

Module IV (8 hours)

Recurrent Neural Networks: Back Propagation, vanishing gradients, exploding gradients, truncated backpropagation through time, Gate Recurrent Units (GRUs), Long Short-Term Memory (LSTM) cells, solving the vanishing gradient problem with LSTMs

Sequence Learning problem:-

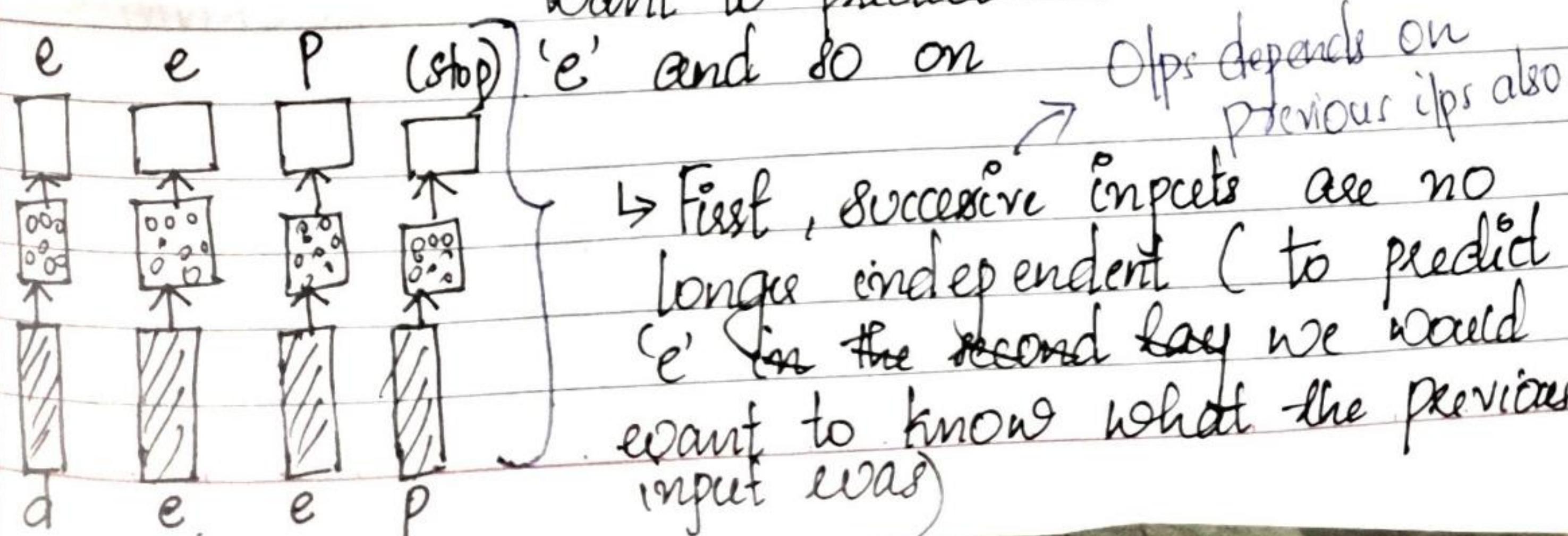
- * In feedforward and convolutional neural nets the size of the input was always fixed

e.g.: fed fixed size (32×32) images to convolutional neural nets for image classification

- * Also, each input to the network was independent of the previous or future inputs

- * In many applications the input is not of a fixed size. Moreover, successive inputs may not be independent of each other

e.g.: Consider the task of auto completion
Given the first character 'd' you want to predict the next character



The length of the input
is not fixed
↳ Second, the length of the inputs
and the number of predictions
you need to make is not fixed
eg (deep, learn, machine, have
different number of characters)

↳ Third, each n/w is performing the
same task
($E/p = \text{character}$
 $O/p = \text{a character}$)

eg:-
↳ Sentimental Analysis
↳ Speech Recognition
↳ Video Labeling

How to model tasks involving sequences

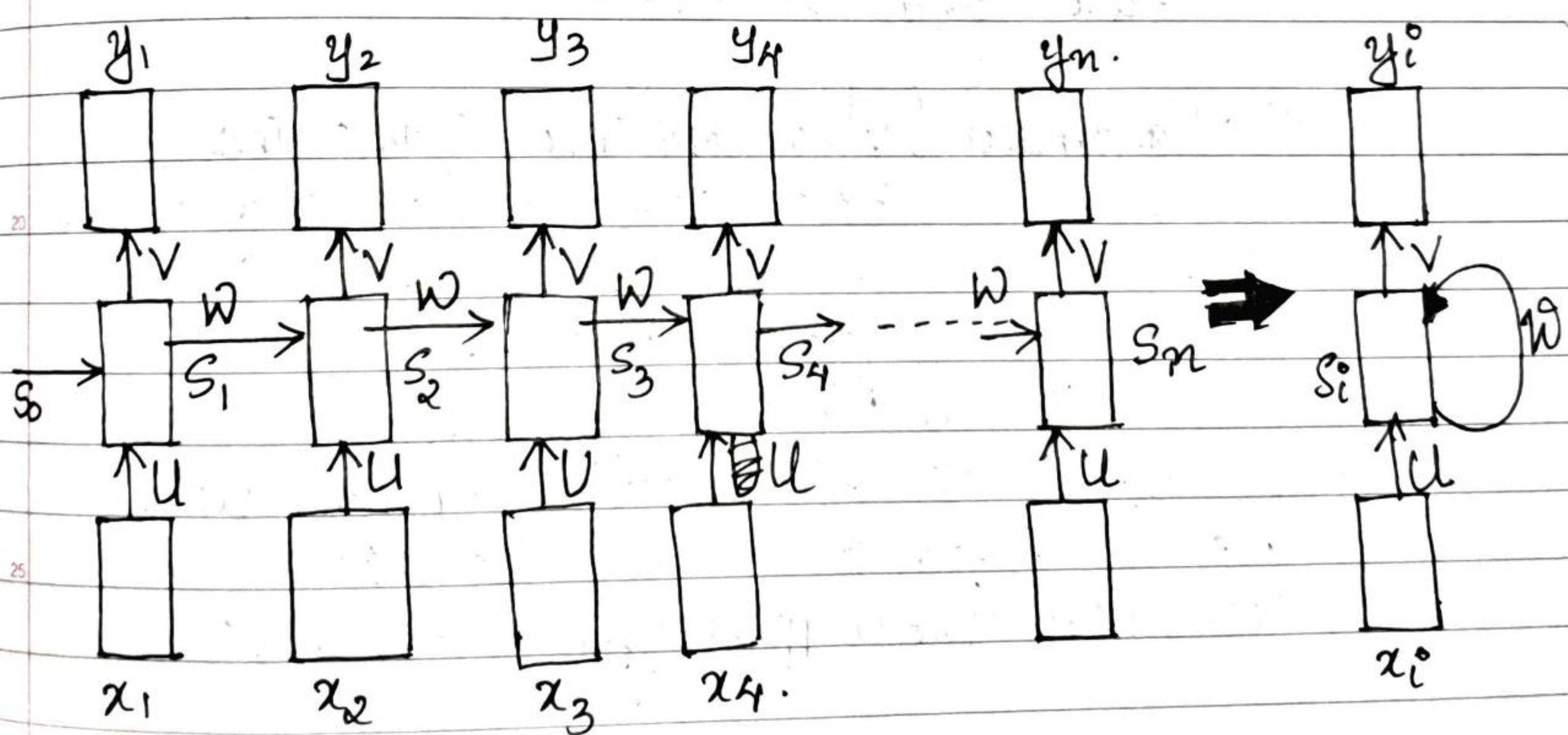
Target points:-

- * Account for dependence b/w inputs
- * Account for variable number of inputs
- * The function executed at each time step is the same.

To train a model that has a kind of sequence
inputs need Recurrent Neural Network (RNN)

Recurrent Neural Network:-

- * RNN is a type of neural network in which the output from the previous step is fed as input to the current step. It is a type of ANN which uses sequential data or time series data.
- * RNN is normally used in NLP and some other domains.
- * RNN's also have a hidden stage which used to capture information about a sentence (Scenario-NLP).
- * RNN's have a memory, which is used to capture information about the calculations made so far.



$$s_i^o = \sigma(Ux_i + Ws_{i-1}^o + b)$$

$$y_i^o = O(Vs_i^o + c)$$

σ

$$y_i^o = f(x_i^o, s_i^o, U, V, W)$$

hidden state at time 'i'
where s_i = state of the network at timestep 'i'.
 c = true class

* The parameters are W, U, V which are shared across timesteps.

* Same network can be used to compute y_1, y_2, \dots, y_{10} or y_n .

Backpropagation.

* Parameter 'U' transformation the input x_t to the state s_t .

* Parameter 'W' transforms the previous state s_{t-1} to the current state s_t .

* Parameter 'V' maps the computed internal state s_t to the output o_t .

s is replaced with h

* If tanh is the activation fn used and softmax is the activation fn used in the o/p layer

$$h_t = \tanh(Ux_t + Wh_{t-1})$$

$$y_t = \text{softmax}(Vh_t)$$

Types of RNN Architectures:

- * The common architectures

 - ↳ One to One

 - ↳ One to Many

 - ↳ Many to One

 - ↳ Many to Many

- * One to One:-

 - ↳ This model is similar to a single layer neural net as it only provides limited predictions

 - ↳ It is mostly used fixed-size ip 'x' and fixed-size op 'y'

eg:- image classification
(e.g. dog vs cat)

y (op layer)

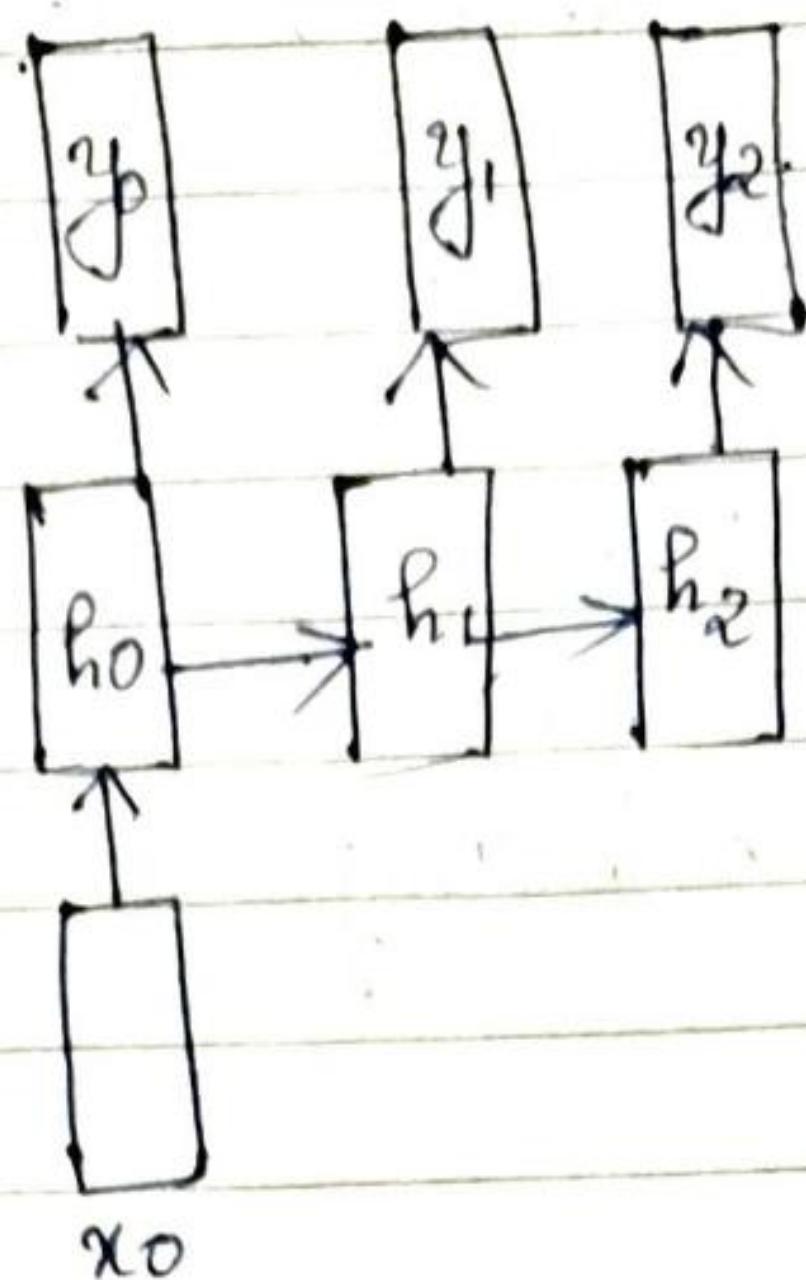
hidden layer

x (ip layer)

- * One to Many:-

 - ↳ This consist of a single ip 'x', activation 'a' and multiple op 'y'.

eg:- Generating an audio stream. If takes a single audio stream as ip and generates new tones or new music based on that stream.



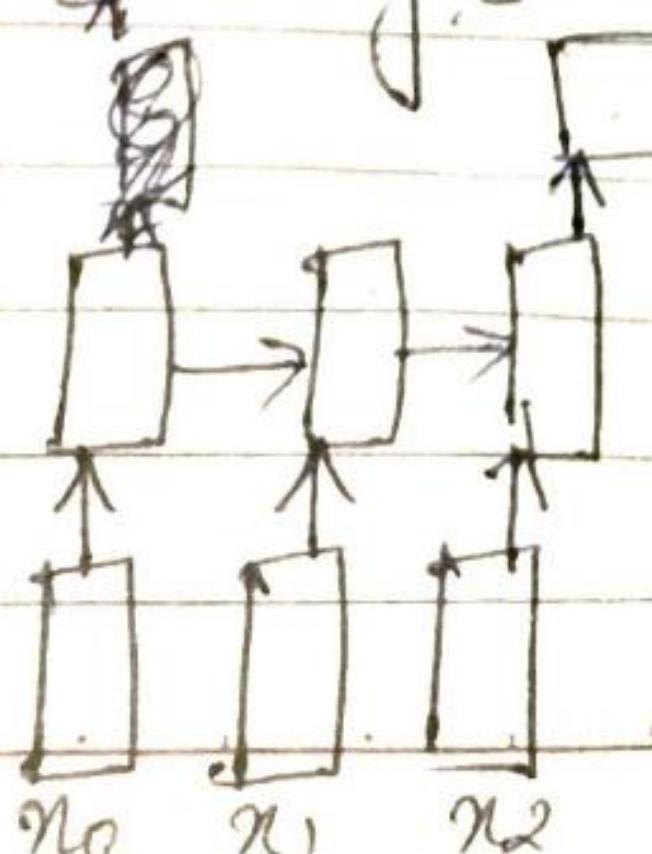
eg:- Image captioning
 (i/p as image and o/p tells the content in the image)

Many to One:-

* This consist of multiple inputs 'x' (such as words or sentence), activation 'a' and produce a single o/p 'y' at the end.

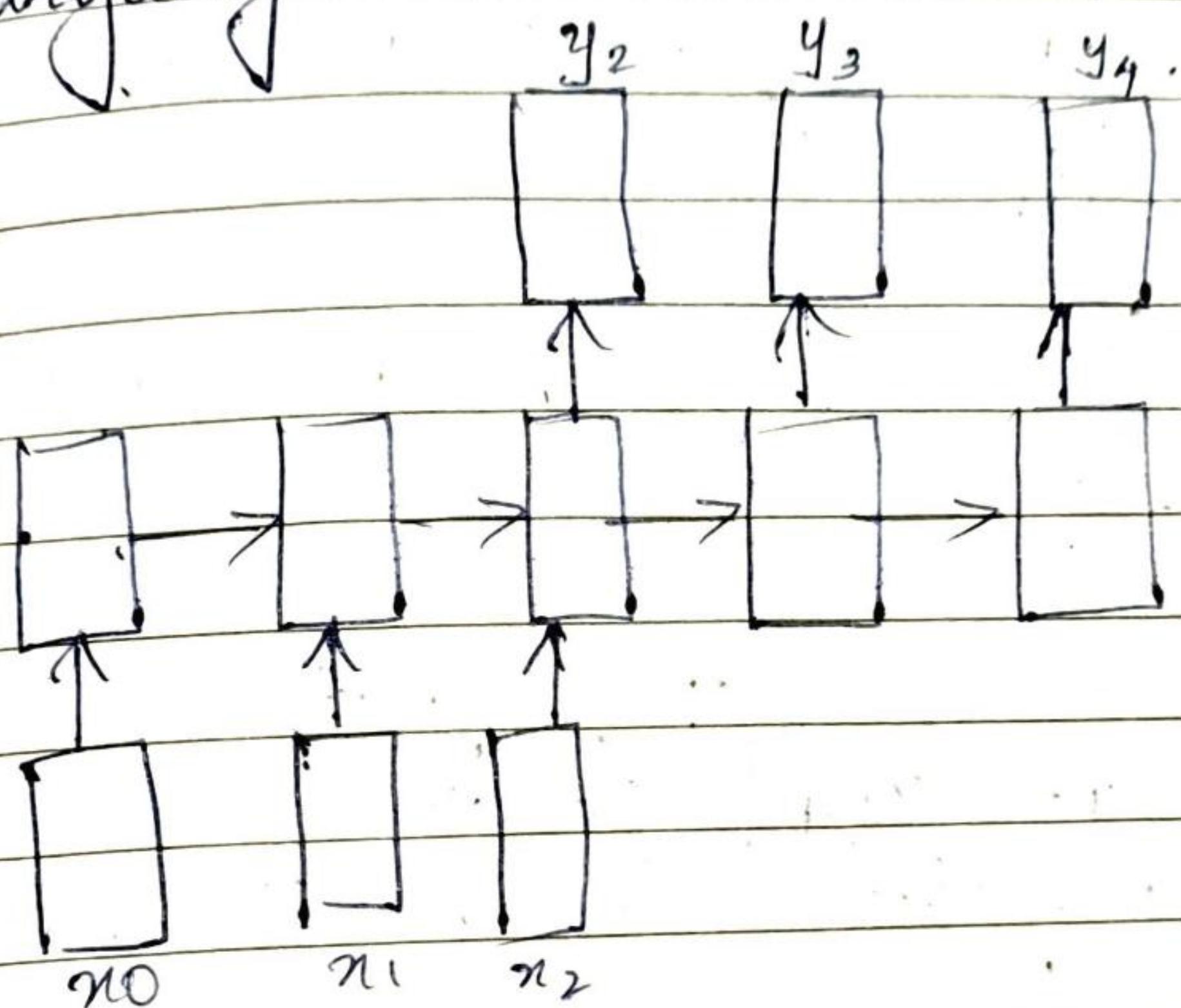
* This type of architecture mostly used to perform sentiment analysis as it processes the entire i/p (collection of words) to produce a single o/p (positive, negative, or neutral sentiment)

* eg:- stock price prediction



Many to Many :-

→ Many-to-Many sequence learning can be used for machine translation where the IP sequence is in some language and the OP sequence is in some other language.



Backpropagation through time (BPTT)

* Dimensions of the parameters

$$x_i \in \mathbb{R}^n \text{ (n-dimensional ip)}$$

$$s_i \in \mathbb{R}^d \text{ (d-dimensional state)}$$

$$y_i \in \mathbb{R}^k \text{ (say 'k' classes)}$$

$$w \in \mathbb{R}^{n \times d}$$

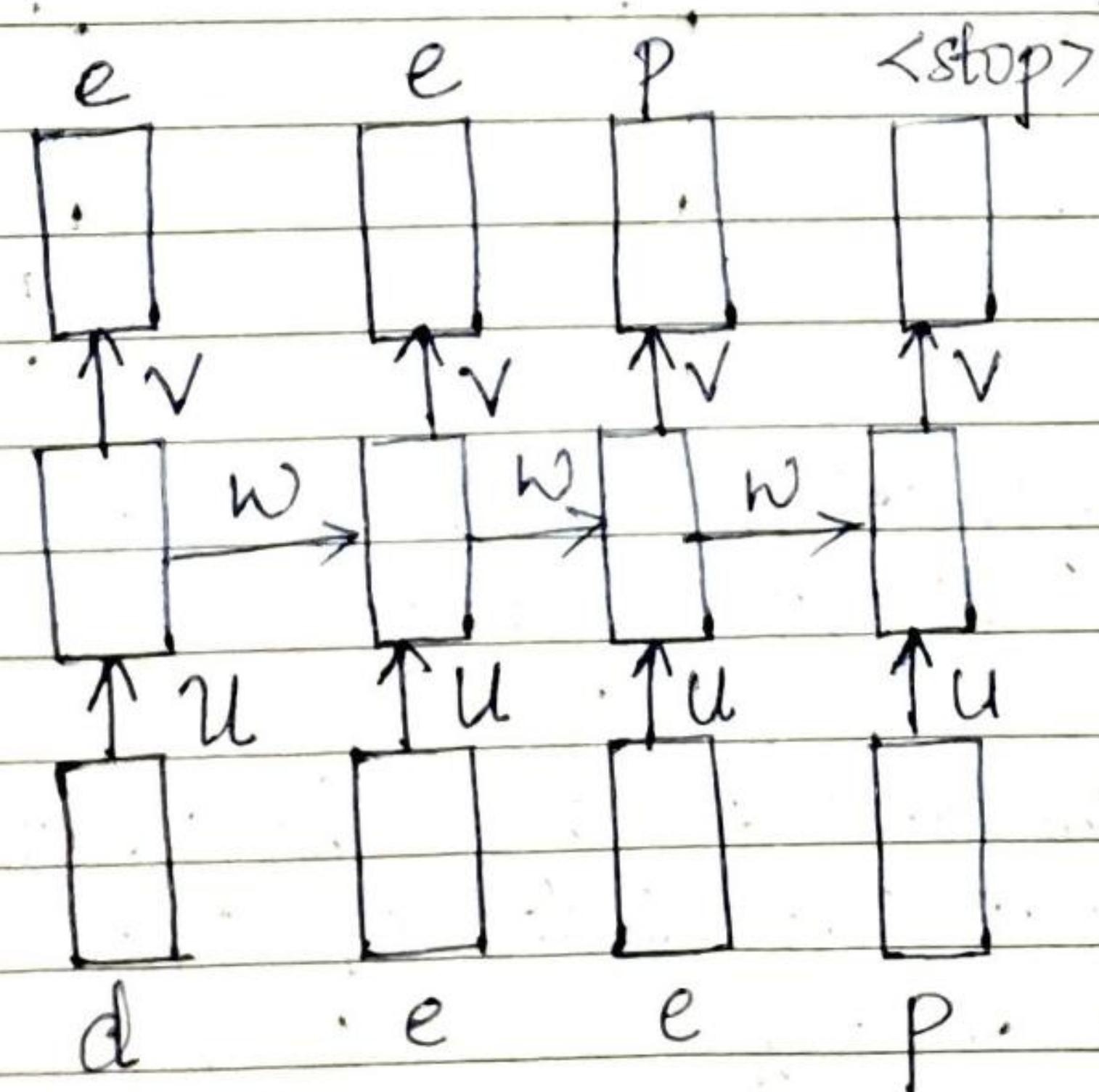
$$v \in \mathbb{R}^{d \times k}$$

$$w \in \mathbb{R}^{d \times d}$$

* let us consider an example "task of auto completion (predicting the next character)

↳ for simplicity we assume that there are only α characters in our vocabulary
· (d, e, p, <stop>)

↳ At each timestep we want to predict one of these α characters



* Suitable op fn for this task is softmax

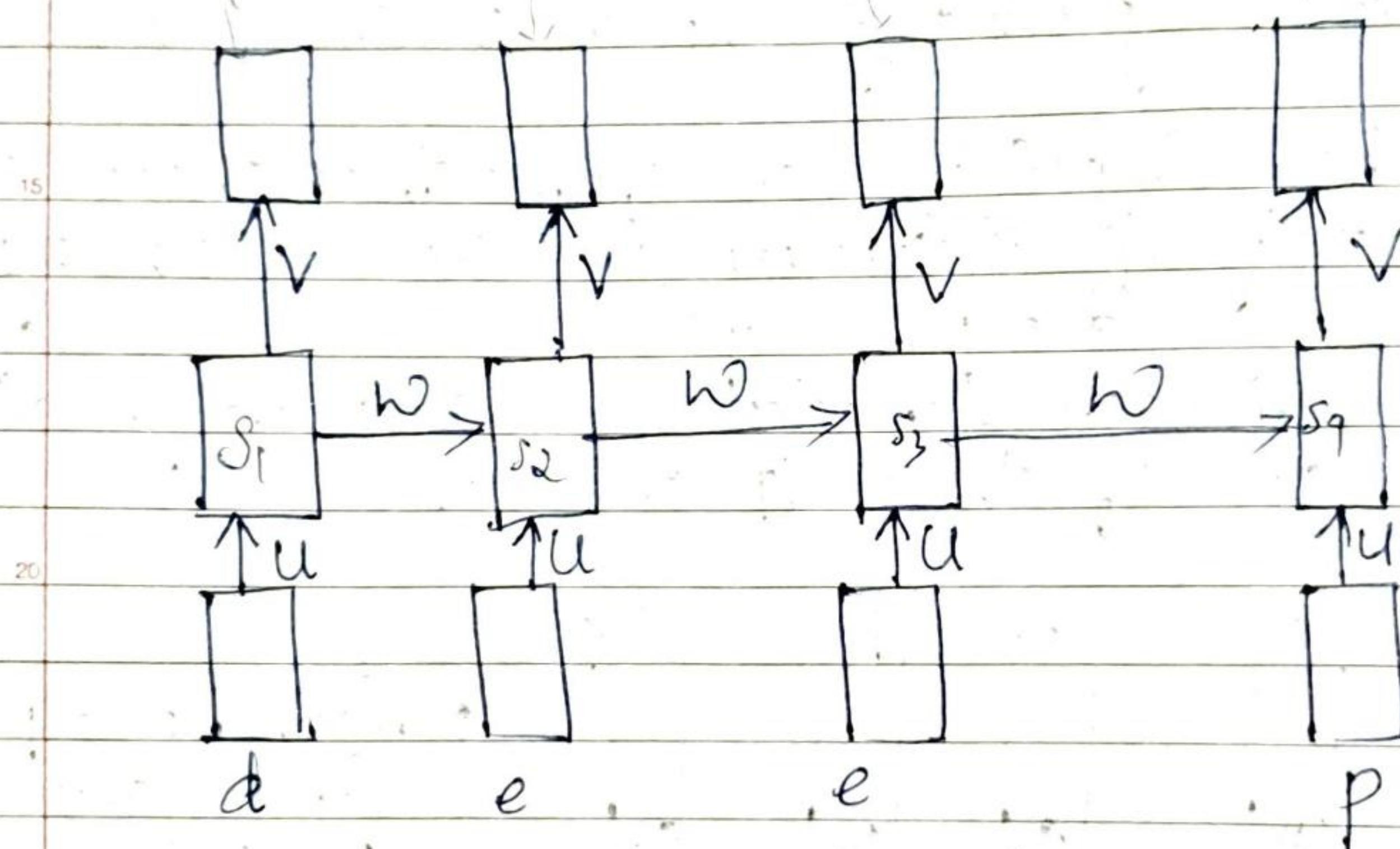
* Suitable loss fn for this task is cross entropy

* Suppose we initialize u, v, w randomly and the net predicts the probabilities

* Initially predict 'd' as input, u, v, w are randomly initialized except matrix which predicts some probability distribution

Predicted	Predicted	Predicted	Predicted
d 0.2	0	d 0.2	0.2
e 0.1	1	e 0.7	0.1
p 0.1	0	p 0.1	0.7
stop 0.1	0	stop 0.1	0.7

true probability true probability



* Objective function is to reduce the error on 'd' loss function

* Total loss made by the model is simply the sum. of the loss over all time steps T

$$L(\theta) = \sum_{t=1}^T L_t(\theta)$$

$$L_f(\theta) = -\log(p_{t^*})$$

↗ e at time step 1
 ↗ e at time step 2
 ↗ e at time step 3
 ↗ e at time step 4
 ↗ true class at time step t
 ↗ loss at time step t

p_{t^*} = Predicted probability of trace character at time step t^* .

T = number of time steps.

* For backpropagation we need to compute the gradient w.r.t w, u, v

→ ① $\frac{\partial L(\theta)}{\partial v}$

$$\frac{\partial L(\theta)}{\partial v} = \sum_{t=1}^T \frac{\partial L_t(\theta)}{\partial v}$$

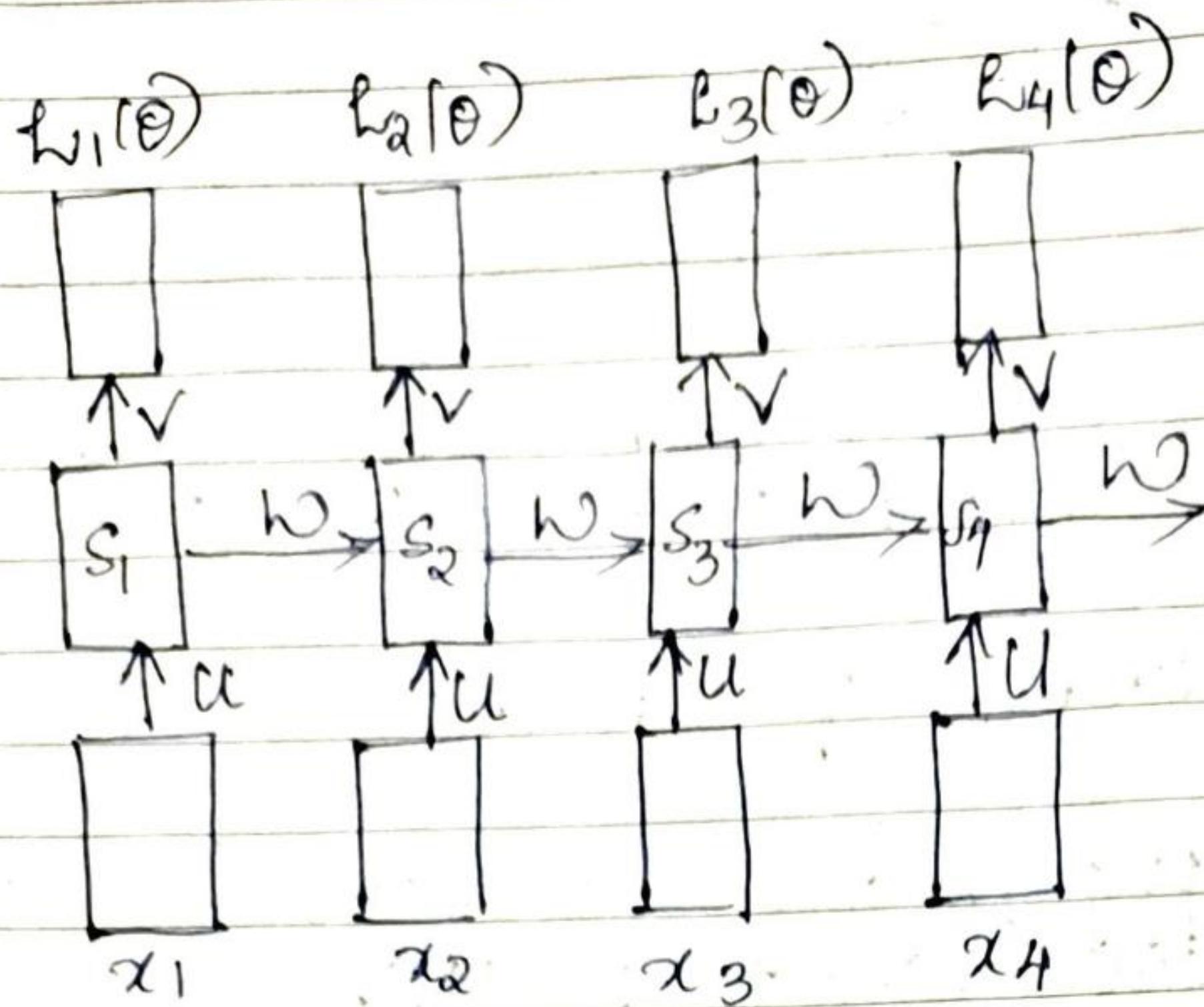
* Each term is the summation is simply the derivative of the loss w.r.t the weights in the OLP layer

→ ② $\frac{\partial L(\theta)}{\partial w}$

$$\frac{\partial L(\theta)}{\partial w} = \sum_{t=1}^T \frac{\partial L_t(\theta)}{\partial w}$$

→ By the chain rule of derivatives we know that $\frac{\partial L_t(\theta)}{\partial w}$ is obtained

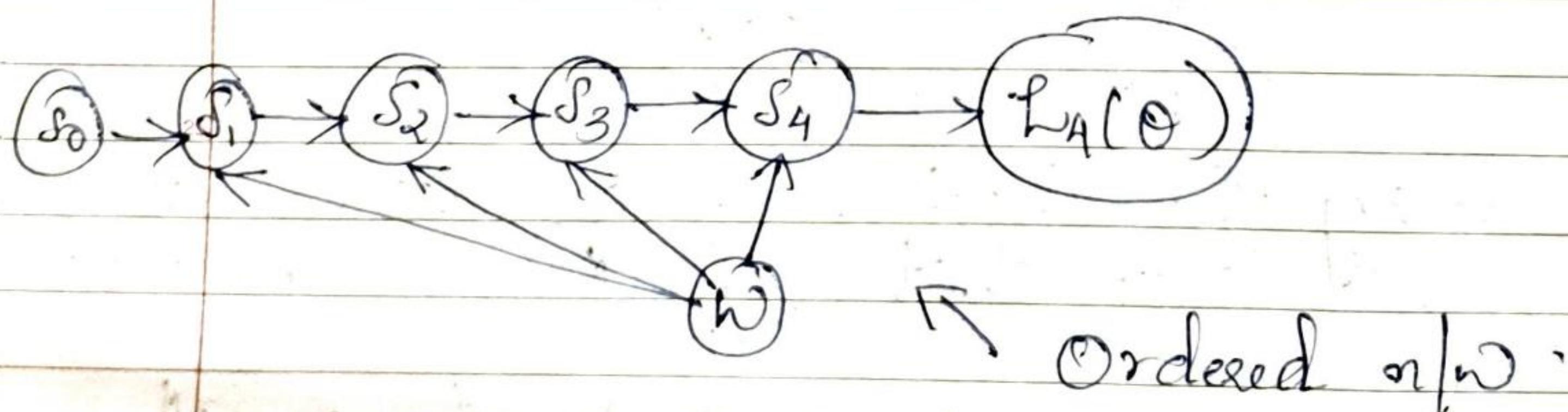
by summing gradients along all the paths from $L_t(\theta)$ to w .



No. of Paths connecting $L_4(\theta)$ to $w = 4$ (according to eq)

- * $L_4(\theta)$ depends on s_4
- * s_4 depends on w and s_3
- * s_3 depends on w and s_2
- * s_2 depends on w and s_1
- * s_1 depends on w and s_0

where s_0 is a constant starting state



- * In an ordered n/w each state variable is computed at a time in a specified order (first s_1 , then s_2 and so on)

* Now we have

$$\frac{\partial h_4(\theta)}{\partial w} = \frac{\partial h_4(\theta)}{\partial s_4} \frac{\partial s_4}{\partial w}$$

Can be computed easily using chain rule since $h_4(\theta)$ depends on V and V depends on s_4

$$\frac{\partial h_4(\theta)}{\partial s_4} = \frac{\partial h_4(\theta)}{\partial V} \times \frac{\partial V}{\partial s_4}$$

* $\frac{\partial s_4}{\partial w} = \frac{\partial}{\partial w} (\sigma(ux_4 + ws_3 + b))$
 (Let's eliminate σ)

$$\frac{\partial s_4}{\partial w} = \frac{\partial}{\partial w} (ux_4 + ws_3 + b)$$

$$= \otimes S_3$$

s_4 depends on s_3 , which depends on s_2 and so on
 (In such an ordered n/w, we can't compute $\frac{\partial s_4}{\partial w}$ by simply treating

s_3 as constant (bcz it also depends on w)

↳ In such ordered n/w the total derivative $\frac{\partial s_4}{\partial w}$ has two parts
 ↳ Explicit and Implicit

Explicit:- $\frac{\partial s_4}{\partial w}$, treating all other inputs as constant

Implicit :- Summing over all implicit paths from s_A to w

$$\frac{ds_A}{dw} = \underbrace{\frac{d^+s_A}{dw}}_{\text{Explicit}} + \underbrace{\frac{ds_A}{ds_3} \frac{ds_3}{dw}}_{\text{Implicit}}$$

$$= \frac{d^+s_A}{dw} + \frac{ds_A}{ds_3} \left[\underbrace{\frac{d^+s_3}{dw} + \frac{ds_3}{ds_2} \frac{ds_2}{dw}}_{\text{Explicit}} \right] \underbrace{\frac{ds_2}{ds_1} \frac{ds_1}{dw}}_{\text{Implicit}}$$

$$= \cancel{\frac{d^+s_A}{dw}} + \cancel{\frac{ds_A}{ds_3}} \frac{d^+s_3}{dw} + \cancel{\frac{ds_3}{ds_2}} \left[\cancel{\frac{d^+s_2}{dw}} + \cancel{\frac{ds_2}{ds_1} \frac{ds_1}{dw}} \right]$$

$$= \frac{d^+s_A}{dw} + \frac{ds_A}{ds_3} \frac{d^+s_3}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \left[\frac{d^+s_2}{dw} + \frac{ds_2}{ds_1} \frac{ds_1}{dw} \right]$$

$$\frac{ds_2}{ds_1}$$

$$= \frac{d^+s_A}{dw} + \frac{ds_A}{ds_3} \frac{d^+s_3}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{d^+s_2}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{ds_2}{ds_1} \frac{ds_1}{dw}$$

$$\left[\frac{ds_1}{dw} \right]$$

$$= \frac{d^+s_A}{dw} + \frac{ds_A}{ds_3} \frac{d^+s_3}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{d^+s_2}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{ds_2}{ds_1} \frac{d^+s_1}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{ds_2}{ds_1} \frac{ds_1}{dw}$$

$$= \frac{d^+s_A}{dw} + \frac{ds_A}{ds_3} \frac{d^+s_3}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{d^+s_2}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{ds_2}{ds_1} \frac{d^+s_1}{dw} + \frac{ds_A}{ds_3} \frac{ds_3}{ds_2} \frac{ds_2}{ds_1} \frac{ds_1}{ds_0} \frac{d^+s_0}{dw}$$

* For simplicity, we will short-circuit some of the paths.

$$\frac{\partial s_4}{\partial w} = \frac{\partial s_4}{\partial s_4} \frac{\partial^+ s_4}{\partial w} + \frac{\partial s_4}{\partial s_1} \frac{\partial^+ s_1}{\partial w} + \frac{\partial s_4}{\partial s_2} \frac{\partial^+ s_2}{\partial w} + \frac{\partial s_4}{\partial s_3} \frac{\partial^+ s_3}{\partial w} + \frac{\partial s_4}{\partial s_0} \frac{\partial^+ s_0}{\partial w}$$

neglected

$$\frac{\partial s_4}{\partial w} = \sum_{k=1}^4 \frac{\partial s_4}{\partial s_k} \frac{\partial^+ s_k}{\partial w}$$

* Finally we have

$$\frac{\partial L_A(0)}{\partial w} = \frac{\partial L_A(0)}{\partial s_4} \frac{\partial s_4}{\partial w}$$

$$\text{we know } \frac{\partial s_4}{\partial w} = \sum_{k=1}^4 \frac{\partial s_4}{\partial s_k} \frac{\partial^+ s_k}{\partial w}$$

In general,

$$\frac{\partial L_t(0)}{\partial w} = \frac{\partial L_A(0)}{\partial s_t} \sum_{k=1}^t \frac{\partial s_t}{\partial s_k} \frac{\partial^+ s_k}{\partial w}$$

* This algorithm is called Backpropagation through time (BPTT) as we backpropagate over all previous time steps

Ques.

→ (3)

$$\boxed{\frac{\partial L(\theta)}{\partial u}}$$

$$\frac{\partial L(\theta)}{\partial u} = \sum_{k=0}^T \frac{\partial L(\theta)}{\partial s_k} \frac{\partial s_k}{\partial u}$$

$$\boxed{\frac{\partial s_k}{\partial u} = x_k}$$

$$\sum_{k=0}^T \frac{\partial L(\theta)}{\partial s_k} = x_k$$

$$\frac{\partial L(\theta)}{\partial s_k} = ?$$

$$s_k = \sigma(u x_k + w s_{k-1} + b)$$

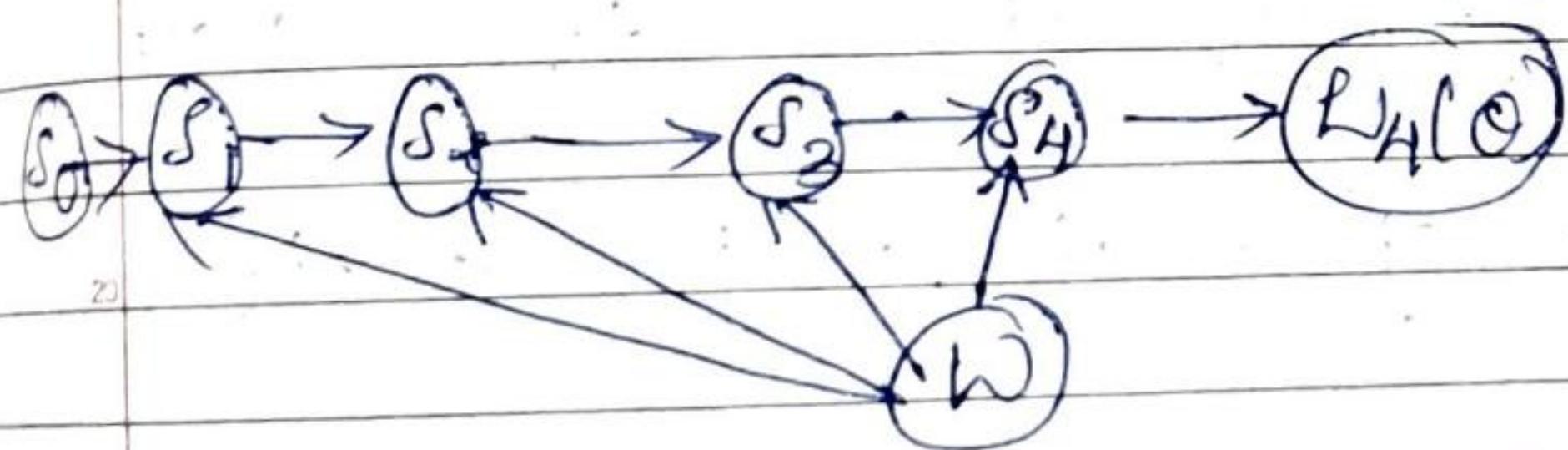
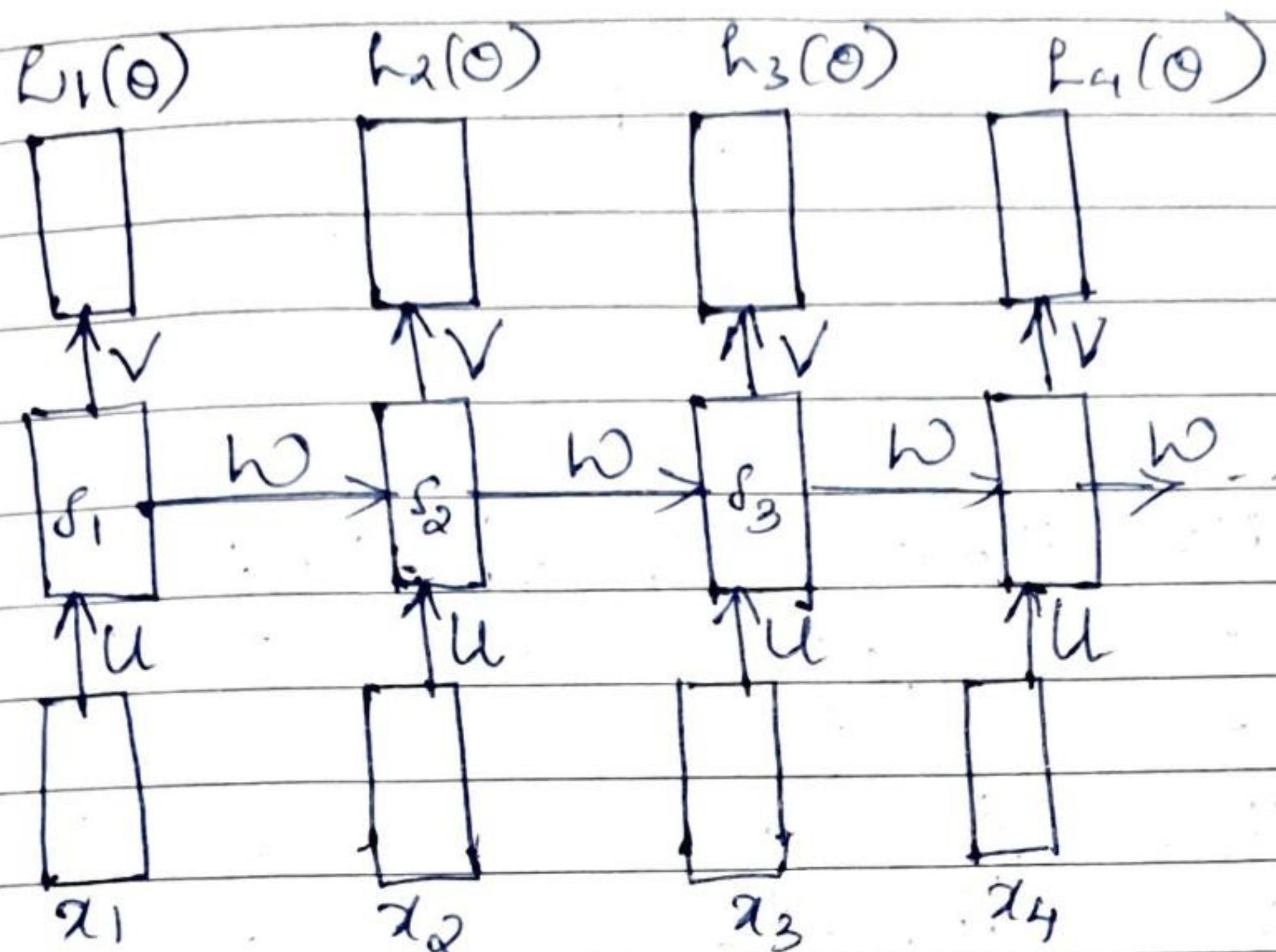
\downarrow
 u contributes to all s_k 's
 up to $k=t$
 s_k depends on s_{k-1}

$$\frac{\partial L(\theta)}{\partial s_{k-1}} = \frac{\partial L(\theta)}{\partial s_k} \frac{\partial s_k}{\partial s_{k-1}}$$

\therefore only need to compute
 $\frac{\partial s_k}{\partial s_{k-1}}$ for each decrement of k
 till the first term so,

which is a
 constant.

The Problem of Exploding and Vanishing Gradients



$$\frac{\partial L_t(t)}{\partial w} = \frac{\partial L_t(t)}{\partial s_t} \leq \prod_{k=1}^t \frac{\partial s_{t-k}}{\partial s_k} \frac{\partial s_k}{\partial w}$$

* We will now focus on $\frac{\partial s_t}{\partial s_k}$ and highlight an important problem in training RNN's using BPTT

$$\begin{aligned} \frac{\partial s_t}{\partial s_k} &= \frac{\partial s_t}{\partial s_{t-1}} \frac{\partial s_{t-1}}{\partial s_{t-2}} \cdots \frac{\partial s_{k+1}}{\partial s_k} \\ &= \prod_{j=k}^{t-1} \frac{\partial s_{j+1}}{\partial s_j} \end{aligned}$$

* Let's explore one such term in the product ie $\frac{ds_{j+1}}{ds_j}$

can also written as

$$\frac{ds_j}{ds_{j-1}}$$

$\rightarrow \boxed{\frac{ds_j}{ds_{j-1}}}$:-

a_j (activation) - $w s_{j-1} + b + l x_j$
(activation) $s_j = \sigma(a_j)$.

$$\therefore \frac{ds_j}{ds_{j-1}} = \frac{ds_j}{daj} \frac{daj}{ds_{j-1}}$$

current previous

ds_j - vector daj - matrix
 daj - vector ds_j - matrix

Let $a_j = [a_{j1}, a_{j2}, a_{j3}, \dots, a_{jd}]$ [d dimension]

$$s_j = [\sigma(a_{j1}), \sigma(a_{j2}), \sigma(a_{j3}), \dots, \sigma(a_{jd})]$$

$$\textcircled{1} \quad \frac{ds_j}{daj} = \begin{bmatrix} \frac{ds_{j1}}{daj_1} & \frac{ds_{j2}}{daj_1} & \frac{ds_{j3}}{daj_1} & \dots \\ \frac{ds_{j1}}{daj_2} & \frac{ds_{j2}}{daj_2} & \dots & \dots \\ \dots & \dots & \frac{ds_{j3}}{daj_3} & \dots \\ \dots & \dots & \dots & \frac{ds_{j4}}{daj_4} \end{bmatrix}$$

$$\begin{aligned}
 &= \begin{bmatrix} \sigma'(\alpha_{j1}) & 0 & 0 & 0 \\ 0 & \sigma'(\alpha_{j2}) & 0 & 0 \\ 0 & 0 & \sigma'(\alpha_{j3}) & 0 \\ 0 & 0 & 0 & \sigma'(\alpha_{j4}) \\ 0 & 0 & 0 & 0 \end{bmatrix} \\
 &= \text{diag}(\sigma'(\alpha_j))
 \end{aligned}$$

$$\frac{ds_j}{ds_{j-1}} = \frac{ds_j}{\sigma'(\alpha_j)} \cdot \frac{\sigma'(\alpha_j)}{ds_{j-1}} \rightarrow \begin{array}{l} \text{right} \\ \text{Matrix} \end{array}$$

$$= \text{diag}(\sigma'(\alpha_j)) \cdot W$$

* Let's find the magnitude of $\frac{ds_j}{ds_{j-1}}$

$$\left\| \frac{ds_j}{ds_{j-1}} \right\| = \left\| \text{diag}(\sigma'(\alpha_j)) \cdot W \right\|$$

$$\begin{aligned}
 \text{Rule: } \|c\| &= \|ab\| \\
 &\leq \|a\| \|b\|
 \end{aligned}$$

$$\leq \|\text{diag}(\sigma'(\alpha_j))\| \|W\|$$

→ Let's see find out $\|\text{diag}(\sigma'(\alpha_j))\|$

$\sigma'(\alpha_j)$ is a bounded fn (Sigmoid, tanh)
 $\sigma'(\alpha_j)$ is also bounded

$$\hookrightarrow \sigma(a_j) \leq \frac{1}{4} = \gamma \text{ (if } \sigma \text{ is logistic)}$$

$$\hookrightarrow \sigma'(a_j) \leq 1 = \gamma \text{ (if } \sigma \text{ is tanh)}$$

$\hookrightarrow \|\nabla\| \leftarrow \text{excegnt matrix is also bounded by some value } \gamma, \text{ then}$

$$\left\| \frac{\partial s_j}{\partial s_{j-1}} \right\| \leq \gamma$$

* Now consider the whole product

$$\left\| \frac{\partial s_t}{\partial s_k} \right\| = \left\| \prod_{j=k+1}^t \frac{\partial s_j}{\partial s_{j-1}} \right\|$$

$$\leq \prod_{j=k+1}^t \gamma \quad [(t-k) \text{ terms}]$$

$$\leq (\gamma)^{t-k}$$

* If γ and t greater than 1, exploding gradient problem arises.

1st case $\cancel{\textcircled{P}}$ If γ and t less than 1 or

$\gamma t < 1$, the gradient will vanish

2nd case \rightarrow If $\gamma t > 1$, the gradient will explode

* To avoid exploding / vanishing gradient problem is to use truncated Backpropagation

Instead of selecting $(t-k)$ terms

$$\downarrow \downarrow \\ T_1$$

We will restrict to some minimum previous states or terms.

↓ occurs vanishing gradient

* In the first case, the term goes to zero exponentially fast, which makes it difficult to learn some long period dependencies. This problem is called the vanishing gradient.

* In the second case, the term goes to infinity exponentially fast, and their value becomes a NaN due to the constant process. This problem is called the exploding gradient.

Truncated Backpropagation Through Time

* Backpropagation is the training algorithm used to update the weights in a neural net in order to minimize the error $J(w)$ between the expected op and the predicted op for a given ip.

* For sequence prediction problems where there is

In order dependence of observations, RNNs are used instead of classical feed forward neural nets.

- * RNNs are trained using a variation of the backpropagation algorithm called BPTT
- * BPTT unrolls the RNN and propagates the error backward over the entire sequence, one timestep at a time. The weights are then updated with the accumulated gradients.
- * BPTT can be slow to train RNNs on problems with very long input sequences.
- * In addition to speed, the accumulation of gradients over so many timesteps can result in a shrinking of values to zero, or a growth of values that eventually overflow, or explode.
- * A modification of BPTT is to limit the no. of timesteps used on the backward pass and in effect estimate the gradient used to update the weights rather than calculate it fully. This variation is called Truncated Backpropagation through time (TBPTT).

Training algorithm

* The TBPTT has two parameters

- k_1 : Defines the no. of timesteps shown to the NW on the forward pass
- k_2 : Defines the no. of timesteps to look at when estimating the gradient on the backward pass.

LSTM Networks (Long Short Term Memory)

* RNN's idea is that they might be able to connect previous information to the present task.

* Sometimes, we need only to look at recent information to perform the present task

only one context { eg:- if we are trying to predict the last word in "the clouds are in the sky"

* In this case, the gap b/w the relevant information and the place that it is needed is small

* But there are also cases where we need more context

eg:- "I grew up in France. I speak fluent French".

* Recent information suggests that the next word is probably the name of a language. But if we want to narrow down which language, we need the context of "Name" from further back.

* Here in this case, the gap b/w the relevant information and the point where it is needed to become is very large.

* As the gap grows, RNN's become unable to learn to connect the information.

* This problem of RNN is called short term Memory Problem.

* This problem is alleviated by the new now called LSTM.

↳ LSTM Network:-

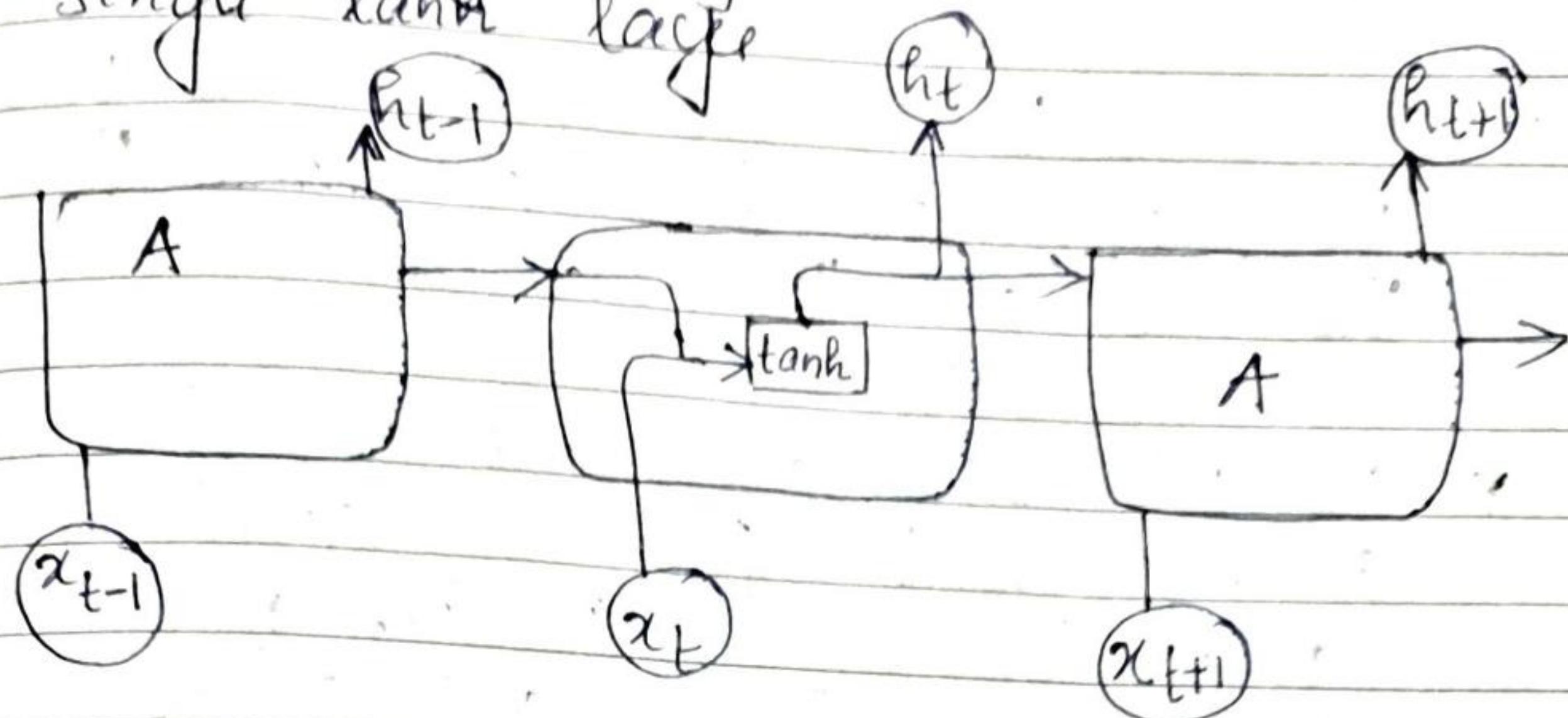
* Long short term memory n/w - are a special kind of RNN, capable of learning long-term dependencies.

* They were introduced by Hochreiter and Schmidhuber in 1997.

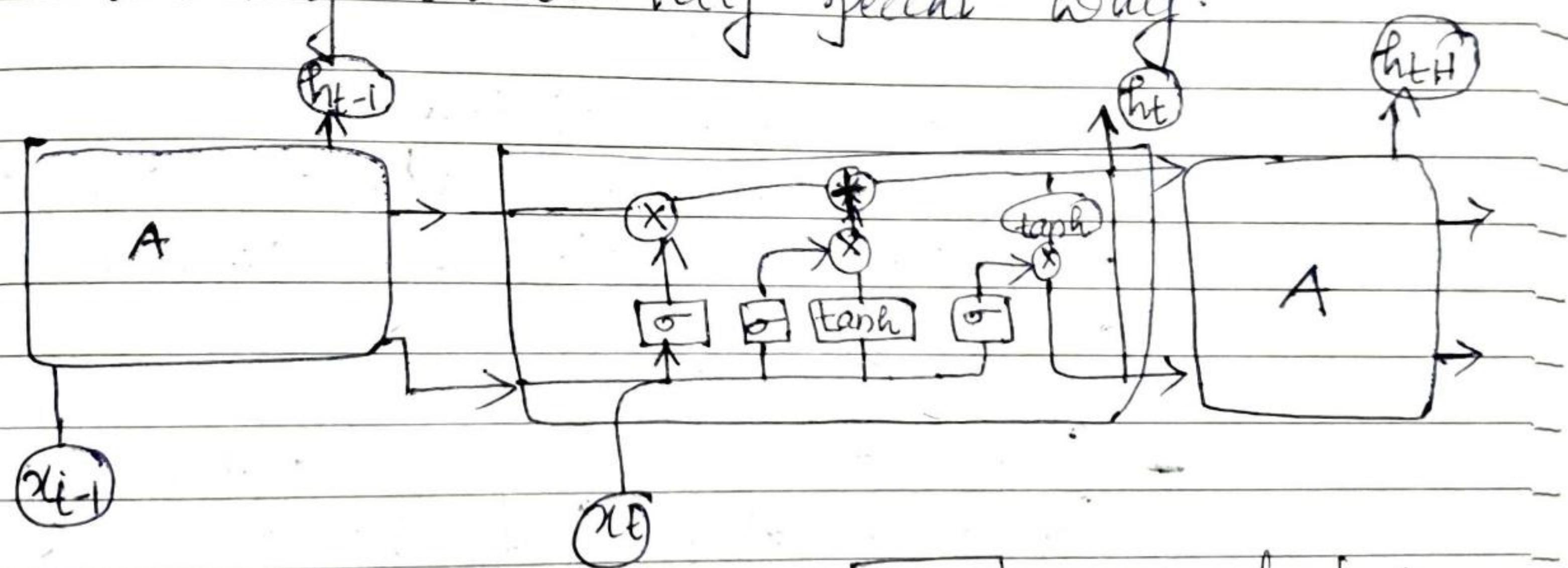
* All RNN have the form of a chain of repeating modules of neural n/w.

In RNN's, this repeating module will

have a
a single
vector
layer simple structure . Such a



* RNNs also have this chain like structure,
but the repeating module has a
different structure. Instead of having
a single neural net layer, there are four
interacting in a very special way.

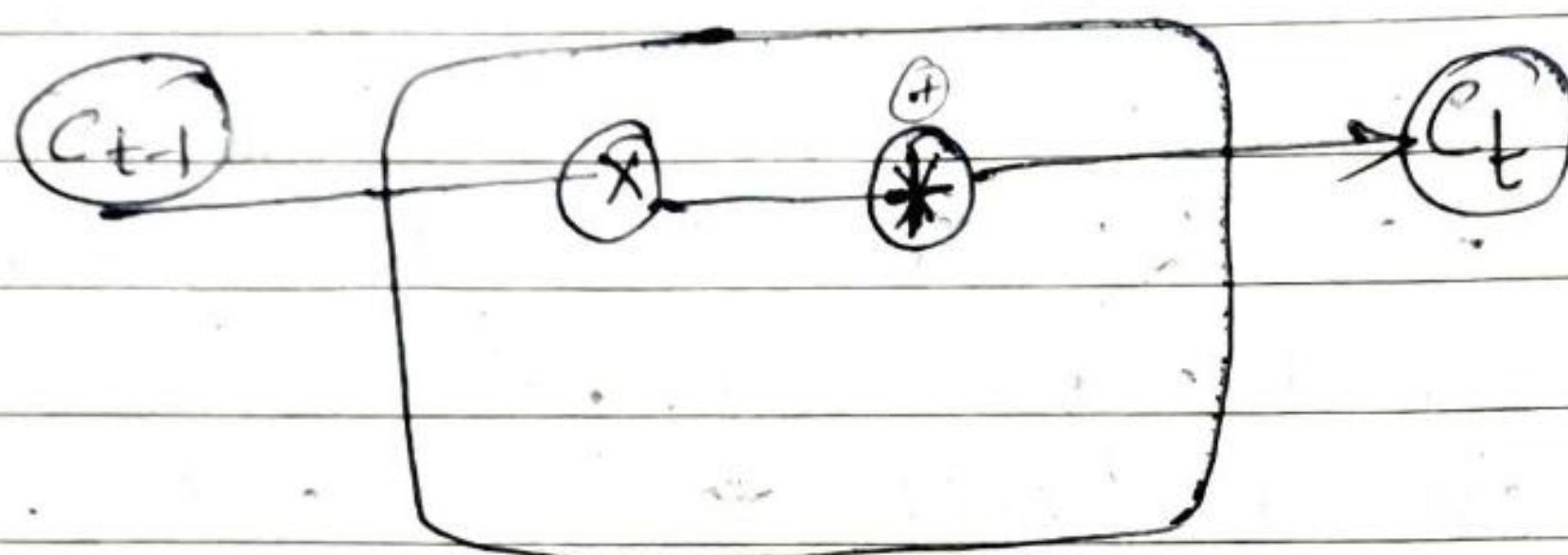


$\boxed{\quad}$ \Rightarrow Neural net layer

- $\circ \rightarrow$ pointwise operation
- \rightarrow Vector transfer
- \rightarrow concatenate
- \rightarrow Copy

Steps :- (cell state)

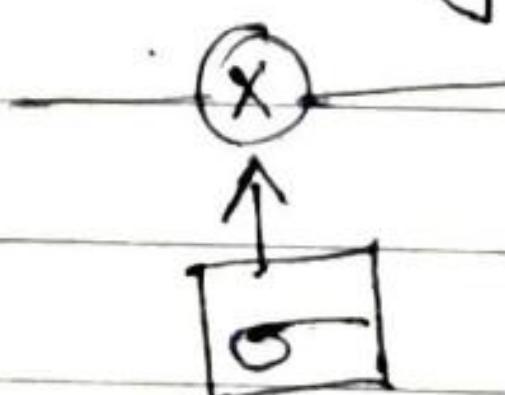
- * The key to LSTMs is the cell state, the horizontal line running through the top of the diagram.
- * The cell state is kind of like a conveyor belt. It runs straight down the entire chain, with only some minor linear interactions.



- * LSTMs have the ability to remove or add information to the cell state, carefully regulated by structures called gates.

→ Gates are a way to optionally let information through.

→ They are composed out of a sigmoid neural net layer and a point-wise multiplication operation.



→ Sigmoid layer ops number b/w 300 and one describing how much of each component should be let through.

→ A value of zero means "lets nothing through"
 → A value of One means "let everything through"

* LSTM has three of these gates to protect
 and control the cell state.

- ① Forget Gate layer
- ② Input Gate layers
- ③ Output Gate layer.

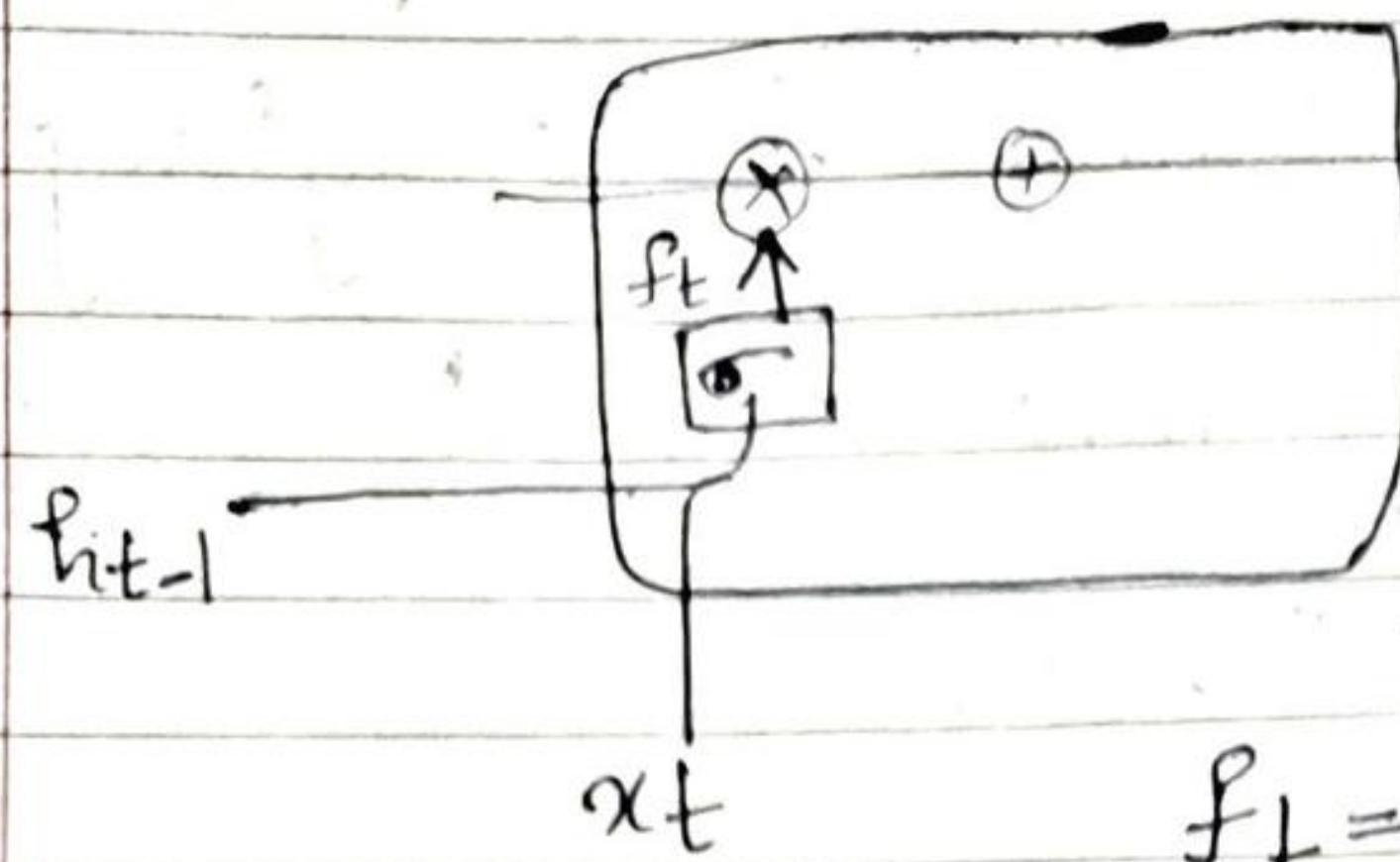
Forget Gate layer:

* Decide what information to throw away I forget from the cell state.
 This decision is made by a sigmoid layer called "forget Gate layer".

* It looks at h_{t-1} and x_t and
 outputs a no: 0 and 1 for
 each no: in the cell state C_{t-1} .

* e.g:- language model trying to predict
 the next word based on all
 the previous ones. In such a
 problem, the cell state might
 include the gender of the present
 state subject, so that correct
 pronouns can be used.

When we see a new subject,
 we want to forget the gender
 of the old subject.



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

Step 2:- Input Gate layer:-

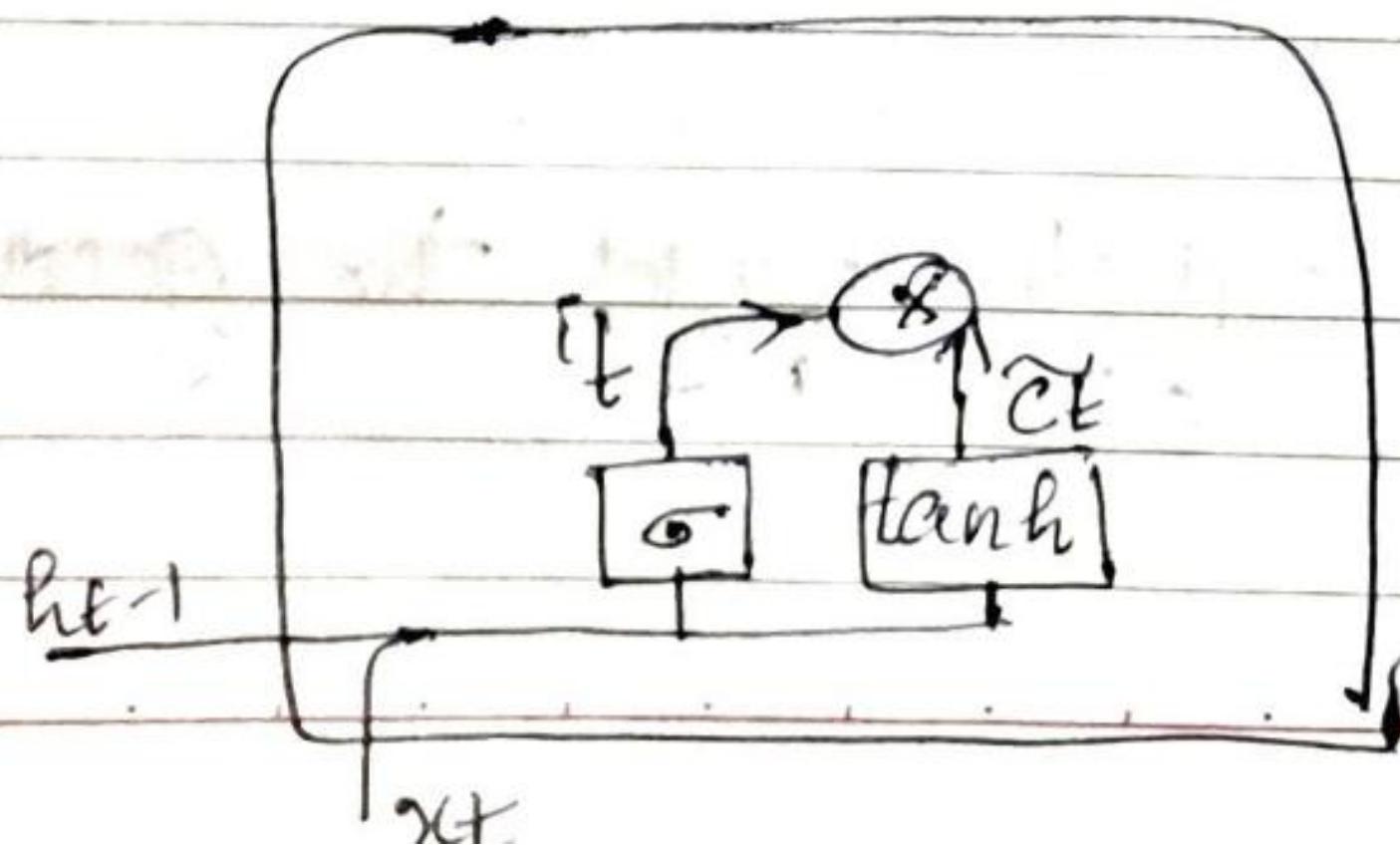
- * This step decides what new information is going to store in the cell state.

- * This has two parts

- ① \hookrightarrow Sigmoid layer called the forget gate layer decides which values will get updated

- ② \hookrightarrow A 'tanh' layer creates a vector of new candidate value \tilde{c}_t , that could be added to the state.

- * eg:- In our language model, want to add the gender of the new subject to the cell to replace the old one we are forgetting



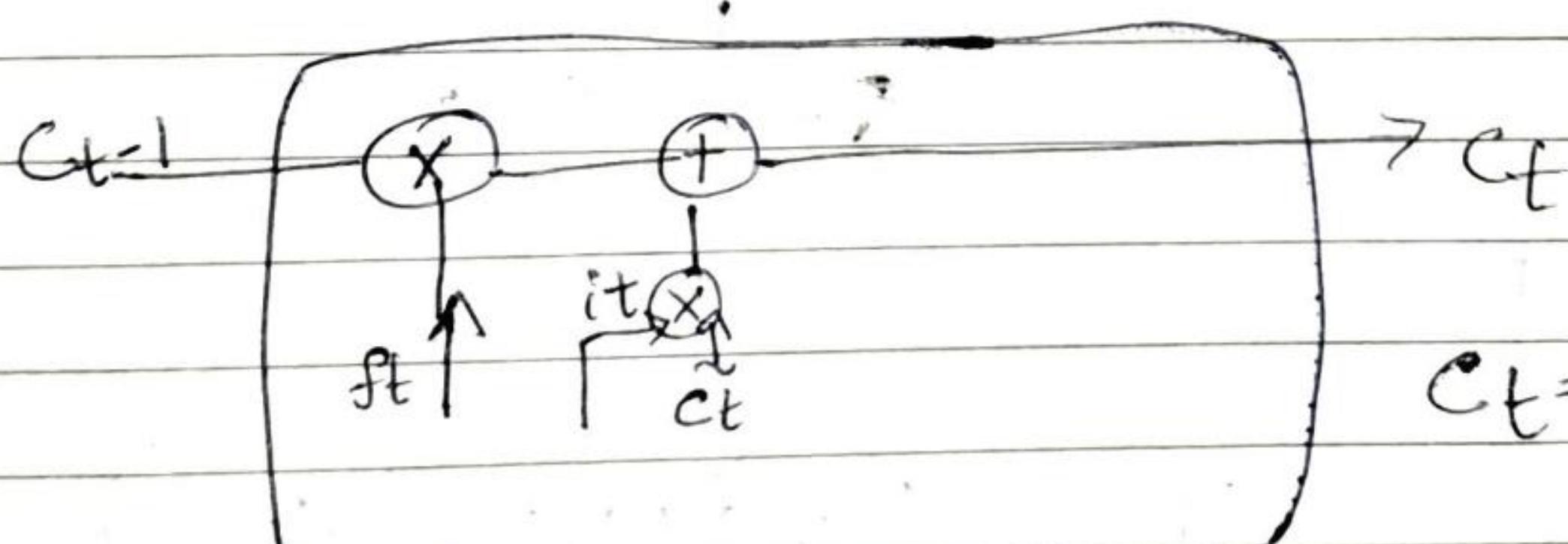
$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$\tilde{c}_t = \tanh(W_c \cdot [h_{t-1}, x_t] + b_c)$$

Steps:- Output Gate layer

- * Update the old cell state C_{t-1} , into the new cell state C_t . Previous steps already decided what to do, in this step, we just need to actually do it.
- * Multiply the old state by f_t , forgetting the things we decided to forget earlier. Then we add $i_t * \tilde{c}_t$.
 - new candidate value to be updated

eg:- In the language model, this is where we had actually drop the information about the object's gender and add the new information.



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

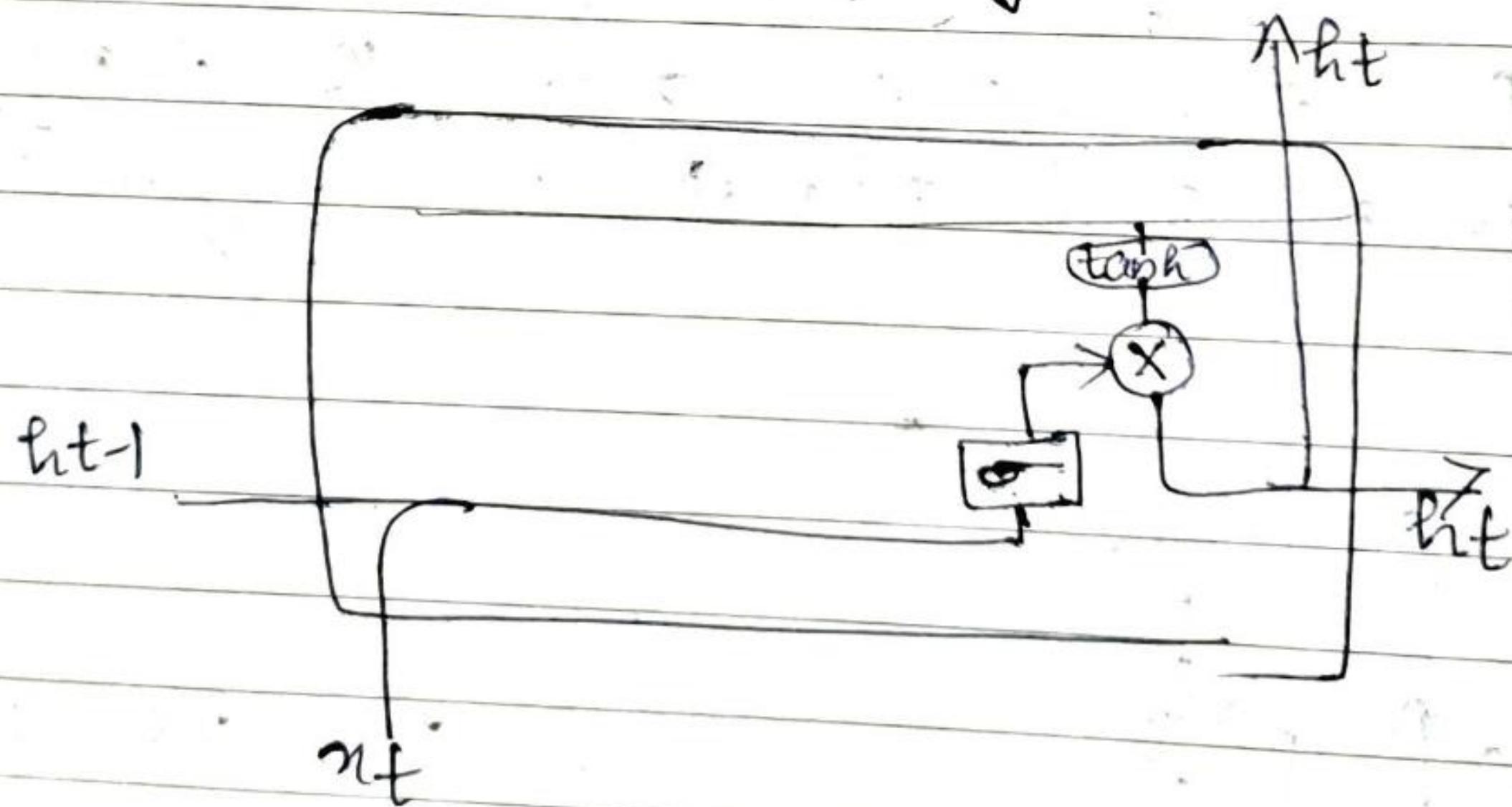
Step 3: Output Gate layer.

* decides what we are going to do

① * ReLU layer which decides what parts of the cell state we are going to output.

② * Then, put the cell state through tanh (ReLU) and multiply it by the o/p of the sigmoid gate, so that we only o/p the parts we decided to

e.g.: In the language model,



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

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