

10/12/18

### SRN

#### Recurrent Neural Networks

ANN's

can't deal with temporal data

ANN's

lack memory

ANN's

have fixed Architecture

RNN's are biologically realistic (recurrent connectivity in the visual cortex of brain)

#### typical applications

Image captioning

(one dp - multiple dps)

Sentiment analysis

(multiple ips - one dps)

Machine translation

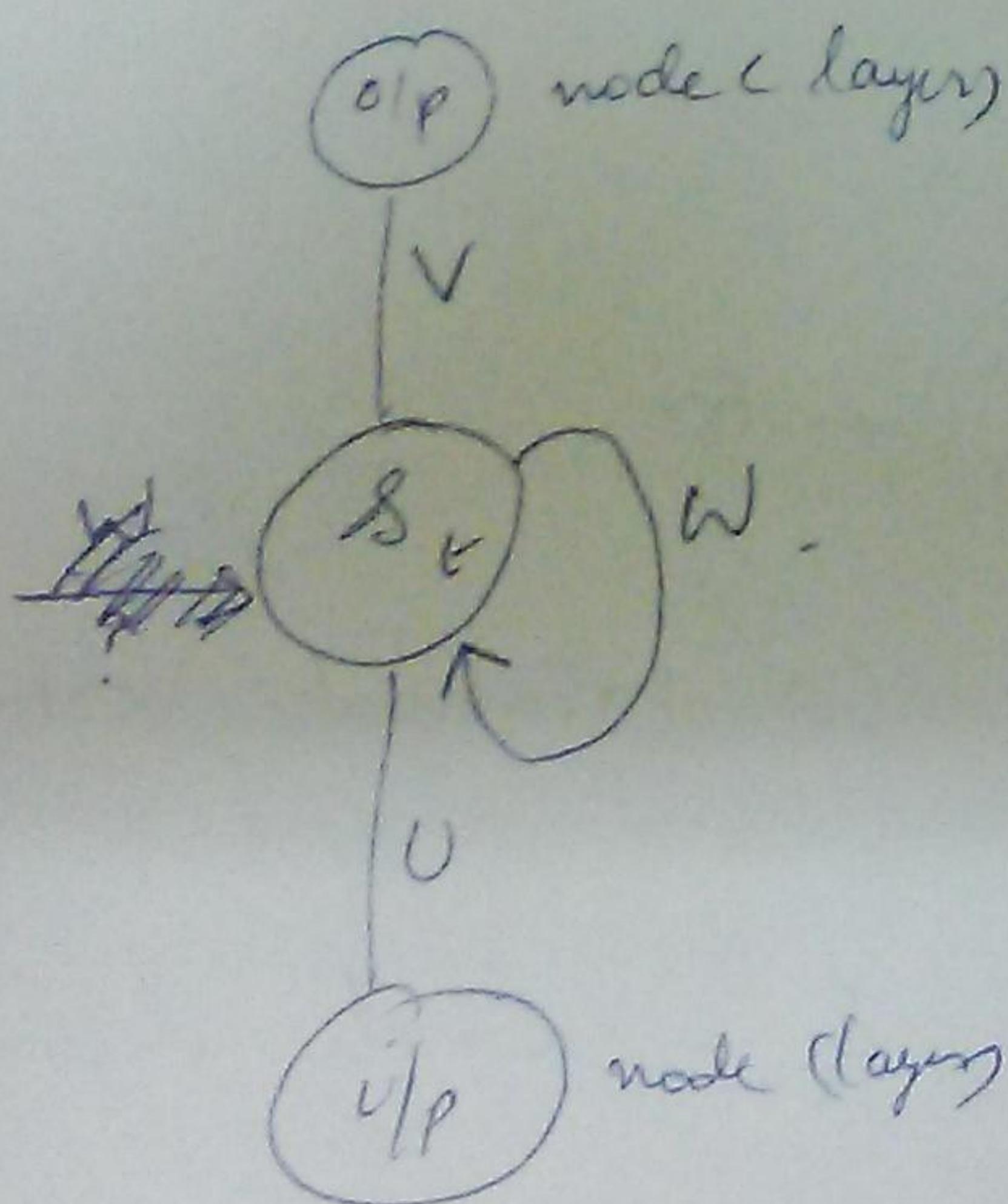
(multiple ips - multiple dps)

Video captioning

(multiple ips - multiple dps)

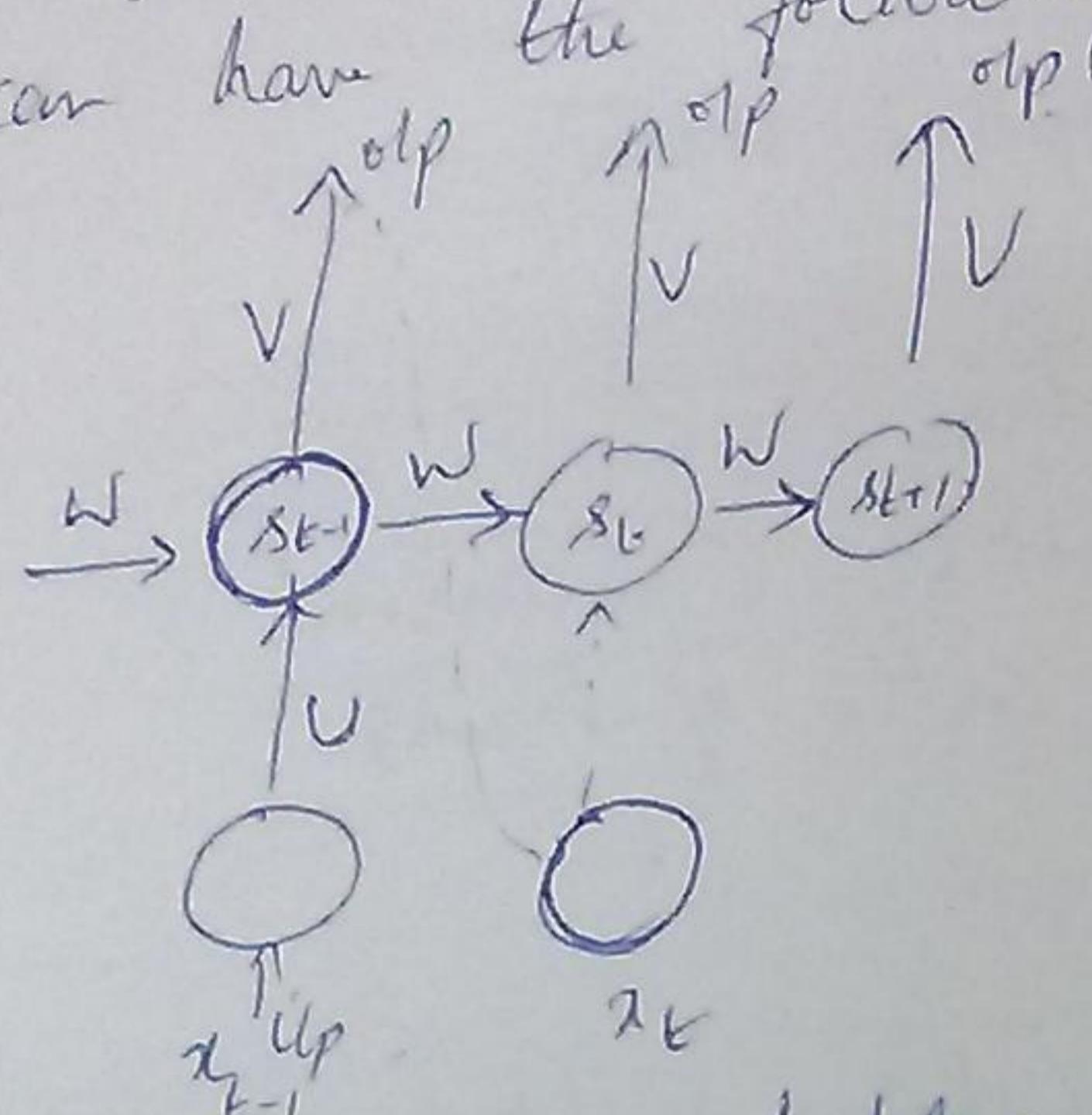
Speech analysis.

#### typical RNN



$s_t$  = state (hidden) of the Recurrent net  
(actually a layer). The state  $s_t$  at  
any point of time it's  $s_t$ . (a

if we unfold or unroll the RNN over time we can have the following



(one to many.

Ex- Image captioning  
or even word / sentence prediction)

There are 3 hidden states  $s_{t-1}, s_t, s_{t+1}$  (This can go on and on)

$$y_t = V s_t$$

$$y_{t-1} = V s_{t-1}$$

$$y_{t+1} = V s_{t+1}$$

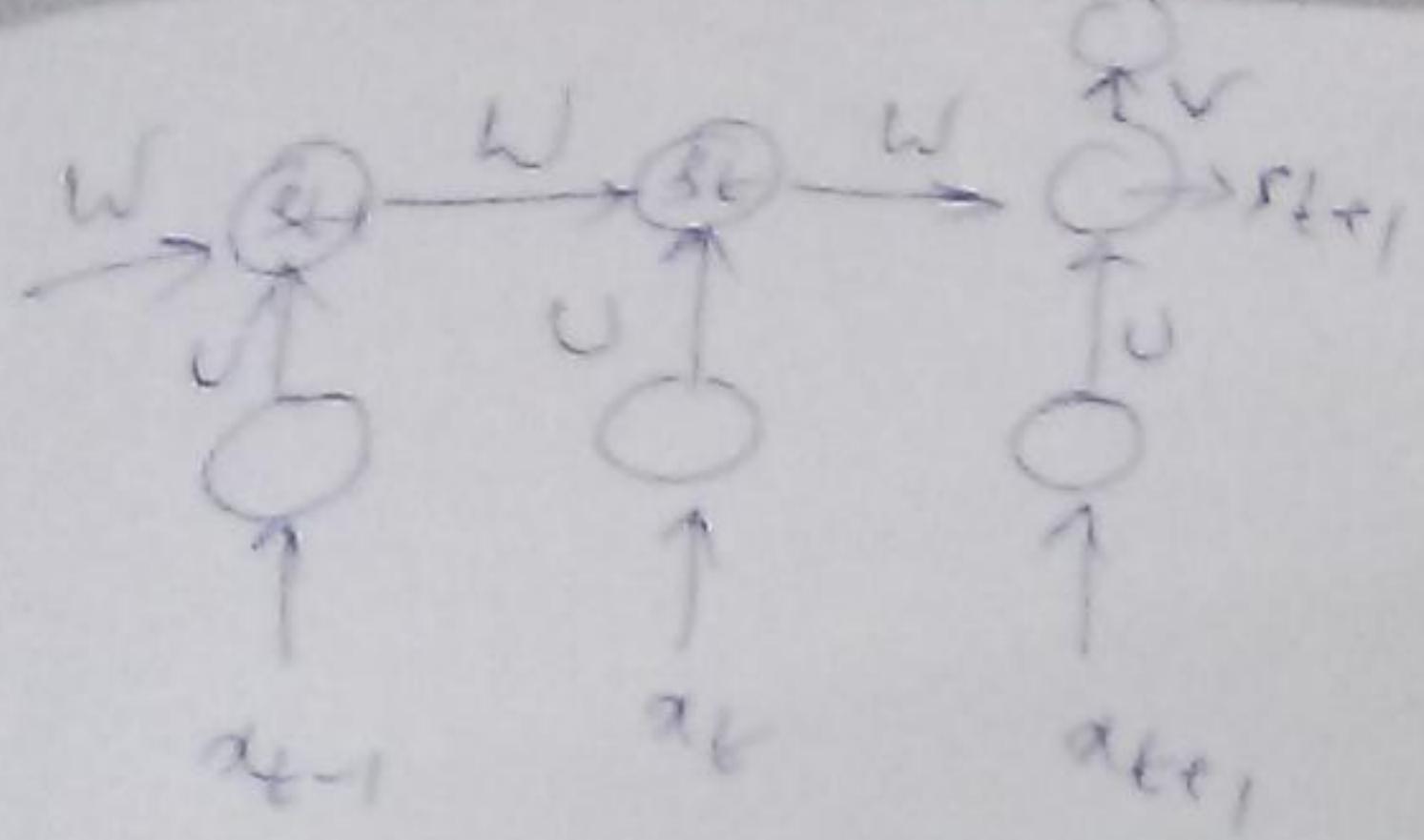
Tanh is a common function:

$$s_t = \tanh(W s_{t-1} + U x_t).$$

The op can have a softmax layer too

so  $y_t = \text{softmax}(V s_t)$ . (a probabilistic vector)

[or at  $t+1$ ,  $y_{t+1}$  depends on  $s_{t+1}$ , which depends on  $s_t$ , which in turn depends on  $s_{t-1}$ , and so on. like eqn ①]

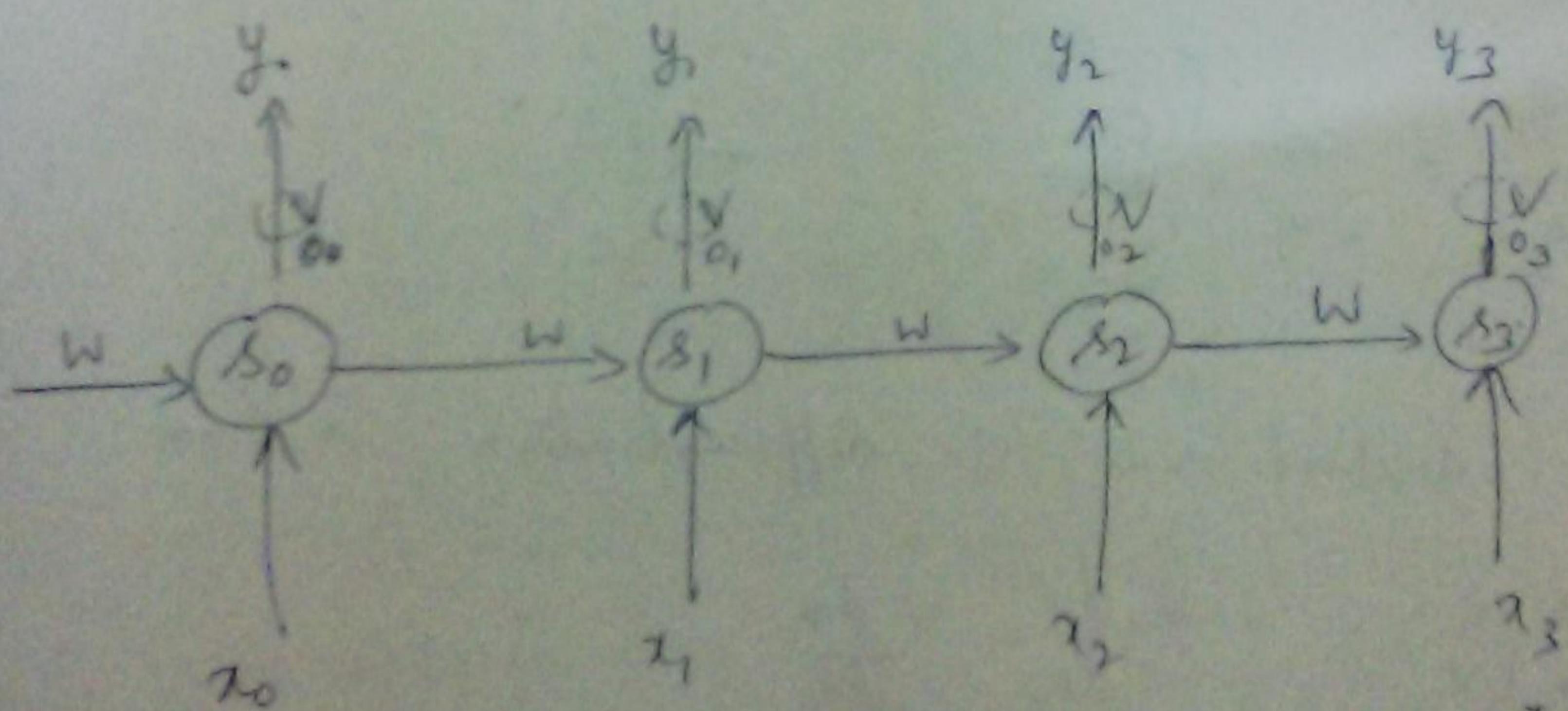


Coming to me.  
Ex. Sentiment classifier)

After all inputs are fed in, the classification is dp. of softmax ( $v s_{t+1}$ ). Here  $s_{t+1}$  has been computed taking all previous hidden states and previous ips.

Makes sense also, because for a given sentence, the sentiment can be guessed only after ALL words of the sentence have been processed.

A typical RNN unfolded over time



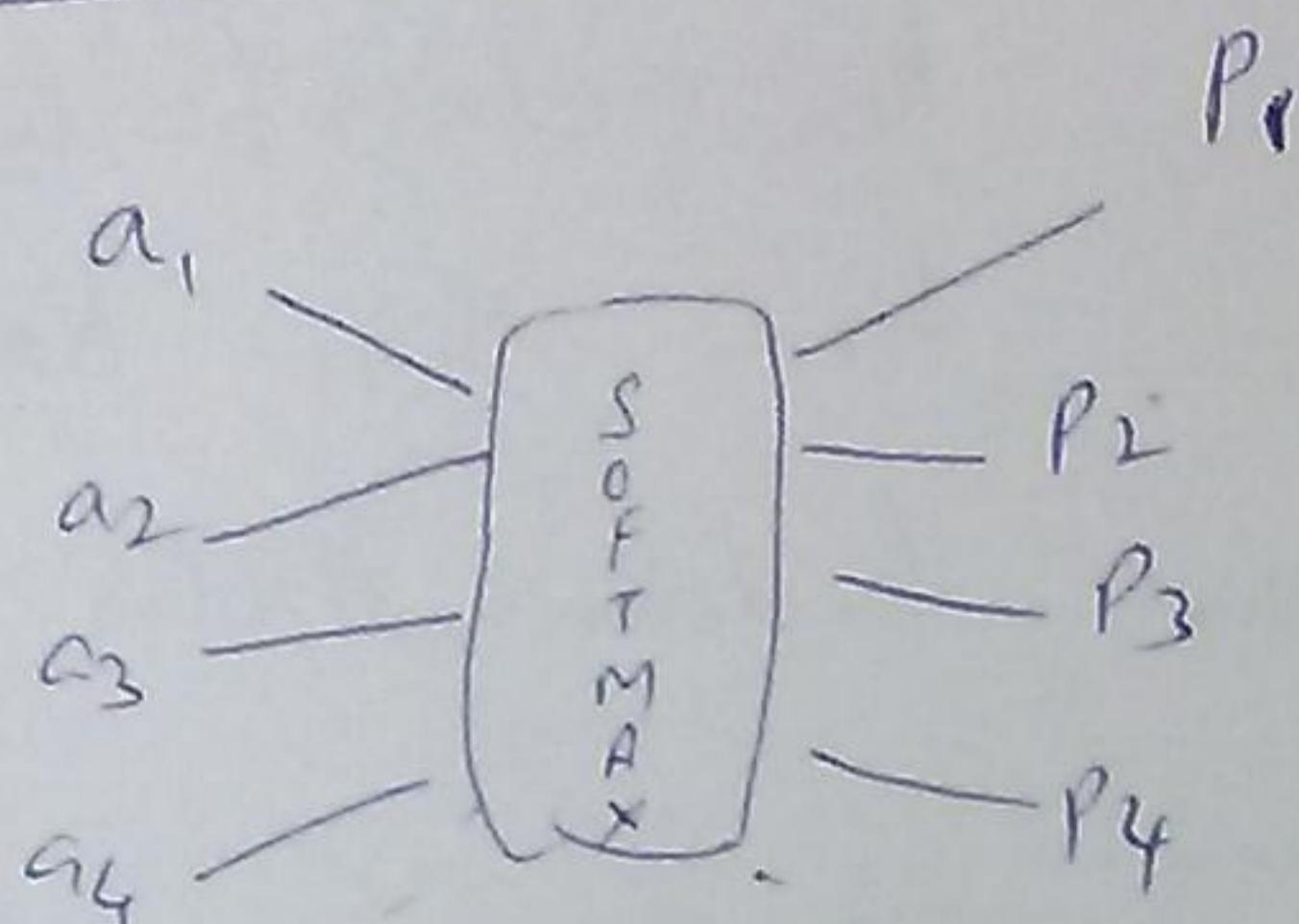
0, 1, 2, 3 are instances of time. So in reality there is only ONE RNN neuron.

$$O_1 = VS_1$$

$$O_2 = VS_2$$

$$y_1 = \text{softmax}(O_1)$$

Softmax detour



$$P_i = \frac{e^{a_i}}{\sum_{j=1}^n e^{a_j}}$$

If we have to find the differential of softmax  
we have to resort to vector calculus)

$$\frac{\partial P_i}{\partial a_j} = \frac{\partial}{\partial a_j} \left\{ \frac{e^{a_i}}{\sum_{j=1}^n e^{a_j}} \right\} \quad -\textcircled{2}$$

From quotient rule of differentiation, we have

$$\frac{d}{dx} \left( \frac{u}{v} \right) = \frac{v \cdot du - u \cdot dv}{v^2}$$

Applying on  $\textcircled{2}$ , we have

$$\frac{\partial p_i}{\partial a_j} = \frac{\sum_{j=1}^n e^{a_j} \cdot \frac{\partial}{\partial a_j} (e^{a_i}) - e^{a_i} \frac{\partial}{\partial a_j} (\sum_{j=1}^n e^{a_j})}{\sum_{j=1}^n e^{a_j} \cdot \sum_{j=1}^n e^{a_j}}$$

now 2 cases. First,  $i \neq j$ .

$$\frac{\partial p_i}{\partial a_j} = \frac{0 - e^{a_i} \cdot e^{a_j}}{\sum_{j=1}^n e^{a_j} \cdot \sum_{j=1}^n e^{a_j}} = - \frac{e^{a_i}}{\sum e^{a_j}} \cdot \frac{e^{a_j}}{\sum e^{a_j}}$$

$$\frac{\partial p_i}{\partial a_j} \text{ when } (i \neq j) = -p_i p_j \quad - \textcircled{3}$$

2nd case,  $i=j$ ,

$$\frac{\partial p_i}{\partial a_j} = \frac{\sum_{j=1}^n e^{a_j} \cdot e^{a_j} - e^{a_j} \cdot e^{a_j}}{\sum_{j=1}^n e^{a_j} \cdot \sum_{j=1}^n e^{a_j}}$$

$$= \frac{e^{a_j} \sum_{j=1}^n e^{a_j}}{\sum e^{a_j} \cdot \sum e^{a_j}} - \frac{e^{a_j} \cdot e^{a_j}}{\sum e^{a_j} \cdot \sum e^{a_j}}$$

$$\frac{\partial p_i}{\partial a_j} (i=j) = p_j - p_j^2 \\ = p_j(1-p_j) \quad - \textcircled{4}$$

Back on our RNN; now

$$E_t(y_t, \hat{y}_t) \rightarrow \text{cross entropy. (at time } t\text{)}$$

$$= -y_t \log(\hat{y}_t)$$

$$E(y, \hat{y}) \quad (\text{sum of all errors from various times})$$

$$E(y, \hat{y}) = -\sum_t y_t \log(\hat{y}_t)$$

No BPPT (w.r.t V)

$$\frac{\partial E_3}{\partial V} = \frac{\partial E_3}{\partial \hat{y}_3} \cdot \frac{\partial \hat{y}_3}{\partial o_3} \cdot \frac{\partial o_3}{\partial V}$$

$$\text{Now } \underline{o_3 = s_3 \cdot V}, \text{ so } \frac{\partial o_3}{\partial V} = s_3.$$

$$\frac{\partial E_3}{\partial V}, \underline{(\hat{y}_3 - y_3) \cdot s_3.} \quad (\text{long proof} \rightarrow \text{not important})$$

Now coming to the impact of W (and also  $\theta$ )

W impacts  $E_3$  via  $s_3$ . But the same W impacts  $E_3$  thru  $s_2$  (which impacts  $s_3$ ) and W also impacts  $E_3$  through  $s_1$  which impacts  $s_2$  which impacts  $s_3$

$$s_0 \quad \frac{\partial E_3}{\partial W} = \sum_{k=0}^3 \frac{\partial E_3}{\partial \hat{y}_3} \cdot \frac{\partial \hat{y}_3}{\partial o_3} \cdot \frac{\partial o_3}{\partial s_k} \cdot \frac{\partial s_k}{\partial W}$$

$$\frac{\partial s_t}{\partial W} = \underline{\frac{\partial s_t}{\partial s_{t-1}}} \cdot \underline{\frac{\partial s_{t-1}}{\partial s_{t-2}}} \cdots \underline{\frac{\partial s_0}{\partial W}}$$

note that  $s_t = \tanh(Ws_{t-1} + ux_t)$

for a tanh function,  $\tanh(x) \quad \frac{\partial}{\partial x} \tanh(x) = \frac{1 - \tanh^2(x)}{1 + \tanh^2(x)}$

so if  $z_t = ws_{t-1} + ux_t$ , then  $\frac{\partial s_t}{\partial z_t} = \frac{1 - s_t^2}{1 + s_t^2}$

$(s_t \tanh(z_t))$

and we have  $\frac{\partial z_t}{\partial W} = s_{t-1}$  (provided  $s_{t-1}$  is not a function of  $s_{t-1}$ )

Vanishing Gradient / Exploding gradient

Due to small derivatives and numerous multiplications, derivatives will be near 0 and updates are impossible.

LSTMs to the Rescue:

Long-short-term Memory:

Recall that in the Vanilla RNN, each hidden neuron is actually a hidden layer (when unfolded)

LSTM:

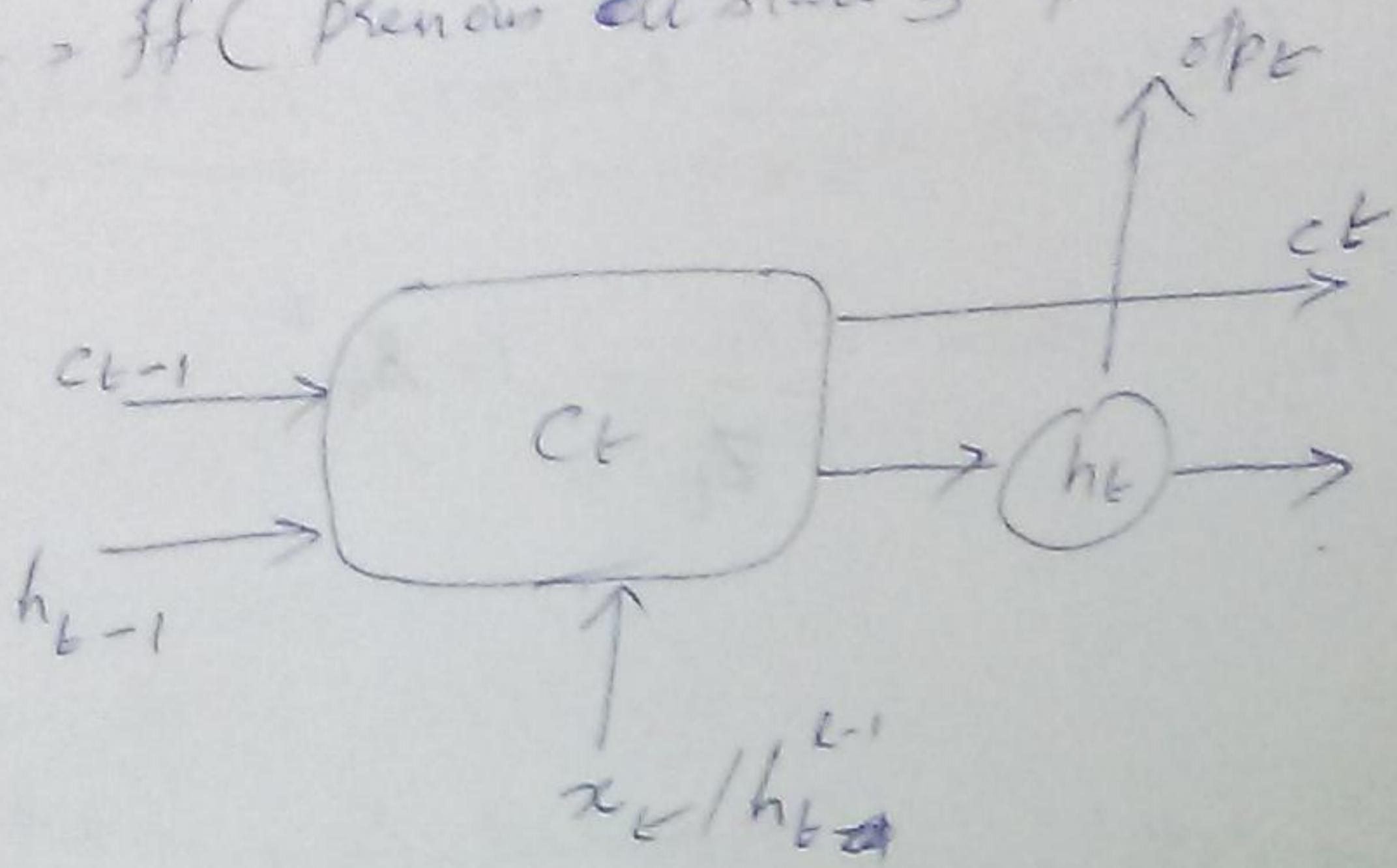
Hidden States still exist,

$h_t$  or  $s_t$  is now a function of something called as the "cell state".

so  $h_t = f(a)$ .

$$c_t = ff(a_{t-1}, x_t)$$

$a_t = ff(\text{previous cell state } s \text{ & previous hidden states})$   
in depth  $+ x_t$ )

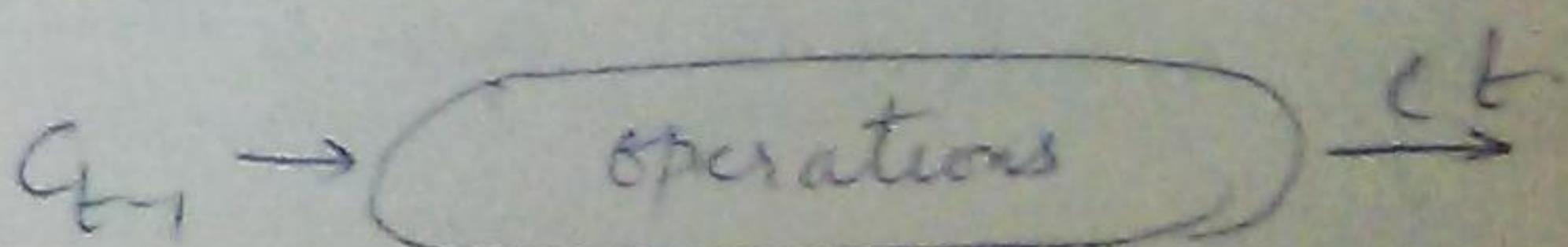


( $t-1$ ) previous depth  
ignore for now

Consider only one hidden neuron.

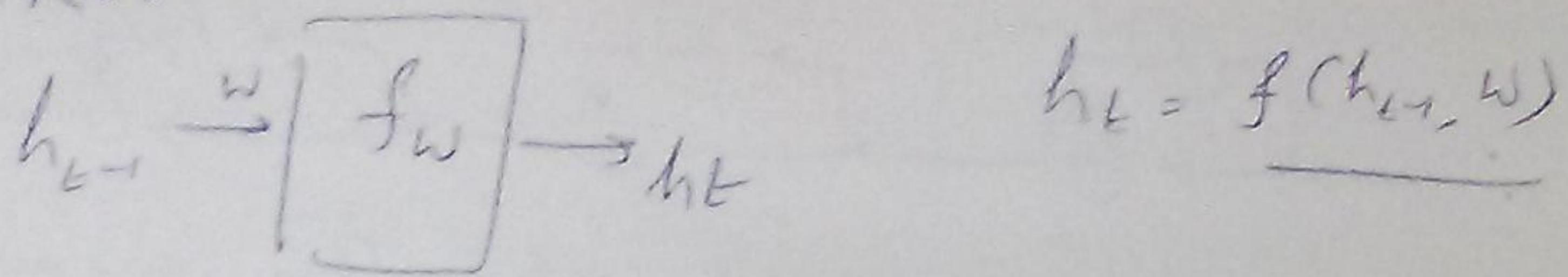
An LSTM cell can be compared with a conveyor belt making a burger (say) -  
there are workers at strategic points on the belt. As the burger moves over the belt, the workers enhance or enrich the burger.

cell state  $\Rightarrow$  conveyor belt



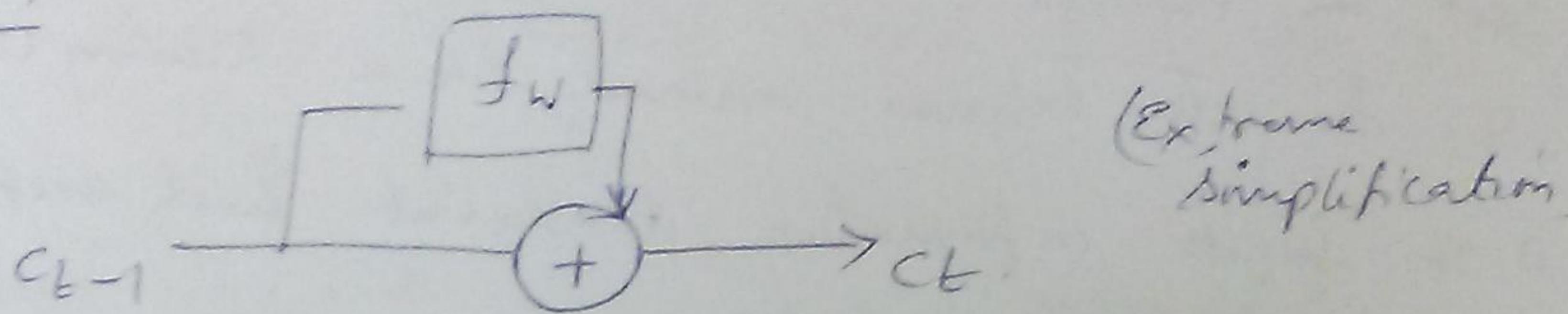
So as information passes through a cell, it may be completely modified or partly modified or may be passed on as is.

Traditional RNN



(derivatives multiply and hence can vanish)

LSTM



$$c_{t-1} \oplus \text{some vector} = c_t$$

$f_w$  = some function  
(not the same as  
vanilla RNN)

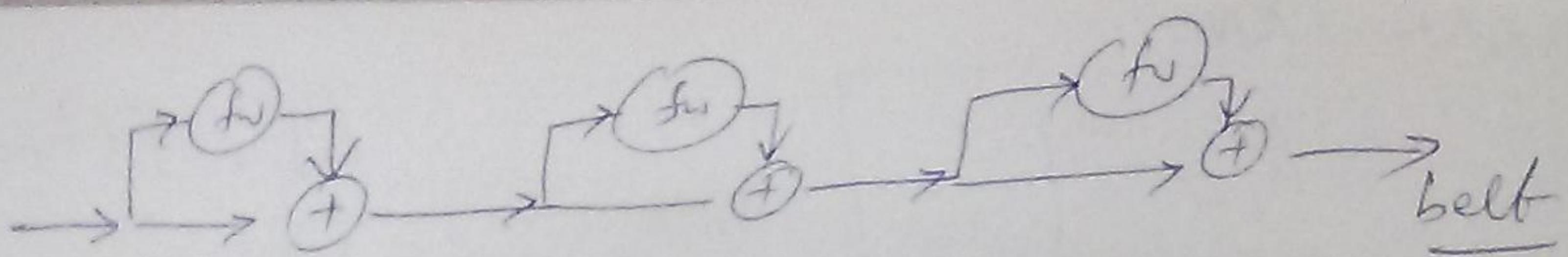
$$c_t = c_{t-1} + f_w(c_{t-1})$$

This can be expanded as

$$c_t = c_0 + f_w(c_1) + f_w(f_w(c_2)) + \dots + f_w(f_w(\dots f_w(c_{t-1}))) - \textcircled{I}$$

Why is this better?

Because when we differentiate, we differentiate sum of terms. So the differentials will now add up (not multiply). So we have a gradient (gradient super-highway). Even if some gradients go zero, it may not result in a zero total gradient.



So, far we have added info. This means it will keep growing and at some point of time we can get an explosion'

- So, the cell behaves instead like a memory unit:
- 1) it 'forgets' information it doesn't find useful from the previous cell-state
  - 2) it writes new information it does find useful from the current input
  - 3) it "reads out" a part of its information to help the computation of the hidden state  
(Reset - write - read)

The summation part  $\rightarrow$  "writing" info.  $\rightarrow$  Add info to the cell

Resetting part is Multiplicative and occurs before writing to memory. The multiplying vector is full of 1's and 0's and we perform an index-wise multiplication on the previous cell-state to eliminate the useless information.

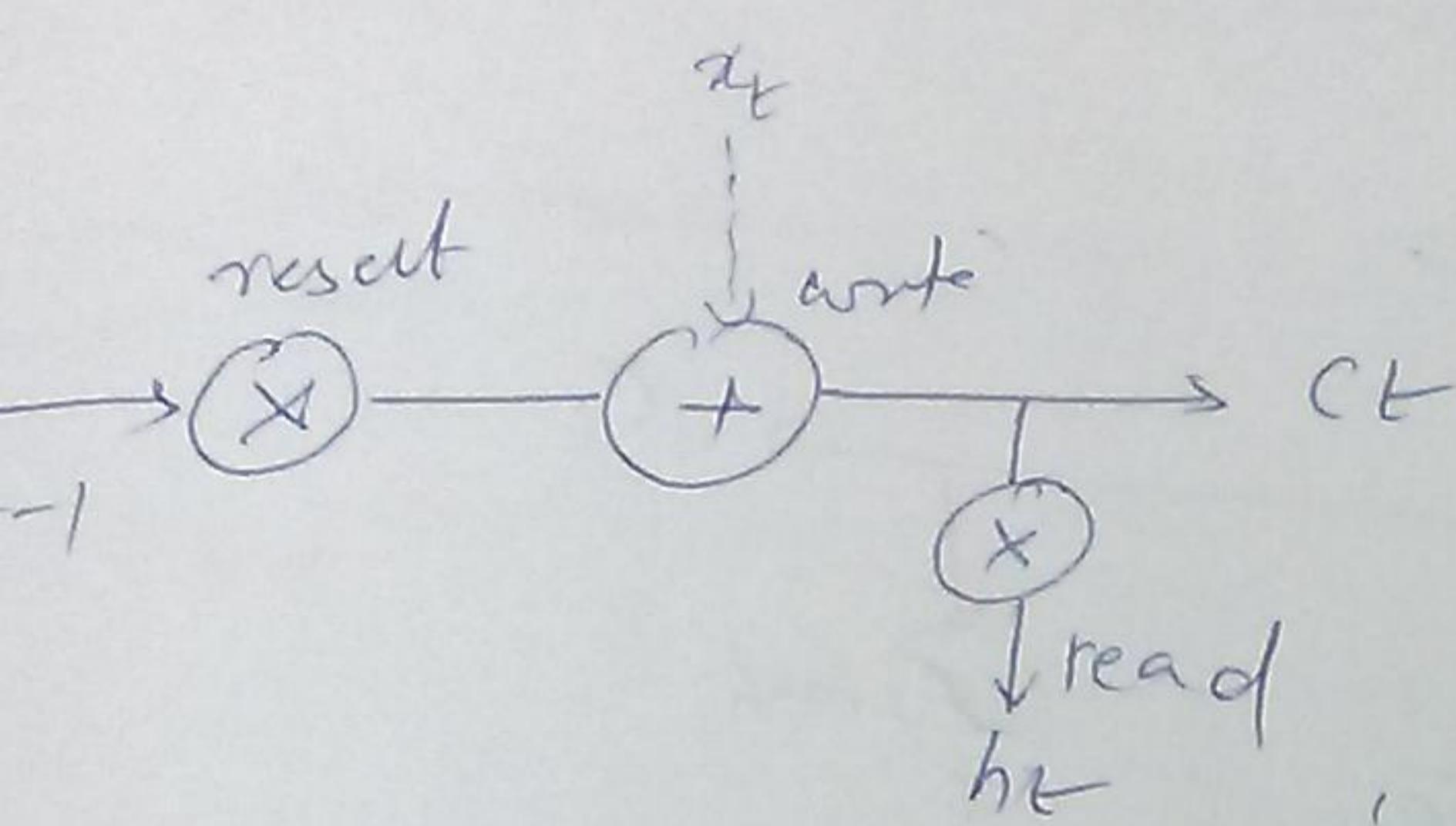
Reading Part is also multiplicative but doesn't modify the cell-state - instead it modifies info flowing into the hidden state

and thereby influences the hidden state. we do element-wise multiplication

$$\begin{bmatrix} a \\ b \end{bmatrix} \odot \begin{bmatrix} c \\ d \end{bmatrix} = \begin{bmatrix} ac \\ bd \end{bmatrix}$$

if  $a=0$ , info in  $c$  is lost.

so a slightly enhanced version of our cell will be



This concept is called 'Gating'

- Gates

1. Forget Gate (actually 'remember gate')

A sigmoid function is used to compute the 'Forget Gate'. Its values are between 0 & 1. If we get '0' we forget. If we get '1', we remember fully.

Write Gate ( $G_w$ )  $\rightarrow$  partly responsible for 'write' process.

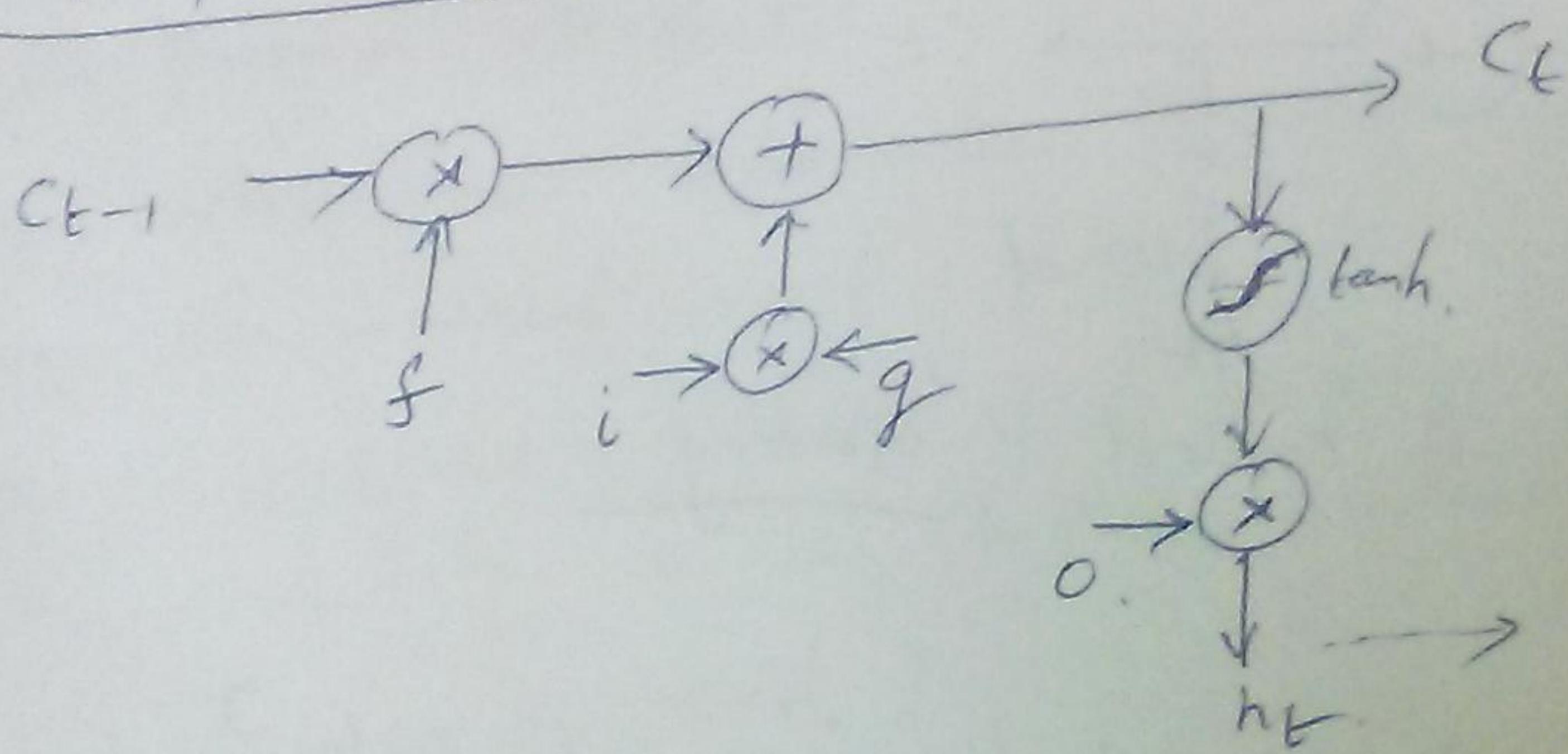
Stores a value between -1 and 1 and represents how much we want to add (input to cell state) by computing with a 'tanh' function.

Computed with an input gate ( $i$ ) - The other gate responsible for the 'write' process. How much of ' $g$ ' is allowed in?

Incrementing or decrementing counter. Computed with a sigmoid function.

4) Output Gate: Passed through sigmoid, and is a number between 0 and 1.  $\rightarrow$  Modulates or modulates which aspects of the cell-state can be drawn by the hidden state for 'Read' operation.  
 $(\tanh)$  also

Slightly enhanced version



$$c_t = f \odot c_{t-1} + i \odot g \\ h_t = o \odot \tanh(c_t)$$

(no activation function.  $c_t$  stays stable by "forgetting" and "writing")

How are Gates calculated?

They have their own weights and are functions of  $(h_{t-1}, x_t)$

$$f_t = \text{sigmoid}(W_{xf} x_t + W_{hf} h_{t-1})$$

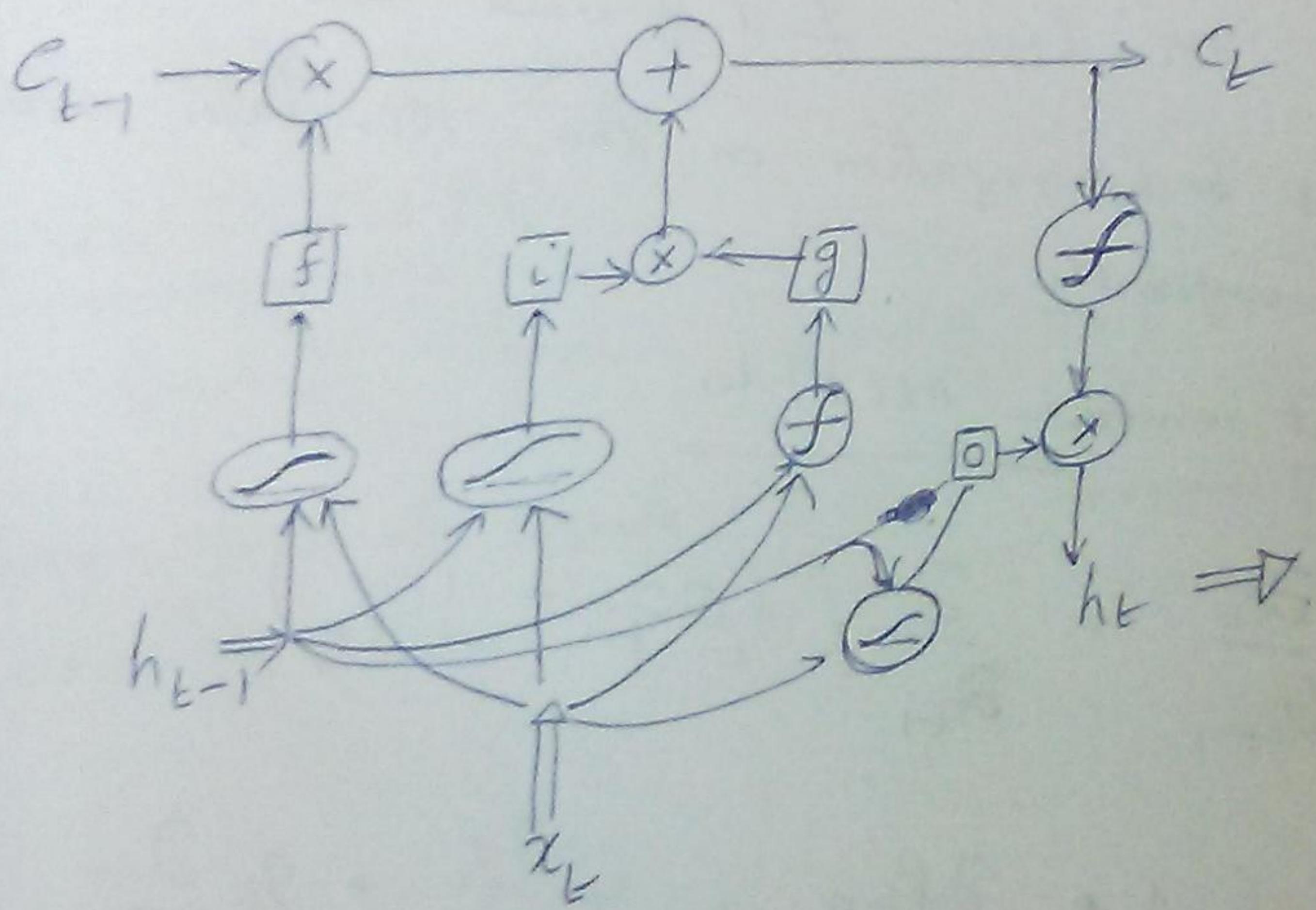
$$g_t = \tanh(W_{xg} x_t + W_{hg} h_{t-1})$$

$$i_t = \text{sigmoid}(W_{xi} x_t + W_{hi} h_{t-1})$$

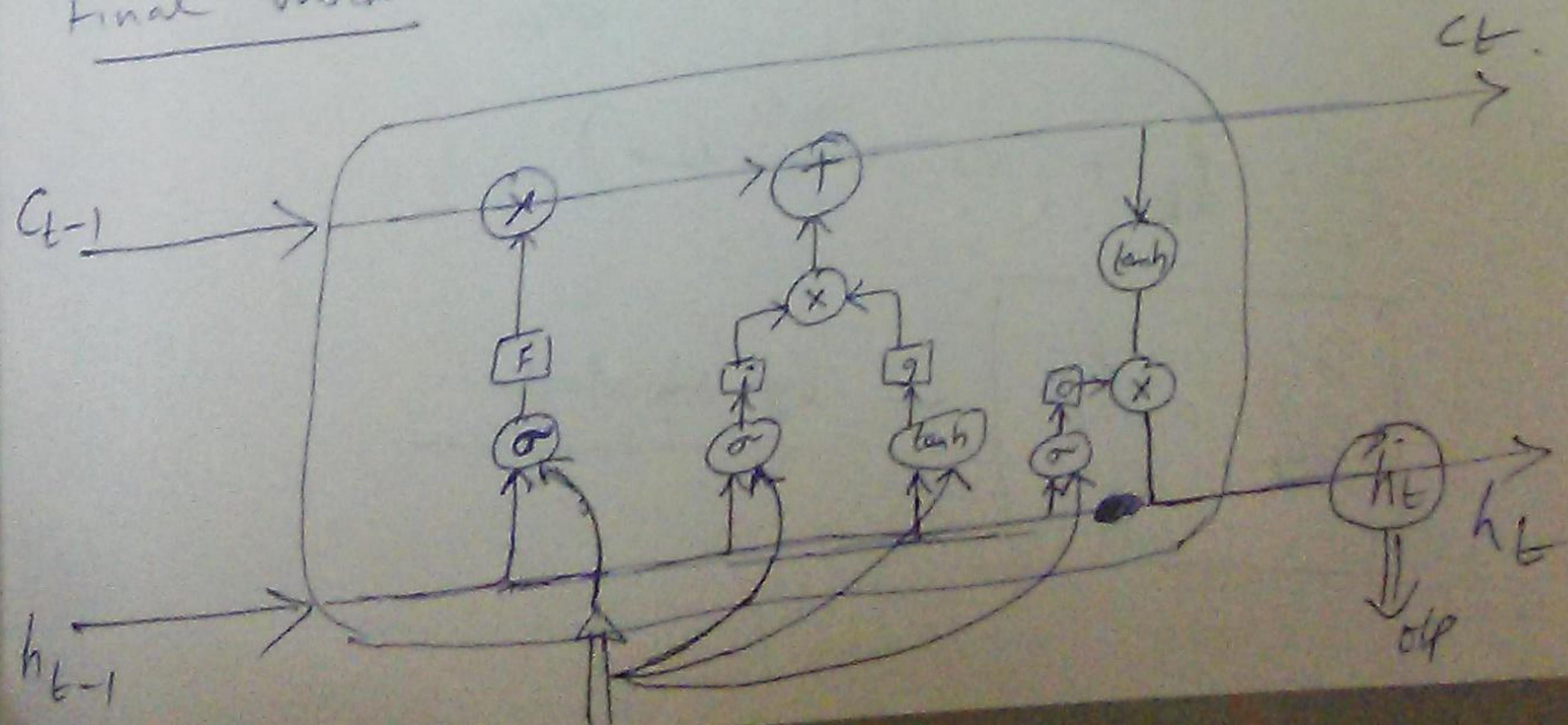
$$o_t = \text{sigmoid}(W_{xo} x_t + W_{ho} h_{t-1})$$

The whole thing is end-end differentiable.

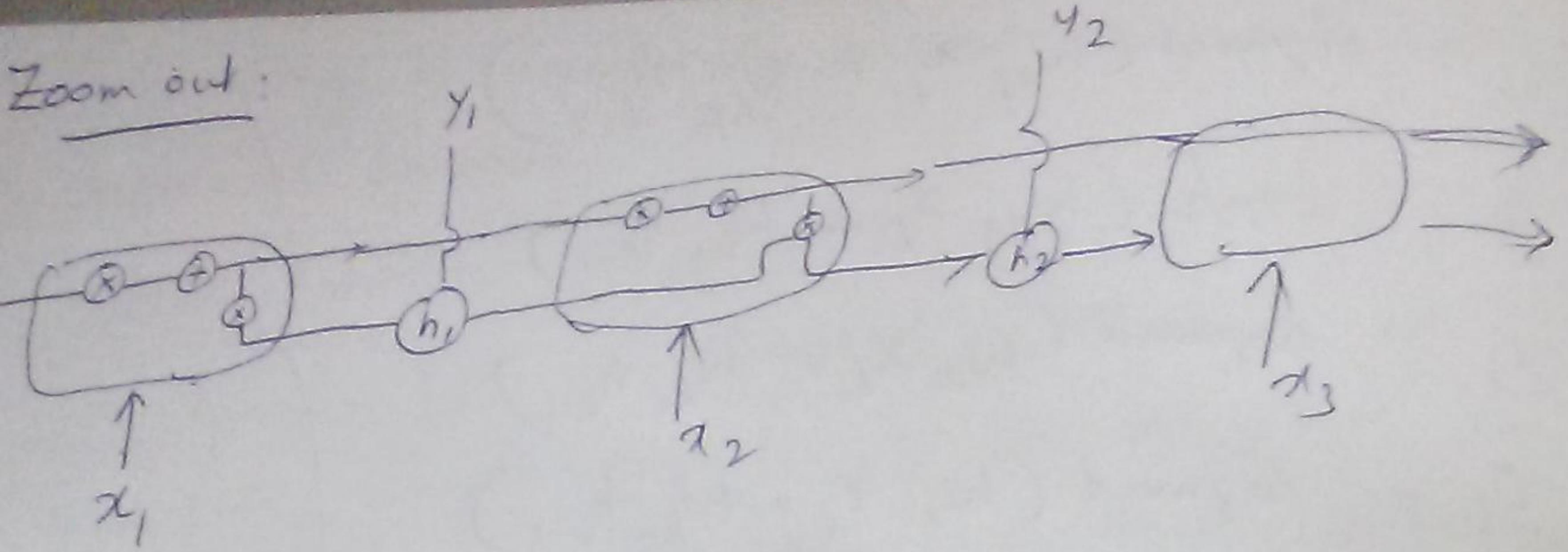
A near-full version



Final Version



Zoom out:



Training:

Truncated BPTT is used here. Ex: if a 1000 word 'data' is sent in, we effectively have a 100-layer feed-forward net. Forward & backpropagation on this 100-layer net can be impossible.

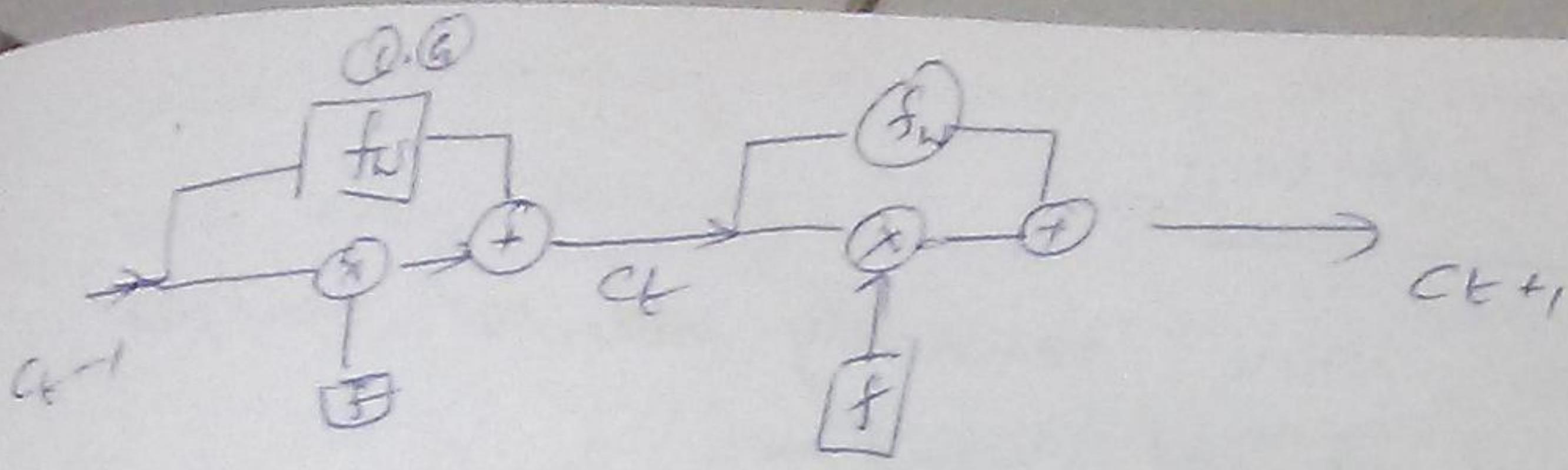
We don't remember ALL states:

$$\frac{\partial c_t}{\partial c_{t-1}} = \frac{\partial}{\partial c_{t-1}} (f_t^0 c_{t-1} + i_t^0 g_t)$$

$$= f_t^0 1 + \frac{\partial f_t^0}{\partial g_{t-1}} + i_t^0 \frac{\partial g_t}{\partial c_{t-1}} + g_t \frac{\partial i_t^0}{\partial c_{t-1}}$$

$$\frac{\partial c_t}{\partial c_{t-1}} = f_t^0 + (\text{various } \overset{\text{may be}}{\text{gradient}})$$

$$\left[ \frac{\partial c_t}{\partial c_{t-1}} = f_t^0 \right] - \underline{\text{Very important}}$$



Usually,  
the gates bias is set very high, so that we  
can start with all 1's in  $f$ .

if a gradient vanished, cell states before that  
state are not contibuting (we need to forget them)  
which might be good.

### Exploding Gradient Pt:

Use Gradient clipping       $\text{grad} = \min(\text{grad}, \text{threshold})$   
Exploding gradients quite possible because we have  
summation of terms.

## Dropout Technique:

A technique where randomly selected neurons are "dropped out" at random. Their contribution to activation of downstream neurons is temporarily removed or the forward pass and any weight<sup>updates</sup> are not applied to the neurons on the backward pass.

This means that if neurons are randomly dropped during training then other neurons must step in and handle representations required to make predictions for the missing neurons. This is believed to result in multiple independent internal representations being learned by the network. This results in the net being less sensitive to the specific weights of the neurons.

Every neuron has a probability  $p$  of being dropped out  $\rightarrow$  dropout rate. Typically set to 0.5

After training, neurons don't get dropped anywhere

[neurons are dropped at random at every training instance. This means a neuron dropped for one instance may come back for the next instance]

would a company do well if its employees forced  
to coin every morning to decide whether to come to  
work or not? So the company cannot rely on ONE  
individual to fill in the coffee machine or manage  
supplies. The expertise is spread across several people.  
The company could be more resilient. If one person  
quit, it couldn't make too much difference.

A different way of understanding this is to realize that  
with dropout, a different neural net is generated  
at each training step. Since each neuron can  
either be present or absent, there are  $2^N$  networks  
possible ( $N = \text{total no. of droppable neurons}$ ). So if we  
went through 10000 training steps, we have effectively  
10000 nets, each with one training instance. The  
resulting neural net can be seen as an averaging  
'ensemble' of these smaller nets.

One small technical detail:

Since  $p=0.5$ , each neuron is connected to twice  
as many input neurons (on avg) as it was during  
training. To compensate, we need to multiply each  
neuron's weight by  $\frac{0.5}{\text{[keep rate]}}$  AFTER training. [Keep rate  
=  $1-p$ ]

(Remember: After training, all neurons are back.  
The net is up to full-strength)

## Binary-crossentropy vs categorical-crossentropy

$$E(y, \hat{y}) = -\sum_t y_t \log(\hat{y}_t) \rightarrow \begin{matrix} \text{Binary crossentropy} \\ \text{log entropy} \end{matrix}$$

## Categorical-crossentropy

For categorical-crossentropy, the targets need to be one-hot vectors, the size of the target vector being equal to no. of classes possible.

## Binary-cross-entropy

$$BCE(t, o) = -(t \log o + (1-t) \log(1-o))$$

## Categorical CE

$$CCE(t, o) = -\sum_x t(x) \log(o(x))$$

## Naive Bayes vs RNNs

Both are classifiers. They can accomplish multi-class classifications.

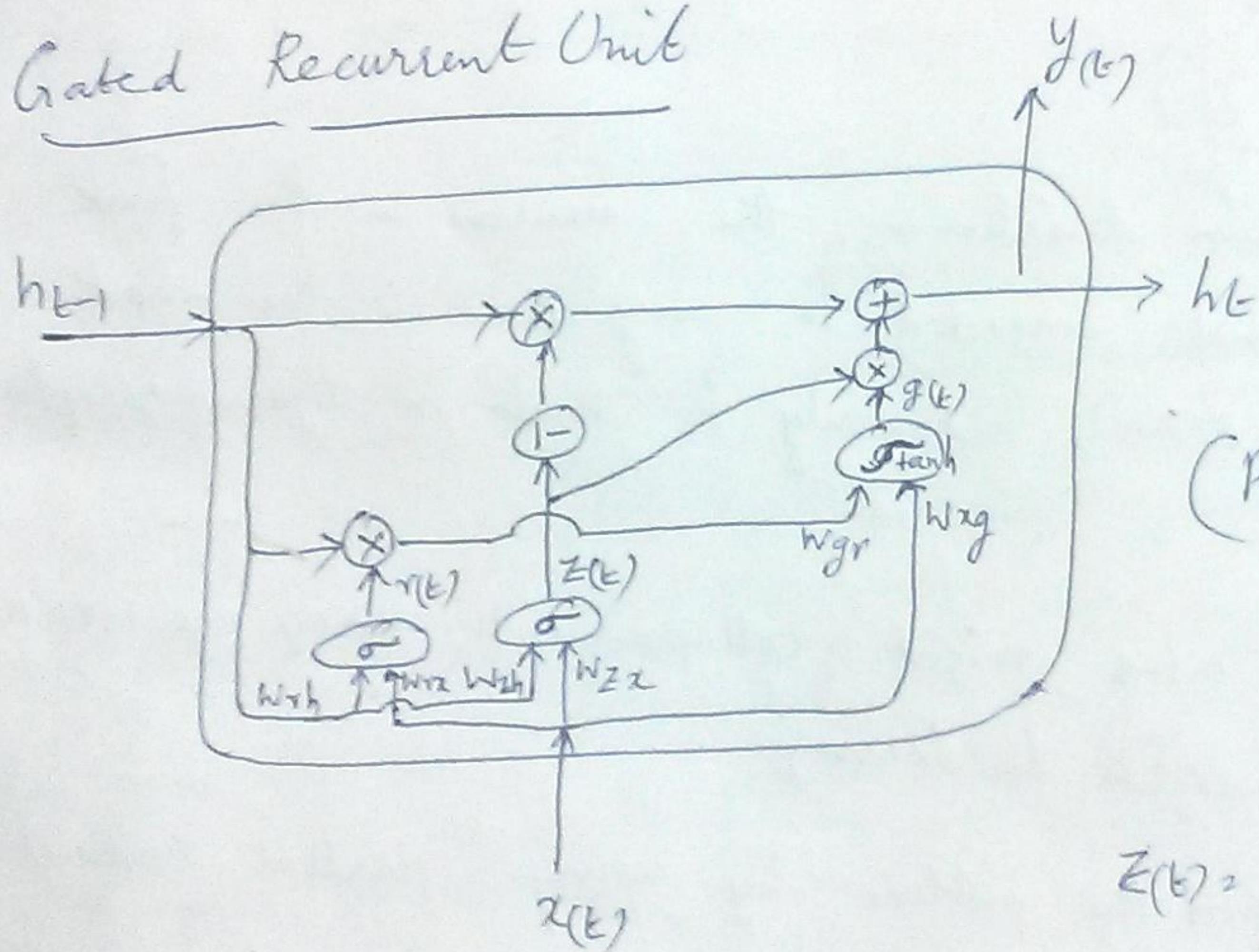
Naive Bayes needs some preprocessing (RNNs too) in terms of (bag of words + tf-idf  $\rightarrow$  n-grams) and assumes i.i.d data. If the assumption holds, it's a very powerful classifier.

RNNs take data sequentially and have memory.  
correlation between words is maintained.

NB → Generative Model - That is, during training, the model tries to see how the 'data' was generated. Try to figure out the underlying distribution that produced the training examples. (HMM, GANs)

RNN - Discriminative: It tries to figure out the diff between the positive and negative samples to learn how to classify. (SVM, decision trees too)

### Gated Recurrent Unit



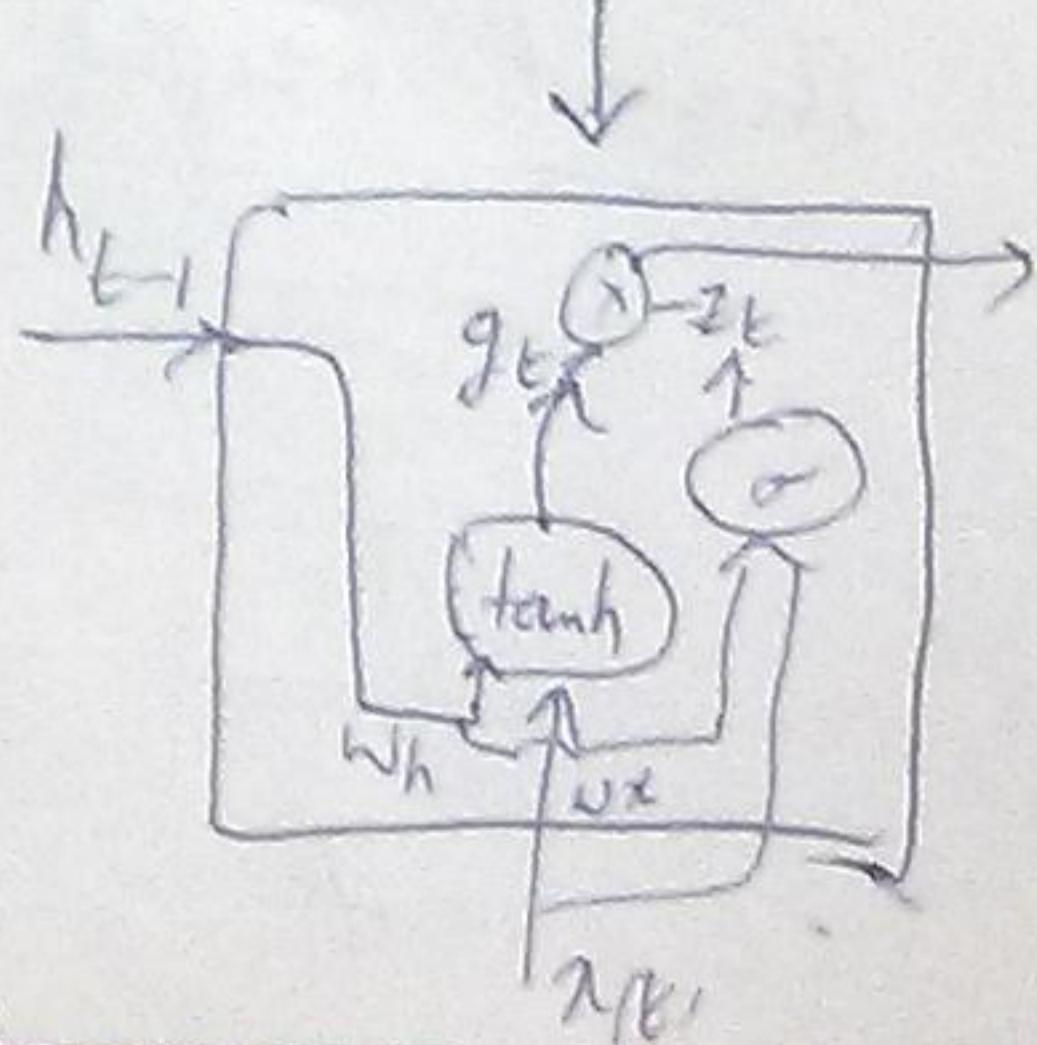
(Performance is similar to LSTM)

$z_t = \sigma(w_{xz} h_{t-1} + w_{zx} x(t))$

$h_t = w_{rh} h_{t-1} + w_{rx} x(t)$

$h_t = h_{t-1} \odot (1 - z_t) + z_t \odot g_t$

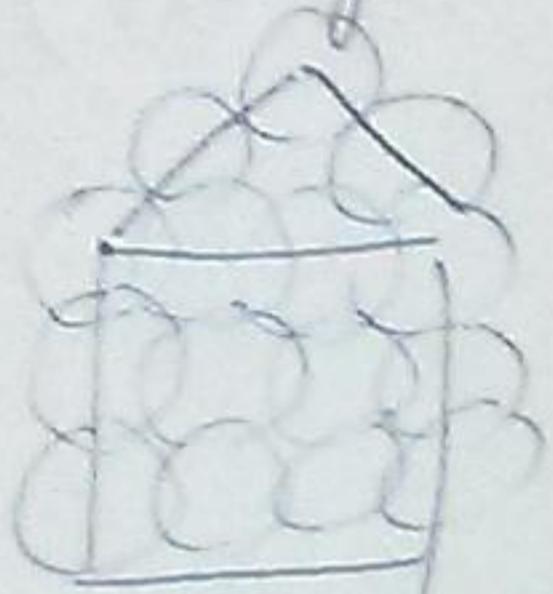
$g_t = \tanh(w_{rg} z_t + w_{gr} (r_t \odot h_{t-1}))$



## Convolutional Neural networks:

- Deep Feed-forward class of Neural nets.
- Visual Imagery - Computer Vision.
- Based on the operation of the visual cortex.

The neurons in the visual cortex have a small receptive field. These fields might overlap and together they tile the whole visual field.

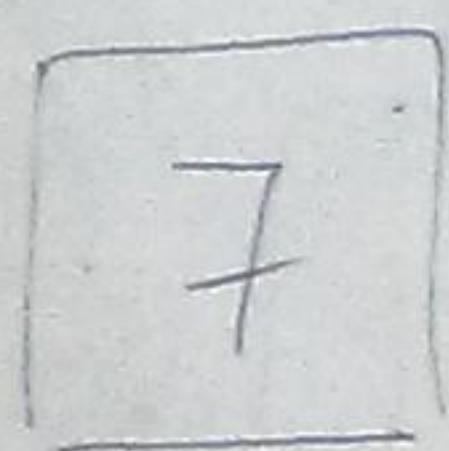


So, the main deviation is, the neurons in the first layer are not connected to every pixel in the input image. (like ANNs), but only to pixels in their receptive areas.

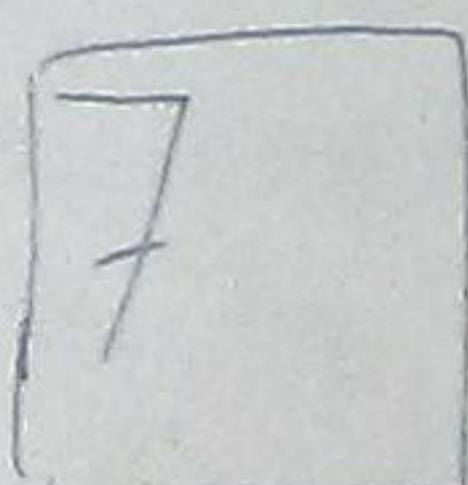
Traditional nets might collapse with image size  $100 \times 100$ .  
(Too many weights to learn)

CNN solves the problem by using partial connections.  
CNN also solves the problem of "Shift variance".

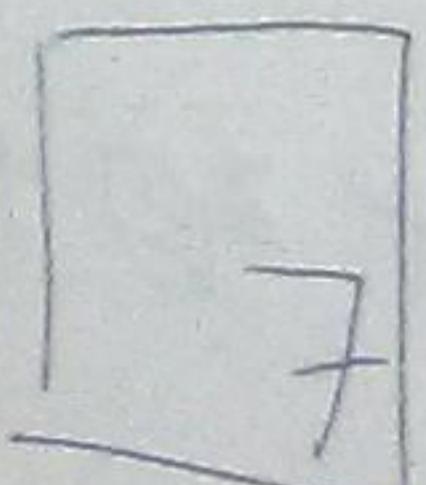
Shared-weight architecture.



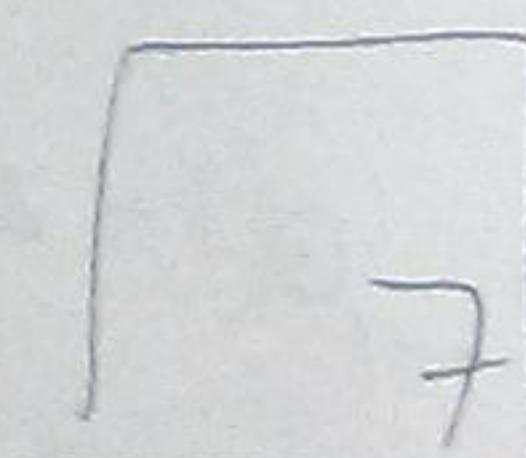
7



?



2



3 ??

Repeating  
grid

We use filters → basically the weights of a neuron  
that need to be learnt. So this effectively means  
the 'features' will be learnt.

### Padding:

- Preserve spatial size. With very deep networks,  
feature-map size can reduce very quickly. Hence padding
- Helps with 'dimension' & math
- Preserves information at the borders  
(Compare to how we see)

$$O = (i - k) + 2p + 1$$

$O$  = output dimension

$k$  = filter dimension

$p$  = padding

$i$  = input dimension

### Filters:

Parameters - Window size? Stride=?  
filter → also called kernel or feature detector  
Output of filter operation on an image part is called  
'Activation Map' or 'feature map'

The more the number of filters, the better the  
classification or learning.

Depth: No. of filters used

Non-padding: if applied → wide convolution  
if not applied → narrow convolution.

ReLU: Non-linearity  $\sigma$  and getting rid of negative values

Spatial pooling: Sub sampling ~ downsampling.

- Reduce dimensionality.
  - a) Max pooling
  - b) Mean pooling
  - c) Sumpooling.
- window-size
  - ↳ stride. Usually equal to window dimension
- Reduces input Representations (manageable)
- Reduces no. of parameters and computations, reducing overfitting.
- Makes the n/w invariant to minor transformations (we take max or mean which reduces the effect)

Fully connected Layer:

Traditional Multi-layer Perceptron (MLP) that uses a softmax activation unit at the output.

This is a fully connected n/w.

The o/p from the pooling layer  $\rightarrow$  high-dimensional features.

Using Existing Models

Training a CNN from scratch can be extremely time-consuming.

We can leverage existing trained Models.  
Popular Models are

- 1) VGG -16/19 - Oxford
- 2) ResNet 50 - Microsoft
- 3) InceptionV3 - Google
- 4) LeNet
- 5) AlexNet etc

We can either use the architecture and train the model on our data from scratch

or we can freeze many layers (initial layers, which learn low-level features like curves, edges etc) and train the subsequent layers only.

(Refer the document (softcopy) for more detail)

GANs :

Please refer last Semester's notes for now !!

## A Quick Note on Vanishing Gradients problem

The deeper a network is, the greater the chance of vanishing gradients due to numerous multiplications.

### Some common solutions

#### ~~Bad idea~~

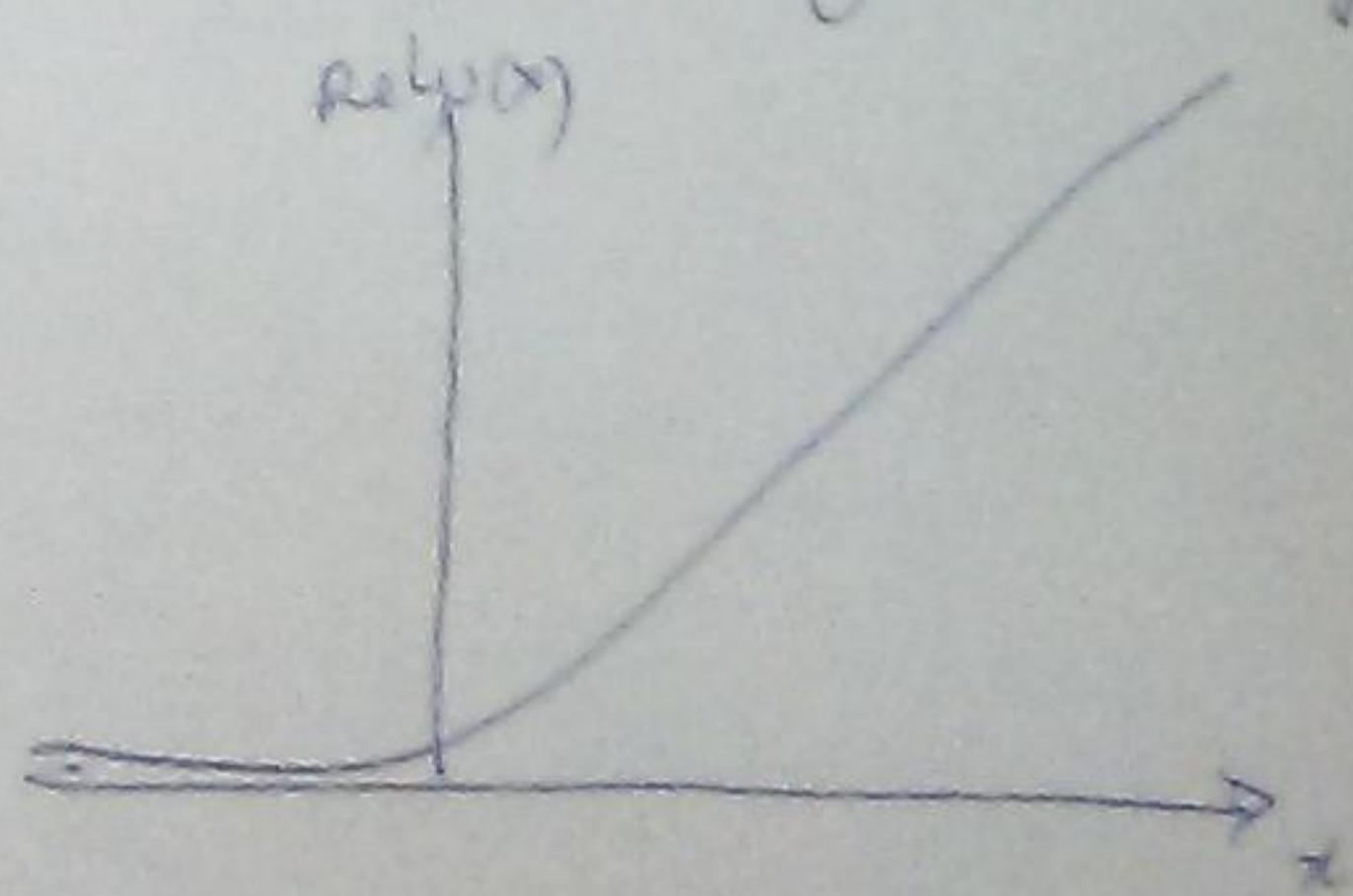
1. Use non-saturating activation functions.  
(with sigmoids, we have the max derivative as  $\sigma(1-\sigma)$  which is  $= 1/4$ . Hence the max derivative we can get for multiplication is 0.25  
with tanh, it is slightly better. we can get a max value of '1').

ReLUs are better.

we can consider

$$\text{ReLU}(x) \cdot \log_e(1+e^x)$$

This ensures the following graph



for positive values of 'x' derivative is significant

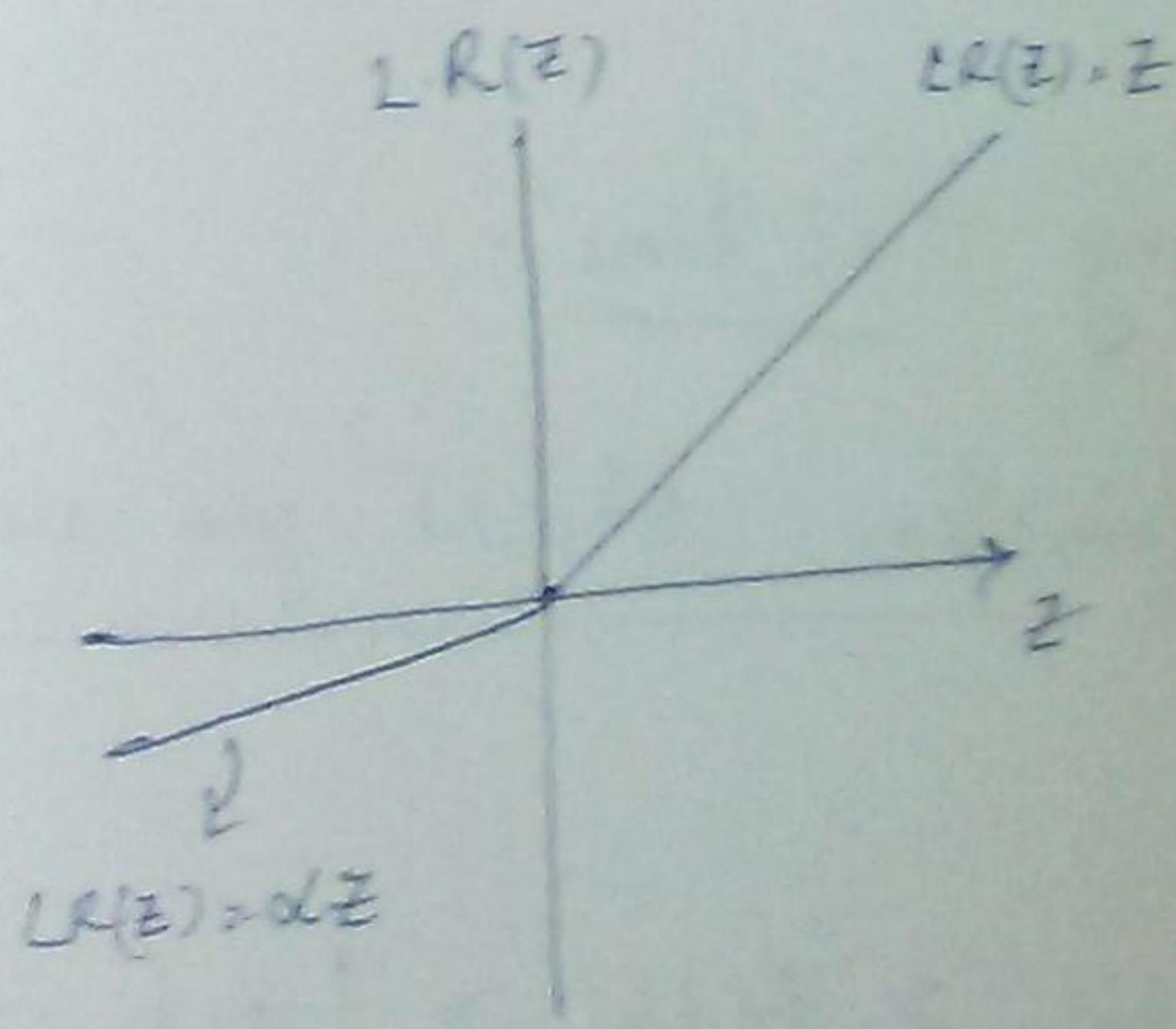
But for negative values of  $z$ , derivative might reach zero, and so we have no weight updates.

This is called 'dying ReLU' problem.

Once a ReLU unit is dead, no op comes from the concerned neurons. Learning ceases.

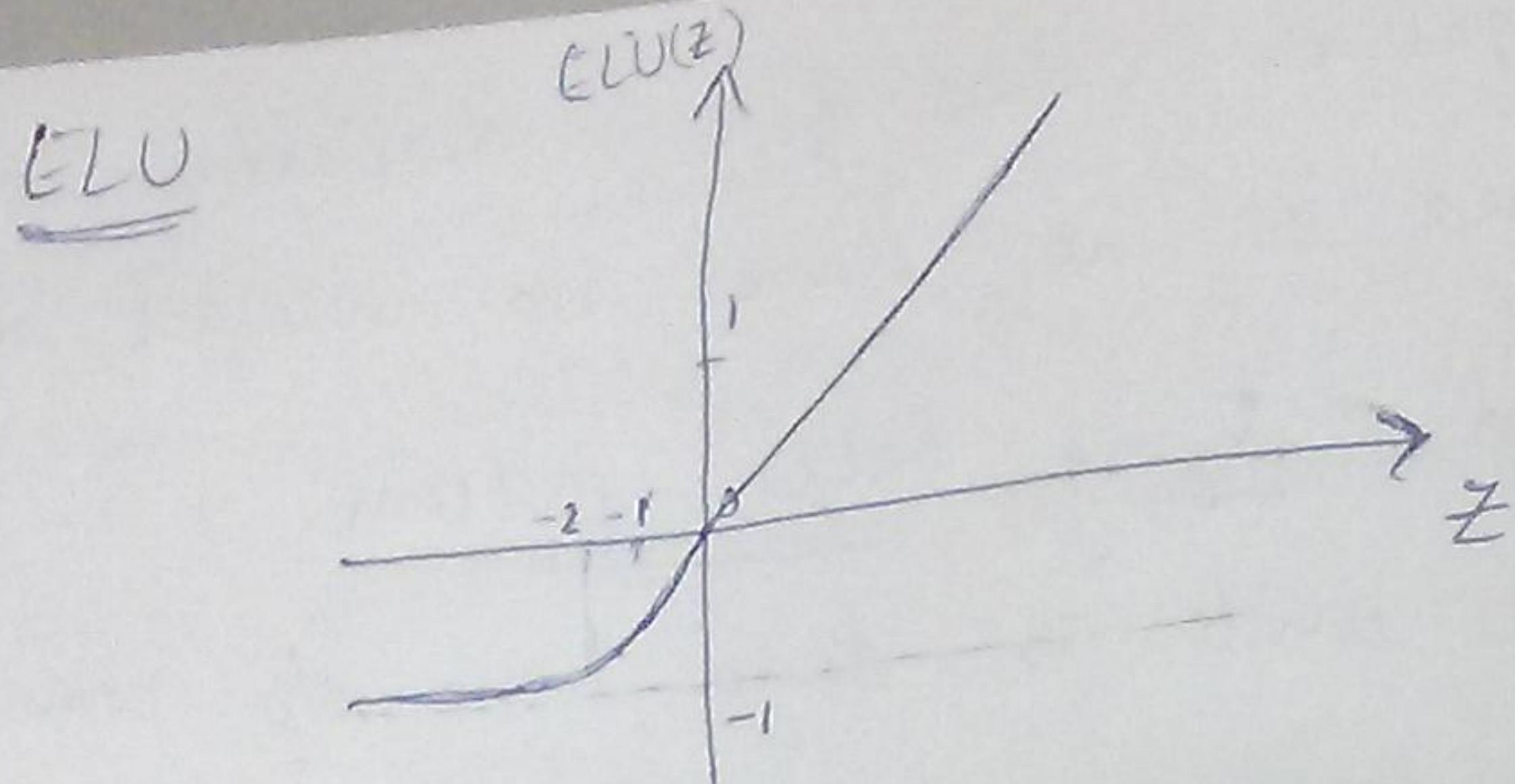
2) Therefore, we can use a modified ReLU, called 'Leaky ReLU'.

$L\text{-ReLU}(z) = \max(\alpha z, z)$ . The parameter  $\alpha$  decides how much the function 'leaks'.  $\alpha$  is the slope of the function for  $z < 0$ . (Usually set to 0.01 or lesser).



The small slope ensures the ReLU never dies.

3) A novel variation of the ReLU is the ELU also called Exponential Linear Unit. This performs best. (training time reduces and performance increases)



this graph is  
for  $\alpha = 1$   
(not to scale)

$$\text{ELU}_\alpha(z) = \begin{cases} \alpha(e^z - 1), & z < 0 \\ z & z \geq 0 \end{cases}$$

This ensures that for negative values of  $z$ , we have a constant off from the ELU. This ensures we don't have the vanishing gradient issue and also don't have the 'dying ReLU' problem.

Drawback:

can be slower during 'test time'.

Which activation function should you use?

Generally speaking,

$$\text{ELU} > \text{Leaky ReLU} > \text{ReLU} > \tanh > \sigma$$

For Exploding Gradients, use Gradient clipping.

(S will cover the concept of Batch Normalization after the Test)

## Convolutional Neural Networks - CNN

chanda's

convolution - winding together

widely popular with image classification Traditional image processing methods are still popular.

The problem:

Identifying or classifying an image. We want the computer to differentiate between all images and learn the unique features that make a dog a dog and a cat a cat.

The computer should be able to perform image-classification by looking for low level features like curves or edges and then building up more abstract concepts through a series of convolutional layers.

Needless to say, the network will be humongous.

The biological inspiration is the visual cortex where each neuron has a limited "local receptive field". This means they react to visual stimuli located only within a limited region of the visual field. Different visual fields of different neurons might overlap.

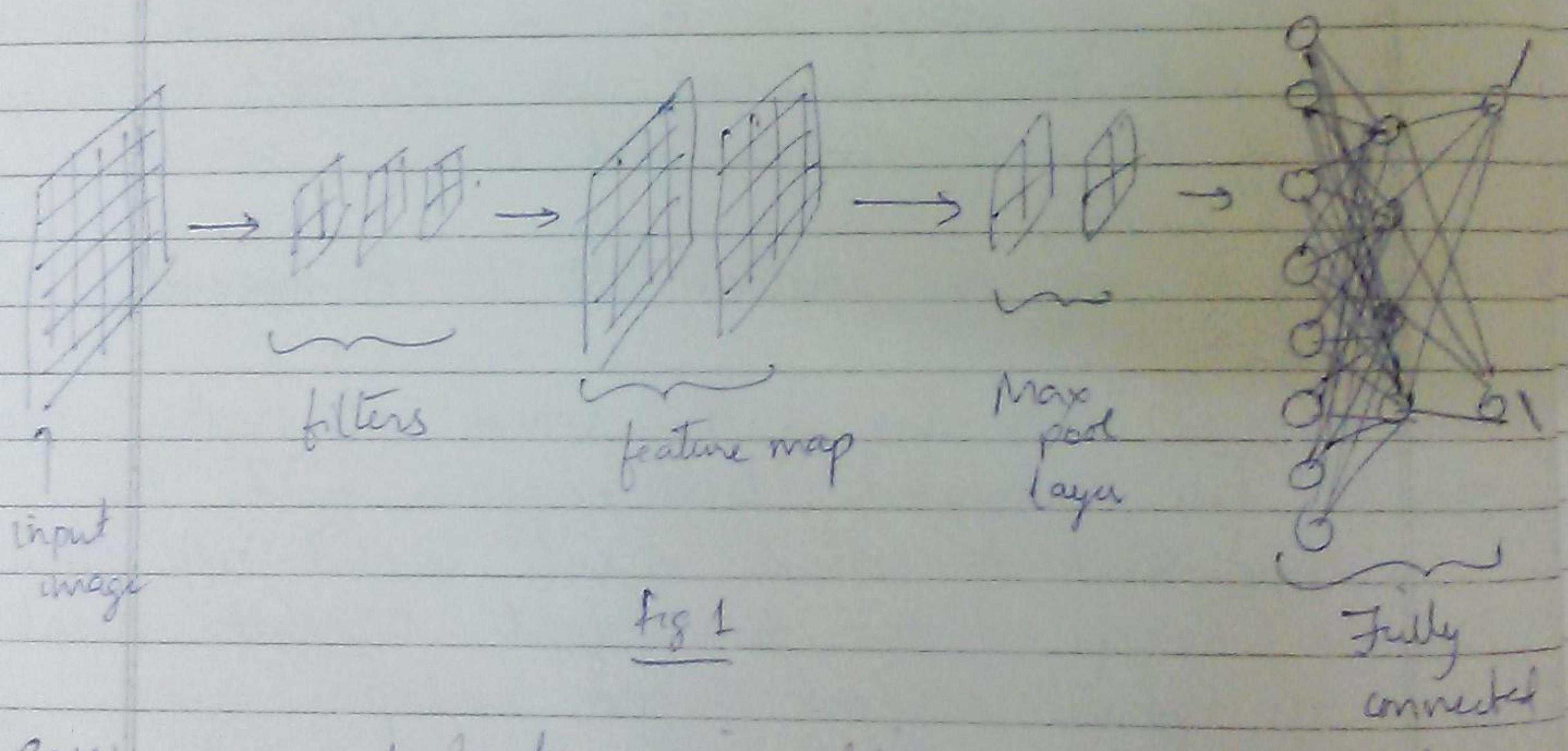


← The circles represent 'receptive fields' of each neuron. One circle per neuron.

Example - Task is to identify forward / or backward slash \ or X.

The convolutional NN contains mainly of 3 main components

- a) The convolutional layer - consists of filters or neurons or kernels which react to a limited visual field.
- b) The Max pool layer (To reduce computation)
- c) The fully connected layer - final classification.



CNNs are robust to noise like translation, scaling, rotation or weighting

The idea is to match tiny parts of the image with a reference we have and then building up to full object detection.

Filter: A neuron with weights. The no of inputs to the neuron = the size of the filter.

Ex: a  $2 \times 2$  filter means a neuron with 4 weights and 4 inputs. The values of the filter will be the weights of the neuron. Each neuron will observe an area = the size of the filter. That is, it observes a  $2 \times 2$  area in the image.

There is a 'stride' factor (stride) which is usually set to 1. This is the no of pixels by which a filter moves over the input image so that the next neuron can react.

Each filter is super-imposed on an area of the input image =  $2 \times 2$ . The off of the neuron observing this part of the image is  $w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4$  (perception)

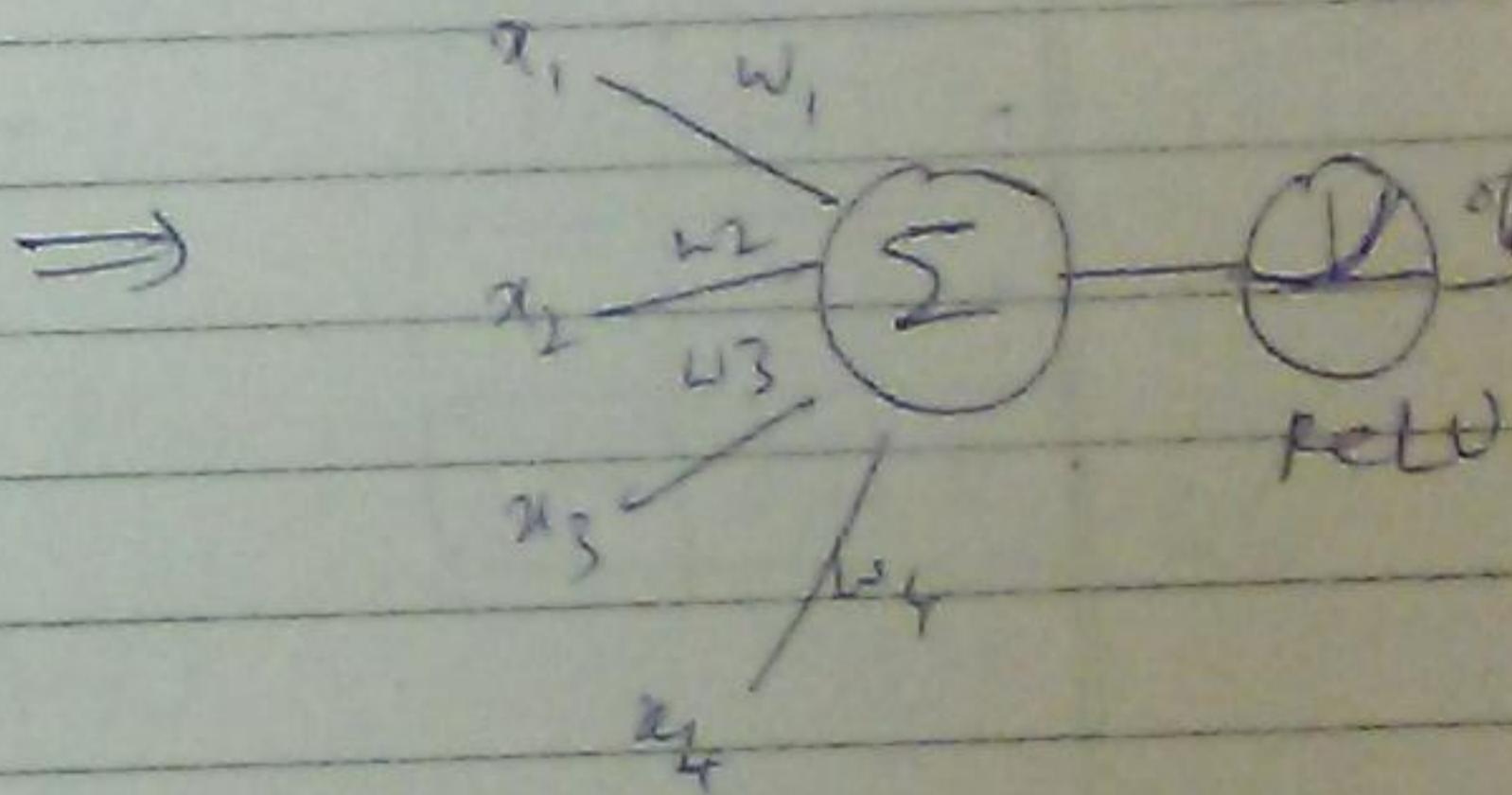
where

$$\begin{bmatrix} w_1 & w_2 \\ w_3 & w_4 \end{bmatrix}$$

$$\begin{array}{|c|c|} \hline x_1 & x_2 \\ \hline x_3 & x_4 \\ \hline \end{array}$$

filter

area of ilp  
image in  
visual field of  
perception



Each perception uses the 'ReLU' as activation unit.

The filters are small parts of the actual (correct, classified) image which give high op when the ilp image part on which it is super-imposed, matches the filter.

Example - A filter of size  $2 \times 2$  representing a '1' can be

$$\begin{array}{|c|c|} \hline \cancel{1} & \cancel{1} \\ \hline \end{array} \Rightarrow \begin{array}{|c|c|} \hline -1 & 1 \\ \hline 1 & -1 \\ \hline \end{array}$$

Similarly a filter representing a ' \ ' is

$$\begin{array}{|c|c|} \hline \cancel{1} & \cancel{-1} \\ \hline \cancel{-1} & \cancel{1} \\ \hline \end{array}$$

which is  $\begin{array}{|c|c|} \hline 1 & -1 \\ \hline -1 & 1 \\ \hline \end{array}$

Now if we present this to an ip image which is actually a forward /, (ip image is  $3 \times 3$  say)

$$\begin{array}{c} a \ b \ c \\ \hline 1 & -1 & -1 & 1 \\ \hline 2 & -1 & 1 & 1 \\ \hline 3 & 1 & -1 & -1 \\ \hline \end{array}$$

The filter starts at top left corner and slides right side, one pixel at a time.

(Diff neurons - 4 to be precise are waiting to operate on the ip area presented to them)

1<sup>st</sup> neuron covers a, b, c

2<sup>nd</sup> neuron covers b, c

3<sup>rd</sup> neuron covers a, b

4<sup>th</sup> neuron covers a, b, c

The ops of the neurons (assuming ReLU) are

(Applying the forward filter  $f \Rightarrow \begin{bmatrix} -1 & 1 \\ 1 & -1 \end{bmatrix}$ )

$$1^{\text{st}} \text{ neuron} \rightarrow -1 \times -1 + 1 \times -1 + 1 \times -1 + -1 \times 1 = -2$$

$$2^{\text{nd}} \text{ neuron} \rightarrow -1 \times -1 + 1 \times 1 + -1 \times -1 + 1 \times 1 = 4$$

$$3^{\text{rd}} \text{ neuron} \rightarrow -1 \times -1 + 1 \times 1 + 1 \times 1 + -1 \times -1 = 4$$

$$4^{\text{th}} \text{ neuron} \rightarrow -1 \times 1 + 1 \times -1 + 1 \times -1 + -1 \times -1 = -2$$

So the feature map (for this filter) is

$$\begin{bmatrix} -2 & 4 \\ 4 & -2 \end{bmatrix} = \begin{bmatrix} 0 & 4 \\ 4 & 0 \end{bmatrix} \quad (\text{See Fig 1 for reference})$$

(likewise there are other feature maps for other filters)

Now to reduce computation, we will feed this to a max pool layer. The Max pool layer also has a window size and stride (usually equal to window size)

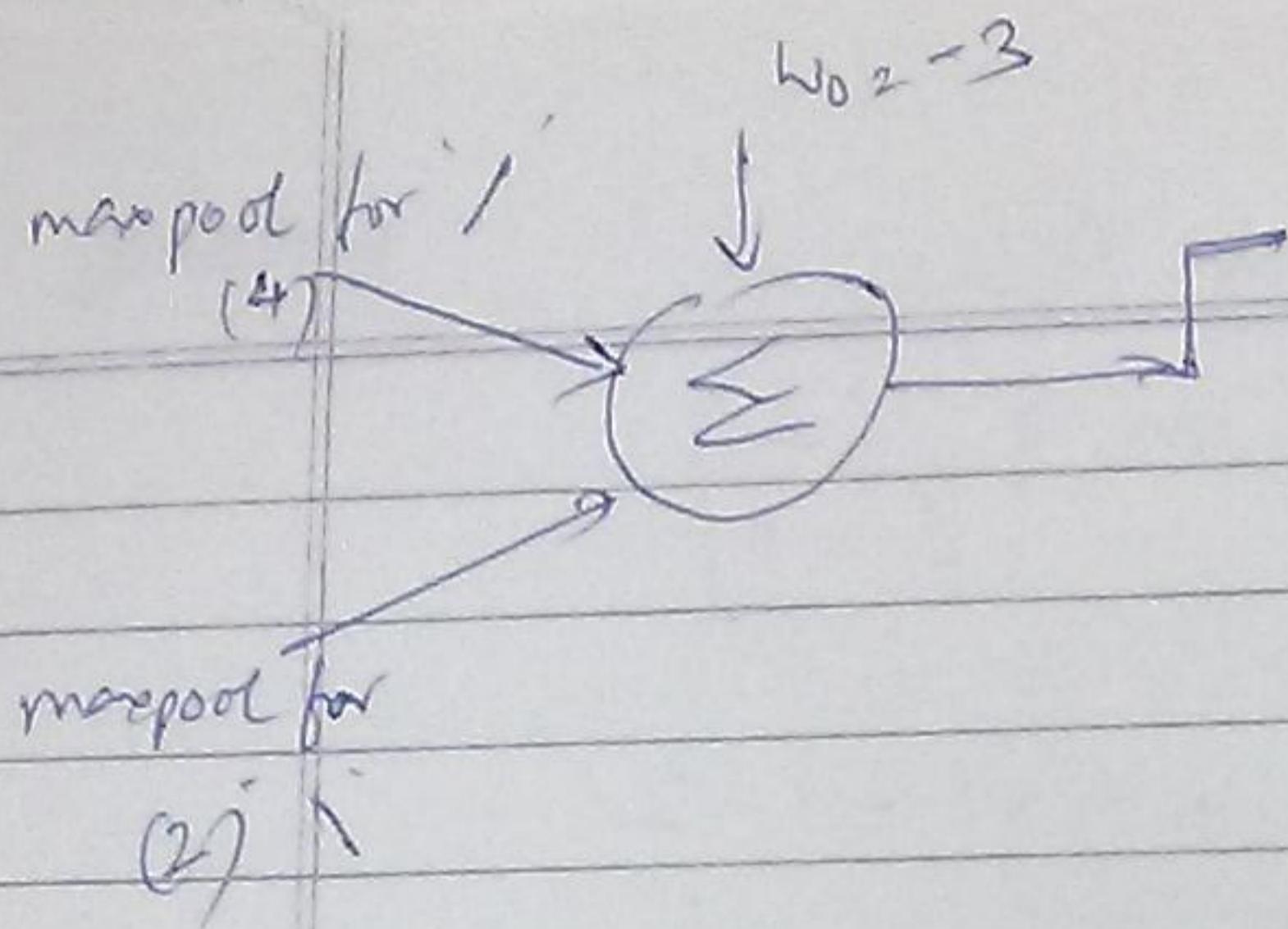
Here let's consider  $2 \times 2$ . The Max pool layer simply takes the max value in the window.

So for the above feature map, Max pool outputs 4

(The second filter for the backslash) would of the following feature map

$$\rightarrow \begin{bmatrix} 2 & -4 \\ -4 & 2 \end{bmatrix} \xrightarrow{\text{ReLU}} \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \quad \begin{array}{l} \text{Max pool w/k } \\ \underline{2} \text{. h/w.} \end{array}$$

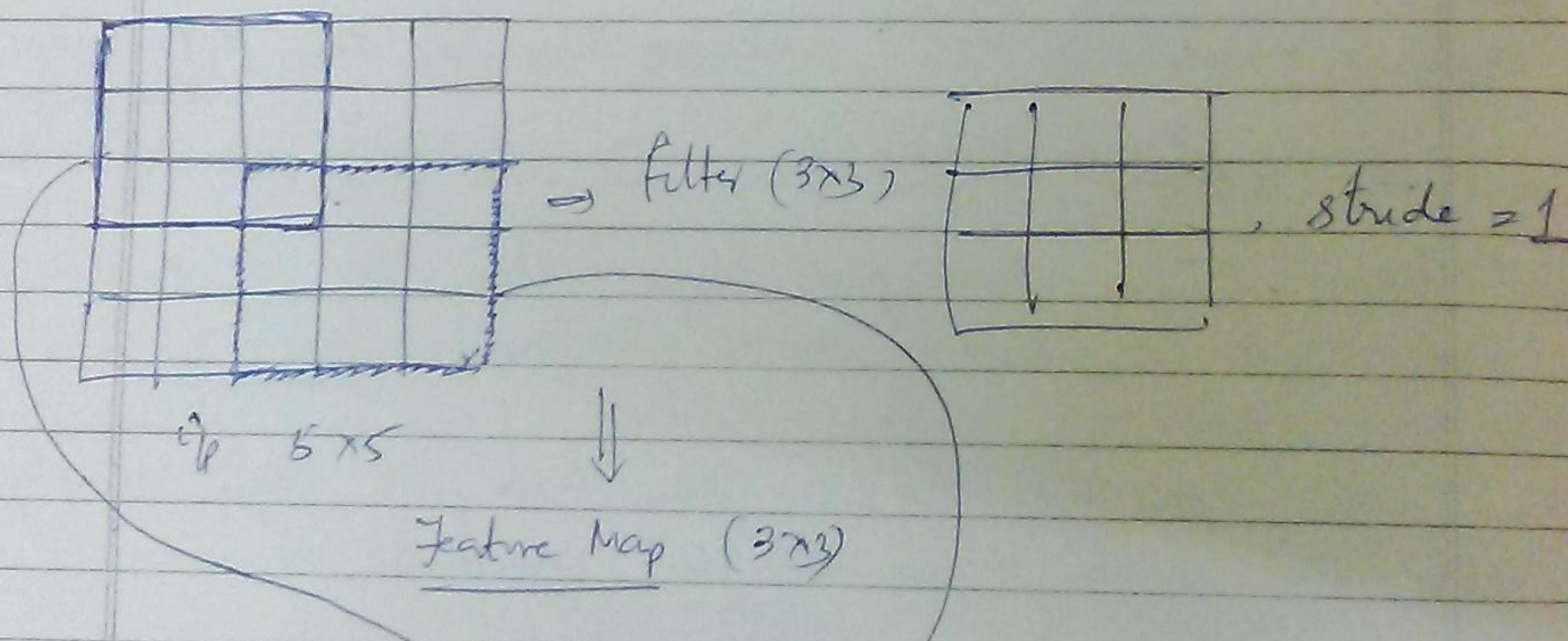
These 2 Max pool ops can now be inputs to a fully connected Neural net. In this case, we only need a single perceptron with  $w_0 = -3$ .



if the neuron triggers it's a forward slash else a backward slash (1 or 0)

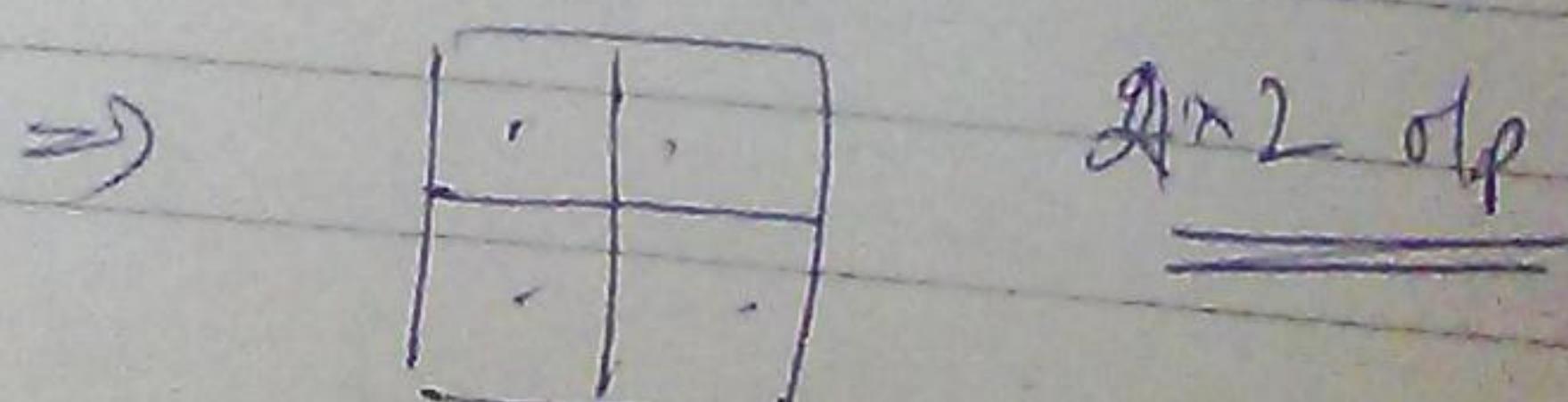
of the same thing is imagined on a bigger scale, we have a real CNN

Training happens by the well-known Backpropagation  
(The filters are also learnt)



↓ here one cell represents the value obtained from a neuron having nine weights

Maxpool (2x2), stride = 2 (placed on the feature map)



Now if there are 2 filters, then there will be 2 o/p from the maxpool layer, each  $2 \times 2$  size - So 8 o/p - <sup>chandra's</sup>

These 8 o/p are fed to a fully connected neural net having 8 ip nodes, one or more hidden layers and an op layer having as many nodes as classifications required

