

ResNets

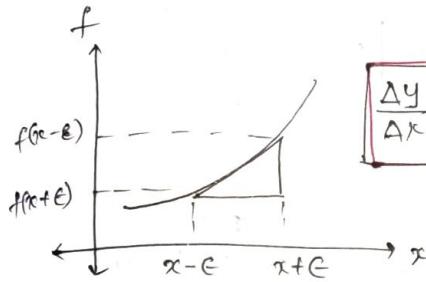
Residual Networks

Gradient Checking

* Numerical approximation = of derivative of f

$$\lim_{\epsilon \rightarrow 0} \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$$

* Actual der \rightarrow calculated by formula or backprop.



$$\frac{\Delta y}{\Delta x} = \frac{f(x+\epsilon) - f(x-\epsilon)}{2\epsilon}$$

because ϵ is small,
difference also is small

$$\text{grad-check} = \frac{\|d\mathbf{w} - d\mathbf{w}_{\text{aprx}}\|_2}{\|d\mathbf{w}\|_2 + \|d\mathbf{w}_{\text{aprx}}\|_2}$$

Euclidean distance

$\|d\mathbf{w}\|_2 < 10^{-7} \rightarrow \text{correct}$

$\|d\mathbf{w}\|_2 > 10^{-3} \rightarrow \text{false/not correct}$

Ex: $\mathbf{w} = w_1, w_2, x = x_1, x_2 \quad f(w_1, w_2, x_1, x_2) = w_1^2 x_1 + w_2^2 x_2$

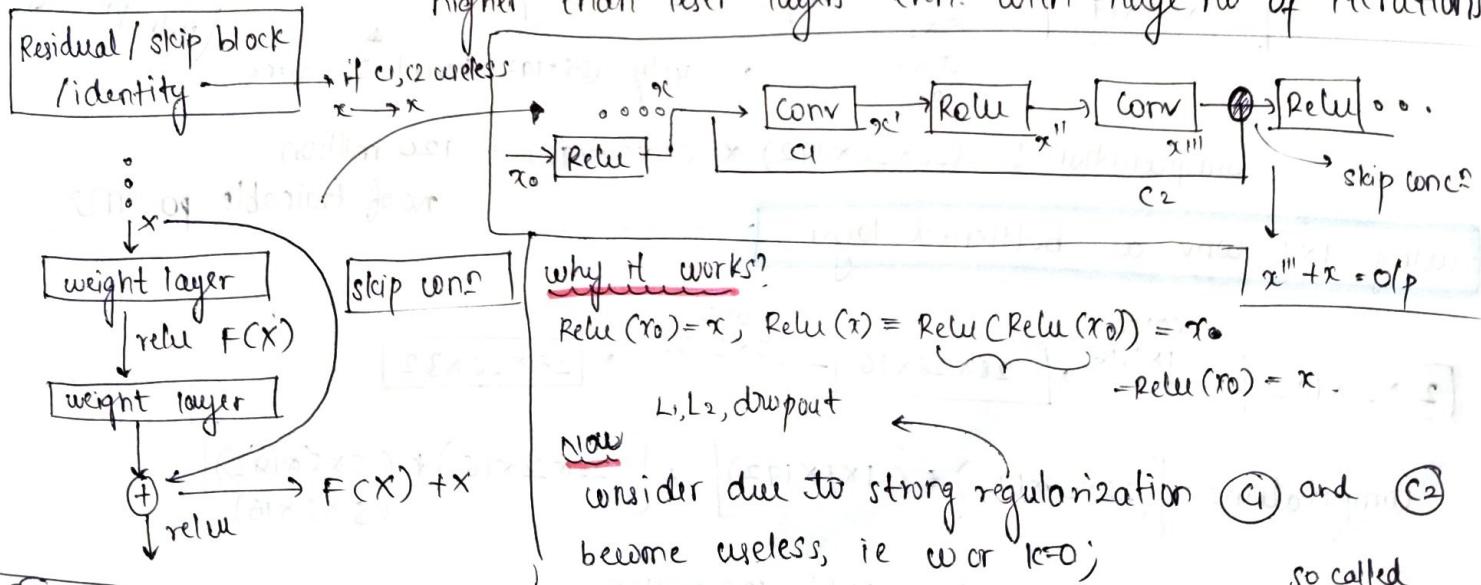
$$\frac{df}{dw_1} = 2 \cdot w_1 x_1 \rightarrow ① \quad \frac{df}{dw_2} = 2 \cdot w_2 x_2 \quad \left(\frac{df}{dw_1} \right)_{\text{aprx}} = \frac{f(w_1 + \epsilon_1, w_2, \dots) - f(w_1 - \epsilon_1, \dots)}{2\epsilon}$$

$$(d\mathbf{w}) = 2w_1x_1 = 6 \rightarrow ① \quad (d\mathbf{w}_1)_{\text{aprx}} = 5.9999999 \rightarrow ②$$

$$\text{grad-check} = \left(\frac{6 - 5.999999}{6 + 5.99999} \right) = 4.251268e^{-13}$$

\therefore as $< e^{-13} \rightarrow$ backprop
is correct

Resnets :- motivation \rightarrow when no of layers were increased, the loss was higher than lesser layers even with huge no of iterations



$$\begin{aligned} \text{ReLU}(\text{ReLU}(x)) &= \text{ReLU}(x) \\ \text{if } +\text{ve}, \text{ReLU}(x) &= \text{ReLU}(x) = x \\ -\text{ve} &= 0 \end{aligned}$$

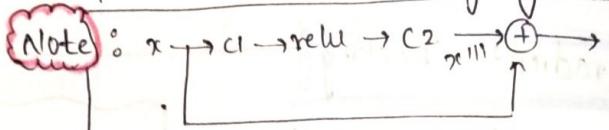
this wont propagate further

$$\text{o/p} = x''' + x = x$$

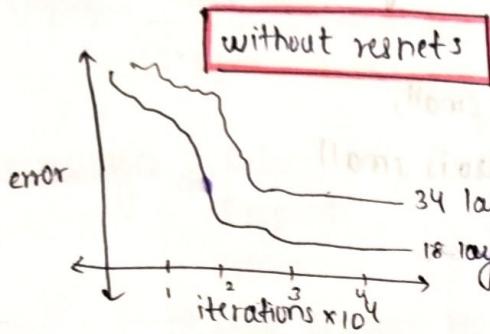
so called skip connection

\therefore C1 and C2 don't matter ie skip

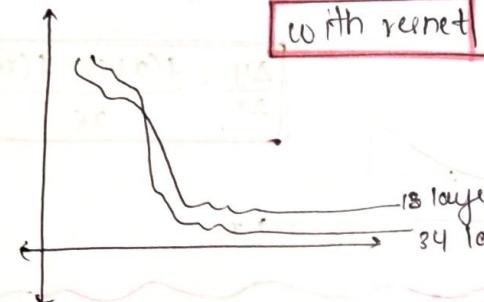
(22) Key takeaway → guaranteed no decrease in performance, small increase
 * adding additional layers would not hurt performance as they get skipped if they are found useless.
 * if new layers are useful reg+weights are non-zero.



make sure $x^{(1)}$ and $x^{(2)}$ have same shape - use padding



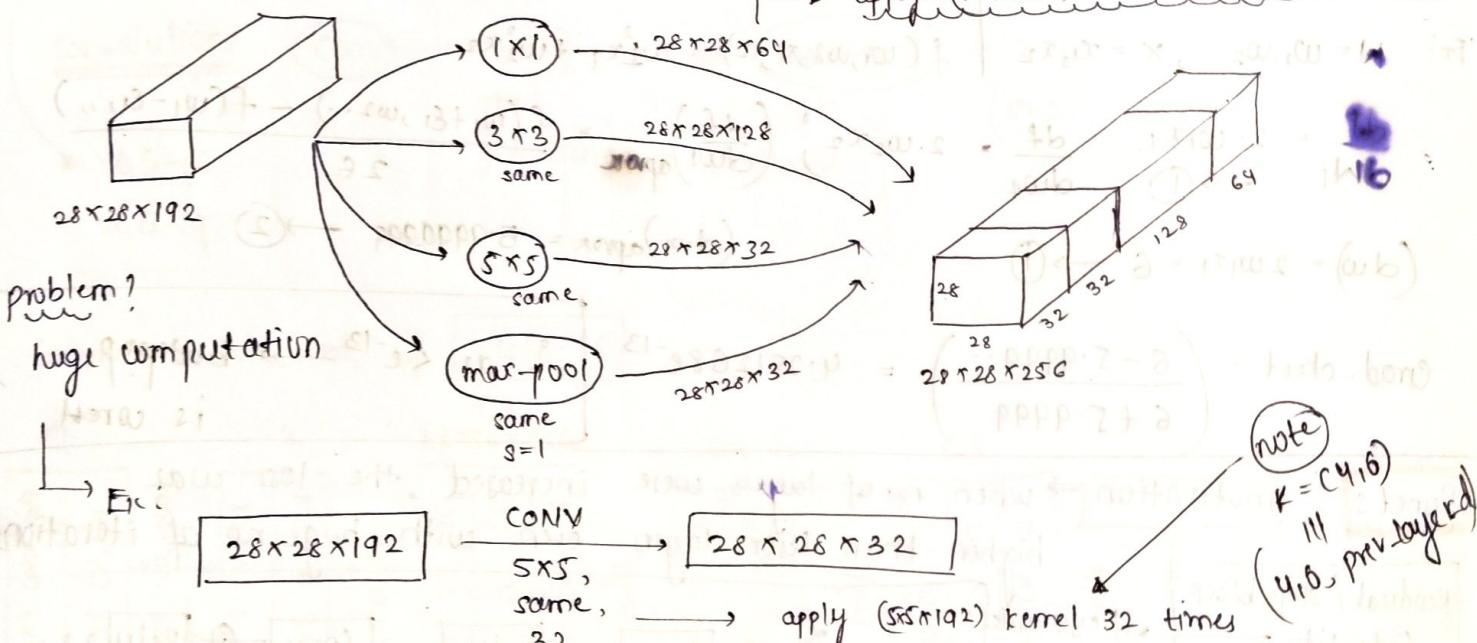
with
res-nets



Inception Network

apply more than 1 type of kernels at once

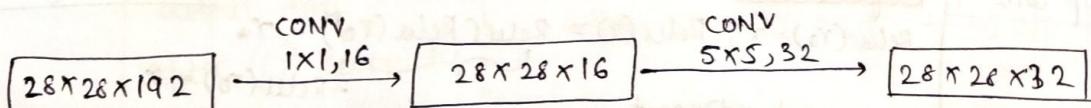
apply $(1 \times 1 \times 192)$ kernels 64 times



$$\text{computation: } (28 \times 28 \times 192) * (5 \times 5 \times 192) = 120 \text{ million.}$$

using 1x1 conv as bottleneck layer

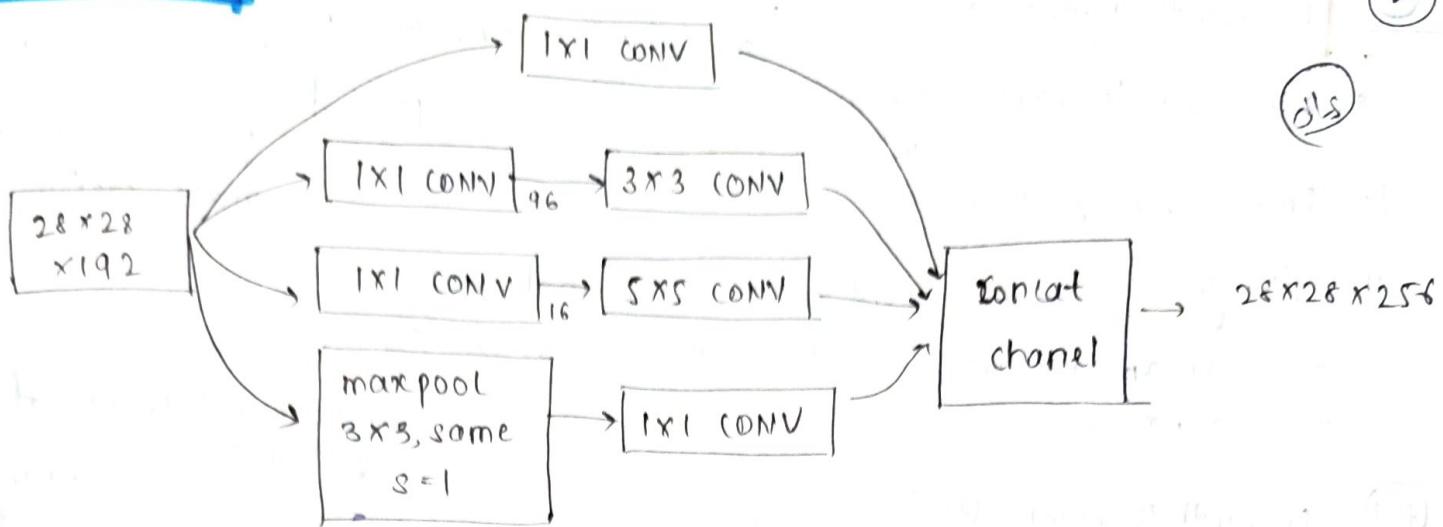
no of trainable params



$$\text{computations: } [(28 \times 28 \times 192) * (1 \times 1 \times 192)] + [(28 \times 28 \times 16) * (5 \times 5 \times 16)] / (5 \times 5 \times 16)$$

$$= 12.4 \text{ million} < 120 \text{ million.}$$

$(\text{inp} \rightarrow \text{conv} \rightarrow \text{conv})$ is better computationally than $(\text{inp} \rightarrow \text{conv})$



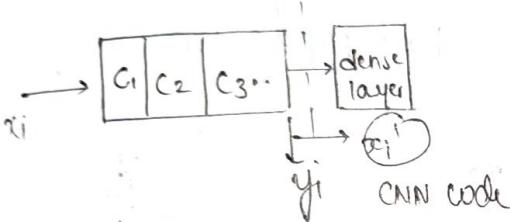
Transfer Learning

→ using pretrained models

Ex: vgg16 trained with imagenet dataset ; New dataset = $D = \{x_i, y_i\}$

they identify basics

case 1 CNN codes



for x_i in D ,

get x_i^1 → vector obtained
from last conv layer

$\therefore D = \{(x_i^1, y_i)\}$

→ train any ml model

Case 2 fine tune last layers

- Freeze most layers → Initial layers are helpful
- fine tune only last few layers → use your dataset to train last layers.

case 3 tune all layers

- use vgg + ImageNet model as a initialization
- train the model using ur dataset

the cases

- where to use?
- factors → similarity with Imagenet and size of dataset

case 1 $D \approx \text{ImageNet}$
IDL small

→ CNN codes

case 3 $D \approx \text{ImgNet}$

IDL medium sized

→ fine tune last few layers

case 5

IDL large
 $D \neq \text{ImgNet}$

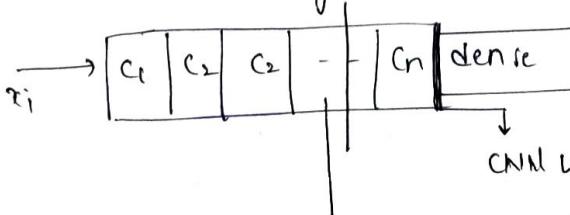
→ use pretrained model as init and train/tune whole mod,

case 2 $D \approx \text{ImgNet}$
IDL large

→ tune full n/w

case 4 IDL is small

$D \neq \text{ImgNet}$



more generalized as

C_n will be useful iff similar

pick from inter same layer but not last

Recurrent Neural Networks

retains and leverages sequence info

Ex: Amazon fine food reviews

D: x_i, y_i

each word is d dimensional vector

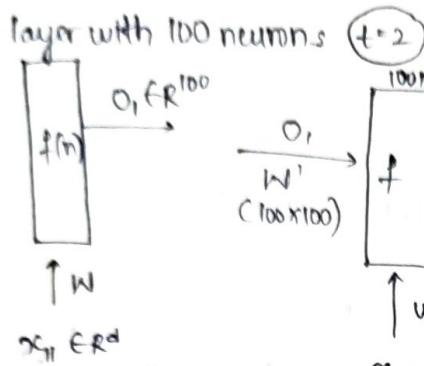
$$x_1 = \langle x_{11}, x_{12}, x_{13}, x_{14} \rangle$$

where $d = \text{size of vocabulary}$

$$x_2 = \langle x_{21}, x_{22}, x_{23} \rangle$$

core idea: sequence of inputs
 machine translation
 time series
 speech recognition
 image captioning

$t=1$

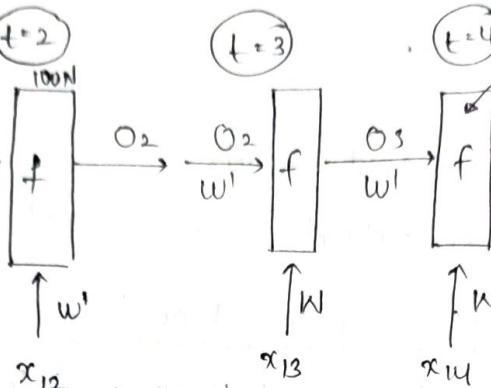


$$N: d \times 100$$

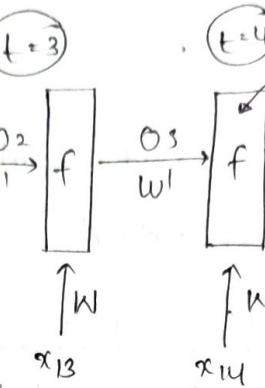
f_m : activation f_m

$$O_2 = f(w' O_1 + w x_{12})$$

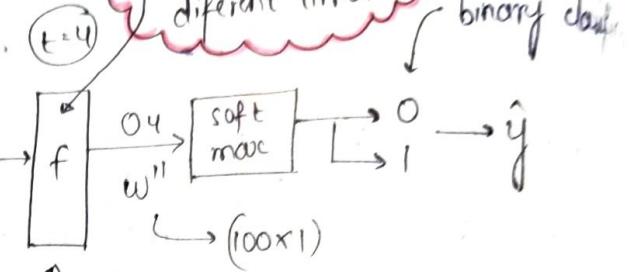
$t=2$



$t=3$



$t=4$



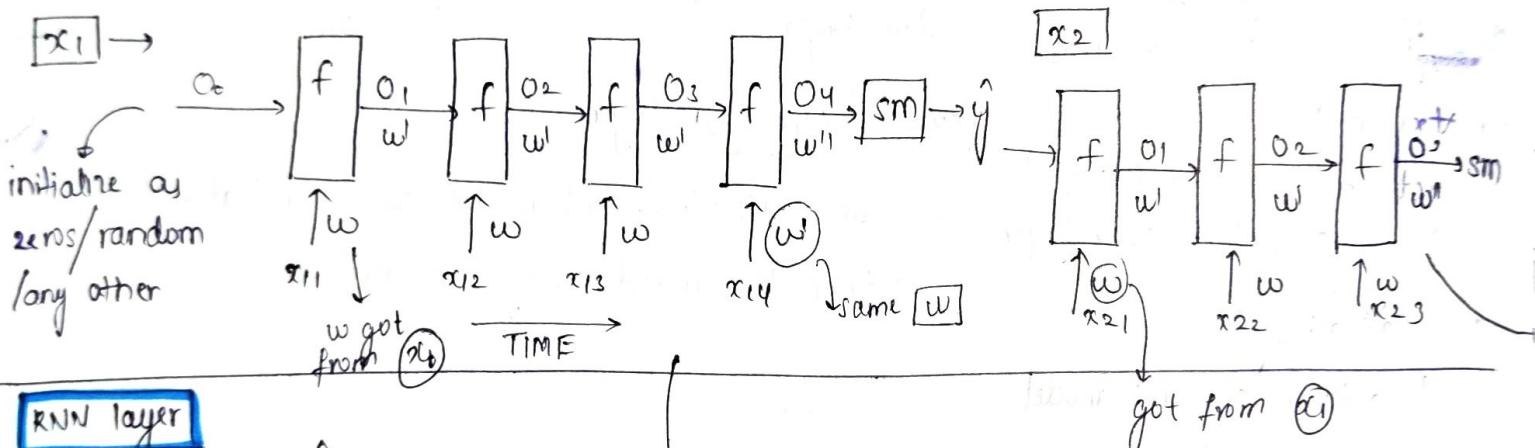
binary class

\hat{y}

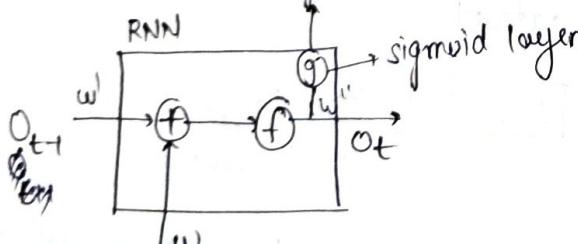
found out by forward and

back propagation

$x_1 \rightarrow$



RNN layer

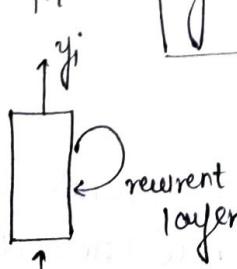


In the whole process we have only 3 matrices, w, w' and w''

→ they are given different inputs according to time t

$$O_t = f(w x_t + O_{t-1} w')$$

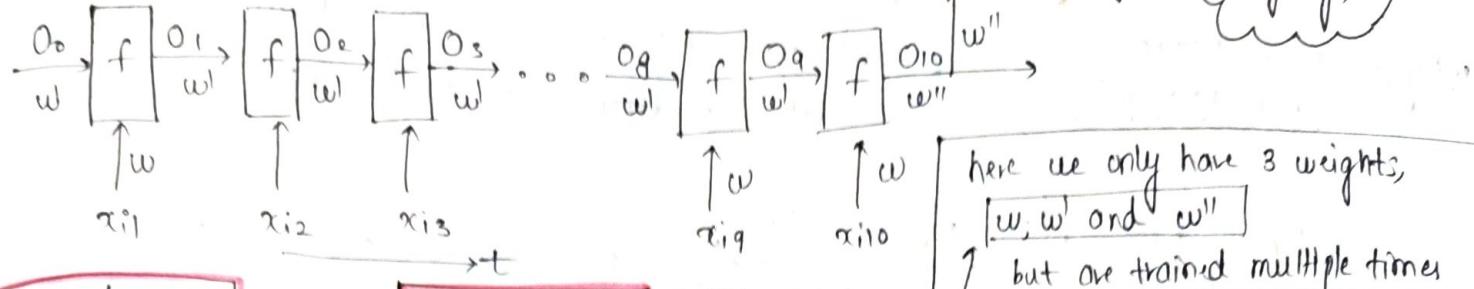
$$y = g(O_t w'')$$



one-to-one



= MLP



Forward prop.

$$o_0 = f(o_0 w' + x_{i1} w)$$

$$o_1 = f(o_0 w' + x_{i2} w)$$

$$\vdots$$

$$o_{10} = f(o_9 w' + x_{i10} w)$$

$$\hat{y}_i = g(w \cdot o_{10})$$

Backward prop.

$$① \frac{\partial L}{\partial \hat{y}_i} \rightarrow ② \frac{\partial L}{\partial o_{10}} = \frac{\partial L}{\partial \hat{y}_i} \cdot \frac{\partial \hat{y}_i}{\partial o_{10}} \rightarrow ③ \frac{\partial L}{\partial w''} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial o_{10}} \frac{\partial o_{10}}{\partial w''}$$

$$w_t = \frac{\partial o_{10}}{\partial w} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial o_{10}} \frac{\partial o_{10}}{\partial w} \rightarrow \left(\frac{\partial L}{\partial w} \right)_{t=1}$$

$$w'_{t=10} \approx \frac{\partial L}{\partial w'} = \frac{\partial L}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial o_{10}} \frac{\partial o_{10}}{\partial w'} \quad \left(\frac{\partial L}{\partial w'} \right)_{t=1}$$

Problems with backprop over time

$\left(\frac{\partial L}{\partial w} \right)_{t=1}$ = multiplication and addition of multiple values → vanishing gradient
Exploding gradient

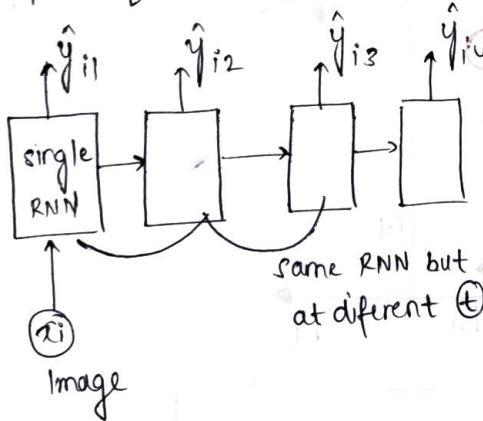
Type of RNN architectures

① many-to-one RNN

usu: sentiment analysis
movie review, rating

② one-to-many RNN

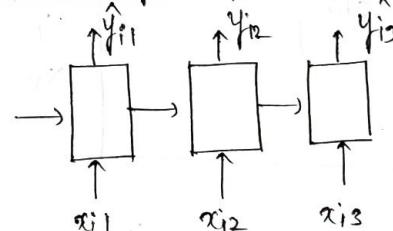
inp: single
out: sequence



Ex: Image captioning

many-to-many [inp] = [out]

① same length : parts of speech



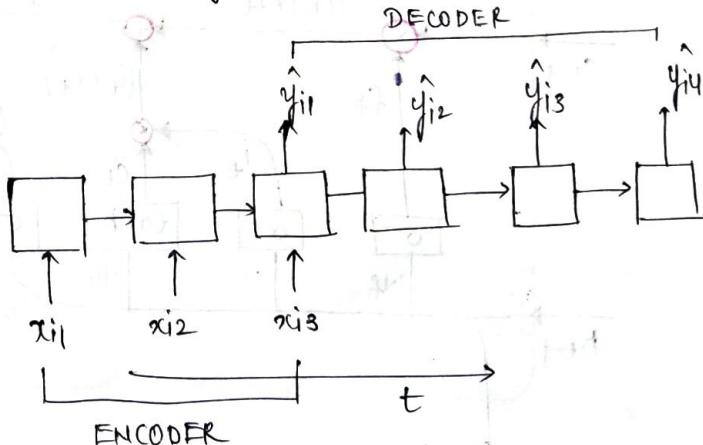
→ word → which part of speech

(seq2seq)

[parallel execution]
can be done

② different length

Encoder - Decoder



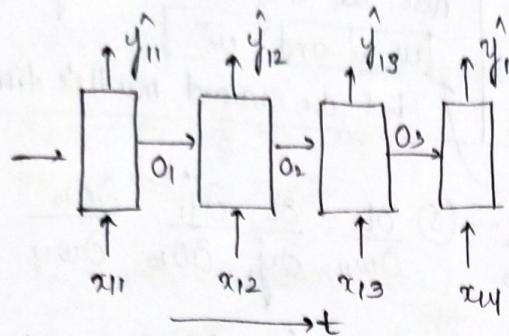
Ex: language translation = Eng → Spanish

a sentence in Eng → may have only 8 words when with 10 words translated into Spanish

LSTM

why we need LSTM when RNN exists?

RNN can't take care of long term dependency
ie later o/p depends on earlier input



why \hat{y}_{14} depends on x_{11} ?

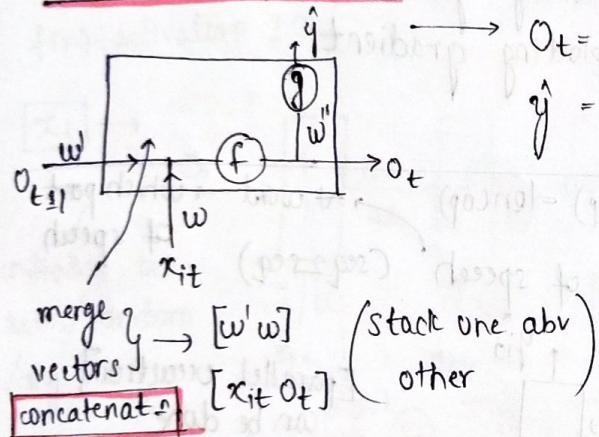
→ there is lot of propagation & time b/w x_{11} and \hat{y}_{14}

→ here,
* \hat{y}_{14} depends more on x_{14} and O_3 and depends very less on O_1 and x_{11}
* let's say we want y_{14} to be more dependent on x_{11} which can't be accomplished

Ex- In some languages starting words might be of more importance when translating

LSTM

Improved simple RNN



$$O_t = f([w, w'] \cdot [x_{it}, O_{t-1}])$$

$$\hat{y} = g(O_t w')$$

similarity

$$C_t = O_t$$

$$X_t = X_{it}$$

$$h_t = \hat{y}_t$$

$$\text{Ex: } [1\ 2\ 3] \times [1\ 1\ 1]$$

Notations

○ → point operation
(apply 2 each)

□ → Neural network layer

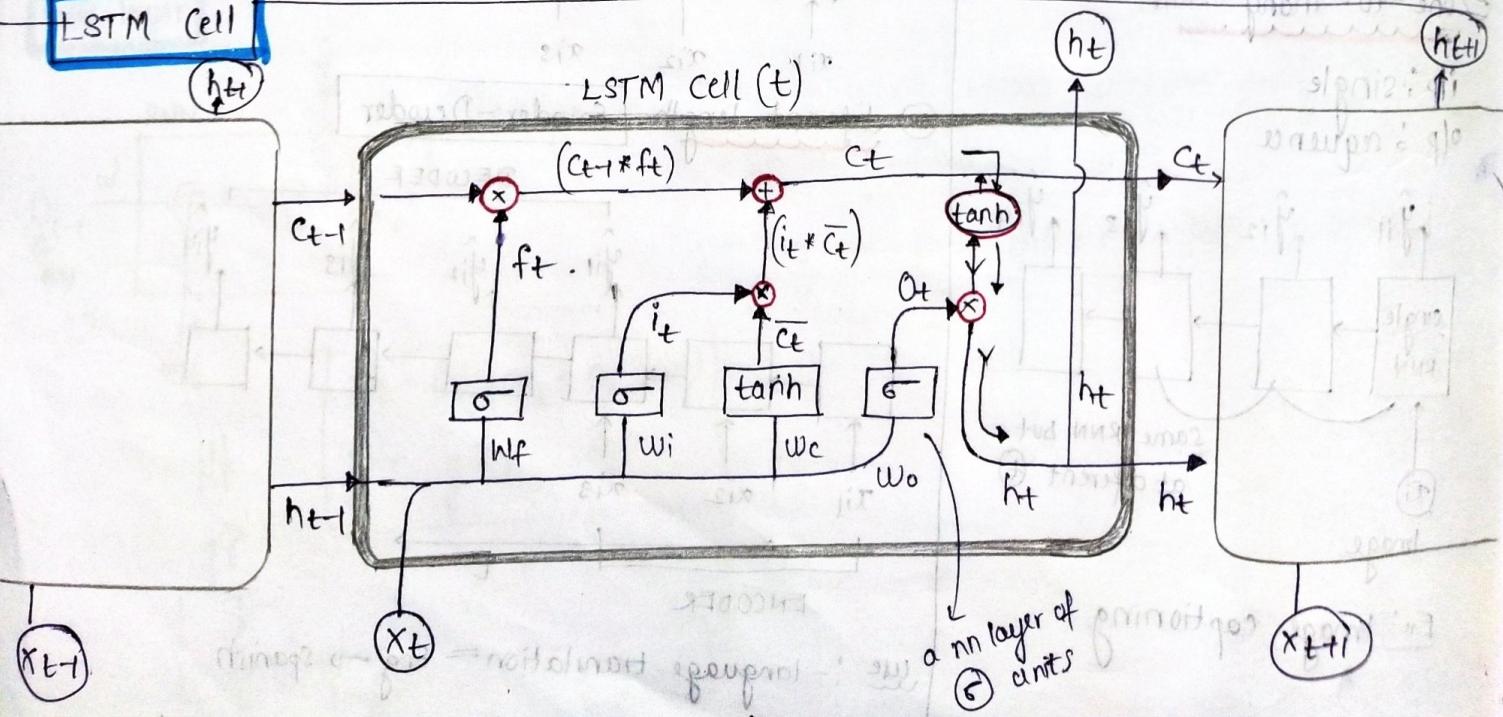
C_t → cell state.

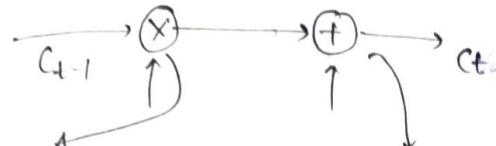
X_t → i/p point

↳ concatenate

↳ copy

LSTM Cell



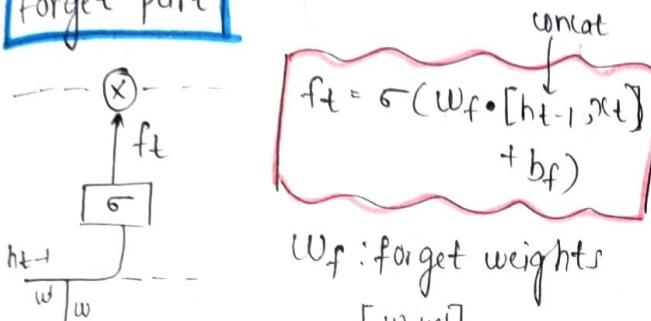


how much to forget? how much more info to be added?

if ip to \otimes = [1, 1, ...] $\rightarrow [c_{t-1} \rightarrow c_t]$
and ip to \oplus = [0, 0, 0, ...]

cell / layer is skipped

Forget part

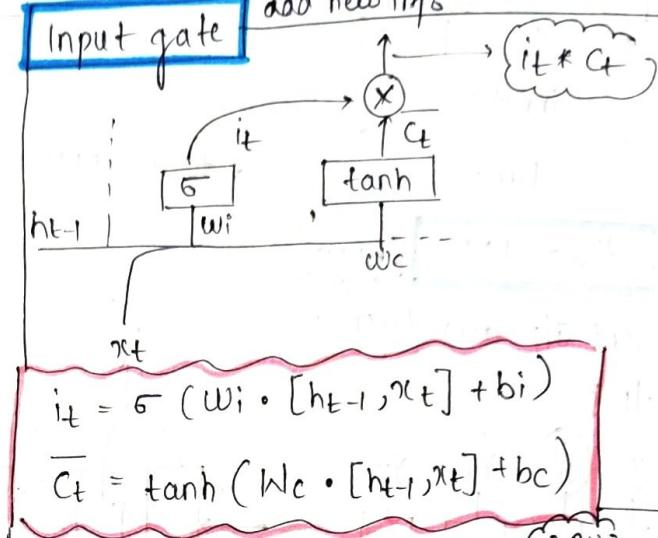


w_f : forget weights
[w, w_i]

b_f : bias factor

$f_t \rightarrow$ how much delta to forget

Input gate



identity + forget + input gate \Rightarrow

$$c_t = (f_t * c_{t-1}) + (i_t * \bar{c}_t)$$

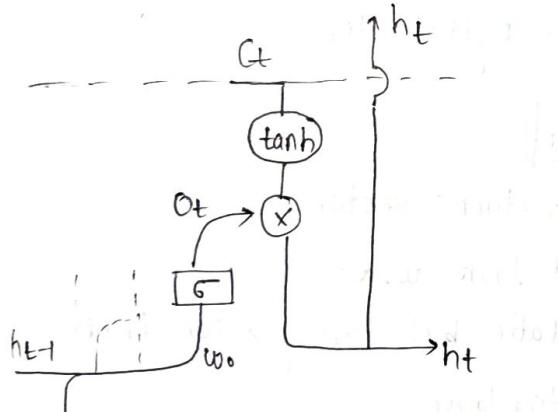
- * few param \rightarrow better with less data
- * not good for long sequence

Output gate

how much op to send?

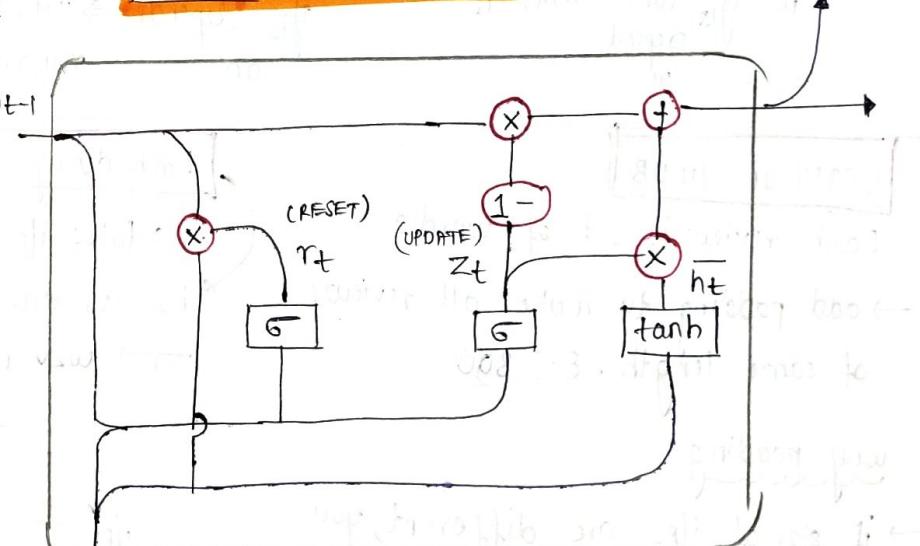
GRU | Gated RNN unit

simplified LSTM



$$\begin{aligned} &\text{GRU vs LSTM} \\ &\text{LSTM have 3 gates} \\ &\text{GRU have 2 gates - reset, update} \end{aligned}$$

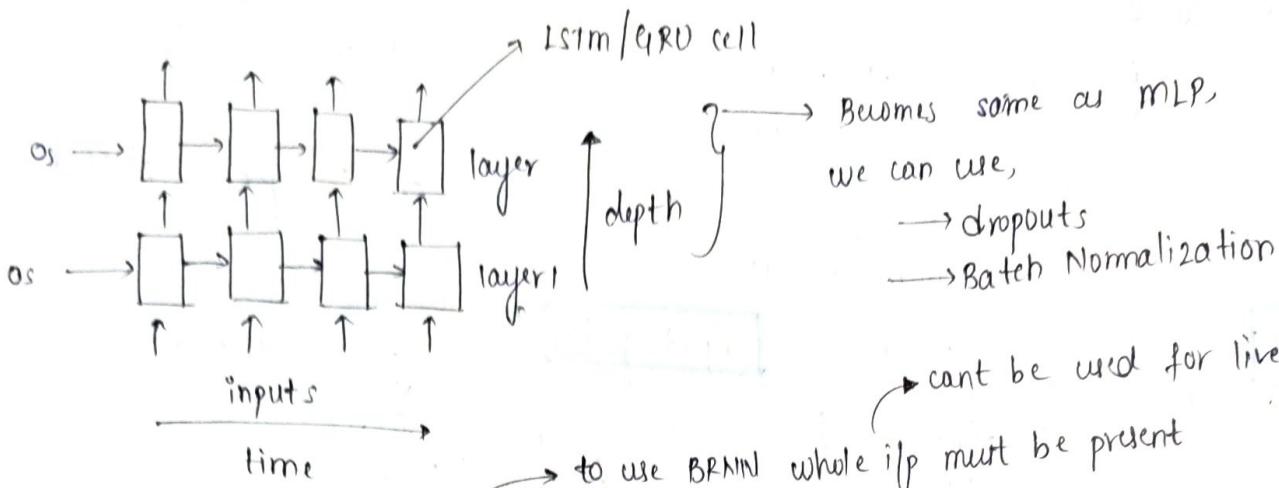
As fewer gates \rightarrow fewer eqn
 \rightarrow faster to train



$$\begin{aligned} &z_t = \sigma(w_z \cdot [h_{t-1}, x_t]) \\ &r_t = \sigma(w_r \cdot [h_{t-1}, x_t]) \\ &h_t = \tanh(w \cdot [r_t * h_{t-1}, x_t]) \\ &h_t = (1 - z_t) * h_{t-1} + (z_t * \bar{h}_t) \end{aligned}$$

$1, 1, \dots$
- [ip vector]

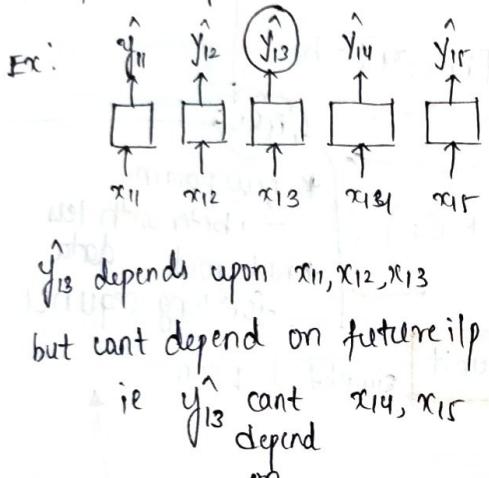
Deep RNN



Bi-Directional RNN

BRNN

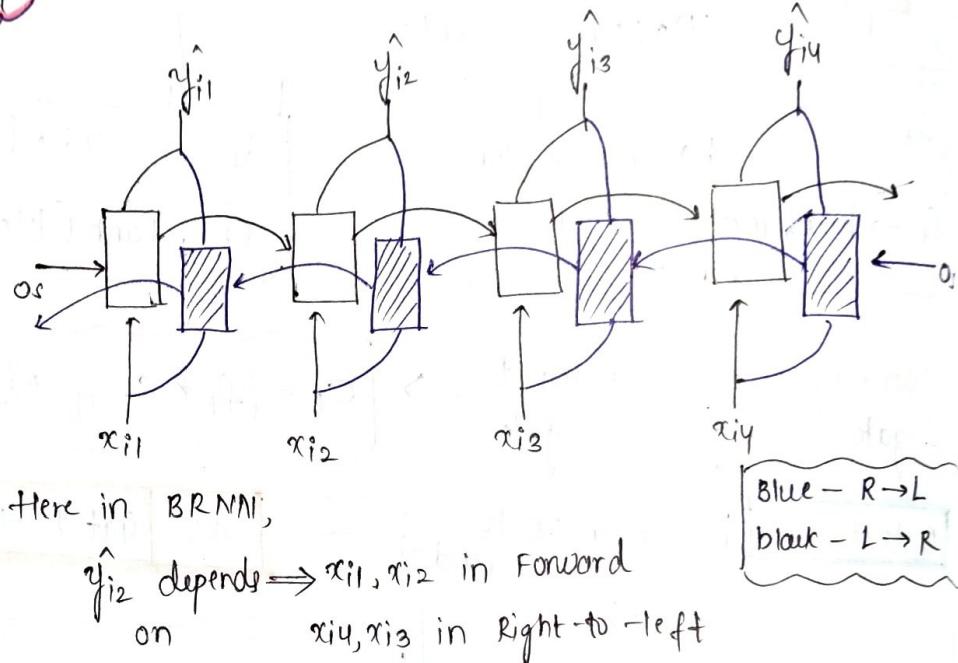
what if current o/p also depends on future ilp?



\hat{y}_{i3} depends upon x_{i1}, x_{i2}, x_{i3}

but can't depend on future ilp

i.e. \hat{y}_{i3} can't depend on x_{i4}, x_{i5}



LSTM on IMDB

Each review is set of words

→ add padding to make all reviews of same length. E.g. 300

why padding?

→ if size of ilps are different, you can only pass one ilp at a time
= stochastic ep → very small

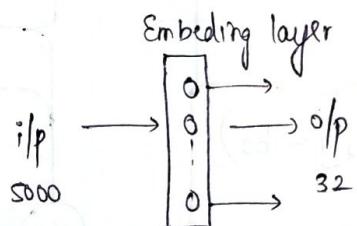
→ if ilp size is same, you can parallelly process inputs like in batch ep → faster

Embedding layer

→ takes ilp → returns vector

→ how is it different from w2v?

→ w2v is static but here you can train



No of trainable params

$$= 5000 \times 32$$

= 160K params

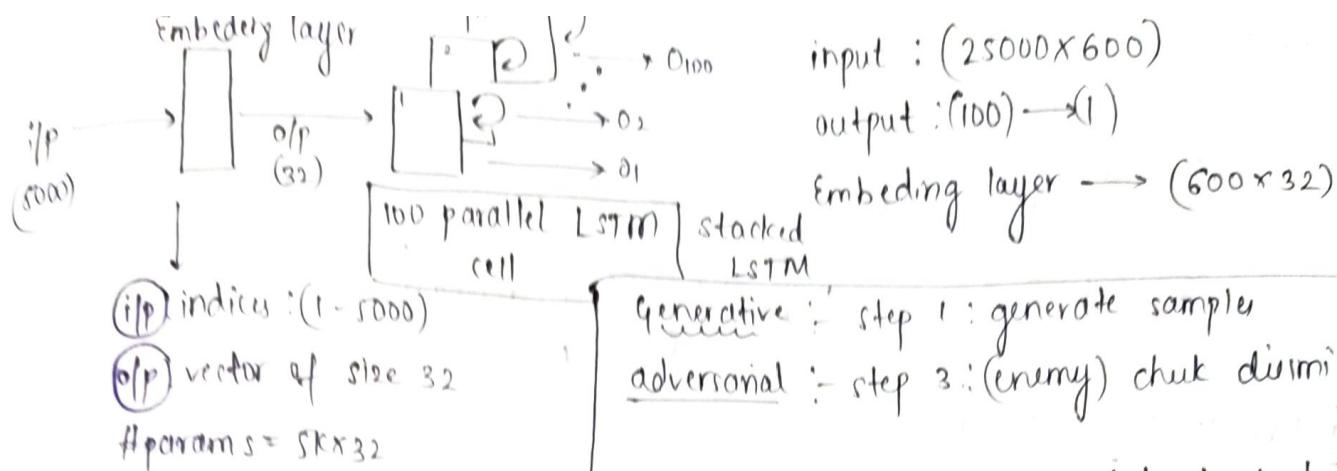
No of train params in LSTM

$$m = \text{no of inputs} = 32$$

$$n = \text{no of o/p} = 108$$

Biases

$$\# \text{params} = 4(nm^2 + n^2 + n)$$



Generative Adversarial Networks

simpler task : 1D scalar

\rightarrow (step 3) check how similar D and D'

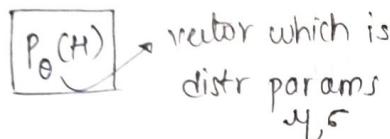
D = set of heights

generate D' similar to D

\rightarrow (s2) generate PDF of D , $P(D)$

* find which distr it is

exit is normal distr



$$P_\theta(H) = N(0 = [y, \sigma])$$

\rightarrow (s2) generate samples randomly

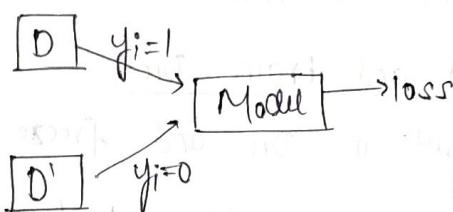
$$\text{from } P_\theta(H) = D'$$

① use statistical tests

(prob model) JS-dist
KL-div

problem: it works only with small dimension

② model approach



\rightarrow if loss is less, then D and D' are separable

$\therefore D$ and D' are not similar

\rightarrow if loss is large ie bad performing model

$\rightarrow D$ and D' are similar

Note : probabilistic models don't work for multivariate analysis. so use ML/DNN

Dataset D'

step ①

generative model (DNN)

②

newly gen dataset (D')

$y_i=0$

Discriminator model

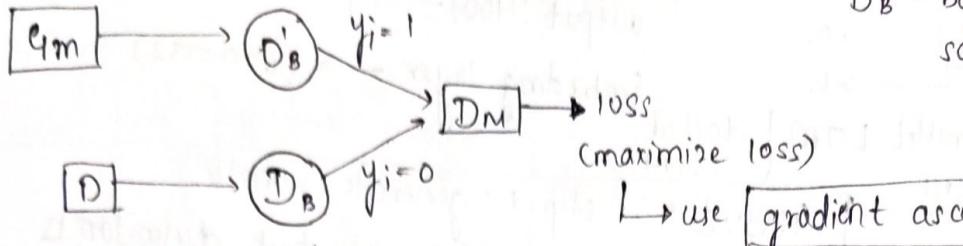
Binary classifier

use for generation of data once architecture is finalized

D

$y_i=1$

(3b) Train GAN



D'_B = batch of dataset
say of 100 datapoints

working : Alternate train optn

step 1 Keep D_m fixed and train G_m

- for D_m initialize we a random sample as train data ex: normal dist? vector
- train the model for it
- * Here G_m does up-sampling; i.e. $(100d \rightarrow 784d)$ MNIST

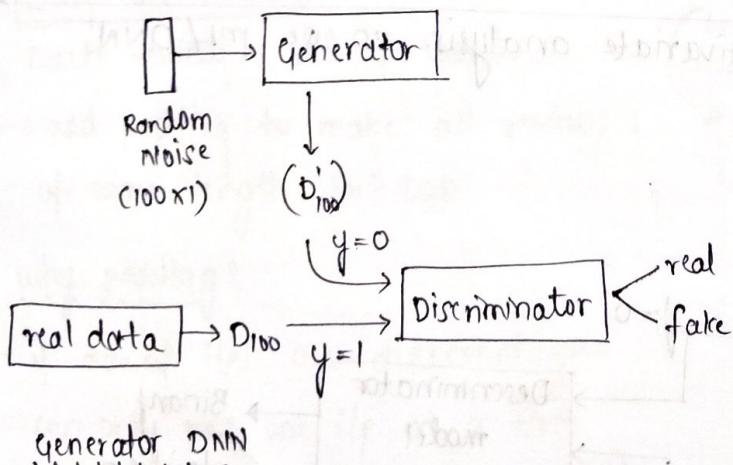
- ① Initially D'_B generated would be very-very random and thus different from D . So loss will be very less which is bad for us
- * we need to find and maximize loss using gradient ascend

step 2 Keep G_m fixed and train D_m

- * only update weights in D_m and freeze weights in G_m

→ Alternatively do ① and ② until satisfactory loss is obtained

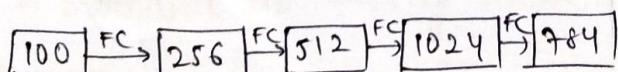
GAN for MNIST



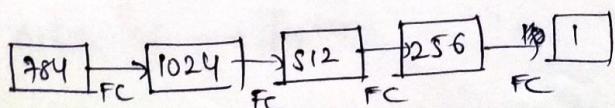
Tips

- ① use normal distribution for sampling
- ② Normalize images b/w 1 and -1
 - use tanh as last layer of generator [$\text{tanh} - (-1, 1)$, ~~sigmoid~~ $- (0, 1)$]
- ③ use batch normalization
 - Each batch must contain only all real or all fake
- ④ avoid ReLU, MaxPool; use Leaky ReLU
 - use adam
- ⑤ use dropouts

generator DNN

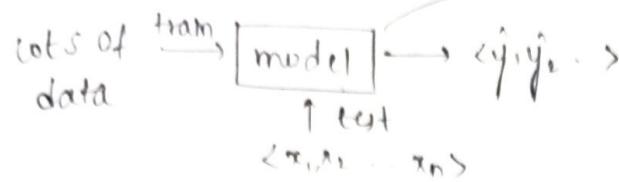


discriminator DNN



Encoder Decoder

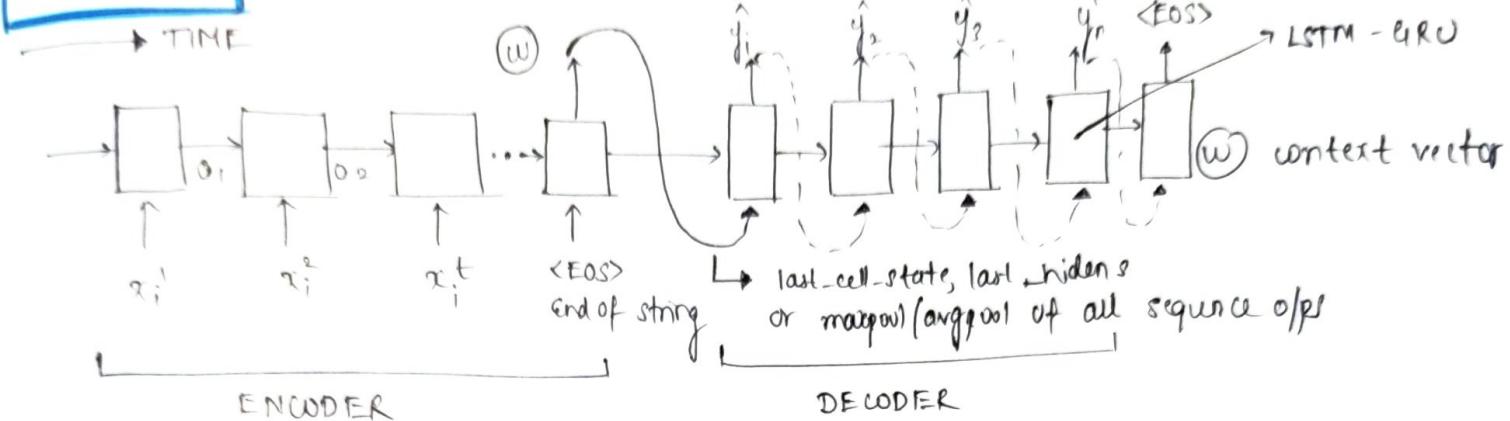
Probabilistic repr of model



$$\text{argmax } p(y_1, y_2, \dots, y_n | z)$$

find y which has max probability wrt data pt z

SEQ-2-SEQ



what is output sentence is long? [cho et al]

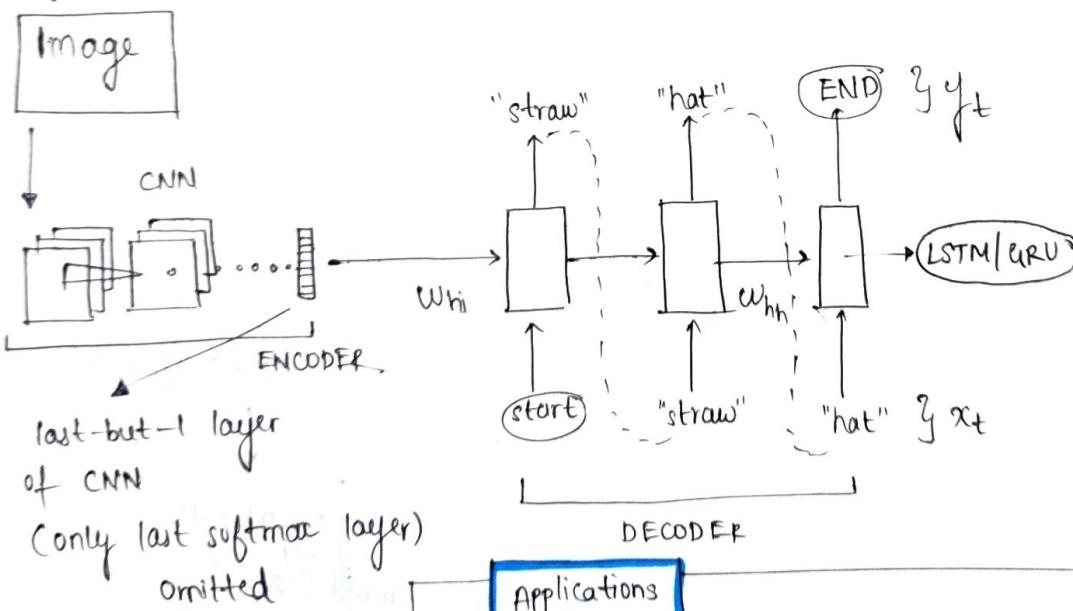
→ context vector is fed to all decoder cells, rather than only first decoder cell

but LSTM can't take 2 if at time stamp

used modified LSTM

Image-2-caption

Karpathy et al



Applications

- ① translation of text
- ② Email auto complete and smart replies
- ③ Image captions and descriptions
- ④ statically finding code errors



Char RNN

- * many-to-many RNN
- * one o/p corresponding to each i/p

(train) let sequence be $c_1, c_2, c_3 \dots c_n$

if i/p is c_2 \rightarrow learn this
o/p should be c_3

given large wrpus of sequence of data
predict which char follows given
char most probably

Word Embedding Layer

(33)

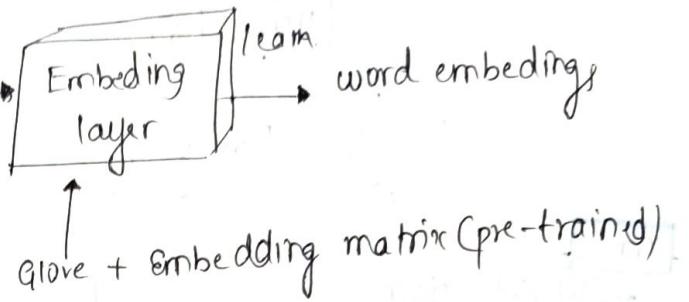
word embedding?

→ represent word with vector → TF-IDF, Bow, Word2Vec etc

word-embedding-layer → let nw learn embeddings by itself

Process

sentence (seq of words) → Integer representation



Ex:

docs = ["sentence is one ... ",
"sentence - 2 ... ".]

① Tokenizer (word → int)

{ "word": index }

use glove-vector / pretrained model

→ t = Tokenizer().fit_on_texts(docs)

x_tr = t.fit

→ x_tr = t.texts_to_sequences(docs)

assume ur glove vector has 200 dim → vector for each weight

create Embedding matrix

	0	1	2	...	200
0					
1					
2					
...					
300					

vocab size
↑ each word d
↑
(300 × 200)

let len(t.word_index) = 300 (300 words)

padding make all sentences have same len

→ maxLen = for \vec{v} in x_tr find largest \vec{v}

→ $x_tr = \text{pad_sequence}(x_tr, padding='post',$
 $\text{dtype}=\text{int}$, $\text{maxlen}=\text{maxLen})$

↑ each word in vocab,
have a v-dim vector

for word, idx in t.word_index:
embed[idx] = glove.get(word)

∴ we have each sentence/data point with len = maxlen

Ex → [1. 0 0 0 0 0 ...]

Embedding Layer

total no of words in vocab/len(em)

e = Embedding (input_dim = len(vocab),

output_dim = 200,

input_length = maxlen,

trainable = false,

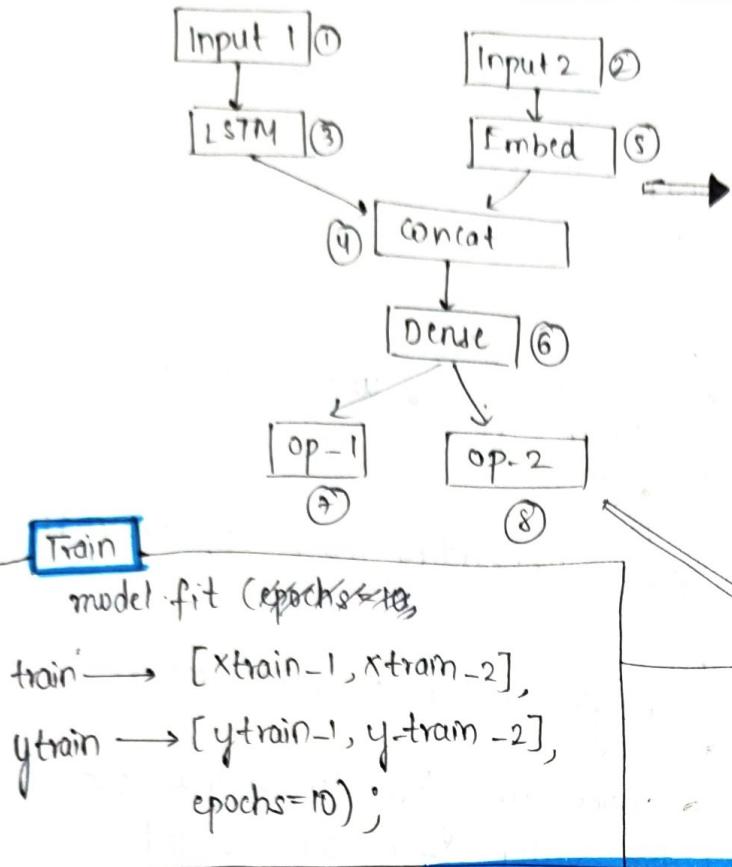
weights = [Embed_matrix])

→ output dimension of each word.
(200 in our case)

length of each training data.

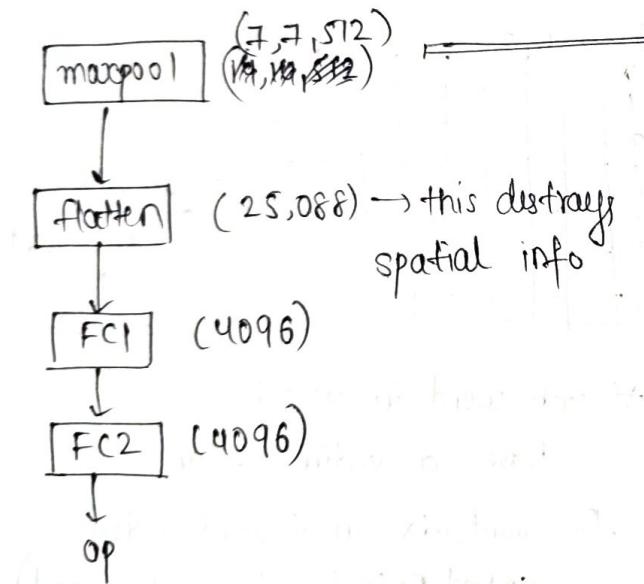
} for pretrained glove vector

Functional API Keras



$\textcircled{1} \text{ Inp-1} = \text{Input}(128, 7)$
 $\textcircled{3} \text{ LSTM} = \text{LSTM}()$ (Inp-1)
 $\textcircled{2} \text{ Inp2} = \text{Input}(10, 7)$
 $\textcircled{5} \text{ Embed-L} = \text{Embedding}()$ (Inp2)
 $\textcircled{4} \text{ concat} = \text{concatenate}()([\text{LSTM}, \text{Embed-L}])$
 $\textcircled{6} \text{ Dense} = \text{Dense}()$ (concat)
 $\textcircled{7} \text{ op-1} = \text{Dense}(10, \text{softmax})$ (Dense)
 $\textcircled{8} \text{ op-2} = \text{Dense}(2, \text{softmax})$ (Dense)

Replace FC layer by CONV 2D

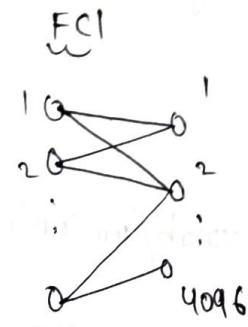
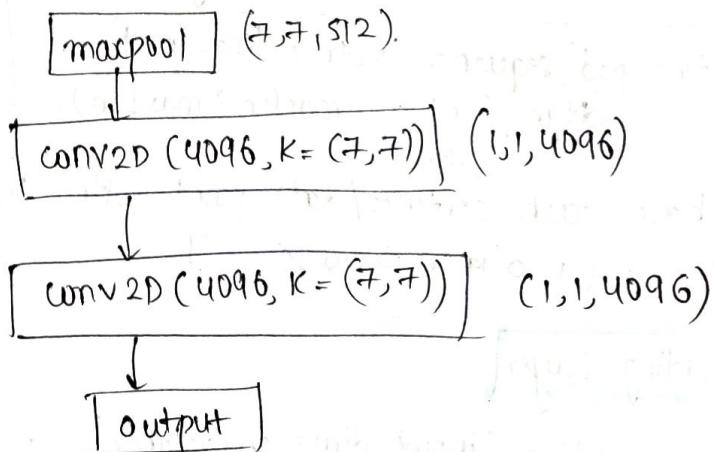


replace FC by CONV 2D

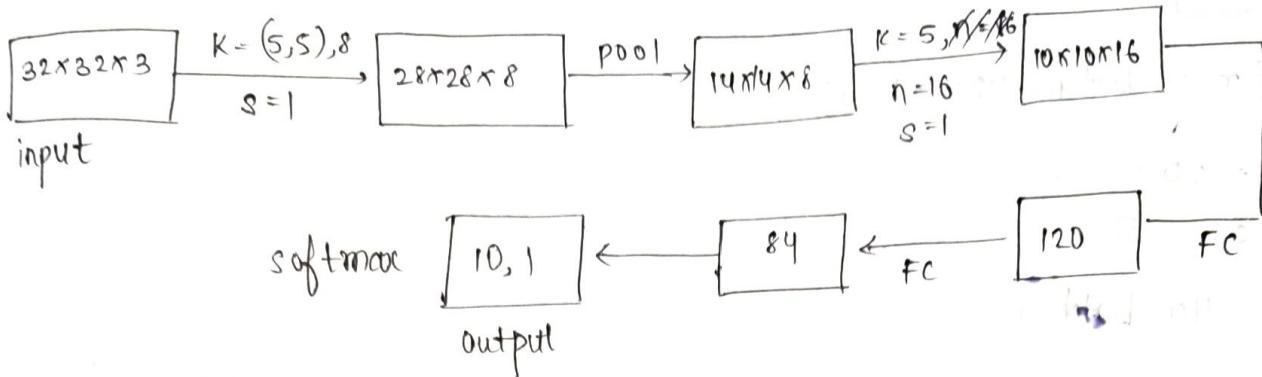
for FC1 we need $4096 = (1, 1, 4096)$

$* (7, 7, 512) \xrightarrow[s=1, fn=4096]{K=(7 \times 7)} (1, 1, 4096)$

$\left(\frac{(7-7+1)}{1}, \frac{(7-7+1)}{1}, 4096 \right)$
 $= 1, 1, 4096$



Calculate CNN params



* conv layers

$$(n, n) \rightarrow \left(\frac{n-k+2p+1}{s}, \frac{n-k+2p+1}{s}, \text{nof filters} \right) \rightarrow \text{Bias term for each filter}$$

$$\# \text{params in conv layer} = \left[\left(K \times K \times \text{nof filter in previous layer} \right) + 1 \right] * \text{nof filters in current layer}$$

$$\# \text{params in FC} = \text{FC(Fcn)} = \# \text{neurons in last layer}$$

* $\# \text{neurons in current FC} + \text{FCn}$

Bias term

Count

layer	shape	# params	activation size	prev. kernel ↑ $[(5 \times 5 \times 3) + 1] * 8$ ↓ bias cur kernel
Input	$(32 \times 32 \times 3)$	0	$32 \times 32 \times 3 = 3072$	
CONV1 $K=5, n=8$	$(28 \times 28 \times 8)$	$\left(\frac{(28 \times 28 \times 3) + 1}{5} \right) 8$ = 608	608×6272	
POOL1	$(14 \times 14 \times 8)$	0	1568	
CONV2 $K=5, n=16$	$(10 \times 10 \times 16)$	$\left(\frac{(14 \times 14 \times 8) + 1}{5} \right) 16$ $\left[(5 \times 5 \times 8) + 1 \right] * 16$ = 3216	1600	
FC1 (120)	$(120, 1)$	$(10 \times 10 \times 16 \times 120)$ + 120 = 192120	120	

* Embedding layer (vocab size, output_dim)

o/p $\rightarrow (\text{None}, \text{ip-size}, \text{output})$

if $\rightarrow \text{vocab-size} * \text{output-dim}$

36

Batch size

```
→ model.fit(x, y, batch_size = b);
```

$$\text{Let } \tan(\tau) = 200.$$

batch-size = 5

\therefore we get 40 batches

training → Epoch 1 → divide data into 40 batches → pass batch 1 & update weight
→ pass batch 2 & update weight.

\therefore In one epoch, weights are

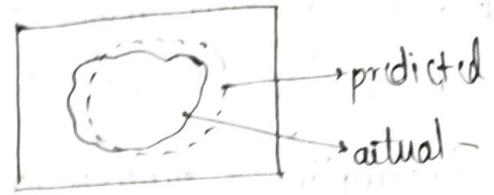
update $\left(\frac{\text{size of dataset}}{\text{batch-size}}\right)$ no of time ie 40 in our case

Image Segmentation

(37)

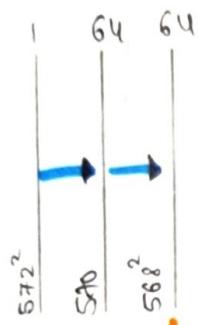
→ divide into areas / segments

measures :- ① jaccard similarity = $\frac{R \cap R'}{R \cup R'}$
 (100)

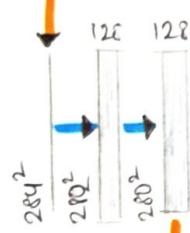


② cross entropy = per pixel comparison

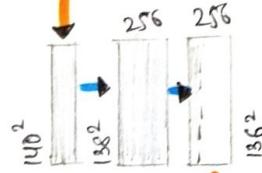
U-net



copy and crop



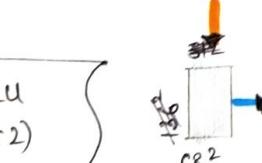
crop (20x20)



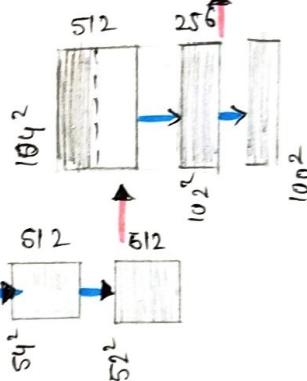
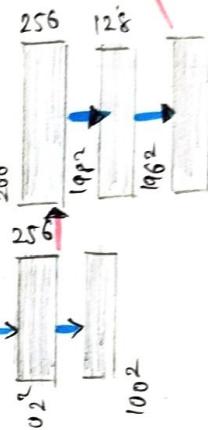
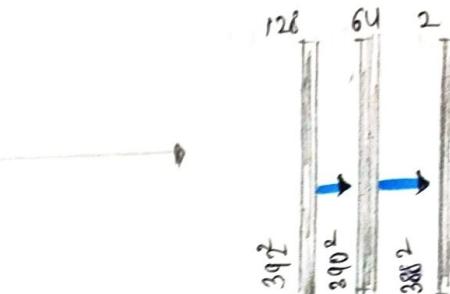
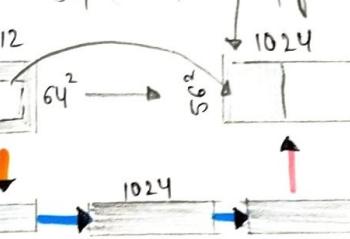
copy and crop

copy (64x1024)

COPY AND CROP



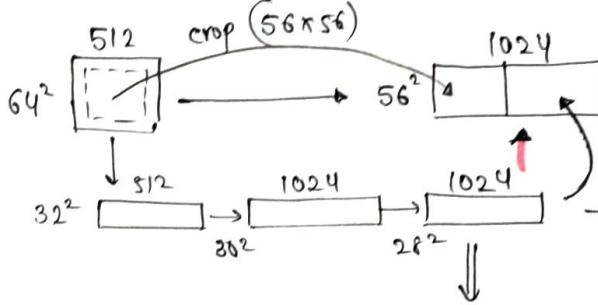
crop (36x36)



concat (upsample, copy-crop)

→ conv 3x3, ReLU
 → max pool (2x2)
 stride = 2
 → upconv (2x2)
 upsample + conv (2x2)

upsample base layer



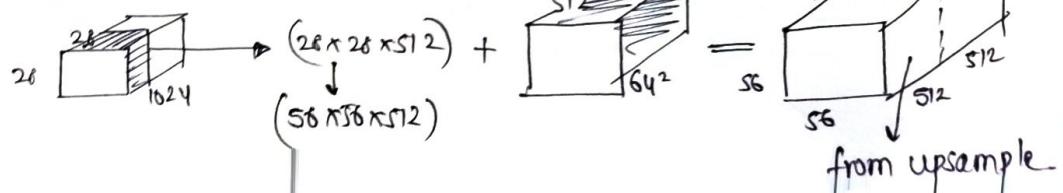
upsample

① upsample : duplicate rows and columns

$$28 \times 28 \times 1024 \rightarrow 56 \times 56 \times 1024$$

② conv2D ($f = 512$, (2x2))
 $(56 \times 56 \times 512)$

crop from copy



10

③ General idea

- * as we go down, image size halves and no. of filters doubles
ie we go into pixel level as we go down
- * (upsampling + copy-n-crop) acts like a residual unit ie get information from lower level as well as from copy-crop unit
- * Here o/p = $(388 \times 388 \times 2)$
i/p - grayscale \rightarrow 2 masks are predicted
ie binary segmentation to identify 2 regions

LSTM + Time Series

moving window average

consider you have data for t days

future data ($t+k$) depends on average of previous K days

other ways

mean avg

median

exponentially weighted

dataset : weather temp each 10 minutes given

Univariate prediction

assume $K=20$ we take window size = 20

* Normalize all

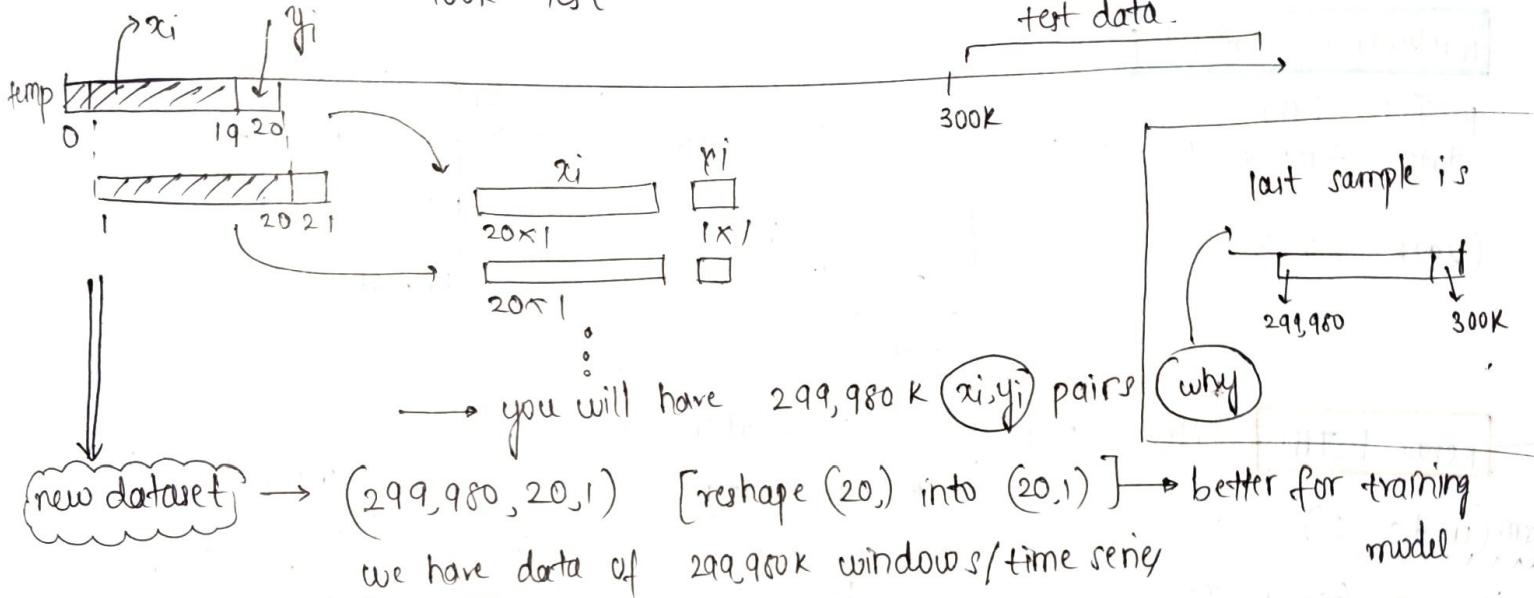
data

only train data

dataset creation

300K - train

100K - test



LSTM training diagram

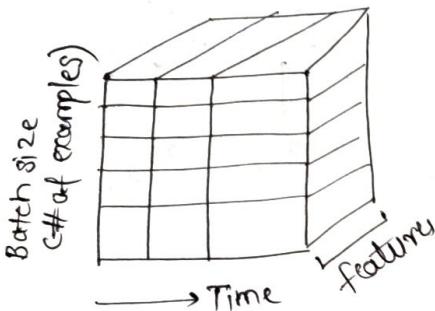
* LSTM expects data in this format \rightarrow [batch, time, features]

height (y) \rightarrow samples

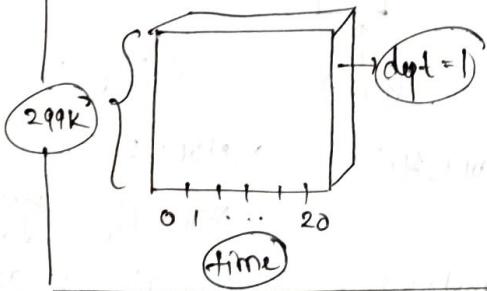
$x \rightarrow$ time series

depth \rightarrow features
(univariate 1)

loss = mae
mean absolute error



our data

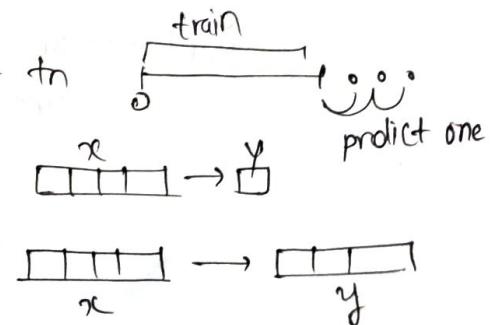


multivariate - single step

→ given a window's predict only one data at t_1, \dots, t_n

* single step → predict only one pt at future

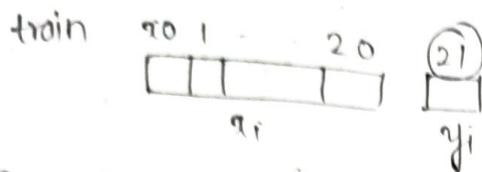
* multi step → predict over multiple points



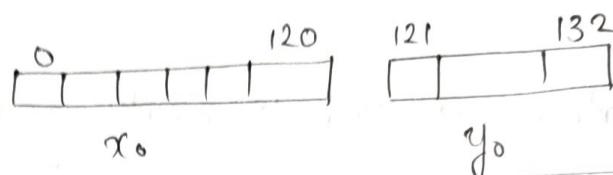
40

Dataset generation hawks

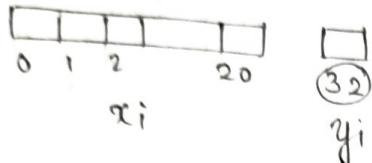
① predict next hour itself



② predict next 12 hours data using past 120 hours → multi step



② predict next 12 hrs



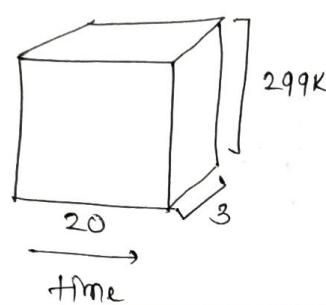
$$\therefore \begin{cases} x_{\text{train}} = (299K, 20, 1) \\ y_{\text{test}} = (299K, 12, 1) \text{ or } (299K, 12) \end{cases}$$

Multivariate dataset

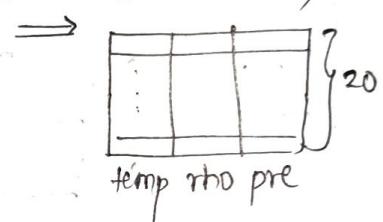
predict - temp

data - temp, rh, pre \Rightarrow

(299K, 20, 3)



one xtrain (20x3)



Keras LSTM

ip = (None, 3, 2) \rightarrow 2 features
 \downarrow
 4 time-steps

dim (units = 2)

ip \rightarrow (x₀, x₁, x₂)

② \rightarrow dim = no of units



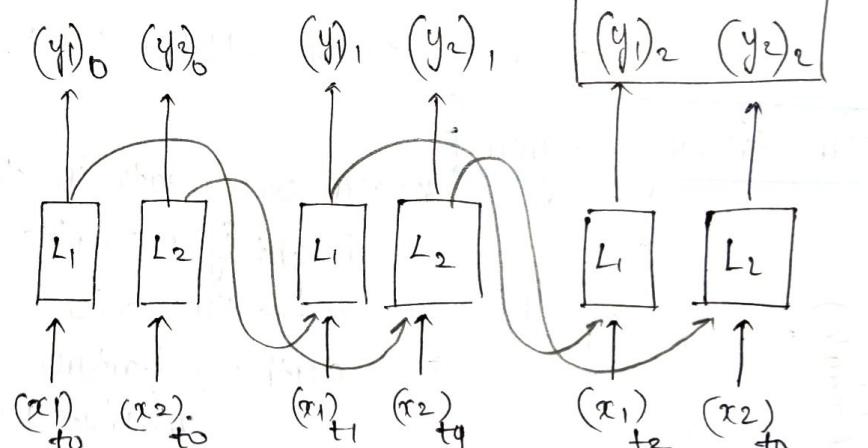
(H_{1,2})

(y₁)₂, (y₂)₂



(None, 3, 2) \rightarrow (None, 2)

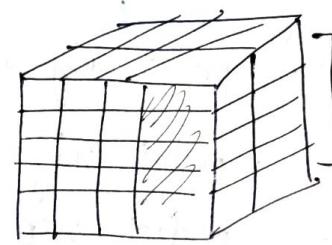
None, 3, 2



units = 2 \rightarrow generate 2 data for each point \rightarrow one o/p from last unit

None, 3, 2

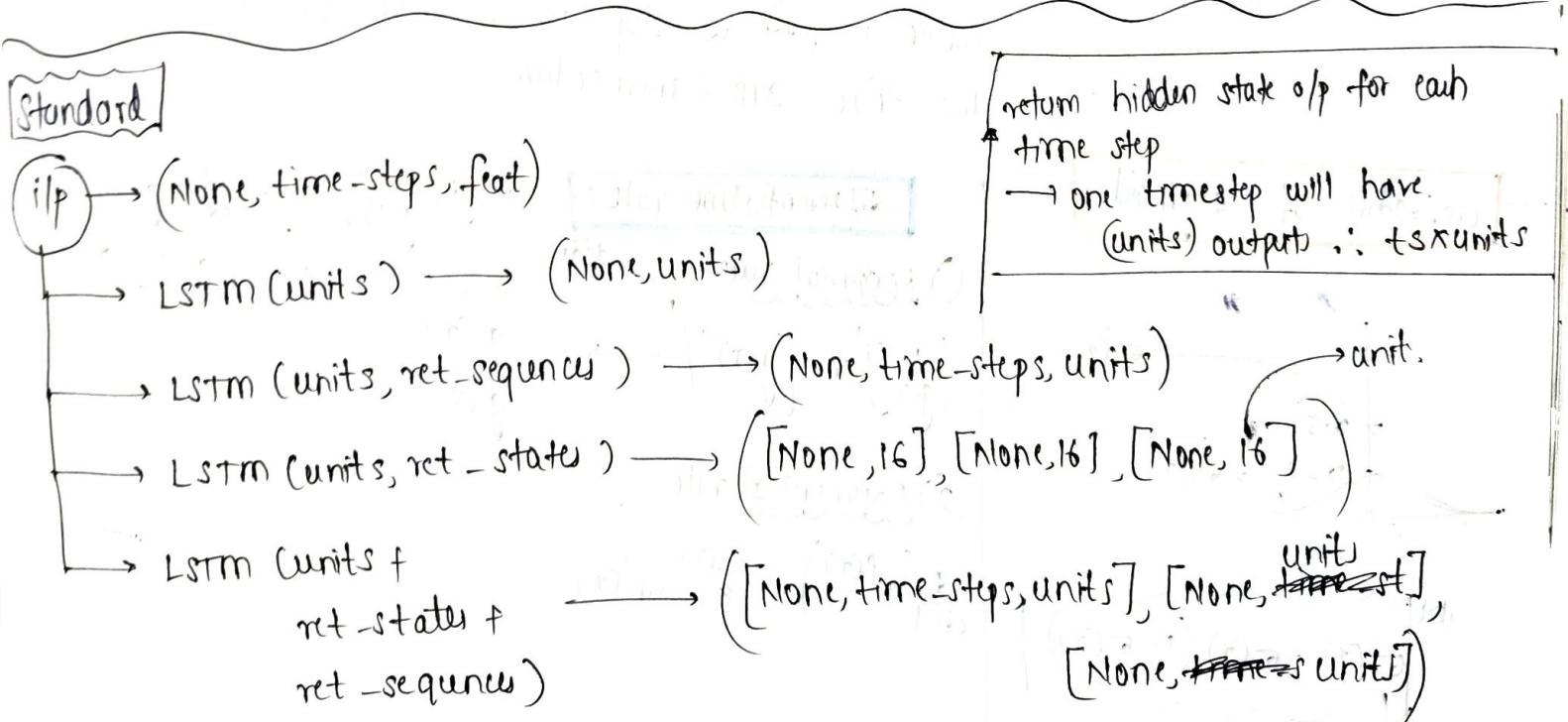
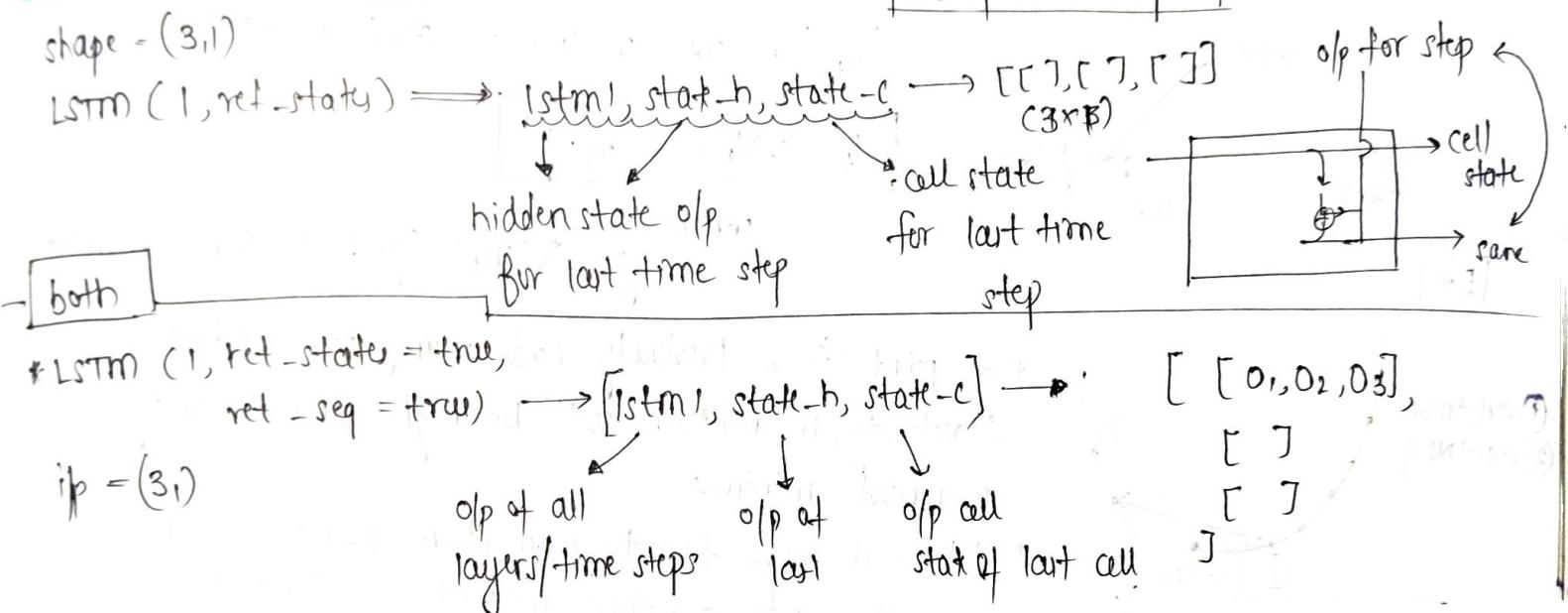
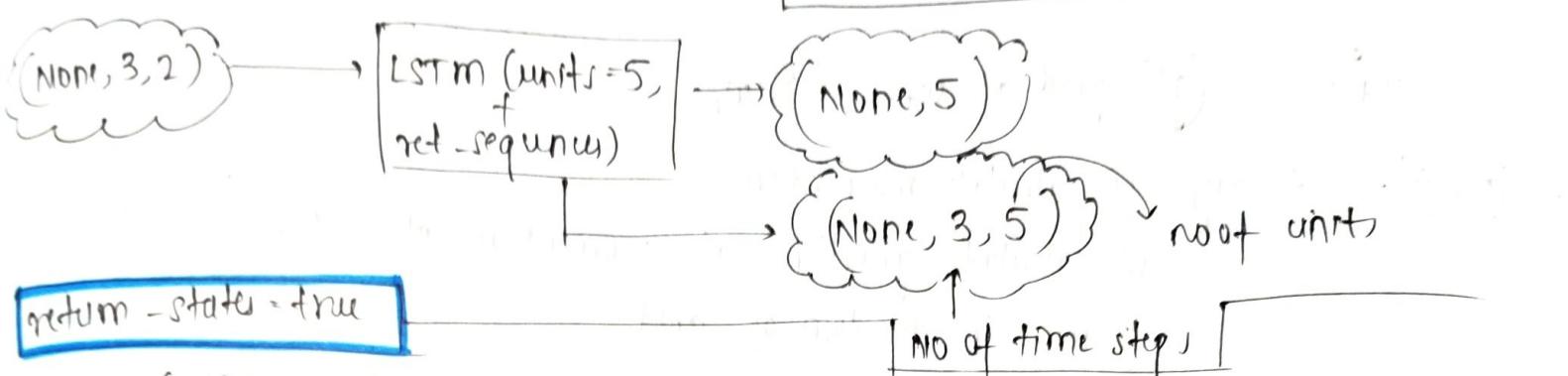
None



2 features

return-sequence = true

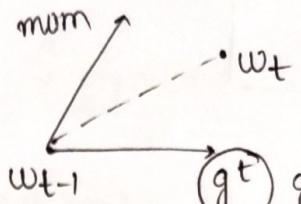
don't only return last-state o/p → $[[y_{10}, y_{20}], [y_{11}, y_{21}], [y_{12}, y_{22}]]$
 but o/p for each state and time step → $(None, 3, 2)$



Nesterov Accelerated gradient

momentum = ewa of previous gradients

$$\rightarrow \text{SGD + momentum} = w_t = w_{t-1} - [\gamma v_{t-1} + \eta g_t]$$

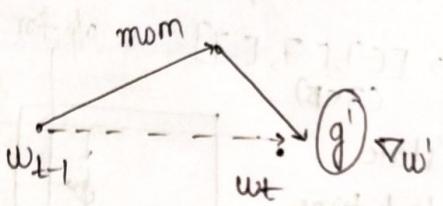


$$\text{grad at } w_{t-1} = \nabla w_t$$

[NAG] → it says ① calculate momentum

② calculate gradient at momentum pt. ($w_{t-1} - \gamma v_{t-1}$)

③ move in direction of grad



$$(w_t)' = w_t - \gamma \cdot v_{t-1} \rightarrow \text{step 1: momentum}$$

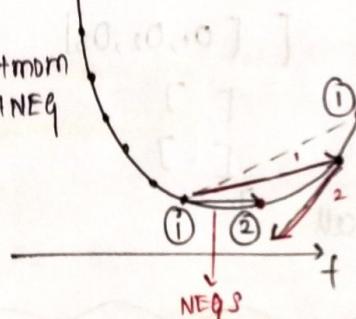
$$v_t = \gamma v_{t-1} + \eta \cdot g'$$

where,

$g' = \nabla w_t'$ (grad at w_{t-1} moment)

[Ex]

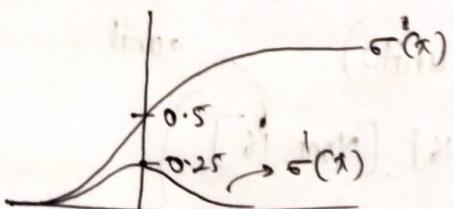
- ① SGD + mom
- ② SGD + NEG



* at point ①, we previously had many +ve updates
 \therefore velocity / momentum will be high and it may overshoot minima

* NEG: - it moves first to momentum point which is on other side, and then gradient is calculated which is **-ve** as slope is -ve \therefore it overshoots less than SGD + momentum.

derivative of sigmoid



$$\sigma'(x) = \sigma(x)(1-\sigma(x))$$

max = 0.25 at $x=0$

min = 0

differentiation rules

① Reciprocal rule

d/dx

$$\left[\frac{1}{g(x)} \right]' \text{ or } \left[g(x)^{-1} \right]' \rightarrow -\frac{g'(x)}{g(x)^2}$$

$$\frac{d}{dx} \left(\frac{1}{g(x)} \right) = -\frac{d}{dx} \left(\frac{1}{g(x)} \right)^2$$

② Exponential rule

$$\frac{d}{dx} \left(e^{f(x)} \right) = e^{f(x)} \cdot f'(x)$$

$$\Rightarrow \sigma'(x) = \frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] = \boxed{\frac{-\frac{d}{dx}(1+e^{-x})}{(1+e^{-x})^2}} \rightarrow ① \text{ using reciprocal rule}$$

solving Numerator of ①

$$-\frac{d}{dx}(1+e^{-x}) = -\left[\frac{d}{dx}(1) + e^{-x} \cdot \frac{d}{dx}(-x) \right] = -\left[0 + (e^{-x})(-1) \right]$$

$$= e^{-x}$$

$$\therefore ① \text{ becomes } \rightarrow \boxed{\frac{e^{-x}}{(1+e^{-x})^2}} \rightarrow ② \quad \therefore \sigma'(x) = \frac{e^{-x}}{(1+e^{-x})^2}$$

simplifying ② further

$$\frac{e^{-x} \cdot 1}{(1+e^{-x})(1+e^{-x})} \rightarrow \frac{(e^{-x}+1-1) \cdot 1}{(1+e^{-x})(1+e^{-x})} \rightarrow \left[\frac{1}{(1+e^{-x})} \right] \cdot \left[\frac{e^{-x}+1-1}{1+e^{-x}} \right]$$

$$\rightarrow \sigma(x) \left[\frac{e^{-x}+1-1}{(1+e^{-x})} \right] \rightarrow \sigma(x) \cdot \left[\frac{e^{-x}+1}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right] \quad \left\{ \begin{array}{l} \text{split into} \\ \text{fractions} \end{array} \right\}$$

$$= \sigma(x) [1 - \sigma(x)]$$

$$\therefore \sigma'(x) = \sigma(x) \cdot (1 - \sigma(x))$$

YOLO Algorithm

object detection

Object Detection

* focus on bounding boxes

and not [on pixels] → use image segmentation (slower and not suitable for real-time apps)

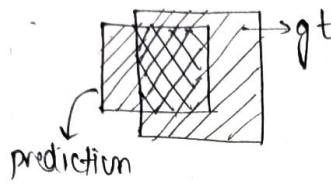
* IP → algo → o/p image

with bounding boxes with confidence level (0-1)

Performance Metrics

$$IoU = \frac{\text{overlap area}}{\text{union area}}$$

gt (ground truth) - human labelled data



* prediction is correct if $IoU \geq 0.5$

Steps

① convert o/p to binary ie True or false using IoU. & bounding boxes

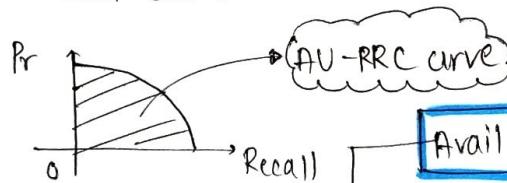
② Image may have multiple bounding boxes,

so,

* for each class (chair, person...)

→ calculate avg-precision = area under pre-recall curve

→ plot Pr-recall curve



Avg-options

* mAP (mean average precision)

* R-CNN (region-CNN) * FPN - FRCN

* SSD (single shot detector)

* RetinaNet

Trade off

→ speed vs mAP

video based

(self-dr cars)
real time face

medical diagnosis (precision is imp)

or
OCR

YOLO → faster and good mAP

→ YOLOv3-320 → input image size

YOLOv3-416

YOLOv3-608

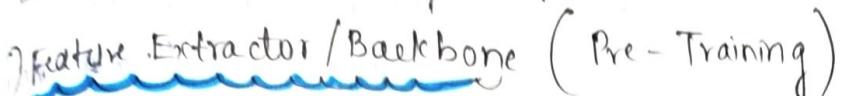
Note: as image size ↑ time ↑

YOLO v3

45

Architecture

→ only used for initialize



YOLOV2
V3

DarkNet-19
DarkNet-53

DarkNet - 53

53 conv-layer (fully convolutional)

→ padding = same

Layer	Filters	Size	Op
conv	32	(3x3)	256×256
conv	64	(3x3)	128×128
		stride=2	
IX	conv	32	(1x1)
	conv	64	(3x3)
	residual		128×128
2X	conv	128	$3 \times 3 / 2$
	conv	64	(1x1)
	conv	128	(3x3)
	residual		64×64
8X	conv	256	$3 \times 3 / 2$
	conv	128	(1x1)
	conv	256	(3x3)
	residual		32×32
8X	conv	512	$3 \times 3 / 2$
	conv	256	(1x1)
	conv	512	(3x3)
	residual		16×16
4X	conv	1024	$3 \times 3 / 2$
	conv	512	(1x1)
	conv	1024	(3x3)
	residual		8×8

Augpool global

Sulfuric

416x416x3 → 13x13x1024

(416/32) last filter size

Characteristics

- ① 53 conv layers
 - ② fully conv neural net
 - ③ Each conv
 - = conv → Batch Norm → Leaky ReLu
 - ④ conv with stride-2 → acts as max-pool
 - ⑤ no maxpooling

* whole model is trained
on ImageNet and save
weights

→ this only detects image
objects but not bounding
boxes

why we DarkNet when we had
Resnets?

Baseline	Top-5	Bops	BFlops/s	FPS
Res-152	77	2930	1090	37
Dark	77.2	18.7	1457	78
Net				

→ darknet almost gives same or
little less accuracy comp to resNet

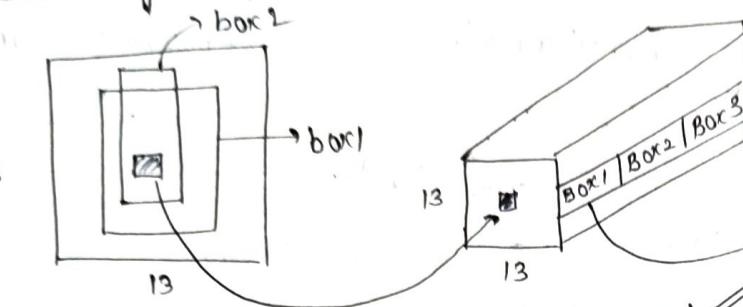
* BN ops is less (billions of ops needed)

- * BFCops/s is high
- * and FPS is high

④ ② Bounding boxes and output

IP $416 \times 416 \times 3 \xrightarrow{\text{DNN}} 13 \times 13 \times 1024 \xrightarrow{\text{Box}} 13 \times 13 \times 425$ O/P

How do you represent boxes?



one image (13×13)

$(K+5) * B$

Ex: consider 80 classes

* box $i = (P_0, x, y, w, h, c_1, \dots, c_{80})$

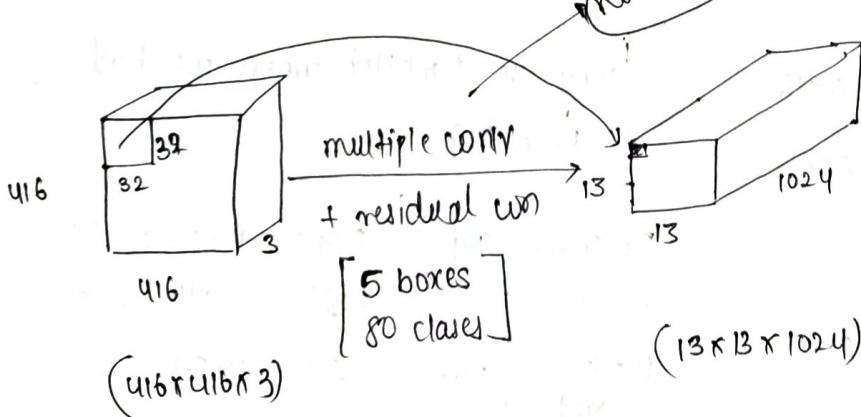
$$\begin{aligned} * \text{ score} &= P_0 * C = \begin{bmatrix} P_0 C_1 \\ P_0 C_2 \\ P_0 C_3 \\ \vdots \\ P_0 C_{80} \end{bmatrix} = \begin{bmatrix} 0.12 \\ 0.21 \\ 0.80 \\ \vdots \\ 0.10 \end{bmatrix} \rightarrow \text{probability that object bounded } \in \text{ class } 2 \\ &\qquad\qquad\qquad P(C=c=k) = P_0 * P_k \text{ or } P_0 * C_k \end{aligned}$$

* assume we get $P_0 C_3 = 0.80$ as max

\therefore the box bounds object

belonging to class = 3

not darkest



assume $K=80$ i.e.
 $B=85$

(Each cell may belong to many boxes) How of boxes

Each bounding box

x	y	tw	th	P_0	C_0, \dots, C_K	$* B$
-----	-----	-------------	-------------	-------	-------------------	-------

(x, y) → centre of box with dimensions $w \times h$

$P_0 \rightarrow$ objective score → (0/1)
whether it contains an obj / not

$C_0 \rightarrow C_K$: class scores for K nof

$$\frac{416}{13} = 32$$

\therefore Each pixel in o/p is of size (13×1024) is captured from $(32 \times 32 \times 3)$ grid input.

Effective Receptive Field

- Anchor Boxes → prior boxes where you highly probable to find
 * Tried to predict (tx, ty, tw, th) directly but couldn't get good results
 So,
- ① Fix some predefined boxes called anchor-boxes : $[Cx, Cy, Pw, Ph]$
 - ② define bounding boxes in terms of anchor boxes →
- Given box coordinates and anchor boxes you can find $[bx, by, bw, bh]$
- Named box + anchor box → actual bounding box
- Objectness score P_o
- Is there any object or not? 0/1
- ↳ have a log-regression for each box $[c_0 \dots c_k]$
- * Tried software — didn't work because of hierarchy presence
- Soln
 Build one logistic regression for each class i.e. u need $P_o + [c_0 \dots c_k]$
- log-regressions
- LOSS
- log-loss for P_o
 + log-loss for $(c_0 \dots c_k)$
 + $\lambda_{co-ord} * \text{sq-loss for } tx, ty, \sqrt{tw}, \sqrt{th}$
 all bounding boxes containing object in y_{train}
- Anchor Boxes → prior boxes where you highly probable to find
 * Tried to predict (tx, ty, tw, th) directly but couldn't get good results
 So,
- (b_x, b_y) = priors
-
- $b_x = \sigma(tx) + C_x$
 $b_y = \sigma(ty) + C_y$
 $b_w = P_w \cdot e^{t_w}$
 $b_h = P_h \cdot e^{t_h}$
- $\sigma = 0-1$
- (b_x, b_y) should be at most 4 pixels away from (Cx, Cy)
- How to find how many anchor boxes and location?
- apply k-means clustering on COCO dataset which is set of all bounding boxes
- $K=5$ was found to be best which indicates that there are 5 anchor boxes have high probability region
- if for any box, if P_o is very small
 → it has very less chance of enclosing any object
 ∴ can be ignored
- $\lambda_{co-ord} = 5 \rightarrow t_x, t_y \rightarrow 0 \rightarrow \infty$
 $\lambda_{noobj} = 0.5 \rightarrow P_o = (0 \rightarrow 1)$
- $\sum \lambda_{noobj} * \text{log-loss for } (P_o)$
 all boxes in predicted which don't have obj
 (predicted boxes not present in y_{train})

④ multiscale prediction

* 13×13 image feature can only detect large objects

Idea

* predict with different grid sizes

↳ How to get?

Extract from layers in darkNet

① last layer $\rightarrow 13 \times 13 \times 1024$

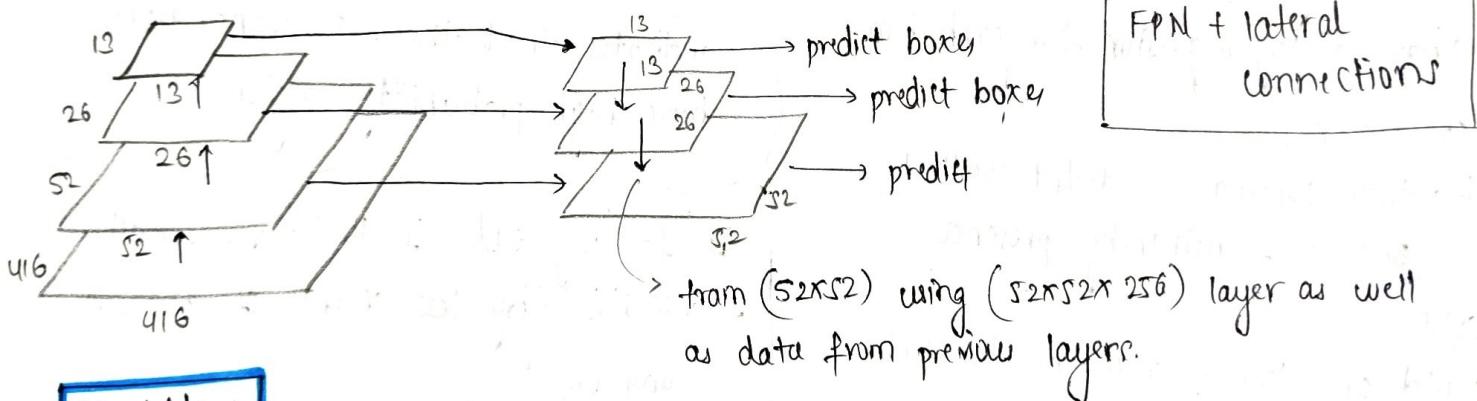
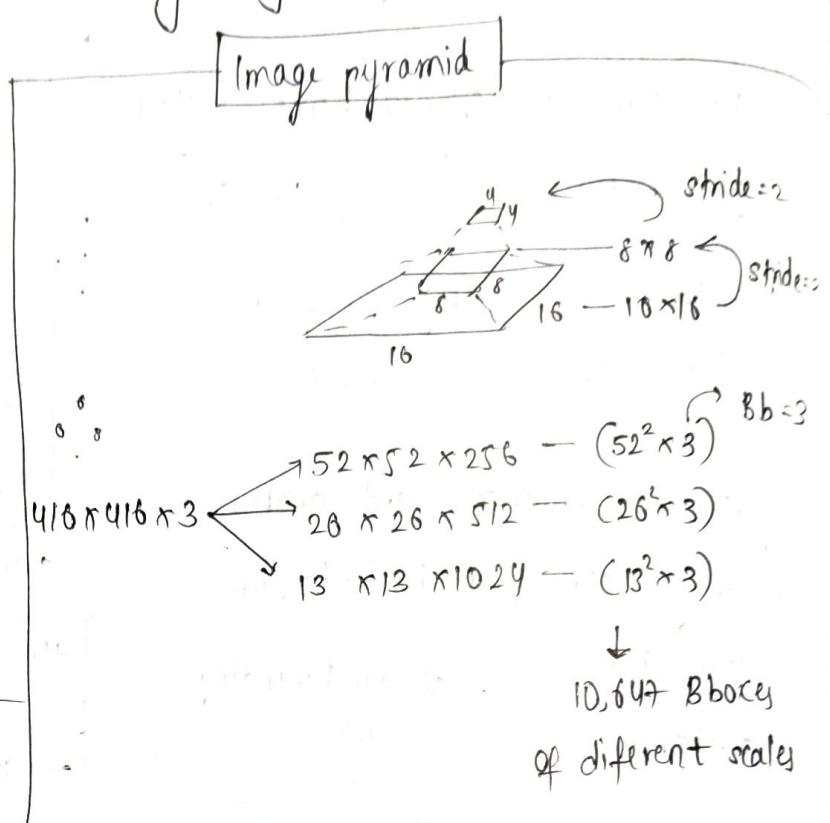
② last but-1 $\rightarrow 26 \times 26 \times 512$
conv layer with $s=2$

③ 2nd last $\rightarrow 52 \times 52 \times 256$
conv layer
 $s=2$

Feature pyramid networks

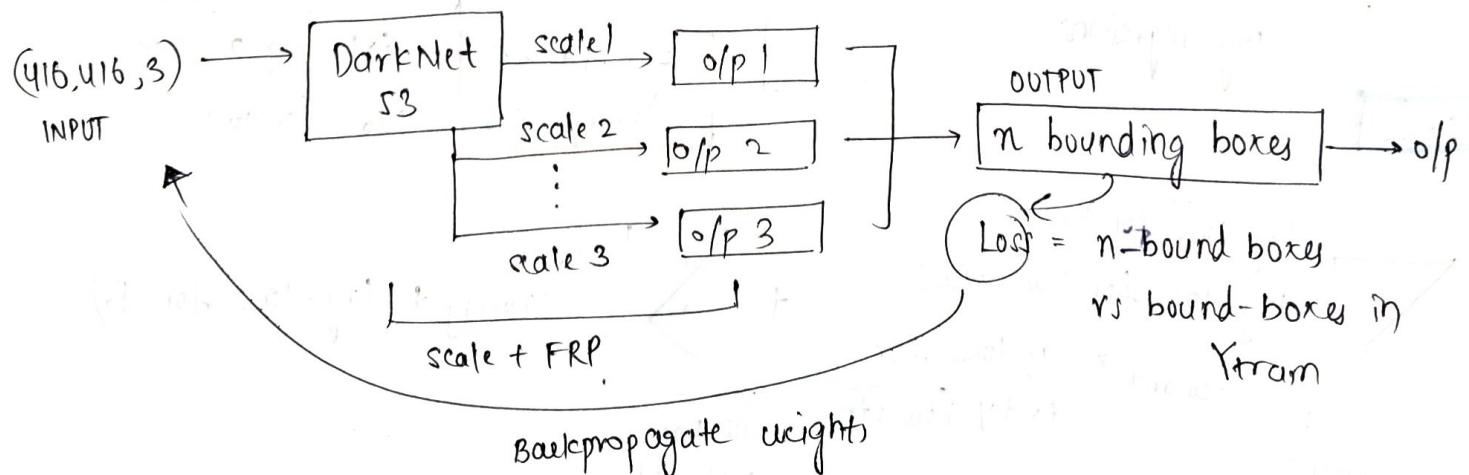
FPN

→ not only use different grid sizes
to predict but also cross-grid information to use

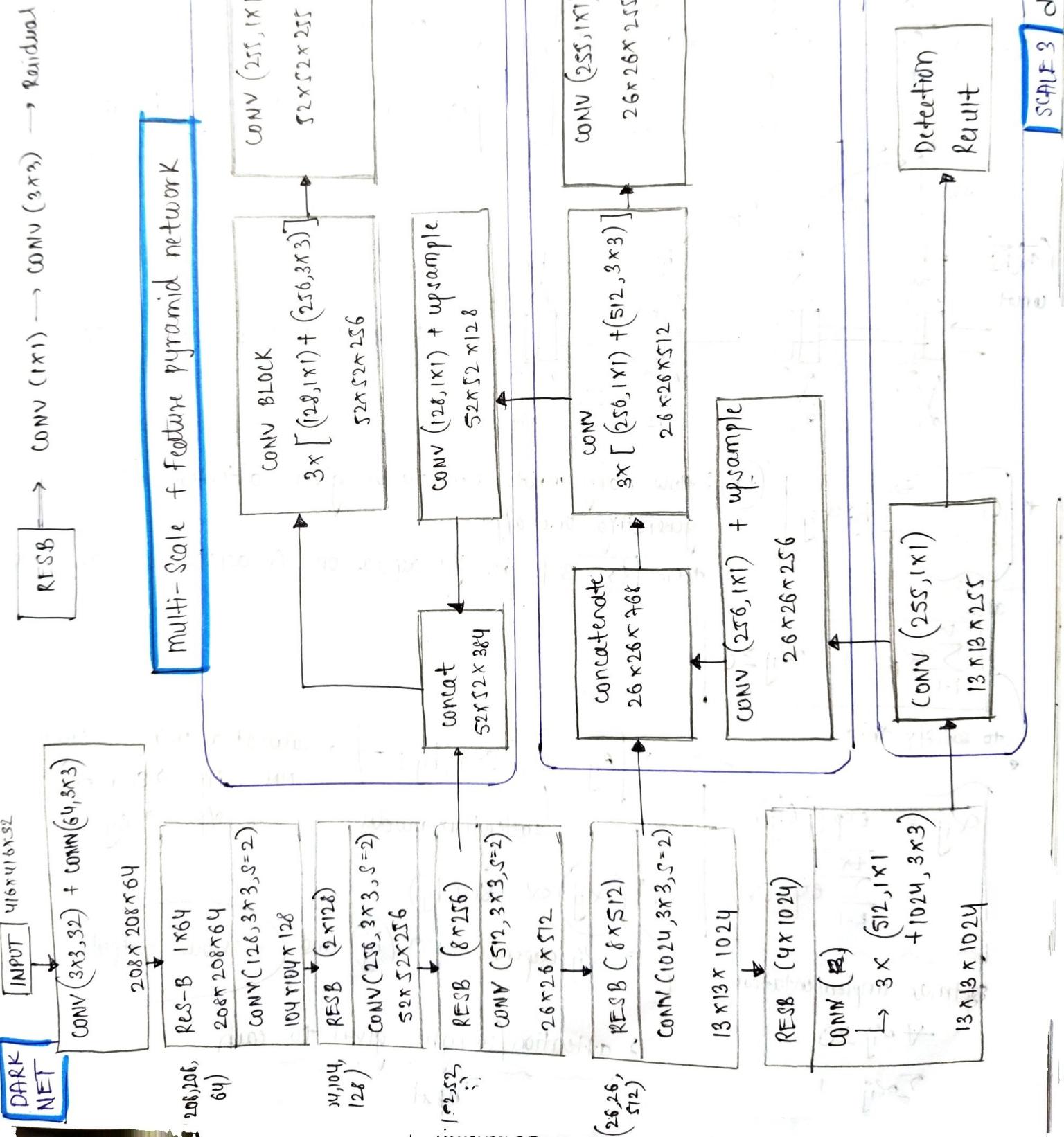


FPN + lateral connections

Workflow



Yolo V3
ARCHITECTURE



50

Attention Models

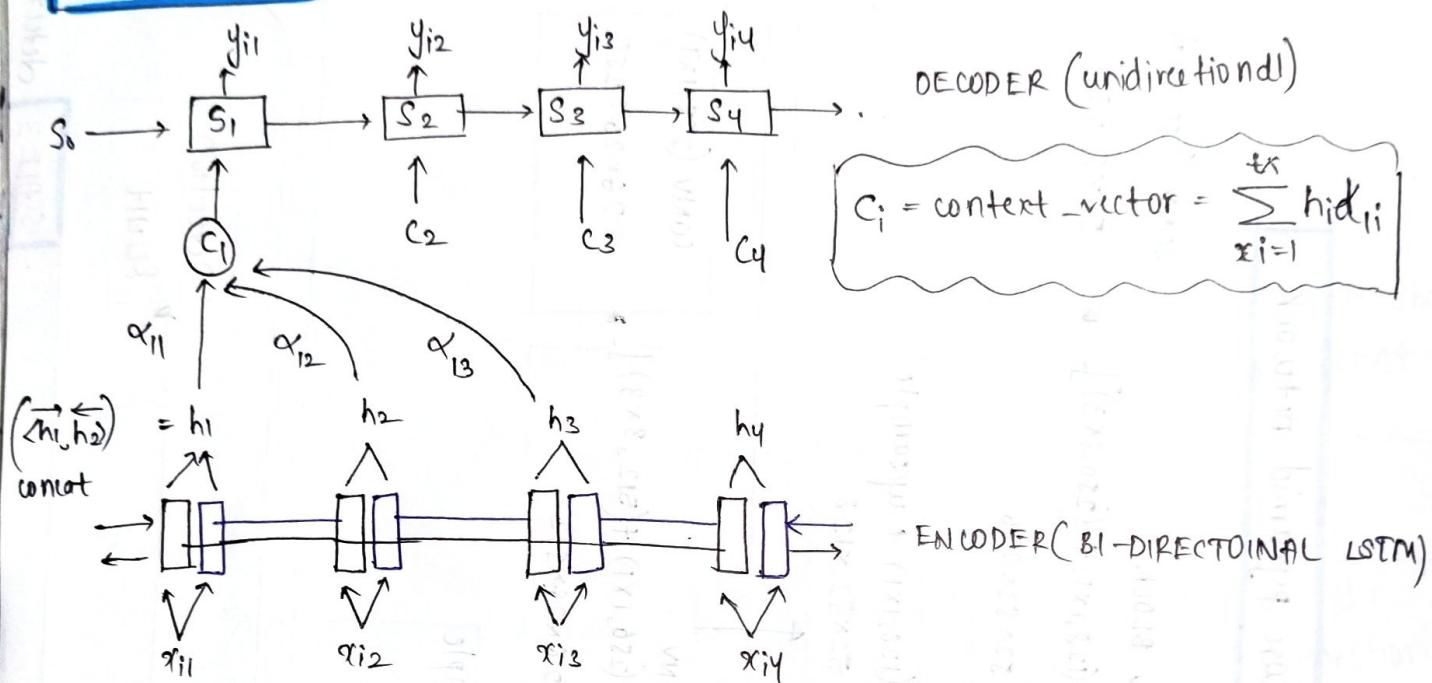
this may not capture the essence of whole sentence

why you need? In encoders we pass only one o/p from encoder

* Encoders decoders don't work well for long sentences

attention? We give attention to a few words (t_x) translate them and shift attention to next set of words (Human characteristic)

Architecture



$$C_i = \sum_{j=1}^{t_x} h_j \alpha_{ij}$$

t_x : How many inputs need to be given attention for generating one o/p.

Here $t_x = 3$ ie y_{i1} depends on C_i depends on x_{i1}, x_{i2}, x_{i3}

also,

$$\sum_{j=1}^{t_x} \alpha_{ij} = 1, \alpha_{ij} \geq 0$$

to satisfy this

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{t_x} \exp(e_{ik})}$$

softmax implementation

$$+ \alpha_{ij} \geq 0$$

$$\sum \alpha_{ij} = 1$$

$$e_{ij} = a(S_{i-1}, h_j) \rightarrow \begin{cases} \text{calculated using 2-layer NN} \\ i/p \rightarrow S_{i-1}, h_j \\ o/p \rightarrow e_{ij} \end{cases}$$

$$\therefore \alpha_{ij} \propto (S_{i-1}, h_j)$$

α_{ij} depends upon (h_j) and previous output

attention / weight given to each input

Time complexity

$K_1 = \text{len}(\text{input})$

$K_2 = \text{len}(\text{output})$

$$TC = O(K_1 \cdot K_2)$$

if $TX = K_2$,

$$\text{then } TC = O(K_1 \cdot K_2)$$

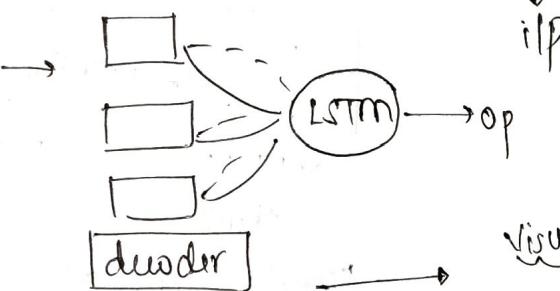
i.e consider whole sentence as
attention to generate o/p.

Image Captioning

i/p → feature extraction using CNN
encoder

generate K d-dim vectors

each vector representing one region of image



α_{ij} importance

Ex: for calculating α_{12} you use $\alpha_1, \alpha_2, \alpha_3$

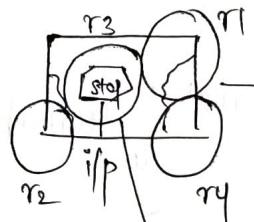
in this case if α_1 is large than
that particular word had more impact
on prediction

Ex: french → english

	you	and	call	→ o/p
ti	0.9	0	1	
ans	0.1	2	0	
me	0	2	10	

to predict $w_3 = \text{call}$
me was most useful.

visualize α_{ij}



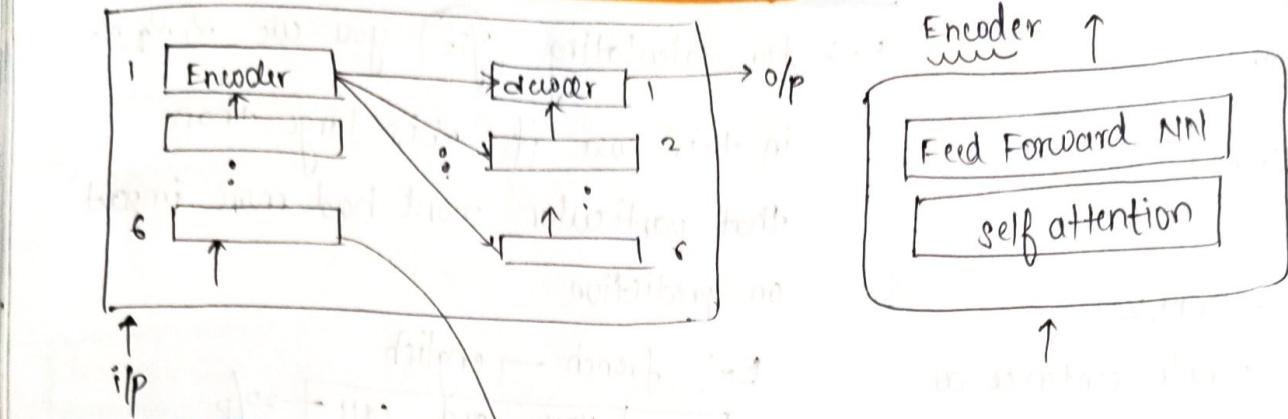
A mountain road with
stop sign

see highest α_{ij}

say highest is α_{23}

∴ region (3) contributed
most

Transformers



ENCODER

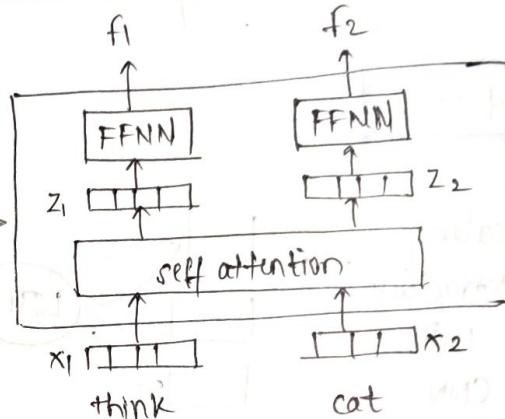
Self attention

the boy crossed the road as
it was empty
what does 'it' refer/connect to

steps

hyper-params :

$d_k = 64$, W^Q - Query mat } learnable
 W^K - keys
 W^V - value



* inputs sent at once

No time-based sending

Input Embedded	x_1 [512] think	x_2 [512] cat
Query	q_1 [64]	q_2
Key	k_1	k_2
Value	v_1	v_2
Score	$q_1 \cdot k_1 = 112$, $q_1 \cdot k_2 = 96$, $q_1 \cdot k_n = ?$	
$\left(\frac{\text{score}}{\sqrt{d_k}}\right)$ stable Δ grad	14	12
softmax	0.08	0.12

$$q_i = x_i \cdot W^Q$$

$$k_i = x_i \cdot W^K$$

$$v_i = x_i \cdot W^V$$

$$(n - \text{no of work}) = q_1 k_1 - n$$

soft x val v_1 [] + v_2 [] + ... → here $v_1 > v_2$ as v_1 is of same word

sum (z_1) []

(53)

matrix computation

$$X \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \times W^Q = Q \rightarrow \text{softmax} \left(\frac{Q \times K^T}{\sqrt{dk}} \right) V$$

$$X \times W^K = K$$

$$X \times W^V = V$$

$$= \boxed{Z} \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}$$

: self-attention power of some words which matter more in its neighborhood (ie same sentence)

one-headed : only 1 set of matrix

multi-headed (8-pairs of W^Q, W^K, W^V) → same steps but with (W_Q^a, W_K^a, W_V^a) to $(W_Q^a, W_K^a, W_V^a) = 8$ (2is)

: for X we get $(z_0, z_1, z_2 \dots z_7)$

$$\boxed{\text{softmax}(z_0, z_1, \dots, z_7) \times W^V = Z}$$

positional Embedding

we pass all words at once, how to maintain time/position

$$\begin{array}{l} \text{input: je seri wuld} \\ x_1 \boxed{} \quad x_2 \boxed{} \quad x_3 \boxed{} \\ + \quad + \quad f \\ p_1 \boxed{} \quad p_2 \boxed{} \quad p_3 \boxed{} \\ \downarrow \quad \downarrow \quad \downarrow \\ x'_1 \boxed{} \quad x'_2 \boxed{} \quad x'_3 \boxed{} \end{array}$$

embeddings

positional encoding

indicates position as well as distance b/w words

* pos vectors are static

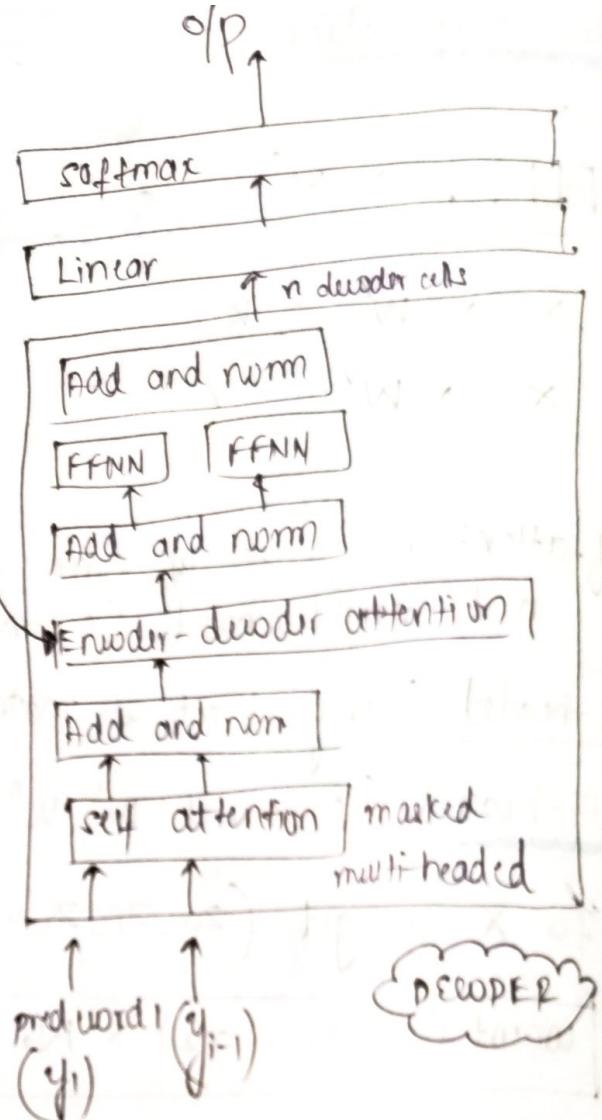
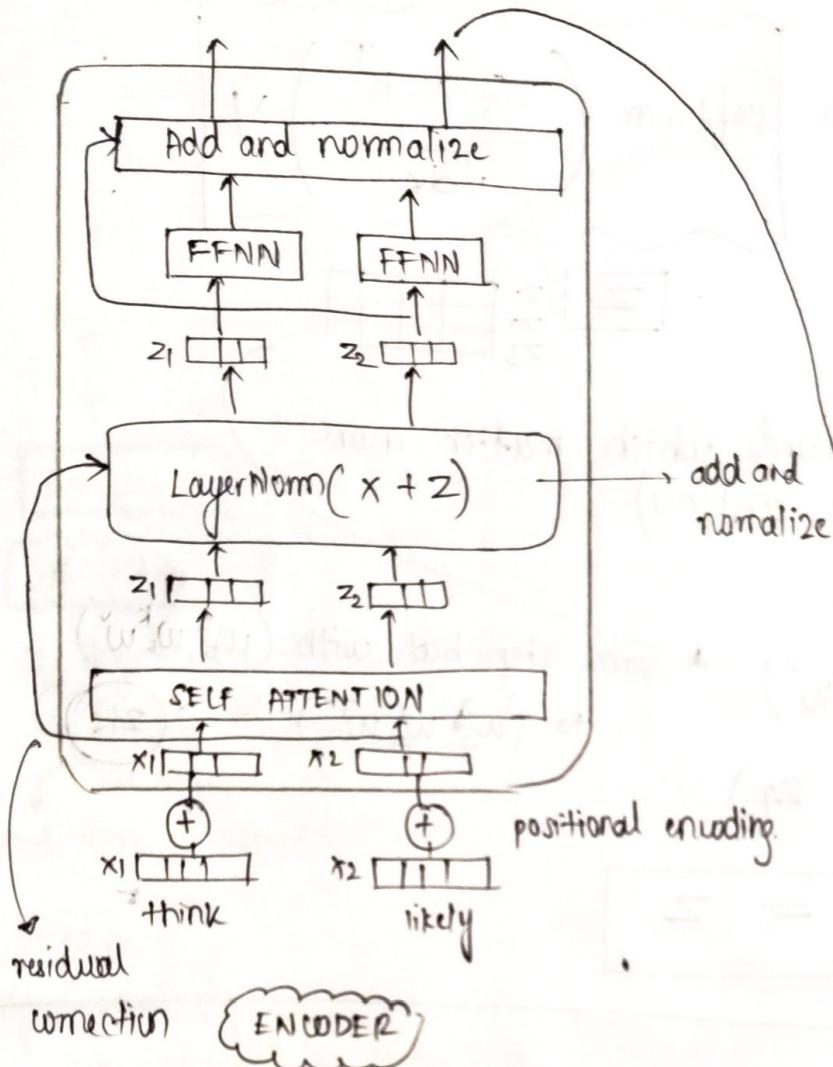
ie p_1, p_2, p_3 are precalculated

st $p_1 > p_2 > p_3$ position

$(p_1 - p_2) < (p_1 - p_3)$ distance

54

ENCODER architecture



Train phase

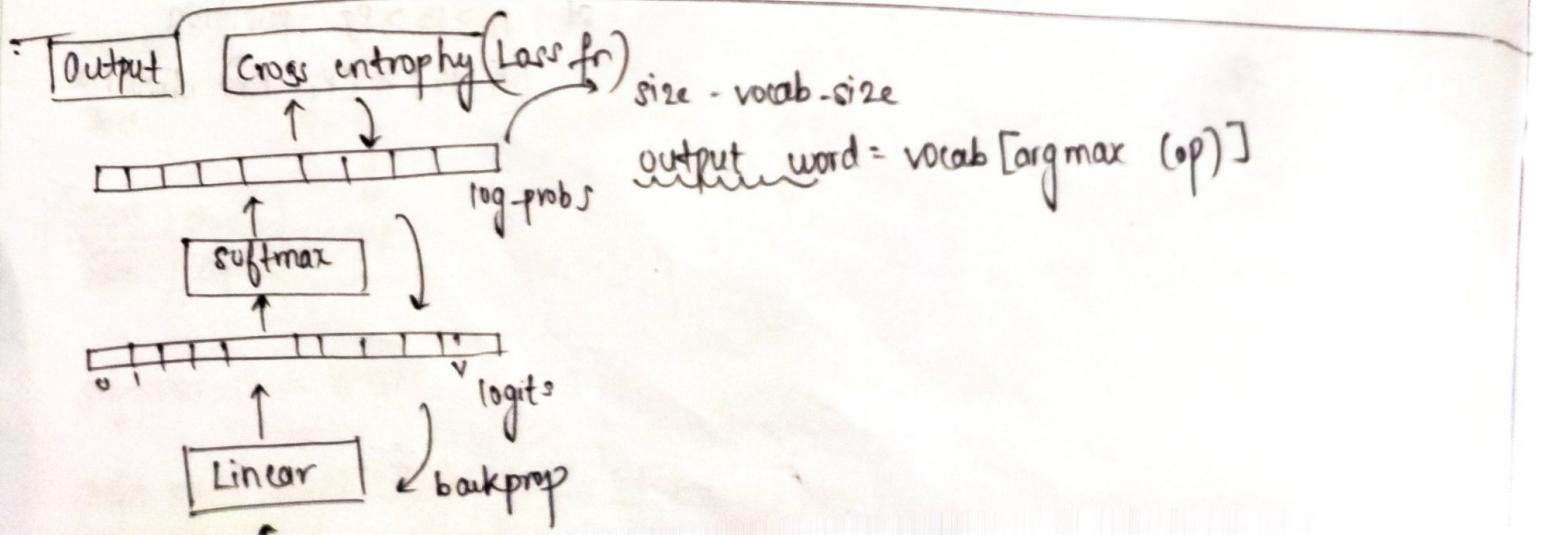
① pass whole sentence to encoder, you get a output

② pass it to decoder and predict first word

second word :- don't touch encoder

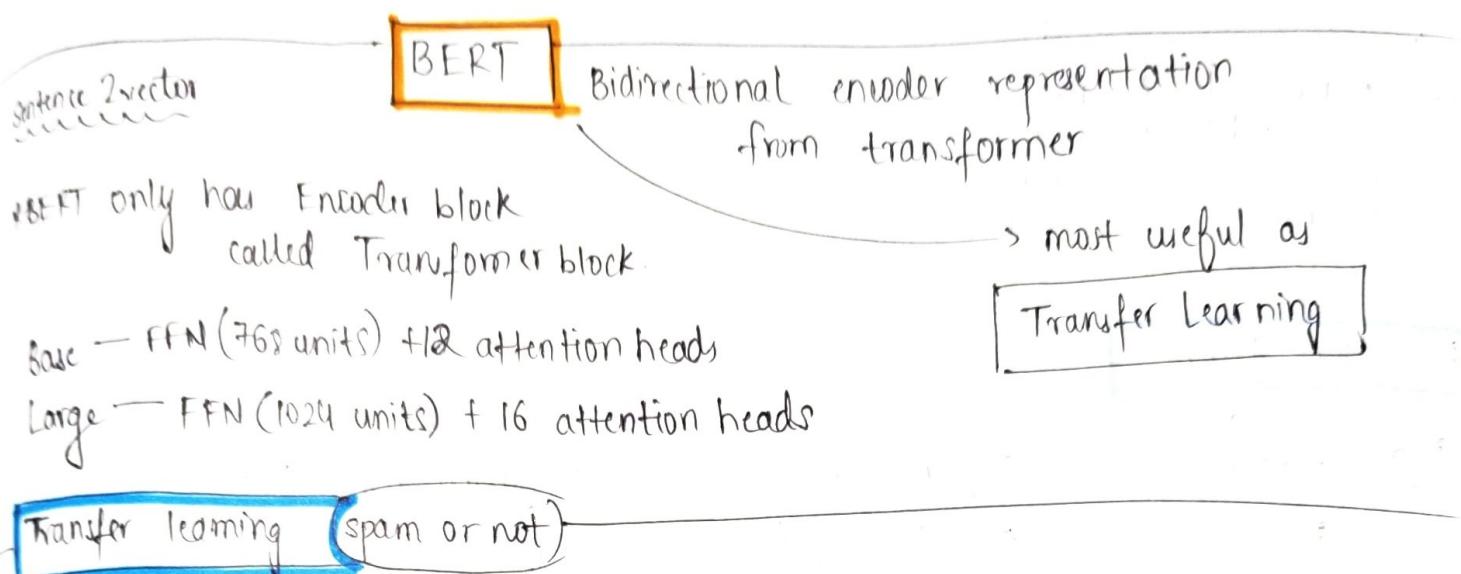
→ pass first word to encoder → o/p second word

third word :- pass 1st and 2nd o/p word → third o/p word

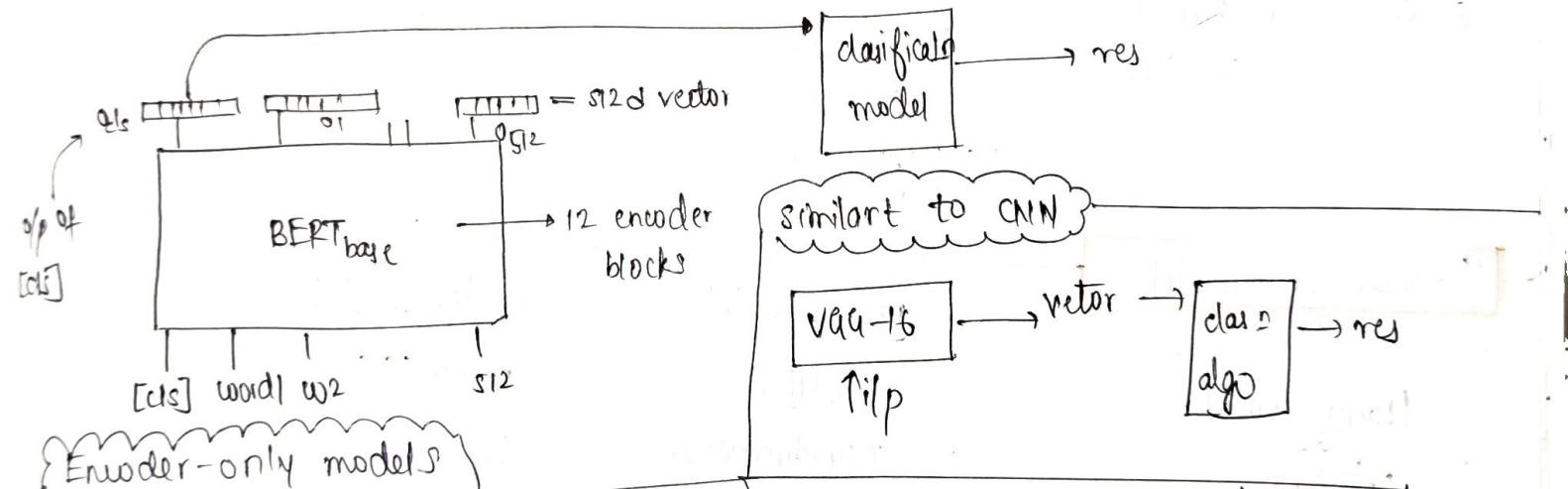


Advantages of Transformers

- * Faster than RNN models as all data is ingested once.
- * state-of-the-art (2018)
- * Inspired from CNN (self-attention/neighborhood \approx convolutions in CNN).



Base Base Bert; can take upto 512 worded sentency and produces 512 o/p's each of dimension (768)



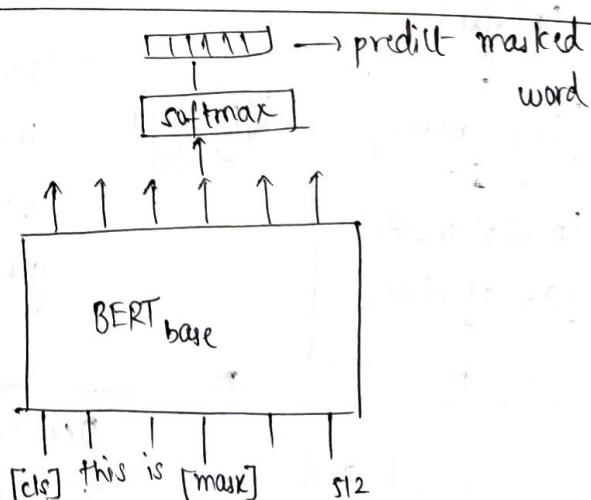
BERT Training \approx word-2-vec training

* Bidirectional models, while training look @ words before and after

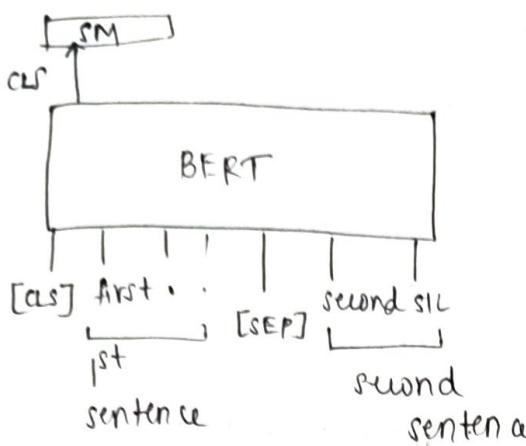
steps

* provide any sentence as ilp and while doing so randomly mask 15% of words/token

* try to predict masked tokens given other words and optimize



BERT with 2 sentences

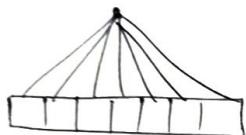


BERT - training

NWP = next word prediction

Attention types

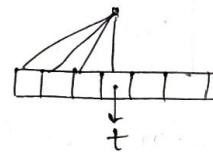
① Self attention



- + depends on all words
- & Encoder (no timestep concept)

② masked attention

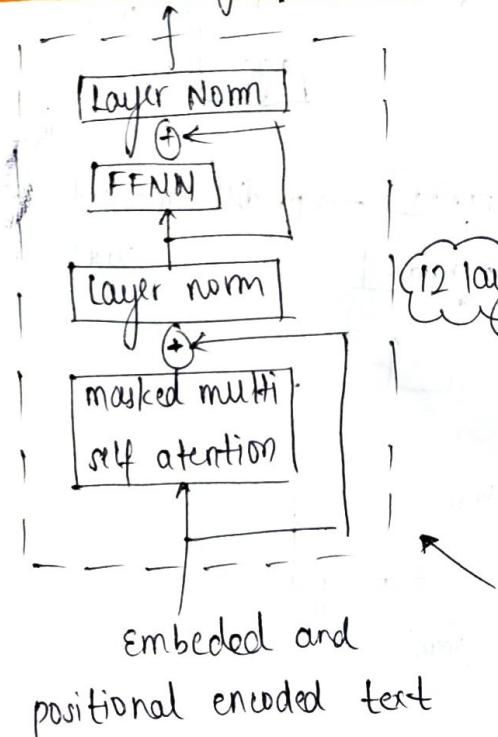
- + at timestep (t_i), self attention calculated on $t_0 \dots t_{i-1}$
- + used in decoders as time based



③ sparse self attention

- + if $\text{len}(\text{sentence}) = N$ choose only \sqrt{N} words for calculating self attention
- + Reduces time complexity.

Decoder only models



Decoder models

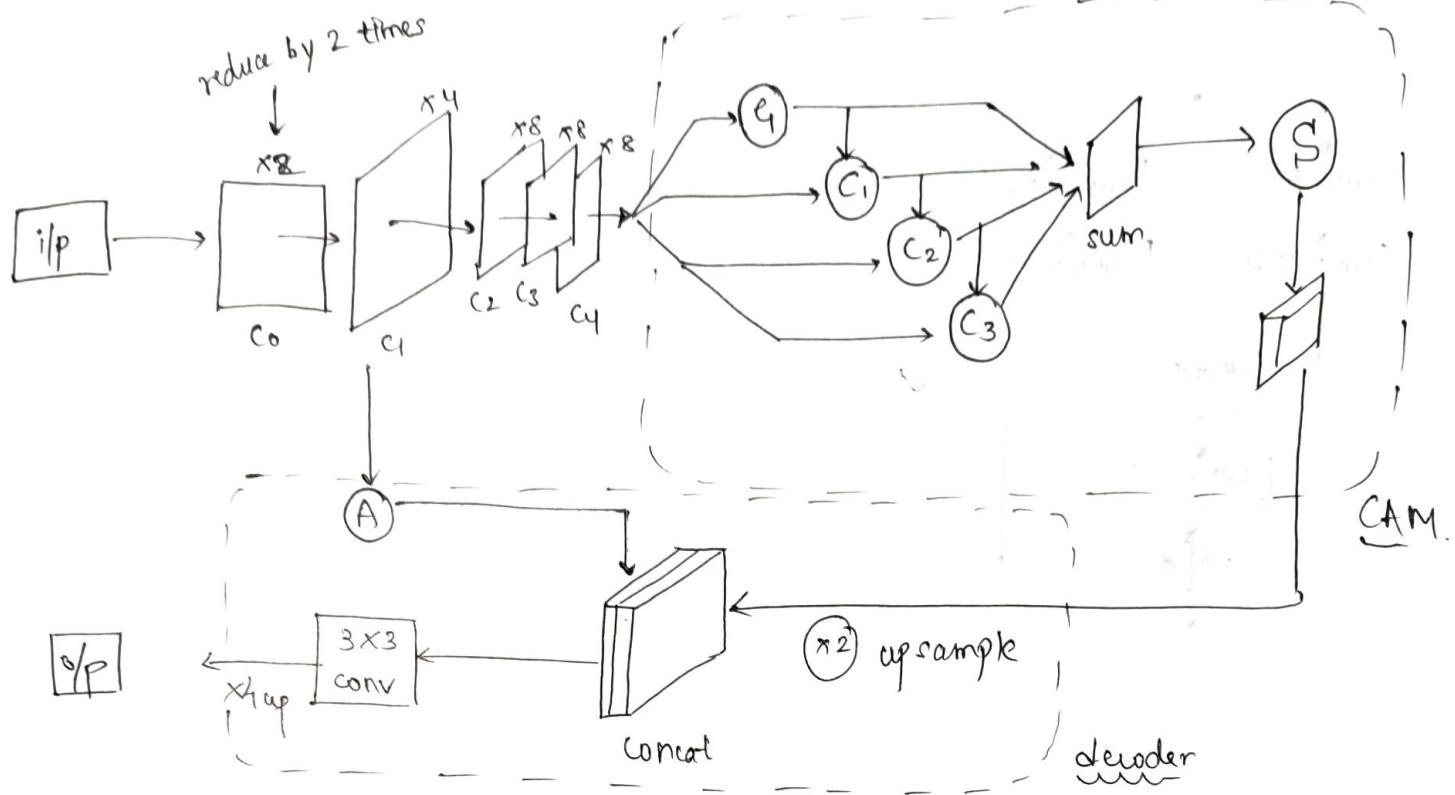
- GPT
- unidirectional
- normally use masked attention
- used more in generative tasks

Encoder models

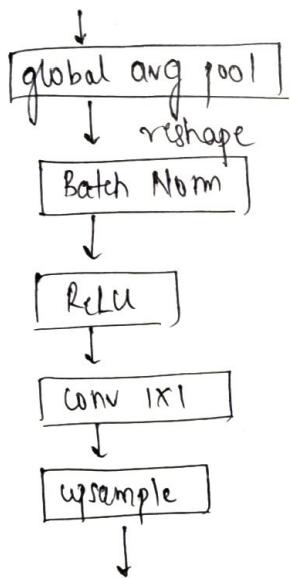
- BERT
- Bidirectional
- simple self attention
- used mostly for text representation/embedding tasks

Decoder only Transformers = GPT-1

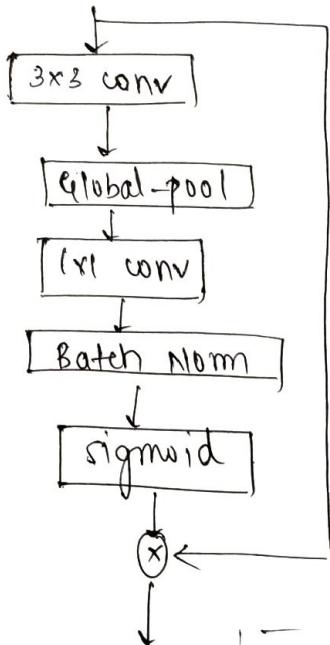
CANET - Image segmentation



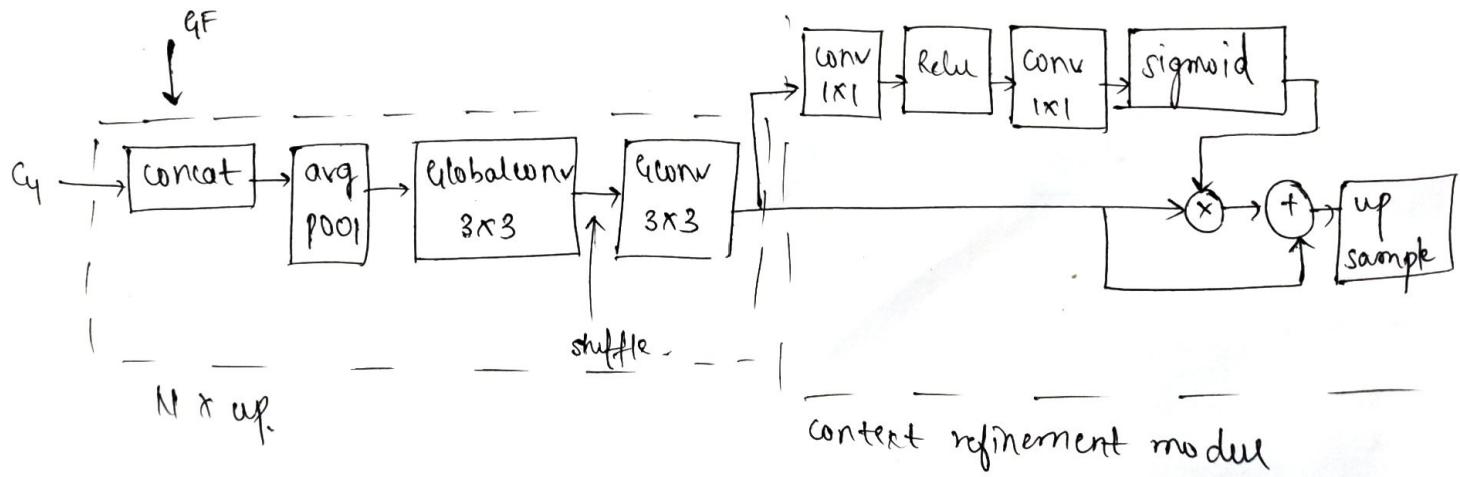
① ⑨ Global flow



③ ⑤ Feature selection module

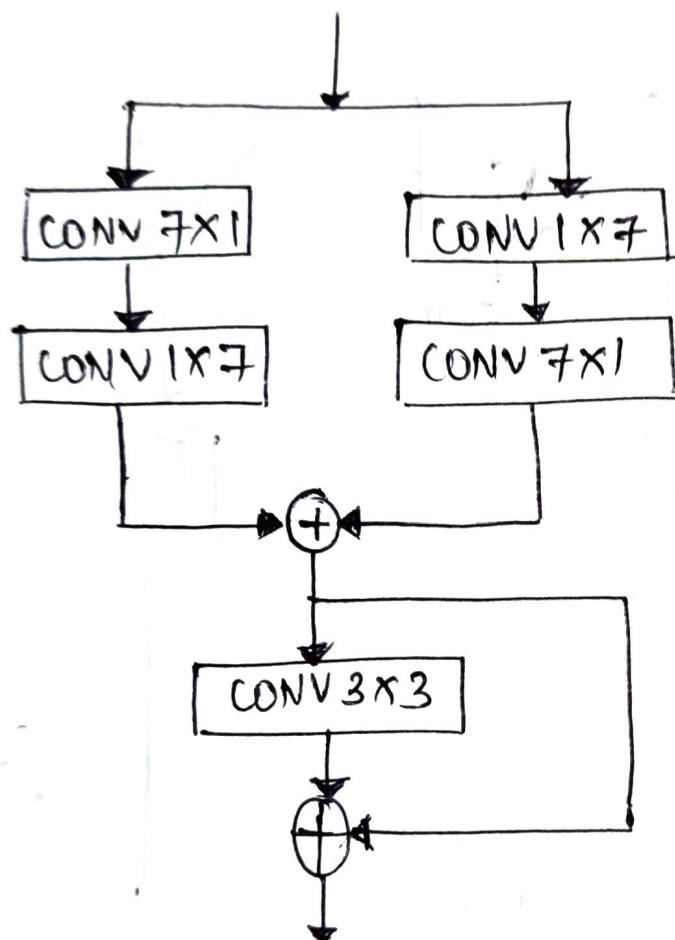


② ⑦ Context fusion module



58

AGCN
Adapted global convolution n/w

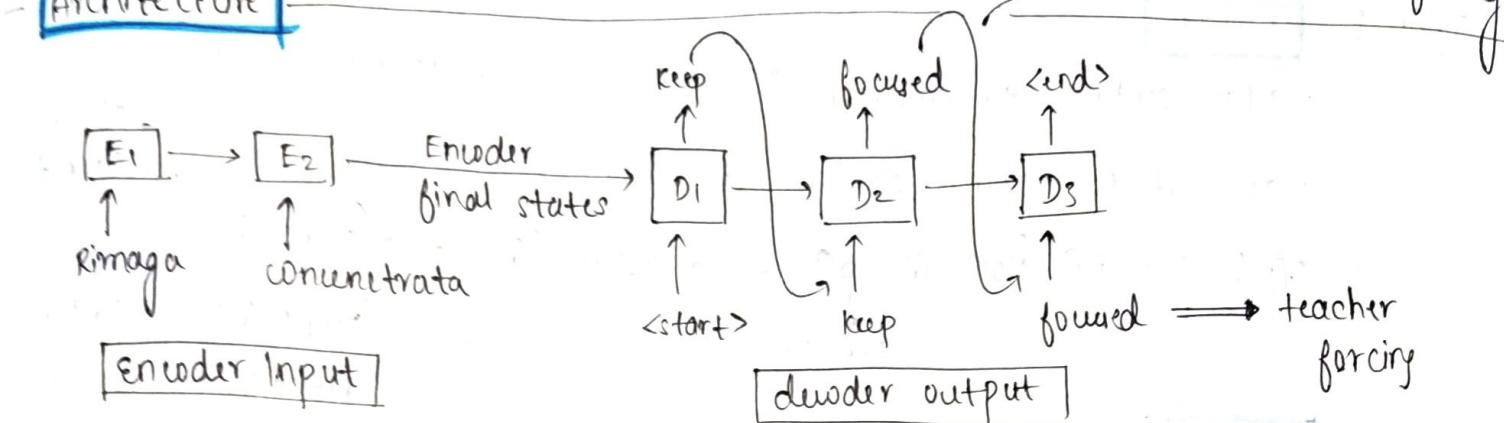


time - <-- sequence

use case

→ Italian to English translation

Architecture



data

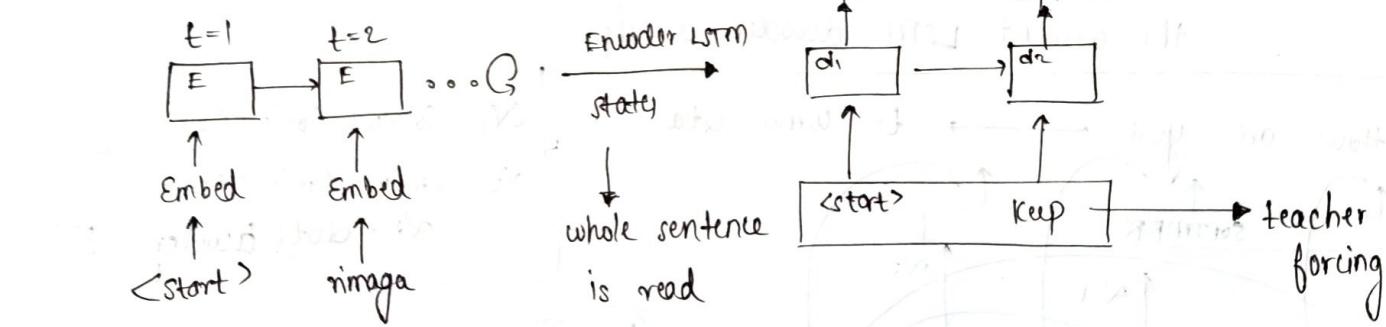
(x) rimaga concentr → Encoder i/p

(y) keep focused → decoder i/p: <start> keep focused
decoder o/p: keep focused <end>

apply tokenizer to convert text into integers.

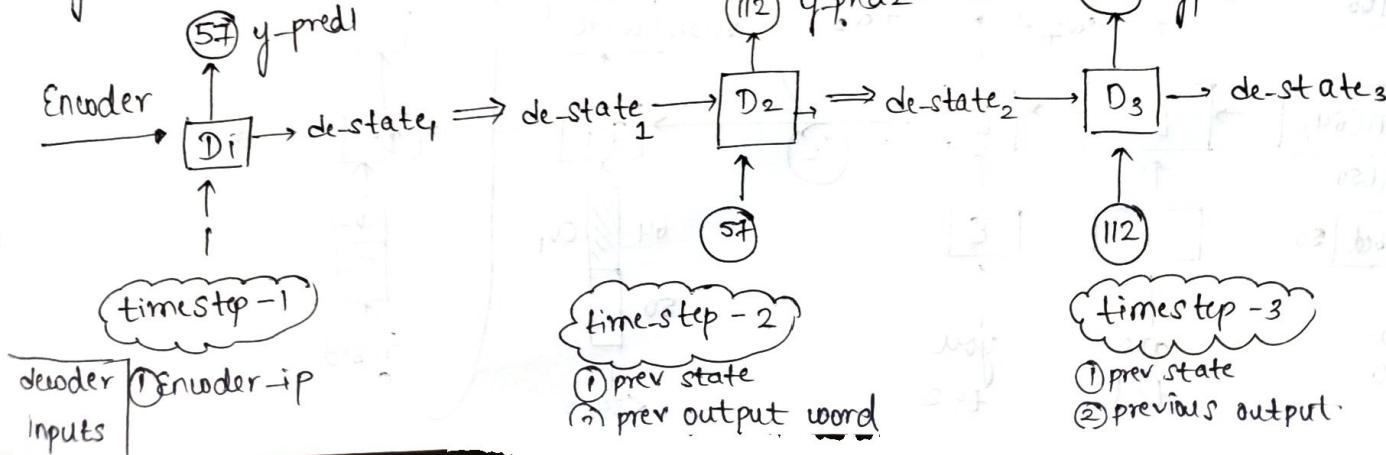
↪ (batch size, sentence-length) $\xrightarrow[\text{layer}]{\text{Embedding}}$ (batch-size, sentence-len, output-dim)

Ex



Inference

↪ you don't have teacher forcing, sentences



Seq-2-Sq Attention

Steps

Train :- Encoder

→ Pass sentence to encoder : words → embedding layer

↳ pass it to LSTM (seq=True, st=True)

O/p: Ex (1) (None, 600, 60) → LSTM (64) → (None, 600, 64)

↓ embeded

o/p. of each ts

(None, 64) → last h_t
(None, 64) → last c_t

Decoder

→ pass desired o/p sentence (+teacher forcing)

at t_s (t)

① using h_{t-1} and h_e find context vector

h_{t-1} : decoder hidden-state
at prev t_s

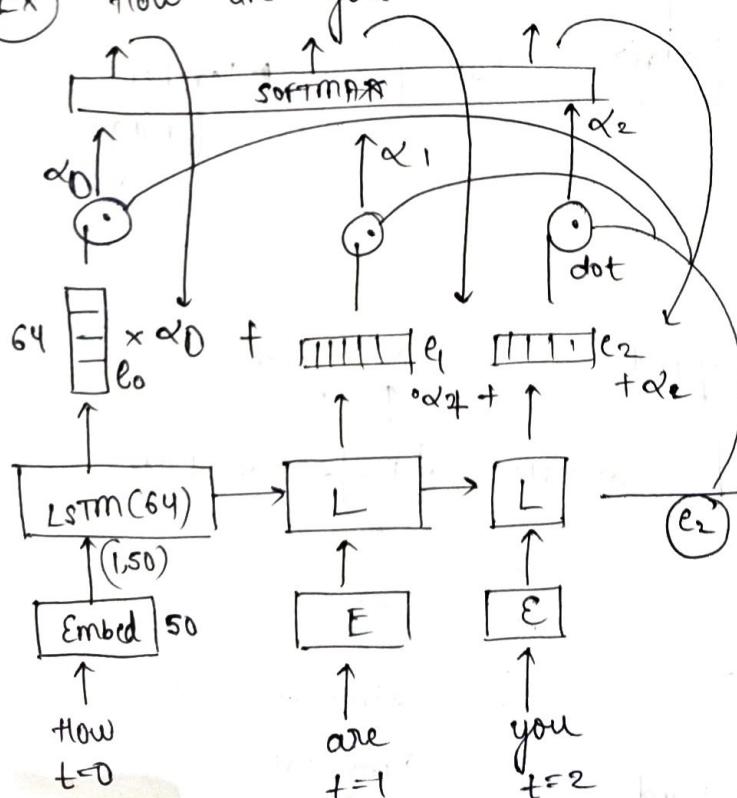
h_e : (h₁, h₂, ..., h_n)
encoder outputs

$$\begin{bmatrix} h_1 \\ h_2 \\ h_n \end{bmatrix} \cdot \begin{bmatrix} \alpha_{t1} \\ \alpha_{t2} \\ \vdots \\ \alpha_{tn} \end{bmatrix} = h_1 \alpha_{t1} + h_2 \alpha_{t2} + \dots + h_n \alpha_{tn}$$

② concat [o/p sentence + context-vector] and pass on

the current LSTM decoder as i/p.

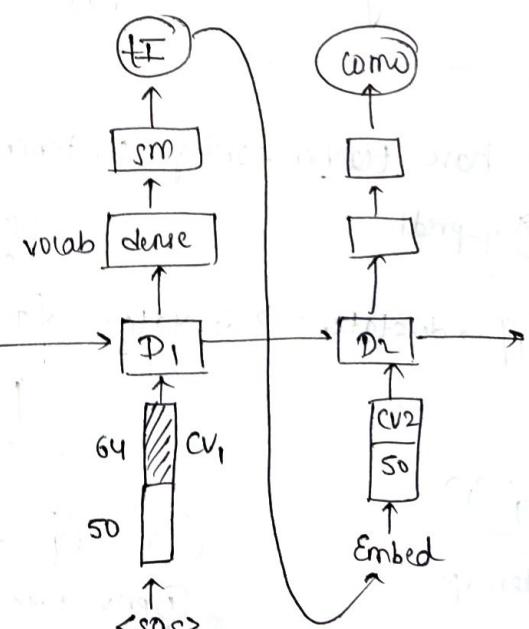
Ex How are you → ti como estas



$$CV_1 = e_0 \alpha_0 + e_1 \alpha_1 + e_2 \alpha_2$$

$$CV_2 = e_0 \alpha_0 + e_1 \alpha_1 + e_2 \alpha_2$$

$$\alpha_i = \text{dot}(\text{decoder}_t, E)$$



How to find α_{it}

Encoder

ip $\rightarrow (\text{None}, 600, 50) \rightarrow (\text{None}, 600, 64)$

Decoder

$LSTM = 64$

at time t

$(600, 64) \rightarrow \text{all hidden states} \rightarrow h_e$
 $(64, 1) \rightarrow \text{last h-s}$
 $(64, 1) \rightarrow \text{last c-s}$

val from prev decoder step $= (h_{t-1}) (64, 1) \rightarrow \text{what if decoder LSTM (128)?}$

now

for first we last
units of encoder

$$h_e \cdot h_{t-1} = (600, 64) \cdot (64, 1)$$

$$= (600, 1)$$

\therefore you get 600 α $\boxed{\alpha_0 - \alpha_{600}}$

\rightarrow pass it to softmax so that $\sum \alpha_i = 1$

$$\boxed{\sum \alpha_i = 1}$$

* use a matrix W of shape
(enc-units \times decoder-units)

$$h_e(k) \cdot h_{t-1} = (600, 64) (64, 128) (128, 1)$$

$$= 600, 1$$

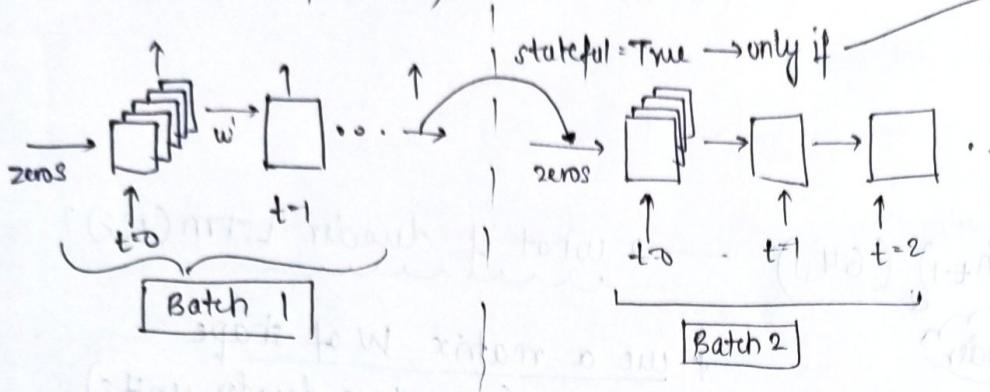
Learned by backprop.

char-RNN concepts

Stateful LSTM

→ used when batches are sequential

→ `LSTM(5, stateful=True)`



Batches must be sequential

1 2 3 4	5 6 7 8
9, 10, 11, 12, 13	14, 15, 16, 17
:	:

Batch 1 Batch 2

Here stateful can be used

* while training batch 1, initially weights are zeros and are gradually updated
They are discarded at last step

+ when batch 2 is trained again initial weights are zeros

↳ Instead we can use the weights which are obtained from last unit of batch 1 ← stateful

[use] if dataset is proper,

we can learn longer RNN sequences

Ex: given ex dataset

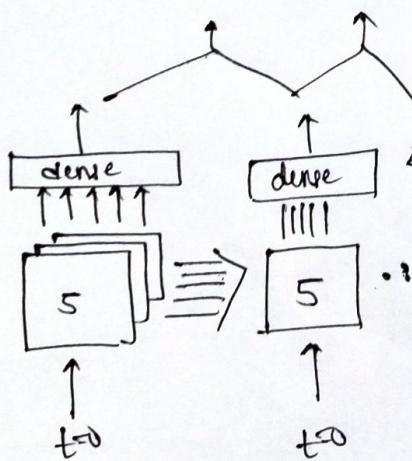
we would train sentence/sequence
of len=4 (0-4, 5-8)

* with stateful we can learn longer
ie (0, 8)
(9, 17)

TimeDistributed Layer

→ apply a dense layer on top of output of each LSTM timestep

* `LSTM(5, ret-seq=True)`



Time distributed dense layer
Code → Time Distributed (Dense(vocab-size))

Ex: `LSTM(64, ret-seq=True)`

$(\text{Batch}, 18) \rightarrow (\text{batch}, \text{time_steps}, 64)$

$(\text{Batch}, 18) \xrightarrow{\text{LSTM}} (\text{batch}, \text{ts}, 64)$

↓ Timedist (dense(16))

$(\text{batch}, \text{ts}, 16)$

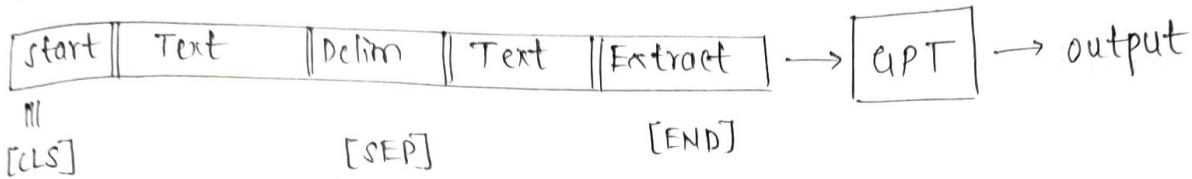
Generative Pre Training

GPT-1

Architecture: ~~Decoder only transformer~~
 ↳ 12 decoder cells (pg 56)

Training: Next word prediction
 → give word (w_i) predict (w_{i+1})

Tokens



GPT-2

Improvements on GPT-1

* max input length = 1024 (GPT1 - 512)

* vocab size = 50,257 (GPT1 - 40K)

* Layer Norm is used @ input of each sub block

* larger batch size : 512 (GPT1 - 64)

* token embedding size = $d_{model} = 768$

* used 10x more data to train

+

Had 10x more parameters

* Zero-shot learning

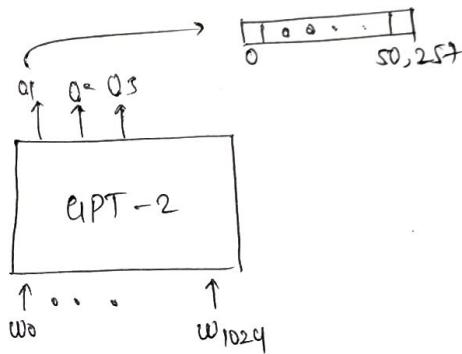
① Pretrain model (NWP)

② fine tune model by giving your dataset

③ make predictions

few-shot learning

① → ② → ③



Params	layer	d_{model}	
117M	12	768	→ GPT2 small
345M	24	1024	
762M	36	1280	
1.542B	48	1600	→ GPT2 large

disadvantages

performs poorly on rare or specialized content

64

GPT-3

→ 100x more params (175B)

→ trained on

$\begin{bmatrix} 300B & \text{common crawl data} \\ \text{text tokens} & (2016-2019) \end{bmatrix}$

→ uses sparse attention

disadvantages

- Cost and compute power
- Interpretability
- Biases in predictions