## W1 - L1

$X$:   $\overset{<1>}{\text{Harry}}$ $\overset{<2>}{\text{Potter}}$ and $\overset{}{\text{Hermione}}$ $\overset{}{\text{Granger}}$ invented $\overset{<7>}{\text{a}}$ new $\overset{<9>}{\text{spell}}$

$\overset{<t>}{X}$

$y$:    1    1    0    1    1     0   0   0   0

    ↓ Name      ↓ NotName

$X^{(i)<t>}$: $t^{th}$ word of $i^{th}$ example

$T_X^{(i)}$: input sequence length of $i^{th}$ training example

$y^{(i)<t>}$: $t^{th}$ element of $i^{th}$ output seq

$T_y^{(i)}$: output sequence length of $i^{th}$ training example

$x^{<t>}$ = feature for $t^{th}$ word

$T_X = 9$ (length of input)

$T_y = 9$ ( " " output sequence)

## Representing Words

Vocabulary

$$\begin{bmatrix} a \\ aaron \\ and \\ harry \\ potter \\ zulu \end{bmatrix} \begin{matrix} 1 \\ 2 \\ \vdots \\ 367 \\ \vdots \\ 4075 \\ 6830 \\ \vdots \\ 10,000 \end{matrix}$$

$X^{<1>} =$ HARRY $(R^{10,000}) \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow 4075$

$X^{<7>} =$ $(R^{10,000}) \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$

## Why not standard NN for NER?

$$\overline{\text{Input}} = \overline{T_x} \times 10,000 \begin{cases} x^{<1>} \quad O \rightarrow \\ x^{<2>} \quad O \rightarrow \\ \vdots \\ x^{<T_x>} \quad O \rightarrow \end{cases}$$

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \rightarrow \begin{array}{l} O \quad y^{<1>} \\ O \quad y^{<2>} \\ \vdots \\ O \\ O \quad y^{<T_y>} \end{array}$$

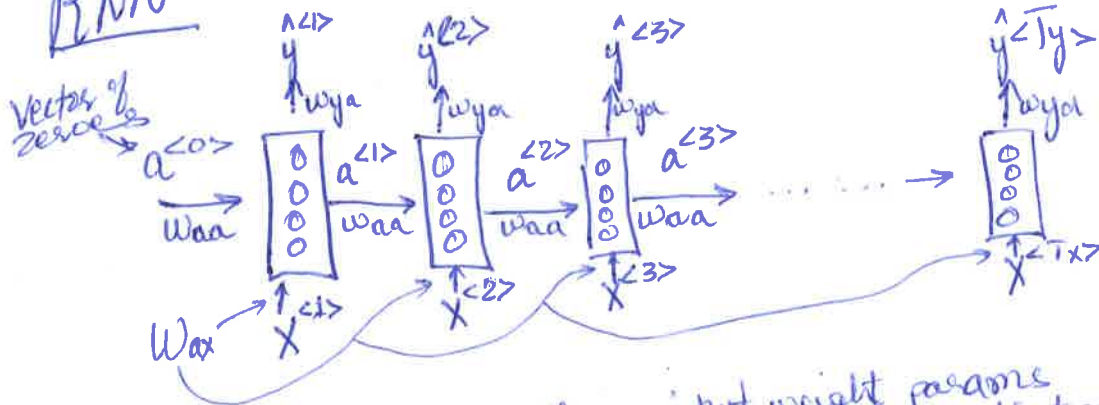### Problems

- Inputs, outputs can be different lengths in different examples
- Doesn't share features learned across different positions of text.

## RNN



Vector of zeros

$a^{<0>}$

$W_{aa}$

$W_{ax}$

Here $\overline{T_x = T_y}$

$W_{ax}$: input weight params
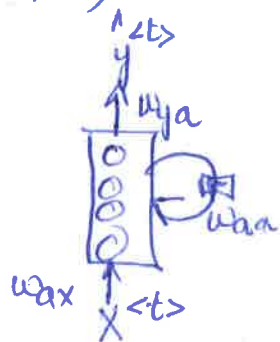$W_{aa}$: params for activations
$W_{ya}$: output params

  ☆ In RNN, parameters are shared throughout the network.

  ☆ One disadvantage of RNN is that it does not take into account future inputs to make predictions on output at "t". Therefore, we will see Bidirectional RNN (BRNN) that does this job.
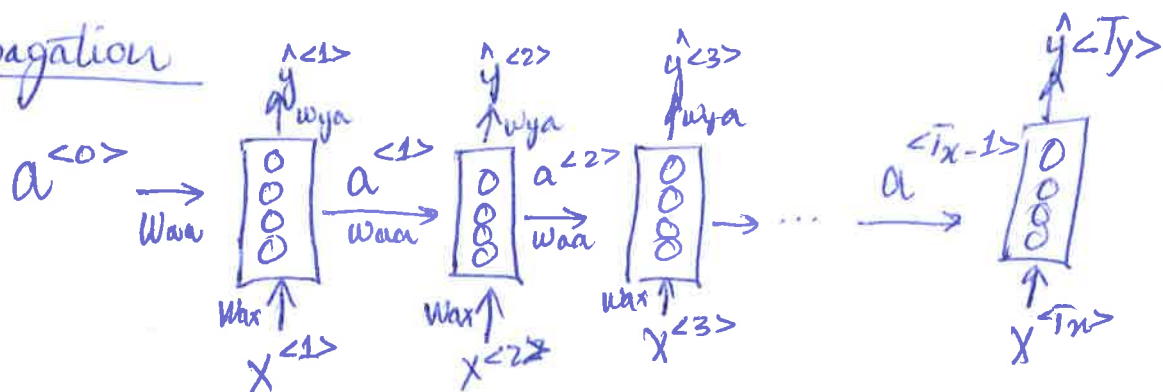
For instance :  He said, "Teddy Roosevelt was a great President" we can predict it as name if we take into account future words.

In some books, RNNs are shown like this:



# Forward Propagation



$$a^{<0>} = \vec{0}$$

$$a^{<1>} = g_1\left(W_{aa}\, a^{<0>} + W_{ax}\, x^{<1>} + b_a\right)$$

$g_1 :$ $\underline{tanh}$/Relu Activation

Common in ANN

$$\hat{y}^{<1>} = g_2\left(W_{ya}\, a^{<1>} + b_y\right)$$

$g_2 :$ Sigmoid activation (binary task)

OR

$g_2 :$ softmax (for K class classification)

$$\Rightarrow a^{<t>} = g\left(W_{aa}\, a^{<t-1>} + W_{ax}\, x^{<t>} + b_a\right)$$

$$\hat{y}^{<t>} = g\left(W_{ya}\, a^{<t>} + b_y\right)$$

$$W_{ax}$$

2nd notation means it will be multiplied by quantity like "x"

1st notation means it will calculate quantity like "a".

$$a^{<t>} = g\left(\underbrace{w_{aa}\, a^{<t-1>} + w_{ax}\, x^{<t>}} + ba\right)$$

$$\hat{y}^{<t>} = g\left(w_{ya}\, a^{<t>} + by\right) \quad \text{simplify this bit using stacking}$$

## Simplifying

$$a^{<t>} = g\left(w_a\left[a^{<t-1>},\, x^{(t)}\right] + ba\right)$$

In our running example:

$$x^{<t>} = R^{10,000} \quad \text{& lets say } a^{<t-1>} = R^{100}$$

$$\Rightarrow \quad a^{<t>} = 100 \times )$$

$$w_{aa} = 100 \times 100$$

$$w_{ax} = 100 \times 10,000$$

And $w_a = \left[\begin{array}{c|c} w_{aa} & w_{ax} \end{array}\right] = \left(100 \times 10,100\right)$

(with dimensions $100$, $100$, $10,000$)

(Stacking $w_{aa}, w_{ax}$)

And $\left[a^{<t-1>},\, x^{<t>}\right] = \left[\begin{array}{c} a^{<t-1>} \rightarrow R^{100} \\ \\ x^{<t>} \rightarrow R^{1000} \end{array}\right]$

Stacking $a^{<t-1>}$ & $x^{<t>}$

$$\Rightarrow \boxed{w_a\left[a^{<t-1>},\, x^{(t)}\right] = \left[w_{aa} \mid w_{ax}\right]\left[\begin{array}{c} a^{<t>} \\ x^{<t>} \end{array}\right] = \begin{array}{l} w_{aa}\, a^{<t-1>} \\ + w_{ax}\, x^{<t>} \end{array}}$$

The advantage of this stacking is that instead of two parameters ($w_{aa}$ & $w_{ax}$), we have just one parameter $w_a$ — which will help us form complex networks.
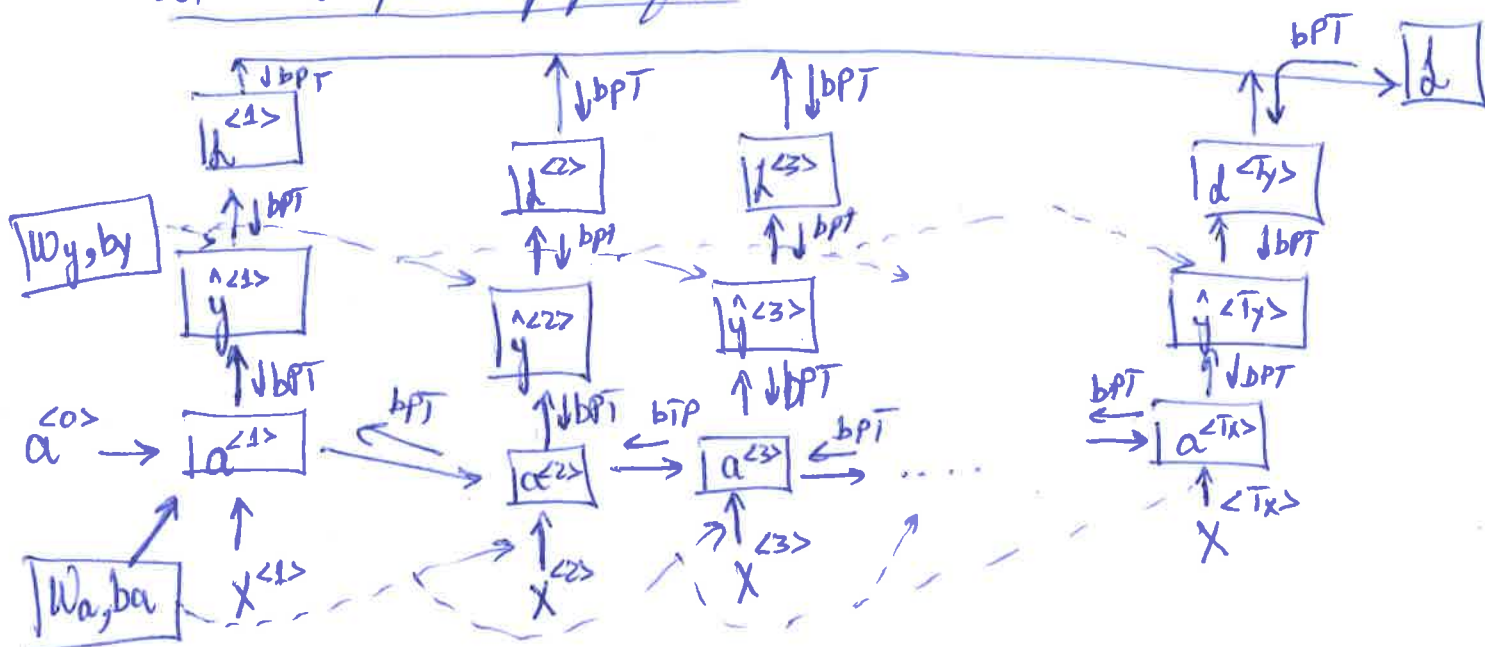
Thus $\hat{y}^{<t>} = g\left(w_y\, a^{<t>} + by\right)$

subscripts now indicates we are calculating "y"

# Backpropagation through time (BPT)

Lets see forwardprop of RNN:



## Loss function

In our example of NER, we will define loss for each prediction of word.

⇒ loss for $t^{th}$ word or $t^{th}$ token (lets define logistic loss)

$$\mathscr{L}^{<t>}\left(\hat{y}^{<t>}, y^{<t>}\right) = -y^{(t)}\log\hat{y}^{(t)} - (1-\hat{y}^{<t>})\log(1-\hat{y}^{<t>})$$

↑ Prediction   ↑ Ground truth

For entire sequence:

$$\mathscr{L}(\hat{y}, y) = \sum_{t=1}^{T_y} \mathscr{L}^{<t>}\left(\hat{y}^{<t>}, y^{<t>}\right)$$

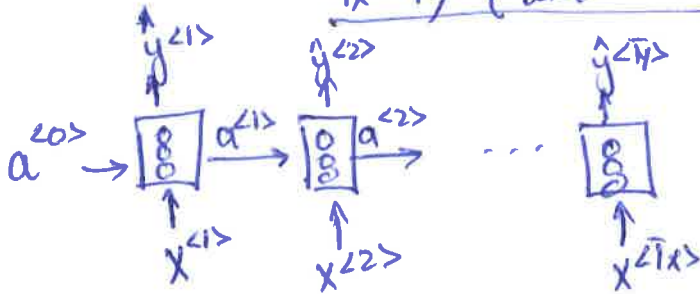* BPT will update $W_a, b_a, W_y, b_y$ by imputing partial derivatives.

# Different types of RNN's

Presentation inspired by blogpost
Andrej Karpathy: "The unreasonable effectiveness
of RNN"

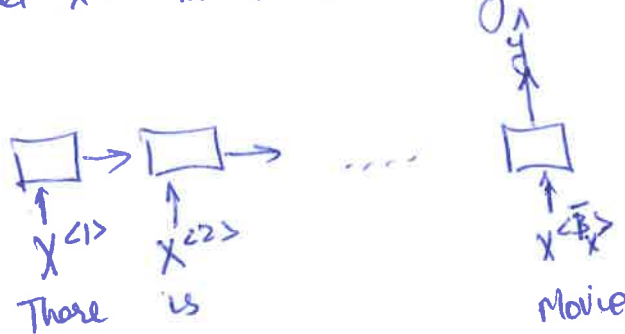## 1) MANY-to-MANY Architecture

$$T_x = T_y \text{ (like our NER)}$$



## 2) MANY-to-ONE Architecture

Consider Sentiment analysis

$X = text$

$y = 0/1$ or five star rating

Let $x =$ There is nothing to like in this movie



There     is                    Movie

## 3) ONE-to-ONE Archietecture



(Normal NN)
like in Course 1 & Course-2

# 4) One-to-Many Architecture

Example is Music generation

$$X \to y^{<1>} \ y^{<2>} \ \cdots \ y^{<T_y>}$$
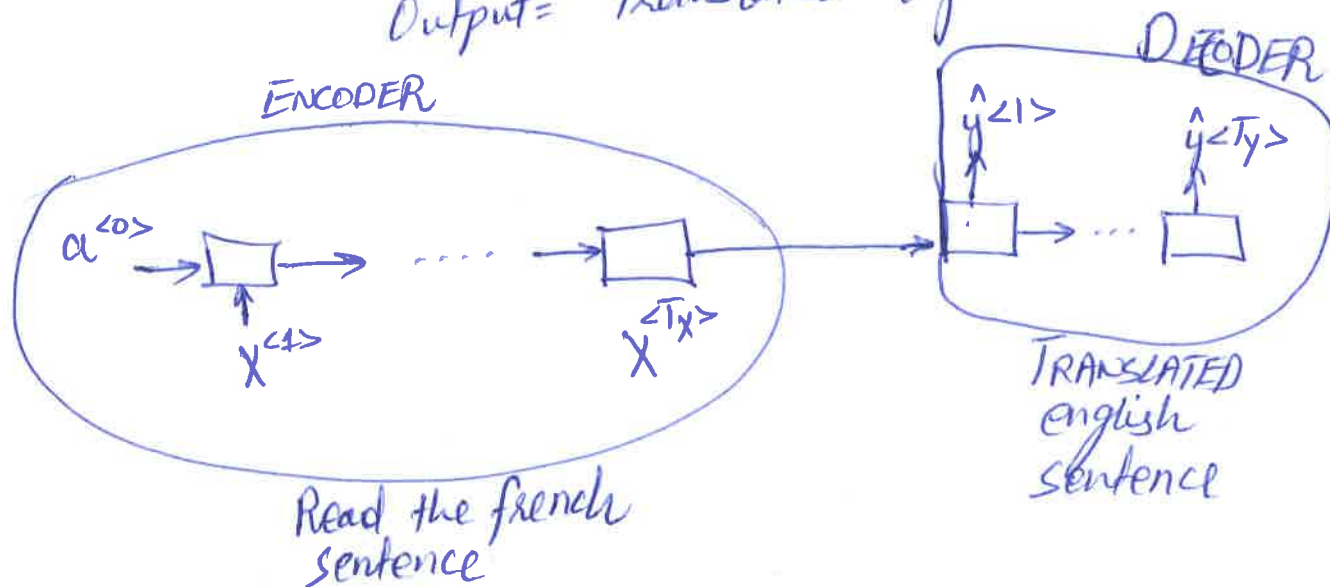


# 5) MANY-to-MANY Architecture

* $T_x \neq T_y$ like in Machine translation

Input = french sentence
Output = translated english sentence

ENCODER

DECODER



Read the french sentence

TRANSLATED english sentence

# LANGUAGE Model

It takes input a text sequence,
Lets say an output of speech recognition $\hat{y} = \hat{y}^{<1>}, \hat{y}^{<2>}, \ldots, \hat{y}^{<T_y>}$

$$\boxed{voice \longrightarrow text (\hat{y})} \begin{pmatrix} \text{Apples and pears are} \\ \text{delicious} \end{pmatrix}$$

then language model will output the probability
of $\hat{y}$ i.e $P(\hat{y})$

## (Learning a Language Model using RNN)
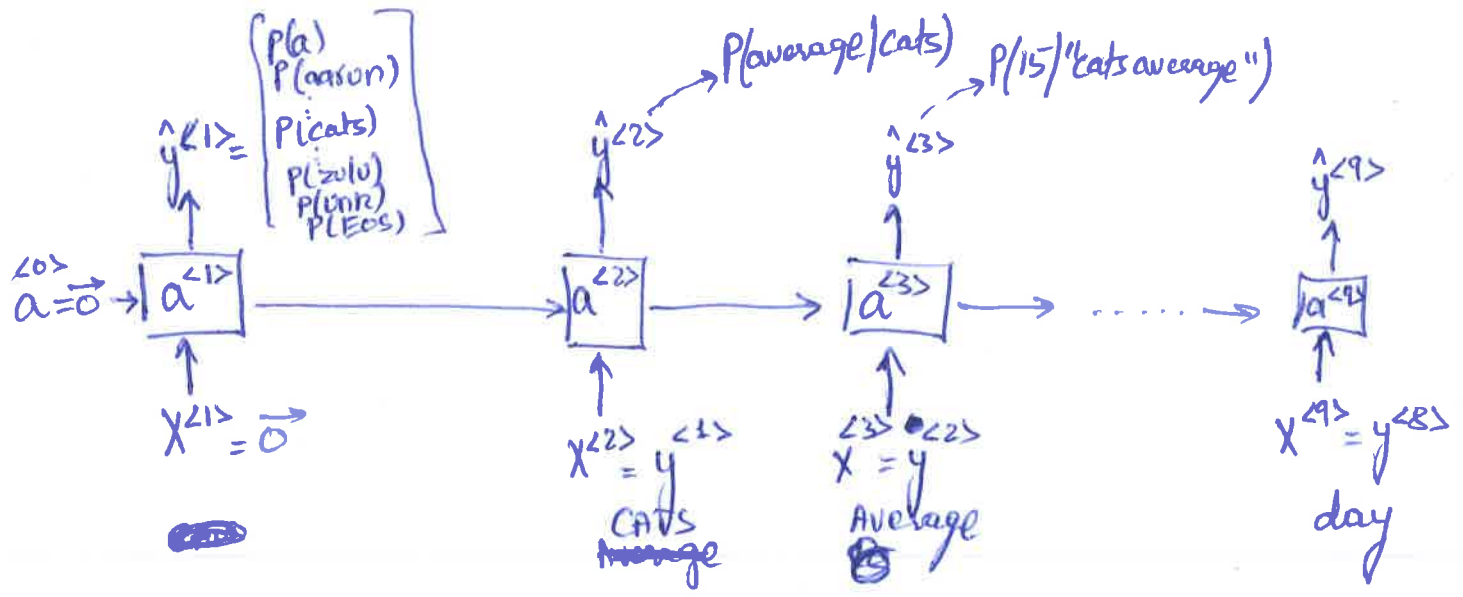
* Training Set: Large corpus of english text
* Running Examples: CATS average 15 hours of sleep a day. <EOS>
  $\quad\quad\quad y^{<1>} \quad\quad y^{<2>} \quad y^{<3>} \quad\quad\quad\quad\quad\quad y^{<8>} \quad y^{<9>}$

* Lets imagine our dictionary is:

$$\begin{bmatrix} a \\ aaron \\ cat \\ zulu \end{bmatrix} \updownarrow \; 10,000 \text{ words}$$

## RNN Model

Initialize with zero vector: $a^{<0>}$ & $x^{<1>}$

$$\hat{y}^{<1>} = \begin{bmatrix} P(a) \\ P(aaron) \\ P(cats) \\ P(zulu) \\ P(unk) \\ P(EOS) \end{bmatrix}$$

$P(average|cats) \quad P(15|"cats average")$

$\hat{y}^{<2>} \quad\quad\quad \hat{y}^{<3>} \quad\quad\quad\quad\quad\quad \hat{y}^{<9>}$

$a^{<0>} = \vec{0} \rightarrow \boxed{a^{<1>}} \longrightarrow \boxed{a^{<2>}} \longrightarrow \boxed{a^{<3>}} \longrightarrow \cdots \longrightarrow \boxed{a^{<4>}}$

$X^{<1>} = \vec{0}$

$X^{<2>} = y^{<1>}$  CATS

$X^{<3>} = y^{<2>}$  Average

$X^{<9>} = y^{<8>}$  day

Loss function (softmax):

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = -\sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

Total loss $L = \sum_t \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$

After training, if we are given a sentence, then we can calculate its probability ~~aftersing~~ using a language model.
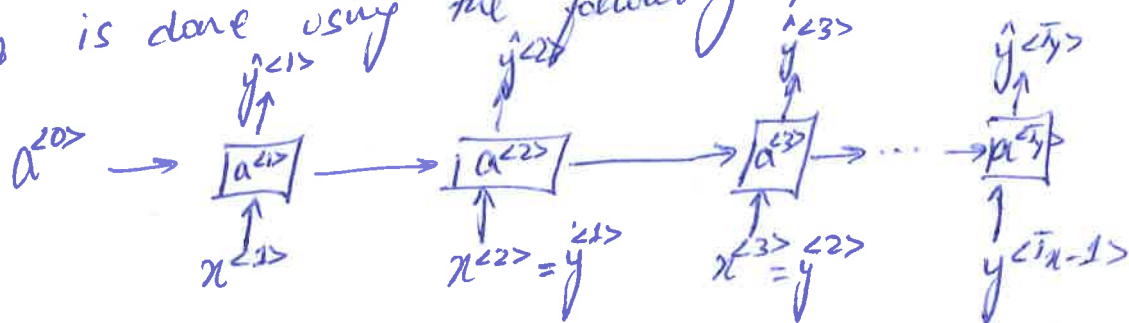
Sentence : $y^{<1>}, y^{<2>}, y^{<3>}$

$\Rightarrow P(y^{<1>}, y^{<2>}, y^{<3>}) = P(y^{<1>}) \, P(y^{<2>}|y^{<1>}) \, P(y^{<3>}|y^{<1>}, y^{<2>})$

The advantage of language model is that we can sample sequences.

# SAMPLING Novel Sequences
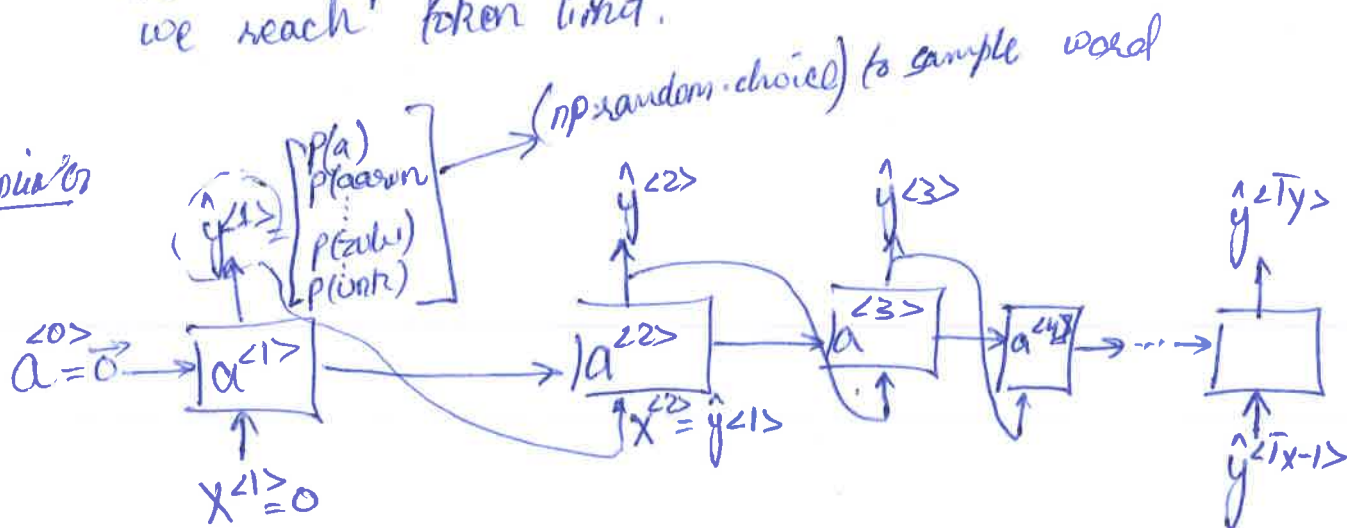
$P(y^{<1>}, \dots y^{<T_y>}) = ??$

TRAINING is done using the following RNN:



For sampling, we will do something different i.e our goal is to generate a new sequence of words. We will keep generating words until we get <EOS> token or we reach token limit.
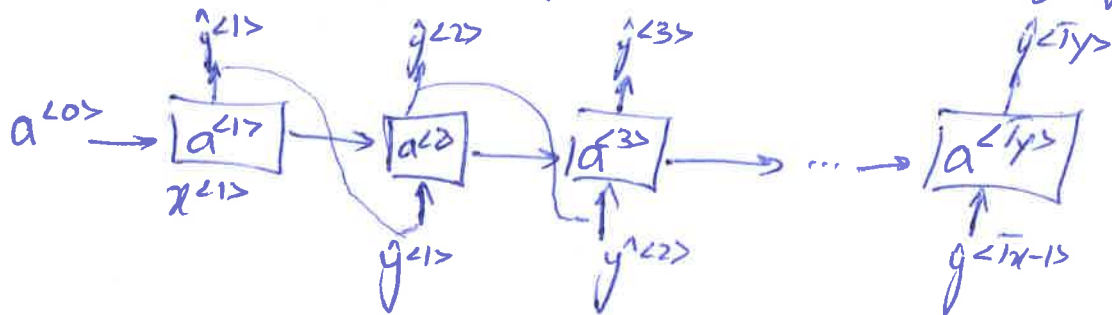
(np.random.choice) to sample word

Sampling

Uptil now, we have seen word-level models
where vocabulary = [a, aaron, ..., zulu, <UNK>]
We can also build character-level model ~~idea~~ where:

$$Vocabulary = [a, b, c, ..., z, \underset{space}{\underline{\quad}}, A, B ... Z, 0, 1 ... 9]$$

And now use RNN to learn character-level language
model.



Here $\hat{y}^{<1>}$ ... $y^{<T_y>}$ are characters,

## W1 - 28

### VANISHING Gradient with RNNs

* RNNs are prone to vanishing gradient problem
  which will be solved using GRUs.
* There is one more issue, exploding gradient which
  is solved using gradient clipping (indication: you will
  see NaN values).

## W1 - 29

### GATED Recurrent Unit (GRU)

It is a modification to RNN's hidden
layer that helps RNN capturing long-range
dependencies

PAPER →  Cho et al., 2014. On the properties of neural machine
         translation. Encoder-decoder approaches

      →  Chung et al. 2014. Empirical Evaluation of Gated
         Recurrent NN on Sequence Modeling

# Intuition for GRUs

RNNs seems to have trouble in learning long-term dependencies.

Sent-1: The <u>cat</u>, which already ate ... <u>was</u> full.

Sent-2: The <u>cats</u>, which already ate ... <u>were</u> full

In both sentences, cat/cats depend on <u>was/were</u>. How to make RNNs learn these type of dependencies.

## Notations (GRU variables)

$c$ = memory cell (it will help memorize ~~senten~~ if cat was singular or plural)

$c^{<t>}$ = memory cell value at location "$t$"

In GRU:

$c^{<t>} = a^{<t>}$ : memory cell value is equal to output activation at time step "$t$". In LSTM, $c^{<t>}$ & $a^{<t>}$ will be different values.

$\tilde{c}^{<t>}$ = candidate to replace $c^{<t>}$

$\Gamma_u$ = u stands for update "$u$" gate

$\Gamma$ " " " gate

# Equations governing GRU unit:

At every timestep "t":

$\tilde{c}^{<t>}$ = candidate to replace $c^{<t>}$

$$\tilde{c}^{<t>} = tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

Important ideas:

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

↖ "u = update"

(eVA) ← $$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$\Gamma_u$ : it is mostly very close to 0 or "1". Most of the times it is 0.

## Illustration:

$\begin{vmatrix} \Gamma_u = 1 \\ c^{<t>} \\ c = 1 \\ cat \end{vmatrix}$  $\Gamma_u = 0$  $\Gamma_u = 0$  - - - - - $\begin{matrix} c^{<t>} \\ \end{matrix}$

Sent: The $\begin{vmatrix} \Gamma_u = 1 \\ c^{<t>} \\ c = 1 \\ cat \end{vmatrix}$, which already ate, ... was full

Lets say we want to learn a concept (Singular subject). Thus we will set $\Gamma_u = 1$ when we have singular noun or singular subject in a sentence. Thus, the moment we detect "cat", we will set $\Gamma_u = 1$ & $c^{<t>} = 1$.

$\Gamma_u = 0$ means don't update $c^{<t>}$ & hang-on its old value. So until the end of sentence, you have successfully memorized that "cat" was singular

## Pictorial View of GRU



$\hat{y}^{<t>}$

softmax

(eVA)

$c^{<t-1>} = a^{<t-1>}$

$\tilde{c}^{<t>}$   $\Gamma_u$

tanh   $\sigma$

$c^{<t>}$
$a^{<t>}$

$x^{<t>}$

Consider eq$(A)$, most of the times $\Gamma_u$ is close to zero (something like 0.000001) so as we scan from left to write, we maintain $c^{<t>} = 1$ as $c^{<t>} = c^{<t-1>}$ when $\Gamma_u \approx 0$. Thus it makes possible to learn a concept even after long-term dependencies.

## Dimensions

Lets assume $c^{<t>} = a^{<t>} = R^{100}$ (i.e 100 hidden units)

Then $\tilde{c}^{<t>}$ & $\Gamma_u$ will also be 100-dimensional

Then in $c^{<t>}$ (eq/A), $*$ will be element-wise multiplication

## Full-GRU Unit

In academic
literature

relevance

$\mathcal{R} \leftarrow \quad \tilde{c}^{<t>} = \tanh\left(w_c\left[\Gamma_r * c^{<t-1>}, x^{<t>}\right] + b_c\right)$

$u \leftarrow \quad \Gamma_u = \sigma\left(w_u\left[c^{<t-1>}, x^{<t>}\right] + b_u\right)$

$r \leftarrow \quad \Gamma_r = \sigma\left(w_r\left[c^{<t-1>}, x^{<t>}\right] + b_r\right)$

$h \leftarrow \quad c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1-\Gamma_u) + c^{<t-1>}$

$*$ GRU makes RNN able to capture long-term dependencies in input. To capture long-term dependencies, GRU & LSTM are the two main ideas that make RNN achieve their goal.

$\underline{\text{LSTM}}$

**Seminal-Paper**
   Hochreiter & Schmidhuber A97. LSTM

**Blog**   Chris OALAH)

___

LSTM lets RNN learn very long-term dependencies (just like GRU)

$$\underline{\text{GRU (2-gates)}}$$

$$\tilde{c}^{<t>} = \dots$$

$$\Gamma_u = \dots$$

$$\Gamma_r = \dots$$

$$c^{<t>} = \dots$$

$$a^{<t>} = c^{<t>}$$

$$\underline{\text{LSTM (3-gates)}}$$

$$\tilde{c}^{<t>} = \tanh\left(W_c[a^{<t-1>}, x^{<t>}] + b_c\right)$$

$$\text{(Update gate)} \quad \Gamma_u = \sigma\left(W_u[a^{<t-1>}, x^{<t>}] + b_u\right)$$

$$\text{(forget gate)} \quad \Gamma_f = \sigma\left(W_f[a^{<t-1>}, x^{<t>}] + b_f\right)$$

$$\text{(output gate)} \quad \Gamma_o = \sigma\left(W_o[a^{<t-1>}, x^{<t>}] + b_o\right)$$

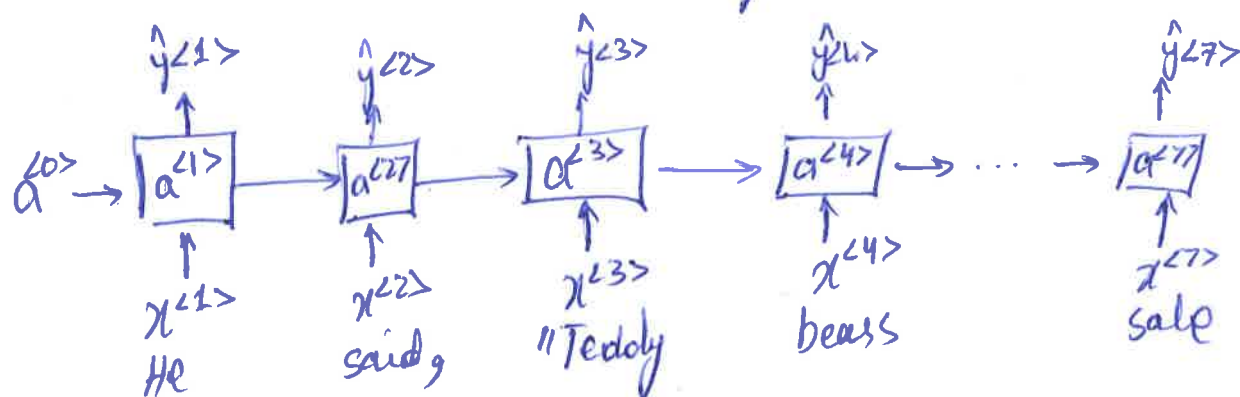$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * c^{<t>}$$

In ~~BRRR~~ Getting information from the future

Sentence-1 : He said, "Teddy bears are on sale!"

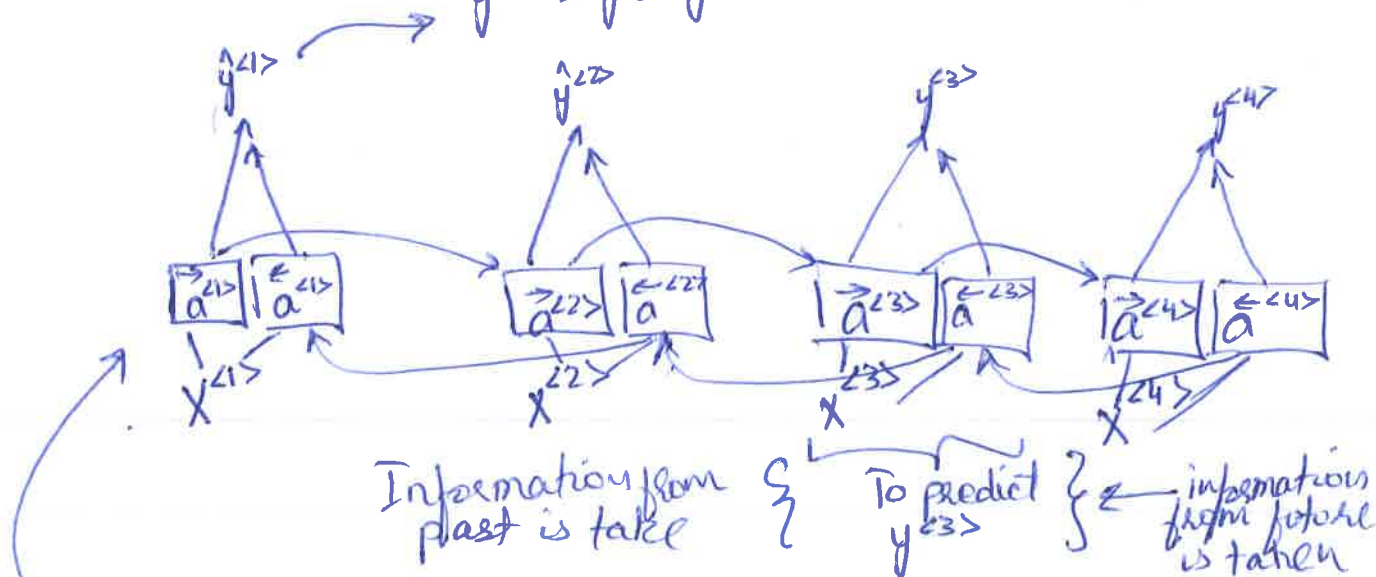Sentence-2 :   "     "  , "Teddy Roosevelt was a great President!"



It is hard to know if $\hat{y}^{<3>} = 1$ or $0$ (i.e Teddy) by just

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\downarrow\quad\quad\downarrow$

$\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad\quad$ Name $\quad$ Not name

having a look on $x^{<1>}$ & $x^{<2>}$. This ~~sd~~ statement is true if above units are ↗RNN, GRU or LSTM units. BRNN helps deciding

$\quad\quad\quad\quad\quad\quad$ simple

whether $\hat{y}^{<3>}$ is name or not by taking into account entire sequence.

## BRNN

### Simplified 4-word sentence

$$\hat{y}^{<t>} = g\left(W_y \left[\overrightarrow{a}^{<t>}, \overleftarrow{a}^{<t>}\right] + b_y\right)$$



Information from { To predict } ← information
plast is take { $y^{<3>}$ } from future
is taken

### Acyclic graph

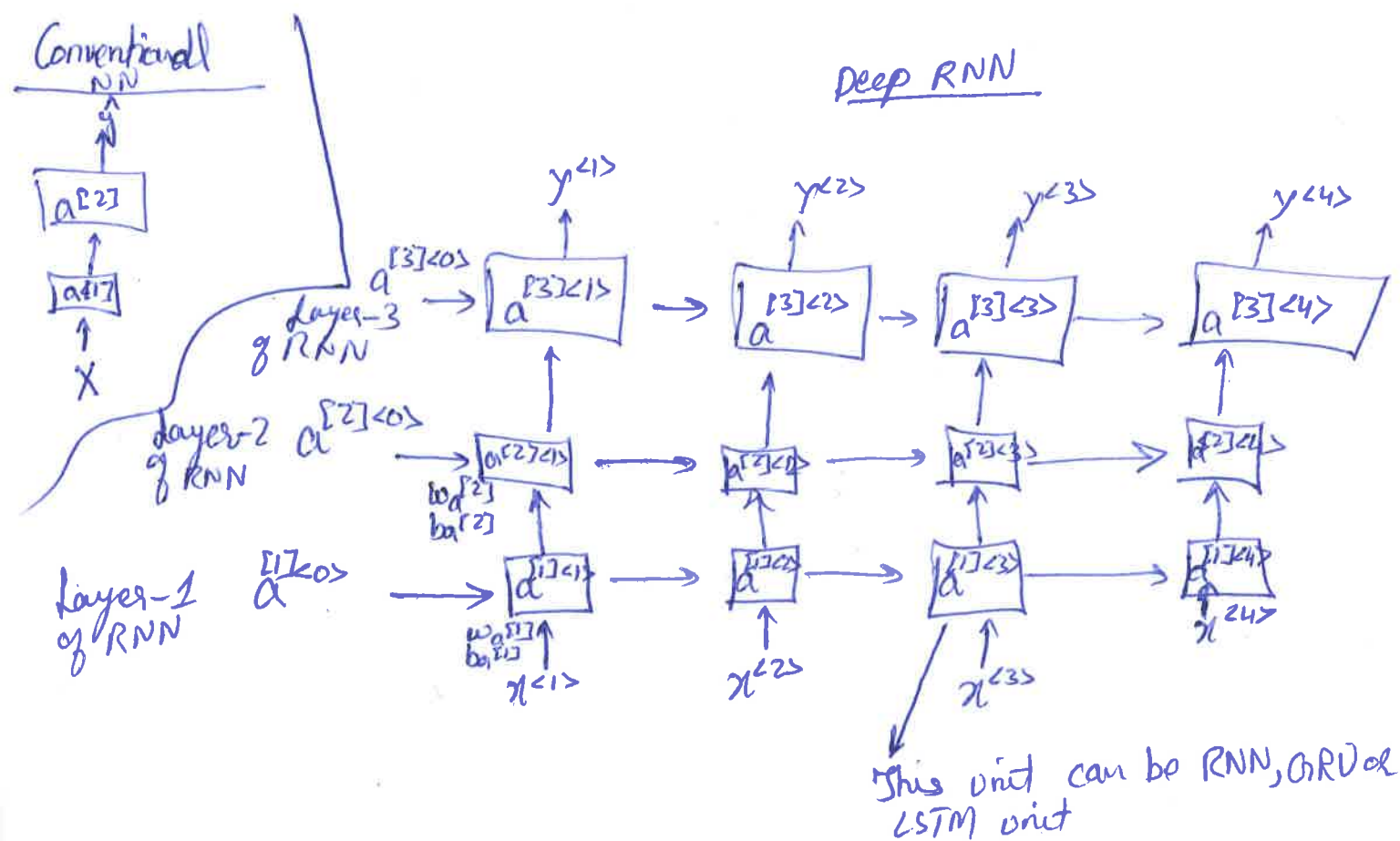For many NLP task, bi-directional RNN with LSTM blocks is most common.

One disadvantage of BRNN is that we need entire sequence to make prediction.

## TAKEAWAY

1) LSTM & GRU help RNN learn long-term dependencies

2) BiRNN helps making prediction by taking into account the entire sequence ~~ie at time "t"~~ ($0 \rightarrow t$ & $t+1 \rightarrow t$ is taken into account for making prediction at "t")

U21-212                          Deep RNN Model

Conventional NN



Deep RNN

$a^{[2]}$

$a^{[1]}$

$X$

Layer-3 of RNN    $a^{[3]<0>}$

Layer-2 of RNN    $a^{[2]<0>}$

Layer-1 of RNN    $a^{[1]<0>}$

$y^{<1>}$   $y^{<2>}$   $y^{<3>}$   $y^{<4>}$

$a^{[3]<1>}$ → $a^{[3]<2>}$ → $a^{[3]<3>}$ → $a^{[3]<4>}$

$a^{[2]<1>}$ → $a^{[2]<2>}$ → $a^{[2]<3>}$ → $a^{[2]<4>}$

$w_a^{[2]}$
$b_a^{[2]}$

$a^{[1]<1>}$ → $a^{[1]<2>}$ → $a^{[1]<3>}$ → $a^{[1]<4>}$

$w_a^{[1]}$
$b_a^{[1]}$

$x^{<1>}$        $x^{<2>}$        $x^{<3>}$        $x^{<4>}$

This unit can be RNN, GRU or LSTM unit

$$a^{[2]<3>} = g\left(w_a^{[2]}\left[a^{[2]<2>}, a^{[1]<3>}\right] + b_a^{[2]}\right)$$

## Modification

We can implement BRNN (for deep RNN) or append another n/w at the output with no-horizontal connection.