

26/7/24

COURSE 5

classmate

Date _____
Page _____

Examples of sequence data:

- Speech recognition
- music generation
- sentiment classification
- DNA seq. analysis
- Video Activity recognition
- Name entity Recognition

Both I/O are seq.

Only input or output is seq.

Example: Name Entity Recognition

x: Harry Potter & Hermoine Granger invented
 $x^{<1>} \quad x^{<2>} \quad x^{<3>} \quad \dots \quad x^{<t>}$
a new spell
 $x^{<q>}$

y: 1 1 0 1 1 0 0 0 0
 $y^{<1>} \quad y^{<2>} \quad \dots \quad y^{<t>} \quad \dots \quad y^{<q>}$

$x^{(i)<t>} \Rightarrow t^{\text{th}}$ element of i^{th} training eg.

$T_x^{(i)} = q$ (here) \Rightarrow length of i^{th} example
 $T_x = T_y$ (here)

* One-Hot Notation

We need a **Vocabulary** for tasks like above, the size of which is generally 30,000 or 50,000.

We'll use 10,000 for learning purpose

Vocabs

 $a^{<1>}$ $a^{<2>}$ $a^{<3>}$

1	a	0	0	0
2	aaron	0	0	0
3	utile	0	0	0
367	aid	0	0	0
4075	harry	1	0	0
6830	potter	0	1	0
10,000	zulu	0	0	0

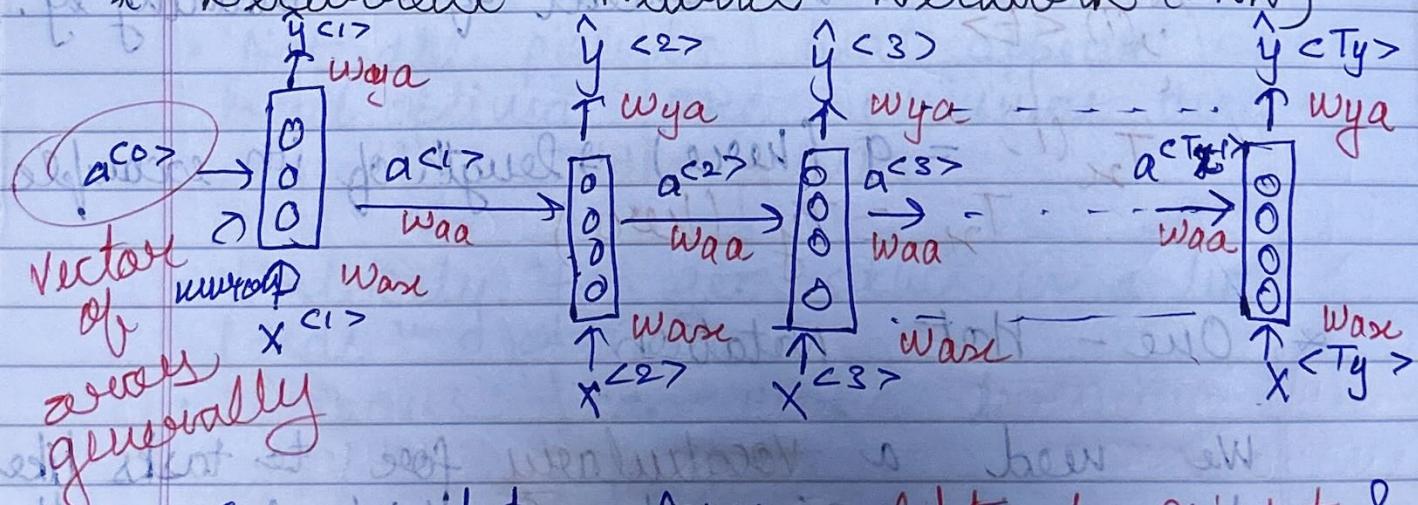
dimension

One-hot, since only one element is 1

* Why not a standard Network?

- I/O can be different lengths in different examples.
- Doesn't share features learned across different positions of text.

* Recurrent Neural Network (RNN)



- Scans data from left to right & the parameters are shared.

- Disadvantage : Considers data only previous to the current element.

Eg :

He said, " Teddy Roosevelt was a President"
 He said, " Teddy Bears are on sale!"
 Here, based off first 2 words cannot decide what Teddy is.

Forward Propagation :

$$a^{(0)} = \vec{0} \quad a^{(1)} = g(w_{00} a^{(0)} + w_{01} x^{(1)} + b_0) \\ \hat{y}^{(1)} = g(w_{y0} a^{(1)} + b_y)$$

\rightarrow tanh / ReLU

$$a^{(t)} = g(w_{00} a^{(t-1)} + w_{01} x^{(t)} + b_0) \\ \hat{y}^{(t)} = g(w_{y0} a^{(t)} + b_y)$$

if classifier, multiplied by a to compute y .
 sigmoid

$$a^{(t)} = g(W_a [a^{(t-1)}, x^{(t)}] + b_a)$$

- compressing a 2 parameter material to a 1 parameter material.

$$W_a = [w_{00}, w_{01}] \in \mathbb{R}^{100 \times 10000}$$

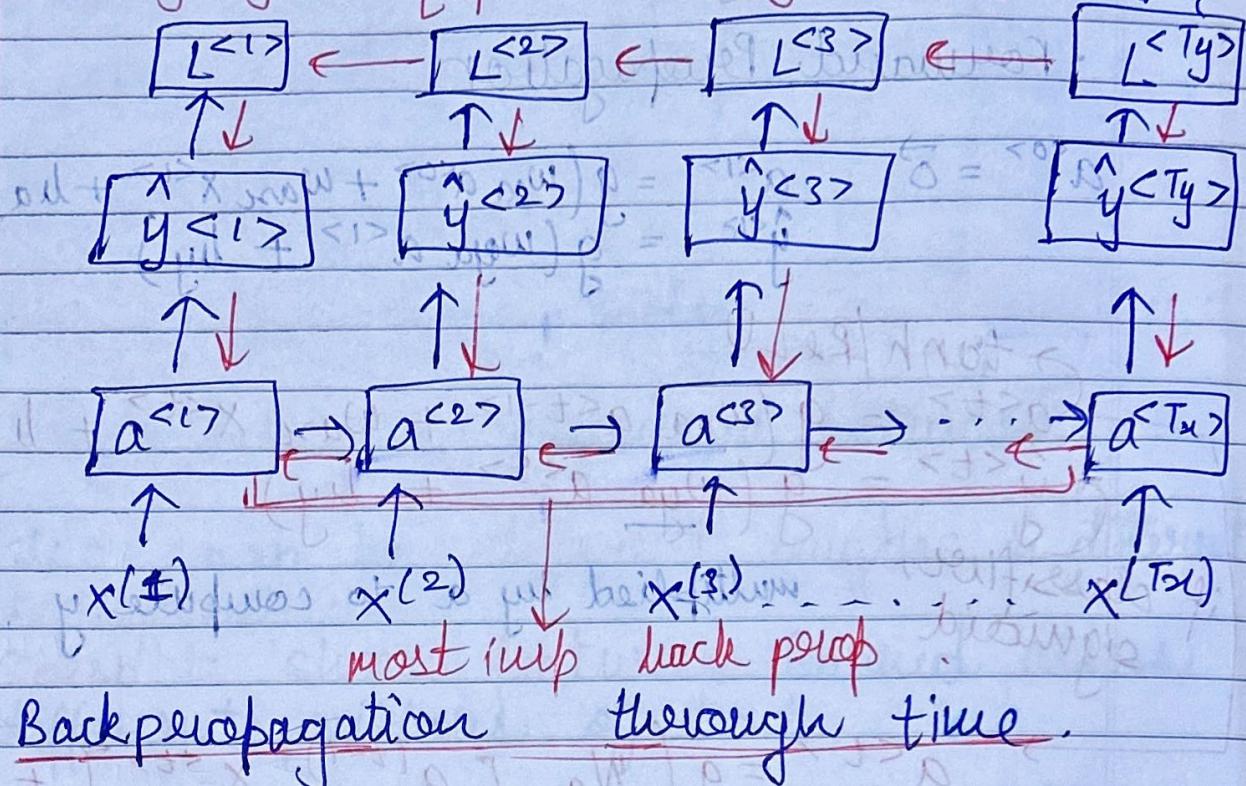
$$[a^{(t-1)}, x^{(t)}] = \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} \in \mathbb{R}^{10000 \times 100}$$

$$[w_{00}, w_{01}] \begin{bmatrix} a^{(t-1)} \\ x^{(t)} \end{bmatrix} = w_{00} a^{(t-1)} + w_{01} x^{(t)}$$

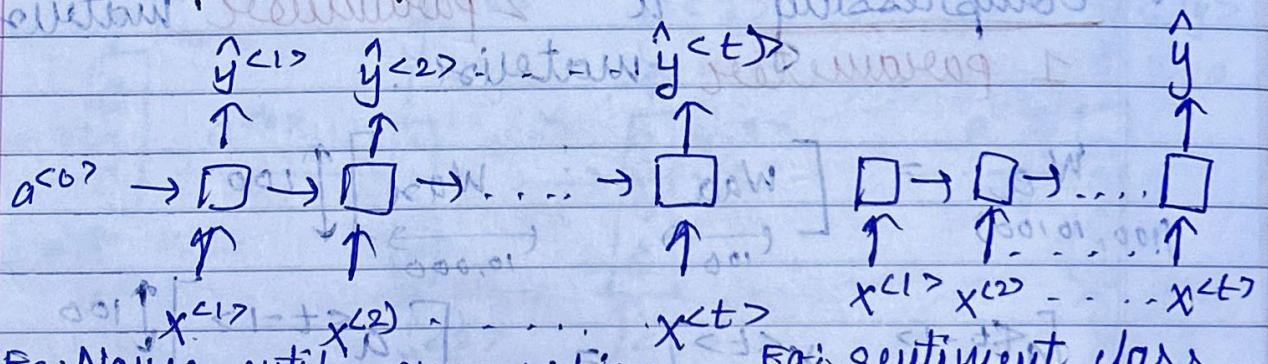
$$L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

$$\text{Junction} = -\hat{y}^{(t)} \log \hat{y}^{(t)} - (1 - \hat{y}^{(t)}) \log (1 - \hat{y}^{(t)})$$

$$L(\hat{y}, y) = \sum_{t=1}^{T_y} L^{(t)}(\hat{y}^{(t)}, y^{(t)})$$



* Examples of RNN architecture:



Eg: Name entity recognition

• Many-to-many
 $T_x = T_y$

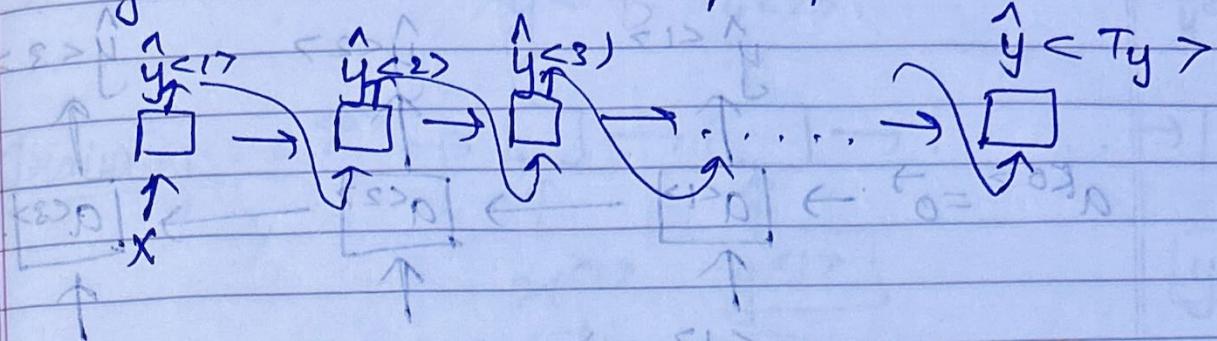
• One-to-one

Eg: sentiment class.

• Many-to-one

- ### • one-to-many

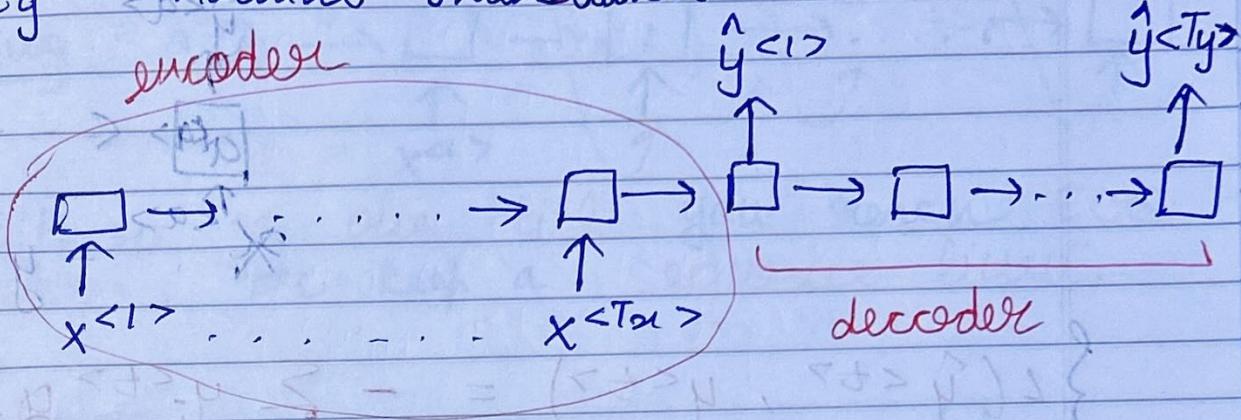
Eg : Music Generation / Seq Generation



- many-to-many $T_{2L} \neq T_Y$

Eg: Machine Translation

encodes



* Language Modelling with an RNN:

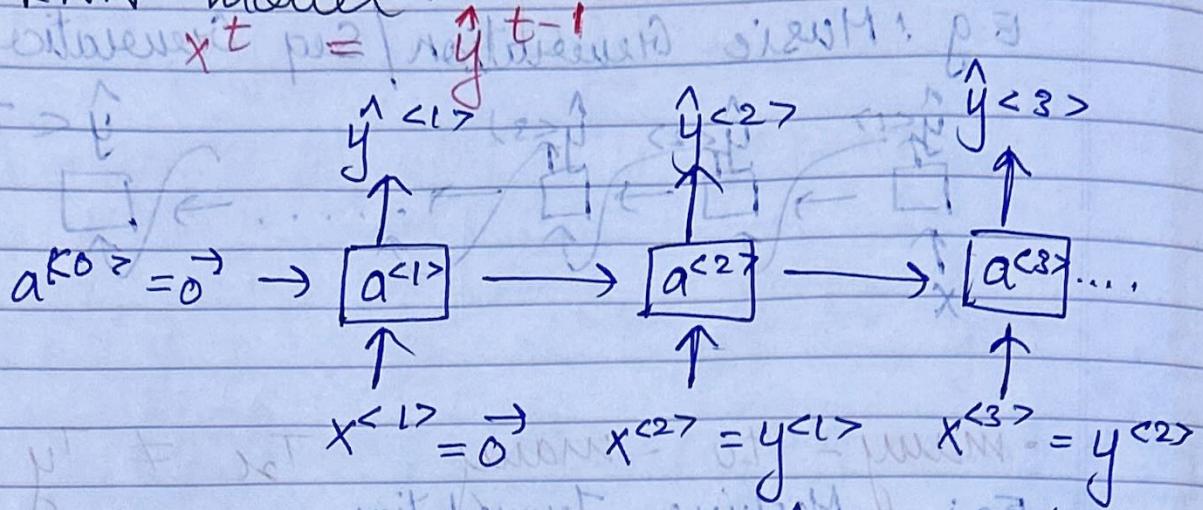
1. Training set: large corpus of English text.
→ Tokenise. (using one-hot or other)

Cats average 15 hours of sleep a day. <EOS>
y<1> | y<2> , y<3> | y<4> | y<5> | y<6>

for word not in dictionary, we give a ^{our} unique token.

~~token~~.
The Egyptian ~~it~~ ^{not in our dict} ~~House~~ is a bread of cat. <EOS>
<UNK> unknown word

RNN model:



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

softmax loss

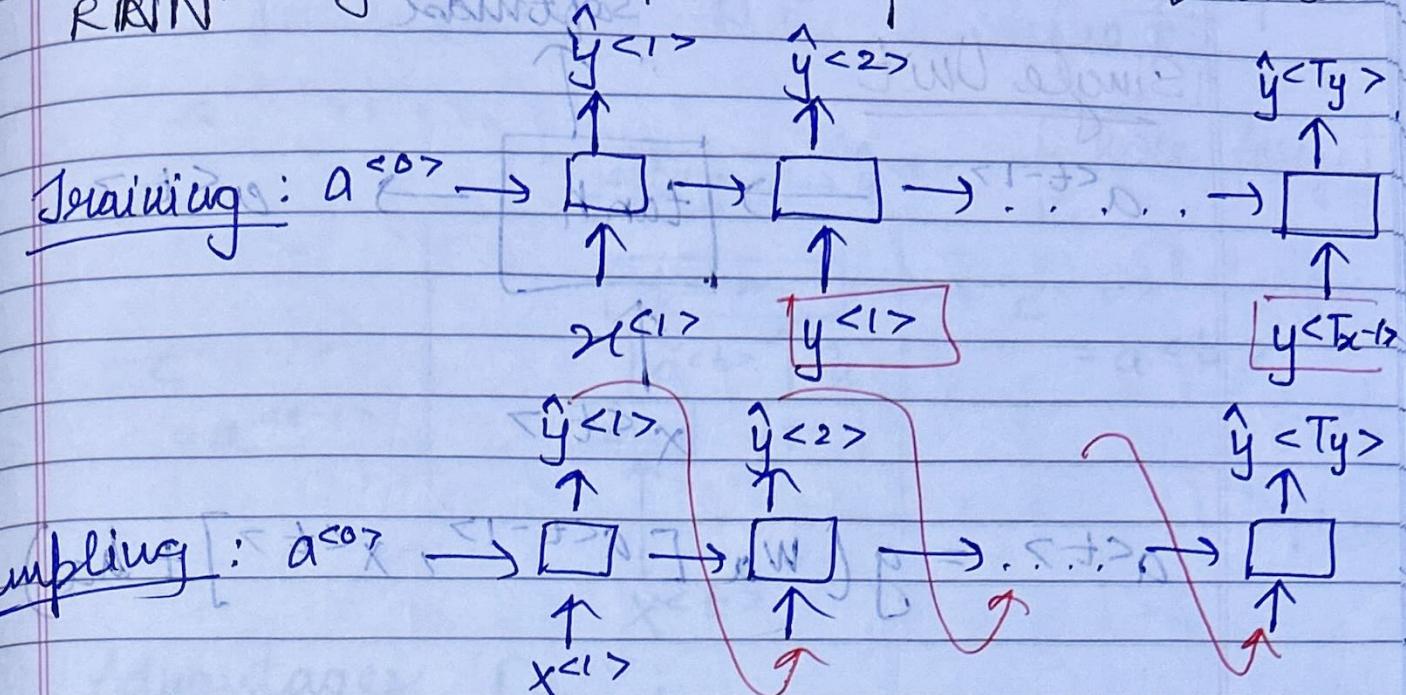
$$L = \sum_t L^{(t)} (\hat{y}^{(t)}, y^{(t)})$$

$$\left\{ \begin{array}{l} P(y^{(1)}, y^{(2)}, y^{(3)}) \\ = P(y^{(1)}) P(y^{(2)} | y^{(1)}) P(y^{(3)} | y^{(1)}, y^{(2)}) \end{array} \right.$$

softmax probability

We continue with this process till we reach the EOS token.

* Sampling a sequence from a trained RNN



Either go ~~out~~ till you reach EOS token, or keep a sequence limit.

* Character - level language model :

Vocabulary: [a, b, ..., z, ,, ., , ..., 9, A, ..., Z]

pros : no worry about unknown word tokens
 cons : can end up with large sequence & it is also computationally expensive

• used only in more specialised applications

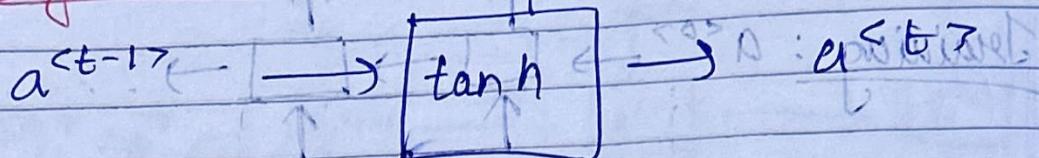
* Vanishing gradients is a problem to be dealt with in RNN

** Exploding gradient $\xrightarrow{\text{apply}} \text{Gradient clipping}
 (not the problem here)$

$y^{<t>}$

Softmax

RNN unit

Single Unit

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

GRU (simplified)

 c = memory cell

$$c^{<t>} = a^{<t>} \quad (\text{in RNN})$$

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

$$= r (w_u [a^{<t-1>}, x^{<t>}] + b_u)$$

uppercase gamma Γ_u update \rightarrow gamma u , always lies b/w 0 & 1
 \rightarrow gamma forget gate

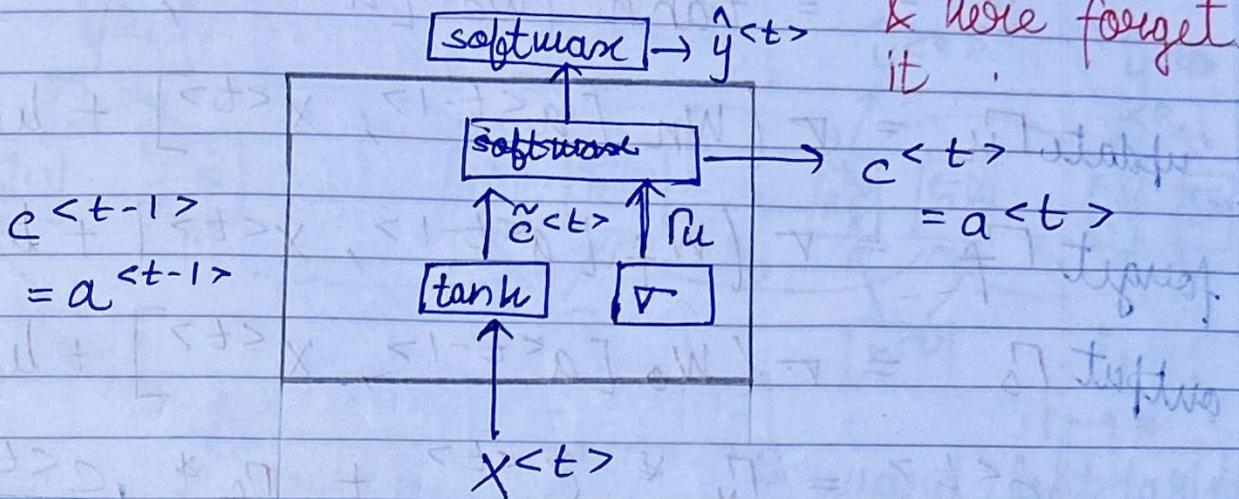
\rightarrow the above func' decides whether or not to update the gate value at the cell at appropriate time.

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

The cat, which already ate.... was full.

$$\text{HT2} \quad F_u = 1 \quad F_u = 0 \quad F_u = 0 \quad F_u = 0 \quad \frac{F_u + 0}{T} = 1$$

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



& here forget it.

Advantages F_u :

- can be very small, thus $c^{<t>} = c^{<t-1>}$
- F_u can be -ive, as long as it is -ive, my round off the update gate will be zero.

$c^{<t>}, \tilde{c}^{<t>}, F_u \Rightarrow$ same dimension

Full GRU:

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

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

$$F_{\text{rel}} = \sigma (W_x [c^{<t-1>} + x^{<t>}] + b_x)$$

(relevance gate)

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

* LS TMs ... to whole slides, too
 → easier & faster versions of LSTM.

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

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

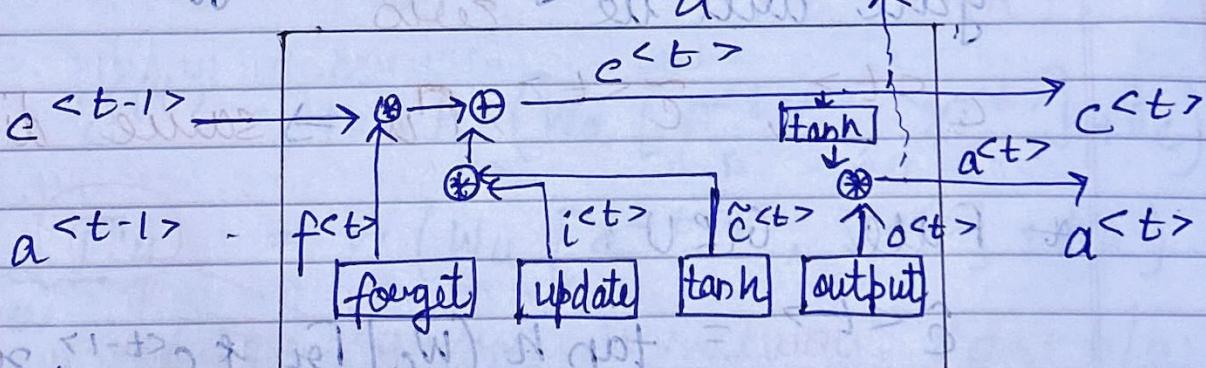
$$\text{forget } \Gamma_f = \sigma(W_f [a^{<t-1>}, x^{<t>}] + b_f)$$

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

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

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

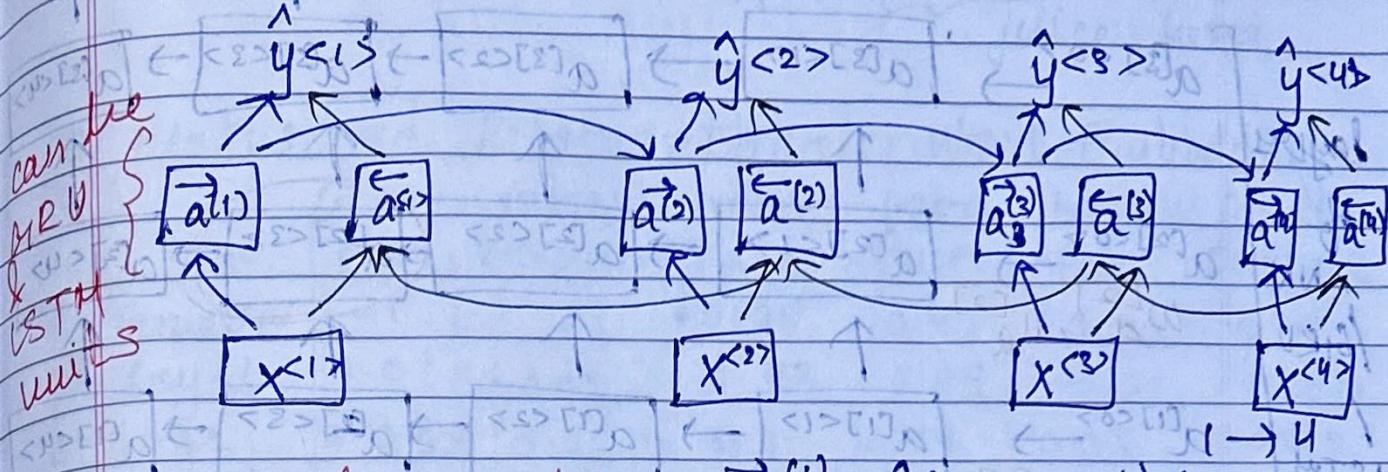
peephole connection



- Re Peephole Connection, $c^{<t-1>}$, shows that it can go for all 3 gates computation as $a^{<t-1>}, x^{<t>}$ determine/affect the value of the gates
- $c^{<t-1>}$, first element affects only first element of gate, thus a one-to-one

connection.

* Bidirectional RNN (BRNN)



- acyclic graph, $\vec{a}^{(1)}$ (forward) left to right
 $\leftarrow \vec{a}^{(1)}$ (backward) right to left

$$\hat{y} = g(w_y [\vec{a}^{(t)}, \leftarrow \vec{a}^{(t)}] + b_y)$$

- At step $\vec{a}^{(3)}$, we can input both info from past & present, which goes into both forward & backward things at this step, & info from future.

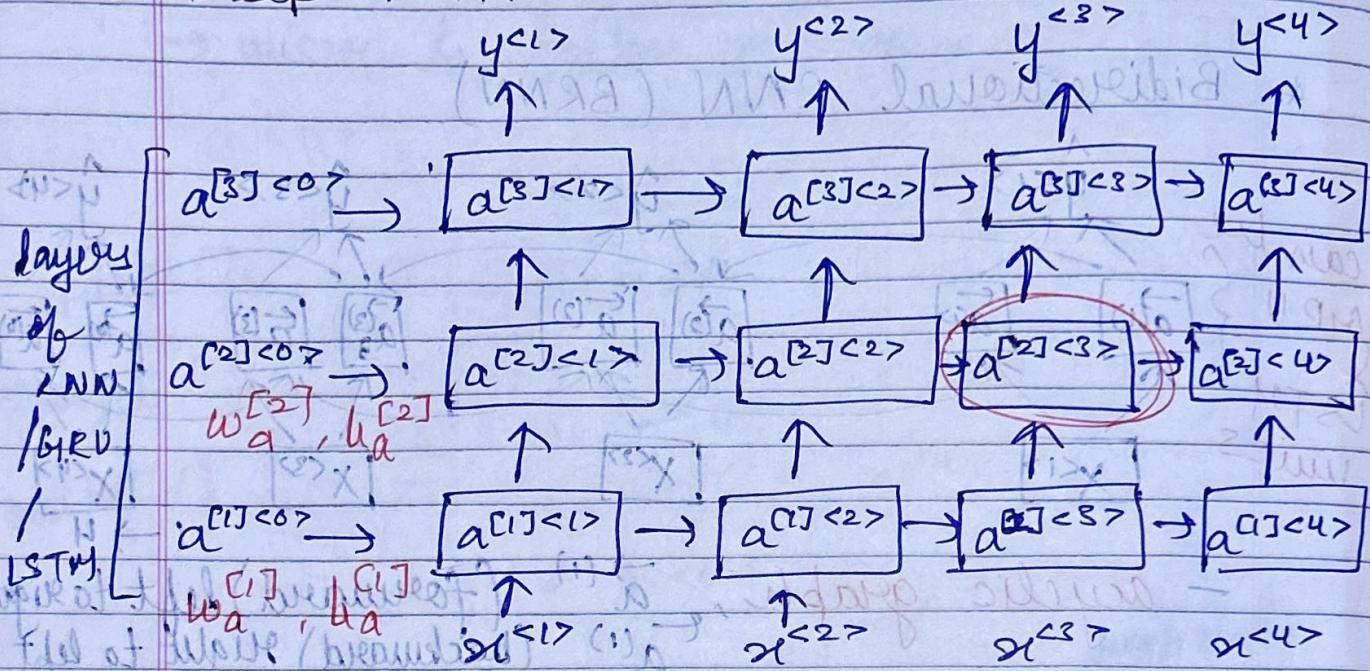
Thus to understand Teddy it takes :

He said, "Teddy Roosevelt . . ."

~~Takes both input to determine what Teddy is,~~

- Disadvantage : You need the entire sequence to make correct predictions
- BRNN w/ LSTM, & other hybrid stacked versions are often used for good efficiency.

* Deep RNN:



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

Computationally expensive.