

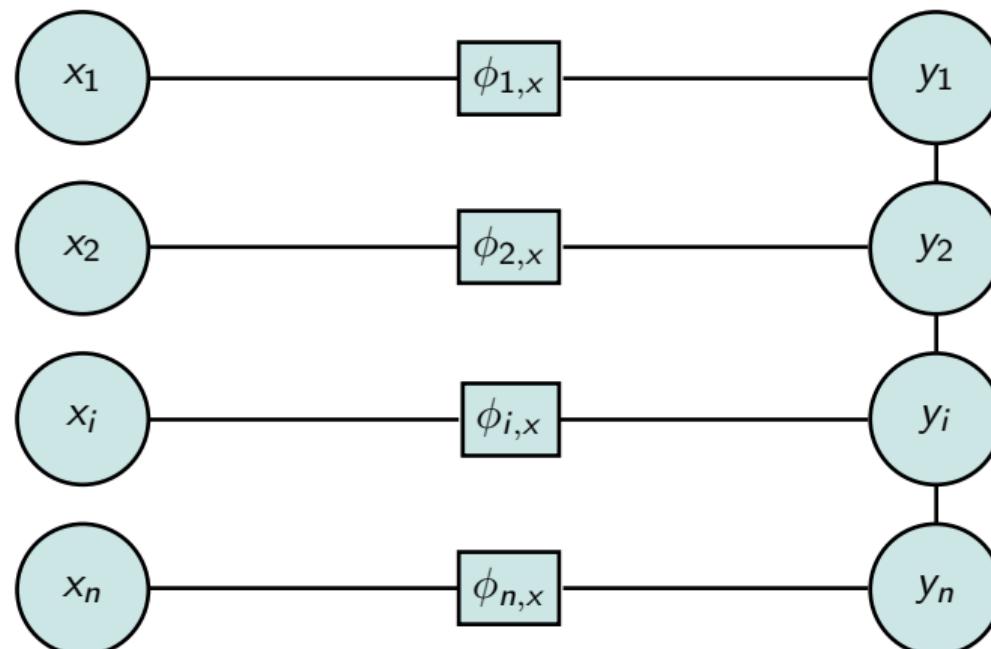
# Lecture 25a: Lego World of Deep Learning, Recurrent Neural Networks, LSTMs, CTC, etc

Instructor: Prof. Ganesh Ramakrishnan

# Recap: Linear Conditional Random Fields (CRF)

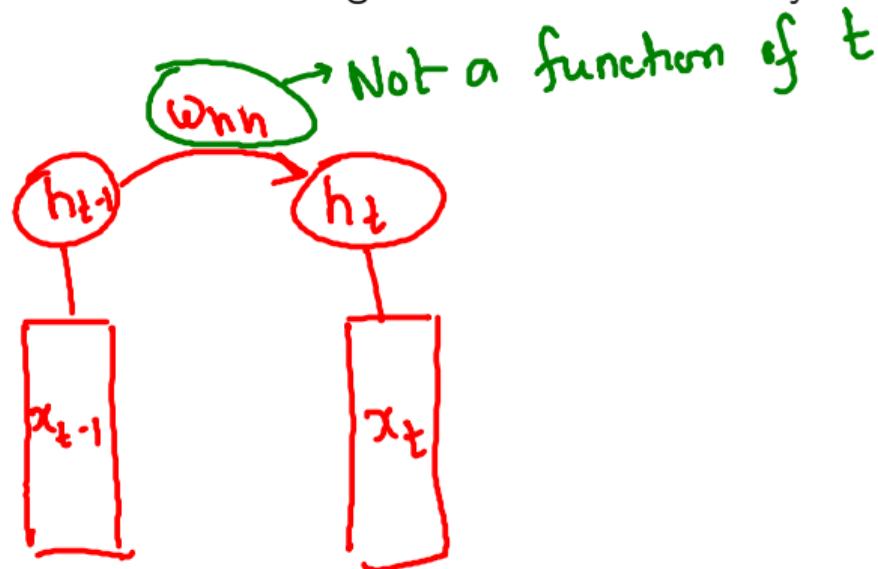
Just as CRF was an extension of Logistic Regression (LR) can Neural Networks (cascade of LRs) be extended to sequential output?

inputs                     $x$ -potentials      classes &  $y$ -potentials  $\phi_{i,y}$



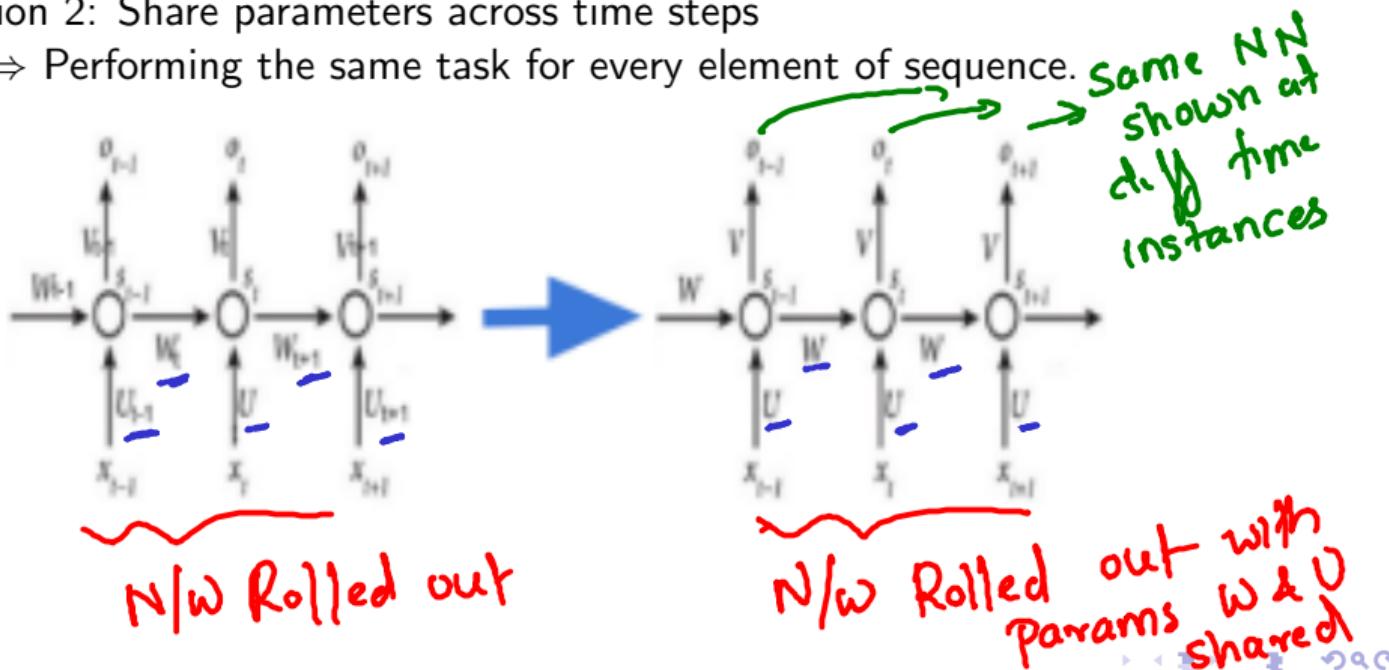
# Recurrent Neural Network (RNN) Intuition

- Recall: In CNN we used the trick of common parameters for many neurons
- RNN intuition 1: We want a neuron's output at time  $t$  to depend on its state  $s$  at time  $t - 1$
- RNN intuition 2: Share parameters across time steps
- Recurrent  $\Rightarrow$  Performing the same task for every element of sequence.



# Recurrent Neural Network (RNN) Intuition

- Recall: In CNN we used the trick of common parameters for many neurons
- RNN intuition 1: We want a neuron's output at time  $t$  to depend on its state  $s$  at time  $t - 1$
- RNN intuition 2: Share parameters across time steps
- Recurrent  $\Rightarrow$  Performing the same task for every element of sequence.

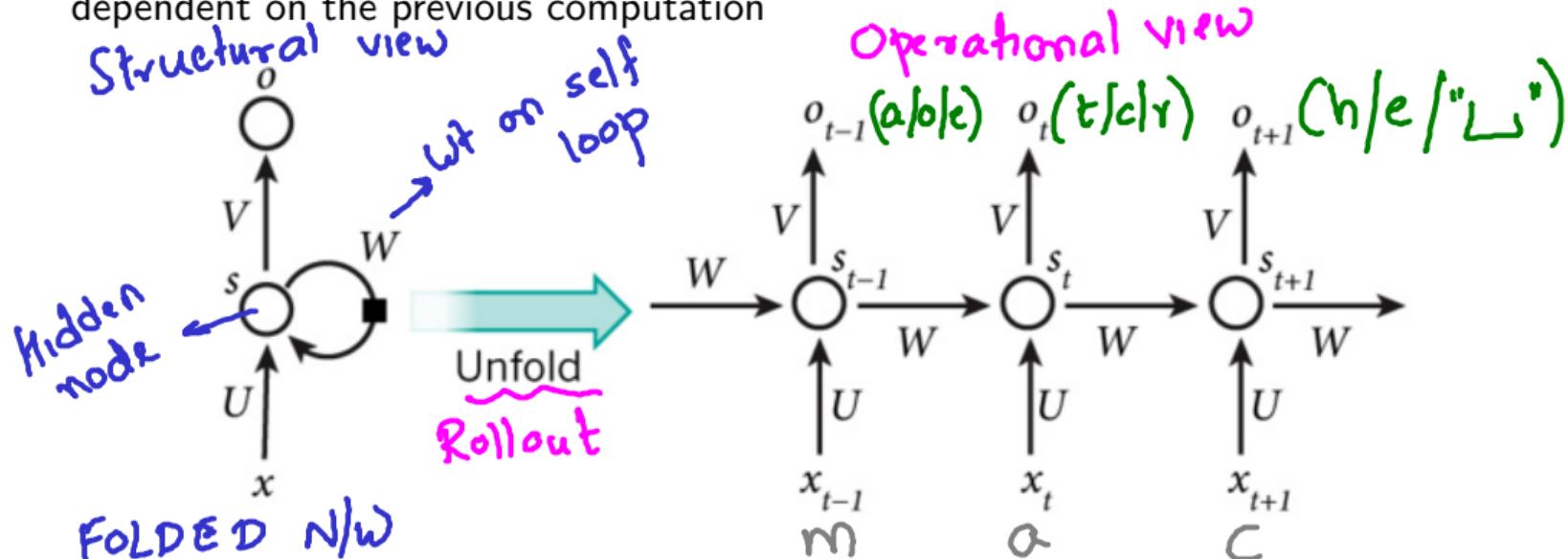


## Tutorial 8: Problem 7

Try the text generation RNN (Recurrent Neural Network) demo at  
<http://www.cs.toronto.edu/~ilya/rnn.html>. State any interesting observations.  
How would you improve the performance of the RNN?

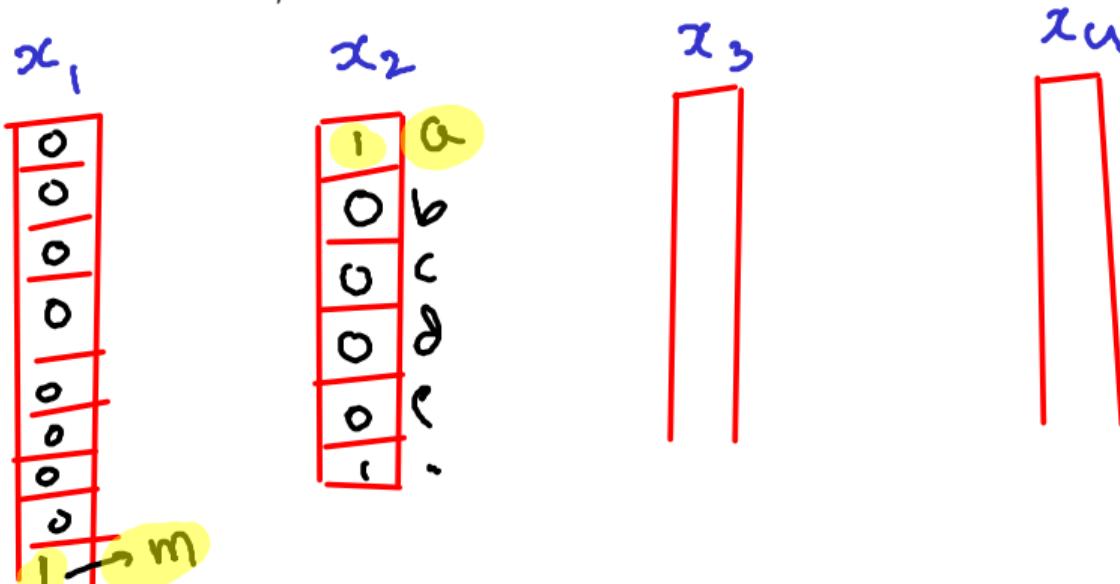
# RNN: Compact representation

- Generalization of Neural networks to Sequential tasks such as *language modeling*, *word prediction*, etc..
- Perform the same *task* for every element of the sequence, with the output being dependent on the previous computation



# RNN: One Hot Encoding for Language Model

- With 3 characters in vocabulary,  $a, b$  and  $c$ , what would be the best encoding to inform each character occurrence to the network?
- One Hot Encoding: Give a unique key  $k$  to each character in alpha-numeric order, and encode each character with a vector of vocabulary size, with a 1 for the  $k^{th}$  element, and 0 for all other elements.

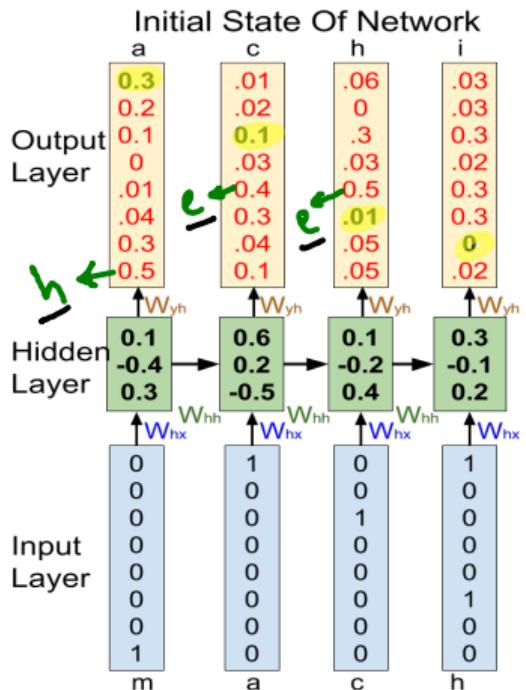


## RNN: One Hot Encoding for Language Model

- With 3 characters in vocabulary,  $a, b$  and  $c$ , what would be the best encoding to inform each character occurrence to the network?
- One Hot Encoding: Give a unique key  $k$  to each character in alpha-numeric order, and encode each character with a vector of vocabulary size, with a 1 for the  $k^{th}$  element, and 0 for all other elements.

a	b	c
1	0	0
0	1	0
0	0	1

# RNN: Language Model Example with one hidden layer of 3 neurons



} Model gets a chance to correct its wts through back propagation through time (BPTT)  
⇒ Correcting the green characters to the highlighted modes

# RNN: Language Model Example with one hidden layer of 3 neurons

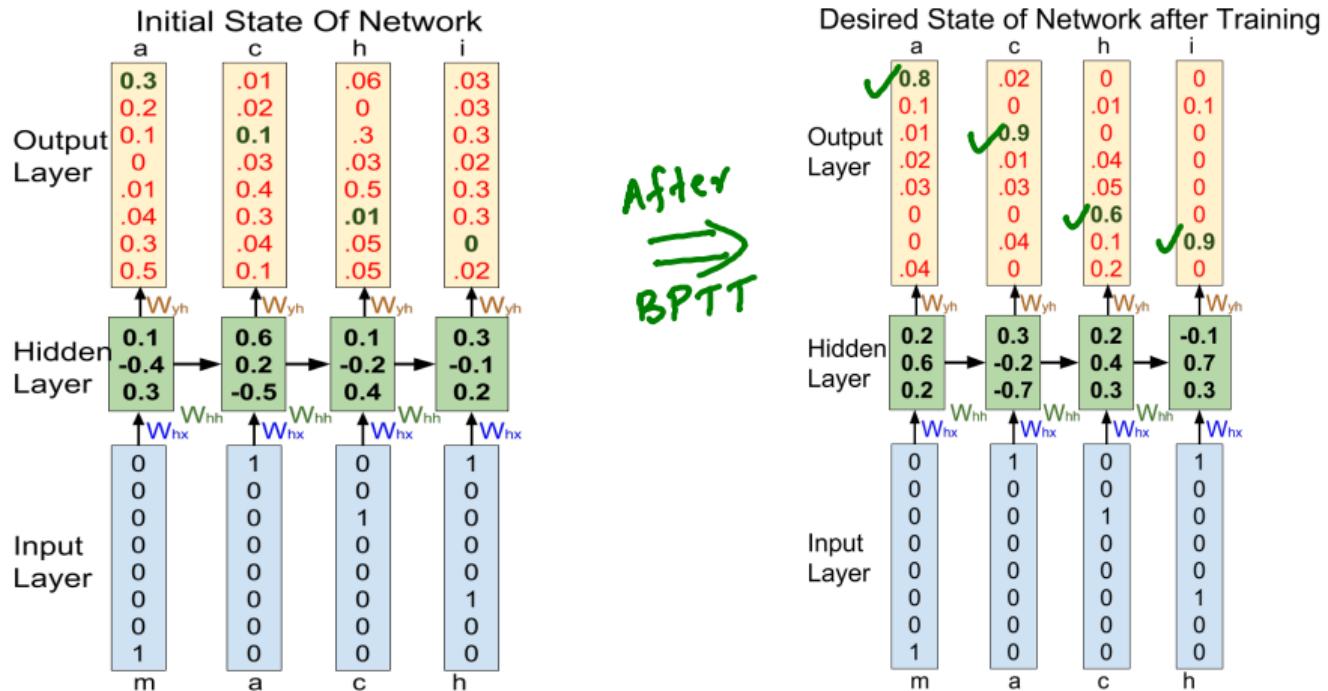


Figure 1: Unfolded RNN for 4 time units

# RNN: Equations (Used for speech to text by youtube)

- $\underline{h_t} = g(W_{hh}h_{t-1} + W_{hx}x_t + b_h)$ 
  - Activation function  $g$  could be sigmoid  $\sigma$  or its extension to multiclass called softmax<sup>1</sup>,  $\tanh\left(\frac{1-e^{-2x}}{1+e^{-2x}}\right)$  which is simply a scaled<sup>2</sup> and shifted version of the sigmoid function
  - A network may have combination of different activation functions<sup>3</sup>
- $\underline{y_t} = \underline{W_{yh} h_t}$
- The new (present) hidden state depends upon the previous hidden state(s) and the present input.
- The present output depends upon present hidden state (and in turn upon previous hidden states).

---

<sup>1</sup>Tutorial 7

<sup>2</sup> $\tanh(x) = 2\sigma(2x) - 1$

<sup>3</sup><http://www.wildml.com/2015/10/>

## Back Propagation Through Time: BPTT Illustration

- $h_1 = g(W_{hh}h_0 + W_{hx}x_0 + b_h)$ , initialize  $h_0$  and  $x_0$  as zero vectors.
- At  $t = 2$  we present  $x_2$  as 'a' at input and desire  $y_2$  as 'c' at output in one hot encoded form as shown previously
- $h_2 = g(W_{hh}h_1 + W_{hx}x_1 + b_h)$
- At  $t = 3$ ,  $x_3 = 'c'$ ,  $y_3$  we desire is 'h'.
- $y_3 = W_{yh} g(W_{hh}h_2 + W_{hx}x_2 + b_h)$

# Back Propagation Through Time: BPTT Illustration

Forward prop.

- $h_1 = g(W_{hh}h_0 + W_{hx}x_0 + b_h)$ , initialize  $h_0$  and  $x_0$  as zero vectors.
- At  $t = 2$  we present  $x_2$  as 'a' at input and desire  $y_2$  as 'c' at output in one hot encoded form as shown previously
- $h_2 = g(W_{hh}h_1 + W_{hx}x_1 + b_h)$
- At  $t = 3$ ,  $x_3 = 'c'$ ,  $y_3$  we desire is 'h'.
- $y_3 = \underline{W_{yh}} g(W_{hh}h_2 + W_{hx}x_2 + b_h)$
- Put  $h_1$  and  $h_2$  in the last equation and then perform gradient descent on weights (i.e., **tune weights** through back propagation) to fit  $y_3$  first corresponding to vectors  $\underline{x_3, x_2}$  and  $\underline{x_1}$ .
- Similarly use  $h_1$  in equation for  $y_2$  and **tune weights** to fit  $y_2$  corresponding to vectors  $x_2$  and  $x_1$ .
- Then **tune weights** for  $y_1$ .

$y_3 \rightarrow h_2 \rightarrow y_2 \rightarrow \dots$

- In previous example, we used
  - ① **sequence length** of 4, (that is, the number of time steps to unroll the RNN)
  - ② **batch size** of 1, (that is, the number of input vectors presented at a single time step to the RNN)
- We can vary these parameters according to the task at hand.
- One hot encoding is the best suited encoding for such tasks while training the neural networks.

BPTT involves lot more multiplication than normal backprop.



- **Vanishing gradient:** When training networks with long range dependencies, the influence of the given input decays exponentially as it cycles around the network's recurrent connections.
- Gradient vanishes especially when we use sigmoid function and several gradient values  $v$  with  $|v| < 1$ , get multiplied during BPTT to give a zero.
- Instead, if we used an alternative function that gives value  $> 1$  as output, we will face the problem of 'exploding gradient'.

# RNN: Vanishing Gradient Problem

The sensitivity(derivative) of network w.r.t input( $x_t = 1$ ) decays exponentially with time, as shown in the unfolded (for 7 time steps) RNN below. Darker the shade, higher is the sensitivity w.r.t to  $x_1$ .

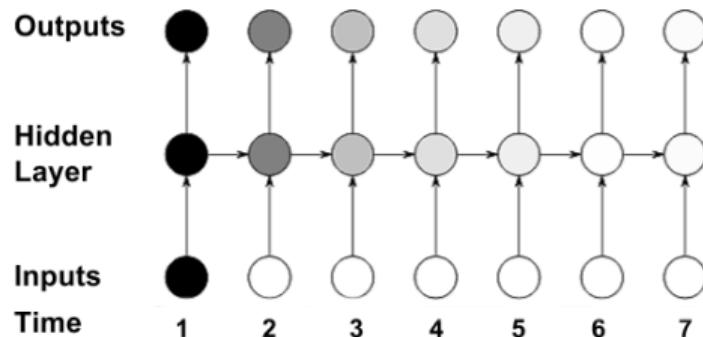


Image reference: Alex Graves 2012.

# Long Short-Term Memory (LSTM) Intuition

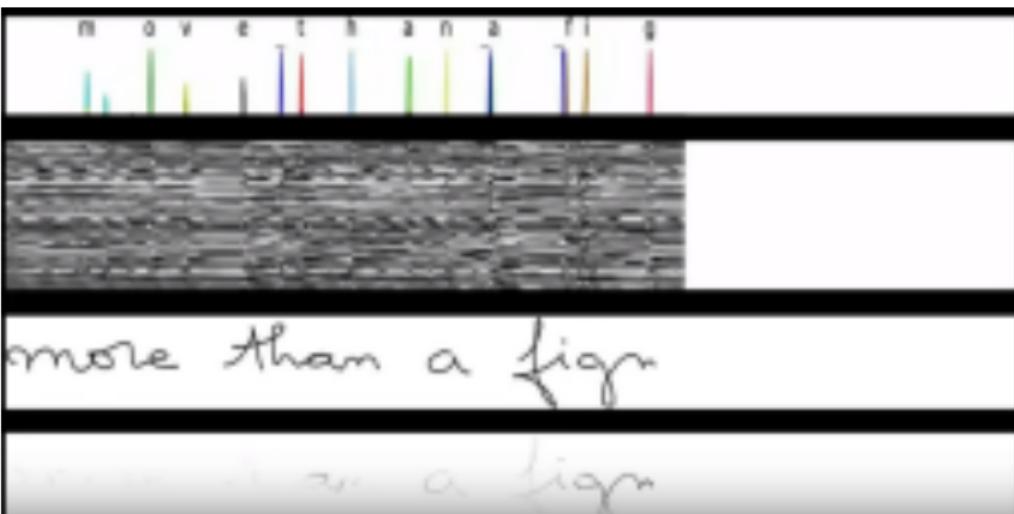


- Learn when to propagate gradients and when not, depending upon the sequences.
- Use the memory cells to store information and reveal it whenever needed.
- I live in **India**.... I visit **Mumbai** regularly.  
    
- For example: Remember the context "India", as it is generally related to many other things like language, region etc. and forget it when the words like "Hindi", "Mumbai" or End of Line/Paragraph appear or get predicted.

# Recap: The Lego Blocks in Modern Deep Learning

- ① Depth/Feature Map
- ② Patches/Kernels (provide for spatial interpolations) - **Filter**
- ③ Strides (enable downsampling)
- ④ Padding (shrinking across layers)
- ⑤ Pooling (More downsampling) - **Filter**
- ⑥ **RNN and LSTM** (Backpropagation through time and Memory cell)
- ⑦ **Connectionist Temporal Classification**
- ⑧ Embeddings (Later, with unsupervised learning)

# Demonstration of Alex Graves's system working on pen coordinates



- ① Top: Characters as recognized, without output delayed but never revised.
- ② Second: States in a subset of the memory cells, that get reset when character recognized.
- ③ Third: Actual writing (input is x and y coordinates of pen-tip and up/down location).
- ④ Fourth: Gradient backpropagated all the way to the xy locations. Notice which

# LSTM Equations

- $f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f)$
- $i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i)$
- We learn the forgetting ( $f_t$ ) of previous cell state and insertion ( $i_t$ ) of present input depending on the present input, previous cell state(s) and hidden state(s).

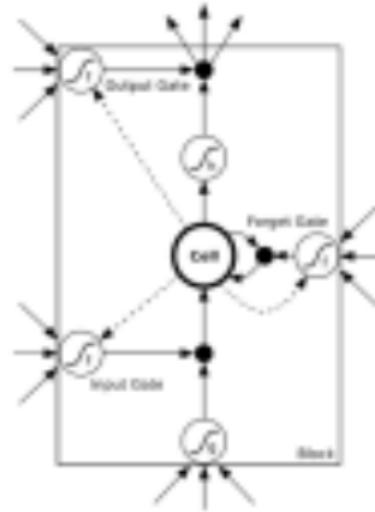


Image reference: Alex Graves 2012

# LSTM Equations

- $c_t = f_t c_{t-1} + i_t \tanh(W_{hc} h_{t-1} + W_{xc} x_t + b_c)$
- The new cell state  $c_t$  is decided according to the firing of  $f_t$  and  $i_t$ .

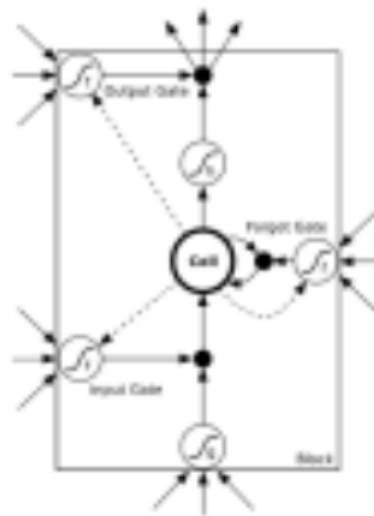


Image reference: Alex Graves 2012.

# LSTM Equations

- $c_t = f_t c_{t-1} + i_t \tanh(W_{hc} h_{t-1} + W_{xc} x_t + b_c)$
- $f_t = \sigma(W_{xf} x_t + W_{hf} h_{t-1} + W_{cf} c_{t-1} + b_f)$
- $i_t = \sigma(W_{xi} x_t + W_{hi} h_{t-1} + W_{ci} c_{t-1} + b_i)$
- Each gate is a vector of cells; keep the constraint of  $W_{c*}$  being diagonal so that each element of LSTM unit acts independently.
- $o_t = \sigma(W_{xo} x_t + W_{ho} h_{t-1} + W_{co} c_{t-1} + b_f)$
- $h_t = o_t \tanh(c_t)$

# LSTM Gradient Information remain preserved

The opening 'O' or closing '-' of input, forget and output gates are shown below, to the left and above the hidden layer respectively.

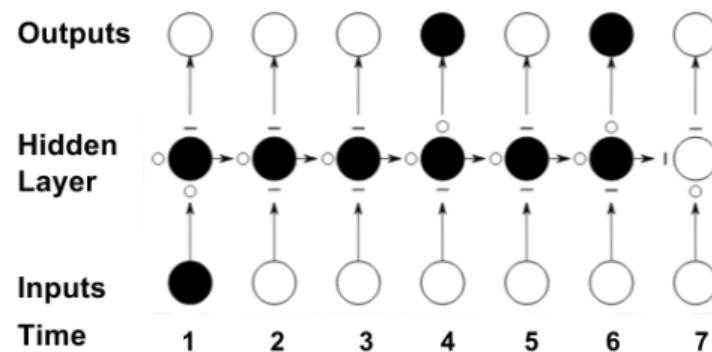


Image reference: Alex Graves 2012.

# LSTM V/S RNN results on Novel writing

A RNN and a LSTM, when trained appropriately with a Shakespeare Novel write the following output (for few time steps) upon random initialization.

RNN's Novel	LSTM's Novel
<p>kentages, whom falled me, Compassion, that the wind of England <u>topased</u> To grief all go.</p> <p>AUPIDOUS: Marry, it were any graven, And high againes.</p> <p>PAULINA: Natus my <u>plateries</u>, this are the sity, hot continue stilen <u>sensable</u> Tybalt.</p> <p>GREMIO: God sues we <u>charnet</u>, Since this is shoulders to the king.</p> <p>KING EDWARD IV Look, myself, Sir Jupe:</p>	<p>Monsur and <u>Parrage</u>, is no newe <u>Desraud</u> is a death barshers: Her much.</p> <p>LEONTES:</p> <p>JULIET: O <u>Wazewat</u> that she were a pitifly, that nothing.</p> <p>Second Lord: They are of my feul me: power and your <u>gallibg</u>.</p> <p>Lore and the <u>glatisfal</u> indeed, stamp; hell, if <u>procing</u> dages.</p> <p>LADY GREY: No moter them <u>belisew</u> the barges That Emrillth these out with her good</p>
...	...

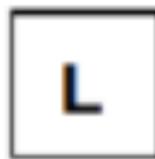
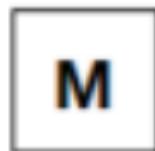
# Sequence Labeling

- The task of labeling the sequence with discrete labels. Examples: Speech recognition, handwriting recognition, part of speech tagging.
- Humans while reading/hearing make use of context much more than individual components. For example:- Yoa can undenstand dis, tough itz an eroneous text.
- The sound or image of individual characters may appear similar and may cause confusion to the network, if the proper context is unknown. For example: "in" and "m" may look similar whereas "dis" and "this" may sound similar.

## Type of Sequence Labeling Tasks

- Sequence Classification: Label sequence is constrained to be of unit length.

Input Images



Output

M

L

## Type of Sequence Labeling Tasks

- Segment Classification: Target sequence consist of multiple labels and the segment locations of the input is known in advance, e.g. the timing where each character ends and another character starts is known in a speech signal.

Input Image Sequence	Output text sequence
Segmented sentence!	Segmented sentence!

- We generally do not have such data available, and segmenting such data is both tiresome and erroneous.

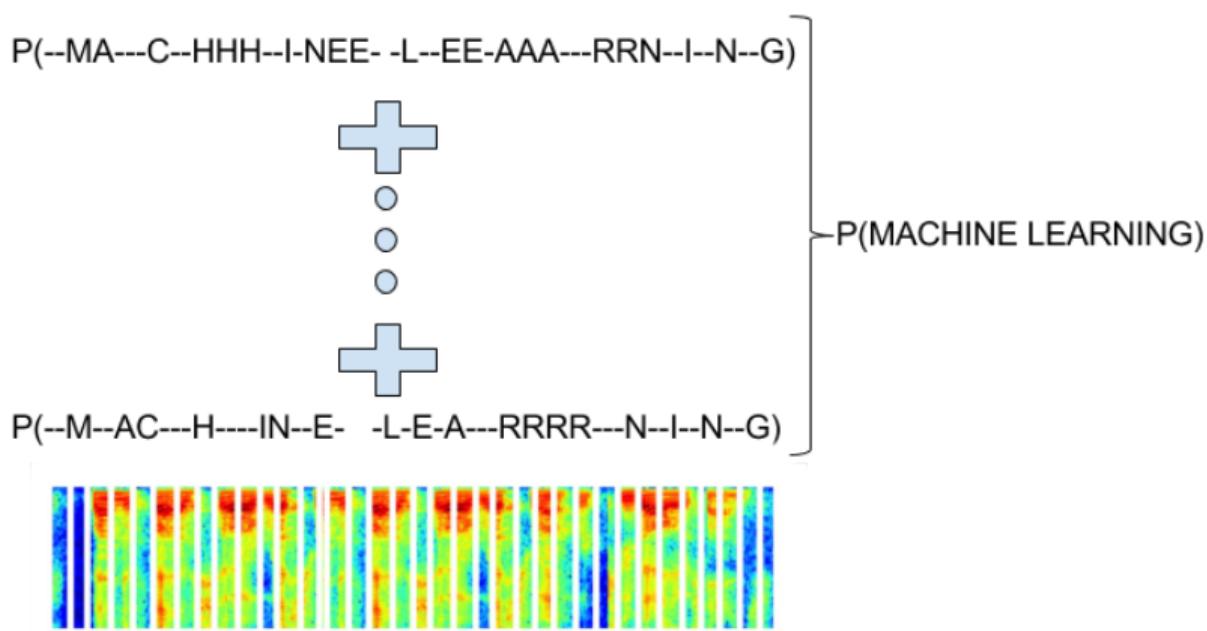
## Type of Sequence Labeling Tasks

- Temporal Classification: Tasks in which temporal location of each label in the input image/signal does not matter.
  - Very useful, as generally we have higher level labeling available for training, e.g. word images and the corresponding strings, or it is much easier to automate the process of segmenting the word images from a line, than to segment the character images from a word.

# Connectionist Temporal Classification (CTC) Layer

- For temporal classification task: length of label sequence  $\neq$  length of input sequence.
- CTC label predictions at any time in input sequence.
- Predict an output at every time instance, and then decode the output using probabilities we get at output layer in vector form.
- e.g. If we get output as sequence "-m-aa-ccch-i-nee- -lle-a-rr-n-iinnn-g", we will decode it to "machine learning".
- While training we may encode " machine learning" to"-m-a-c-h-i-n-e- -l-e-a-r-n-i-n-g-" via C.T.C. Layer.

# CTC Intuition [Optional]



## CTC Intuition [Optional]

NN function :  $f(x_T) = y_T$

- $x_T$ : input image/signal  $x$  of length  $T$ .
  - For image: each element of  $x_T$  is a column(or its feature) of the image.
- $y_T$ : output sequence  $y$  of length  $T$ .
  - each element of  $y_T$  is a vector of length  $|A'|$  (where  $A' = A \cup \{"\text{-}\"$  i.e. alphabet set  $\cup$  blank label).

$\ell_U$  : Label of length  $U(jT)$ .

Intuition behind CTC:

- generate a PDF at every time-step  $t \in 1, 2, \dots, T$ .
- Train NN with objective function that forces Max. Likelihood to decode  $x_T$  to  $\ell_U$ (desired label).

## CTC Layer: PDF [Optional]

$$P(\pi \rightarrow x) = \prod_{t=1}^T y_t(\pi_t)$$

- path  $\pi$  : a possible string sequence of length T, that we expect to lead to  $\ell$ . For example: -p-a-t-h-, if  $\ell$  = "path".
- $y_i(n)$ : probability assigned by NN when character  $n (\in A')$  is seen at time i. "-" is symbol for blank label.
- $\pi_t$  :  $t^{th}$  element of path  $\pi$ .

$$P(\ell \rightarrow x) = \sum_{label(\pi)=\ell} P(\pi \rightarrow x)$$

## CTC Layer: PDF [Optional]

$$P(\ell|x) = \sum_{label(\pi)=\ell} P(\pi|x) = \sum_{label(\pi)=\ell} \prod_{t=1}^T y_t(\pi_t)$$

- Question: What could be possible paths of length  $T = 9$  that lead to  $\ell = "path"$ ?
- Answer:
- Question: How do we take care of cases like  $\ell = "Mongoose"$ ?
- Answer:

## CTC Layer: PDF [Optional]

$$P(\ell \rightarrow x) = \sum_{label(\pi)=\ell} P(\pi \rightarrow x) = \sum_{label(\pi)=\ell} \prod_{t=1}^T y_t(\pi_t)$$

- Question: What could be possible paths of length  $T = 9$  that lead to  $\ell = "path"$ ?
- Answer: -p-a-t-h-", pp-a-t-h-", -paa-t-h-", -ppa-t-h-", -p-aat-h- etc.
- Question: How do we take care of cases like  $\ell = "Mongoose"$ ?
- Answer: We change  $\ell = "Mongoose"$  to  $\ell = "Mongo-ose"$ .
- Question: During training  $\ell$  is known, what to do at testing stage?
- Answer: Next Slide.

## CTC Layer: Forward Pass Decoding [Optional]

$$P(\ell \rightarrow x) = \sum_{label(\pi)=\ell} P(\pi \rightarrow x) = \sum_{label(\pi)=\ell} \prod_{t=1}^T y_t(\pi_t)$$

Question: During training  $\ell$  is known, what to do at testing stage?

- ① Brute force: try all possible  $\ell$ s, all possible  $\pi$ s for each  $\ell$  to get  $P(\ell \rightarrow x)$  and choose best  $\ell$ .
  - Rejected as expensive.
- ② Best Path Decoding - most likely path corresponds to the most likely label.
  - $P(A1) = 0.1$ , where A1 is the only path corresponding to label A.
  - $P(B1) = P(B2) = \dots = P(B10) = 0.05$ , where B1..B10 are the 10 paths corresponding to label B.
  - Clearly B is preferable over A as  $P(B \rightarrow x) = 0.5$ .
  - But Best Path Decoding will select A.
  - Rejected as inaccurate.
- ③ Prefix Search Decoding - NEXT

## CTC Layer: Prefix Search Decoding [Optional]

- ① Initialize prefix  $p^* = \phi$ ; //  $p^*$  is prefix of /ell
- ②  $P(\phi) = 1$ ; // as  $\phi$  is prefix of every word.
- ③ try  $p_{new} = p^* + k$ , for all  $k \in A \cup \{eos\}$ ; // + represent concatenation.
- ④ Maintain  $L_p$ : list of growing prefixes; // —A— + 1 new values per iteration.
- ⑤ Maintain  $P_p$ : list of probabilities of corresponding elements in  $L_p$ ; // How to find  $P(p^*+k)$ ? Next Slide.
- ⑥ if  $P(p^*+eos) \geq \max(P_p)$ : stop and go to step 8;
- ⑦ else: update  $p^*$  with prefix having max prob. and repeat from step 3;
- ⑧  $p^*$  is the required prefix.

In practice beam-search is used to limit the exponentially growing  $L_p$  and make the decoding faster.

# CTC Layer: Prefix Search Decoding

Consider the DAG shown below with  $A = \{x,y\}$  and  $e$  representing the end of string.  
The steps a-f represent the path followed by Prefix Search Decoding Algorithm.

- What  $\ell$  would the Best Path Decoding Produce?
- What  $\ell$  would the Prefix Search Decoding Produce?

