



- I) Why do we need non-linear?
  - Complex decision boundary
  - Single neuron can only compute a linear function even if nonlinearity activation is used. We need nets to get complete non-linear decision boundary
  - Bias term tells us how active neuron is w/o input
  - Loss function is the important one to optimize that helps to actually learn
  - Universal approximator tells us network has capacity to learn any function but doesn't tell us about how to learn them?
- II) No activation :- Piece wise linear  $\rightarrow$  Take max of linear functions at that space.
  - w/o activation we just get linear units, all multiple linear units collapse to single linear function. But using linear units is used in model compression or matrix decomposition  $\xrightarrow{\text{PCA}}$   $1000 \times 1000 \Rightarrow 1000 \times 100 + 100 \times 1000$
- III) Bias need :- Some systems need to output some defined value as default state. To keep deep default state of n/w. It is considered as another input hence could be moved
- IV) MLP capacity :- Empirical tuning, while compressing output to lower dim we tend
  - To remove unnecessary variance like noise and include variance inherent to data's true nature
- V) softmax :- Represent probability distribution, exponential makes sure o/p is always  $0 < p_i < 1$ , normalization applied. To add up to one.
  - It is like sigmoid with multiple classes. This is most like softmax rather than softmax. To get softmax we must multiply them with final output
- VI) Loss function :- Maximum likelihood :- Learn a model under which data is more likely.
  - Maximizing  $P(Y=1|z)$ . Use logarithm to avoid underflow/overflow.
  - issues :-  $\log \text{softmax}(z)_i = z_i - \log \sum_j \text{exp}(z_j)$   
 $\xrightarrow{\text{if } z_j \rightarrow -\infty}$  Due to presence of log, softmax works well for  $\rightarrow$  log likelihood

Disclaimer: Since it's a student notes, there can be many errors. Please don't use notes as is

## MLP Basics - Chapter 1

I) Why do we need non-linear?

→ complex decision boundary

→ Single neuron can only compute a linear function even if non-linearity activation is used. Use neural nets To get complete non linear decision boundary

→ Bias Term Tells us how active neuron is w/o input

→ Loss function is the important one To optimize that helps to actually learn

→ Universal approximator Tells us network has capacity to learn any function but doesn't tell us about how to learn them?

II) Max out Activation :- Piece wise linear → Take max of linear func at that space.

→ w/o activation we just get linear units, all multiple linear units collapse to single layer neuron. But using linear units is used in model compression or matrix decomposition

→ PCA  $1000 \times 1000 \Rightarrow 1000 \times 100 + 100 \times 1000$

III) Bias need :- Some systems need to output some defined value as default state. To keep default state of m/w. It is can be treated as another input hence could be moved

IV) MLP capacity :- Empirical Tuning, while compressing output to lower dim we tend to remove unnecessary variance like noise and include variance inherent to data's true nature

V) softmax :- Represent probability distribution, exponential makes sure o/p is always we, Normalization applied To add up to one.

→ It is like sigmoid with multiple classes. This is most like soft sigmoid rather than softmax. To get softmax we must multiply them with final output

VI) Loss function :- Maximum likelihood :- Learn a model under which data is more likely :- Maximizing  $P(Y=1|z)$ . Use logarithm To avoid underflow/overflow issues.

$$\log \text{softmax}(z)_i = z_i - \log \sum_j \exp(z_j)$$

→ Due to presence of log, softmax works well for  $\rightarrow$  log likelihood

→ Loss function like MSE doesn't work well with softmax or softmax squishes data, the gradients near 1 or 0 become less & less hence learning cannot happen

### Which loss function To use?

- i) Multi class classification :- crossentropy / negative loglikelihood + softmax
- ii) Mean square error :- Linear o/p layer → used for regression problems like house prediction

### Learning by Gradient descent :- Why do we need To compute gradients?

→ Because we don't have the info of network parameters v/s error plot beforehand, so that we can visually handpick best parameters.

→ Knowing this beforehand we need to plot for every combo of  $w_1, w_2$  v/s error curve

→ For a network with millions of parameters  
This is not possible

→ Hence find gradient that tells us maximum change of error. Traverse in direction where error decreases

\* \* \* - Problem :- If network has lots of parameters then it becomes expensive to compute derivatives directly. grad descent is a learning algorithm we use to pollute our network's parameters.

→ Use backprop approach to help calculate gradients efficiently

$$w_i = w_i - \alpha \Delta w_i \rightarrow \text{gradient descent algorithm}$$

$$\Delta w_i = \Delta w_j \cdot \Delta w_k \cdot \Delta o_n \rightarrow \text{back propagation algorithm}$$

grad descent v/s Back prop

→ Backprop is chain rule backwards. Use computation graph to do this in software

### Types of Training

- a) Batch Training :- smooth descent, but can take really long for huge batch  
→ Also complete data may not be fitting memory  
→ The experience / loss of previous examples can help in learning next examples

- b) Stochastic learning : Faster learning, but can be a bit jumpy and can miss minima or stuck in local minima. Even noise in data is unnecessary
- captured. Cannot take advantage of GPU
  - c) Trade off - Minibatch Training :

## Chapter - 2

### CNN - Block 1

Motivation :- Recognize objects, audio independent from its position on image/spectrogram

Idea :- Use shared weights  $\rightarrow$  Reusing same weight applied at different positions

Implementation :- same filter kernel applied throughout image/spectrogram.

Activation map :- Output after convolving filter with input. Some weights are applied throughout different places at image

How it's learnt ? Weights & Biases are learnt through gradient descent + backprop.

$\rightarrow$  Each time filter is applied on different position of image we may receive different gradient for same parameter

$\rightarrow$  So take all these gradients, average them up and then update the param



Receives 25 gradients as filter moves throughout image

Average all 25 of them & update

$\rightarrow$  To apply convolution one must mostly do it on linear space  
Example : spectrogram [frequency (log scale) v/s Time (linear scale)]

### Properties of CNN

- a) Local connectivity :- Unlike dense fully connected n/w where every neuron is connected to all neurons from previous layer. In convolution n/w receptive field for convolution n/w receptive field is not complete, we say

To increase receptive field of CNN as we go deeper

b) shared weights :- Unlike FC, the same filter/param is used to apply to all position, hence weights are shared.

c) Translation equivariance :- No rotation / scale / shear, just shifting. invariance

→ if we move image to another position, convolution can work with almost same filters w/o much change  $f(g(x)) = g(f(x))$

Filters and Activation : Each filter when applied results in a single off channel

Multichannel computation :- Convolution only on axis that has semantic hence X, Y

This means we cannot convolve across channels as There is no ordering. no ordering

$$R \text{ } g \text{ } B = g \text{ } R \text{ } B = B \text{ } g \text{ } R = B \text{ } R \text{ } g$$

Filter output size :- Activation size decreases as convolved.

Padding : To ensure output activation have desired weight as input.

a) Valid padding

Don't want to add info if it's not there

b) Same padding

Simplifying giving edge pixels equal priority

c) Full padding

Convolution to text processing  
→ paying equal attention to all words

→ odd sized filter makes sense as center pixel can look at all local pixels

Striding : Controls how often we apply convolution on input. How much conv you apply before moving again. Downscale / Decrease resolution

Dilation : We skip values like stride but also skip on input when applying

→ Improves receptive field w/o having filter matrix huge.

→ Bigger filters are more specific whilst smaller filters are more general

Wavenet : 16kHz per second is audio, nothing much interesting happens at millisecond scale. Need a filter that has receptive field of 16k samples

→ We cannot use a big filter, hence use stacked small filters with dilation

→ Really less parameters to train, and resolution is not reduced

Pooling :- Differs from dilation and striding. Aim of pooling is not downsample but it's a side effect of pooling when used along with strides.

- Indilation, increase in receptive field but not reducing total resolution but it's not the case in pooling
- Approximate invariance to small translations when combined with maxout
- No trainable parameters
- But loses some info about location of pixel, helpful to detect whether obj exists and not where?



Even with slight rotation  
maxout can always take max pixel

Strong prior :- Conv filters are local, hence have a bias. This doesn't work well with training tasks where we need to pay attention arbitrarily as local connectivity might not help. This can cause overfitting, underfitting



filter is local hence strong prior only to certain positions (non-zero values mostly)  
(all zeros no prior)

Subsampling → (convolution → subsampling)

+ bias

↓

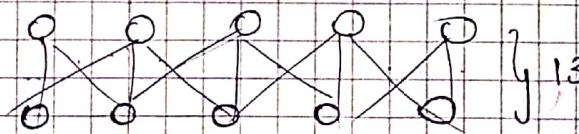
output

→ Bias parameter inside subsamples / pooling

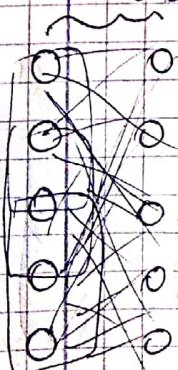
→ pooling inside subsamples (strides = 2)

→ No flattening and dense layer at end but used a direct convolution. It's degraded as we use 4x4 filter alphabet character spanning 16 channels. Eats up entire subsample feature map

25 connections



How can we encourage learning complex filters?  
encourage filter sparsity



\* → Pooling is a cheap and fixed operation, that controls amount of invariance to input translation

\* → Striding is a costly operation compared to pooling which controls (non-linear) (subsampling rate) by amount of skipping

## CNN - block 2 (Important architectures - only)

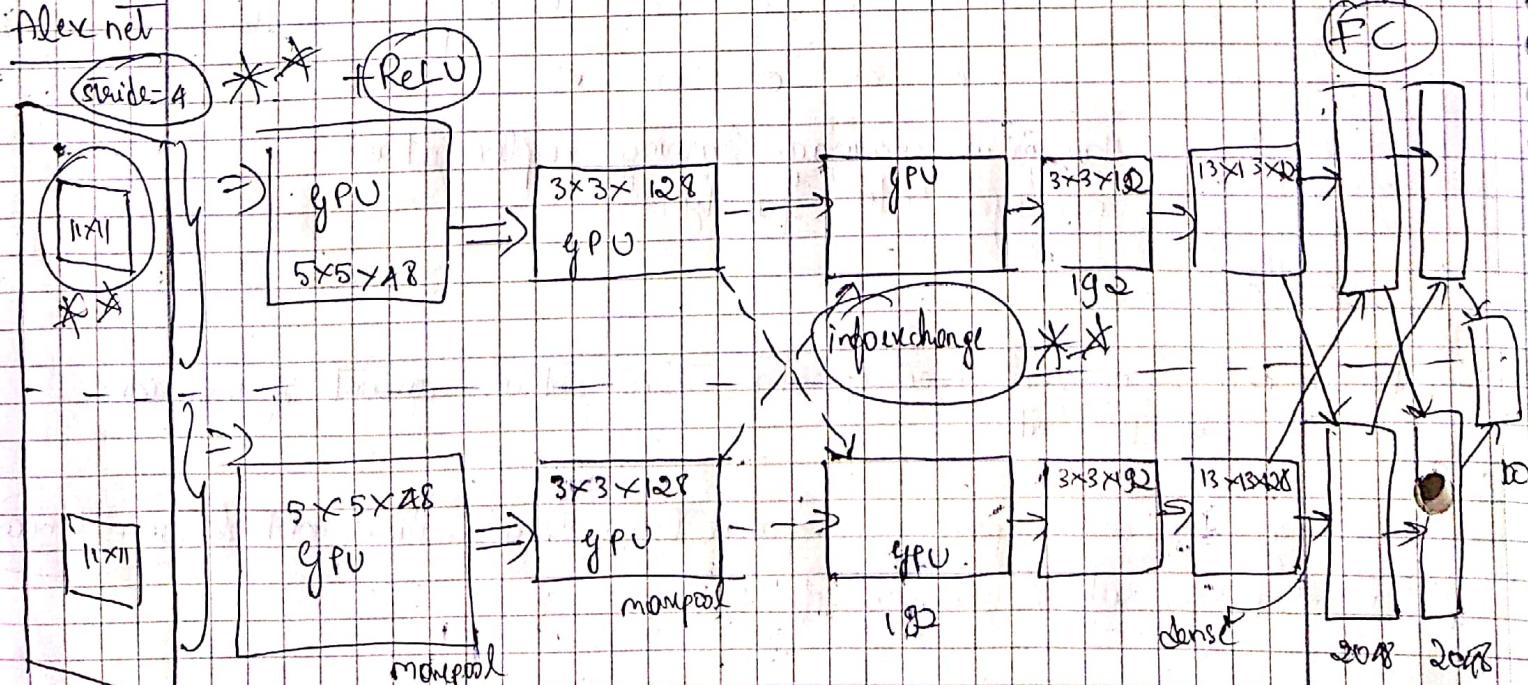
### Alexnet

- Distributed on Two GPUs, Two streams → One looking at Top part and another at bottom part of image. \*
- 11x11 filter size right at start → doesn't help in generalization can only learn specific patterns. Skipping lot of relevant info
- stride of 4 → can lead to Too much downsampling \* so
- Gross linking of activations between Two streams of two gpus → lower speed  
But important as one needs to have info on complete image to take decision
- Usage of ReLU To boost Training

### ZF net

- Remove or keep filter size To 7x7 to retain more info \*
- How did they find it? They used deconvolutional layer to view input from activation. and also unpooling for maxpooling
- Using unpooling is difficult as maxpooling leads to loss of info, hence save the info of which position max element occurred & use it. Other values are set to zero
- Can see what do individual layers detect and recognise?
- first effort to know how convnet learns
- Deconvolution results in checker board patterns

### Alex net



\* Alexnet  
 → Remember about 1x1 filter with stride, Exchange of info b/w two GPU  
 , Max Pooling and ReLU

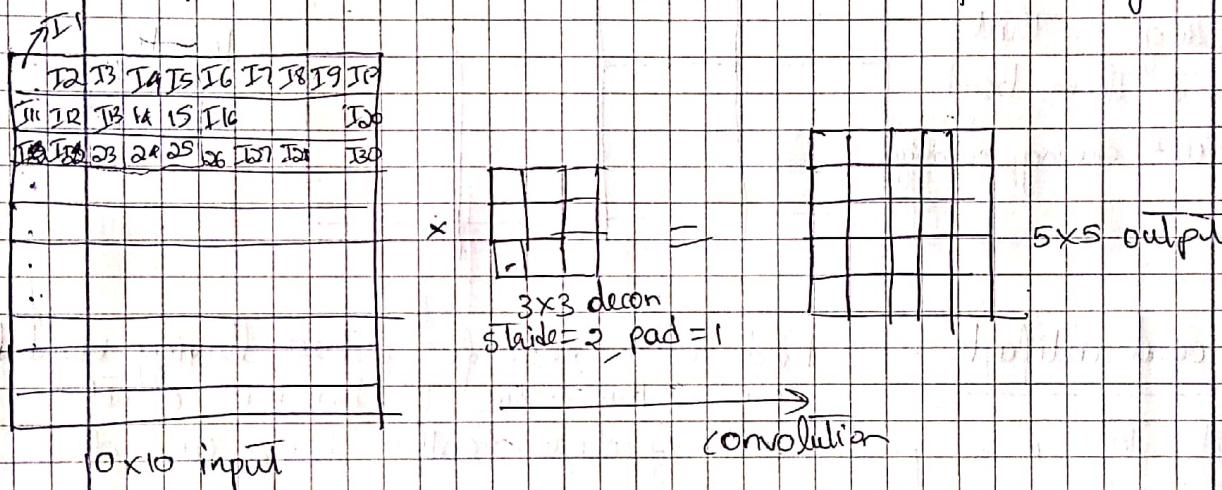
• ZFnet → Remember using 7x7 filters with stride=2 to retain more info  
 that was getting lost in Alexnet. Deconvolution/unpooling layers were first  
 introduced to visualize feature maps.

Q) How does deconvolution work and would it result in info loss?

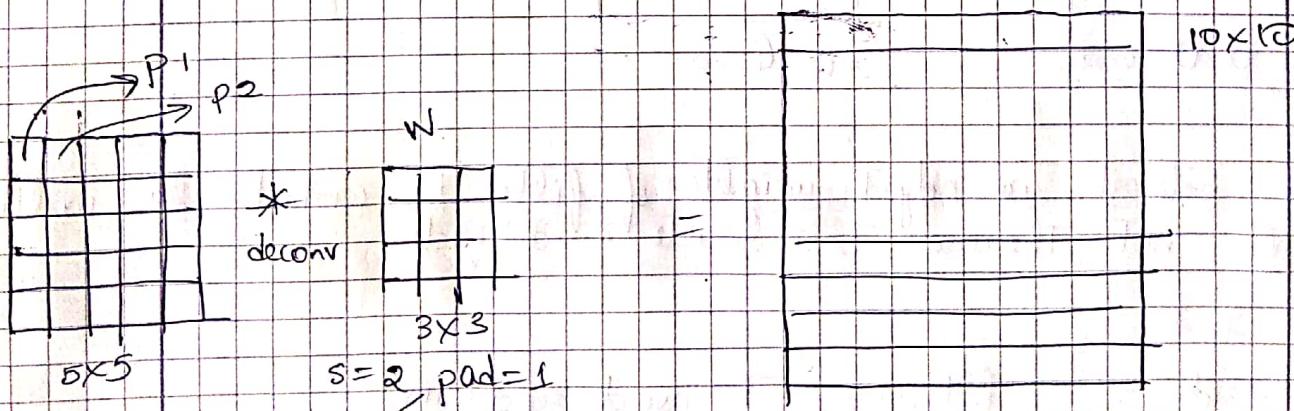
No, it doesn't end with information loss. It can be approximately reversed

as convolution is an operation based on system of linear equations;

ReLU too helps in this as there cannot be -ve pixel range.



$$\text{formula } \frac{(n+2p-f)}{s} + 1 = \text{floor} \left( \frac{10*2+1-3}{2} \right) + 1 = 5 \times 5$$



$$p_1 = w_1 * I_1 + w_2 * I_2 + \dots + w_{35} * I_{35}$$

$$p_2 = w_1 * I_3 + w_2 * I_4 + \dots + w_{36} * I_{36}$$

$$p_3 = w_1 * I_5 + w_2 * I_6 + \dots + w_{37} * I_{37}$$

$$p_4 = w_1 * I_7 + w_2 * I_8 + \dots + w_{38} * I_{38}$$

$$p_5 = w_1 * I_9 + w_2 * I_{10} + \dots + w_{39} * I_{39}$$

$$p_6 = w_1 * I_{11} + w_2 * I_{12} + \dots + w_{40} * I_{40}$$

$$p_7 = w_1 * I_{13} + w_2 * I_{14} + \dots + w_{31} * I_{31}$$

$$p_8 = w_1 * I_{15} + w_2 * I_{16} + \dots + w_{32} * I_{32}$$

$$p_9 = w_1 * I_{17} + w_2 * I_{18} + \dots + w_{33} * I_{33}$$

$$p_{10} = w_1 * I_{19} + w_2 * I_{20} + \dots + w_{34} * I_{34}$$

$$p_{11} = w_1 * I_{21} + w_2 * I_{22} + \dots + w_{35} * I_{35}$$

$$p_{12} = w_1 * I_{23} + w_2 * I_{24} + \dots + w_{36} * I_{36}$$

$$p_{13} = w_1 * I_{25} + w_2 * I_{26} + \dots + w_{37} * I_{37}$$

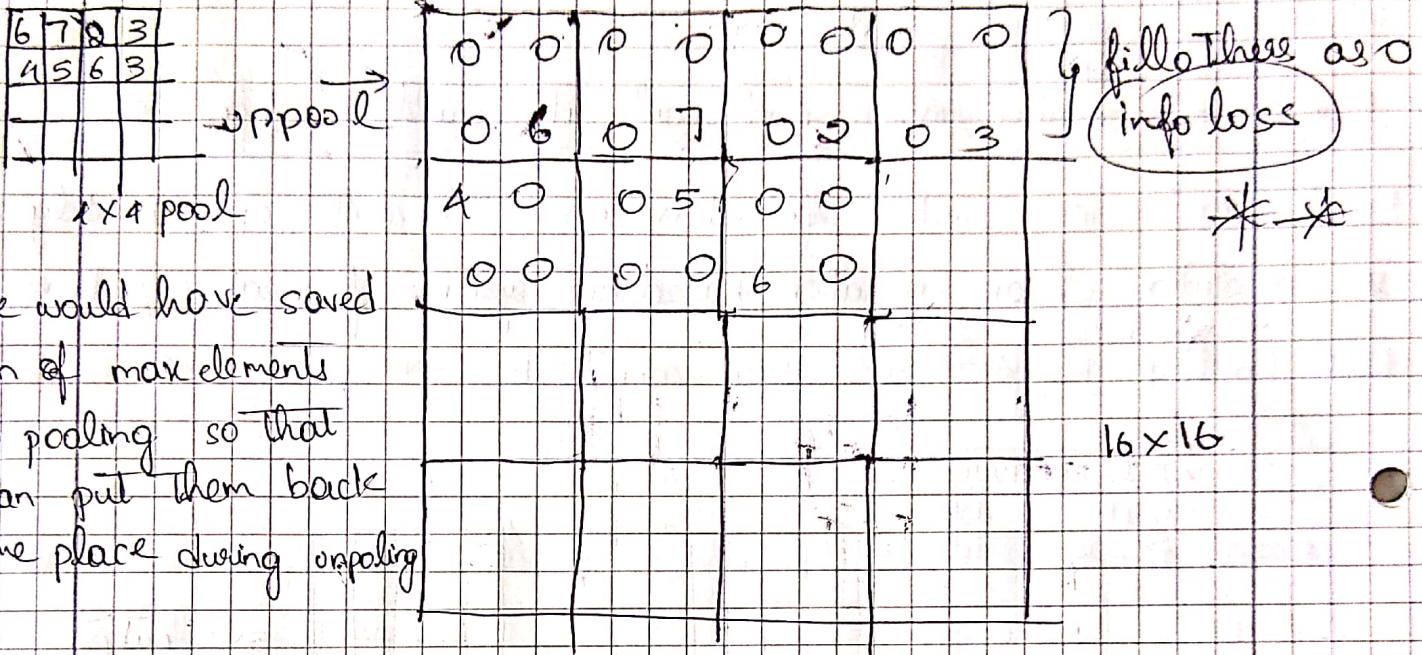
$$p_{14} = w_1 * I_{27} + w_2 * I_{28} + \dots + w_{38} * I_{38}$$

$$p_{15} = w_1 * I_{29} + w_2 * I_{30} + \dots + w_{39} * I_{39}$$

$$p_{16} = w_1 * I_{31} + w_2 * I_{32} + \dots + w_{40} * I_{40}$$

→ As we can see that we have a system of linear eqn which can be solved to finally get  $(I_1, I_2, \dots, I_{100})$  back

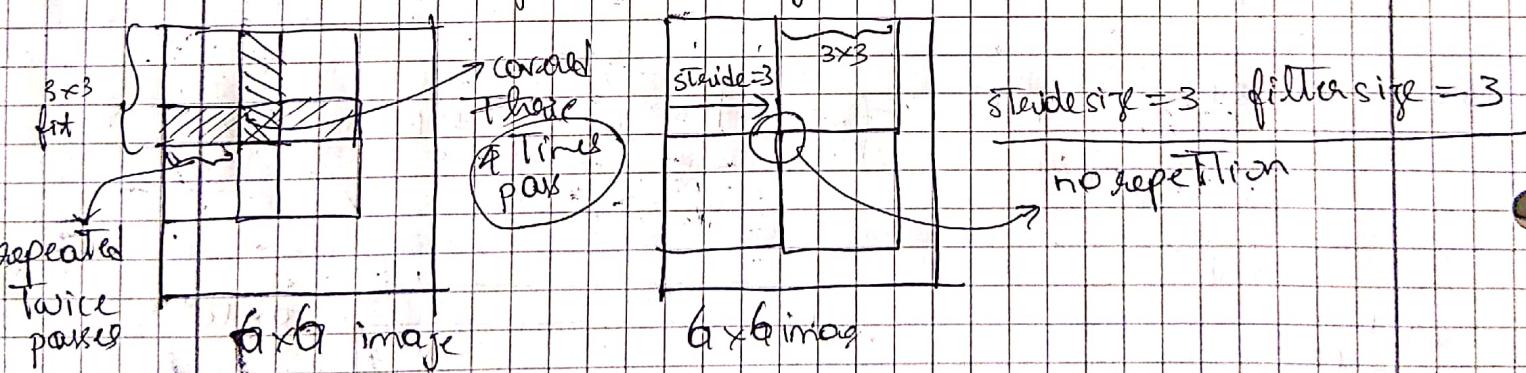
Unpooling: cause info loss as same as pooling



→ We would have saved position of max elements during pooling so that we can put them back to same place during unpooling

Checker board artifacts : tend to be more for pictures having vivid bright colors. This is due to bias term of net.

Bias tends to more produce avg color rather than rare colors



→ Although network can adjust weights of filter to combat this multi-pass. This is often not the case. This results in artifacts

Solution: \*

- i) Make stride size = filter size to avoid repetition
- ii) Better upsampling Techniques

U-net :- High level image  $\rightarrow$  low level image  $\rightarrow$  features  $\rightarrow$  low level image

$\rightarrow$  Often times used for segmentation, where we predict each pixel high res image or group of pixels. Basically finds segmentation mask

$\rightarrow$  Pooling switches are not saved during up sampling as data at same levels are different

$\rightarrow$  Activations at same level are concatenated. This what gives segmentation

$\rightarrow$  At the end have 2 channels  $\rightarrow$  one neuron per pixel per channel (one for foreground and another for background)

$\rightarrow$  U-net is a special architecture used for segmentation semantics.

Object detection :- Put bounding box around required objects & classify them

Semantic segmentation :- Paint a specific class or classes pixel with single color ignore the rest (group of people)

Instance segmentation :- Paint individual person with different pixel colors & ignore rest of classes

What's different in this architecture? Before they were using sliding window technique which are slower to perform semantic segmentation. This architecture changes the

Tasks and Problems that paper is solving

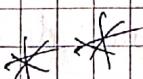
$\rightarrow$  Very few training examples in biomedical training

$\rightarrow$  Usage of Data augmentation To create a generic model

$\rightarrow$  Overlap tile strategy

$\rightarrow$  Need for per pixel classification.

Idea



$\rightarrow$  Down sample To extract features and remove unnecessary variance.

$\rightarrow$  Upsample To get precise locations for these features in high res layers

$\rightarrow$  Concatenating downsampled feature maps with corresponding upsampling layer

$\rightarrow$  To assemble a precise output

$\rightarrow$  Overlap Tile strategy :- To predict pixels at border region, missing context is extrapolated by mirroring

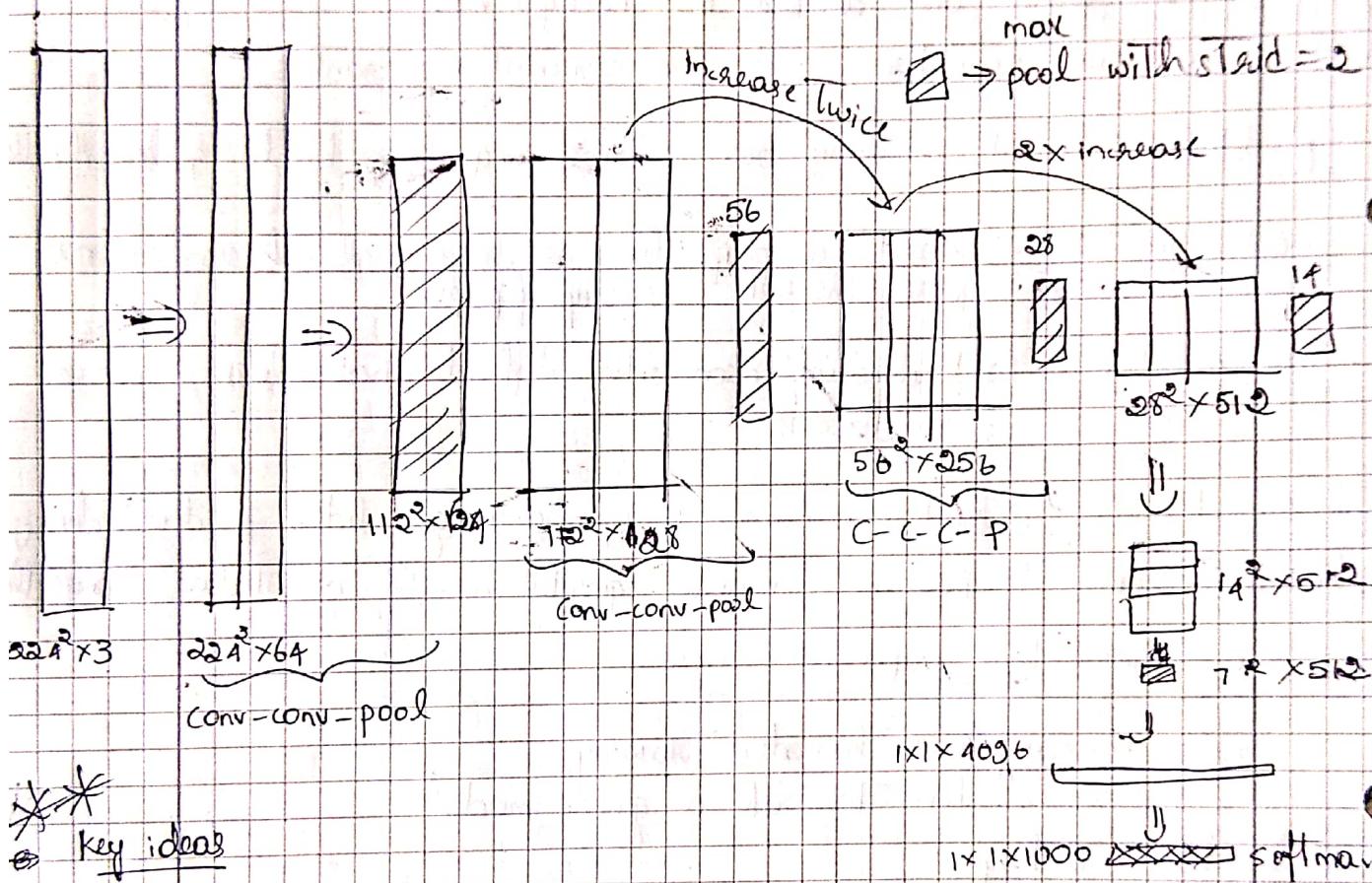
$\rightarrow$  give more weight to background pixels in loss func That helps separating objects of same class close to each other.

## Probabilistic V-net for ambiguous images

- i) Image data annotated by medical students and inexperienced professionals differently
- ii) Need for a model which can handle borderline stuff

VGG net ( $\text{VGG-16}$ ) ( $2-\text{M}; 2-\text{M}; 3-\text{M}; 3-\text{M}; 3-\text{M}$ )

Design principle :- Simplicity & depth

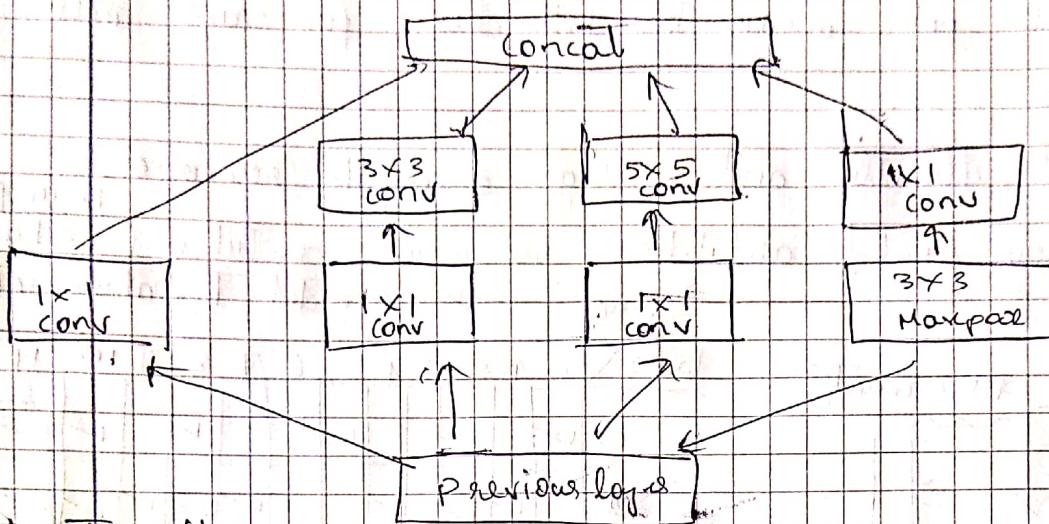


### ~~Key ideas~~

- Simple and elegant design
- Usage of more filters with lower dim size of  $3 \times 3$  instead of  $5 \times 5$  or  $7 \times 7$  as neural network can essentially use two  $3 \times 3$  or three  $3 \times 3$  to makeup for  $5 \times 5$  or  $7 \times 7$  filters. This makes filters more discriminative
- Usage of  $1 \times 1$  convolution in some configurations to increase non-linearity
- Elegant design by usage of same padding and double channel size after each maxpool downsampling
- Lesser kernel size leads to less parinfo among edge pixels even if no of filters are more
- This preserves resolution and

Inception Module :- Apply convolutions of different sizes in parallel and concatenated.

- Due to concatenation, filters/channels may become too large a number.
- Use  $1 \times 1$  convolution as dimensional reduction to reduce no. of channels



Inception V1

GoogleNet :- 9 inception modules  
less parameters in alexnet, use average pool at last layer instead of FC layer to reduce params

Inception V2 :- Adding batch norm. In intermediate layers more changes are present

- What comes into hidden layer is a feature map of previous layer which are constantly changing. All inputs except if image changes, hence learning is much harder. (Ex - Input has pixel range of 0-255 but layers can have 0-500 & layers b/w 0-350 etc and each iteration this may change)
- Normalize intermediate input maps to make learning faster as feature maps will be in defined range

$$y_i = \gamma x_i + \beta \quad x_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} \quad \text{small delta to avoid } \sigma_B = 0$$

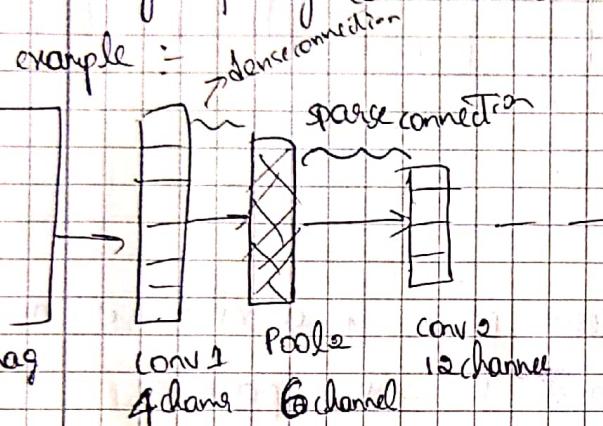
Inception V3 :-

- Avoid representational bottlenecks in early part of network
- Higher dimensional representations are easier to process locally within n/w
- Spatial aggregation can be done over lower dimensional embeddings w/o loss of representational power
- Balance width and depth of n/w

Inception network :- general idea & problems its solving (concrete own expression  
 Problems its Tackling :- Inception net - Version 1  
 Network with too much capacity is prone to overfit (like VGG). Hence  
 bigger n/w might not always be better choice for cases with  
 less training data

→ Parameters not used efficiently and wastage of compute resources

Ideas :- Using sparsely connected architectures even inside convolution (See LeNet paper)

For example :-   
 Imag → conv1 (4dans) → Pool2 (6 channel) → sparse connection → conv2 (12 channel)

	0	1	2	3	4	5	6	7	8	9	10	11
Pool2	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
conv2	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
channel	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

### \* Motivation:

→ As we can see that not every output of pool2 is utilized by conv2 and selection is mostly random (3, 4, all 6)

Why? This helps to break symmetry, if everything is connected to everything then features learnt by conv2 filters may look the same & this is not useful.

→ Breaking this symmetry helps learn complementary patterns see for example (0, 0) and (1, 1) one may learn horizontal and another vertical patch

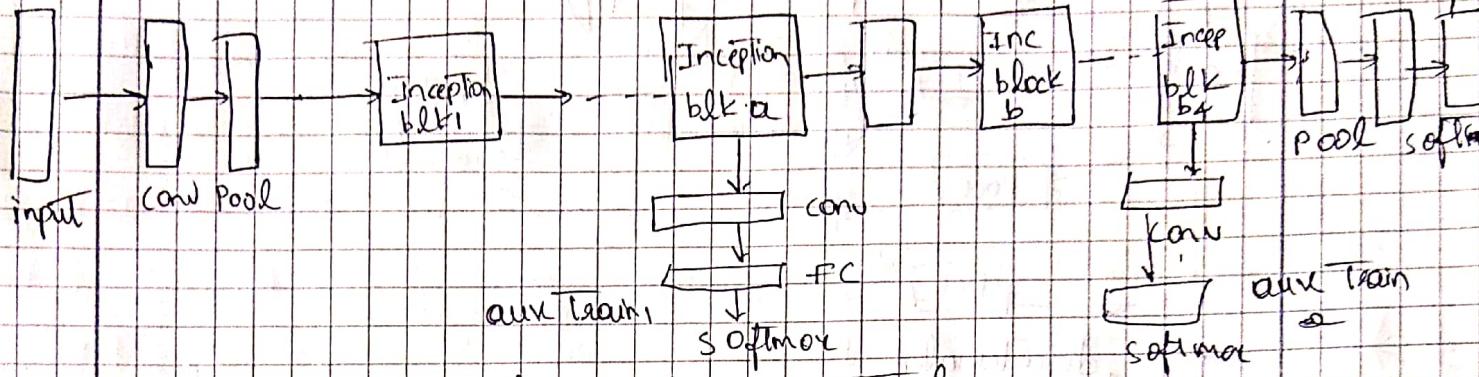
→ we can also combine all maps of previous pool layer and use it in one combo of (ii) This provides an overview

Problems with this approach :- GPUs and Tensor libraries works well and optimized for dense matrix operations

Solutions in inception n/w :- Use combination of filters and concatenate their output to mimic this sparse connectivity. This also takes care of hardware utilization problems as all operations are done on dense stuff.

→ Usage of  $1 \times 1$  convolutions so that the param size won't blow up when feature maps are passed from one Inception block to another. They also add additional non-linearity through ReLU

Training Inception net - After stacking inception modules they also had a couple of auxiliary outlets opened for training. This was removed during inference. This acted as regularizers.



- 9 Inception modules with 100 layers in total

Inception net-V2:

→ has fewer parameters than Alexnet

→ Use average pool at the end instead of an fully connected (FC) layer to keep no. of total parameters down

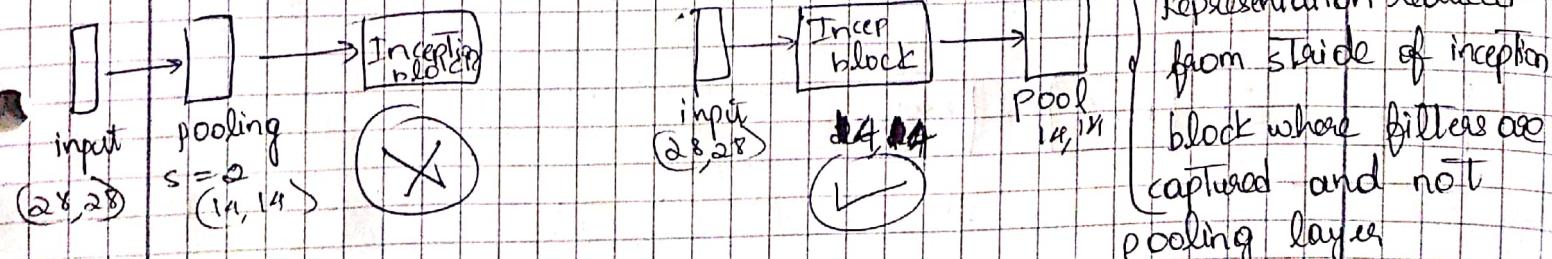
Inception net-V2: Solving the problems during Training by using Batch normalization. This makes input of intermediate layers to be in pre-expected range of near  $(0, 1)$

problems in solving

Inception net-V3: No concrete explanations on general design rules of inception net-V1 was provided in earlier paper. In this we discuss on These principles and see how can we customize them to our needs

Design principles

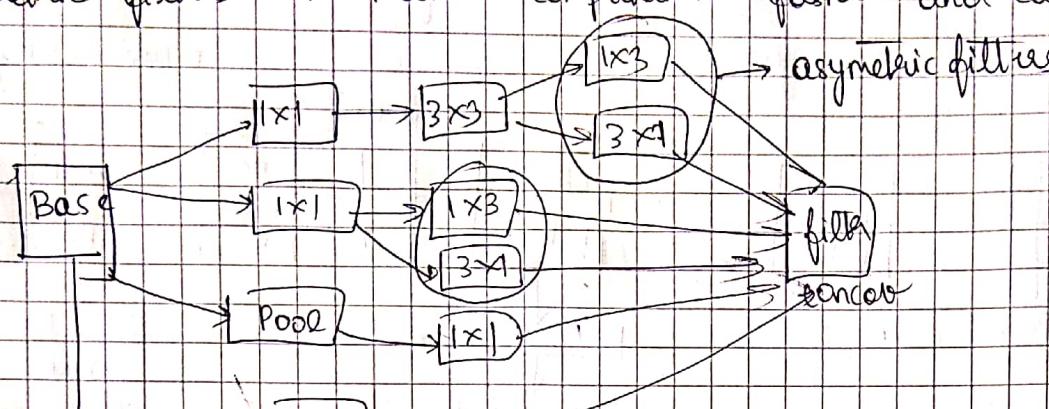
i) Avoid representational bottlenecks in early layers: Simply put it don't do Too much pooling or any other info reduction where spatial info ( $28 \times 28$ ) is suddenly decreased or halved ( $14 \times 14$ ) in early layers. Give n/w time to learn filters on high dim structure and decrease gently.



2) Higher dimensional representations are easier to process locally within n/w :-

→ Simply put make or process any spatial related info inside Inception block and not outside  $\rightarrow n \times n \rightarrow (1 \times n)(n \times 1)$

→ Use asymmetric filters To make computation faster and eat less memory



→ Capturing more activations using different filter structure leads to more discriminative featurmaps and distinct ; Helpful to add variety & train

Ex:  $1 \times 3 \rightarrow$  detect and output feature map that is horizontal

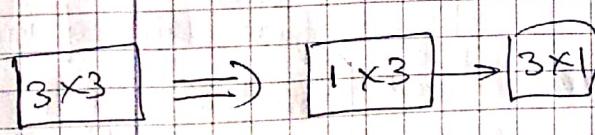
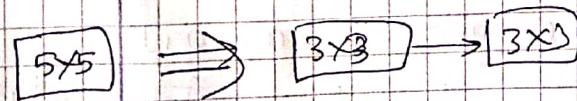
$3 \times 1 \rightarrow$  — " — " — vertical

$3 \times 3 \rightarrow 1 \times 3$  } Already  $3 \times 3$  would have subsampled and lowered resolution  
 $\rightarrow 3 \times 1$  } so applying another set of filters would provide blurred image  
lower res featurmaps That may be helpful in some case

3) Spatial aggregation can be done with over lower dim embeddings w/o much loss in representation power

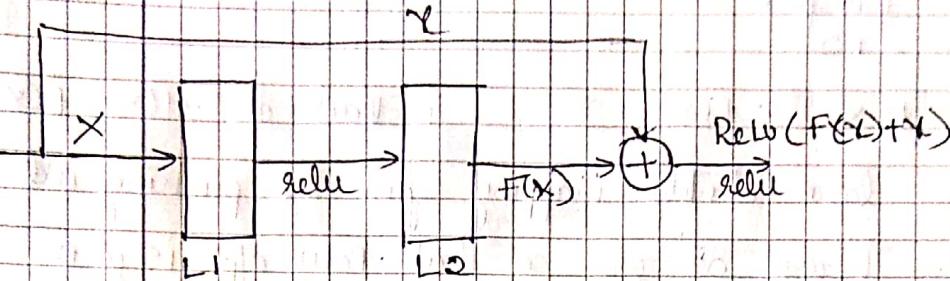
→ Simply put using Two  $(3 \times 3)$  filters instead of one  $(5 \times 5)$  would not harm the process

→ To take a step further apply general factorization and factorize  $3 \times 3 \rightarrow (1 \times 3) \rightarrow (3 \times 1)$  This reduces parameters drastically and doesn't impact learning



A) Balance depth and width of network : This is simply number of inception blocks, filter size and concatenated maps etc. like hyper parameters should be chosen wisely

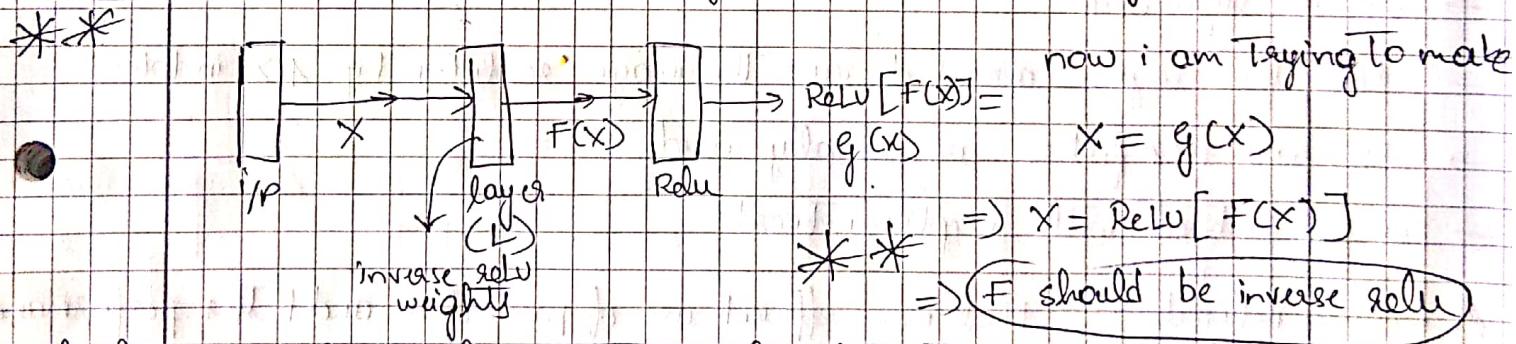
### Residual networks



Problems in tackling : Surprisingly deeper networks that are too deep resulted in higher training error.<sup>1</sup> This was because, in extreme deep networks, some intermediate layers may not learn anything new, but try to pass the info it got from previous layer unadulterated.

→ Here lies the problem, once non-linearity is applied in layer, it has to learn weights such that it reverses the non-linearity as we just need to pass the signal and not learn anything. This process of adjusting weights are difficult.

For example - let's say i want my final output to be just  $x$  i.e.  $g(x) = x$



→ if layer (L) learns the weights such that it can invert feature ReLU operation Then signal  $x$  would be passed as is w/o much changes

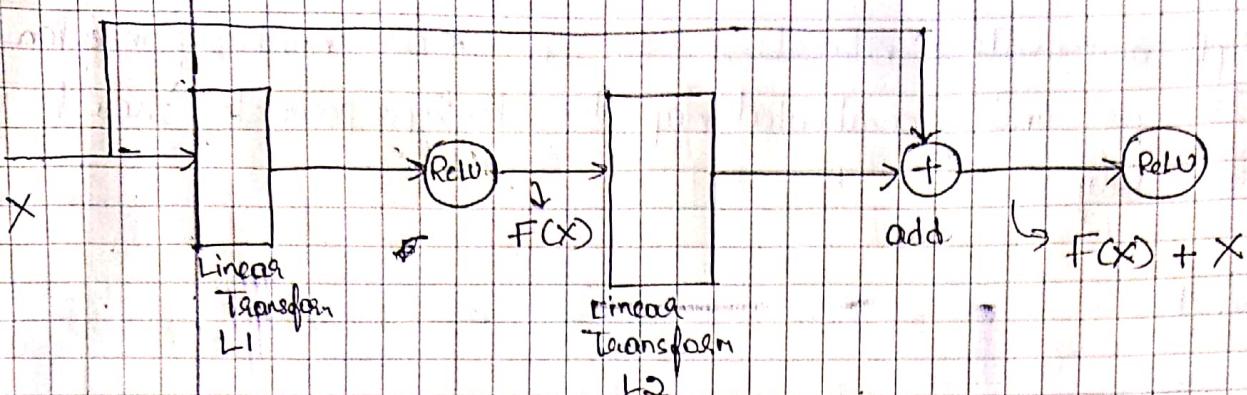
$$\Rightarrow x = \text{ReLU}[F(x)] = \text{ReLU}[\text{ReLU}^{-1}(x)] = x$$

→ Often in reality this is difficult and we need different structures

\* Solution by ResNet : Reformulate the problem, instead of learning weights

To make inverse ReLU, learn weights to approach a constant value so learning during these identity phase could be better.

→ Also provide an option to learn other complex func instead of identity when necessary



→ Learn residual function  $h(x) = F(x) - x$  instead of  $h(x) = f(x)$

This way, it is easier to learn identity mapping just by pushing all weights of  $L1 \neq$  close to zero or zero so that ReLU o/p is zero hence final o/p is just  $X$ . In case more complex func needs to be learnt then learn different non-zero weights in residual layer "L1"

\* → The above is easily realized by providing a short cut connection from input to addition layer. As a side effect, this also improves backprop as there is no challenge/obstacle for gradients to travel in the short cut path, hence even initial layers can be trained easily with less gradient vanishing problem

### Typical practical ResNet

→ 152 layer deep

→ First two layers almost decrease the image resolution by 4x in total

→  $3 \times 3$  conv + 64 filters are normally used

→ Can be combined with inception network

Inception network : Much more efficient use of parameters and better performance

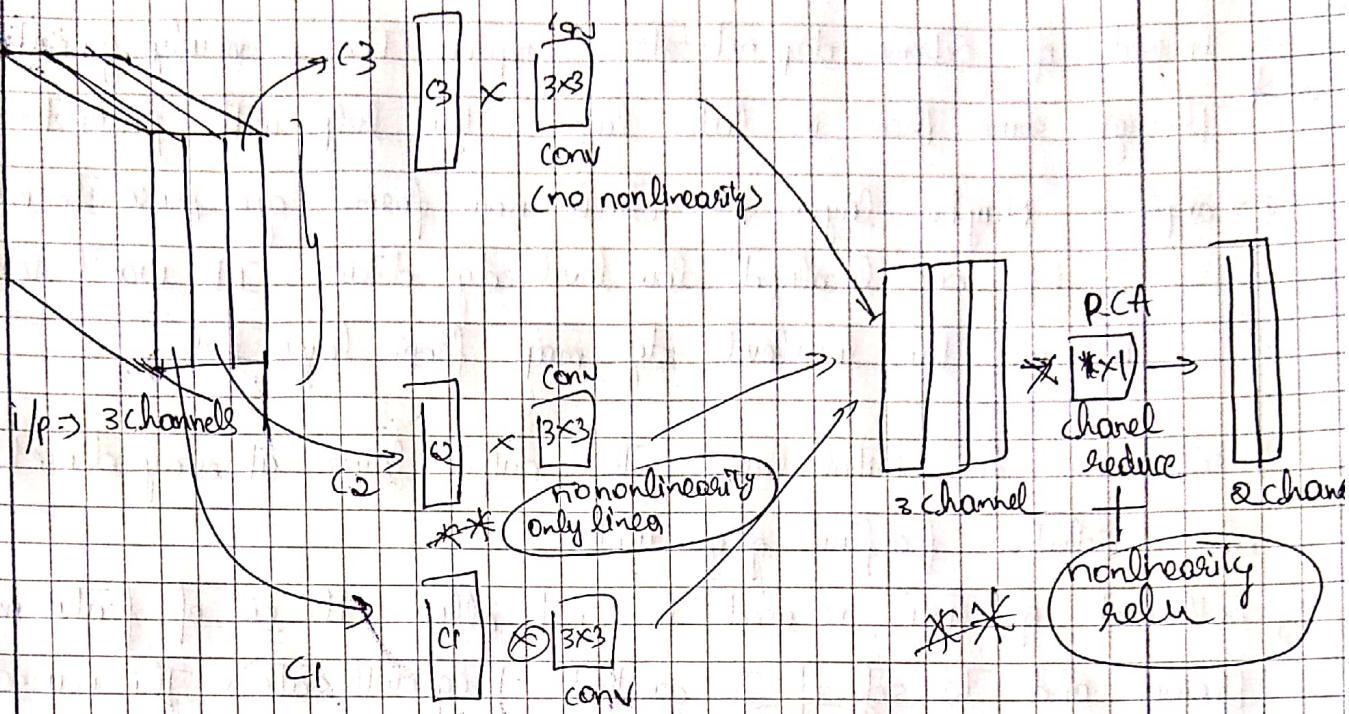
Than Inception-V3. Extreme version of inception network

Principle: In inception net we tried to somewhat decouple them by using multiple filters in  $L1, L2$  and concatenating all output. Here we take this to extreme and do channel by channel convolution using separate kernels of same size

→ ReLU is applied to the output of these operations and concatenated and reduced by  $(1 \times 1)$  PCA type of operation. (Later we apply ReLU activation).

(and not in between doing depthwise convolution)

out -

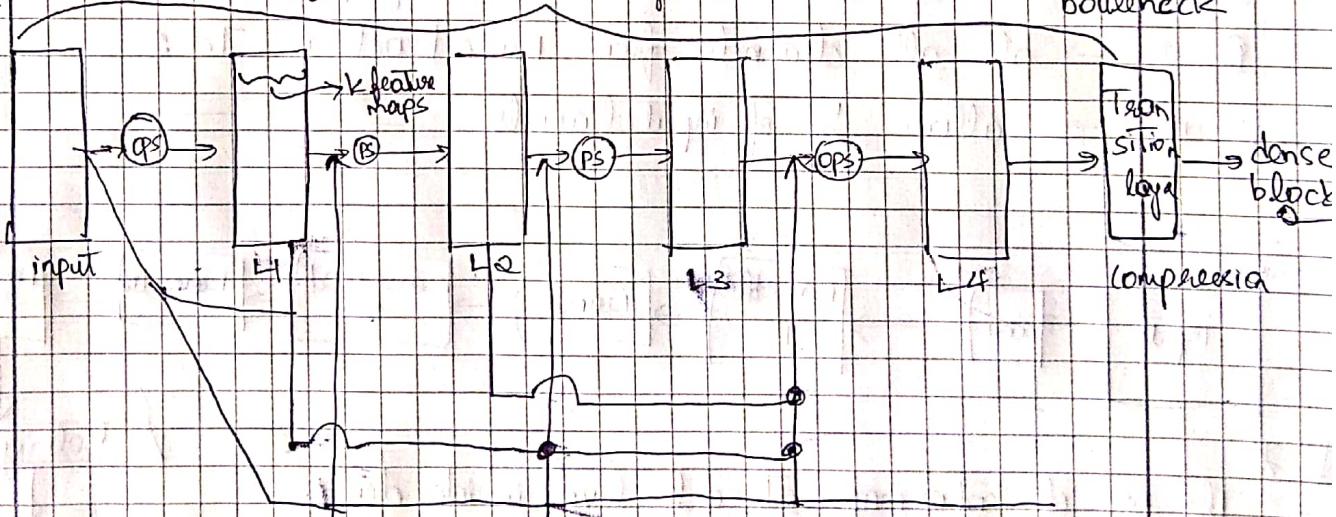


DenseNet : Prob

Problems it's tackling : Nothing serious in particular, but more optimization of network capacity and feature reuse promotion and increase info flow w/o much adulteration.

Solutions :

Dense block  $\leftarrow$   $k=4$  (growth rate)  $\text{Ops} = \text{Convolution} + \text{ReLU} + \text{bottleneck}$



$$\rightarrow \text{No of direct connections} = \frac{L * (L+1)}{2} \quad \text{where } L \text{ is no of layers inside dense block}$$

$$\rightarrow \text{Input To } l^{\text{th}} \text{ layer} = X_l = f(l)([X_0, X_1, \dots, X_{l-1}])$$

where each layer produces  $K$  feature maps

$$\rightarrow \text{Total feature maps} = \underbrace{K_0}_{\text{Input maps}} + \underbrace{K * (l-1)}_{\text{growth rate}} \underbrace{\text{no of layers}}$$

## Highlights

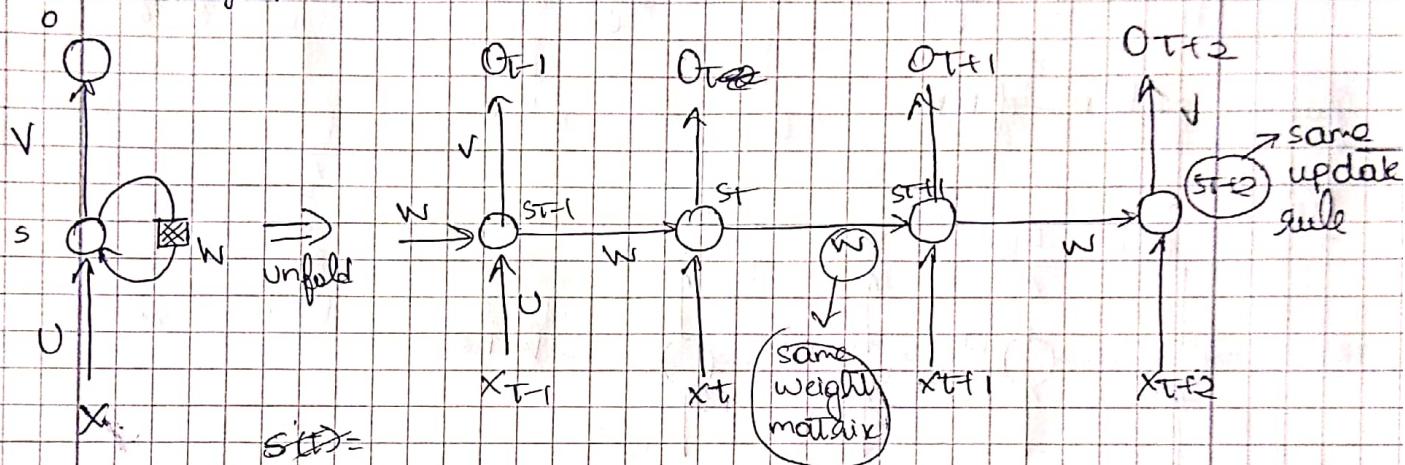
- Instead of extreme deep networks, exploit the previously generated features through reuse them in later layers. This helps with parameter efficiency.
  - \* Example: Maybe layer L3 which now forms eyes/nose etc need some info on localized low level edge details. It won't recreate this but use this low level edge maps from layer L1.
  - Direct access for initial layer to output layer at every dense block allows relative free flow of gradients.
  - Bottleneck of  $(\times 1)$  is used to not allow the no. of feature maps to blowup and to set it to constant  $k$  (growth rate). This is used <sup>inside</sup> between each layer <sub>within</sub> dense block.
  - Transition layers consist of convolution + compression + Pooling to control amount of info exchanged b/w consecutive dense blocks. Again use dynamic number of  $(\times 1)$  conv channels using parameter  $\theta$ 
    - $\theta = 1 \rightarrow$  all all channels from Dense block 1 to 2
    - $\theta = 0.5 \rightarrow$  allow  $\frac{1}{2}$  of input channels
  - \* The no. of feature maps at each level inside dense block is controlled by growth rate  $k$  which is kept less, around typically  $k=12$ . Hence less params are required.
- 
- The diagram illustrates the architecture of a DenseNet. It starts with an 'Input' block, followed by a 'conv' layer. This leads to a sequence of layers: 'Dense block 1', 'Transition', 'Dense block 2', 'Transition', and finally 'Dense block 3'. A 'Pooling' layer is shown below 'Dense block 3', and a 'Linear' layer is shown below that, both with downward arrows indicating their respective outputs. The connections between layers are dense, with skip connections from every layer to every subsequent layer in the sequence.

## RNN - Block A

Motivation: Process sequential data and also data for whom order is important

- Model for which historical states are required
- Output at current step is a function applied to previous step off
- produced using some update rule applied in previous step

### Computational graph



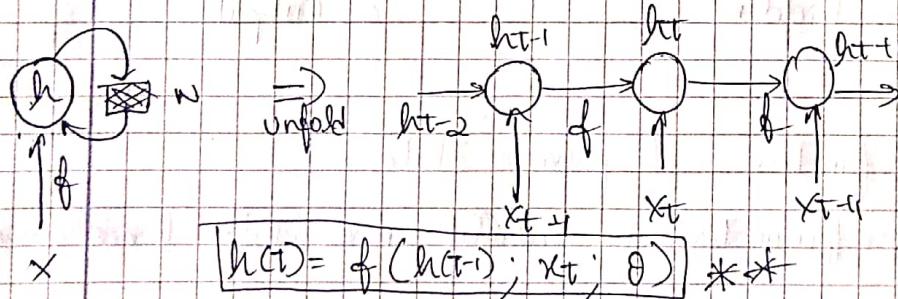
→ Regardless of sequence length dimensionality of input remains constant

### Back Propagation Through Time

Advantage  
→ gradient computation in unfolded n/w w/ mat parameters are expensive

Disadvantage  
→ NO parallelization

→  $O(T)$  where  $T$  is history length complexity



### 1-d Convolution v/s RNN

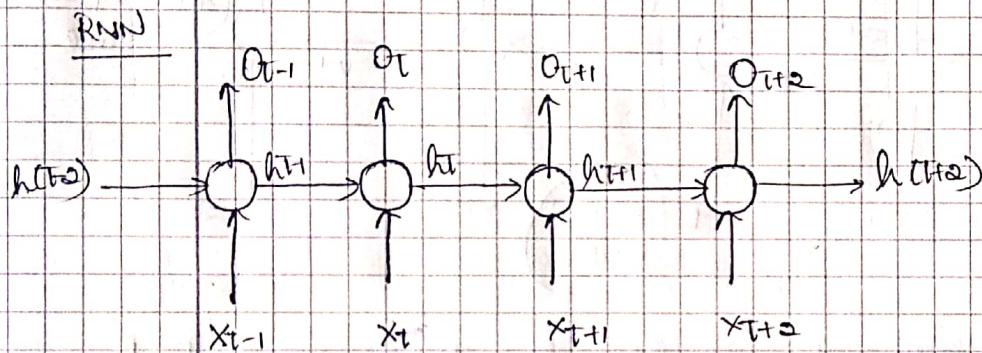
→ 1-D convolution can be easily used to model Timeseries data or speech recognition, wherever there is a fixed length analysis is needed

For example: In Time series, CNN can learn seasonal customer patterns since customer sales patterns repeat year round (like sales in Christmas, summer, Easter etc). Fixed number of filters and fixed length window may be sufficient. Customer behaviour from summer will change in winter hence unlike RNN where all history is used to make decision we can make a quick pattern based matching decision in CNN

similar is the case for audio recording as always, some phenomena tend to occur together. This feature can be position invariant and a pattern.

Where it doesn't work? Even if we increase receptive field of one-d CNN to include large history by stacking & dilation. The amount of receptive field is still finite. Technically, also CNN is a pattern matcher where it tries to sense pattern filters in any position/order. RNN will still work even if there is no visible pattern or the patterns are extremely unique & non repetitive & might not occur again (like some language modelling task).

### gradient problems of RNN



Recurrent neural networks  $\vdash$  A class of neural networks that are mainly used to process sequential data. It consists of three states:

- a) Input
- b) Hidden
- c) Output

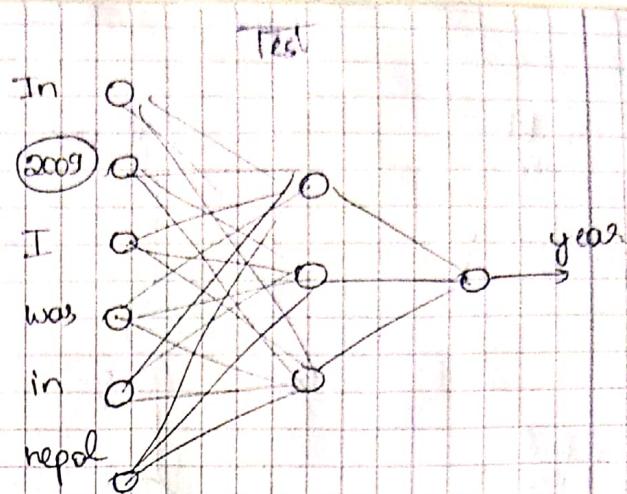
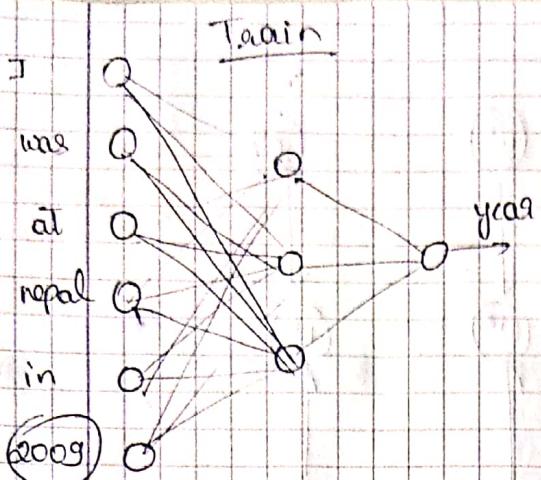
### Properties of RNN

- a) Current state is a function of previous state
- b) We use the same function (parameters) to calculate next states, hence parameters are shared
- c) Regardless of input size, the hidden state dimensions are fixed and they keep a lossy summary of history upto that current time  $T$ .  
We can also ask to present complete i/p sequence history as it's ingested as with case of auto encoders

Parameter sharing and applying same function 'f' at each time steps

### Advantages and Disadvantages of RNN

$\rightarrow$  With just feed forward n/w, let's say we are trying to extract date

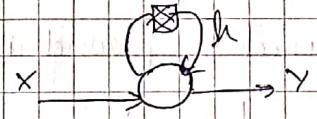


→ The position of the date in sentence matters as weights are learnt separately for each Token and if position changes in Test, it doesn't generalize well as during Train we focused on last Token.

→ To remedy this we need to Train on all possible combos of sentence so that network figures out all rules. This is a disadvantage

→ RNN's parameter sharing & some function helps so that you don't need to learn separate rules of language applied to different position. Everything is encapsulated in hidden state which is kind of summary

→ Once only with folded RNN



There is a disadvantage that the input cannot be defined over constant length of history

$$h(t) = g(t)(x(1), x(2), \dots, x(t-1))$$

$$h(t-1) = g(t-1)(x(1), x(2), \dots, x(t-2))$$

$$h(t+1) = g(t)(x(1), x(2), \dots, x(t))$$

Problem: each Time step  $g(t)$  is looking on variable length of input

This may lead to  $g(t)$  trying to learn a different rule at each time step completely

With unfolding we can define current state using only its fixed previous states. This allows for more generalization & faster training

$$h(t) = f(h(t-1), x(t); \theta)$$

$$h(t-1) = f(h(t-2), x(t-1); \theta)$$

$$h(T+1) = f(h(T), x(T+1); \theta)$$

One i/p & one hidden state

→ Since we have chosen to represent current state only with respect to some defined previous state and not looking at all inputs at every time, gradient propagation becomes problematic and expensive. Even forward prop takes more time ( $O(2) \approx \text{length of sequence}$ )

For ex: compare two cases below

$$h(t-1) = g_1(t)(x^{(1)}, x^{(2)}, \dots, x^{(t-2)})$$

$$h(t) = g_2(t)(x^{(1)}, x^{(2)}, \dots, x^{(t-1)})$$

$$h(t+1) = g_3(t)(x^{(1)}, x^{(2)}, \dots, x^{(t)})$$

Rolled n/w

$$h(t-1) = h(t)f(h(t), x(t), \theta)$$

$$h(t) = f(h(t), x(t), \theta)$$

$$h(t+1) = f(h(t), x(t+1), \theta)$$

unrolled n/w

→ Direct access to inputs, so better gradient propagation

→ No visible dependency of  $g_3(t)$  to  $g_1(t)$  or  $g_2(t)$ , since we would already have complete input sequence we can compute all  $h(t-1), h(t), h(t+1)$  in parallel

→ Need to learn different  $g_1, g_2, g_3$  for different steps, which is really difficult

→ No direct visibility from  $h(t)$  to  $h(t-1)$  or  $h(t)$  should always go through  $h(t-1)$

→ Sequential processing of info on both forward and backward pass as  $h(t)$  cannot see  $h(t-2), h(t-3)$  etc & can only see  $h(t-1)$  as per defined model

→ Just need to learn  $f, \theta$  same is applied at every step

### Different gradient problems in RNN

Consider change of error in final layer  $\approx$  w.r.t first hidden state

$$\frac{\partial E^{[r]}}{\partial h^{[0]}} = \frac{\partial E^{[r]}}{\partial h^{[r]}} \prod_{t=1}^T \frac{\partial h^{[t]}}{\partial h^{[t-1]}}$$

$r \rightarrow$  depth of layers

$E \rightarrow$  error

$h \rightarrow$  hidden state

$$\text{WKT } h^{[t]} = \tanh(a^{[t]})$$

$$a^{[t]} = Wx^{[t]} + Wh^{[t-1]}$$

$$\frac{\partial \tanh \theta}{\partial \theta} = 1 - \tanh^2 \theta$$

↓  
Vanishing  
gradient

Prerequisite

### 1) Bounding of derivatives of $\sigma$ , tanh

WKT range of  $\sigma \in [0, 1]$  &  $\tanh \in [-1, 1]$

$\frac{d\sigma(x)}{dx}$  range would be  $\approx [0, \frac{1}{4}]$

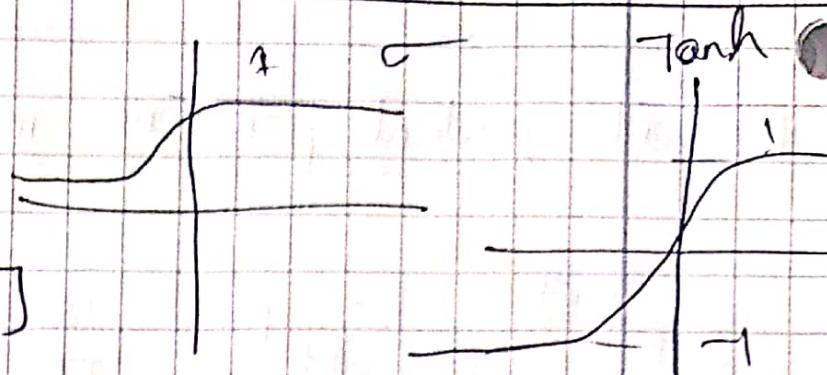
why?  $\frac{d\sigma(x)}{dx} = \sigma(x)(1 - \sigma(x))$  is max at  $\sigma(x) = \frac{1}{2}$

$$\Rightarrow \left( \frac{d\sigma(x)}{dx} \right)_{\max} = \frac{1}{4} \quad * * *$$

Why for tanh

$$\left| \frac{d\tanh x}{dx} \right|_{\max} = 1 \Rightarrow \text{rangebound} = 1$$

why:  $\frac{d\tanh x}{dx} = 1 - \tanh^2 x = 1$



## II) Eigenvalue decomposition and Powerlaw



→ Any weight matrix  $W$  or  $U$  or  $V$  can be decomposed into eigen vectors & values

$$W = M \text{ diag}(\lambda) M^{-1} \quad \text{let } \lambda = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & -3 \end{bmatrix} \text{ for example}$$

→  $-3$ ' is dominant eigen vector by magnitude

$$W^2 = M (\text{diag}(\lambda))^2 M^{-1}, \quad W^3 = M (\text{diag}(\lambda))^3 M^{-1} \Rightarrow$$

$W^T = M (\text{diag}(\lambda))^T M^{-1} \Rightarrow$  Why?  $\lambda$  eigen value act as scaling factor to keep eigen vector same

Hence  $W^T \lambda^T = \begin{bmatrix} 1^T & 0 & 0 \\ 0 & 2^T & 0 \\ 0 & 0 & -3^T \end{bmatrix} \Rightarrow$  we can ignore  $1 \& 2$  & can get away by only considering principal vector  $(3)^T$

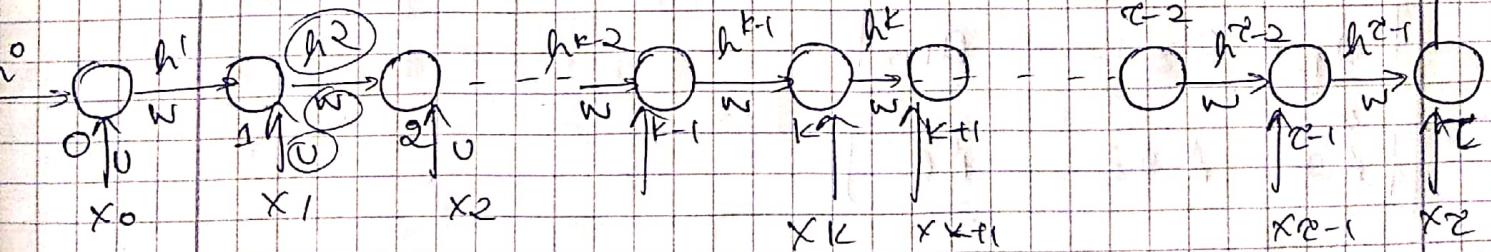
## III) Chain rule and gradient descent

We always try to learn parameters through loss function, because this is how we learn stuff  $\Rightarrow \frac{\partial \text{Out}}{\partial W}$  or  $\frac{\partial \text{Out}}{\partial \text{Bias}}$

## IV) Deriving update rule for BPTT

dependency from  $\epsilon$  to  $k-1$

The quick brown fox jumps over a lazy dog



→ let's say Task is to remember animals Then we need to back propagate from last layer  $\epsilon$  to required layer ' $k-1$ ', it can be  $0^{\text{th}}$  layer if you want too.

My requirement :  $\frac{\partial \text{Loss}}{\partial \mathbf{N}} = ?$

$$\hookrightarrow \text{Apply chain rule} \Rightarrow \frac{\partial \text{Loss}}{\partial h^T} \cdot \frac{\partial h^T}{\partial h^{T-1}} \cdot \frac{\partial h^{T-1}}{\partial h^{T-2}} \cdots \frac{\partial h^k}{\partial h^{k-1}} \cdot \frac{\partial h^{k-1}}{\partial h^{k-2}} \cdots$$

wrt  $h^{k-1} = W h^{k-2} + U x^{k-1}$

$$\frac{\partial h^{k-1}}{\partial W} = \text{constant}$$

$$\Rightarrow \frac{\partial \text{Loss}}{\partial h^T} = \frac{\partial \text{Loss}}{\partial h^2} \cdot \frac{\partial h^2}{\partial h^{2-1}} \cdots \frac{\partial h^{k-1}}{\partial h^{k-2}}$$

\*\*\*

$$\frac{\partial h^T}{\partial h^{T-1}} = \frac{\partial \text{Loss}}{\partial h^2} \frac{\partial h^2}{\partial h^{2-1}} \cdots \frac{\partial h^{T-1}}{\partial h^{T-2}}$$

Consider one step of backprop  $\frac{\partial h^T}{\partial h^{T-1}}$  or  $\frac{\partial h^{T-1}}{\partial h^{T-2}} = ?$

wrt  $h^T = \text{activation } (\mathbf{a}^T) = \text{ReLU } A(\mathbf{a}^T)$

$$\mathbf{a}^T = W h^{T-1} + U x^T$$

$$\Rightarrow \frac{\partial h^T}{\partial h^{T-1}} = \frac{\partial h^T}{\partial \mathbf{a}^T} \cdot \frac{\partial \mathbf{a}^T}{\partial h^{T-1}} \quad \left| \begin{array}{l} \text{if by } \frac{\partial h^{T-1}}{\partial h^{T-2}} = \frac{\partial h^{T-1}}{\partial \mathbf{a}^{T-1}} \cdot \frac{\partial \mathbf{a}^{T-1}}{\partial h^{T-2}} \\ \dots \end{array} \right.$$

activations are vectors hence (Remember, we ~~set~~ set hidden state to size 1024 size during assignment)

$$\mathbf{a}^T = [a_{T1}, a_{T2}, \dots, a_{Tr}]$$

$$\text{if by } h^T = [A(a_{T1}), A(a_{T2}), \dots, A(a_{Tr})]$$

We need to find rate of change of each individual  $h^T$  v/s all vector  $\mathbf{a}^T$   
hence Take Jacobian

For example :

$$\frac{\partial h^T_1}{\partial a^T} = \left[ \frac{\partial A(a_{11})}{\partial a_{11}}, \frac{\partial A(a_{12})}{\partial a_{11}}, \frac{\partial A(a_{13})}{\partial a_{11}}, \dots, \frac{\partial A(a_{1n})}{\partial a_{11}} \right]$$

zeros

$a_{13} \neq a_{11}$  (not dependent / related to denominator, so constant numerator)

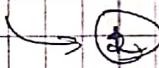
$$\frac{\partial h^T_2}{\partial a^T} = \left[ \frac{\partial A(a_{21})}{\partial a_{21}}, \frac{\partial A(a_{22})}{\partial a_{21}}, \frac{\partial A(a_{23})}{\partial a_{21}}, \dots, \frac{\partial A(a_{2n})}{\partial a_{21}} \right]$$

zero

zero

$$\Rightarrow \begin{bmatrix} \frac{\partial h^T_1}{\partial a^T} \\ \frac{\partial h^T_2}{\partial a^T} \\ \vdots \\ \frac{\partial h^T_n}{\partial a^T} \end{bmatrix} = \begin{bmatrix} \frac{\partial A(a_{11})}{\partial a_{11}} & 0 & 0 & \dots & 0 \\ 0 & \frac{\partial A(a_{22})}{\partial a_{22}} & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \frac{\partial A(a_{nn})}{\partial a_{nn}} \end{bmatrix}$$

diagonal matrix



$$\text{Hence for } \frac{\partial a^T}{\partial h^{T-1}} = \text{diagonal } \frac{\partial \{W h^{T-1} + U x^T y\}}{\partial h^{T-1}} = W - \textcircled{2}$$

partial derivative short notation

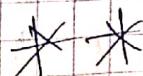
$$\Rightarrow \frac{\partial h^T}{\partial h^{T-1}} = \text{diagonal } \underbrace{\text{Act}^T(a^T)}_{\text{Act}^T(a^T) y^T W}$$

$$\left\| \frac{\partial h^T}{\partial h^{T-1}} \right\| = \left\| \text{diagonal } \{ \text{Act}^T(a^T) y^T W \} \right\|$$

magnitude

from  $\textcircled{II}$  if activation  $\text{Act} = \sigma$  then  $\text{Act}^T \approx [0, \frac{1}{4}]$   $= V$   
 activation =  $\text{Tanh}$  then  $\text{Act}^T \approx [0, 1]$

$$\| \text{grad} \| = \| V W \|$$



$$\text{wkt } \|a \cdot b\| \leq \|a\| \|b\| \Rightarrow \| \text{grad} \| \leq \|V\| \|W\|$$

Use ① and replace  $\|w\|$  by its eigen equivalent

$$\Rightarrow \|\text{grad}\| = \|V\| \|M \lambda N^{-1}\| = \frac{\partial h^T}{\partial h^{T-1}}$$

Now  $\frac{\partial h^{T-1}}{\partial h^{T-2}} = \|V\| \|M \lambda N^{-1}\|$  dominant eigen vector

$$\Rightarrow \left[ \begin{array}{c} \vdots \\ \frac{\partial h^T}{\partial h^{T-1}} \\ \vdots \\ \frac{\partial h^{T-K}}{\partial h^{T-1}} \end{array} \right] = (\lambda \cdot V)^T \|N \lambda^{-1}\|$$
 upper bound of activation derivative

### Vanishing gradient problem

if  $\lambda \cdot V < 1$  Then for significantly large  $T \ll K$   
gradients may vanish

$$\text{ex: } (0.4 * 0.9)^{100-10} = (0.36)^{90} \approx 0$$

final step      intermediate

if  $\lambda \cdot V \geq 1$  Then for significantly large  $T \gg K$

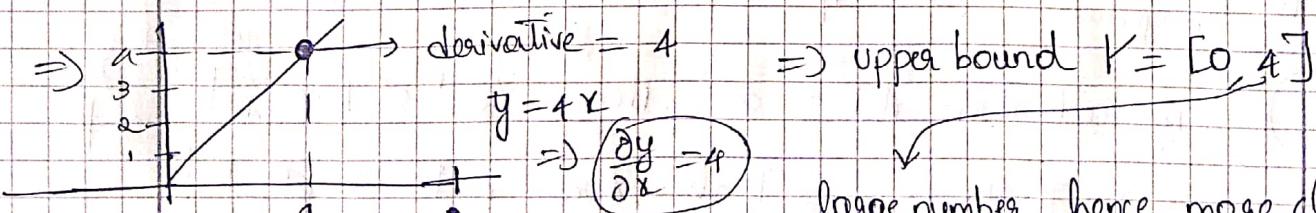
gradients explode

$$(0.4 * 4)^{100-10} = (1.6)^{90} \gg 1$$

15<sup>90</sup>

### Why it's difficult to use ReLU in RNN?

ReLU scales linearly for positive quadrant hence its derivative = scaling factor



$(\lambda \cdot 1)^{T-K}$  becoming huge.

large number, hence more chances of

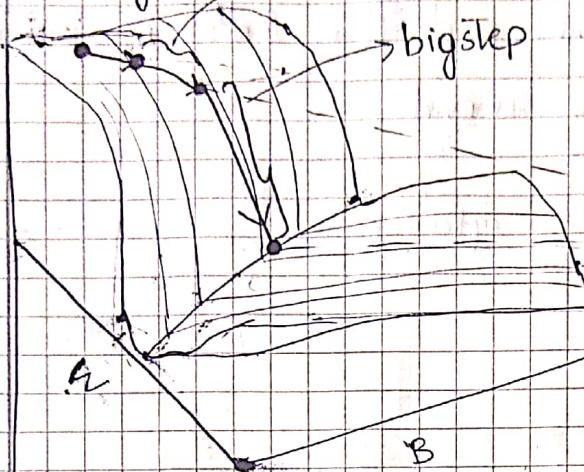
\* Here notice that gradient problem is due to two reasons

- a) choice of activation func
- b) Repeated multiplication of weight matrix with itself

## Remedies

Exploding gradient :- Gradient clipping, activation function, no of layers

Gradient clipping :- small step



Loss  $J(w, b)$  → Due to big step loss decreased suddenly  
→ caused large gradients and caused parameters to jump out of space during gradient descent

new weight after a gradient update  $w_i = w_i - \alpha \nabla w_i$  → Makes learning unstable and (out of space) we lose out on our efforts we put for optimization earlier

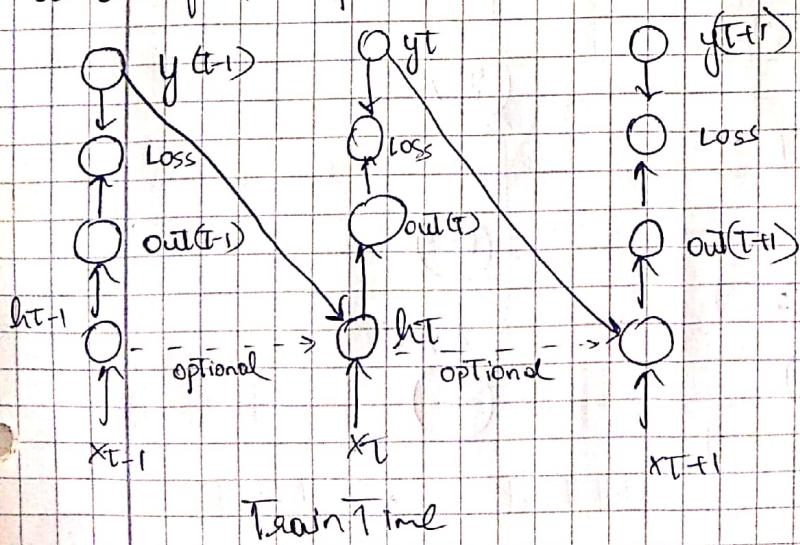
Vanishing gradient :- Activation function, gating, initialization

$w_i = w_i - \alpha \nabla w_i$   $\nabla w_i$  is too small hence no learning takes place

→ Better to have vanishing than exploding gradient as in the latter optimization is not possible and for vanishing gradient it's painfully slow.

$$h^{(T)} = \underbrace{\Theta^T}_{\text{Transpose}} \underbrace{h^{(0)}}_{\text{eigen value}}$$

Teacher forcing :- arises from NLLC criterion. Receive expected output as one of the input to hidden state

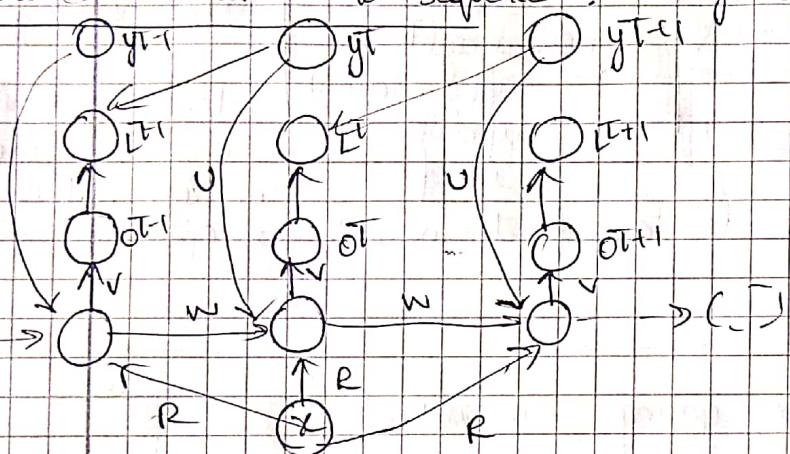


→ if no hidden-hidden connection Then BPTT not required, if connection exists then use BPTT to train  
→ Works well if "out" is high dim & rich  
→ Allows arbitrary distribution over y given sequence of x of same length

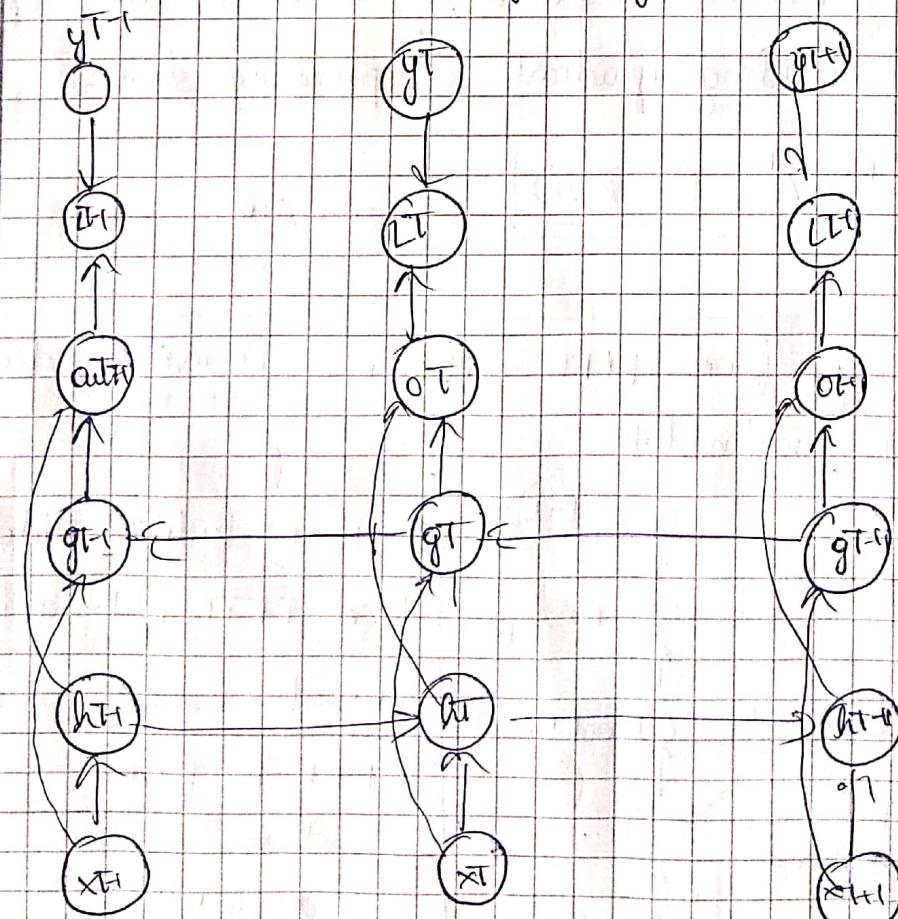
Disadvantages :- Teacher forcing was learned using expected output directly, hence it hasn't learnt much of useful things in its memory state.  $\rightarrow$  Hence during test it doesn't generalize well on unseen instances.

Remedy :- Use a type of curriculum learning, where the network is teacher forced for first few epochs and then learnt in normal RNN way with its  $O/P \rightarrow$  hidden connection removed.

Fixed context vector To sequence :- Image captioning



Bidirectional RNN :- Hand writing recognition, seq-seq models



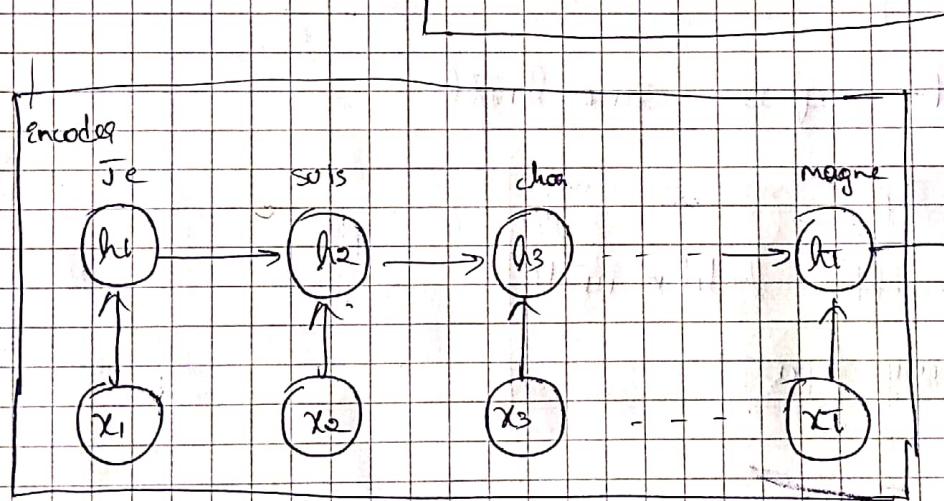
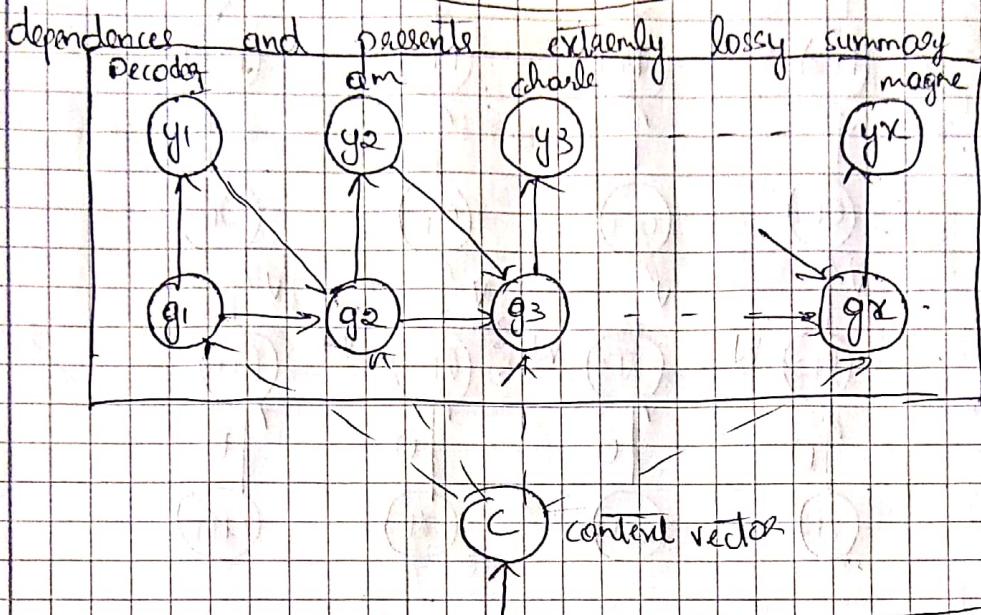
Take into account both past & future info  
 $\rightarrow$  helpful in predicting middle word of sentence

## Encoder-Decoder : Machine Language Translation

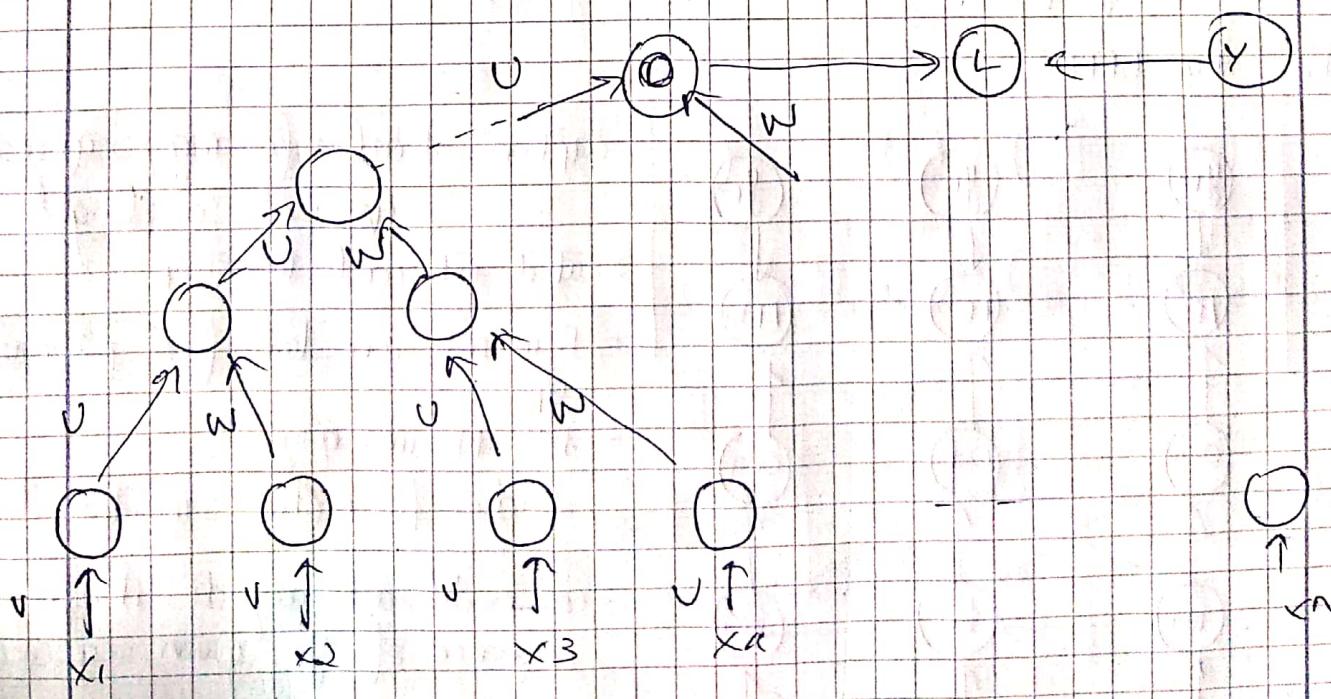
→ Variable length ip and op

bottleneck

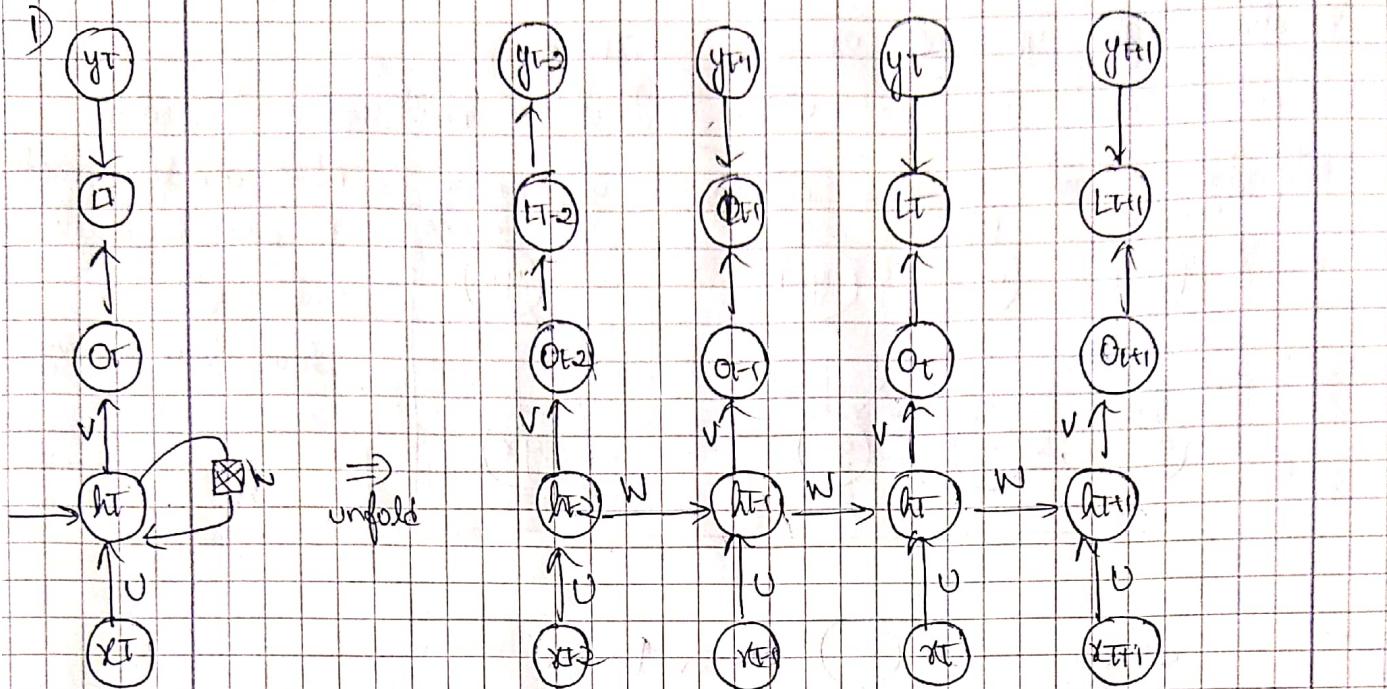
→ Problems faced in context vector (1) as it cannot capture long range dependence and parent child extremely lossy summary. Also gradient needs to flow through context vector slows down learning



Recursive Tree → Sentiment analysis, Sentence parser



Match the following : Architectures with captions



Caption : Sequence To sequence same length

→ output at each time step

→ Recurrent connection b/wn hidden units

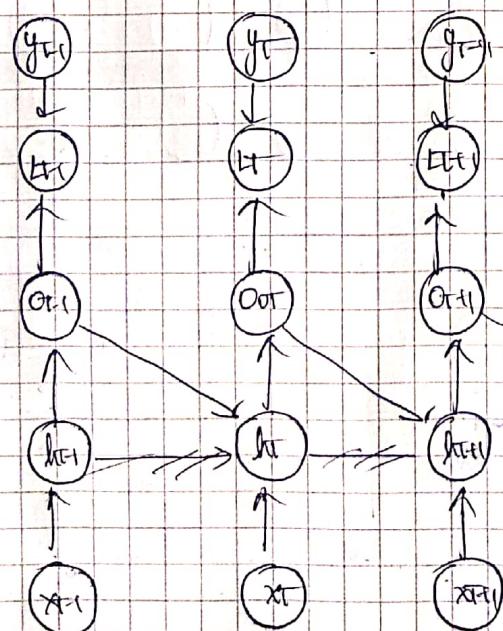
→ Only relevant summary of past

→ Sequential Training

→ can compute any func computable by Turing m/c (UAC)

→ Trainable by BPTT

E)



Caption: Teacher forcing / Sequence To sequence same length

→ output at each time step

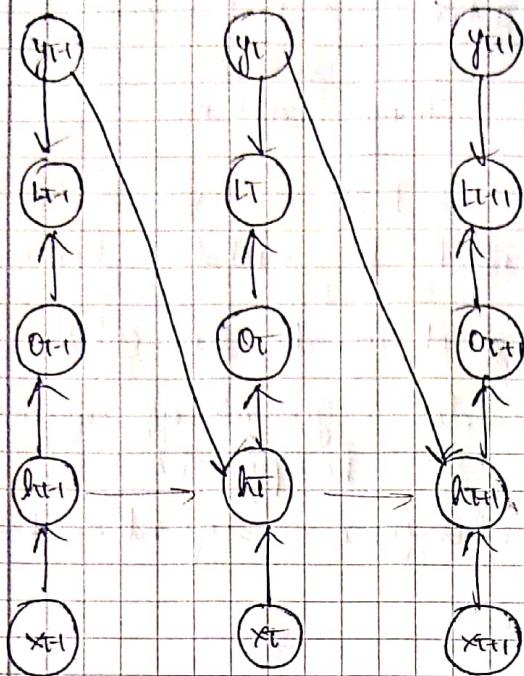
→ Recurrent connection from previous o/p

→ ~~Relevant~~ No summary

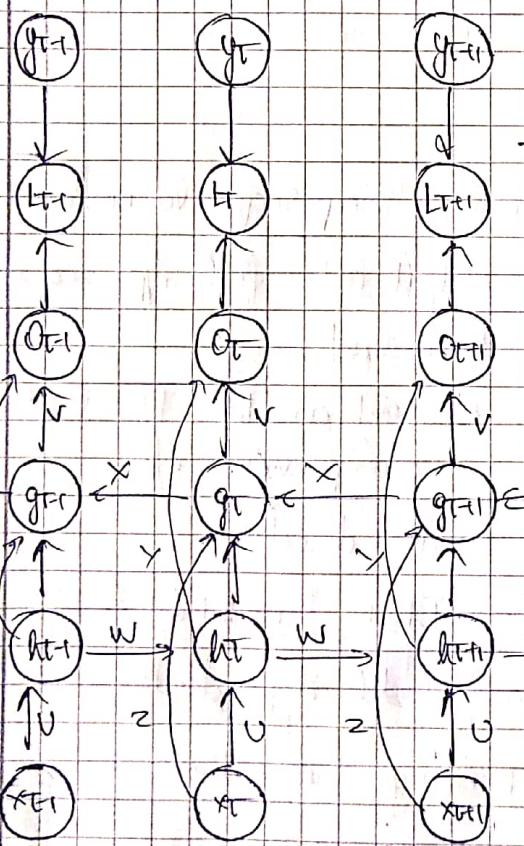
→ Teacher force (but sequential Training)

→ can model arbitrary distribution over sequences of  $y$  given sequence of  $x$

→ lacks important info from past unless  $o$  is high dim & rich



- 3) *Caption:* Teacher forced / seq-seq same length  
 → Output at each timestep  
 → Recurrent connection from both hidden & previous o/p  
 → Summary from past  
 → Parallelized Training / Teacher forced  
 → Can model arbitrary distribution over sequences of  $y$  given sequences of  $x$



- 4) *Caption:* Seq-seq bidirectional same length  
 → Output at each time step

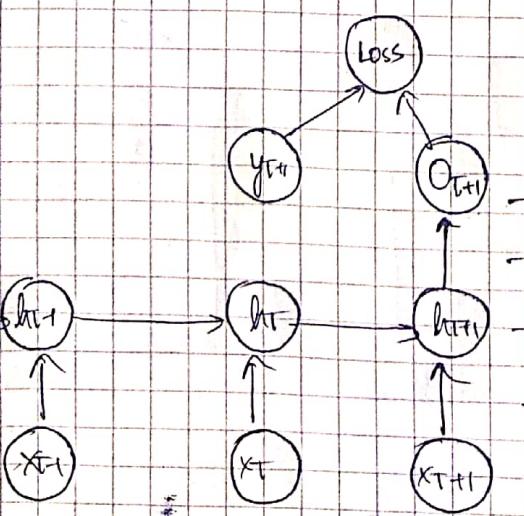
→ Recurrent connection b/wn hidden units

→ Both past and future summary

→ Sequential BPTT Training

compute any function computable by Turing M/C

past summary



- 5) *Caption:* Complex structure To fixed size vector

Text summarization

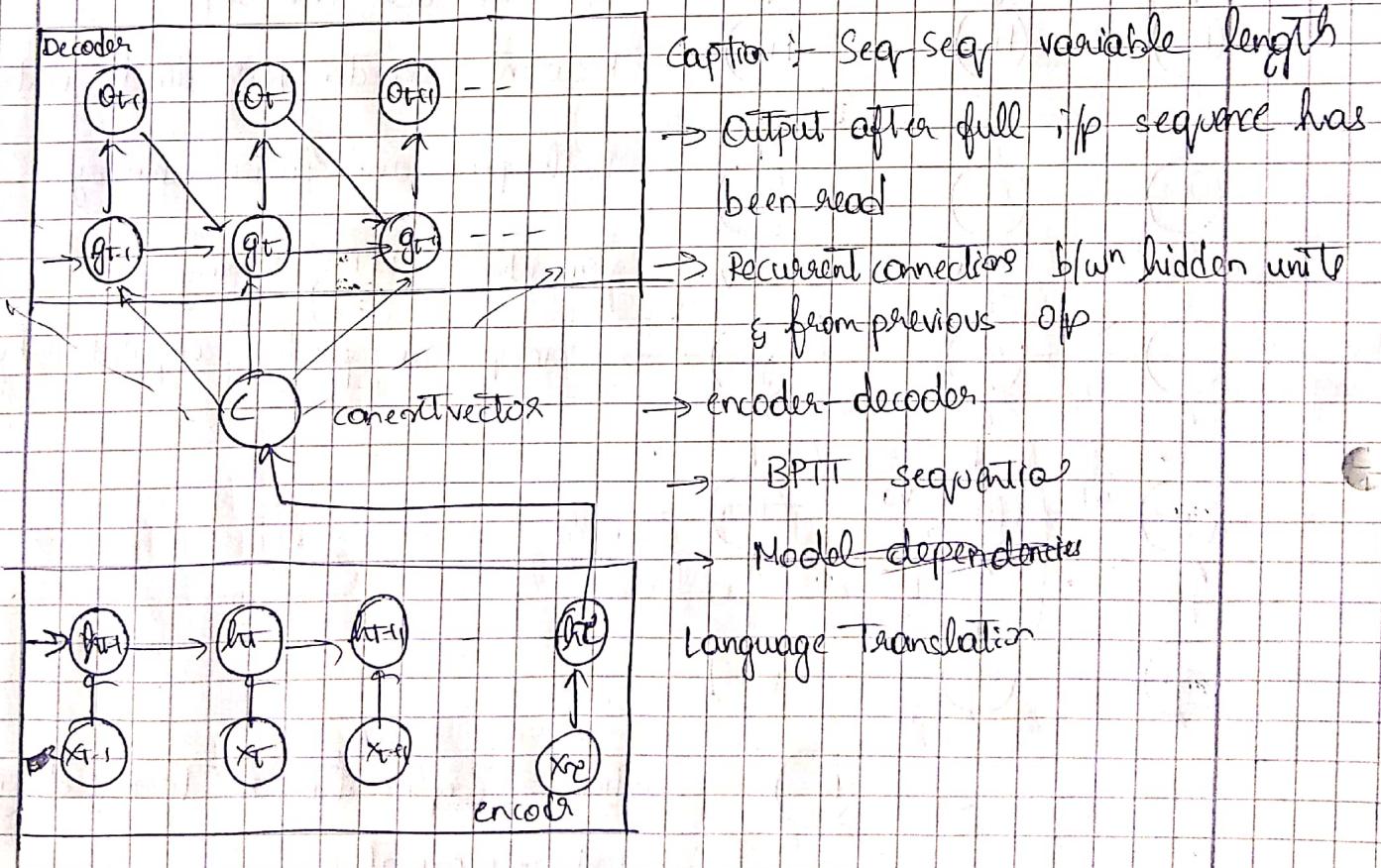
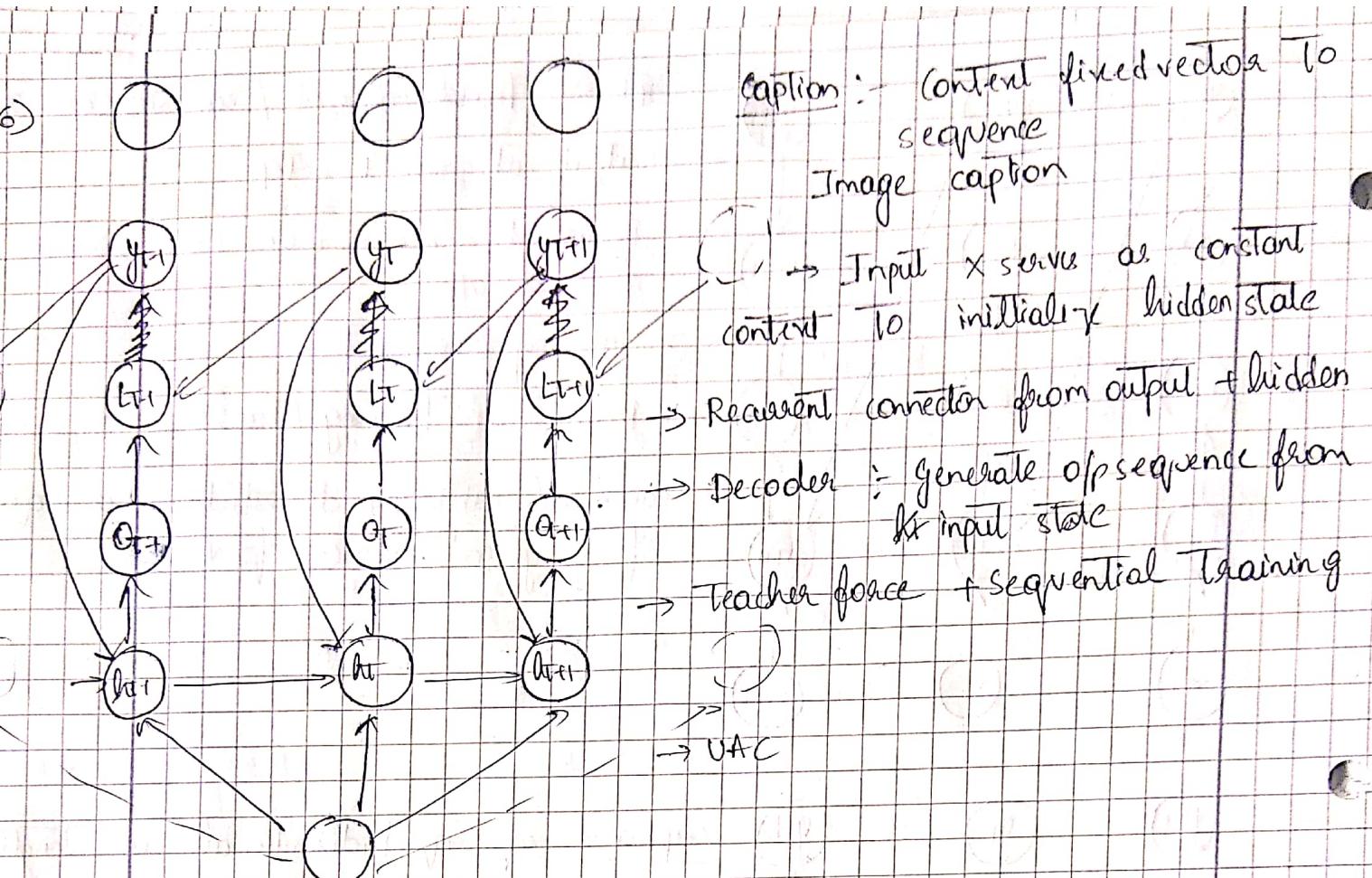
→ Output after all i/p has been read

→ Recurrent connection b/wn hidden units

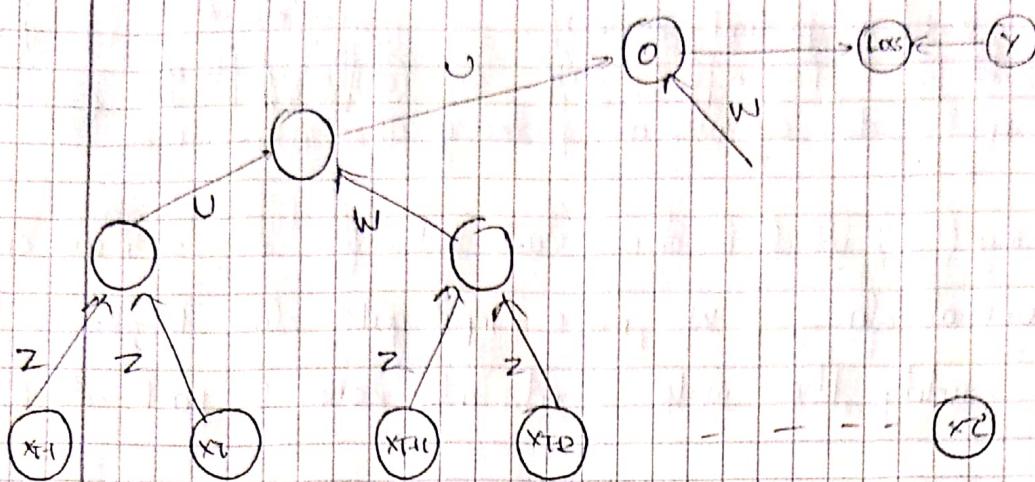
→ Only past summary

→ Sequential Training BPTT

→ Model Encoder :- generate content vector



8)



(caption): Seq - fixed size vector  $\Rightarrow$  Sentiment analysis / Sentence Pairs

$\rightarrow$  Compute graph structured as deep tree

$\rightarrow$  BPIT sequential Training

$\rightarrow$  Encoder Type architecture

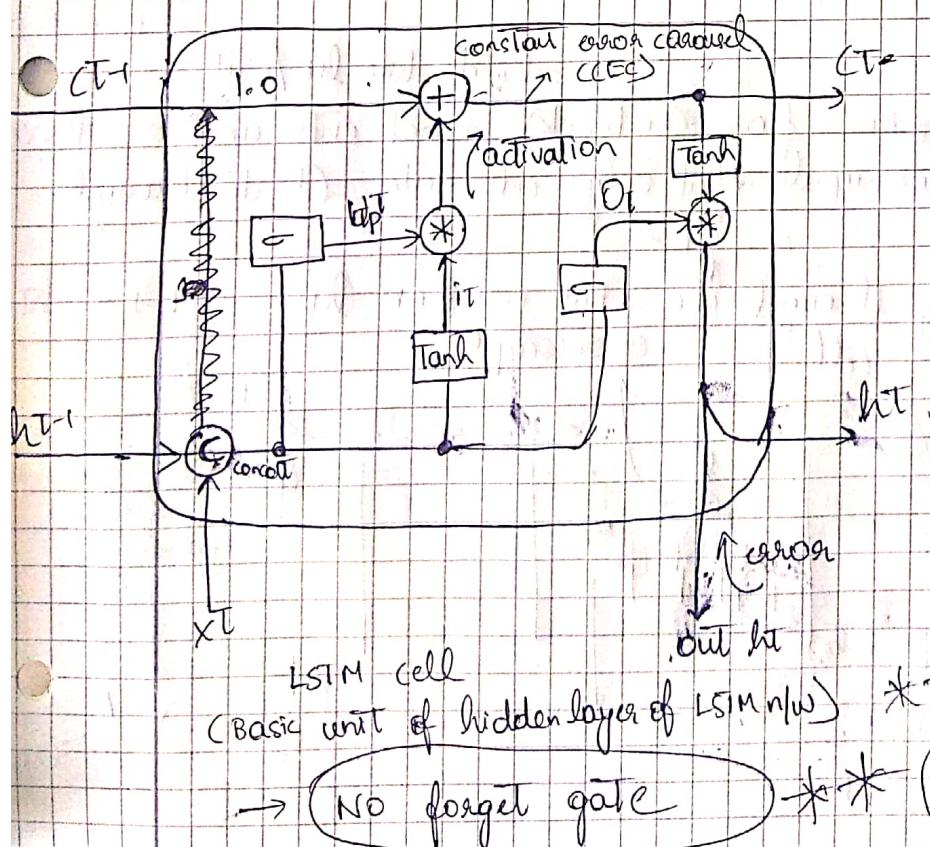
$\rightarrow$  Output after complete sequence  $X$  has been read.

LSTM and O

gated Recurrent Neural Nets

LSTM - Long Short Term Memory

Basic 1997 LSTM cell (w/o forget gate)



$\rightarrow$  New Idea: Instead of utilizing gradient at every step or make update at every step, allow another parallel pathway w/o much obstructions where error can flow back and we can add our new current activation info only if required

$\rightarrow$  Control the info flow through gates

$\rightarrow$  Don't use any heavy technique on error path To make it flow almost constant

How does it solve vanishing gradient problem?

irrelevant ← "The quick brown fox jumps over a lazy dog" irrelevant → relevant  
During forward prop: Let's say we are processing the above string

→ While processing "a" which is an irrelevant token for me, as I am only extracting "animal" class. We make input gate close to zero so that this doesn't muddy the waters and add noise to past cell history.

→ Since past cell history  $C_{t-1}$  is on parallel path even if  $C_t \neq C_{t-1}$  is

$$C_t = C_{t-1} * 1.0 + \begin{cases} \text{input gate} * C_{t-1} & \text{if input gate } C_t \neq 0 \\ \text{no forget gate} & \text{if input gate } C_t = 0 \end{cases}$$

zero  $C_t$  won't become zero. It would still remain  $C_t = C_{t-1}$

→ Input gate protects cell state from being distorted by noise like token "a"  
During backward prop:

→ Errors path through LSTM cell is through output gate and this protects it.

Example:- Propagate error from dog → quick. (Only update near time steps that has (dog, lazy, fox, quick) ignore the rest.)

→ Once the error escape through a valid LSTM output gate, at let's say token "dog" position, Due thinking is, it can ride CEC and reach "lazy", "fox" and "quick"

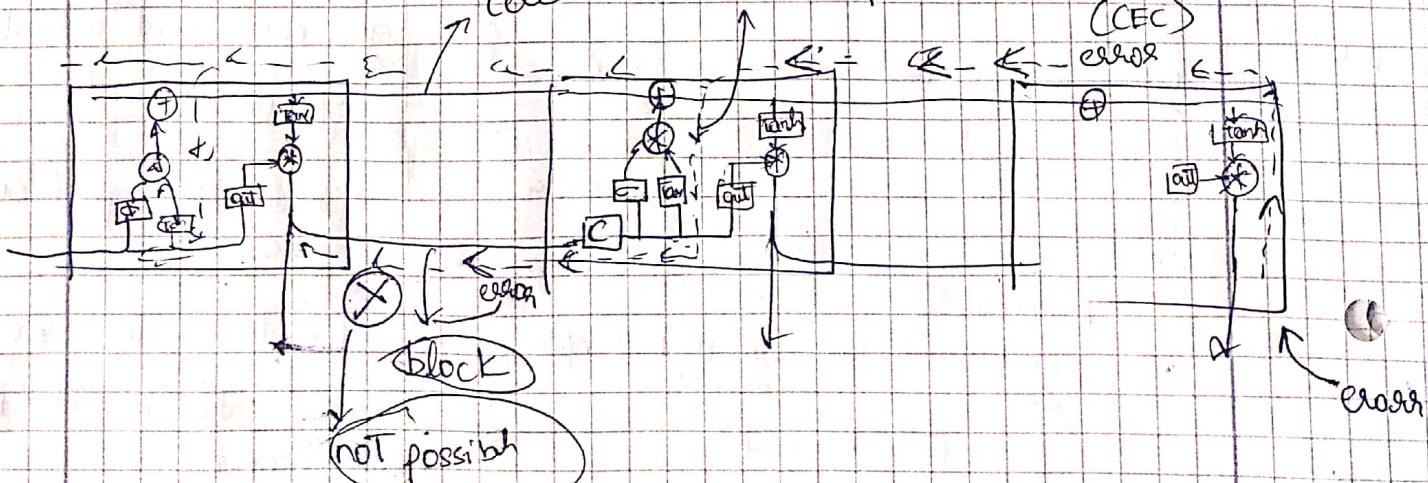
↓  
constant error carried / parallel path

\* \* \*

→ Due to truncated BPTT error from one block cannot enter another. Hence blocks don't exchange error signals but only cell states  $C_t$  do through CEC.

$$(C_t \rightarrow C_{t-1} \rightarrow C_{t-2})$$

Example:- Error passes only through CEC. But not from  $h_t \rightarrow h_{t-1} \rightarrow h_{t-2}$



→ First error signal arrives at memory cell output which gets scaled by output gate  $O_g$  and nonlinearity squasher tank. Then it enters Carousel CFC.

→ It can only escape and update weights of other cell only if their input gates allow so (if they are relevant for us).

→ Let's say they escape through input gate and again scales by Tanh and updates weights and bit (block/output).

→ Now this residual error from  $h_t$  cannot travel to  $h_{t-1}$  due to truncated back prop; but error can still flow through Cell Stack  $C_t$  to  $C_{t-1}$  as there is no obstruction.

\* Problem \*

→ LSTM cannot work where one block bit exclusively serve others like  $h_{t-1}$ ,  $h_{t-2}$  etc as all error are dependent on error propagated from final output to carousel rather than blocks.

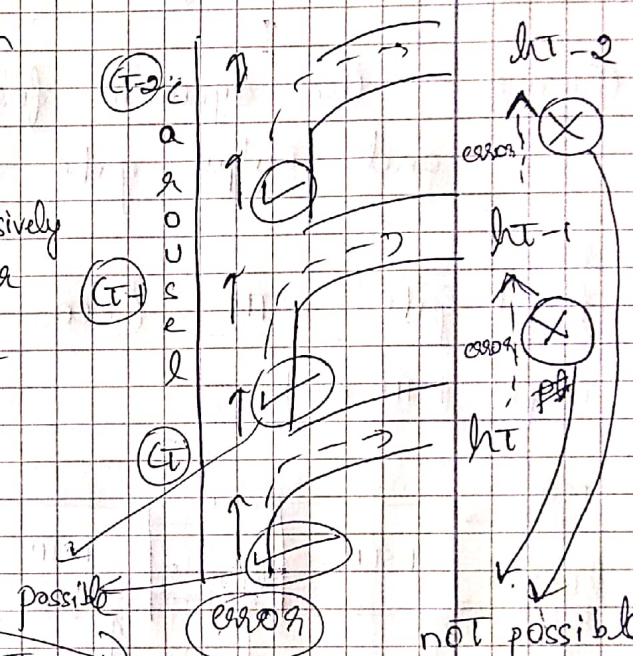
\* This also means that info in cell state gets built up linearly over time

$$C_t = 1.0 * C_{t-1} + I_g * g(\text{NetInput})$$

$$Y_1 = X_1 + mX_2$$

$$Y_2 = 1.0 * Y_1 + mX_3 = X_1 + mX_2 + mX_3$$

$$Y_3 = Y_2 + mX_4 = X_1 + mX_2 + mX_3 + mX_4$$



\* Since we use squashing function Tanh, over time all the inputs output gets dangerously close to +1 or -1 as magnitude of  $t$  grows linearly (if  $t$  becomes 100 or -100 then its Tanh value would be close to  $\pm 1$ .) \*

\* This causes saturation and usually gradients don't backpropagate well enough near these saturation points. Hence we need to (methodically forget some past) stuff so that we move away from

saturation region of Tanh. This can be done if Cell state is not allowed to blowup linearly. How to do this?

Option 1 :- Use <sup>constant</sup> annealing factor to forget stuff of past and make cell not saturable

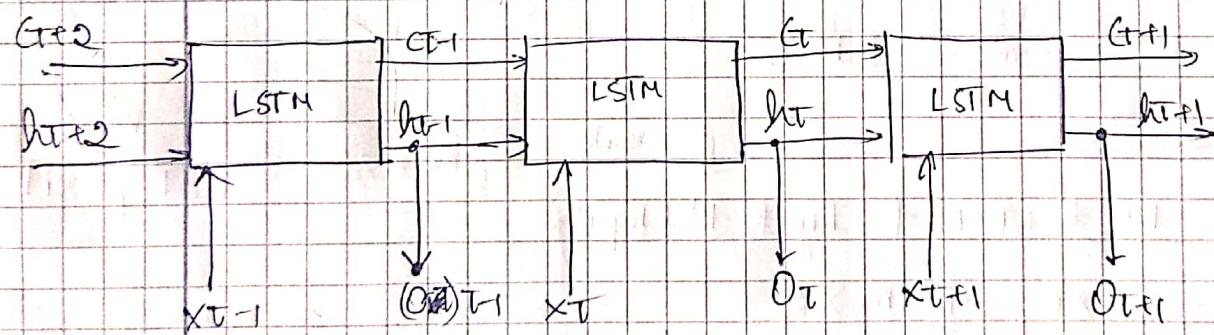
$$G_t = \underbrace{0.9 *}_{\text{annealer (non-learnable / constant)}} (I_{t-1} + I_g * \text{Net}$$

problem :- We should not forget relevant tokens/history/info but only forget irrelevant ones. Having a constant annealer for every LSTM cell makes every cell forget some stuff of history whether it's relevant or not.

Option 2 :- Use dynamic forget gate which is similar to input and output gate & learnable

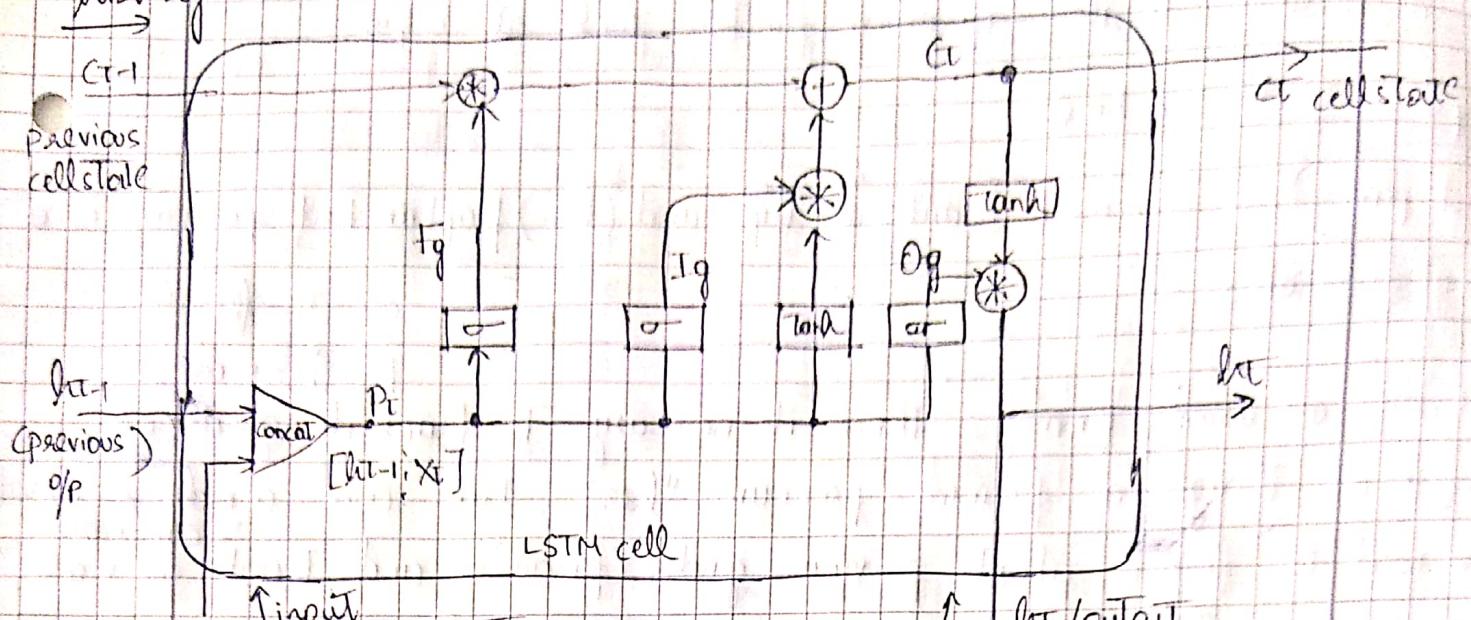
Forget gate LSTM (2003 - Ges Phillips)

Complete LSTM blocks



- Along with hidden state we have a cell state and cell state pathway forms constant cushion of errors to prevent errors from blowing/dying.
- Here, actually cell state carries past summary and history, whilst hidden state acts as output and also links to next state. Let's look inside a cell

history



$$P_t = \text{concat} [h_{t-1}, x_t]$$

$$\text{Forget gate} = \sigma (W_F P_t + b_F) \quad - (1)$$

$$\text{Ingest gate} = \sigma (W_I P_t + b_I) \quad - (2)$$

$$\text{Output gate} = \sigma (W_O P_t + b_O) \quad - (3)$$

$$\tilde{C}_t = \text{tanh} (W_C P_t + b_C) \quad - (4)$$

$$\text{updated cell state} = F_{\text{gate}} * c_{t-1} + I_{\text{gate}} * \tilde{C}_t \quad - (5)$$

$$O_t = \text{tanh} (\tilde{C}_t) \quad - (6)$$

$$h_t = \text{Output gate} * O_t \quad - (7)$$

Explanation

→ We concatenate both  $h_{t-1}$  &  $x_t$  else we may require two weight matrix to train instead of one for every gate.

Forget gate : Learn the gate values, which are from sigmoid distribution to selectively allow only parts of history to move forward. Remove unwanted part of history.

→ This also helps to keep linear sum  $Ct$  within certain limits so that it doesn't saturate squashing function tanh at output.

Example: gender tracking

Black pearl is the most beautiful ship, and her Captain Jack Sparrow is a King of pirates.

 → female

→ WKT ships are usually female. Hence it's necessary to track its gender or confirm it through a feminine pronoun "her". But once a male noun/noun pronoun is introduced we need to forget about previous gender (female) and track for male noun/pronoun like "King". If we don't forget then this may add irrelevant history and noise thereby diluting the prediction process.

Example :-

female

femore

male

famous

Queen Anne's Revenge is a fearsome ship and Captain Teach is her captain, Captain Teach is dad of his first mate Angelica.

[teach] is [dad] of his first mate anglica

female

\* → Without forget gate we have transition from female → Male → Female  
We end up doing wrong prediction & predict may predict all as female

Input gate + Updatable cell state value

I  
I gate

C(CT)

a)  $E_t$ : Updatable cell state values calculates the new values that it should output based on previous output  $H_{t-1}$  & current input  $X_t$

b) Input gate :- Controls the flow of new updates to cell state, w/o which we may add unwanted activations/noise and screw up cell state having past history.

Example! From previous gender tracking example

Black pearl  
relevant

is The most  
↓  
adjective

tracking example  
beautiful ship  
relevant activation

→ Adding activations of "is the most" for my task makes previous cell state "black pearl" noisy. It's like shooting my own foot. Don't update cell state for the activations of those tokens using Input gate.

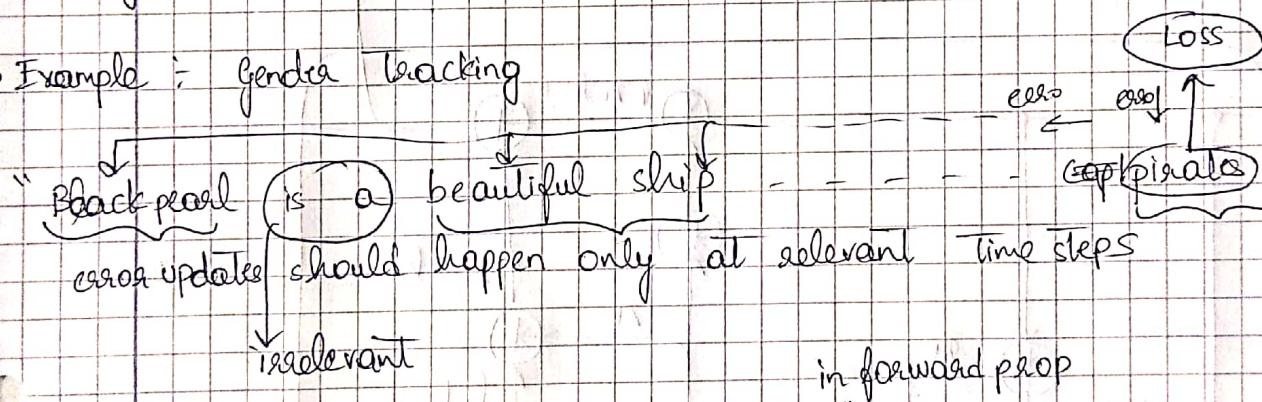
Og Output gate ~~and~~ / New cell state / Output state  $h_t$

a) → Armed with the information of what to forget about past and what relevant info we must add for current step we update existing cell state using this combo

$$C_t = \underbrace{f_{\text{gate}} * C_{t-1}}_{\text{what to forget about past}} + \underbrace{i_{\text{gate}} * \tilde{C}_t}_{\substack{\text{what relevant current} \\ \text{info should be added}}}$$

b) → Output gate :- Output gate helps in shielding against irrelevant errors during back prop. This also helps to output current steps data ( $h_t$ )

→ Example :- Gender tracking



\* → "is a" did not add any gender info, hence useless to be updated during back prop. This also shields against unnecessary updates and errors can travel much longer.

c) Output state :- The output of current step LSTM cell based on current cell state and output gate.

why Tanh? Tanh is centered around zero & squashed b/w [-1, 1] hence its gradient is a bit stronger than sigmoid.

why sigmoid ingates? Gates are similar to boolean logic (0, 1), Sigmoid fits this description.

## Other LSTM variants

Peephole : Include cell state and hidden states from previous step with current i/p  $x_t$  to calculate all gates & update

Gated Recurrent Unit : Simple design, with a single gate for both update and forget, No separate cell state, CTC unchanged

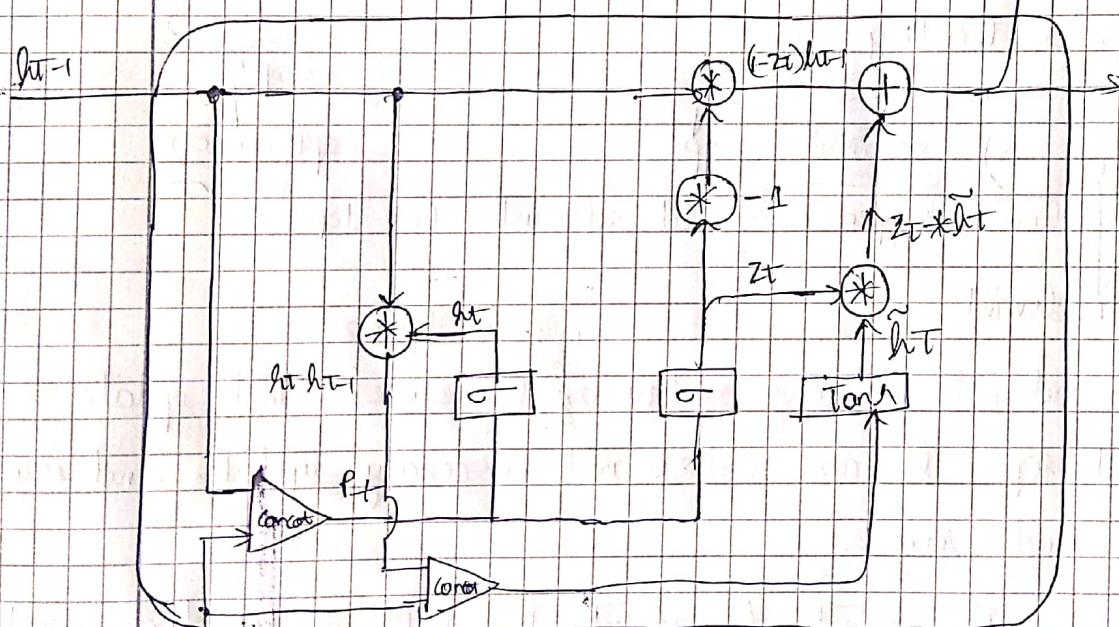
$$P_t = \text{concat} [x_t, h_{t-1}]$$

$$\text{(update)} z_t = \sigma (W_z P_t + b_z)$$

$$\text{(reset)} r_t = \sigma (W_r P_t + b_r)$$

$$\tilde{h}_t = \text{Tanh} (W_h [r_t * \tilde{h}_{t-1}] + b_h)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$



How to initialize gates?

Forget gate : Initialize to near one, so that system can gradually learn what to forget

Input gate : Initialize near zero, so system can increase weights for relevant steps later

Output gate : Initialize near zero, relevant steps later