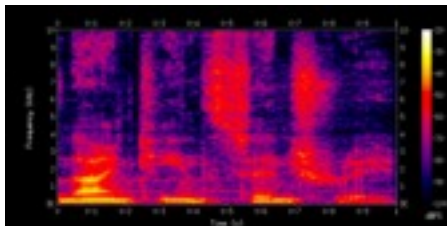
# Two common tasks with sequential data

## 1. Sequence-to-sequence

### 1. Speech recognition using a sound spectrogram

- decompose sound waves into frequency, amplitude using Fourier transforms

→ “Nineteenth Century”



Frequencies increase up the vertical axis, and time on the horizontal axis. The lower frequencies are more dense because it is a male voice. Legend on right shows that the color intensity increases with density

### 1. NLP: Named Entity Recognition

- **Input:** Jim bought 300 shares of Acme Corp. in 2006
- **NER:** [Jim]<sub>Person</sub> bought 300 shares of [Acme Corp.]<sub>Organization</sub> in [2006]<sub>Time</sub>

### 2. Machine Translation: Echte dicke kiste → Awesome sauce

## 2. Sequence-to-symbol

### 1. Sentiment:

- *Best movie ever* → Positive

### 2. Speaker recognition

- Sound spectrogram → Harry

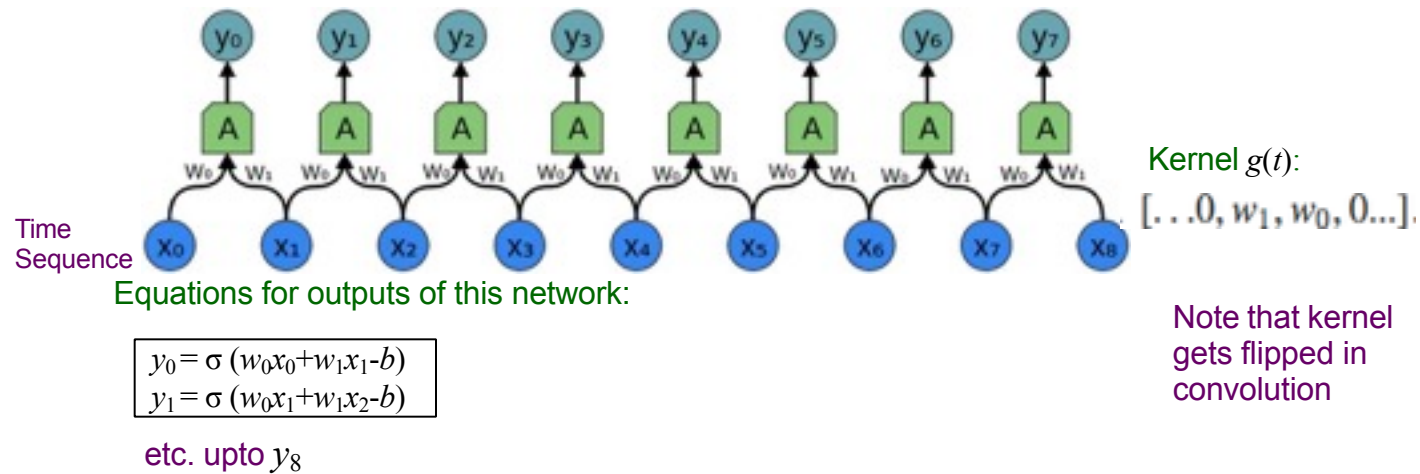
## Recurrent Neural Networks process sequential data

- RNNs are a family of neural nets for sequential data
- Analogy with Convolutional Neural Networks
  - Specialized architectures
    - CNN is specialized for grid of values, e.g., image
    - RNN is specialized for a sequence of values  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}$
  - Scaling & Variable length
    - CNNs readily scale to images with large width/height
      - CNNs can process variable size images
    - RNNs scale to longer sequences than would be practical for networks without sequence-based specialization
      - RNNs can also process variable-length sequences

## RNNs share same weights across Time Steps

- To go from multi-layer networks to RNNs:
  - Need to share parameters across different parts of a model
  - Separate parameters for each value of cannot generalize to sequence lengths not seen during training
  - Share statistical strength across different sequence lengths and across different positions in time
- Sharing important when information can occur at multiple positions in the sequence
  - Given “*I went to Nepal in 1999*” and “*In 1999, I went to Nepal*”, an ML method to extract year, should extract 1999 whether in position 6 or 2
  - A feed-forward network that processes sentences of fixed length would have to learn all of the rules of language separately at each position
  - An RNN shares the same weights across several time steps

# 1-D CNN used with a time sequence



We can also write the equations in terms of elements of a general  $8 \times 8$  weight matrix  $W$  as:

$$\begin{aligned} y_0 &= \sigma(W_{0,0}x_0 + W_{0,1}x_1 + W_{0,2}x_2\dots) \\ y_1 &= \sigma(W_{1,0}x_0 + W_{1,1}x_1 + W_{1,2}x_2\dots) \end{aligned}$$

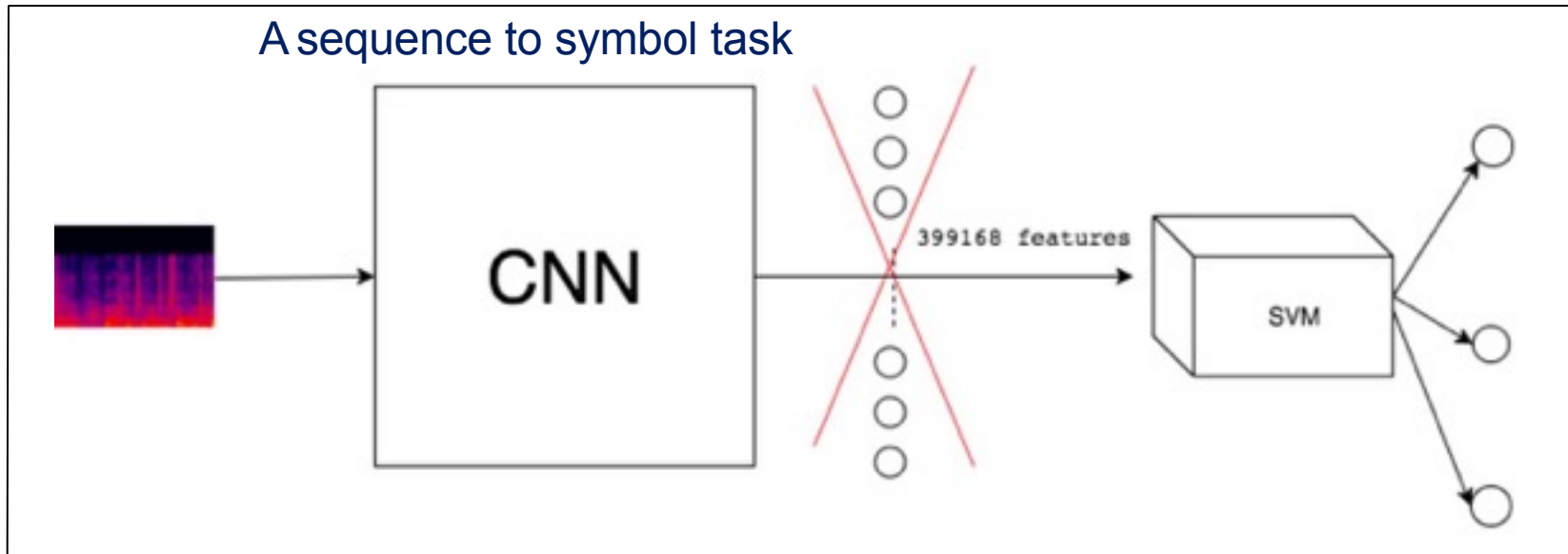
etc. upto  $y_8$

where

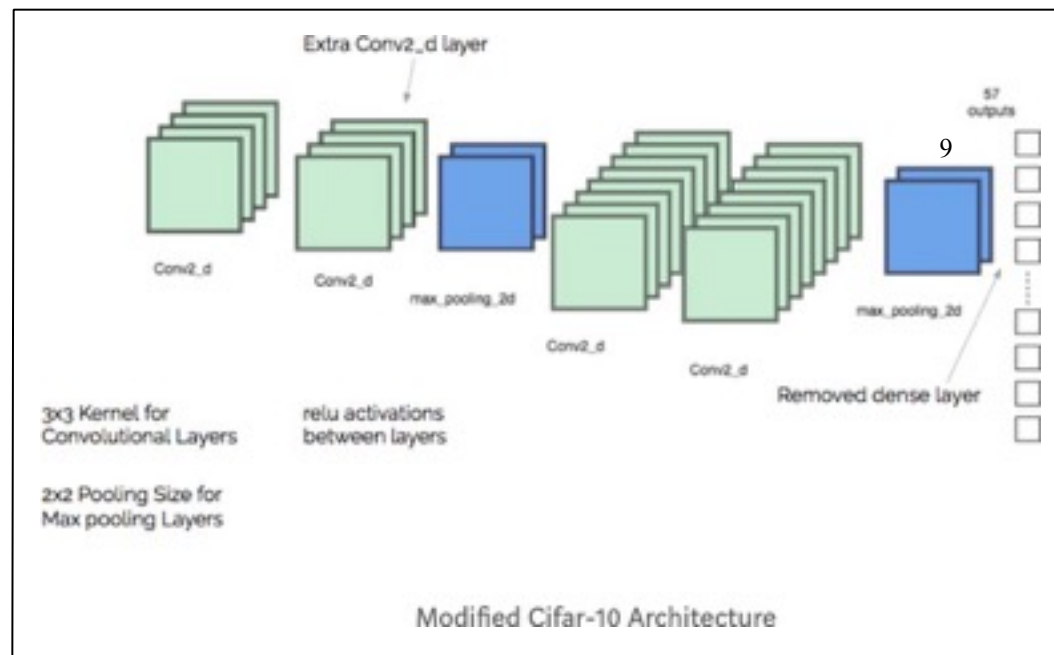
$$W = \begin{bmatrix} w_0 & w_1 & 0 & 0 & \dots \\ 0 & w_0 & w_1 & 0 & \dots \\ 0 & 0 & w_0 & w_1 & \dots \\ 0 & 0 & 0 & w_0 & \dots \\ \dots & \dots & \dots & \dots & \dots \end{bmatrix}$$



# Speaker Recognition with CNN

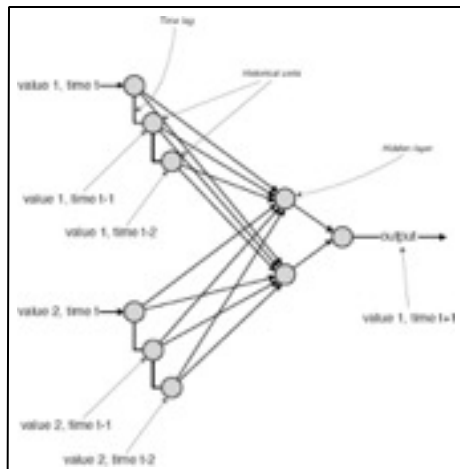


CNN  
Feature Extractor  
(Transfer Learning)



## Time Delay Neural Networks (TDNNs)

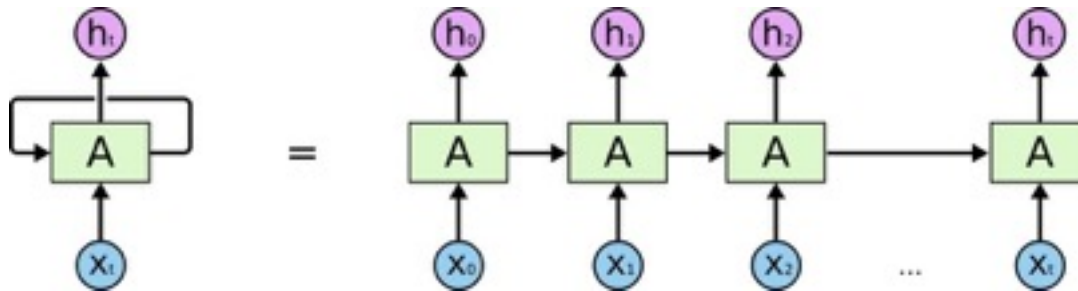
- TDNNs use convolution for a 1-D temporal sequence
  - Convolution allows shared parameters across time, but is shallow
    - Each output is dependent upon a small no. of neighboring inputs
    - Parameter sharing manifests in the application of the same convolutional kernel at each time step



A TDNN remembers the previous few training examples and uses them as input into the network. The network then works like a feed-forward, back propagation network.

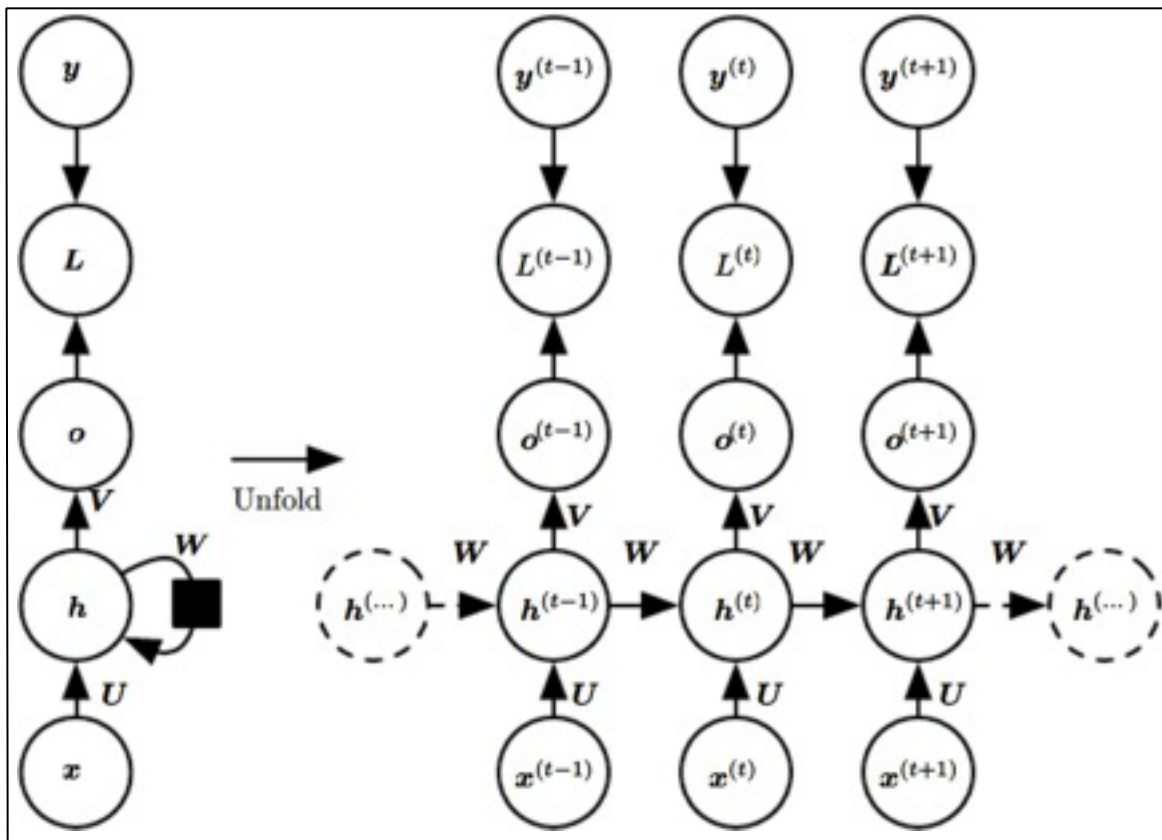
## RNN vs. TDNN

- RNNs share parameters in a different way than TDNN
  - Each member of output is a function of previous members of output
  - Each output produced using same update rule applied to previous outputs
  - This recurrent formulation results in sharing of parameters through a very deep computational graph
- An unrolled RNN



# Computational Graphs for RNNs

- We extend computational graphs to include cycles
  - Cycles represent the influence of the present value of a variable on its own value at a future time step
  - In a Computational graph nodes are variables/operations
  - RNN to map input sequence of  $x$  values to output sequence of  $o$  values
    - Loss  $L$  measures how far each output  $o$  is from the training target  $y$



• Forward propagation is given as follows:

For each time step  $t$ ,  $t=1$  to  $t=\tau$

Apply the following equations

$$o(t) = c + Vh(t)$$

$$h(t) = \tanh^{12}(a(t))$$

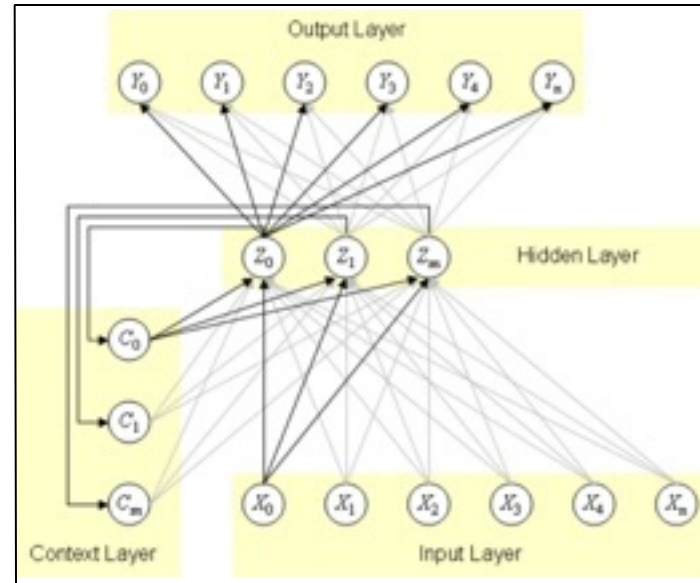
$$a(t) = b + Wh(t-1) + Ux(t)$$

## RNN operating on a sequence

- RNNs operate on a sequence that contain vector  $\mathbf{x}^{(t)}$  with time step index  $t$ , ranging from 1 to  $\tau$ 
  - Sequence:  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)}$
  - RNNs operate on minibatches of sequences of length  $\tau$
- Some remarks about sequences
  - The steps need not refer to passage of time in the real world
  - RNNs can be applied in two-dimensions across spatial data such as image
  - Even when applied to time sequences, network may have connections going backwards in time, provided entire sequence is observed before it is provided to network

## RNN as a network with cycles

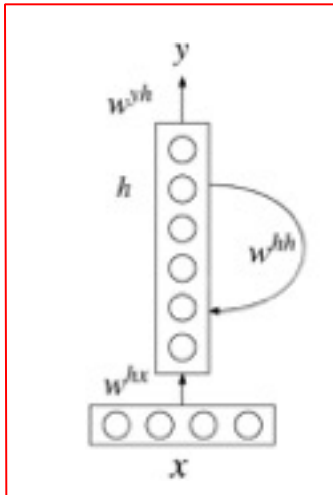
- An RNN is a class of neural networks where connections between units form a directed cycle
- This creates an internal state of the network which allows it to exhibit dynamic temporal behavior
- The internal memory can be used to process arbitrary sequences of inputs



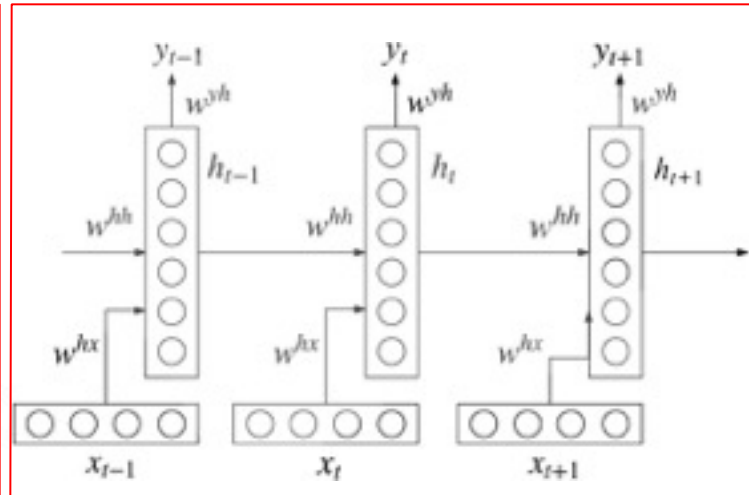
Three layer network with input  $\mathbf{x}$ , hidden layer  $\mathbf{z}$  and output  $\mathbf{y}$ .  
Context units  $\mathbf{c}$  maintain a copy of the previous value of the hidden units

# RNN parameters

Folded network  
with cycles



Unfolded sequence network with three time steps

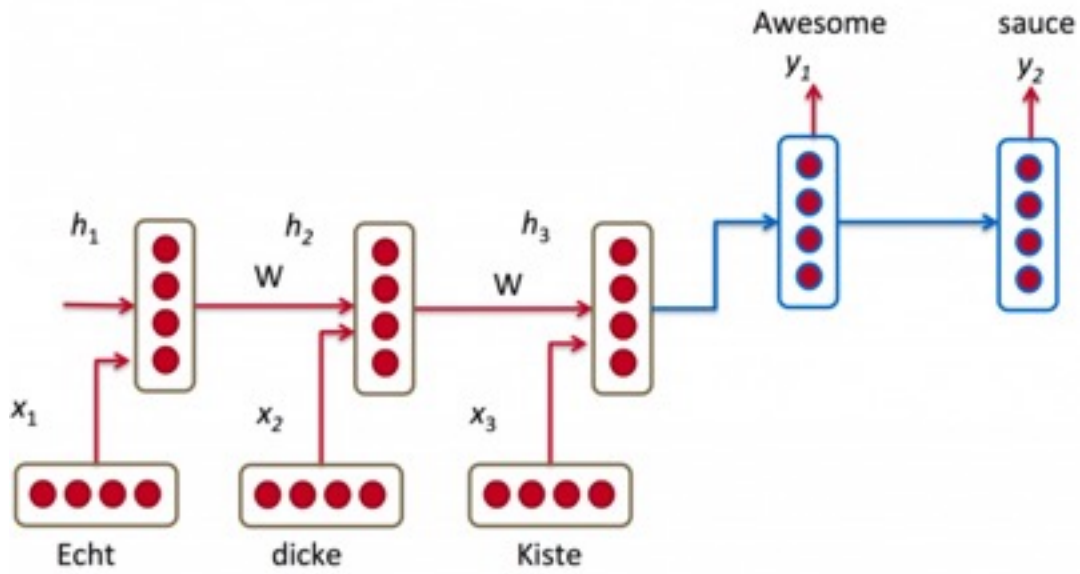


$$h_t = f(w^{hh}h_{t-1} + w^{hx}x_t)$$

$$y_t = \text{softmax}(w^{yh}h_t)$$

Unlike a feedforward neural network, which uses different parameters at each layer, RNN shares the same parameters ( $w^{hx}$ ,  $w^{hh}$ ,  $w^{yh}$ ) across all steps

## RNN for Machine Translation



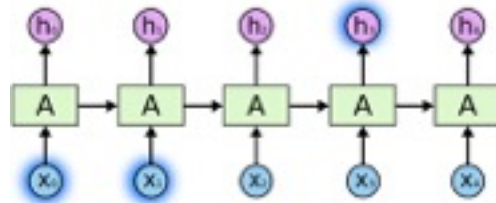


## Limitation of length of time steps

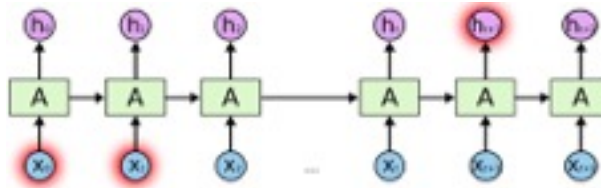
- Length of time steps is determined by length of input
  - e.g., if word sequence to be processed is a sentence of six words, RNN would be unfolded into a neural net with six time steps or layers.
    - One layer corresponds to a word
- Theoretically, RNN can make use of the information in arbitrarily long sequences
  - In practice, RNN is limited to looking back only a few steps due to the vanishing gradient or exploding gradient problem

## Problem of Long-Term Dependencies

- Easy to predict last word in “the clouds are in the *sky*,”
  - When gap between relevant information and place that it’s needed is small, RNNs can learn to use the past information

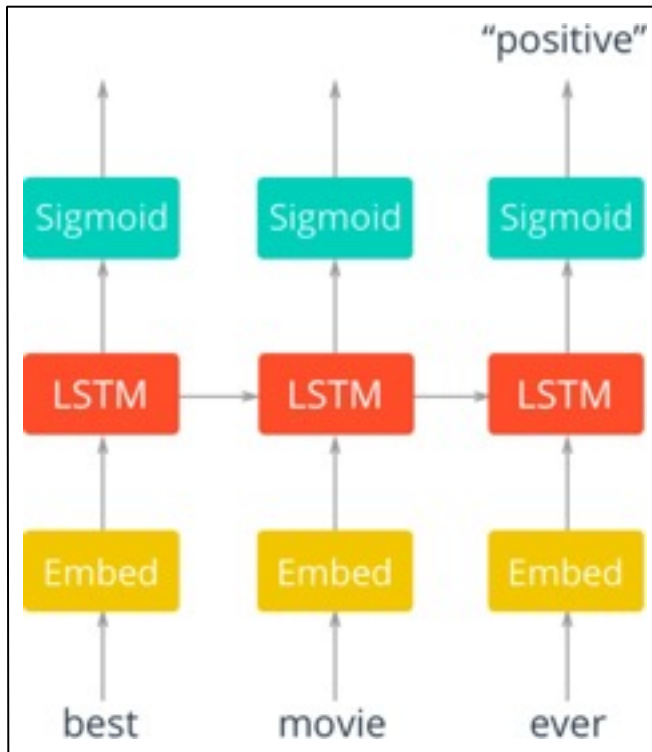


- “I grew up in France... I speak fluent *French*.”
  - We need the context of France, from further back.
  - Large gap between relevant information and point where it is needed



- In principle RNNs can handle it, but fail in practice
  - LSTMs offer a solution

## RNN for sentiment analysis



### Embedding Layer

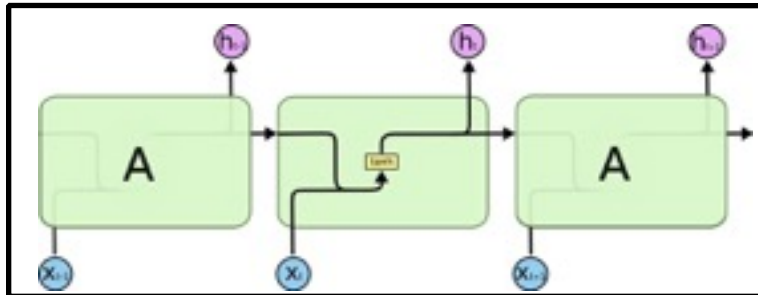
We pass in words to an embedding layer

### Options

1. Actually train up an embedding with Word2vec
2. It is good enough to just have an embedding layer and let network learn the embedding table on its own.

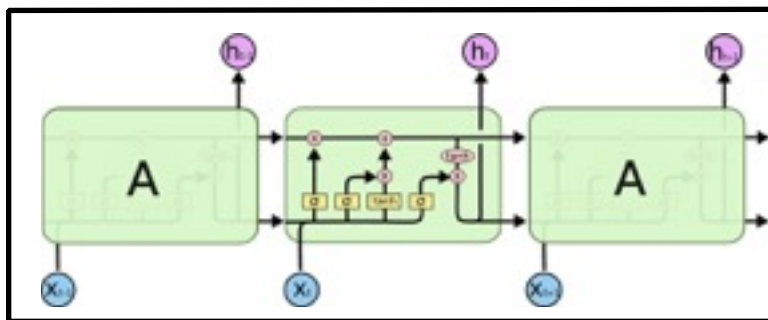
LSTM:

# RNN vs LSTM



## RNN

Repeating module has a simple structure Such as a tanh layer

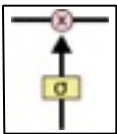


## LSTM

Repeating module has four interacting Layers, with notation:

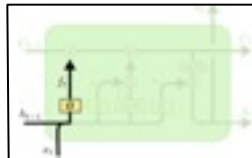


### Three gates of the form



Sigmoid is 0 or 1  
Allows input to go through or not

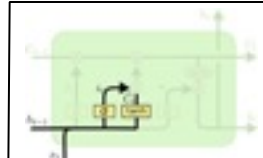
### Forget gate



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

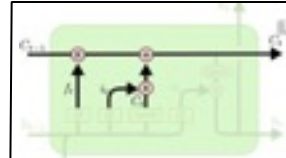
Forget gender of old subject

### Input gate



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

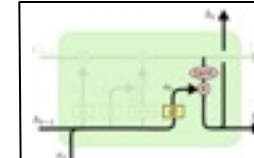
Input gender of new subject



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

Actually drop old  
Add new

### Output gate



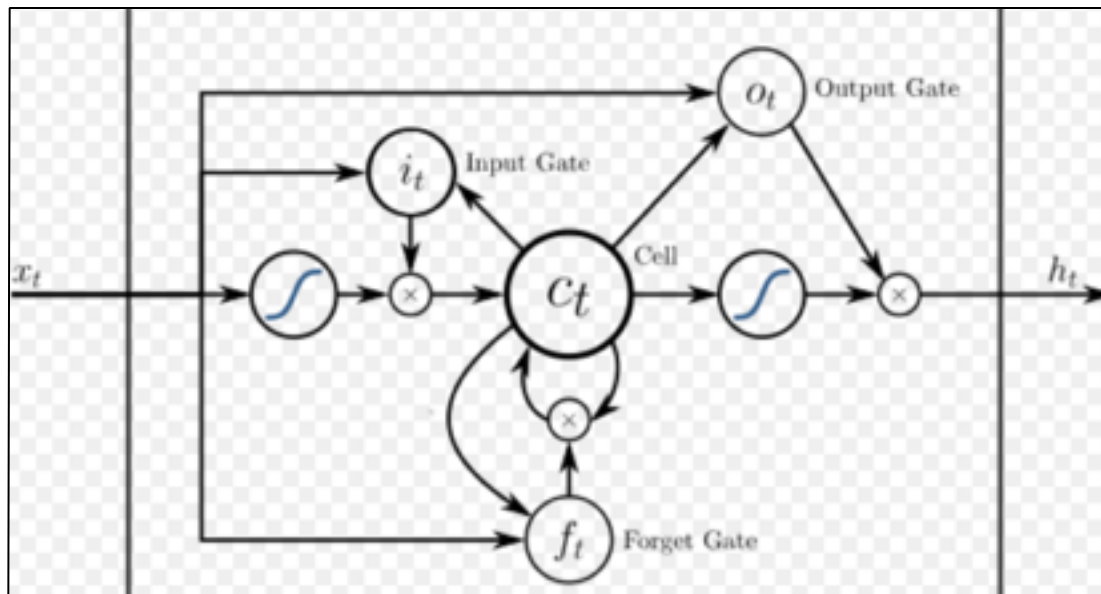
$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

Whether subject is singular  
or plural

## An LSTM variant

- A common LSTM unit is composed of
  - a **cell**, an **input gate**, an **output gate** and a **forget gate**

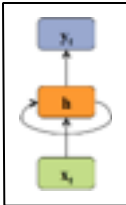
The cell remembers values over arbitrary time intervals and the three *gates* regulate the flow of information into and out of the cell.
- A peephole LSTM is shown below



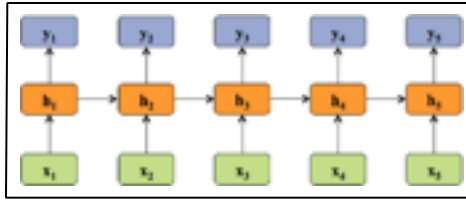
# Summary of Neural Sequential Models

## Recurrent Neural Network

RNN



Unrolled RNN



Definition

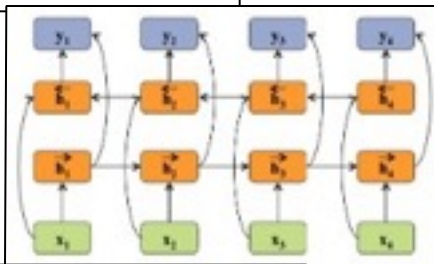
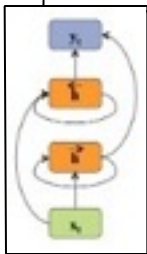
inputs :  $x = (x_1, x_2, \dots, x_T), x_i \in \mathbb{R}^I$   
 hidden units :  $h = (h_1, h_2, \dots, h_T), h_i \in \mathbb{R}^J$   
 outputs :  $y = (y_1, y_2, \dots, y_T), y_i \in \mathbb{R}^K$   
 nonlinearity :  $\mathcal{H}$

Activation Functions

$$h_t = \mathcal{H}(W_{xh}x_t + W_{hh}h_{t-1} + b_h)$$

$$y_t = W_{hy}h_t + b_y$$

## Bidirectional RNN



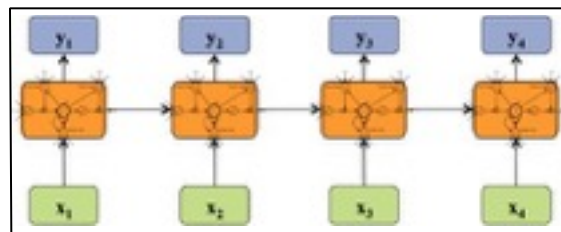
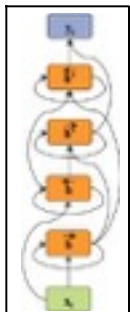
inputs :  $x = (x_1, x_2, \dots, x_T), x_i \in \mathbb{R}^I$   
 hidden units :  $\vec{h}$  and  $\overleftarrow{h}$   
 outputs :  $y = (y_1, y_2, \dots, y_T), y_i \in \mathbb{R}^K$   
 nonlinearity :  $\mathcal{H}$

$$\vec{h}_t = \mathcal{H}(W_{x\vec{h}}x_t + W_{\vec{h}\vec{h}}\vec{h}_{t-1} + b_{\vec{h}})$$

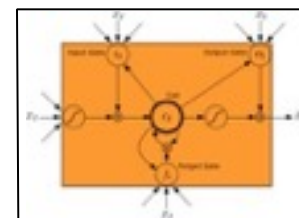
$$\overleftarrow{h}_t = \mathcal{H}(W_{x\overleftarrow{h}}x_t + W_{\overleftarrow{h}\overleftarrow{h}}\overleftarrow{h}_{t-1} + b_{\overleftarrow{h}})$$

$$y_t = W_{\overleftarrow{h}y}\overleftarrow{h}_t + W_{\vec{h}y}\vec{h}_t + b_y$$

## Deep Bidirectional RNN



## LSTM

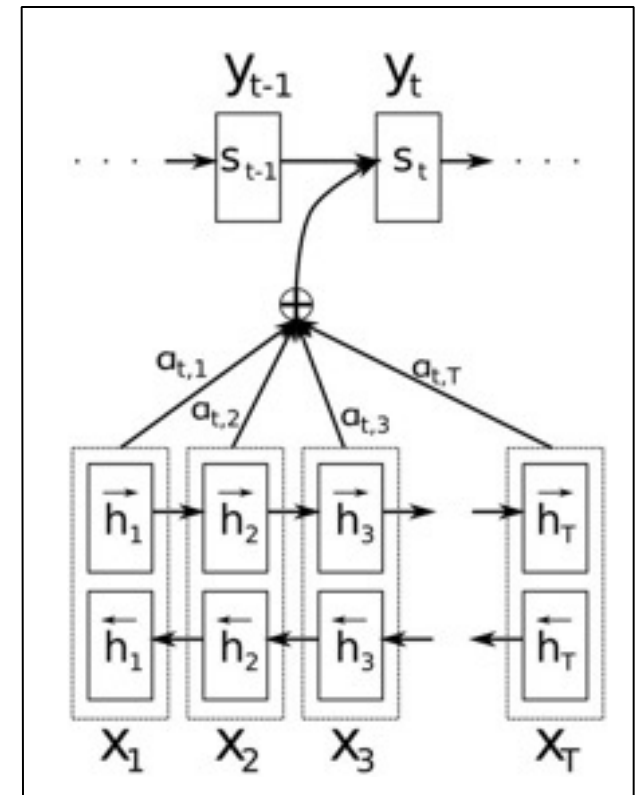


# Attention Mechanisms

- Long range dependencies are still a problem with LSTMs
- With an attention mechanism we no longer try encode the full source sentence into a fixed-length vector
  - Rather, we allow the decoder to “attend” to different parts of the source sentence at each step of the output generation
- Attention is simply a vector, often the outputs of dense layer using softmax function.

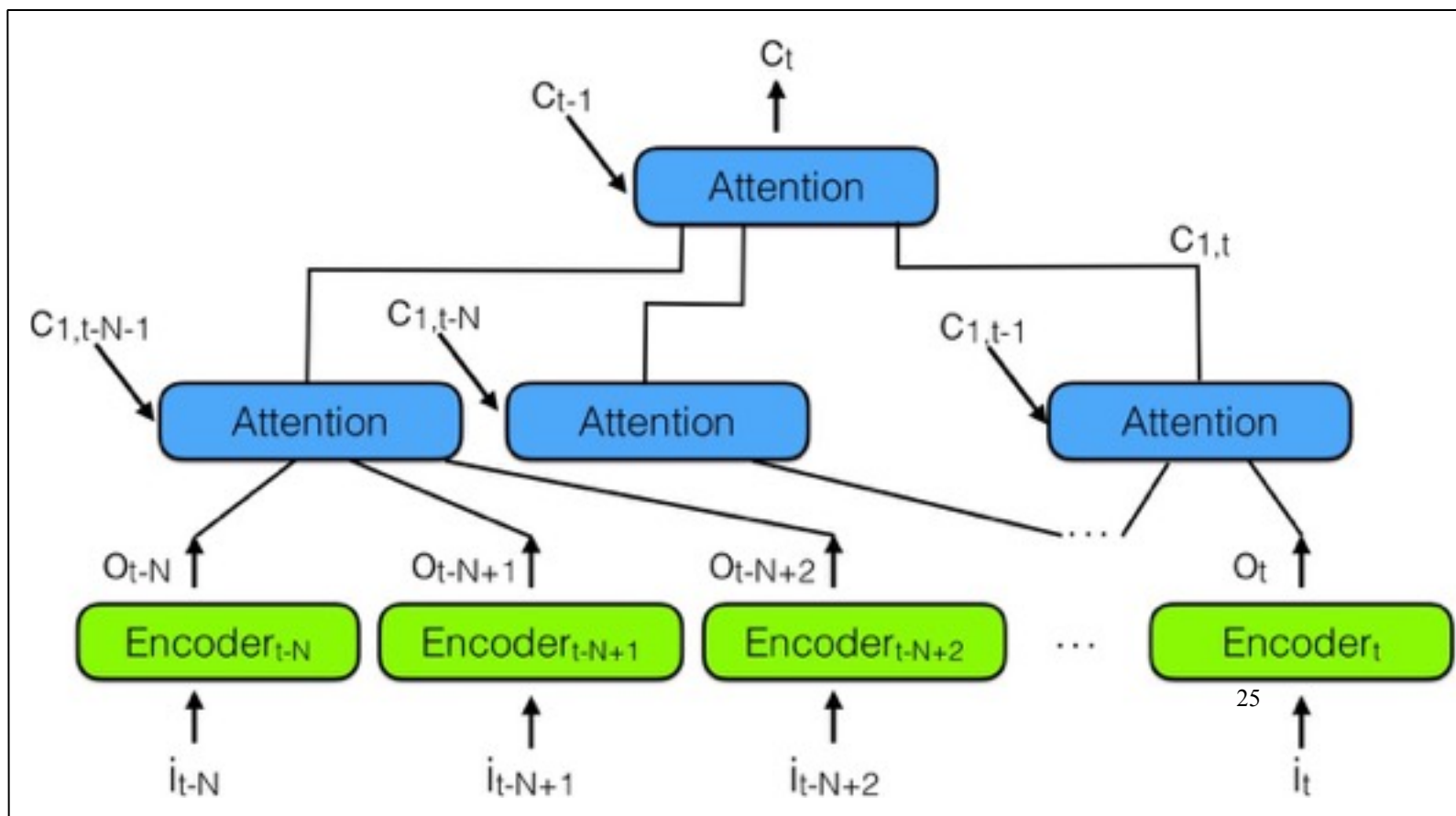
## Attention Mechanism in Translation

- $y$  : translated words
- $x$  : source words.
  - Use of bidirectional RNN is unimportant
- Each decoder output word now depends on
  - weighted combination of all the input states,
    - not just the last state.
  - $a$ 's are weights for each input state
    - if  $a_{3,2}$  is large, decoder pays attention to second state in source sentence while producing third word of target
    - $a$ 's are normalized to sum to 1





## Hierarchical neural attention encoder

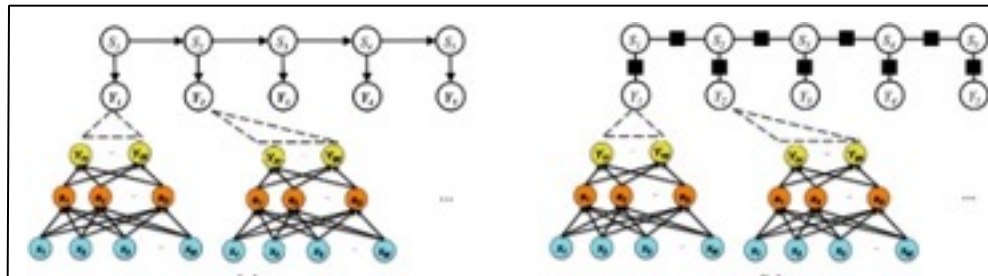


# Deep Learning and Graphical Models

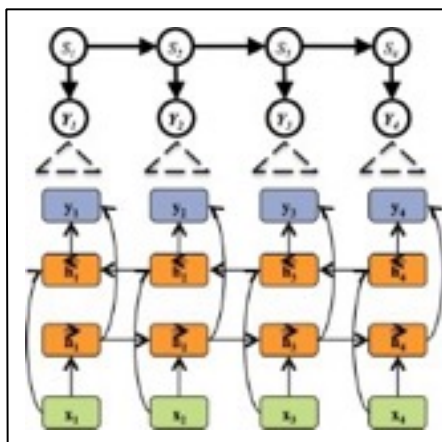
- In deep learning:
  - Tasks of interest:
    - Classification
      - Feature learning
  - Method of learning
    - Back propagation and gradient descent
- In graphical models:
  - Tasks of interest:
    - Transfer learning
    - Latent variable inference
  - Methods of learning
    - Parameter learning methods
    - Structure learning methods
- Hybrid graphical models combine the two types of models
  - They are trained using backpropagation

# Hybrid Graphical Models and Neural Networks

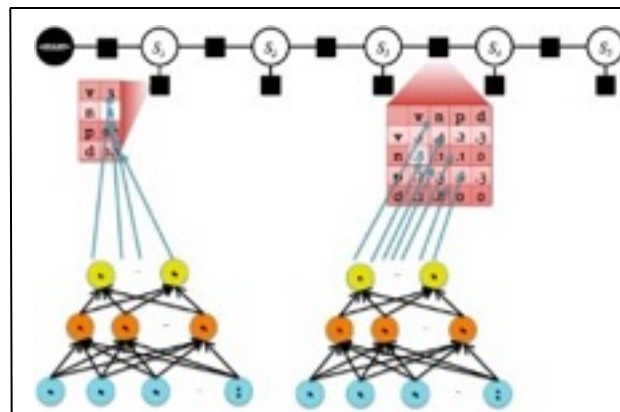
- Hybrid NN and HMM



- Hybrid RNN+HMM



- Hybrid CNN+CRF



# Unfolding Computational Graphs

## Topics in Sequence Modeling

- Recurrent Neural Networks
  1. Unfolding Computational Graphs
  2. Recurrent Neural Networks
  3. Bidirectional RNNs
  4. Encoder-Decoder Sequence-to-Sequence Architectures
  5. Deep Recurrent Networks
  6. Recursive Neural Networks
  7. The Challenge of Long-Term Dependencies
  8. Echo-State Networks
  9. Leaky Units and Other Strategies for Multiple Time Scales
  10. LSTM and Other Gated RNNs
  11. Optimization for Long-Term Dependencies
  12. Explicit Memory

## Unfolding Computational Graphs

- A Computational Graph is a way to formalize the structure of a set of computations
  - Such as mapping inputs and parameters to outputs and loss
- We can unfold a recursive or recurrent computation into a computational graph that has a repetitive structure
  - Corresponding to a chain of events
- Unfolding this graph results in sharing of parameters across a deep network structure

## Example of unfolding a recurrent equation

- Classical form of a dynamical system is

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

- where  $s^{(t)}$  is called the state of the system
  - Equation is recurrent because the definition of  $s$  at time  $t$  refers back to the same definition at time  $t-1$
- For a finite no. of time steps  $\tau$ , the graph can be unfolded by applying the definition  $\tau-1$  times
  - E.g, for  $\tau = 3$  time steps we get

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

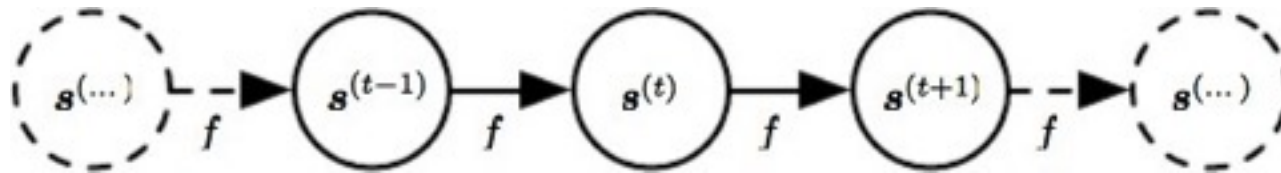
- Unfolding equation by repeatedly applying the definition in this way has yielded expression without recurrence
  - $s^{(1)}$  is ground state and  $s^{(2)}$  computed by applying  $f$
- Such an expression can be represented by a traditional acyclic computational graph (as shown next)

## Unfolded dynamical system

- The classical dynamical system described by

$$s^{(t)} = f(s^{(t-1)}; \theta) \text{ and } s^{(3)} = f(f(s^{(1)}; \theta); \theta)$$

is illustrated as an unfolded computational graph



- Each node represents state at some time  $t$
- Function  $f$  maps state at time  $t$  to the state at  $t+1$
- The same parameters ( the same value of  $\theta$  used to parameterize  $f$  ) are used for all time steps



## Dynamical system driven by external signal

- As another example, consider a dynamical system driven by external (input) signal  $\mathbf{x}^{(t)}$

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

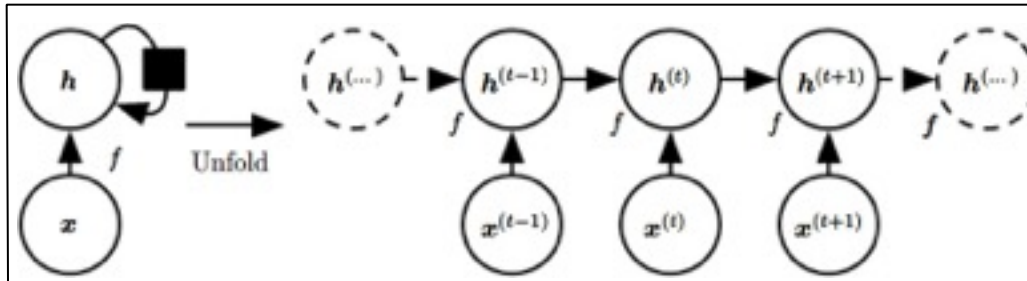
- State now contains information about the whole past input sequence
- Note that the previous dynamic system was simply  $\mathbf{s}^{(t)} = f(\mathbf{s}^{(t-1)}; \boldsymbol{\theta})$
- Recurrent neural networks can be built in many ways
  - Much as almost any function is a feedforward neural network, any function involving recurrence can be considered to be a recurrent neural network

## Defining values of hidden units in RNNs

- Many recurrent neural nets use same equation (as dynamical system with external input) to define values of hidden units
  - To indicate that the state is hidden rewrite using variable  $h$  for state:

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

- Illustrated below:



- Typical RNNs have extra architectural features such as output layers that read information out of state  $h$  to make predictions<sup>34</sup>

## A recurrent network with no outputs

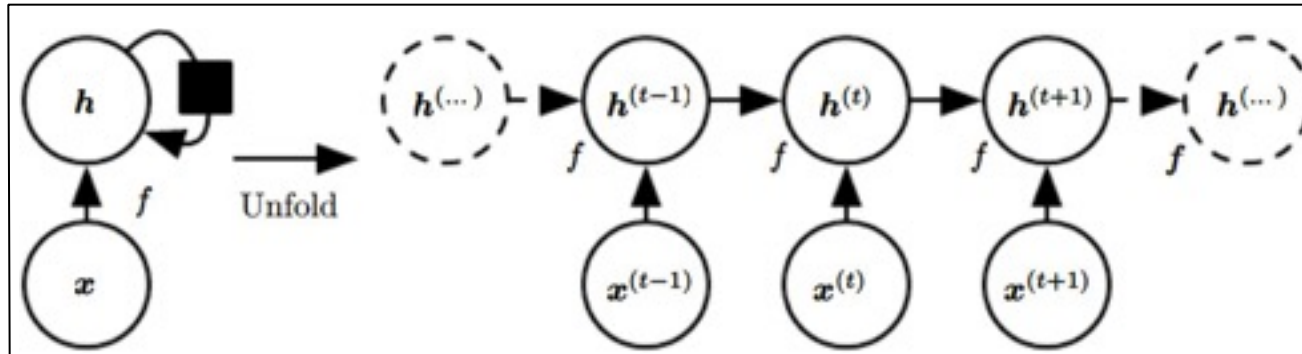
- This network just processes information from input  $x$  by incorporating it into state  $h$  that is passed forward through time

### Circuit diagram:

Black square indicates  
Delay of one time step

### Unfolded computational graph:

each node is now  
associated with one time instance



Typical RNNs will add extra architectural features such as output layers to read information out of the state  $h$  to make predictions

## Predicting the Future from the Past

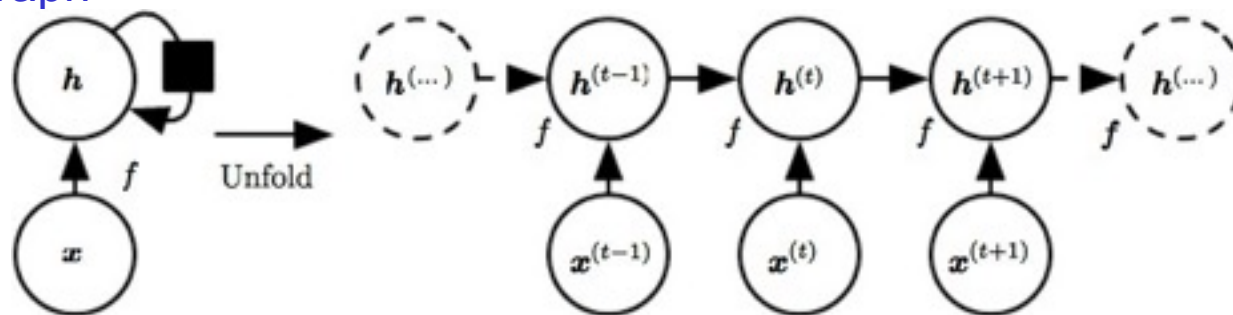
- When RNN is required to perform a task of predicting the future from the past, network typically learns to use  $\mathbf{h}^{(t)}$  as a lossy summary of the task-relevant aspects of the past sequence of inputs upto  $t$
- The summary is in general lossy since it maps a sequence of arbitrary length  $(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  to a fixed length vector  $\mathbf{h}^{(t)}$

## Information Contained in Summary

- Depending on criterion, summary keeps some aspects of past sequence more precisely than other aspects
- Examples:
  - RNN used in statistical language modeling, typically to predict next word from past words
    - it may not be necessary to store all information upto time  $t$  but only enough information to predict rest of sentence
  - Most demanding situation: we ask  $\mathbf{h}^{(t)}$  to be rich enough to allow one to approximately recover the input sequence as in autoencoders

## Unfolding: from circuit diagram to computational graph

- Equation  $h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta)$  can be written in two different ways:  
circuit diagram or an unfolded computational graph



- Unfolding is the operation that maps a circuit to a computational graph with repeated pieces
- The unfolded graph has a size dependent on the sequence length

## Process of Unfolding

- We can represent unfolded recurrence after  $t$  steps with a function  $g^{(t)}$ :

$$\mathbf{h}^{(t)} = g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$$

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

- Function  $g^{(t)}$  takes in whole past sequence  $(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  as input and produces the current state but the unfolded recurrent structure allows us to factorize  $g^{(t)}$  into repeated application of a function  $f$
- The unfolding process introduces two major advantages as discussed next.

## Unfolding process allows learning a single model

- The unfolding process introduces two major advantages:
  1. Regardless of sequence length, learned model has same input size
    - because it is specified in terms of transition from one state to another state rather than specified in terms of a variable length history of states
  2. Possible to use same function  $f$  with same parameters at every step
- These two factors make it possible to learn a single model  $f$ 
  - that operates on all time steps and all sequence lengths
  - rather than needing separate model  $g^{(t)}$  for all possible time steps
- Learning a single shared model allows:
  - Generalization to sequence lengths that did not appear in the training
  - Allows model to be estimated with far fewer training examples than would be required without parameter sharing

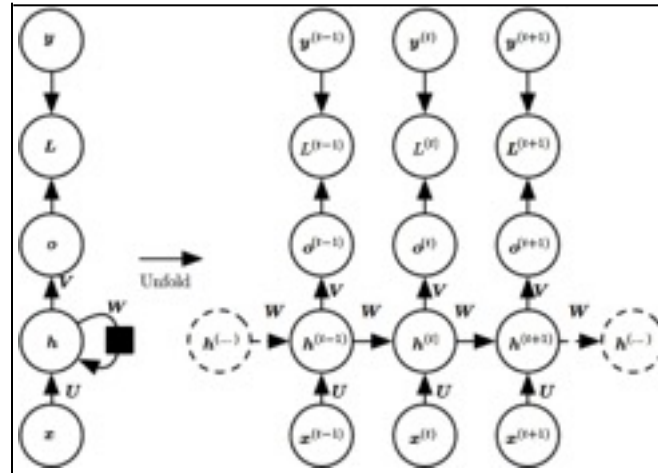


## Both recurrent graph and unrolled graph are useful

- Recurrent graph is succinct
- Unrolled graph provides explicit description of what computations to perform
  - Helps illustrate the idea of information flow forward in time
    - Computing outputs and losses
  - And backwards in time
    - Computing gradients
  - By explicitly showing path along which information flows

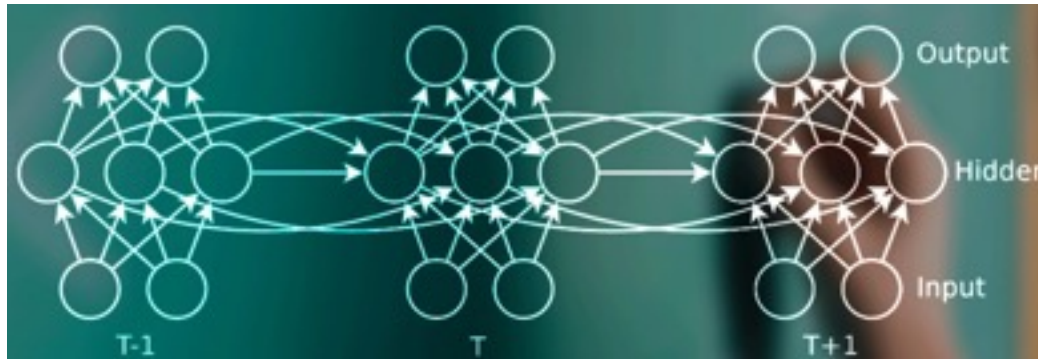
# Unfolding Computational Graphs

- A Computational Graph is a way to formalize the structure of a set of computations
  - Such as mapping inputs and parameters to outputs and loss
- We can unfold a recursive or recurrent computation into a computational graph that has a repetitive structure
  - Corresponding to a chain of events
- Unfolding this graph results in sharing of parameters across a deep network structure



$$\begin{aligned} o^{(t)} &= c + Vh^{(t)} \\ h^{(t)} &= \tanh(a^{(t)}) \\ a^{(t)} &= b + Wh^{(t-1)} + Ux^{(t)} \end{aligned}$$

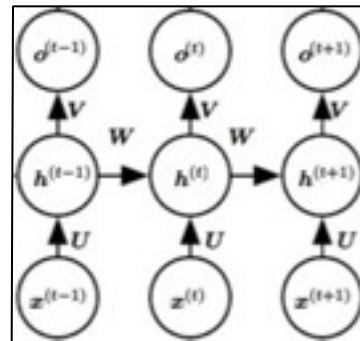
## A more complex unfolded computational graph



Source: Indico corporation

- Two units in input layer, instead of 1
- Three units in hidden layer, instead of 1
- Two units in output layer, instead of 1

Previous one:



# Recurrent Neural Networks

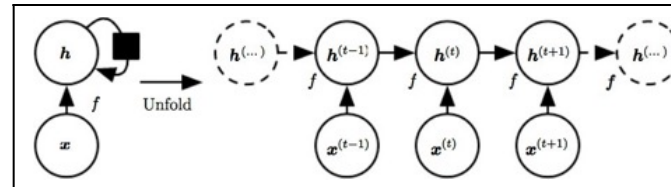
## Topics in Sequence Modeling

- Overview
- 1. Unfolding Computational Graphs
- 2. Recurrent Neural Networks
- 3. Bidirectional RNNs
- 4. Encoder-Decoder Sequence-to-Sequence Architectures
- 5. Deep Recurrent Networks
- 6. Recursive Neural Networks
- 7. The Challenge of Long-Term Dependencies
- 8. Echo-State Networks
- 9. Leaky Units and Other Strategies for Multiple Time Scales
- 10. LSTM and Other Gated RNNs
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory

## Graph unrolling/parameter sharing in RNN design

- Process of graph unrolling:

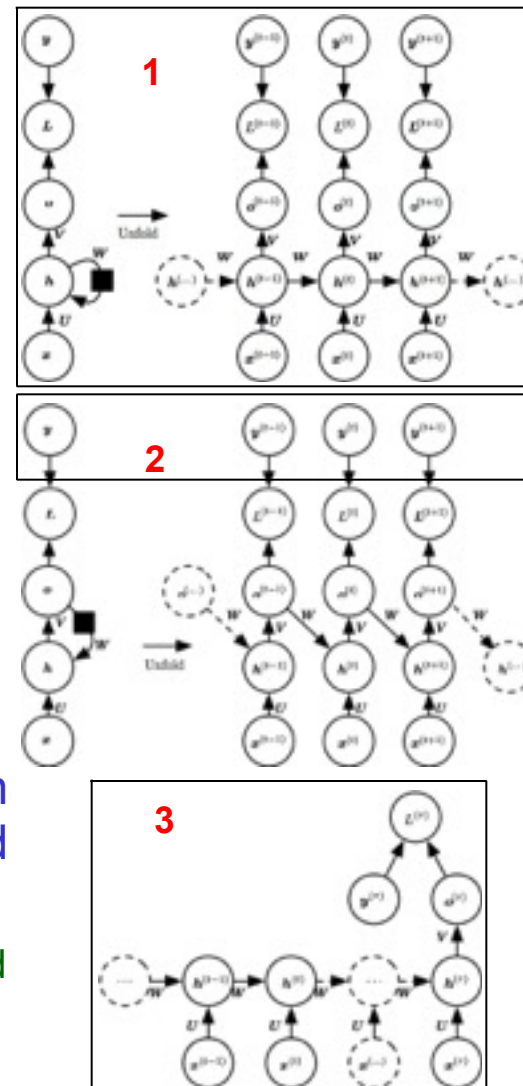
$$\begin{aligned} \mathbf{h}^{(t)} &= g^{(t)}(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)}) \\ &= f(\mathbf{h}^{(t-1)}, \mathbf{x}^{(t)}; \boldsymbol{\theta}) \end{aligned}$$



- Function  $g^{(t)}$  takes in whole past sequence  $(\mathbf{x}^{(t)}, \mathbf{x}^{(t-1)}, \dots, \mathbf{x}^{(2)}, \mathbf{x}^{(1)})$  as input and produces the current state but the unfolded recurrent structure allows us to factorize  $g^{(t)}$  into repeated application of a function  $f$
- It does not need a separate model  $g^{(t)}$  for all possible time steps
- Process of unrolling and parameter sharing allows us to design a wide variety of recurrent neural networks
- Examples of important design patterns shown next

## Three design patterns of RNNs

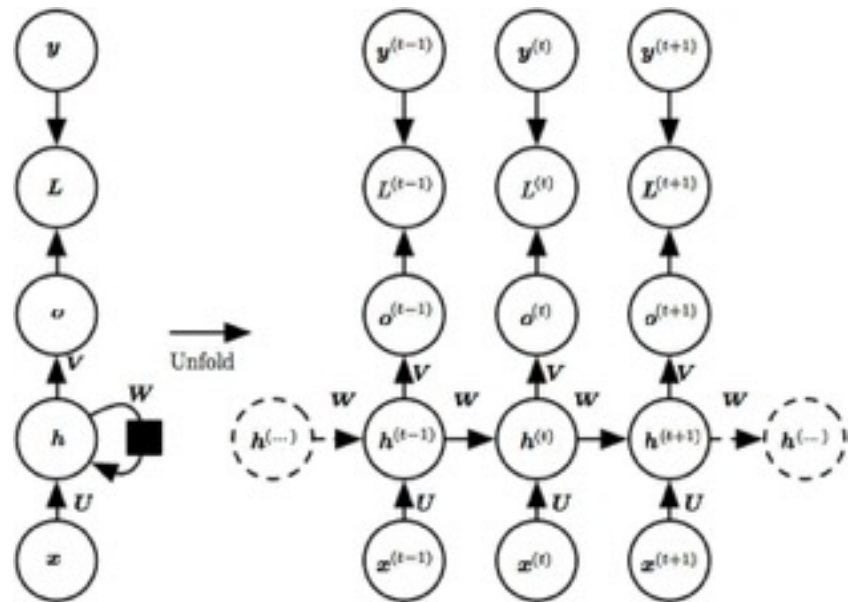
1. Produce output at each time step and have recurrent connections between hidden units
2. Produce output at each time step and have recurrent connections only from output at one time step to hidden units at next time step
3. Recurrent connections between hidden units to read entire input sequence and produce a single output
  - Can summarize sequence to produce a fixed size representation for further processing



## Design 1: RNN with recurrence between hidden units

- Maps input sequence  $\mathbf{x}$  to output  $\mathbf{o}$  values
  - Loss  $L$  measures how far each  $\mathbf{o}$  is from the corresponding target  $\mathbf{y}$ 
    - With softmax outputs we assume  $\mathbf{o}$  is the unnormalized log probabilities
    - Loss  $L$  internally computes  $\mathbf{y} = \text{softmax}(\mathbf{o})$  and compares to target  $\mathbf{y}$
- Update equation

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





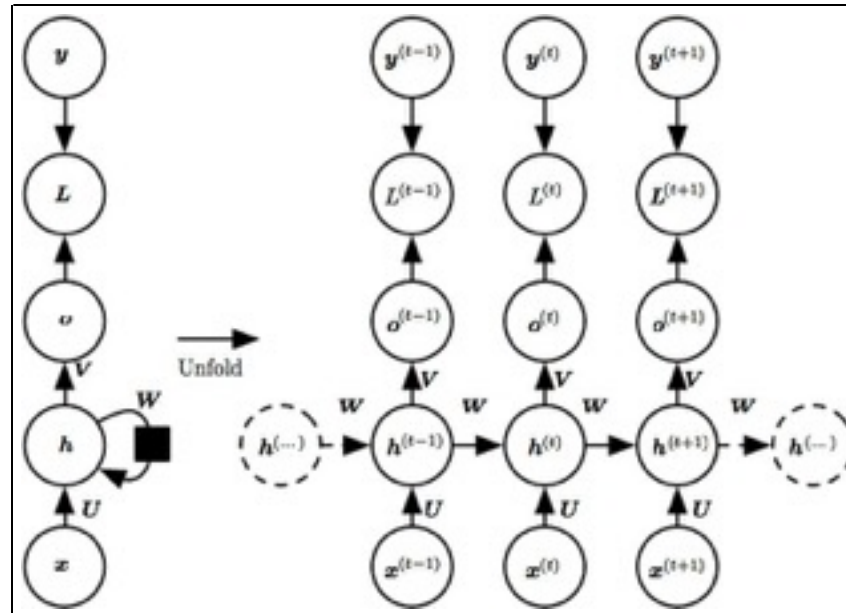
## RNN with hidden unit recurrence is a Universal TM

- RNN with hidden unit connections together with

$$a^{(t)} = b + Wh^{(t-1)} + Ux^{(t)}$$

is Universal, i.e.,

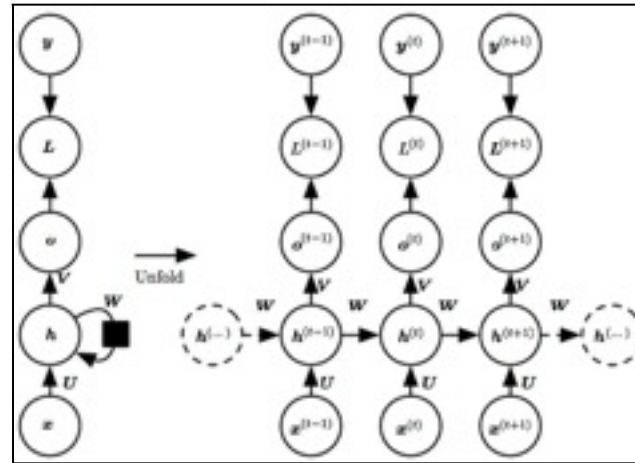
- Any function computable by a TM is computable by such an RNN of finite size



- Output can be read from the RNN after no. of steps asymptotically linear in no. of steps used by the TM and asymptotically linear in length of input

## Forward prop for RNN with hidden unit recurrence

- This design does not specify
  - activation functions for hidden units,
  - form of output and
  - loss function
- Assume for this graph:
  - Hyperbolic tangent activation function
  - Output is discrete



- E.g., RNN predicts words/characters
- Natural way to represent discrete vars:
  - Output  $\mathbf{o}$  gives unnormalized log probabilities of each possible value  $\hat{\mathbf{y}}$
  - Can apply softmax as a postprocessing step to obtain vector of normalized probabilities over the output
- Forward prop begins with specification of initial state  $\mathbf{h}^{(0)}$

## Forward Prop Equations for RNN with hidden recurrence

- Begins with initial specification of  $\mathbf{h}^{(0)}$
- Then for each time step from  $t=1$  to  $t=\tau$  we apply the following update equations

$$\begin{aligned}\mathbf{o}^{(t)} &= \mathbf{c} + V\mathbf{h}^{(t)} \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}) \\ \mathbf{a}^{(t)} &= \mathbf{b} + W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)}\end{aligned}$$

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

- Where the parameters are:
  - bias vectors  $\mathbf{b}$  and  $\mathbf{c}$
  - weight matrices  $U$  (input-to-hidden),  $V$  (hidden-to-output) and  $W$  (hidden-to-hidden) connections

## Loss function for a given sequence

- This is an example of an RNN that maps an input sequence to an output sequence of the same length
- Total loss for a given sequence for a given sequence of  $\mathbf{x}$  values with a sequence of  $\mathbf{y}$  values is the sum of the losses over the time steps
- If  $L^{(t)}$  is the negative log-likelihood of  $\mathbf{y}^{(t)}$  given  $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}$  then

$$L(\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\} | \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(t)}\}) = \sum_t L^{(t)}$$

$$= - \sum_t \log p_{\text{model}}(\mathbf{y}^{(t)} | \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)}\})$$

- where  $p_{\text{model}}$  is given by reading the entry for  $\mathbf{y}^{(t)}$  from the model's output vector  $\mathbf{y}^{\theta}$

## Computing Gradient of Loss Function

- Computing gradient of this loss function wrt parameters is expensive
  - It involves performing a forward propagation pass moving left to right through the unrolled graph followed by a backward propagation moving right to left through the graph
- Run time is  $O(\tau)$  and cannot be reduced by parallelization
  - Because forward propagation graph is inherently sequential
    - Each time step computable only after previous step
  - States computed during forward pass must be stored until reused in the backward pass
    - So memory cost is also  $O(\tau)$
- Backpropagation applied to the unrolled graph with  $O(\tau)$  cost is called *Backward Propagation through time* (BPTT)

## Backward Propagation through Time (BPTT)

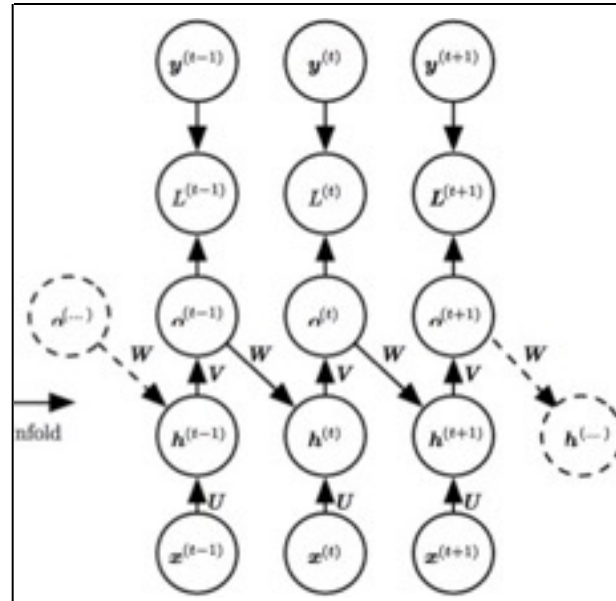
- RNN with hidden unit recurrence is very powerful but also expensive to train
- Is there an alternative?

## Topics in Recurrent Neural Networks

- Overview
  - Design patterns for RNNs
  - RNN with recurrence between hidden units
  - Forward propagation equations
  - Loss function for a sequence
  - RNN with recurrence from output units to hidden units
- 1. Teacher forcing and Networks with Output Recurrences
- 2. Computing the gradient in a RNN
- 3. Recurrent Networks as Graphical Models
- 4. Modeling Sequences Conditioned on Context with RNNs

## Models with recurrence from output to hidden

- Less powerful than with hidden-to-hidden recurrent connections
  - It cannot simulate a universal TM
  - It requires that the output capture all information of past to predict future
- Advantage
  - In comparing loss function to output all time steps are decoupled
    - Each step can be trained in isolation
  - Training can be parallelized
    - Gradient for each step  $t$  computed in isolation
    - No need to compute output for the previous step first, because training set provides ideal value of output
- Can be trained with Teacher Forcing
  - Described next





## Teacher Forcing

- Used for training models with recurrent connections from their outputs
- Teacher forcing during training means
  - Instead of summing activations from incoming units (possibly erroneous)
  - Each unit sums correct teacher activations as input for the next iteration

## Teacher Forcing Procedure

- Teacher forcing is a procedure for training RNNs with output-to-hidden recurrence
- It emerges from the maximum likelihood criterion
- During training time model receives ground truth output  $\mathbf{y}^{(t)}$  as input at time  $t+1$ 
  - We can see this by examining a sequence with two time steps (next)

# Maximum Likelihood in Teacher Forcing

- Consider a sequence with two time steps
- The conditional maximum likelihood criterion is

$$\log (p(\mathbf{y}^{(1)}, \mathbf{y}^{(2)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)})) = \log p(\mathbf{y}^{(2)} | \mathbf{y}^{(1)}, \mathbf{x}^{(1)}, \mathbf{x}^{(2)}) + \log p(\mathbf{y}^{(1)} | \mathbf{x}^{(1)}, \mathbf{x}^{(2)})$$

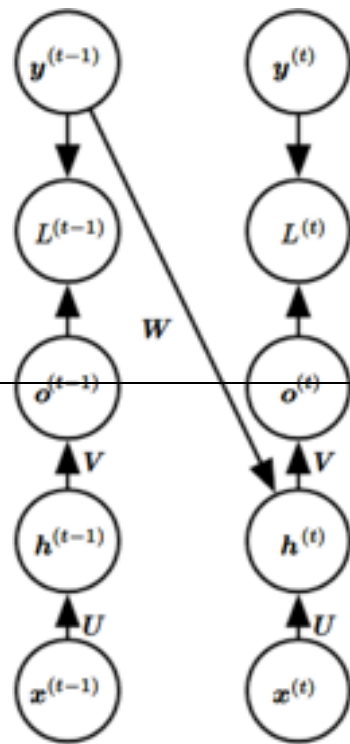
- At time  $t=2$ , model is trained to maximize conditional probability of  $\mathbf{y}^{(2)}$  given both the  $\mathbf{x}$  sequence so far and the previous  $\mathbf{y}$  value from the training set.
  - We are using  $\mathbf{y}^{(1)}$  as teacher forcing, rather than only  $\mathbf{x}^{(i)}$ .
- Maximum likelihood specifies that during training, rather than feeding the model's own output back to itself, target values should specify what the correct output should be

# Illustration of Teacher Forcing

- Teacher Forcing is a training technique applicable to RNNs that have connections from output to hidden states at next time step

## Train time:

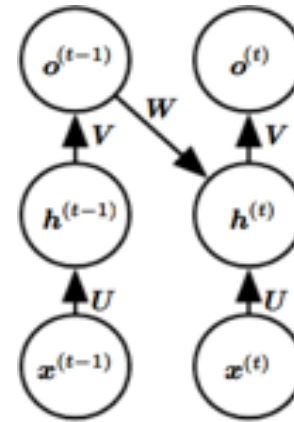
We feed the correct output  $y^{(t)}$  (from teacher) drawn from the training set as input to  $h^{(t+1)}$



Train time

## Test time:

True output is not known. We approximate the correct output  $y^{(t)}$  with the model's output  $o^{(t)}$  and feed the output back to the model



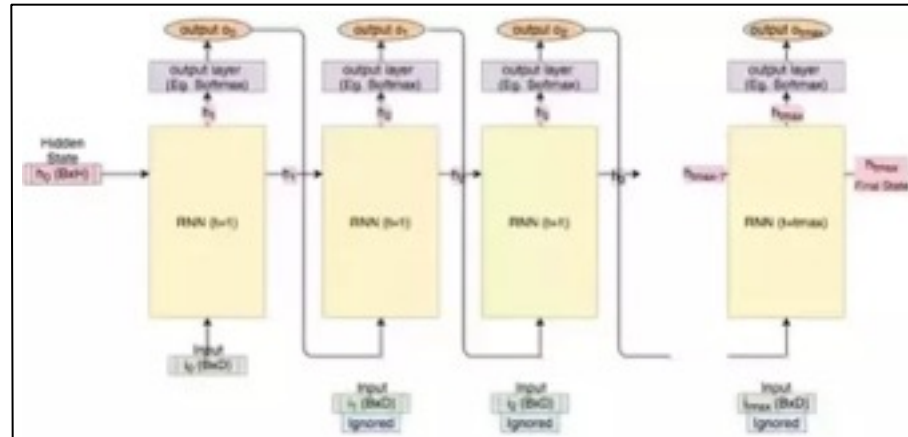
Test time

# RNNs without and with teacher forcing

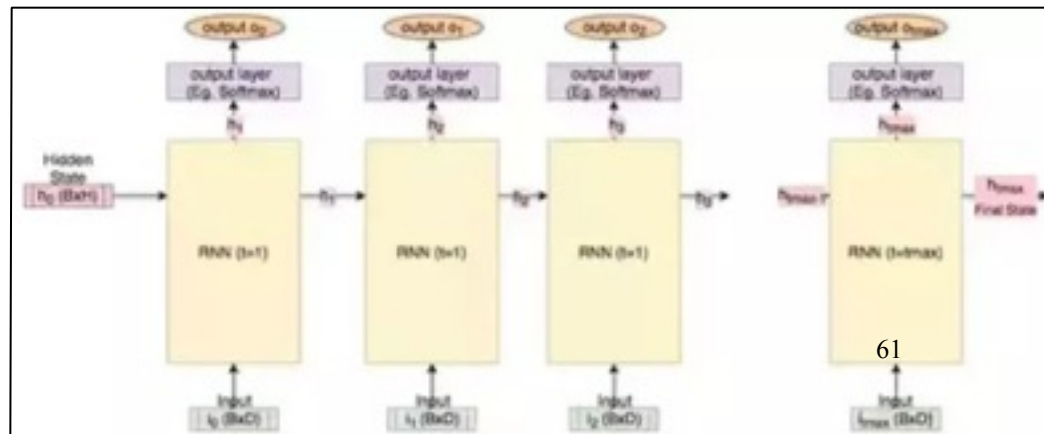
## Without Teacher forcing:

Output of previous time step acts as input during RNN training

RNN learns to generate next output in a sequence based upon the previous value



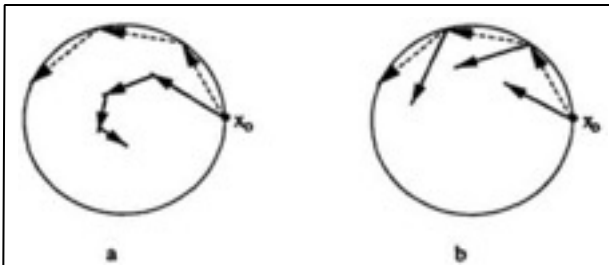
## With Teacher forcing:



Source: <https://www.quora.com/What-is-the-teacher-forcing-in-RNN>

## Visualizing Teacher Forcing

- Imagine that the network is learning to follow a trajectory
- It goes astray (because the weights are wrong) but teacher forcing puts the net back on its trajectory
  - By setting the state of all the units to that of teacher's.



(a) Without teacher forcing, trajectory runs astray (solid lines) while the correct trajectory are the dotted lines

(b) With teacher forcing trajectory corrected at each step

## Teacher forcing for hidden-to-hidden

- We originally motivated teacher forcing as allowing us to avoid backpropagation through time in models that lack hidden-to-hidden connections
- Teacher forcing can still be applied to models that have hidden-to-hidden connections
  - As long as they have connections from the output at one time step to values computed at the next time step
  - As soon as the hidden units become functions of earlier time steps, BPTT algorithm is necessary

## Training with both Teacher Forcing and BPTT

- Some models may be trained with both Teacher forcing and Backward Propagation through time (BPTT)
  - When there are both hidden-to-hidden recurrences as well as output-to-hidden recurrences



## Disadvantage of Teacher Forcing

- If network is to be used in an *open-loop* mode with network outputs (or samples from the output distribution) fed back as inputs
  - In this case the kind of inputs that it will see during training time could be quite different from that it will see at test time
- Solutions to mitigate this problem:
  1. Train with both teacher-forced inputs and free running inputs
    - E.g., predicting the correct target a no of steps in the future through the unfolded recurrent output-to-input paths
    - Thus network can learn to take into account input conditions not seen during training
  2. Mitigate the gap between inputs seen at training time and test time by *generating* values as input
    - This approach exploits a *curriculum learning* strategy to gradually use more of the generated values as input

# Professor Forcing

- Teacher Forcing trains RNNs by supplying observed sequence values as inputs during training and using the network's own one-step-ahead predictions to do multi-step sampling
- Professor Forcing algorithm
  - Uses adversarial domain adaptation to encourage the dynamics of the RNN to be the same when training the network and when sampling from the network over multiple time steps.
  - Applied to
    - language modeling
    - vocal synthesis on raw waveforms
    - handwriting generation
    - image generation
- Professor Forcing is an adversarial method for learning generative models
  - It is closely related to Generative Adversarial Networks

# Computing the Gradient in an RNN

## Topics in Recurrent Neural Networks

1. Overview
2. Teacher forcing and Networks with Output Recurrence
3. Computing the gradient in a RNN
4. RNNs as Directed Graphical Models
5. Modeling Sequences Conditioned on Context with RNNs

# General Backpropagation to compute gradient

To compute gradients  $\mathbb{T}$  of variable  $z$  wrt variables in computational graph  $\mathcal{G}$

---

**Algorithm 6.5** The outermost skeleton of the back-propagation algorithm. This portion does simple setup and cleanup work. Most of the important work happens in the `build_grad` subroutine of algorithm 6.6

---

**Require:**  $\mathbb{T}$ , the target set of variables whose gradients must be computed.

**Require:**  $\mathcal{G}$ , the computational graph

**Require:**  $z$ , the variable to be differentiated

Let  $\mathcal{G}'$  be  $\mathcal{G}$  pruned to contain only nodes that are ancestors of  $z$  and descendants of nodes in  $\mathbb{T}$ .

Initialize `grad_table`, a data structure associating tensors to their gradients

`grad_table[z] ← 1`

**for**  $V$  in  $\mathbb{T}$  **do**

`build_grad(V,  $\mathcal{G}$ ,  $\mathcal{G}'$ , grad_table)`

**end for**

Return `grad_table` restricted to  $\mathbb{T}$

---

# Build-grad function of Generalized Backprop

**Algorithm 6.6** The inner loop subroutine `build_grad(V, G, G', grad_table)` of the back-propagation algorithm, called by the back-propagation algorithm defined in algorithm 6.5.

**Require:** `V`, the variable whose gradient should be added to `G` and `grad_table`.

**Require:** `G`, the graph to modify.

**Require:** `G'`, the restriction of `G` to nodes that participate in the gradient.

**Require:** `grad_table`, a data structure mapping nodes to their gradients

if `V` is in `grad_table` then

Return `grad_table[V]`

end if

$i \leftarrow 1$

for `C` in `get_consumers(V, G')` do

op  $\leftarrow$  `get_operation(C)`

`D`  $\leftarrow$  `build_grad(C, G, G', grad_table)`

$G^{(i)} \leftarrow$  op.bprop(`get_inputs(C, G')`, `V`, `D`)

$i \leftarrow i + 1$

end for

$G \leftarrow \sum_i G^{(i)}$

`grad_table[V] = G`

Insert `G` and the operations creating it into `G`

Return `G`

Ob.bprop(inputs, X, G)

returns

$$\sum (\nabla_x \text{op.f}(\text{inputs})_i) G_i$$

Which is just an implementation of chain rule

If  $y = g(x)$  and  $z = f(y)$  from chain rule

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

Equivalently, in vector notation

where  $\frac{\partial y}{\partial x}$  is the  $n \times n$  Jacobian matrix of  $g$

$$\nabla_x z = \left( \frac{\partial y}{\partial x} \right)^T \nabla_y z$$

## Example for how BPTT algorithm behaves

- BPTT: Back Propagation Through Time
- We provide an example of how to compute gradients by BPTT for the RNN equations

$$\begin{aligned} \mathbf{a}^{(t)} &= \mathbf{b} + \mathbf{W}\mathbf{h}^{(t-1)} + \mathbf{U}\mathbf{x}^{(t)} \\ \mathbf{h}^{(t)} &= \tanh(\mathbf{a}^{(t)}) \\ \mathbf{o}^{(t)} &= \mathbf{c} + \mathbf{V}\mathbf{h}^{(t)} \\ \hat{\mathbf{y}}^{(t)} &= \text{softmax}(\mathbf{o}^{(t)}) \end{aligned}$$

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

## Intuition of how BPTT algorithm behaves

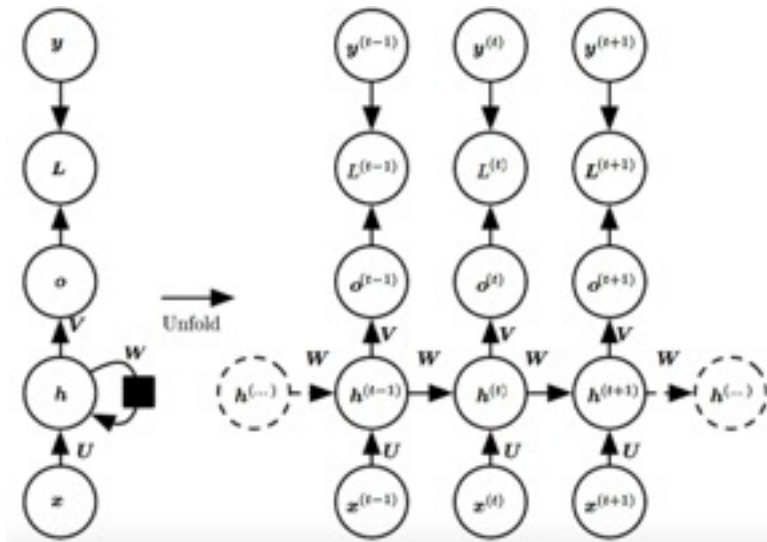
- Ex: Gradients by BPTT for RNN eqns given above

Nodes of our computational graph include the parameters

$U, V, W, b, c$  as well as the sequence of nodes

indexed by  $t$  for  $\mathbf{x}^{(t)}, \mathbf{h}^{(t)}, \mathbf{o}^{(t)}$  and  $L^{(t)}$

For each node  $\mathbf{N}$  we need to compute the gradient  $\nabla_{\mathbf{N}} L$  recursively





# Computing BPTT gradients recursively

- Start recursion with nodes immediately preceding final loss

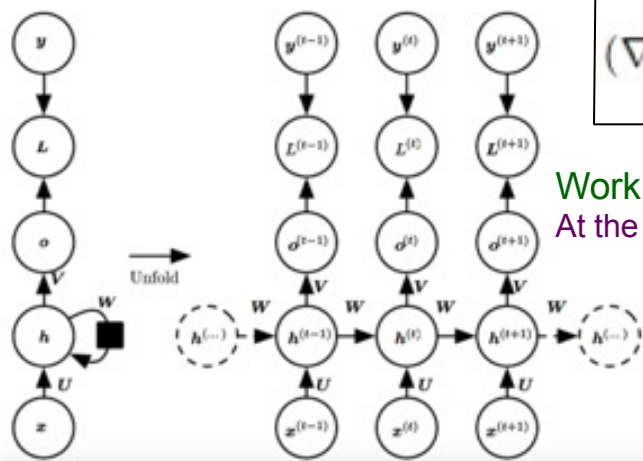
$$\frac{\partial L}{\partial L^{(t)}} = 1$$

- We assume that the outputs  $\mathbf{o}^{(t)}$  are used as arguments to softmax to obtain vector  $\hat{\mathbf{y}}$  of probabilities over the output. Also that loss is the negative log-likelihood of the true target  $\mathbf{y}^{(t)}$  given the input so far.
- The gradient on the outputs at time step  $t$  is

$$(\nabla_{\mathbf{o}^{(t)}} L)_i = \frac{\partial L}{\partial o_i^{(t)}} = \frac{\partial L}{\partial L^{(t)}} \frac{\partial L^{(t)}}{\partial o_i^{(t)}} = \hat{y}_i^{(t)} - \mathbf{1}_{i, y^{(t)}}$$

Work backwards starting from the end of the sequence.  
At the final time step  $\tau$ ,  $\mathbf{h}^{(\tau)}$  has only  $\mathbf{o}^{(\tau)}$  as a dependent. So

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



## Gradients over hidden units

- Iterate backwards in time iterating to backpropagate gradients through time from  $t = \tau - 1$  down to  $t = 1$

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

- where  $\text{diag} \left( 1 - \left( \mathbf{h}^{(t+1)} \right)^2 \right)$  is a matrix with elements within parentheses. This is the jacobian of the hyperbolic tangent associated with the hidden unit  $i$  at time  $t+1$

## Gradients over parameters

- Once gradients on internal nodes of the computational graph are obtained, we can obtain gradients on the parameters
  - Because parameters are shared across time steps some care is needed

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

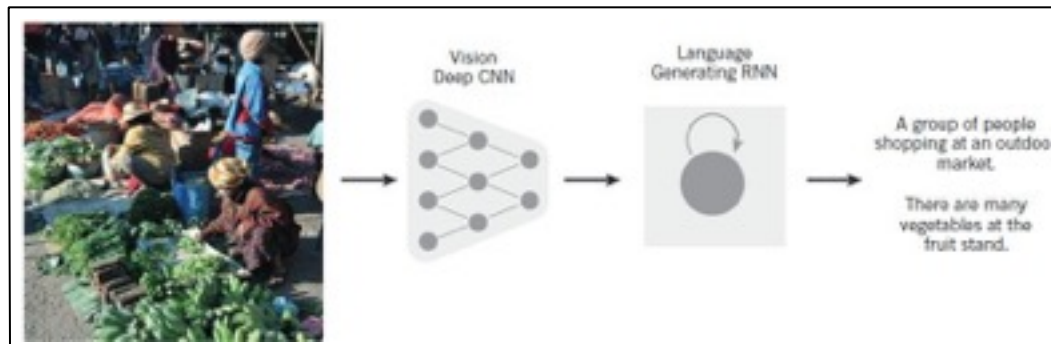
# RNNs as Directed Graphical Models

## Topics in RNNs as Directed Graphical Models

- Motivation
- What are directed graphical models?
- RNNs as Directed Graphical Models
  - Conditioning Predictions on Past Values
  - Fully connected graphical model for a sequence
  - Hidden units as random variables
  - Efficiency of RNN parameterization
- How to draw samples from the model

## Motivation: From images to text

- Combining ConvNets and Recurrent Net Modules
- Caption generated by a recurrent neural network (RNN) taking as input:
  1. Representation generated by a deep CNN
  2. RNN trained to translate high-level representations of images into captions



## Better translation of images to captions

- Different focus (lighter patches given more attention)



A woman is throwing a **frisbee** in a park.



A **dog** is standing on a hardwood floor.



A **stop** sign is on a road with a mountain in the background



A little girl sitting on a bed with a **teddy bear**.



A group of **people** sitting on a boat in the water.



A giraffe standing in a forest with **trees** in the background.

- As it generates each word (**bold**), it exploits it to achieve better translation of images to captions

## What is a Directed graphical model?

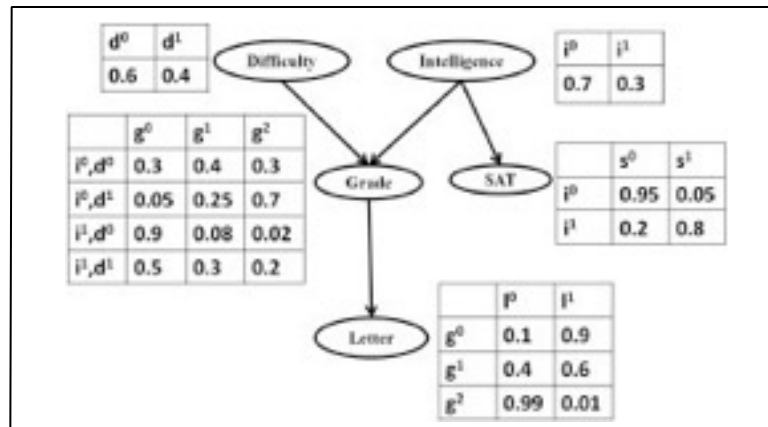
- Specifies a joint distribution  $p(\mathbf{x}) = p(x_1, \dots, x_n)$  in terms of conditional distributions (called a Bayesian Network)

$$p(\mathbf{x}) = \prod_{i=1}^N p(x_i | pa(x_i))$$

- where  $p(x_i | pa(x_i))$  are conditional distributions
  - $pa$  represents parents in a directed graph
  - Computational efficiency is obtained by omitting edges that do not correspond to strong interactions
  - Model is generative:
    - We can easily sample from a BN (ancestral sampling)
- We will see here that RNNs have an interpretation as an efficient directed graphical model



## Example of a Directed Graphical Model (Bayesian Net)

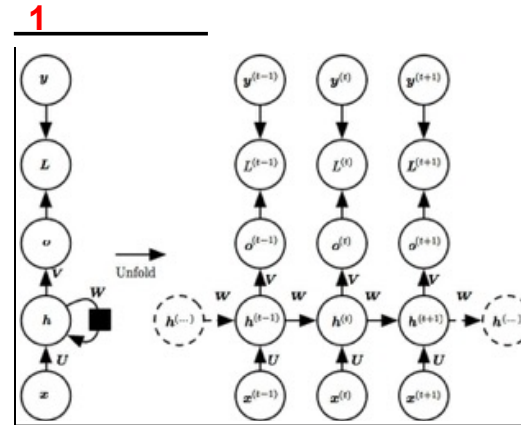


Instead of  $2 \times 2 \times 3 \times 2 \times 2 = 48$  parameters we need 18 parameters using the directed graphical model (which makes certain condition independence assumptions)

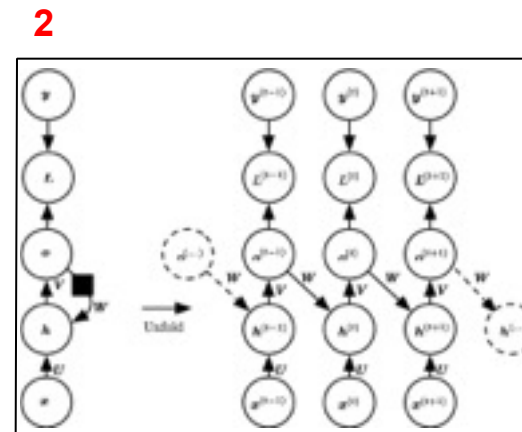
$$\begin{aligned}
 P(D, I, G, S, L) &= P(D)P(I)P(G|D, I)P(S|I)P(L|G) \\
 P(i^1, d^0, g^2, s^1, l^0) &= P(i^1)P(d^0)P(g^2|i^1, d^0)P(s^1|i^1)P(l^0|g^2) \\
 &= 0.3 \cdot 0.6 \cdot 0.08 \cdot 0.8 \cdot 0.4 = 0.004608
 \end{aligned}$$

# CPDs of model depend on RNN Design pattern

1. Recurrent connections between hidden units



2. Recurrent connections only from output at one time step to hidden units at next time step



# Loss function for RNNs

- In feed-forward networks our goal is to minimize

$$-E_{x \sim p^{\text{data}}} [\log p_{\text{model}}(x)]$$

- Which is the cross-entropy between distribution of training set (data) and probability distribution defined by model
  - *Definition of cross entropy* between distributions  $p$  and  $q$  is  
 $H(p, q) = E_p[-\log q] = H(p) + D_{\text{KL}}(p || q)$
  - For discrete distributions  $H(p, q) = -\sum_x p(x) \log q(x)$
- As with a feedforward network, we wish to interpret the output of the RNN as a probability distribution
  - With RNNs the losses  $L^{(t)}$  are cross entropies between training targets  $\mathbf{y}^{(t)}$  and outputs  $\mathbf{o}^{(t)}$ 
    - Mean squared error is the cross-entropy loss associated with an output distribution that is a unit Gaussian

## Types of CPDs in an RNN

1. With a predictive log-likelihood training objective:

$$L(\{x^{(1)}, \dots, x^{(t)}\}, \{y^{(1)}, \dots, y^{(t)}\}) = \sum_t L_t = - \sum_t \log p_{\text{model}}(y^{(t)} | \{x^{(1)}, \dots, x^{(t)}\})$$

we train the RNN to estimate the conditional distribution of the next sequence element  $y^{(t)}$  given past inputs. This means that we maximize log-likelihood

$$\log p(y^{(t)} | x^{(1)}, \dots, x^{(t)})$$

2. Alternatively if the model includes connections from output at one time step to the next

$$\log p(y^{(t)} | x^{(1)}, \dots, x^{(t)}, y^{(1)}, \dots, y^{(t-1)})$$

## Conditioning predictions on past $y$ values

- Decomposing joint probability over sequence of  $y$  values as a series of one-step probabilistic predictions  
is one way to capture the full joint distribution across the whole sequence.
- 1. Do not feed past  $y$  values as inputs for next step prediction, directed model contains no edges from any  $y^{(i)}$  in the past to current  $y^{(t)}$ . In this case outputs  $y^{(i)}$  are conditionally independent given the  $x$  values
- 2. Do feed the actual  $y$  values (not their prediction, but the actual observed or generated values) back into the network  
the directed PGM contains edges from all  $y^{(i)}$  values in the past to the current  $y^{(t)}$  value.

# A simple example

- RNN models a sequence of scalar random variables  $Y = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}$  with no additional inputs  $\mathbf{x}$ .
  - Input at time step  $t$  is simply the output at time step  $t-1$
  - RNN defines a directed graphical model over  $\mathbf{y}$  variables
- Parameterize joint distribution of observations using the chain rule for conditional probabilities:

$$P(Y) = P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}) = \prod_{t=1}^{\tau} P(\mathbf{y}^{(t)} \mid \mathbf{y}^{(t-1)}, \mathbf{y}^{(t-2)}, \dots, \mathbf{y}^{(1)})$$

where the right hand side of the bar  $|$  is empty for  $t=1$

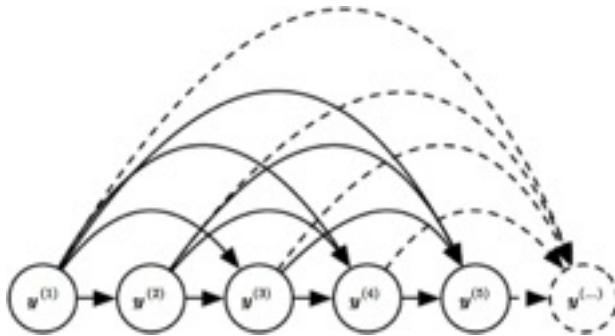
- Hence the negative log-likelihood of a set of values  $\{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)}\}$  according to such a model is

$$L = \sum_t L^{(t)}$$

$$\text{where } L^{(t)} = -\sum \log P(\mathbf{y}^{(t)} = \mathbf{y}^{(t)} \mid \mathbf{y}^{(t-1)}, \mathbf{y}^{(t-2)}, \dots, \mathbf{y}^{(1)})$$

## Fully connected graphical model for a sequence

- Every past observation  $y^{(i)}$  may influence the conditional distribution of some  $y^{(t)}$  for  $t > 1$ 
  - Parameterizing this way is inefficient.



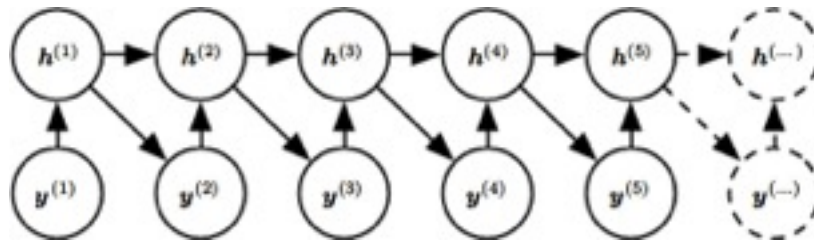
- RNNs obtain the same full connectivity but efficient parameterization as shown next

## Introducing the state variable in the PGM of an RNN

### Introduce the state variable in the PGM of RNN

even though it is a deterministic function of its inputs, helps see we get efficient parameterization

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



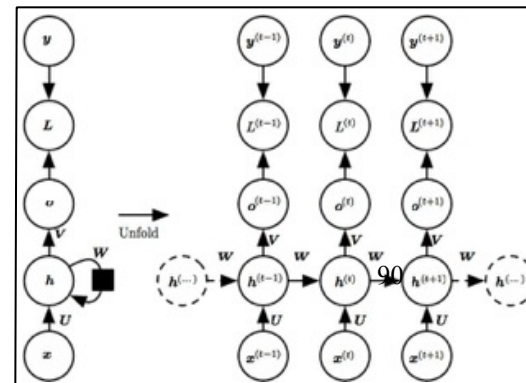
Every stage in the sequence for  $h^{(t)}$  and  $y^{(t)}$  involves the same structure and can share the parameters with other stages



- If  $y$  can take on  $k$  different values, the tabular representation would have  $O(k^\tau)$  parameters
  - By comparison because of parameter sharing the no of parameters in the RNN is  $O(1)$  as a function of sequence length
- Even with efficient parameterization of PGM, some operations are computationally challenging
  - Difficult to predict missing values in the middle of the sequence
- Price RNNs pay for reduced no of parameters is that optimizing parameters may be difficult

## Stationarity in RNNs

- Parameter sharing in RNNs relies on assuming same parameters can be used in different time steps
- Equivalently CPD over variables at time  $t+1$  given the variables at time  $t$  is stationary
- In principle it is possible to use  $t$  as an extra input at each time step
  - Learner discovers time dependence



## How to draw samples from the model

- Simply sample from the conditional distribution at each time step
- One additional complication
  - RNN must have some mechanism for determining the length of the sequence

## Methods for predicting end of sequence

- This can be achieved in various ways
  1. Adding a special symbol for end of sequence
    - When the symbol is generated, the sampling process stops
  2. Introduce an extra Bernoulli output
    - It represents the decision to either continue generation or halt generation at each time step
      - More general than previous method as it can be added to any RNN rather than only RNNs that output a sequence of symbols, e.g., RNN that emits sequence of real numbers
  3. Add an extra output to the model to predict  $\tau$ 
    - Model can sample value of  $\tau$  and then sample  $\tau$  steps worth of data

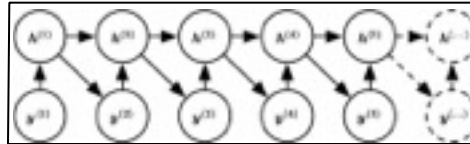
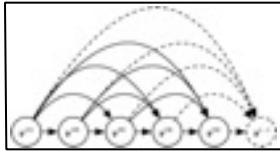
# Modeling Sequences Conditioned on Context with RNNs

# Graphical models of RNNs *without/with* inputs

## 1. Directed graphical models of RNNs *without* inputs

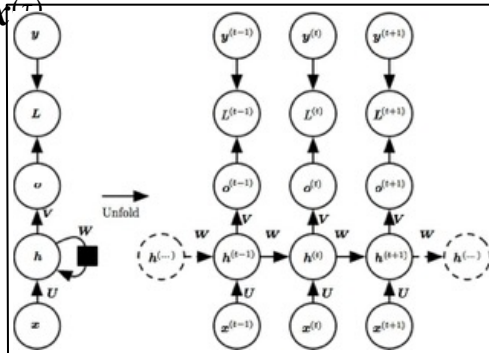
- having a set of random variables  $y^{(t)}$ :

Fully connected graphical model Efficient parameterization based on  $h^{(t)} = f(h^{(t-1)}, x^{(t)}; \theta)$



## 2. RNNs do include a sequence of inputs $x^{(1)},$

$x^{(2)}, \dots, x^{(T)}$

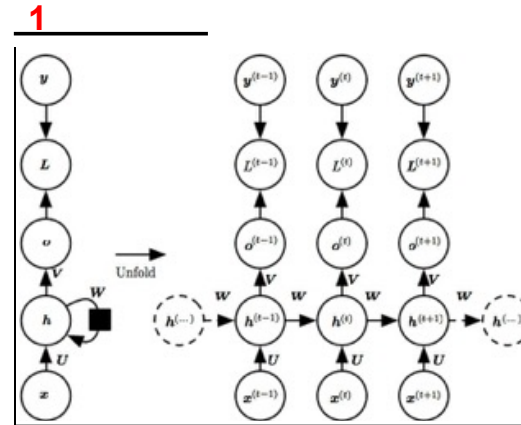


$$\begin{aligned} o^{(t)} &= c + Vh^{(t)} \\ h^{(t)} &= \tanh(a^{(t)}) \\ a^{(t)} &= b + Wh^{(t-1)} + Ux^{(t)} \end{aligned}$$

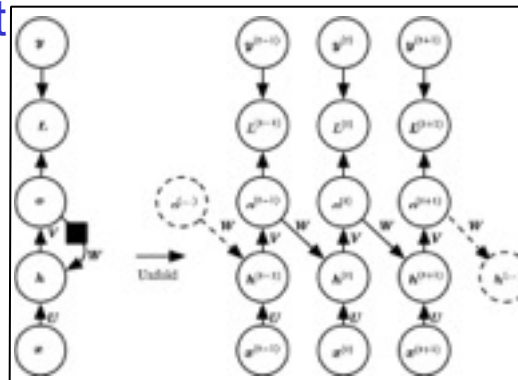
- RNNs allow the extension of the graphical model view to represent not only the joint distribution view over  $y$  variables but also a conditional distribution over  $y$  given  $x$

# CPDs of model depend on RNN Design pattern

1. Recurrent connections between hidden units



2. Recurrent connections only from output at one time step to hidden units at next time step



## Extending RNNs to represent conditional $P(\mathbf{y}|\mathbf{x})$

- A model representing a variable  $P(\mathbf{y}; \boldsymbol{\theta})$  can be reinterpreted as a model representing a conditional distribution  $P(\mathbf{y}|\boldsymbol{\omega})$  with  $\boldsymbol{\omega}=\boldsymbol{\theta}$
- We can extend such a model to represent a distribution  $P(\mathbf{y}|\mathbf{x})$  by using the same  $P(\mathbf{y}|\boldsymbol{\omega})$  as before but making  $\boldsymbol{\omega}$  a function of  $\mathbf{x}$
- In the case of an RNN this can be achieved in several ways
  - Most common choices are described next



## Taking a single vector $\mathbf{x}$ as an extra input

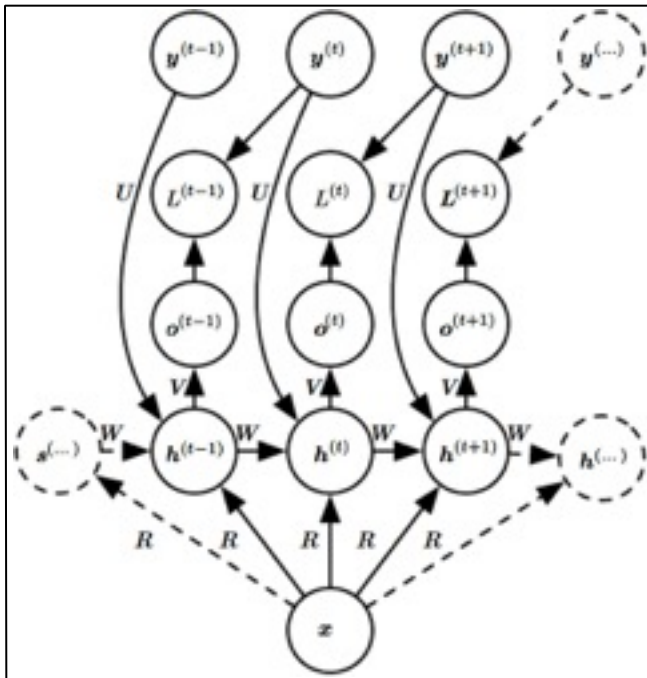
- Instead of taking a sequence  $\mathbf{x}^{(t)}$ ,  $t = 1, \dots, \tau$  as input we can take a single vector  $\mathbf{x}$  as input
- When  $\mathbf{x}$  is a fixed-size vector we can simply make it an extra input of the RNN that generates the  $\mathbf{y}$  sequence
- Common ways of providing an extra input to RNN are
  1. An extra input at each time step, or
  2. As the initial state  $\mathbf{h}^{(0)}$ , or
  3. Both
- The first and common approach is illustrated next

## Mapping vector $x$ into distribution over sequences $Y$

### An extra input at each time step

Interaction between input  $x$  and hidden unit vector  $h^{(t)}$

is parameterized by a newly introduced weight matrix  $R$   
that was absent from the model with only  $y$  values



Appropriate for tasks such as  
image captioning

where a single image  $x$  is input which  
produces a sequence of words describing  
the image

Each element of the observed output  
 $y^{(t)}$  of the observed output sequence  
serves both as input (for the current  
time step) and during training as target

## RNN to receive a sequence of vectors $\mathbf{x}^{(t)}$ as input

- RNN described by  $\mathbf{a}^{(t)} = \mathbf{b} + W\mathbf{h}^{(t-1)} + U\mathbf{x}^{(t)}$  corresponds to a conditional distribution  $P(\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(\tau)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(\tau)})$
- It makes a conditional independence assumption that this distribution factorizes as

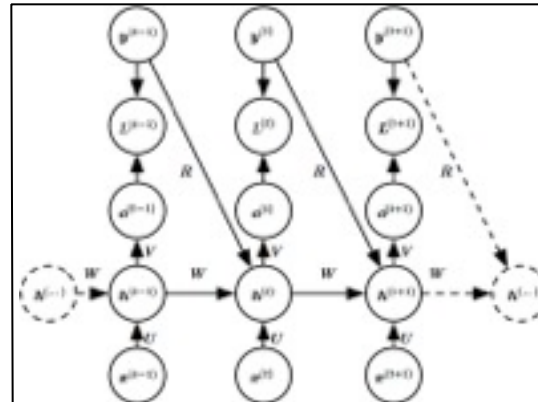
$$\prod_t P(\mathbf{y}^{(t)} | \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(t)})$$

- To remove the conditional independence assumption, we can add connections from the output at time  $t$  to the hidden unit at time  $t+1$  (see next slide)
  - The model can then represent arbitrary probability distributions over the  $\mathbf{y}$  sequence
- Limitation: both sequences must be of same length
  - Removing this restriction is discussed in Section 10.4

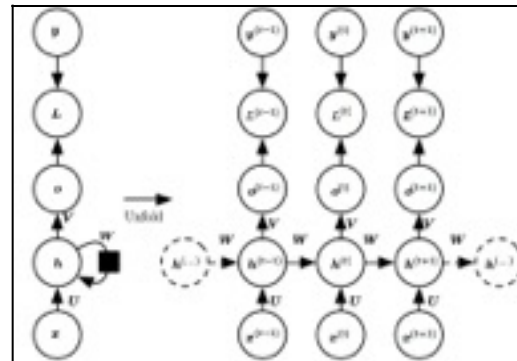
## Removing conditional independence assumption

Connections from previous output to current state allow RNN to model arbitrary distribution over sequences of  $y$

Conditional RNN mapping a variable-length sequence of  $x$  values into a distribution over sequences of  $y$  values of same length



Compare it to model that is only able to represent distributions in which the  $y$  values are conditionally independent from each other given  $x$  values



# BidirectionalRNNs

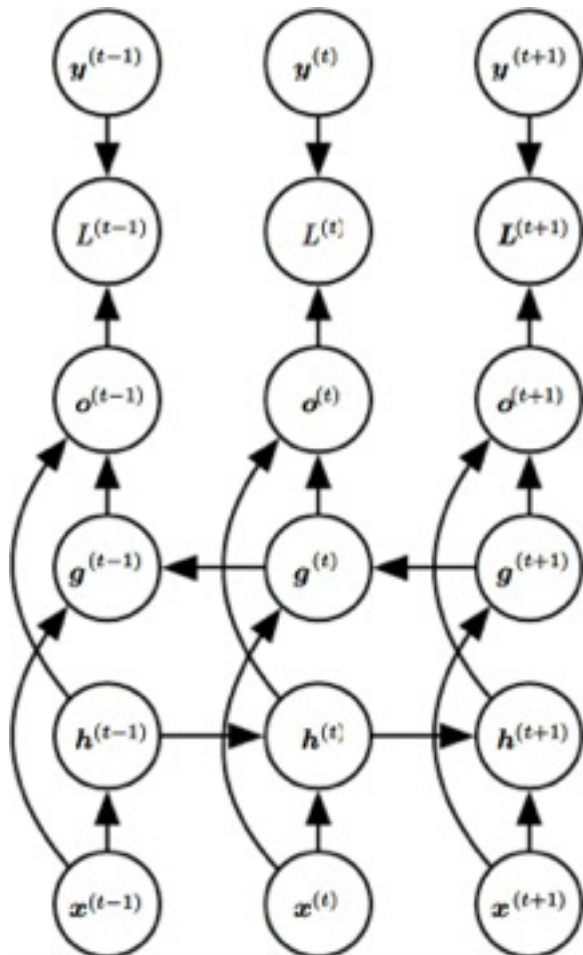
## Need for bidirectionality

- In speech recognition, the correct interpretation of the current sound may depend on the next few phonemes because of co-articulation and the next few words because of linguistic dependencies
- Also true of handwriting recognition

## A bidirectional RNN

- Combine an RNN that moves forward through time from the start of the sequence
- Another RNN that moves backward through time beginning from the end of the sequence
- A bidirectional RNN consists of two RNNs which are stacked on the top of each other.
  - The one that processes the input in its original order and the one that processes the reversed input sequence.
  - The output is then computed based on the hidden state of both RNNs.

## A typical bidirectional RNN



Maps input sequences  $x$  to target sequences  $y$  with loss  $L^{(t)}$  at each step  $t$

$h$  recurrence propagates to the right  
 $g$  recurrence propagates to the left.

This allows output units  $o^{(t)}$  to compute a representation that depends both the past and the future

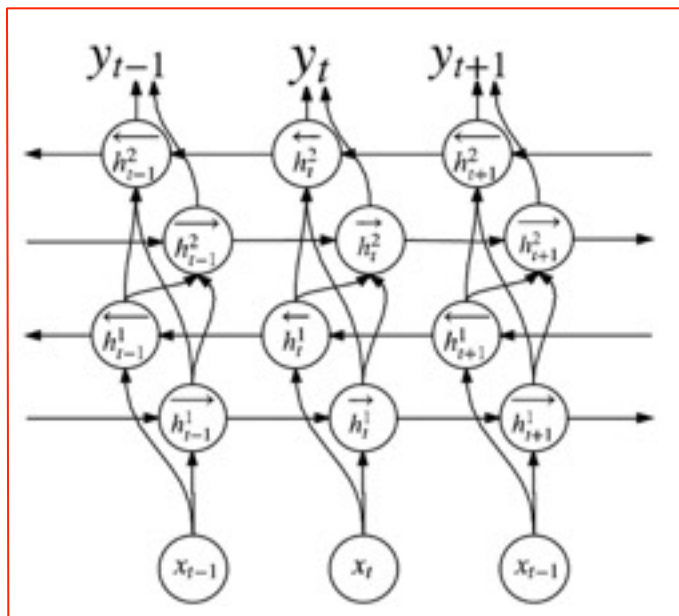
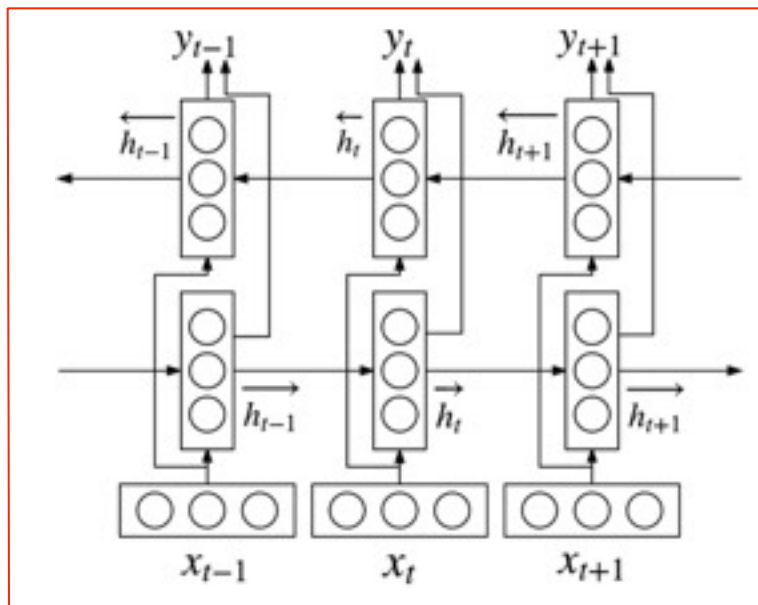


## Parameters of a Bidirectional RNN

Bidirectional RNN

Deep Bidirectional RNN

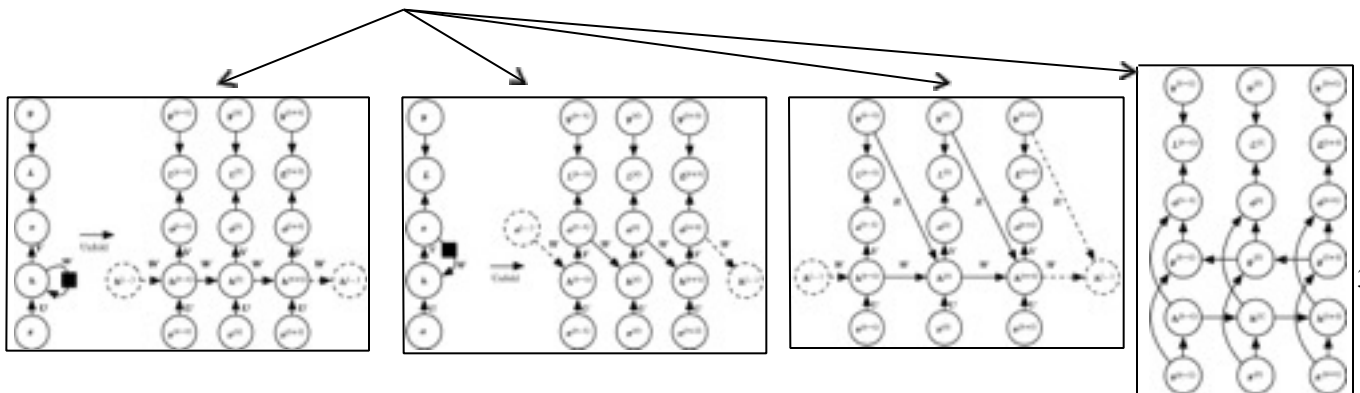
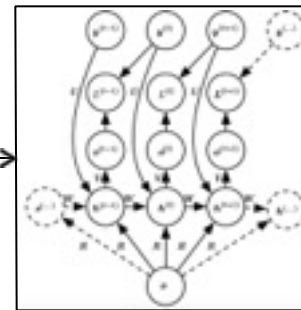
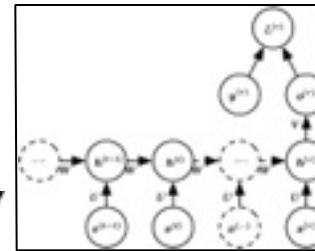
Has multiple layers per time step



# Encoder-Decoder Sequence-to-Sequence Architectures

## Previously seen RNNs

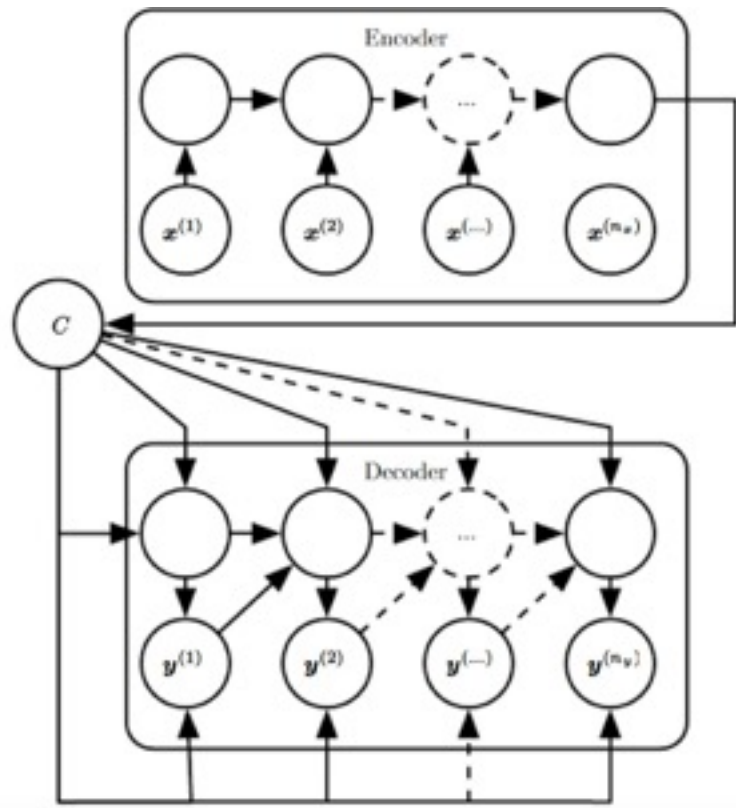
1. An RNN to map an input sequence to a fixed-size vector
2. RNN to map a fixed-size vector to a sequence
3. How an RNN can map an input sequence to an output sequence of the same length



## RNN when input/output are not same length

- Here we consider how an RNN can be trained to map an input sequence to an output sequence which is not necessarily the same length
- Comes up in applications such as speech recognition, machine translation or question-answering where the input and output sequences in the training set are generally not of the same length (although their lengths might be related)

## An Encoder-Decoder or Sequence-to-Sequence RNN



Learns to generate an output sequence  
 $(y^{(1)}, \dots, y^{(n_y)})$

given an input sequence  
 $(x^{(1)}, \dots, x^{(n_x)})$

It consists of an encoder RNN that reads an input sequence and a decoder RNN that generates the output sequence or computes the probability of a given output sequence)

The final hidden state of the encoder RNN is used to compute a fixed size context  $C$  which represents a semantic summary of the input sequence and is given as input to the decoder

# Deep Recurrent Networks

- Recurrent Neural Networks
  1. Unfolding Computational Graphs
  2. Recurrent Neural Networks
  3. Bidirectional RNNs
  4. Encoder-Decoder Sequence-to-Sequence Architectures
  5. Deep Recurrent Networks
  6. Recursive Neural Networks
  7. The Challenge of Long-Term Dependencies
  8. Echo-State Networks
  9. Leaky Units and Other Strategies for Multiple Time Scales
  10. LSTM and Other Gated RNNs
  
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory

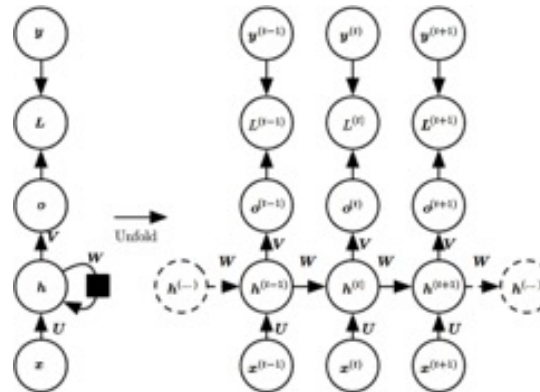
## Computation in RNNs: parameter blocks

- The computation in most recurrent neural networks can be decomposed into three blocks of parameters and associated transformations:
  1. From the input to the hidden state
  2. From the previous hidden state to the next hidden state
  3. From the hidden state to the output



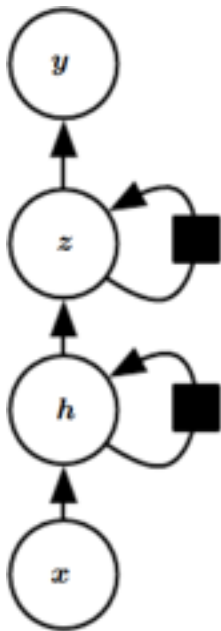
## Blocks of parameters as a shallow transformation

- With the RNN architecture shown each of these three blocks is associated with a single weight matrix, i.e.,
  - When the network is unfolded, each of these corresponds to a shallow transformation.
  - By a shallow Transformation we mean a transformation that would be represented a single layer within a deep MLP.
  - Typically this is a transformation represented by a learned affine transformation followed by a fixed nonlinearity
- Would it be advantageous to introduce depth into each of these operations?
  - Experimental evidence strongly suggests so.
    - That we need enough depth in order to perform the required transformations

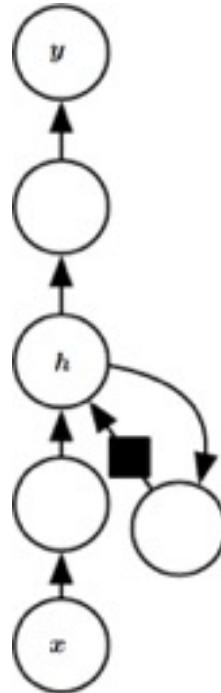


## Ways of making an RNN deep

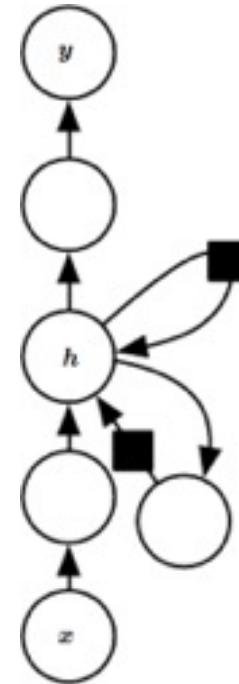
1. Hidden recurrent state can be broken down into groups organized hierarchically



2. Deeper computation can be introduced in the input-hidden, hidden-hidden and hidden-output parts. This may lengthen the shortest path linking different time steps

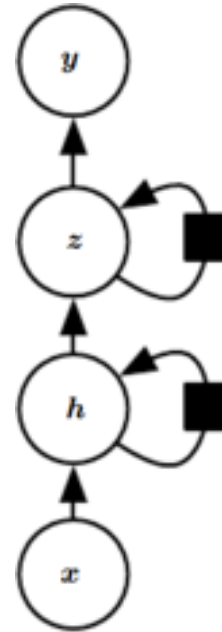


3. The path-lengthening effect can be mitigated by introducing skip connections.



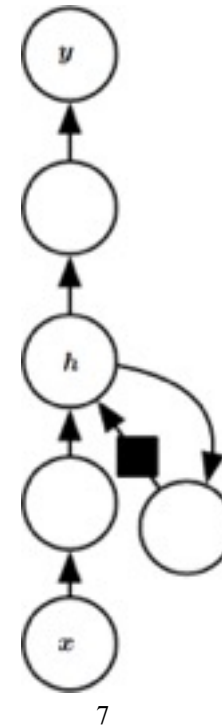
## 1. Recurrent states broken down into groups

We can think of lower levels of the hierarchy play a role of transforming the raw input into a representation that is more appropriate at the higher levels of the hidden state



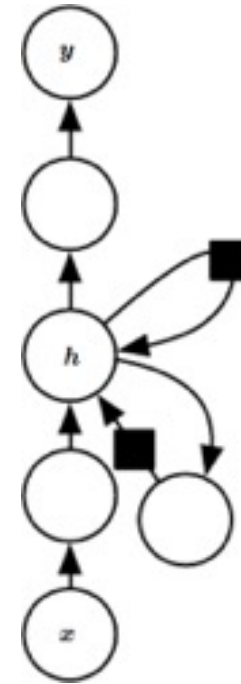
## 2. Deeper computation in hidden-to-hidden

- Go a step further and propose to have a separate MLP (possibly deep) for each of the three blocks:
  1. From the input to the hidden state
  2. From the previous hidden state to the next hidden state
  3. From the hidden state to the output
- Considerations of representational capacity suggest that to allocate enough capacity in each of these three steps
  - But doing so by adding depth may hurt learning by making optimization difficult
  - In general it is easier to optimize shallower architectures
  - Adding the extra depth makes the shortest time of a variable from time step  $t$  to a variable in time step  $t+1$  become longer



### 3. Introducing skip connections

- For example, if an MLP with a single hidden layer is used for the state-to-state transition, we have doubled the length of the shortest path between variables in any two different time steps compared with the ordinary RNN.
- This can be mitigated by introducing skip connections in the hidden-to-hidden path as illustrated here



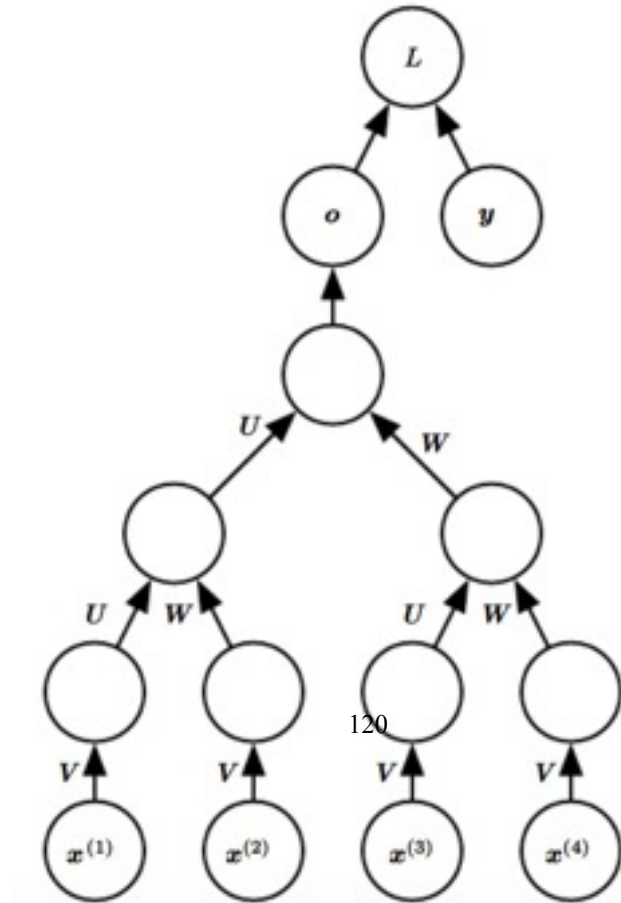
# Recursive Neural Networks

## Recursive Neural Networks

- They are yet another generalization of recurrent networks with a different kind of computational graph
- It is structured as a deep tree, rather than the chain structure of RNNs
- The typical computational graph for a recursive network is shown next

## Computational graph of a Recursive Network

- It generalizes a recurrent network from a chain to a tree
- A variable sequence  $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(t)}$  can be mapped to a fixed size representation (the output  $\mathbf{o}$ ), with a fixed set of parameters (the weight matrices  $U, V, W$ )
- Figure illustrates supervised learning case in which target  $\mathbf{y}$  is provided that is associated with the whole sequence





## Advantage of Recursive over Recurrent Nets

- For a sequence of the same length  $\tau$ , the depth (measured as the no. of compositions of nonlinear operations) can be reduced from  $\tau$  to  $O(\log \tau)$ , which might help deal with long-term dependencies
- An open question is how best to structure the tree

## Need for Recursive nets in NLP

- Deep learning based methods learn low-dimensional, real-valued vectors for word tokens, mostly from a large data corpus, successfully capturing syntactic and semantic aspects of text
- For tasks where the inputs are larger text units, e.g., phrases, sentences or documents, a compositional model is first needed to aggregate tokens into a vector with fixed dimensionality that can be used for other NLP tasks
- Models for achieving this fall into two categories: recurrent models and recursive models

## Recurrent Model for NLP

- Recurrent models deal successfully with time series data
- The recurrent models generally consider no linguistic structure aside from the word order
- They were applied early on to NLP by modeling a sentence as tokens processed sequentially and at each step combining the current token with previously built embeddings
- Recurrent models can be extended to bidirectional ones from both left to right and right to left
- These models consider no linguistic structure aside from word order

## Recursive Models for NLP

- Recursive neural models (also referred to as tree models) by contrast are structured by syntactic parse trees
- Instead of considering tokens sequentially, recursive models combine neighbors based on the recursive structure of parse trees, starting from the leaves and proceeding recursively in a bottom-up fashion until the root of the parse tree is reached
  - Ex: for the phrase the food is delicious, following the operation sequence ((the food) (is delicious)) rather than the sequential order (((the food) is) delicious)

## Advantage of Recursive Model for NLP

- They have the potential of capturing long-distance dependencies
- Two tokens may be structurally closer to each other even though they are far away in word sequence
- Ex: a verb and its corresponding direct object can be far away in terms of tokens if many adjectives lie inbetween, but they are adjacent in the parse tree
- However parsing is slow and domain dependent
- See performance comparison with LSTM on four NLP tasks at [https://nlp.stanford.edu/pubemnlp2015\\_2\\_jiwei.pdf](https://nlp.stanford.edu/pubemnlp2015_2_jiwei.pdf)

## Structure of the Tree

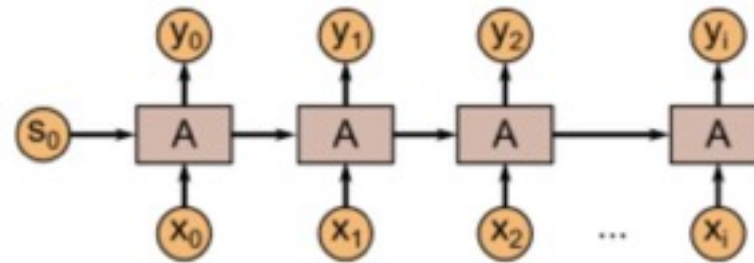
- One option is to have a tree structure that does not depend on the data, such as a balanced binary tree
- In some application domains, external methods can suggest the appropriate tree structure
  - Ex: when processing natural language sentences, the tree structure for the recursive network can be fixed to the structure of the parse tree of the sentence provided by a natural language parser
- Ideally, one would like the learner itself to discover and infer the tree structure that is appropriate for any given input

## Variants of Recursive Net idea

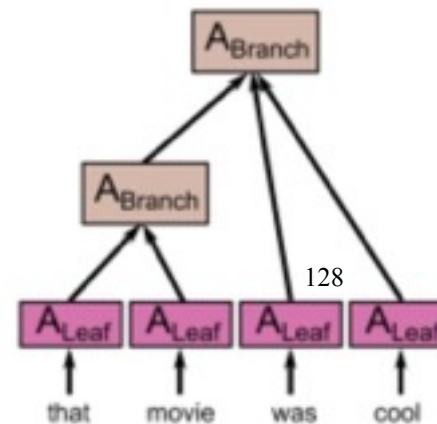
- Associate data with a tree structure and associate inputs and targets with individual nodes of the tree
  - The computation performed for each node does not have to be the artificial neuron computation (affine transformation of all inputs followed by a monotone nonlinearity)
  - Can use a tensor operations of bilinear forms
    - Previously found useful to model linear relationships between concepts when the concepts are represented by continuous vectors

# Recursive Neural Networks

- Recursive neural networks are also called Tree Nets
  - Useful for learning tree-like structures
  - They are highly useful for parsing natural scenes and language



Recurrent Neural Net



Recursive Neural Net



## Unrolling Recurrent and Tree Nets

- In RNNs, at each time step the network takes as input its previous state  $s^{(t-1)}$  and its current input  $x^{(t)}$  and produces an output  $y^{(t)}$  and a new hidden state  $s^{(t)}$ .
- TreeNets, on the other hand, don't have a simple linear structure like that.
- With RNNs, you can 'unroll' the net and think of it as a large feedforward net with inputs  $x^{(0)}, x^{(1)}, \dots, x^{(T)}$ , initial state  $s^{(0)}$ , and outputs  $y^{(0)}, y^{(1)}, \dots, y^{(T)}$ , with  $T$  varying depending on the input data stream, and the weights in each of the cells tied with each other.
- You can also think of TreeNets by unrolling them – the weights in each branch node are tied with each other, and the weights in each leaf node are tied with each other.

## Advantage of Recursive Nets

- The advantage of Recursive Nets is that they can be very powerful in learning hierarchical, tree-like structure.
- The disadvantages are, firstly, that the tree structure of every input sample must be known at training time.

# Long-Term Dependencies: The Challenge

## Topics in Sequence Modeling

- Recurrent Neural Networks
  1. Unfolding Computational Graphs
  2. Recurrent Neural Networks
  3. Bidirectional RNNs
  4. Encoder-Decoder Sequence-to-Sequence Architectures
  5. Deep Recurrent Networks
  6. Recursive Neural Networks
  7. The Challenge of Long-Term Dependencies
  8. Echo-State Networks
  9. Leaky Units and Other Strategies for Multiple Time Scales
  10. LSTM and Other Gated RNNs
  
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory

## Challenge of Long-Term Dependencies

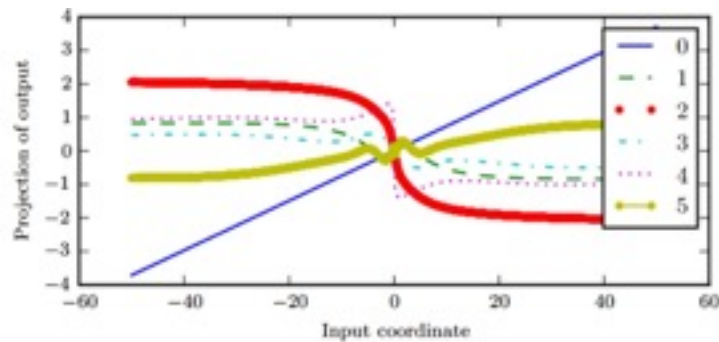
- Neural network optimization face a difficulty when computational graphs become deep, e.g.,
  - Feedforward networks with many layers
  - RNNs that repeatedly apply the same operation at each time step of a long temporal sequence
- Gradients propagated over many stages tend to either vanish (most of the time) or explode (damaging optimization)
- The difficulty with long-term dependencies arise from exponentially smaller weights given to long-term interactions (involving multiplication of many Jacobians)

## Vanishing and Exploding Gradient Problem

- Suppose a computational graph consists of repeatedly multiplying by a matrix  $W$
- After  $t$  steps this is equivalent to multiplying by  $W^t$
- Suppose  $W$  has an eigendecomposition  $W = V \text{diag}(\lambda) V^{-1}$ 
  - In this case it is straightforward to see that
$$W^t = (V \text{diag}(\lambda) V^{-1})^t = V \text{diag}(\lambda)^t V^{-1}$$
  - Any eigenvalues  $\lambda_i$  that are not near an absolute value of 1 will either explode if they are greater than 1 in magnitude and vanish if they are less than 1 in magnitude
- Vanishing gradients make it difficult to know which direction the parameters should move to improve cost
- Exploding gradients make learning unstable

## Function Composition in RNNs

- RNNs involve composition of the same function multiple times, one per step
- These compositions can result in extremely nonlinear behavior



Composing many nonlinear functions:

$\tanh$  here

$\mathbf{h}$  has 100 dimensions  
mapped to a single  
dimension

Most of the space it has a  
small derivative and highly  
nonlinear elsewhere

- Problem particular to RNNs

# Leaky Units and Multiple Time Scales



- Recurrent Neural Networks
  1. Unfolding Computational Graphs
  2. Recurrent Neural Networks
  3. Bidirectional RNNs
  4. Encoder-Decoder Sequence-to-Sequence Architectures
  5. Deep Recurrent Networks
  6. Recursive Neural Networks
  7. The Challenge of Long-Term Dependencies
  8. Echo-State Networks
  9. Leaky Units and Other Strategies for Multiple Time Scales
  
- 10. LSTM and Other Gated RNNs
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory

## Multiple Time Scales

- Goal: to deal with long-term dependencies
- Design model so that
  1. Some model parts operate at fine-grained time scales
    - Handle small details
  2. Other parts operate at coarse time scales
    - Transfer information from the distant past to the present more efficiently

## Building Fine and Coarse Time Scales

- Strategies to build fine and coarse time scales
  1. Adding skip connections through time
  2. Leaky units and a spectrum of different time scales
    - To integrate signals with different time constants
  3. Removal of some connections to model fine-grained time scales

## Adding skip connections through time

- Add direct connections from variables in the distant past to variables in the present
- In an ordinary RNN, recurrent connection goes from time  $t$  to time  $t+1$ . Can construct RNNs with longer delays
- Gradients can vanish/explode exponentially wrt no. of time steps
- Introduce time delay of  $d$  to mitigate this problem
- Gradients diminish as a function of  $\tau/d$  rather than  $\tau$
- Allows learning algorithm to capture longer dependencies
  - Not all long-term dependencies can be captured this way

## Leaky units and a spectrum of time scales

- Rather than an integer skip of  $d$  time steps, the effect can be obtained smoothly by adjusting a real-valued  $\alpha$
- Running Average
  - Running average  $\mu^{(t)}$  of some value  $v^{(t)}$  is  $\mu^{(t)} \leftarrow \alpha \mu^{(t-1)} + (1-\alpha)v^{(t)}$
  - Called a linear self-correction
  - When  $\alpha$  is close to 1, running average remembers information from the past for a long time and when it is close to 0, information is rapidly discarded.
- Hidden units with linear self connections behave similar to running average. They are called *leaky units*.
- Can obtain product of derivatives close to 1 by having *linear* self-connections and a weight near 1 on those connections<sup>14</sup>

## Removing Connections

- Another approach to handle long-term dependencies
- Organize state of the RNN at multiple time scales
  - Information flowing more easily through long distances at the slower time scales
- It involves actively removing length one connections and replacing them with longer connections
- Skip connections add edges

# Long-Short Term Memory and Other Gated RNNs

## Topics in Sequence Modeling

- Recurrent Neural Networks
  1. Unfolding Computational Graphs
  2. Recurrent Neural Networks
  3. Bidirectional RNNs
  4. Encoder-Decoder Sequence-to-Sequence Architectures
  5. Deep Recurrent Networks
  6. Recursive Neural Networks
  7. The Challenge of Long-Term Dependencies
  8. Echo-State Networks
  9. Leaky Units and Other Strategies for Multiple Time Scales
- 10. LSTM and Other Gated RNNs
- 11. Optimization for Long-Term Dependencies
- 12. Explicit Memory

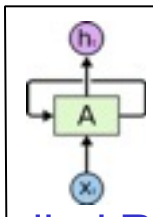


## Topics in LSTM and Other Gated RNNs

- From RNN to LSTM
- Problem of Long-Term Dependency
- BPTT and Vanishing Gradients
- Key intuition of LSTM as “state”
- LSTM Cell (three gates and two activations)
- LSTM usage example
- Equations for: forget gate, state update, output gate
- Step-by-Step LSTM Walk-through
- Gated RNNs

# Recurrent Neural Networks

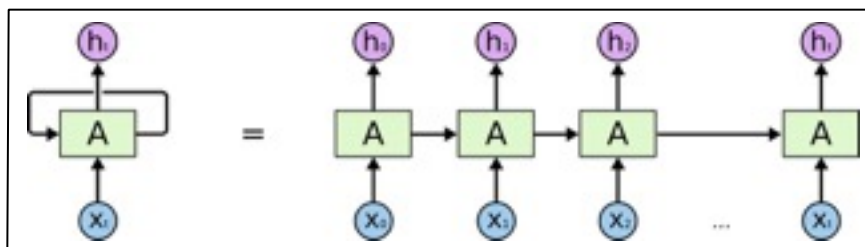
- RNNs have loops



A chunk of neural network  $A$  looks at some input  $x_t$  and outputs a value  $h_t$

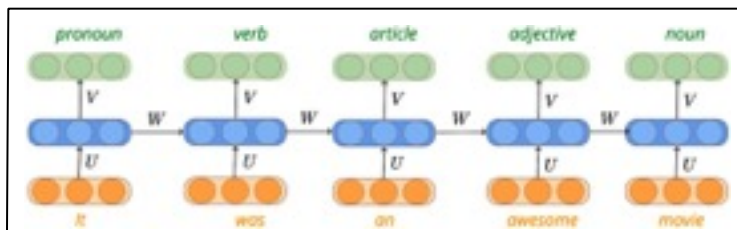
A loop allows information to be passed from one step of the network to the next

- An unrolled RNN

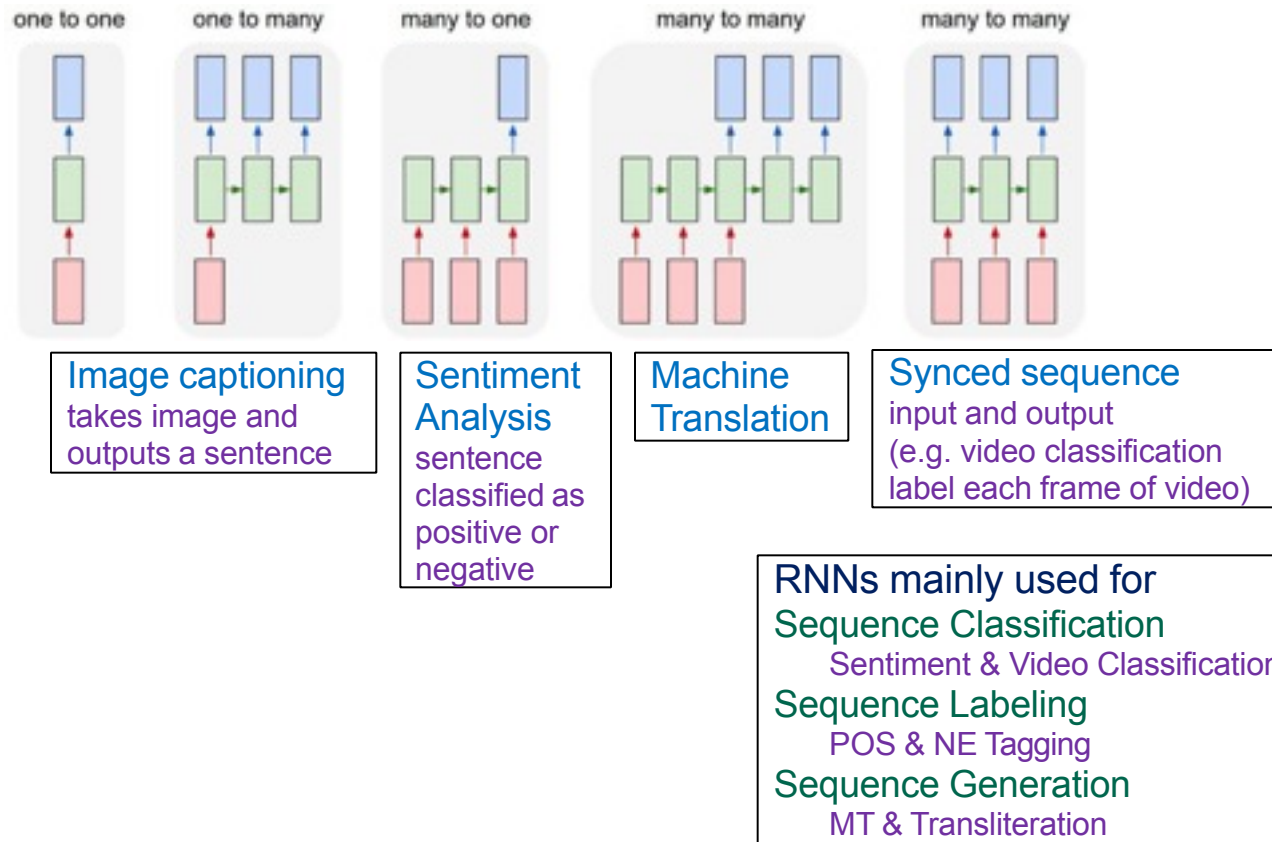


Chain-like structure reveals that RNNs are intimately related to sequences and lists

- Application to Part of Speech (POS) Tagging

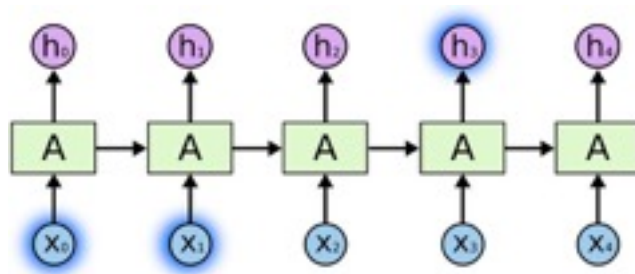


# Different Types of RNNs



## Prediction with only recent previous information

- RNNs connect previous information to the present task
  - Previous video frames may help understand present frame
- Sometimes we need only look at recent previous information to predict
  - To predict the last word of “The clouds are in the *sky*” we don’t need any further context. It is obvious that the word is “sky”



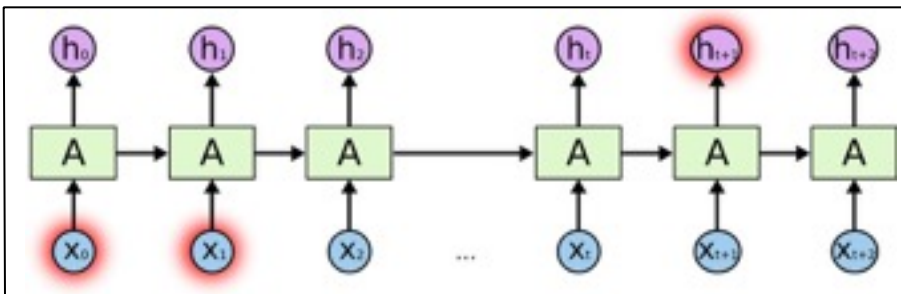
$h_3$  depends on  $x_0$  and  $x_1$

## Long-Term Dependency

- RNNs work upon the fact that the result of an information is dependent on its previous state or previous  $n$  time steps.
- They have difficulty in learning long range dependencies.
- **Example sentence:** The man who ate my pizza has purple hair
  - **Note:** purple hair **is for the** man **and not** the pizza.
  - **So this is a long dependency**

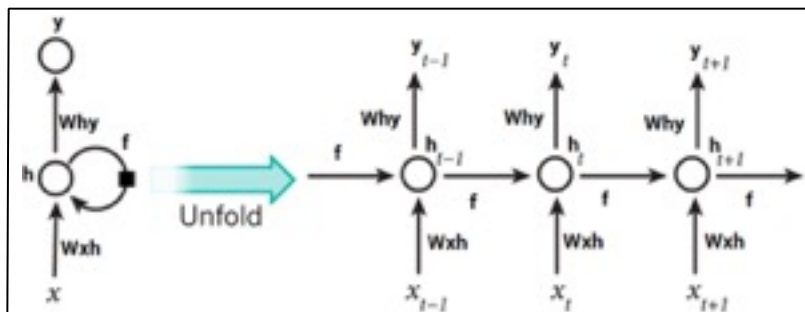
## Long-term dependency and Backprop

- There are cases where we need even more context
  - To predict last word in I grew up in France.....I speak *French*
  - Using only recent information suggests that the last word is the name of a language. But more distant past indicates that it is *French*
- Gap between relevant information and where it is needed is large



- Learning to store information over extended time intervals via RNN backpropagation takes a very long time due to decaying error backflow (discussed next)

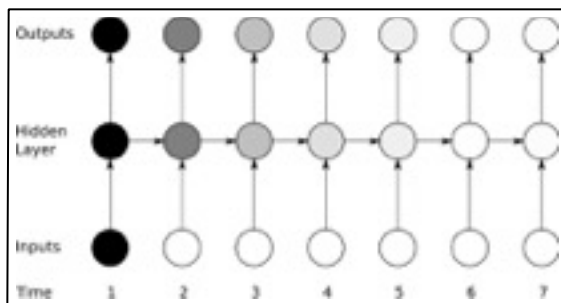
# Backpropagation Through Time (BPTT)



If  $y_t$  is predicted and  $\bar{y}_t$  is training value the cross entropy loss is

$$E_t(\bar{y}_t, y_t) = -\bar{y}_t \log(y_t)$$

Total error is the sum

$$E(\bar{y}, y) = -\sum \bar{y}_t \log(y_t)$$


To backpropagate error, apply the chain rule. The error after the third time step wrt the first:

$$\partial E / \partial W = \partial E / \partial y_3 * \partial y_3 / \partial h_3 * \partial h_3 / \partial y_2 * \partial y_2 / \partial h_1 \dots$$

and there is a long dependency.

## Vanishing gradient problem

If a single gradient approached 0, all gradients rush to zero exponentially fast due to the multiplication.

Such states would no longer help the network to learn anything.

# Vanishing and Exploding Gradients

- Vanishing gradient problem



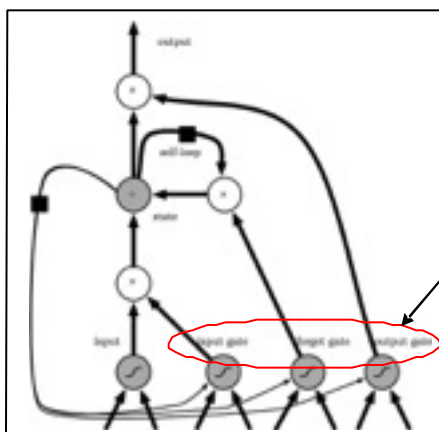
Contribution from the earlier steps becomes insignificant in the gradient

- Exploding gradient problem:
  - Gradients become very very large due to a single or multiple gradient values becoming very high.
- Vanishing gradient is more concerning
  - Exploding gradient easily solved by clipping the gradients at a predefined threshold value.
- Solution: gated RNNs
  - LSTM, GRU (Gated Recurrent Units), a variant of LSTM



# Key intuition of LSTM is “State”

- A persistent module called the cell-state
- Note that “State” is a representation of past history
- It comprises a common thread through time
- Cells are connected recurrently to each other
  - Replacing hidden units of ordinary recurrent networks

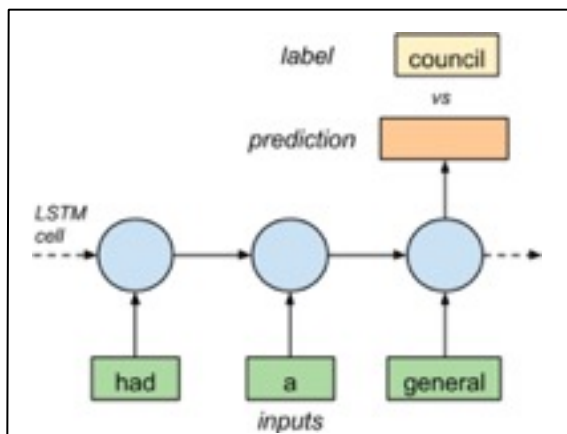


There are three sigmoid gates:  
An input gate (i),  
A forget gate (f)  
An output gate (o)

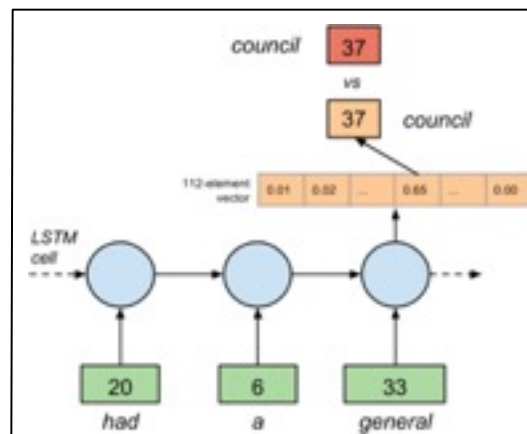
## LSTM usage example

- “long ago , the mice had a general council to consider what measures they could take to outwit their common enemy , the cat.....” Aesop’s fables story with 112 words

Training



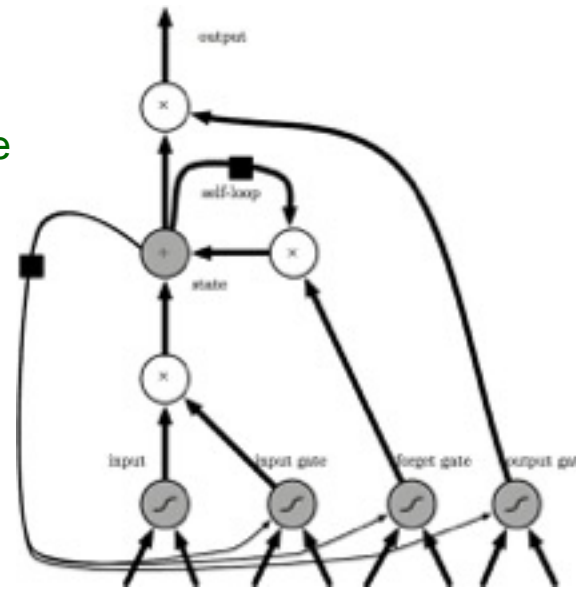
Prediction



154

## LSTM Recurrent Network Cell

- Input feature computed with regular artificial neuron unit
  - Its value can be accumulated into the state
  - If the sigmoidal input gate allows it
- State unit has linear self-loop
- Weight controlled by forget gate
- Output of cell can be shut off by output gate
- All units: sigmoid nonlinearity
  - Input gate can be any squashing
- State unit can also be used as an extra input to gating units
- Black square indicates delay of single time step

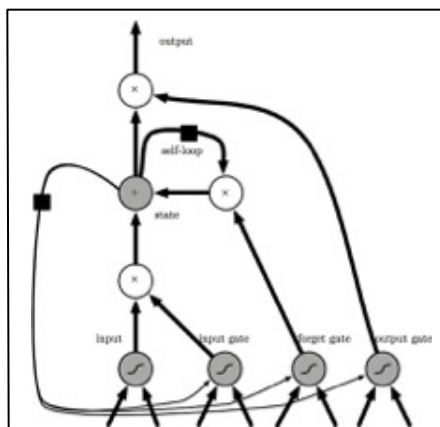


## Core contribution of LSTM

- Clever idea of introducing self-loops to produce paths where the gradient can flow for long durations
- A crucial addition: make weight on this self-loop conditioned on the context, rather than fixed
  - By making weight of this self-loop gated (controlled by another hidden unit), time-scale can be changed dynamically
  - Even for an LSTM with fixed parameters, time scale of integration can change based on the input sequence
    - Because time constants are output by the model itself

# State perturbed by linear operations

- Cell-state perturbed only by a few linear operations at each time step



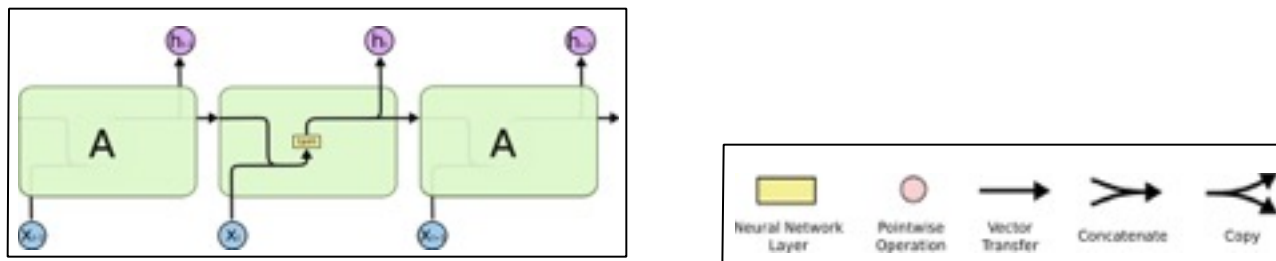
- Since cell state connection to previous cell states is interrupted only by the linear operations of multiplication and addition, LSTMs can remember short-term memories (*i.e.* activity belonging to the same “episode”) for a very long time

## LSTM success

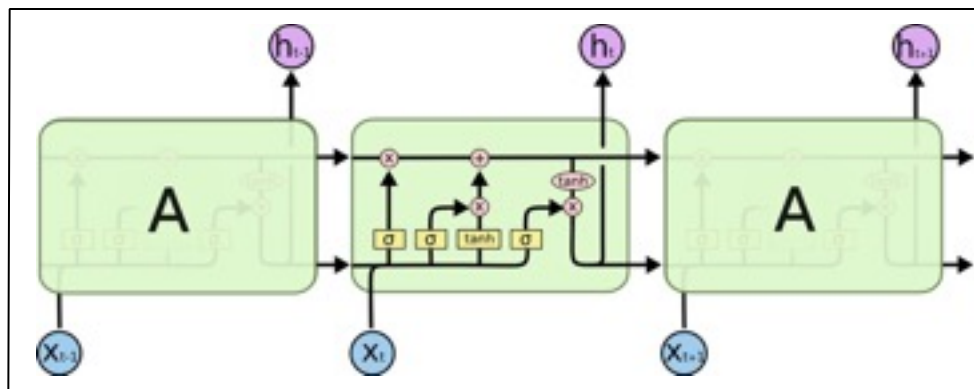
- LSTM found extremely successful in:
  - Unconstrained handwriting recognition
  - Speech recognition
  - Handwriting generation, Machine Translation
  - Image Captioning, Parsing

## RNN vs LSTM cell

- RNNs have the form of a repeating chain structure
- The repeating module has a simple structure such as tanh

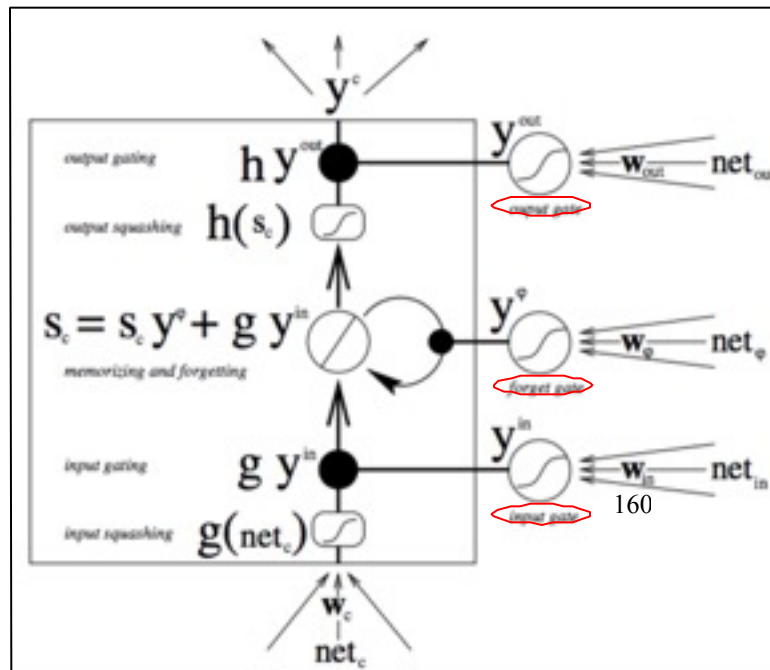
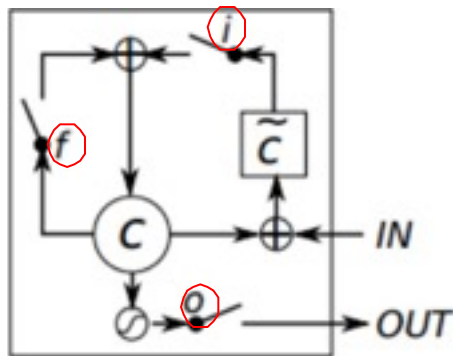


- LSTMs also have a chain structure
  - but the repeating module has a different structure



# LSTM cell

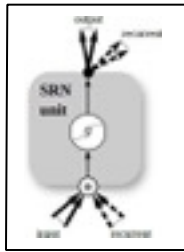
- Like RNNs , LSTMs also have a chain like structure
- But repeating module has a slightly different structure
- Instead of a single layer, multiple layers interact
  - They have three sigmoid gates:
    - An input gate (i),
    - A forget gate (f)
    - An output gate (o)





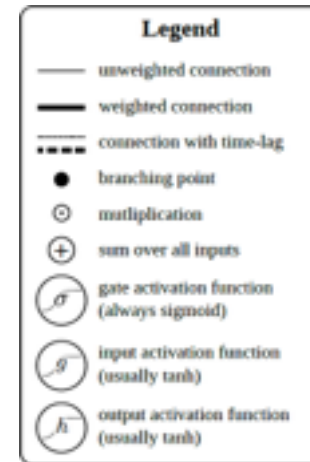
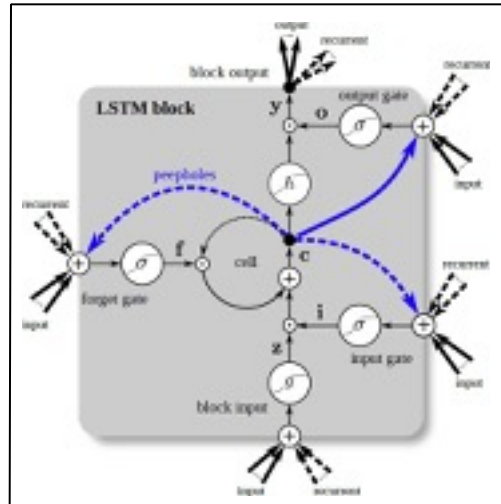
# LSTM Cell equations

Simple  
RNN  
unit



Weights for an LSTM layer:

- Input weights:  $W_z, W_i, W_f, W_o \in \mathbb{R}^{N \times M}$
- Recurrent weights:  $R_z, R_i, R_f, R_o \in \mathbb{R}^{N \times N}$
- Peephole weights:  $p_i, p_f, p_o \in \mathbb{R}^N$
- Bias weights:  $b_z, b_i, b_f, b_o \in \mathbb{R}^N$



Vector formulas for LSTM forward pass:

$$\begin{aligned}
 \tilde{z}^t &= W_z x^t + R_z y^{t-1} + b_z && \text{block input} \\
 z^t &= g(\tilde{z}^t) \\
 \tilde{i}^t &= W_i x^t + R_i y^{t-1} + p_i \odot c^{t-1} + b_i && \text{input gate} \\
 i^t &= \sigma(\tilde{i}^t) \\
 \tilde{f}^t &= W_f x^t + R_f y^{t-1} + p_f \odot c^{t-1} + b_f && \text{forget gate} \\
 f^t &= \sigma(\tilde{f}^t) \\
 c^t &= z^t \odot i^t + c^{t-1} \odot f^t && \text{cell} \\
 \tilde{o}^t &= W_o x^t + R_o y^{t-1} + p_o \odot c^t + b_o \\
 o^t &= \sigma(\tilde{o}^t) && \text{output gate} \\
 y^t &= h(c^t) \odot o^t && \text{block output}
 \end{aligned}$$

BPTT deltas inside the LSTM block:

$$\begin{aligned}
 \delta y^t &= \Delta^t + R_z^T \delta z^{t+1} + R_i^T \delta i^{t+1} + R_f^T \delta f^{t+1} + R_o^T \delta o^{t+1} \\
 \delta \tilde{o}^t &= \delta y^t \odot h(c^t) \odot \sigma'(\tilde{o}^t) \\
 \delta c^t &= \delta y^t \odot o^t \odot h'(c^t) + p_o \odot \delta \tilde{o}^t + p_i \odot \delta \tilde{i}^{t+1} \\
 &\quad + p_f \odot \delta \tilde{f}^{t+1} + \delta c^{t+1} \odot f^{t+1} \\
 \delta \tilde{f}^t &= \delta c^t \odot c^{t-1} \odot \sigma'(\tilde{f}^t) \\
 \delta \tilde{i}^t &= \delta c^t \odot z^t \odot \sigma'(\tilde{i}^t) \\
 \delta \tilde{z}^t &= \delta c^t \odot i^t \odot g'(\tilde{z}^t)
 \end{aligned}$$

Sigmoid used for gate activation, tanh used as input and output activation, point-wise multiplication of vectors is  $\odot$

Source: Greff et.al., LSTM: a search space odyssey, Transactions on Neural networks and Learning Systems, arXiv 2017

## Accumulating Information over Longer Duration

- Leaky Units Allow the network to accumulate information
  - Such as evidence for a particular feature or category
  - Over a long duration
- However, once the information has been used, it might be useful to forget the old state
  - E.g., if a sequence is made of sub-sequences and we want a leaky unit to accumulate evidence inside each sub-sequence, we need a mechanism to forget the old state by setting it to zero
- Gated RNN:
  - Instead of manually deciding when to clear the state, we want the neural network to learn to decide when to do it

## Description of the LSTM cell

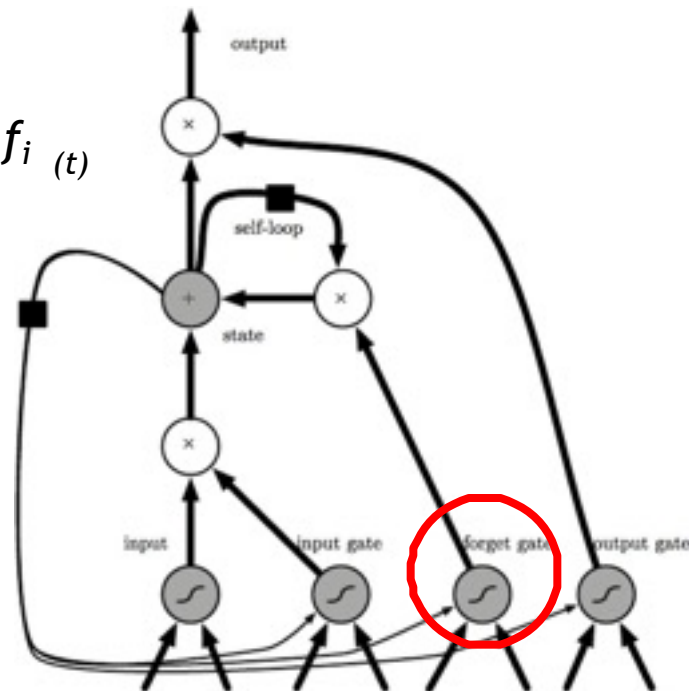
- Forward propagation equations for the LSTM cell are given in the next slide
  - In the case of a shallow recurrent network architecture
  - Deeper architectures have also been successively used
- Instead of a unit that simply applies an elementwise nonlinearity to the affine transformation of inputs and recurrent units LSTM recurrent units have *LSTM cells* that have an internal recurrence (a self-loop) in addition to the outer recurrence of the RNN
- Each cell has the same inputs and outputs as an ordinary RNN
  - But has more parameters and a system of gating units that controls the flow of information

## Equation for LSTM forget gate

- The most important component is the state unit  $s_i(t)$  that has a linear self-loop similar to the leaky units
- However the self-loop weight (or the associated time constant) is controlled by a forget gate unit  $f_i(t)$

(for time step  $t$  and cell  $i$ )  
that sets this weight to a value between 0 and 1 via a sigmoid unit

$$f_i^{(t)} = \sigma \left( b^f + \sum_j U_{ij}^f x_j^t + \sum_j W_{ij}^f h_j^{(t-1)} \right)$$



where  $\mathbf{x}^{(t)}$  is the current input vector and  $\mathbf{h}^{(t)}$  is the current hidden layer vector containing the outputs of all the LSTM cells, and  $\mathbf{b}^f$ ,  $\mathbf{u}^f$  and  $\mathbf{W}^f$  are respectively biases, input weights and recurrent weights for forget gates

## Equation for LSTM internal state update

- The LSTM cell internal state is updated as follows
  - But with conditional self-loop weight

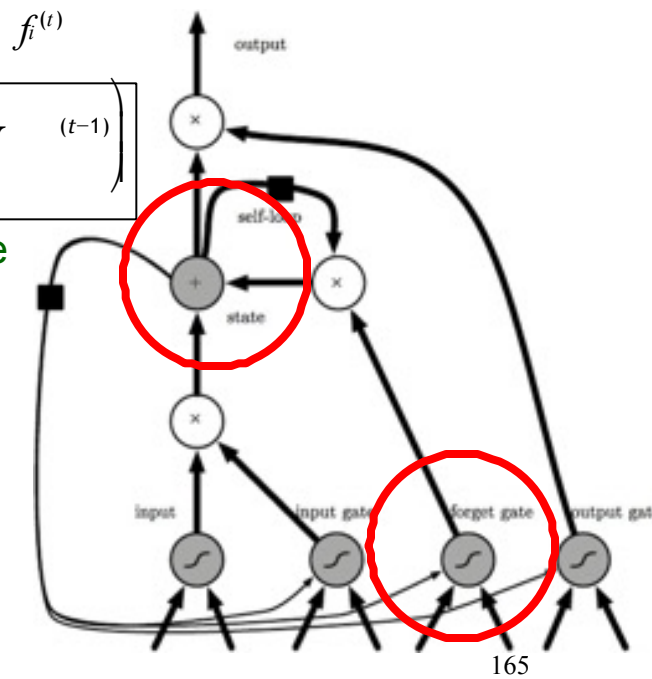
$$s_i^{(t)} = f_i^{(t)} \left( s_i^{(t-1)} + g_i^{(t)} \left( \sigma \left( \sum_j U_{ij}^f x_j^{(t)} + \sum_j W_{ij}^f h_j^{(t-1)} \right) \right) \right)$$

where  $b$ ,  $U$  and  $W$  respectively denote the biases, input weights and recurrent weights into the LSTM cell

- External input gate unit  $g_i^{(t)}$  is
  - computed similar to forget gate
  - With a sigmoid unit to obtain a gating value between 0 and 1 but with its own parameters

$$g_i^{(t)} = \sigma \left( b + \sum_j U_{ij}^g x_j^{(t)} + \sum_j W_{ij}^g h_j^{(t-1)} \right)$$

$h$



## Equation for output of LSTM cell

- Output  $h_i^{(t)}$  of the LSTM cell can be shut off via the output gate  $q_i^{(t)}$  which also uses a sigmoid unit for gating

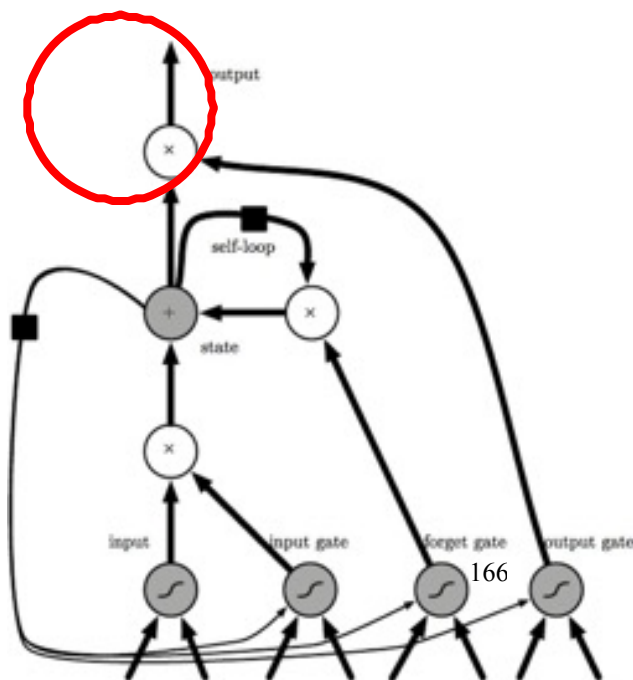
$$h_i^{(t)} = \tanh_s \left( \frac{U_i}{q_i} \right)$$

$$q_i^{(t)} = \left( b^0 + \sum_j U_j^0 + \sum_h W_h^0 \right)$$

$b^0$ ,  $U^0$  and  $W^0$  are biases, input weights and recurrent weights

- Among the variants one can choose to use the cell state  $s_i(t)$

as an extra input (with its weight) into the three gates of the  $i$ -th unit  
This would require three additional parameters



## Power of LSTMs

- LSTM networks have been shown to learn long-term dependencies more easily than simple recurrent architectures
  - First on artificial data sets
  - Then on challenging sequential tasks
- Variants and alternatives to LSTM have been studied and used and are discussed next

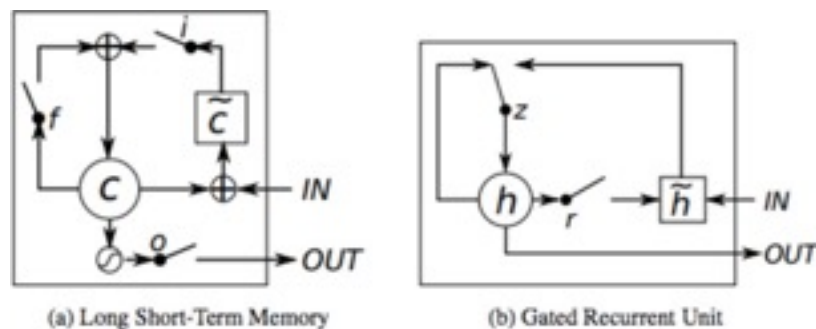
## Other gated RNNs

- What pieces of the LSTM architecture are actually necessary?
- What other successful architectures could be designed that allow the network to dynamically control the time scale and forgetting behavior of different units?
- Some answers to these questions are given with the recent work on gated RNNs, whose units are also known as gated recurrent units or GRUs
- Main difference with LSTM
  - Single gating unit simultaneously controls the forgetting factor and decision to update state unit



## GRUs

- GRUs have simpler structure and are easier to train.
- Their success is primarily due to the gating network signals that control how the present input and previous memory are used, to update the current activation and produce the current state.
- These gates have their own sets of weights that are adaptively updated in the learning phase.
- We have just two gates here, the reset and the update gate



## Idea of Gated RNNs

- Like leaky units, gated RNNs are based on the idea of creating paths through time that have derivatives that neither vanish nor explode
- Leaky units do this with connection weights that are manually chosen or were parameters  $\alpha$ 
  - generalizing the concept of discrete skipped connections
- Gated RNNs generalize this to connection weights that may change with each time step

## Update equations for Gated RNNs

- Update equations are

$$h_i^{(t)} = u_i^{(t)} h_i^{(t-1)} + (1 - u_i^{(t)}) \sigma \left( b_i + \sum_j U_{ij} x_j^{(t)} + \sum_j W_{ij}^{(t)} h_j^{(t-1)} \right)$$

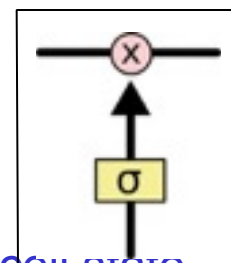
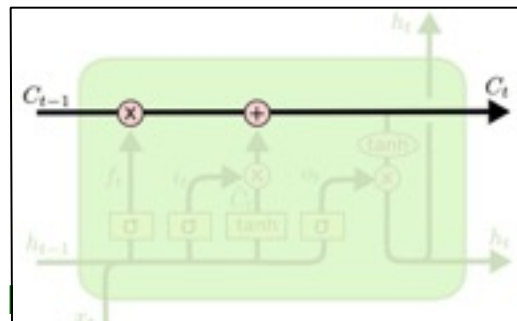
- Where  $u$  stands for the update gate and  $r$  for reset gate. Their value is defined as usual:

$$u_i^{(t)} = \sigma \left( b_i + \sum_j U_{ij}^u x_j^{(t)} + \sum_j W_{ij}^u h_j^{(t-1)} \right) \quad \text{and} \quad r_i^{(t)} = \sigma \left( b_i + \sum_j U_{ij}^r x_j^{(t)} + \sum_j W_{ij}^r h_j^{(t-1)} \right)$$

- Reset and update gates can individually ignore parts of the state vector  $h$ 
  - Update gates act conditional leaky integrators that can linearly gate any dimension thus choosing to copy it (at one extreme of sigmoid) or completely ignore it (at the other extreme) by replacing it by the new target state value (towards which the leaky integrator converges)
  - Reset gates control which parts of the state get used to compute next target state, introducing an additional nonlinear effect in the relationship between past state and future state

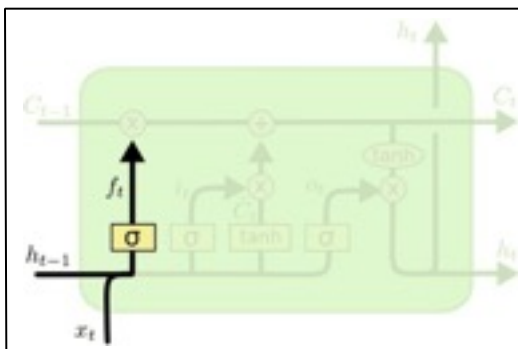
## Core idea behind LSTM

- The key to LSTM is the cell state,  $C_t$ , the horizontal line running through the top of the diagram
- Like a conveyor belt
  - Runs through entire chain with minor interactions
  - LSTM does have the ability to remove/add information to cell state regulated
- Gates are an optional way to let information through
- Consist of a sigmoid and a multiplication operation
- Sigmoid outputs a value between 0 and 1
  - 0 means let nothing through
  - 1 means let everything through
- LSTM has three of these gates, to protect and control cell state



## Step-by-step LSM walk through

- Example of language model: predict next word based on previous ones
  - Cell state may include the gender of the present subject
- First step: information to throw away from cell state



Called *forget gate layer*

It looks at  $h_{t-1}$  and  $x_t$  and outputs a number between 0 and 1 for each member of  $C_{t-1}$  for whether to forget

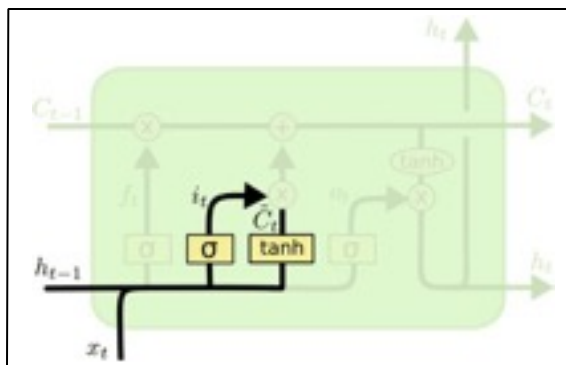
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

### In language model

consider trying to predict the next word based on all previous ones  
 The cell state may include the gender of the present subject  
 so that the proper pronouns can be used  
 When we see a new subject we want to forget old subject

## LSM walk through: Second step

- Next step is to decide as to what new information we're going to store in the cell state



In the Language model,  
we'd want to add the gender  
of the new subject to the  
cell state, to replace the old  
One we are forgetting

This has two parts:

first a sigmoid layer called *Input gate layer*.  
decides which values we will update  
Next a tanh layer creates a vector of  
new candidate values  $\tilde{C}_t$  that could be  
added to the state.

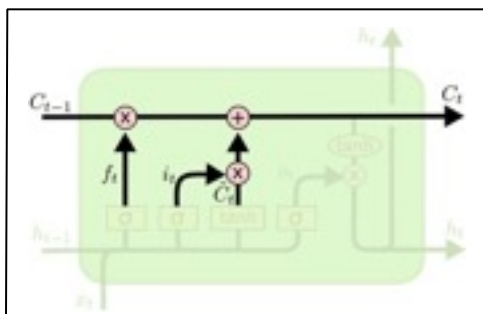
In the third step we will combine these two  
to create an update to the state

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

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

## LSTM walk-through: Third Step

- It's now time to update old cell state  $C_{t-1}$  into new cell state  $C_t$ 
  - The previous step decided what we need to do
  - We just need to do it



We multiply the old state by  $f_t$ , forgetting the things we decided to forget earlier.

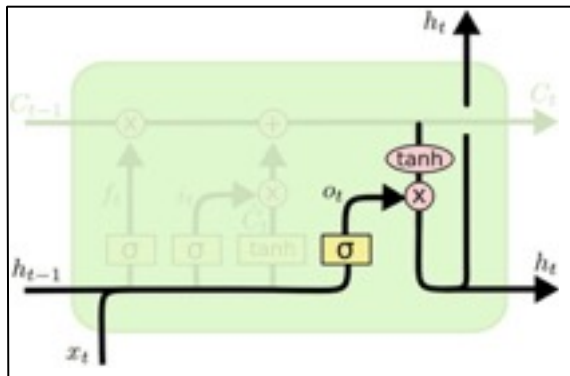
Then we add  
This is the new candidate values, scaled by  
how much we decided to update each  
state value

In the Language model,  
this is where we'd actually  
drop the information about the  
old subject's gender and add  
the new information,  
as we decided in previous steps

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

## LSTM walk-through: Fourth Step

- Finally we decide what we are going to output



This output will be based on our cell state, but will be a filtered version.

First we run a sigmoid layer which decides what parts of cell state we're going to output. Then we put the cell state through tanh (to push values to be between -1 and 1) and multiply it by the output of the sigmoid gate, so that we output only the parts we decided to

For the Language model,

Since it just saw a subject it might want to output

relevant to a verb, in case that

is what is coming next, e.g., it might output whether the subject is singular or plural so that we know what form a verb should be conjugated into if that's what follows next.

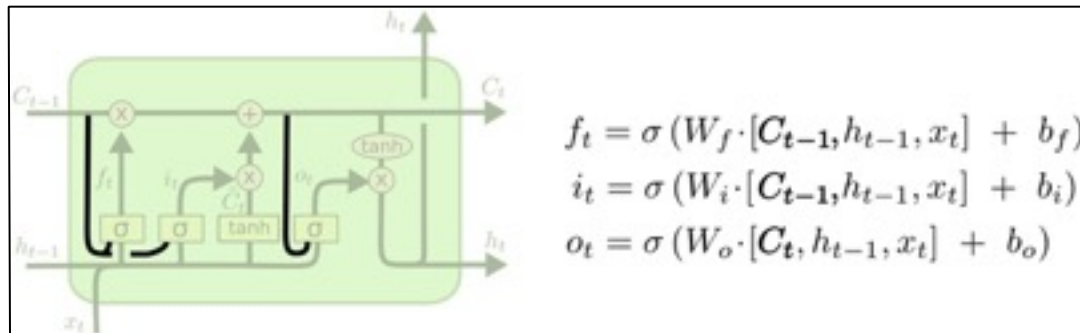
$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

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

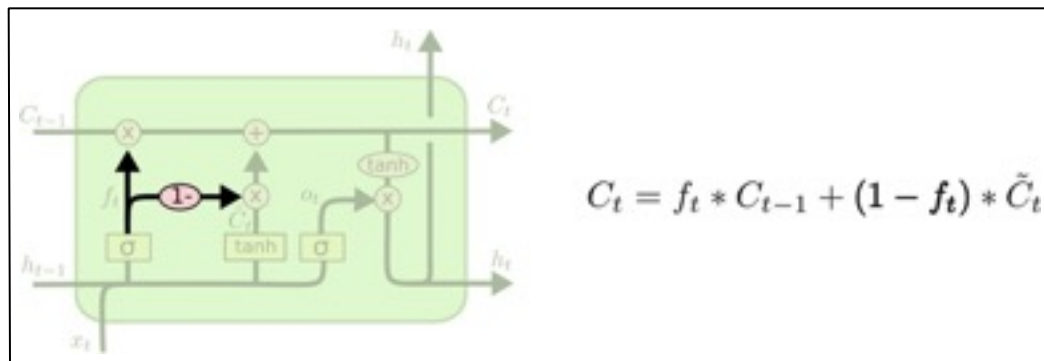


## Variants of LSTM

- There are several minor variants of LSTM
  - LSTM with “peephole” connections



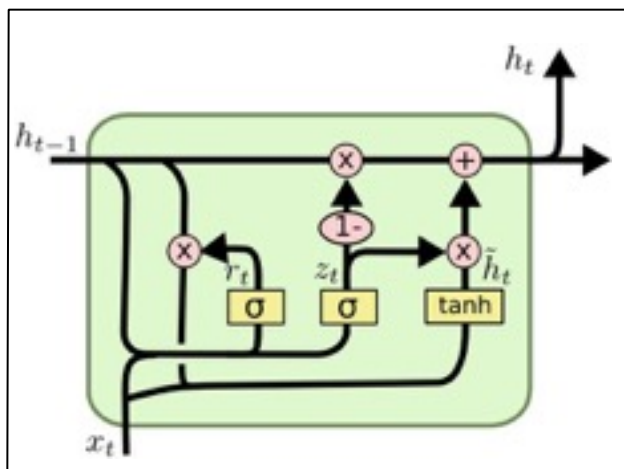
- Coupled forget and input gates



## Gated Recurrent Unit (GRU)

- A dramatic variant of LSTM

- It combines the forget and input gates into a single update gate
- It also merges the cell state and hidden state, and makes some other changes
- The resulting model is simpler than LSTM models
- Has become increasingly popular



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

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

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

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