

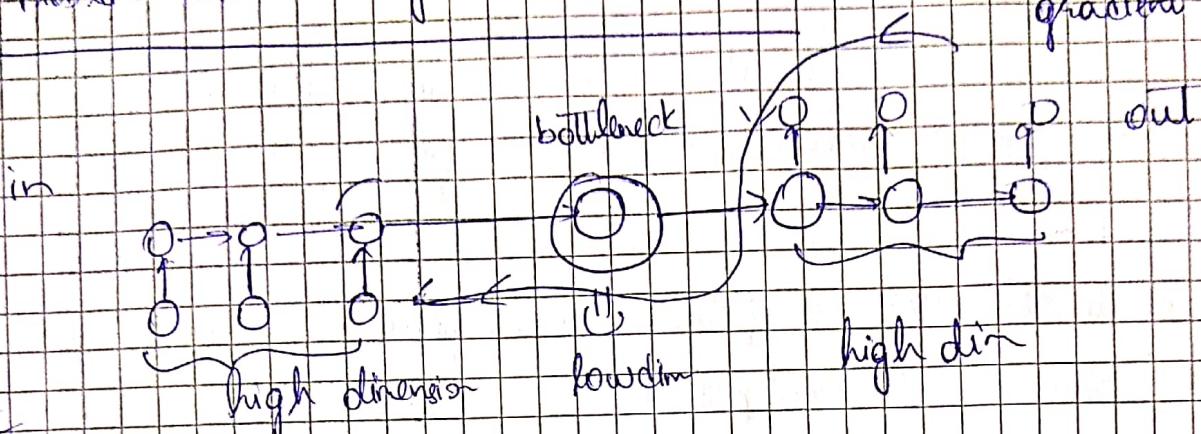
7/06/20

3.1.1

Stable

Attention - 1

Problems with existing encoder-decoder (Need for attention)



* Difficult to capture long term dependencies of long sentences

as everything is mapped to one 512/6124 low dim vector in bottleneck

→ Both long & short sentences have same no of nodes at bottleneck

* Gradient should pass from o/p \rightarrow bottleneck \rightarrow i/p hence encode
difficult to learn

Soln: Why can't decoder look at bottleneck for all the time
* This also

→ This allows for decoder to focus on relevant hidden state
at that particular timestep (No bottleneck)

→ Last bottleneck is connected to all encoders with some weights (attention)
Hence allows for creation of shortcuts dynamically making
gradient update more localized to some encoders

For example: The quick brown fox jumped over a lazy dog

* P_1 P_2 P_3 P_4
(shorter dependency)

→ When decoding P_1 part of sentence P_1 must have higher gradient update & maybe some for P_2 , but P_3/P_4 should receive less prominence

order of prominence

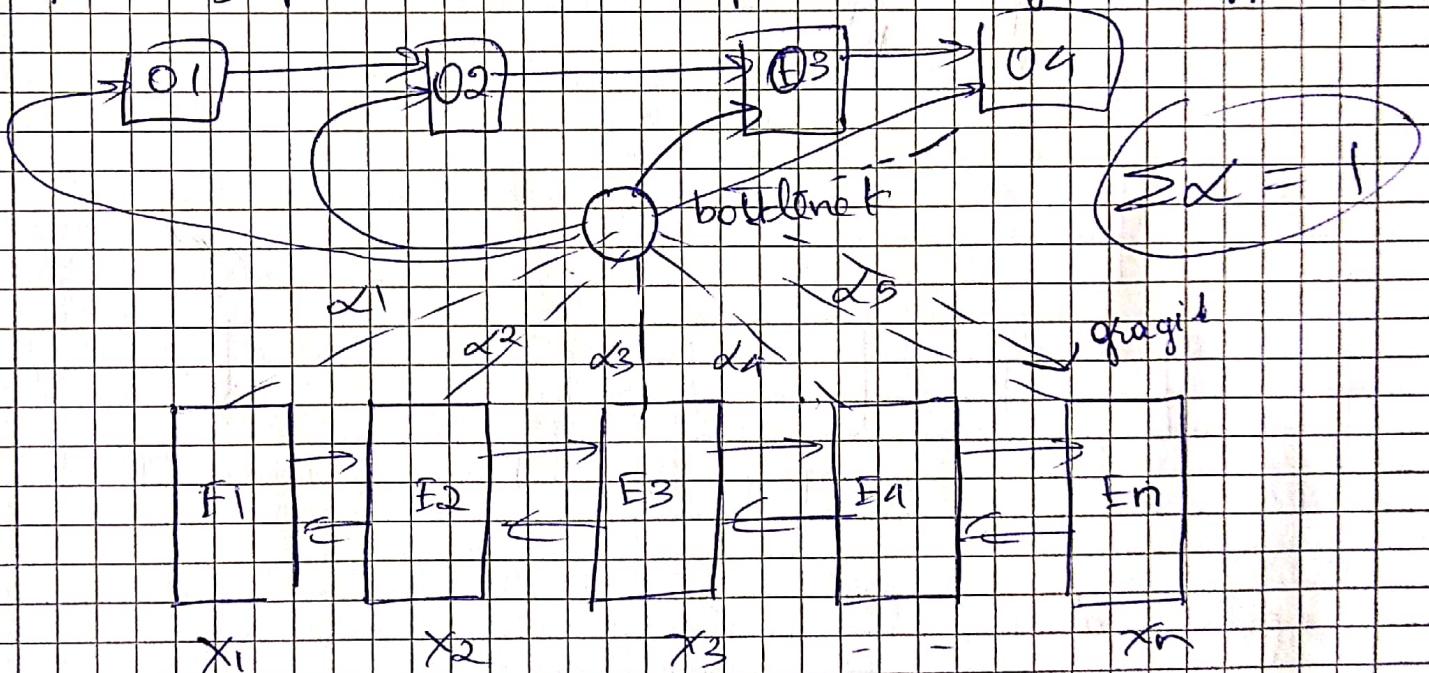
Why for decoding

$$\left\{ \begin{array}{l} P_2 \rightarrow P_2, P_3, P_1, P_4 \\ P_3 \rightarrow P_3, P_2, P_4, P_1 \\ P_4 \rightarrow P_4, P_3, P_2, P_1 \end{array} \right.$$

(Use α)

→ This is dynamic hence allows for better gradient update

→ Encoder can be bidirectional To receive info both from past & future ' α ' is for this dynamic effect



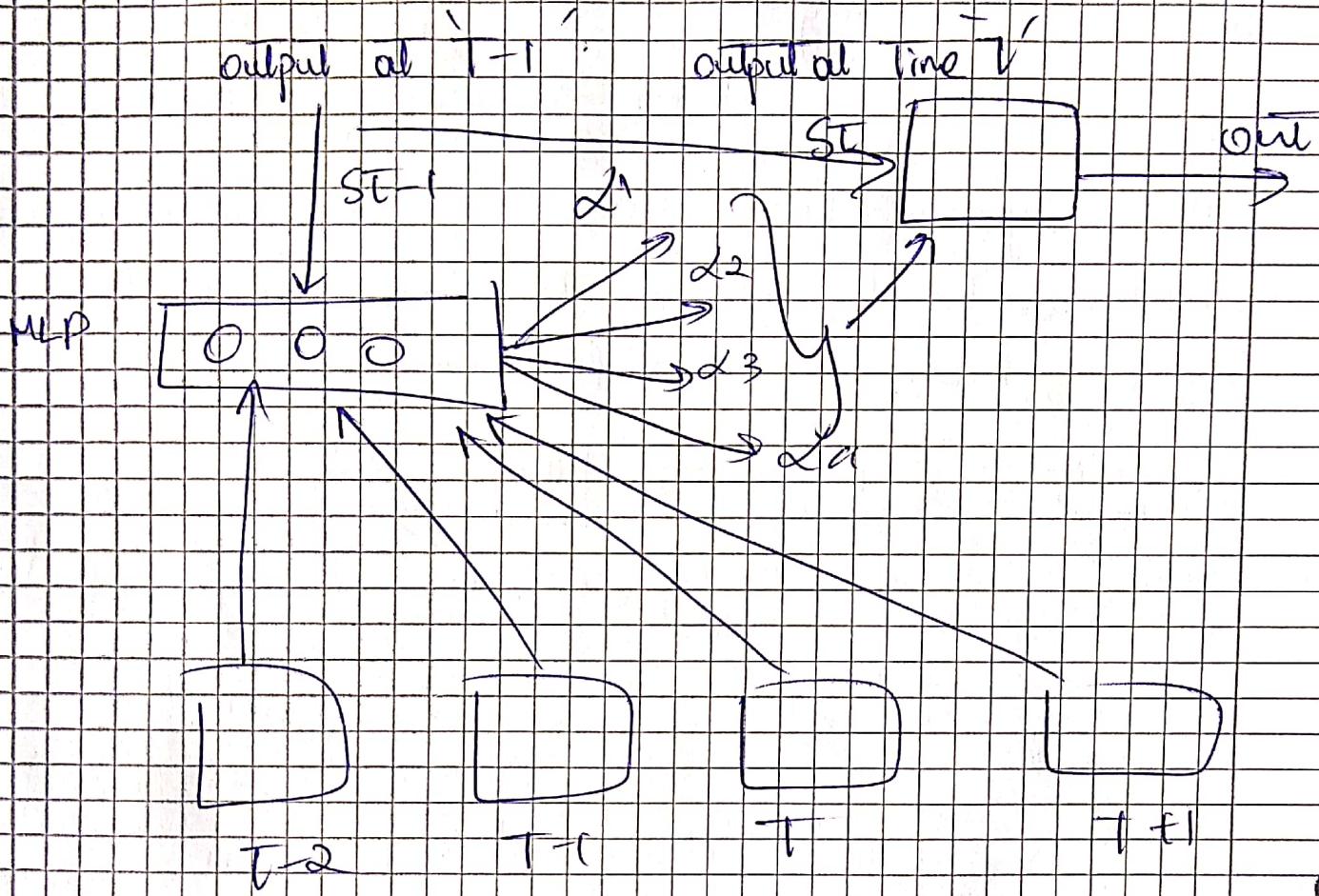
→ At each Time step of decoding we get different weights α , hence gradients flow/update are dynamic

→ At each step decoder consumes a different combination of encoder states. Glimpse of all parts of sentence/sequence available

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

$$e_{ij} = a(s_{i-1}, h_j)$$

$$a(s_{i-1}, h_j) = \underbrace{v_a^T}_{\text{Scale}} \tanh(w_a s_{i-1} + v_a h_j) \underbrace{y_{MLP}}_{\text{Weight}}$$



Computational aspects

- Quadratic $O(n^2)$ computation if we have 'm' decoder steps & 'n' input then $O(mn)$ operations
- for each pair of i, j $\xrightarrow{\text{o/p}} \xrightarrow{\text{i/p}}$
- * * Sum Two vectors \rightarrow apply Tanh \rightarrow dot product
- (can be done in $1/\text{el}$) for all inputs j given i is fixed
(o/p is fixed time)
- Due to vectorization, $O(mn)$ could be computationally fast

* Important in context function

Attention is not just simple MLP layer mostly defines a sort of similarity

Dot Product Attention

- To give attention to current time step decoder, we take output of previous time step & take dot product with every encoder state $\xrightarrow{\text{MLP}}$
- Hence there are no parameters / to be trained. It's just a simple series of dot product (similarity)
- Think this in context of Info retrieval. Which decoder state has more similarity to which encoder state
For example from previous page: P4 is less similar to P1 than P3

→ Downside: Since there are no weights involved both decoder - encoder dimension matrix must be matrix multiplication compatible.

$$W_{\text{Dec}} S_{\text{dec}} + W_{\text{enc}} H_{\text{enc}} \quad \times \quad (\text{additive})$$

Alignment matrix

| Sentence - english sentence
| sentence \rightarrow german sentence

Location based attention

Neural Turing M/C

Why To see all encoders (global)?
Check some defined neighbourhood
based on previous o/p. Use This
To calculate next step

(content) based attention

Dot product
additive

Looks at neighbouring regions for some localised
region based (global) *

Hard attention (Visual domain)

→ Consider human vision. Human vision / receptive field even though is large, can only focus on couple of objects at a time for rest of them would be blurred

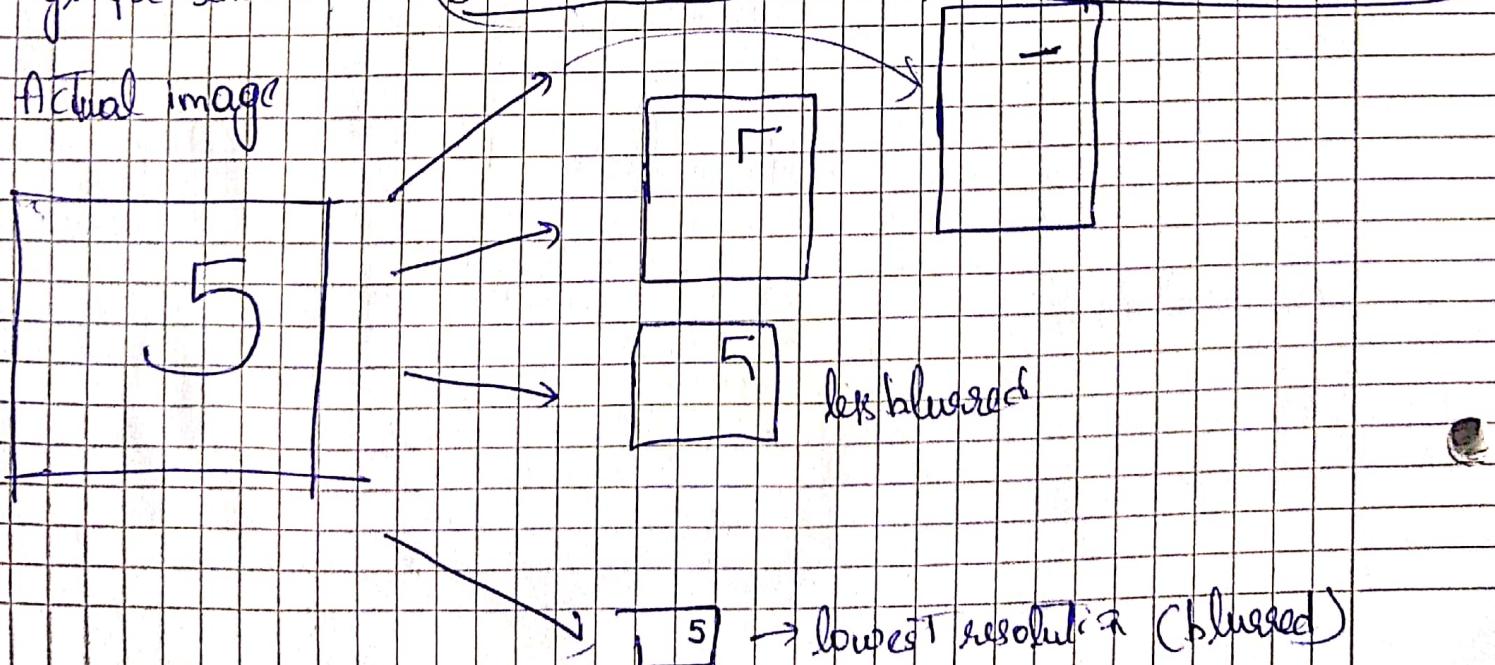
* Hence humans use a hard attention mechanism, To see other objects we need to turn our eyes else not

How To do it?

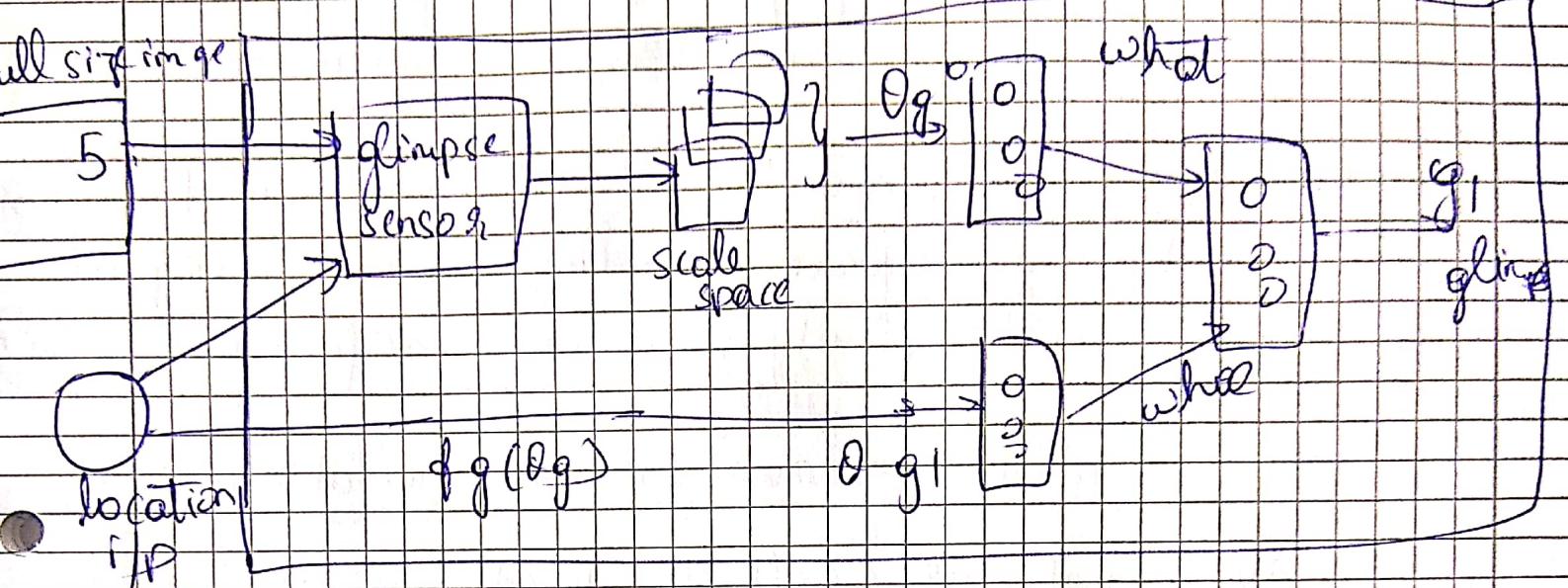
Recurrent attention model (RAM):

glimpse sensor - Builds a scale space of three-4 levels

Actual image



Concatenate scale space like in dense-nell & pass through n/f

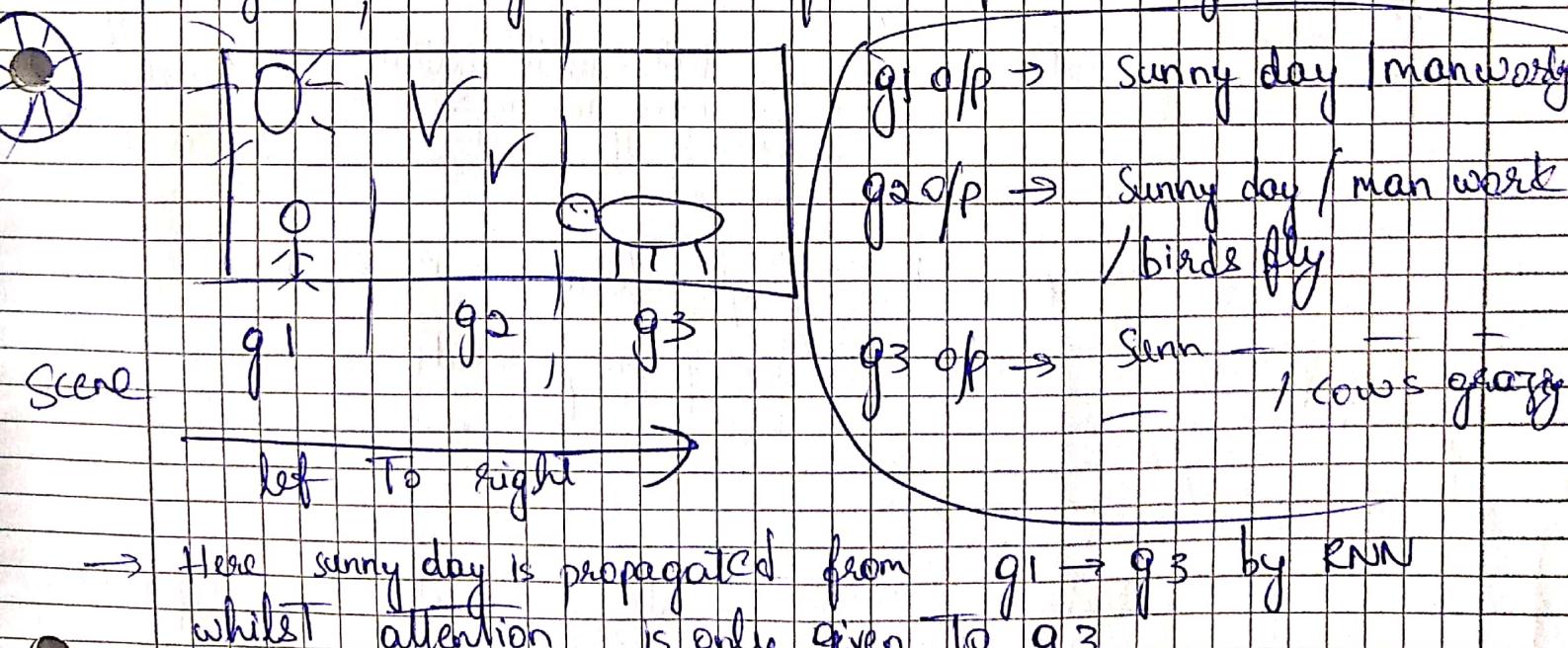


Activation n/w : what do you want to see / on seeing now

Location n/w : Predict where should we look next step

Recurrent n/w : Propagate state fwd

Imagine watching a scene from left To right



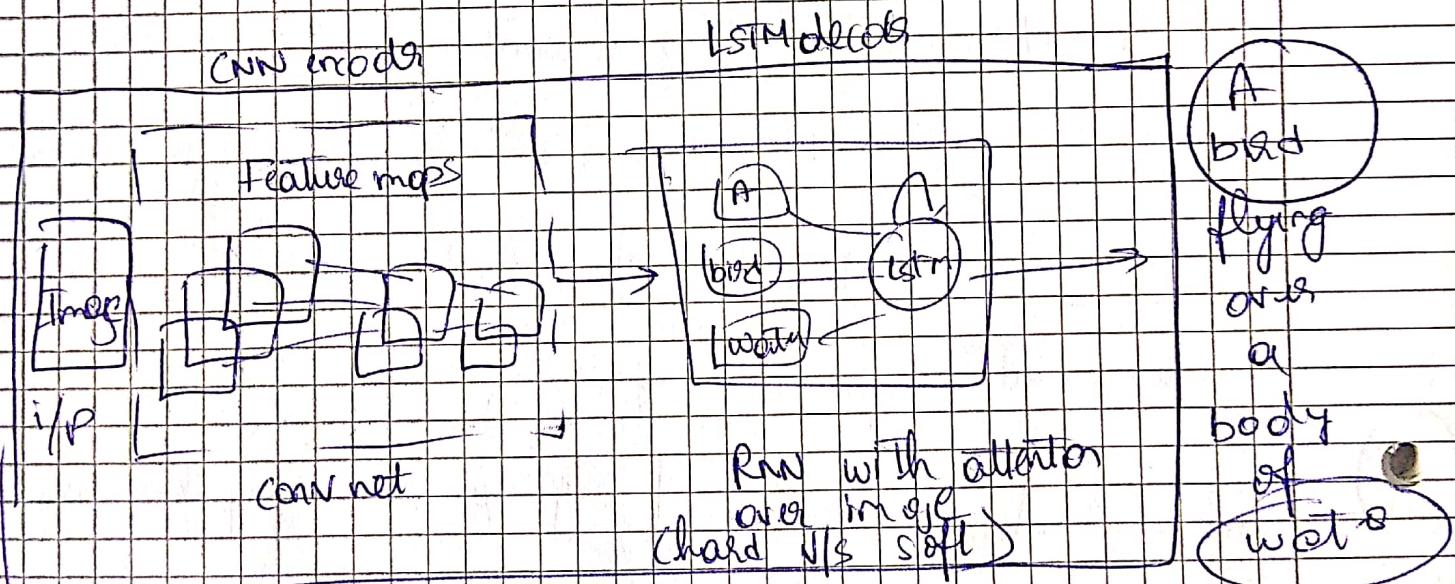
→ Here sunny day is propagated from $g_1 \rightarrow g_3$ by RNN whilst attention is only given to g_3

* Problems : Scale space is gaussian sampling & random, hence cannot back propagate (sampling \rightarrow random \rightarrow non deterministic) easily. Hence train it using reinforcement learning. gaussian filter changes

* Deep RAM : Provide a content vector. To provide a good starting point. Just using reinforcement learning ends up starting at random state. $O(1)$ complexity or $O(k)$ as focus on only one state

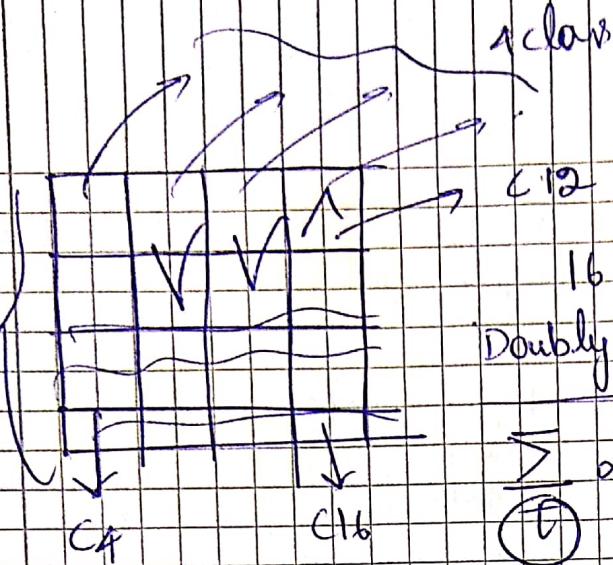
* CTC is a hard attention mechanism but with tractable gradient

* Show, Attend and Tell (Soft v/s Hard visual attention comparison)



A bird \rightarrow high focus \rightarrow more output
water \rightarrow low focus / resolution \rightarrow more blur less attention

* With soft attention they create a grid over image and each grid cell becomes a softmax over attention



16 class softmax for soft attention
Doubly stochastic regularization

$$\sum_{i=1}^n \alpha_i = 1 \approx 1 \quad \text{Index now is } y$$

* * At each Time step consider α we need to pay attention to all $C_1 - C_{16}$ gridcells equally

Aggregated over Time for single grid cell

$$\text{loss: } L_d = -\log(P(y|a)) + \lambda \left(\sum_{\text{Location}} (\hat{P}_c - \sum_i \alpha_i) \right)$$

Typical crossentropy

sum over location

all

for each location

compute deviation

from \hat{P} across Timesteps

Memory based n/w

Von neumann architecture : Separate memory & program

→ Neural n/w Takes care of computation

→ RAM Takes care of memory

- * * Can operate on any knowledge base, KB can be swapped in/out
- Doesn't need neural n/w To represent data, (DBpedia)
- We can scale up/down /swap KB

concept

RNN, MLP



- Neural n/w acts as controller, which controls where & when to read on memory based on external/internal input
- Heads select positions of memory to read or write on them
- Memory is realvalued matrix & not binary *
- Everything differentiable, which means we cannot make (discrete) hard decisions, only soft based
- access only one memory cell, not different one

Memory networks in practice

- Strong supervision needed To access a specific memory cell To address respective query
- direct / reward based
- babb dataset

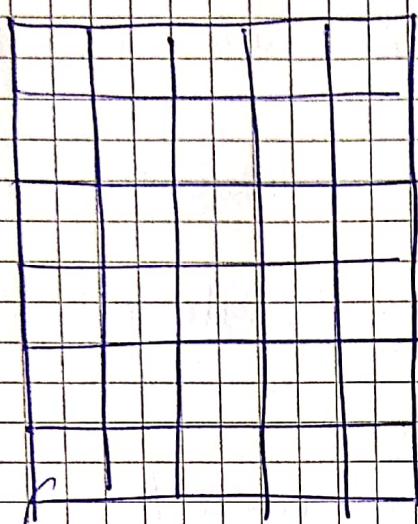
7/06/20

Neural Turing m/c

21:50

External I/O

External o/p



~~each cell n-dim vector
real value~~
memory explicit
~~N finite memory position~~

External I/O

Controllers

Read head

Write head

Memory

(Explicit memory why) we may need to access large amount of info if we need n/w to store & represent actual data/pattern in terms of weights. If its already stored in explicit memory, we can swap in/out similar KB

NTM Read/write :- Train by gradient descent, differentiable w.r.t.

Soft read (Content - location attention) :- Read all memory cells every time but with a softmax mask for ex:-

0.01	0.05	0.03
0.2	0.3	0.1
0.3	0.05	0.05

we read from every cell but focus on (21), (22), (31)

soft write

$$gt \leftarrow \sum_i w_t(i) M_t(i)$$

soft erase :- Basically we write everywhere but with again a softmax mask (attention mask)

$$M_t(i) \leftarrow M_t(i) + w_t(i) \delta_{at}$$

~~forget gate~~

Memory matrix

Soft erase:

$M_U(i)$

$M_U(i)$

$[1 - M_U(i)]$

Memory at previous time

Big α value leads to more forget

* At y Erase & Write control vectors comes from controller

* May lead to slow training as we need to read/write external memory
can be alleviated by use of GPU

* Important Questions: How to know where to read/write
How much to read/write? This is where controller comes in.
We use attention mechanism (Content-based addressing)

i) Content based (Associative Memory)

Retrieve by partial query / input — similar to search engine

ii) Interpolation

Based on previous memory cell accessed, Use previous attention
to determine next access point

iii) Location based addressing (Much more powerful than others)

→ Needed for variable (varying content) (By their nature variables hold variable values & their values are not constant)

Since vars keep on changing, we cannot query by their value we need
to address the location they are stored

234 // AB12 Content based can emulate location based if content has some part of location embedded along with content

query q_T

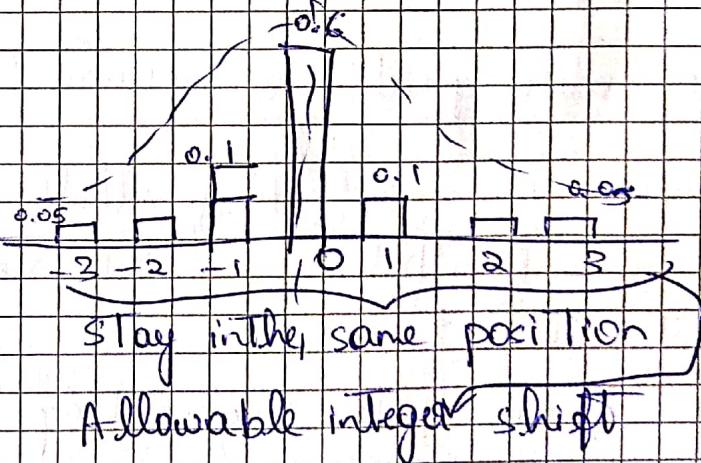
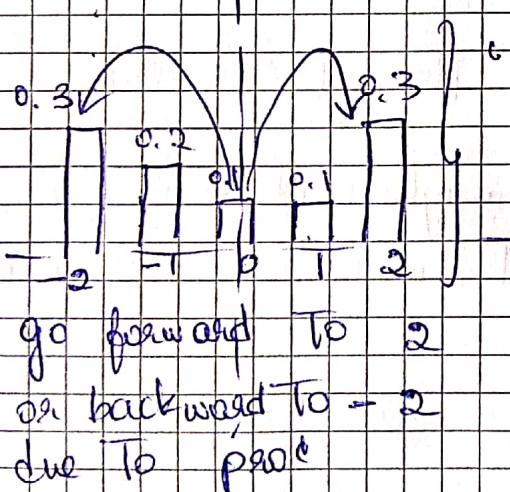
= [1 0.2 0.1 0.05 0.1 0.2]

content info

location info

- a) content based :- Use cosine similarity, key \rightarrow query, query compared To Memory cell $y \rightarrow$ apply softmax
- b) Interpolation :- Use old attention To update new attention
No need for content
- $$W_t^g = g_t W_t^c + (1 - g_t) W_{t-1}^c$$
- How much current info should we keep? \rightarrow NTM should learn this from last atten
- c) NTM location based addressing :- Once you are at the position location of memory cell, you need to read discrete number of memory cells. We can only work with continuous due to different How To deal with This?

* * get a distro over allowable integer shift. Offset from where we are now. \rightarrow How much we are gonna go further to do next read



Why we need modulo :- Consider memory size as 16, Since controller is a neural net it can give any value, like 100/500/1-20 etc which are physically not present on cell. Use modulo to convert them back

Another option is To make network learn a regression model which takes a real valued i/p & converts into distribution

$$\text{eg: } I/P \rightarrow 6.7$$

$$O/P \Rightarrow (6 \times 0.3) + (7 \times 0.7)$$

regression or

Take some real number resolve it into distribution

memory cell 6

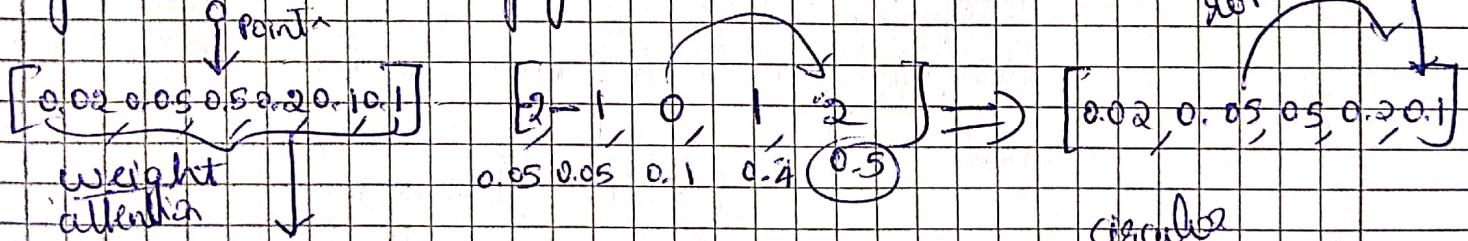
memory cell

How to shift attention based on distribution?

Use weight vector on distribution

length of weight vector is same as no of memory cell

If we want to attend to next / previous location we rotate weight vector accordingly



This module operatⁿ is implemented through convolution

(ii) sharpening: Soft fuzzy attention/weight masks lead to dispense of actual attention (Attention mask which is spiky, begins to flatten out)

→ Use softmax with higher Temp → lower dispersion

NTRM Addressing flow diagram

- We start with previous state, memory of last attention vector
- We start with off of controller (recall v_{lw}) \rightarrow query
- Step 1: Do content based addressing by using key vector and key strength To get content based attention / weight vector
- Step 2: Do interpolation using gates from controller. That tells whether we need to only pay attention To last memory cell or The new content vector / memory cell That was generated by step 1 as part of answer To query

Step 3: Once we have (gated attention vector) Use location based To determine shift parameter from current location

Step 4: Once we have (location based attention) from step 3 Do the sharpening To remedy dispersion. Use final weight vector To operate on data.

Pointer Network

Encoder - Decoder : Output are (positions / indices of input sequence)

NN can learn (sorting / permutations) problems

X Problems

* Convex hull : Network has To output a shape which goes around all points & includes others

Fix



Convex hull

(Output ordering points of hull)

\rightarrow All points contained in hull

Schedule
Delaunay Triangulation, Travelling Salesman, Graph Traversal,
are problems for point-to-point

Differential Neural Computer (DNC) (NTM 2.0)

Changes / Enhancements To NTM

- a) Added Memory management :- Helps when scaling up, Machine doesn't waste memory, if memory becomes free it should be reused for writing again
- b) No Temporal link matrix :- Allows To Track how memory was written from previous Time steps (Are there some memory that are read/written again and again?)
 - Facilitate sequential access to data not written to contiguous block
 - Facilitate retracing steps through Time
 - ↓ Use address pointers To allow writing to non-contiguous block (fragmentation problem)
 - It has a free list To know where are free memory location To work new info
 - Gates keep track and control how data is accessed

DNC Searching by Time :- Change attention based on the architecture

One write, 2 read heads for controller
read multiple things

Attention Transformer notes

vanilla

- Why Attention? Disadvantages of LSTM

→ In vanilla RNNs the output of the encoder is carried out till the last encoder which acts as bottleneck! Obviously for long data there are chances that the last encoder could have forgot about data at beginning.

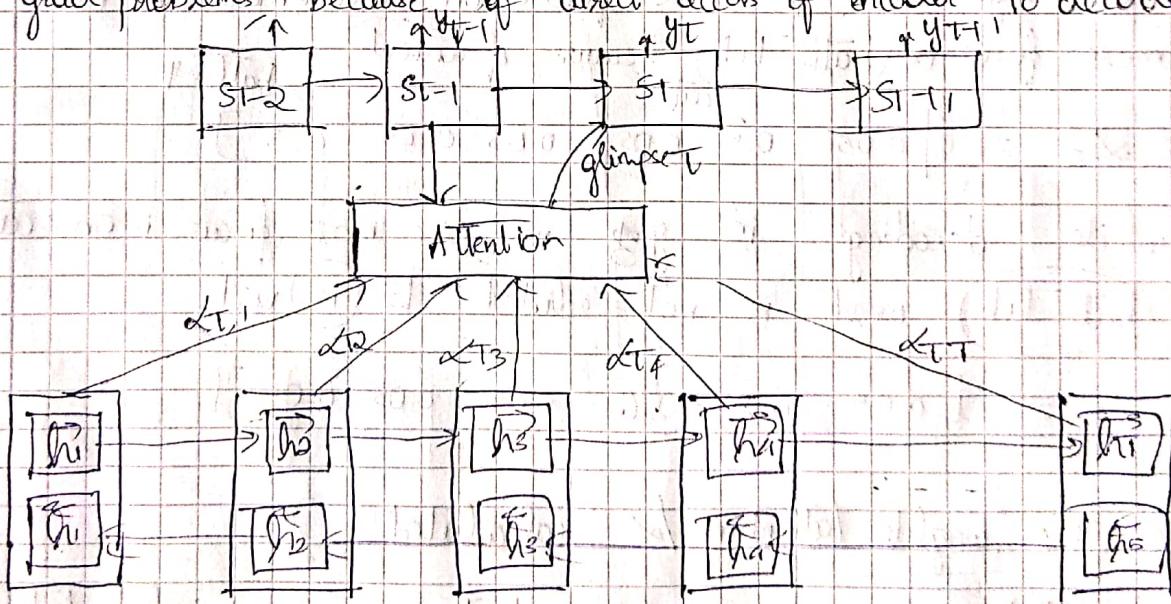
→ Gradients should flow from last encoder to first causing problems in Training (Vanishing grad)

Extra attention :- here all encoder states are fully connected to decoder state.

- state. Each encoder consists a kind of summary of entire sequence with focus on some relevant parts. Hence

→ No bottleneck, because of direct connection to decoder state & summary at

→ No grad problems, because of direct access of encoder to decode



→ At each Time step decoder consumes a different summary with different focus from encoder state. Hence it has glimpse over entire sequence (glimpse vector)

How to interpret / calculate output at Time 't' (y_t) given input X ?

W.K.T, y_t is dependent on the current encoder state s_t and glimpse / context vector at $t - G$

⇒ probability of predicting y_t given its past output & inputs are a func of decoder state & glimpse

$$P(y_t | y_1, y_2, \dots, y_{t-1}, x) = g(y_t, s_t, c_t)$$

How to calculate c_t ? (glimpse vector)

→ Intuitively, glimpse will take a weighted input from all encoder states in parallel. Hence

$$c_t = \sum_{j=1}^J \alpha_{tj} h_j$$

↓
weighted sum

so for every timestep of decoder there is a different weight / focus on encoder state $s_1, s_2, s_3, s_4, s_5, s_6$

Ex: The cat is in the bath tub.

$$\begin{matrix} & \uparrow & \uparrow & \uparrow & \uparrow & \uparrow \\ s_1 & s_2 & s_3 & s_4 & s_5 & s_6 \end{matrix}$$

→ while decoding cat (y_2) it focuses on cat more than others. It also gives more focus on bath tub because of context.

$$\alpha_2 = [0.05, 0.5, 0.05, 0.05, 0.05, 0.3]$$

→ While decoding is y_3 we see more focus is on cat & is because bath tub word is not related to is verb.

$$\alpha_3 = [0.05, 0.2, 0.6, 0.05, 0.05, 0.05]$$

How are the weights / attention α calculated?

- 1) Additive attention
- 2) Dot product
- 3) Scaled dot product

→ These types of attention, each calculated differently & have their own advantage & dis-

Additive attention: An MLP model would take in previous decoder states and one encoder state & calculate attention for this pair.

$$e_{ij} = \text{Va}^T \text{Tanh}(\text{Wa}_{\text{scale}} s_{i-1} + \text{Va}_{\text{scale}} h_j)$$

decoder Time step encoder Timestep scale previous decoder state for every encoder state

→ fix $i = T-1$ To calculate off off

→ iterate through all j 's (encoder state).

→ increment i To calculate attention for next decoder state

* Problems

→ Need to learn weights of MLP which is an added burden

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_j \exp(e_{ij})}$$

Attention → $O(n^2)$ operation is fast as done on gpu, but still quadratic

Dot product attention :- Instead of MLP take a dot product to find cosine similarity. $e_{ij} = s_{i-1} \cdot h_j$

Advantages :- No weights like that of additive, hence simple

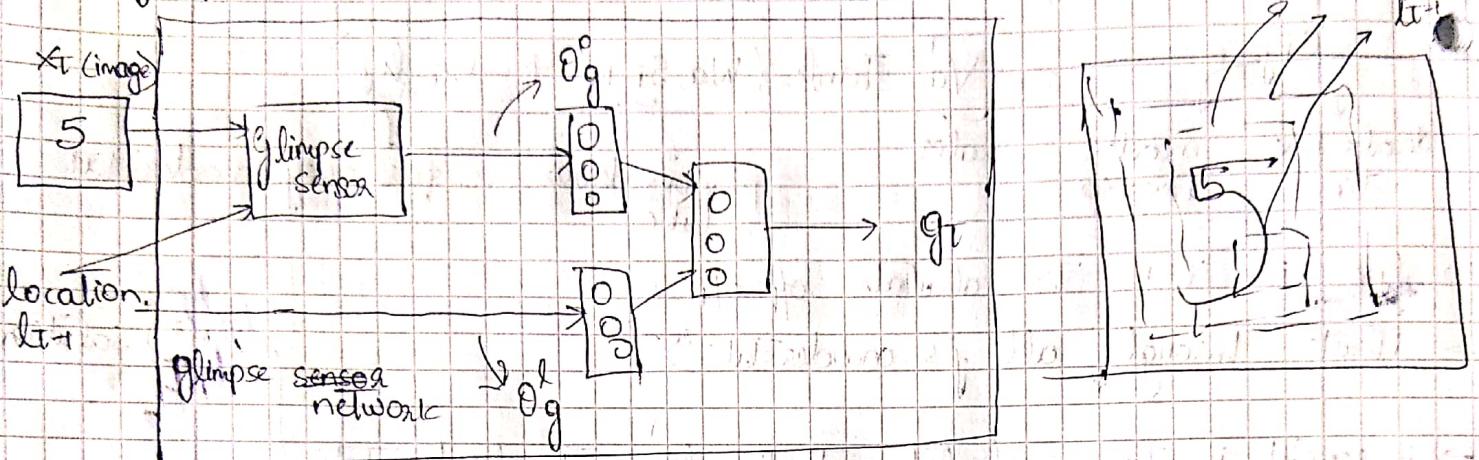
Problems :- s_{i-1} & h_j should have compatible dimensions as we are directly multiplying them. We cannot adjust weight dims using weights like we did in additive

DRAM

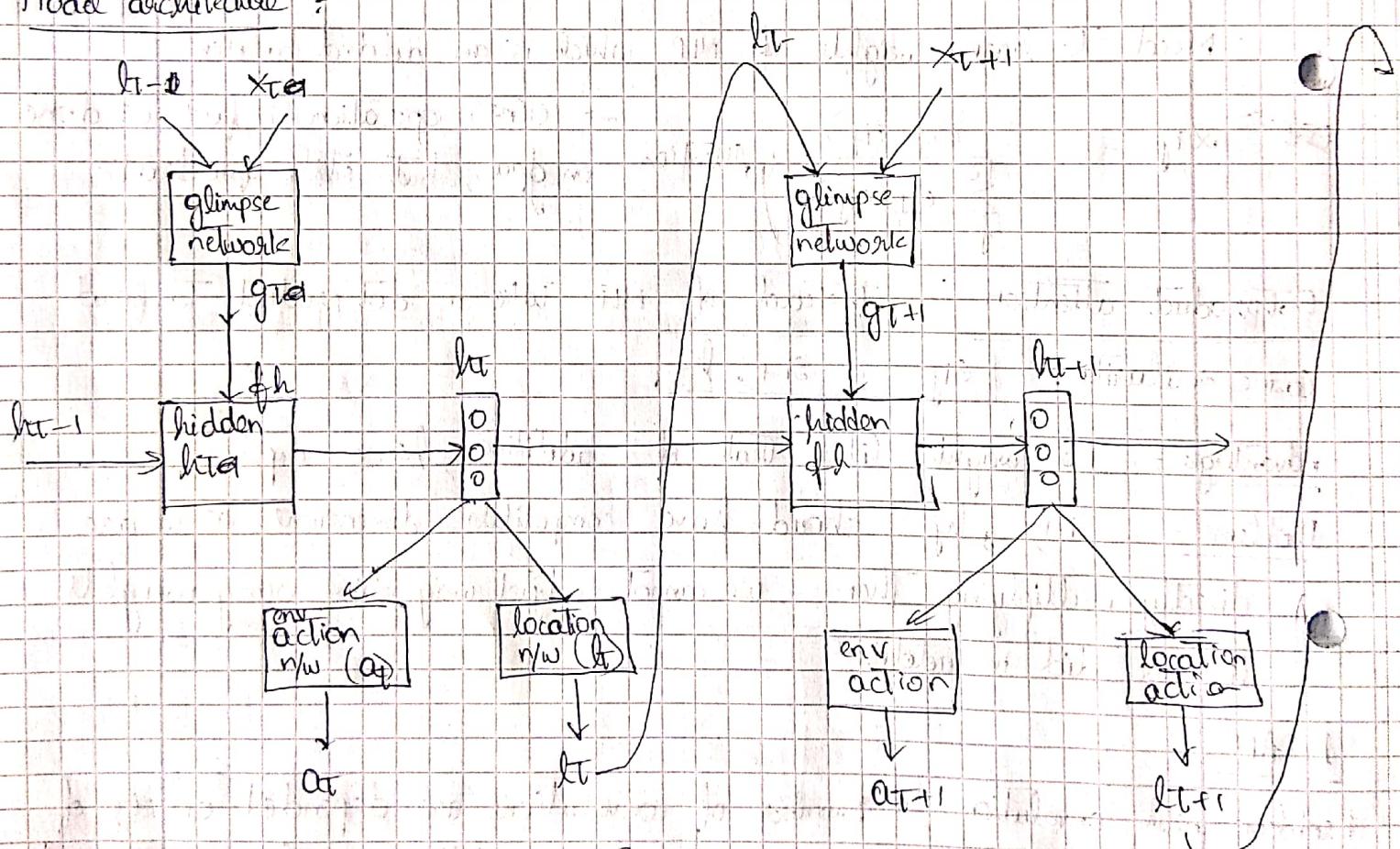
Problems with convolution :- Number of convolutions are depended on size of image. Or atleast even if we increase filter kernel size to make some amount of conv, increasing is still a function of image size.

Solution :- See an image like human does, humans take some initial context about image and only focus on required part of image, dull the rest. This dulling causes less info need to classify / action on any image.

glimpse DREAM paper : Taking inspiration from human vision we can build a glimpse sensor that imitates them.



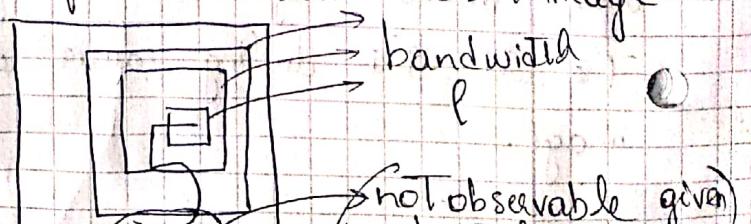
Model architecture :



How does glimpse network work?

→ At each Time step t given an image x_t glimpse sensor builds up a scale space focussing on location l_{t-1} . Hence we observe same image at different scales, but we won't have full bandwidth access to image (States are partially observable)

* Can learn Translation invariance like CNN, but will take more glimpses



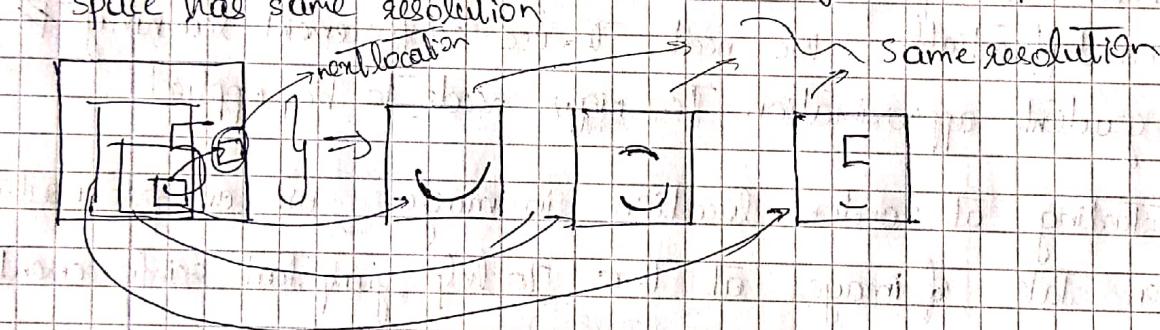
→ This function $f(x_t | t-1)$ helps encode gait-like representation in vicinity of $t-1$ and dim other places. Hence use high resolution at pixels near $t-1$ & low resolution elsewhere.

→ Both location $t-1$ & bandwidth are combined together by bunch of Neural net layers $\{\theta_1, \theta_2 \rightarrow \theta_g\}$ & final output layer forms g_t a glimpse vector having both location info & info of image at different scales.

How does the model work?

i) Sensor Network :- Takes x_t and h_{t-1} as input. Apply bandwidth function at location specified

→ Builds a scale space representation of image each representation of scale space has same resolution



→ Use the location h_{t-1} and linear layers to provide glimpse vector g_t which encodes high res near $t-1$ and low res elsewhere.

ii) Hidden State :- RNN hidden state comprises of info about past observations, if combined with current glimpse g_t , it can encode what are the actions need to be taken next $h_t = (g_t, h_{t-1})$

iii) Action networks :- Two Types i) location decider ii) environment action decider

a) Environment action :- $a_t = f_a(g_t)$

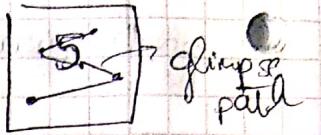
→ Decides what action to take next (classification, moving left/right etc.) based on hidden state

b) location action :- $l_t = f_l(h_t)$

→ Decides which particular pixel/location to go next / focus on next based on hidden state

iv) Reward: If agent's action / classification leads to success then it is rewarded with max reward. Over Time T we sum up the rewards accumulated.

$$R = \sum_{t=1}^T r_t. \text{ We try to maximize this } R$$



Why can't agent cheat and keep going to same location again and again since classification is right?

Soft vs hard attention and problems with DRAM

→ DRAM's attention is specific to a location, hence discrete. Discrete errors are non-differentiable. We need to use reinforcement learning to and some gradient approximation technique needs to be applied.

→ Instead of staring at random location in image, we can pre-load some context / hidden state of image at $t=0$ to help jumpstart reinforcement learning process in location we need.

↳ vicinity

Soft

Deterministic

Exact gradient

$O(\text{input size})$ → all context considered

Typically easy to train

Hard

Stochastic

Approx gradient

$O(1)$ → single or constant num of context/loc
hard to train

Show Attend and Tell

Image captioning is a difficult problem because we need to understand relevant parts of the image and relationship between them. After that this should be translated to text describing objects and relationship.

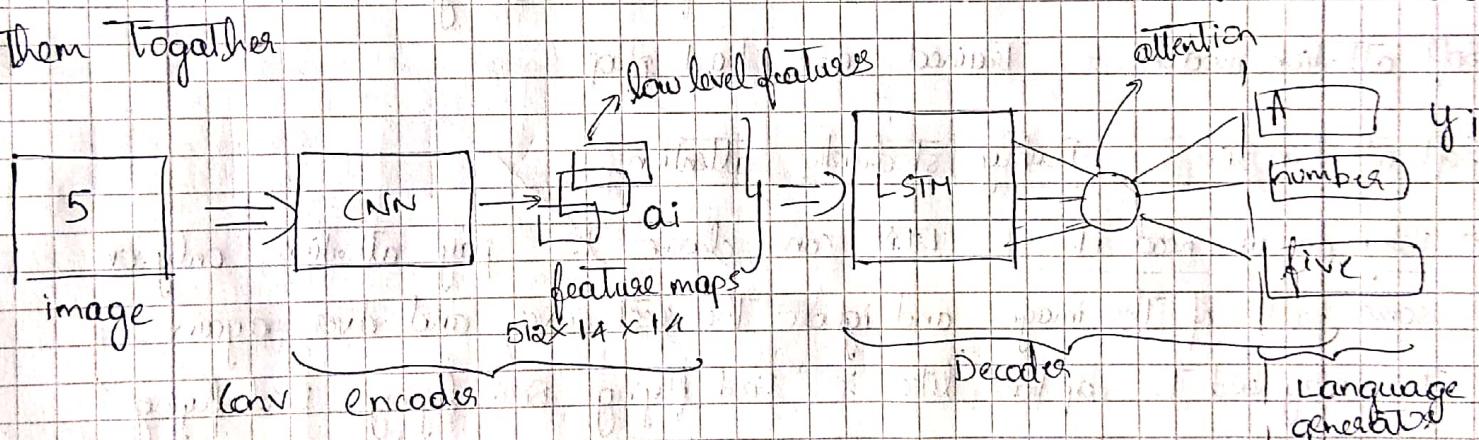
Problems with older approaches

- Training Conv net and Language model separately & combining them
- Throwing away lot of low level representations that are important to caption model and taking one whole image at each step

Current approach : Both Convnet which produce low level features and LSTM

attention decoder which utilizes these low level features to produce captions

Hence model learns abstract concepts rather than detecting objects & stringing them together



Implementation : Both hard and soft versions of attentions were used

i) CNN Encoder :

- Image is encoded into a lower dim transformed to lower dim feature maps. Each map has some part of image represented
- Each map is straightened out and put in a 2-D embedding matrix of 512×196
- Which represents a vector (a_i)
- This allows decoder to selectively use only relevant parts of image from this matrix / feature space.

ii) Decoder + LSTM :

- Decoder tries to predict next word conditioned on previous words generated, previous decodes hidden state (h_{t-1}) and context vector (z_t)

- Context vector is a function on Embedding matrix and attention
- $z_i = \phi[a_i, \alpha_i]$. This allows vector to select suitable subset of low dim features and weigh them accordingly.
- Ex:- Bird flying over water \Rightarrow (Water, over) is not important hence any subset of (α_i) having this image features are weighted less. Since we are focusing on bird, Bird, flying related lower dim representations, are weighted higher.
- Attention is calculated using additive type of model (Bahdanu) conditioned on lower representation a_i and previous hidden state h_{t-1}
- attention $\alpha = f_{nc}(a_i, h_{t-1})$. This can be thought as probability that location i is the place to focus for producing next word.

Training :-

Hard attention based models :- Trained on Reinforcement learning

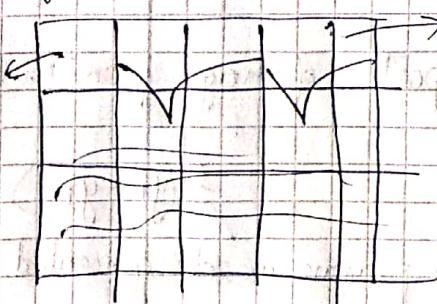
Soft attention models :- Trained on back prop

Soft attention loss :- Doubly stochastic attention

→ Why do we need it? LSTM can choose to pay attention only on some parts of the image and ignore the rest over and over again.

This may lead to captions like :- "Bird flying Bird flying. Bird flying" because it focuses only on some aspects of image.

→ Solution:- Force LSTM to cover every part of image with great focus equally over time ($\sum_i \alpha_{ti} \approx \mathcal{C} \approx 1$) Over time the attention at location i should add up to \mathcal{C} , if $\mathcal{C} = 2$ Then we will provide great attention twice over all parts of image



① location

$$\text{loss} : L_d = -\underbrace{\log[p(y|a)]}_{\text{cross entropy}} + \lambda \sum_i \left(\mathcal{C} - \sum_t \alpha_{ti} \right)^2$$

scalar

force LSTM to pay
good focus on every part of image over
time

Transformers and Self attention

LSTM with self attention

Intra v/s Inter attention

Self attention := Encoding current Token by considering other relevant token in the sequence

→ Correlation b/wn Tokens in The sequence

→ Attending To self

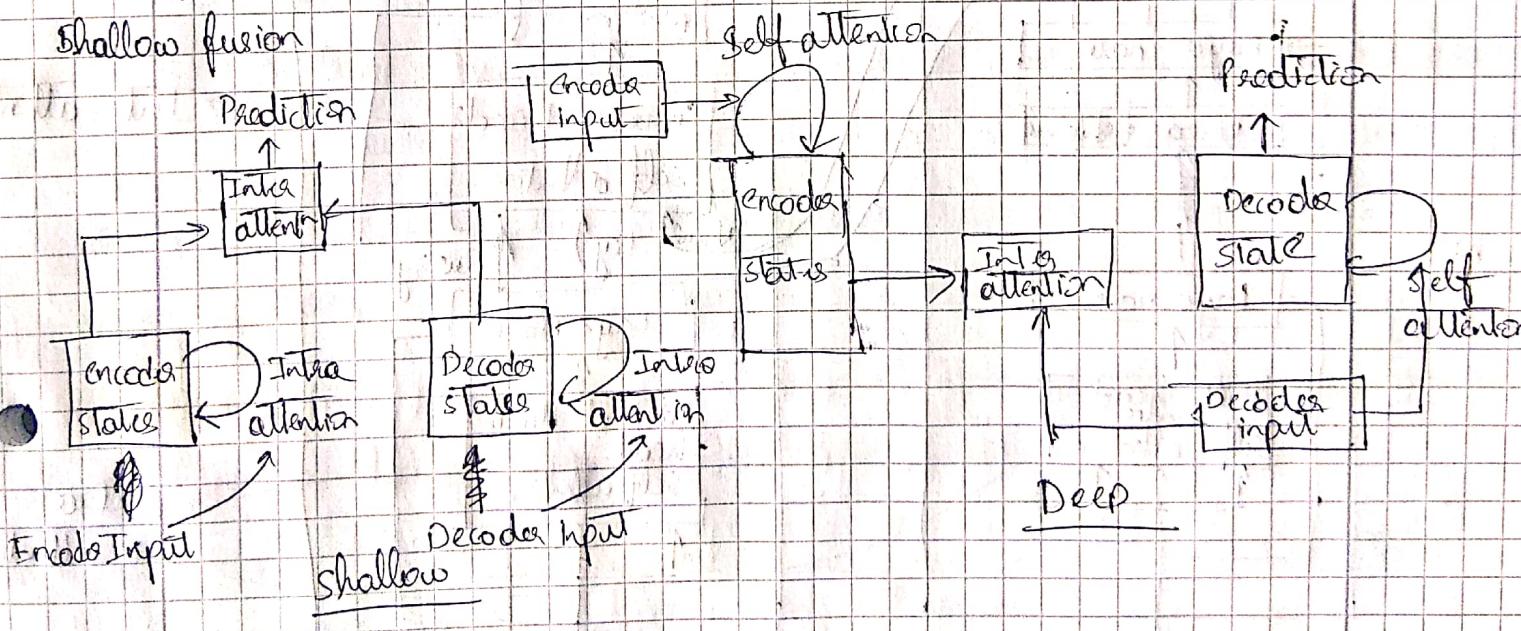
Inter attention := Decoder attends To encoder state

Shallow fusion := Fusion happens only when we do prediction

Deep fusion := Fusion happens at every stage irrespective of prediction, Hence encoder-decoder Tightly bound

→ Intra/Self attention is present in each encoder and decoder part

Shallow fusion



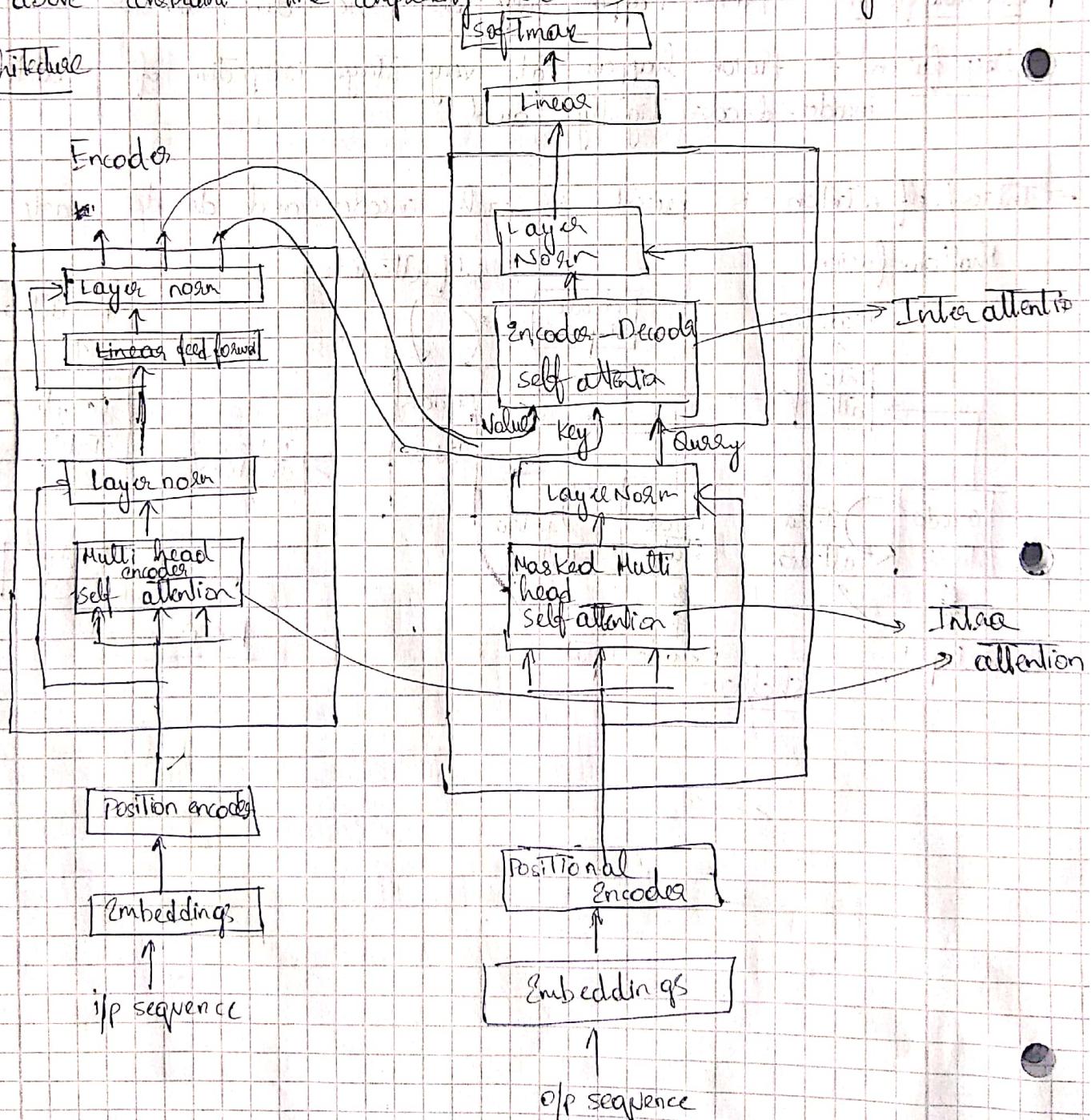
Transformer: can do sequence-sequence modelling w/o recurrence or convolution.

Process sequence parallelly

Need for new architecture: Even with advent of attention and skip connections of gradient paths, the processing still has to be done in sequence as all encoder states must be populated by respective positions of sequence, then only we can parallelly apply attention over them.

- hidden states of encoders are often huge with dimensions of 512/1024 unit
- Usually sequence length \ll no. of hidden dimensions (512/1024). But often in case like poems we have a long sequence. This causes sequential access still, with above constraint time complexity $O(n \cdot d^2)$, seems to be huge. Enter Transformer

Model Architecture



Peeling The architecture

AT Encoder side

- High level Explanation: The input sequence is tokenized and represented as a vector using some embeddings.

→ Another vector describing position of token in the sequence is modulated with embedding. Each embedding has same dimension, $d_{embed} = 512$.

↳ positionvector

→ Each token is then converted to unique trio of (key, value, query) vectors having same dimensions ($d_{key} = d_{val} = d_{query} = 64$). This is done through a simple matrix multiplication with some learnable weight matrices W_k, W_v, W_q .

→ For each Token in encoder self attention block, we find its most similar counterpart in sequence by cosine similarity of Query Token v/s keys of sequence. We use

- This similarity score to find weighted sum of values in sequence as off. This is done ~~multi~~ with multiple blocks in parallel, hence multi headed.

→ The weighted values of each token, of each attention block is concatenated and reprojected back to same embedding dimension space $d_{embed} = 512$ by a feed forward network. These meta/intermediate embeddings serve as input for next stack encoder.

→ After repeating the process over $N_E = 6$ encoder stacks we extract final key, value and query Trio from encoders last stack.

- At decoder side: We kick off with "start" Token, during the beginning and calculate the embedding positionally modulated as usual.

→ During attention calculation at decoder end we mask the future states that are not predicted to avoid copying behaviour.

→ Query of the decoder is issued upon key / values of encoder to in the encoder-decoder intra attention blocks. Again there are multiple attention heads. Aim is to find next key, value token from input given current word's query.

→ After reprojecting the ~~to~~ weighted values of this block back to embedding space, process is repeated $N_E = 6$ Times till last decoder.

→ At the end a linear + Softmax layer is applied to predict next token which would be used as input to decoder in next cycle.

Lower level details (Pooling architecture)

Embeddings :- We cannot use words as it is, hence we need to project each word as a fixed dim vector ($d=512$) in this case to be used as input to transformer.

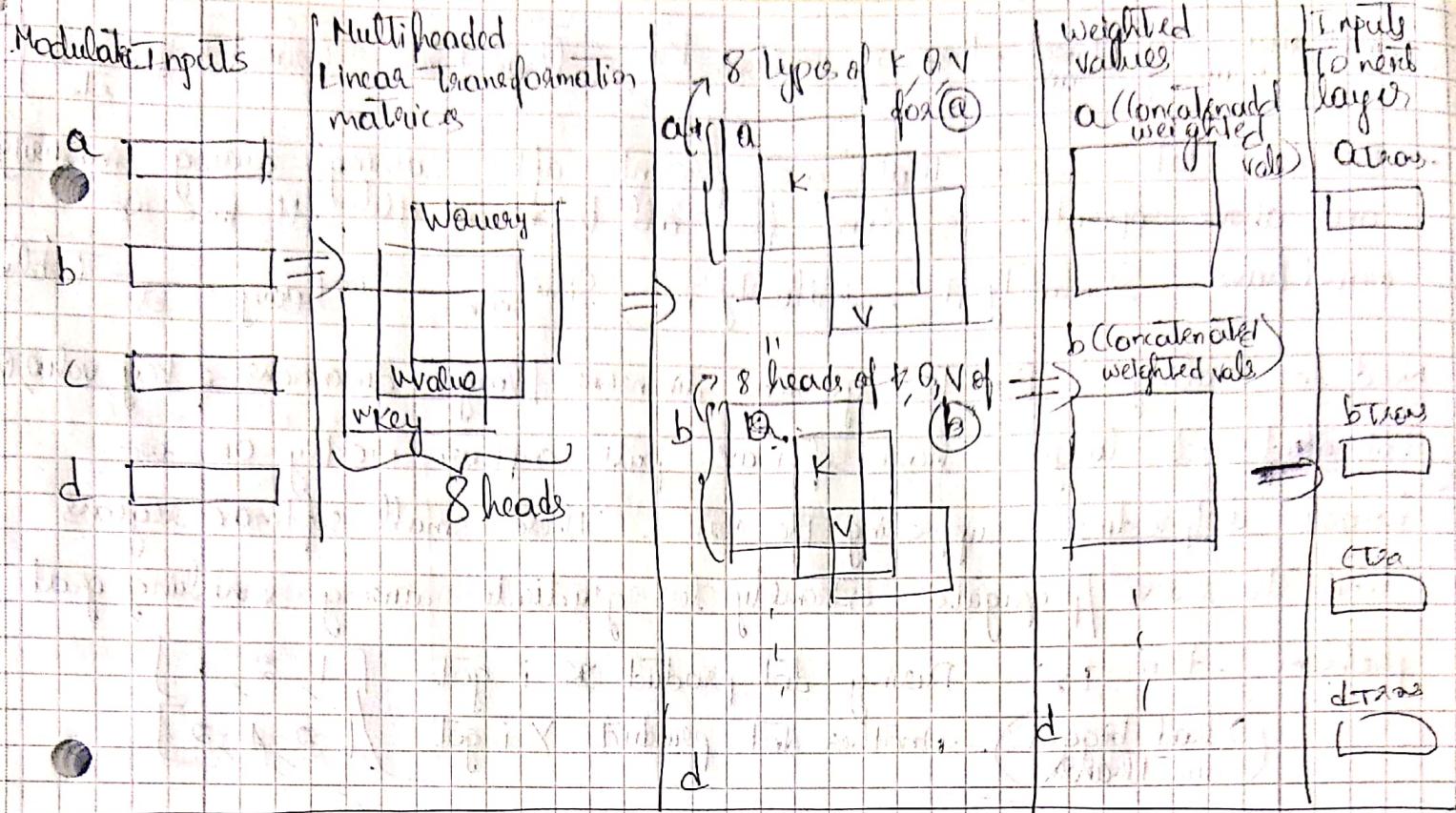
Positional Encoder :- We provide and process embeddings inside encoder & decoder block parallelly, hence we have lost the inherent order that was present during sequential processing of sequence in RNN based networks. This becomes a sort of like bag of words model. Hence we modulate a couple of sinusoidal waves which are continuous (hence can handle arbitrary sequence lengths) instead of using a fixed length positional codes like (000, 001, 010, 011 etc).

$$PE_{(pos, ai)} = \sin(pos / 10000^{ai/dmodel})$$

$$PE_{(pos, ai+1)} = \cos(pos / 10000^{ai/dmodel})$$

Encoder Self attention :- Why self attention here?

- Computational complexity per-layer = $O(n^2d)$ v/s RNN $O(nd^2)$
 → (for most cases) sequence length (n) << dimensions of hidden state or K, Q, V , hence Transformer performs much better!
- Parallelizability :- Even though in attention layer we find similarity for each key with entire sequence, this $O(n^2)$ operation is just matrix multiplication and can run parallel + efficiently on GPU's, whilst RNN cannot leave the sequence order in encoder state. $O(n)$ sequential operations
- Interpretability :- Self attention brings the interpretability aspect of our prediction, as we can easily visualize the reason for prediction through attention weights.



What happens in a single attention head?

I) Calculating key, query & value for each embedding

→ Once we receive an embedding of $d=512$, we use learnable linear matrices W_k, W_q, W_v to transform the embedding to a trio of key, query, value of 64 dims each.

→ Here (key, value) are closely related. They are like answers to questions posed by query. Key acts as an index of database which is actually pointing to value.

Example: - It is raining heavily on Monday?

Consider raining: Value, key mostly represents word rain itself and maybe its neighbor words by small amount. But for Query there are two possibilities how? → Heavily / When → Monday. Hence we see Query would be influenced by those two words rather than "raining".

II) Calculating similarity

→ Now once we have this trio. Let's take a Token query and compare with every other word in sequence including itself through dot product

(key)

$$\text{Similarity}_{\text{rain}} = \left[\vec{\theta}_{\text{word}} \cdot \vec{k}_{\text{rain}} \right], \quad \vec{\theta}_{\text{rain}} \text{ is } \vec{\theta}_{\text{rain, rain}}, \vec{\theta}_{\text{rain, rain}}, \vec{\theta}_{\text{rain, heavy}}, \dots \text{etc}$$

→ We also take tokens that are present after "rain" during similarity and above operation is done for all tokens in file. Hence we would have Similarity it, Similarity is, Sim_{rain}, Sim_{heavy} - etc. Totally

Need for scaling :- Due to growing dimensions (larger dimensions of key, value, a dotproduct gets larger, hence softmax puts emphasis mostly on the largest dotproduct suppressing the rest). These small softmax regions tend to receive / propagate extremely less gradients causing vanishing grad problem.

Ex:- During dot product $\vec{x} \cdot \vec{i}$ got $[1, 2, 1]$

(Y had bigger dim than X) → Another dot product $\vec{x} \cdot \vec{i}$ got $[-2, 4, 2]$

$$\begin{bmatrix} 0.21 \\ 0.57 \\ 0.2 \end{bmatrix} \rightarrow \begin{bmatrix} 0.1 \\ 0.78 \\ 0.1 \end{bmatrix}$$

less grad flow in '0.1' softmax

Divide each similarity by $\sqrt{\dim_{\text{key}}} = \sqrt{64} = 8$

$$\Rightarrow \text{Sim}_{\text{rain}} = \text{Sim}_{\text{rain}} / \sqrt{64}$$

III) Normalizing Using Softmax :- As expected and seen before, we will normalize the similarity scores to add up to one using softmax.

$$(\text{Sim})_{\text{rain}} = [0.1 \ 0.1 \ 0.5 \ 0.1 \ -0.1] \quad (\text{Sim})_{\text{heavy}} = [0.1 \ 0.05 \ 0.2 \ 0.3 \ 0.2 \ 0.1]$$

Finding answer for query :- Unlike databases fixed results, we would return a sort of fuzzy / probabilistic weighted (values) of matching keys.

$$\Rightarrow \text{ValueWeighted Answer}_{\text{rain}} = W_{\text{IT}} \times V_{\text{IT}} + W_{\text{IS}} \times V_{\text{IS}} + W_{\text{rain}} \times V_{\text{rain}}$$

$$\text{Answer}_{\text{rain}} = 0.1 \times \frac{[-]}{64 \dim} + 0.1 \times \frac{[64]}{-}$$

Formulas :- Scaled similarity per token per head

$$\text{Sim}_{ij} = \frac{[\mathbf{Q}_i \cdot \mathbf{K}_j^T]}{\sqrt{d_{\text{key}}}}$$

Soft max normalized sim

$$S_{\text{Sim}} = \frac{e^{\text{Sim}_{ij}}}{\sum_{j=1}^N e^{\text{Sim}_{ij}}}$$

Weighted Attention $(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_{\text{key}}}}\right)\mathbf{V}$

→ So at the end for each token of each attention head we have weighted values of entire sequence

Multiple Attention heads :

Why is it needed? This allows model to attend to information from different representation subspaces at different positions

What does this mean? Consider the sentence below

" Ram and sita went to forest and she saw a golden deer."

→ If we have a single attention head then it could probably predict only subject (Ram, sita), verb (went, saw) object (golden deer)

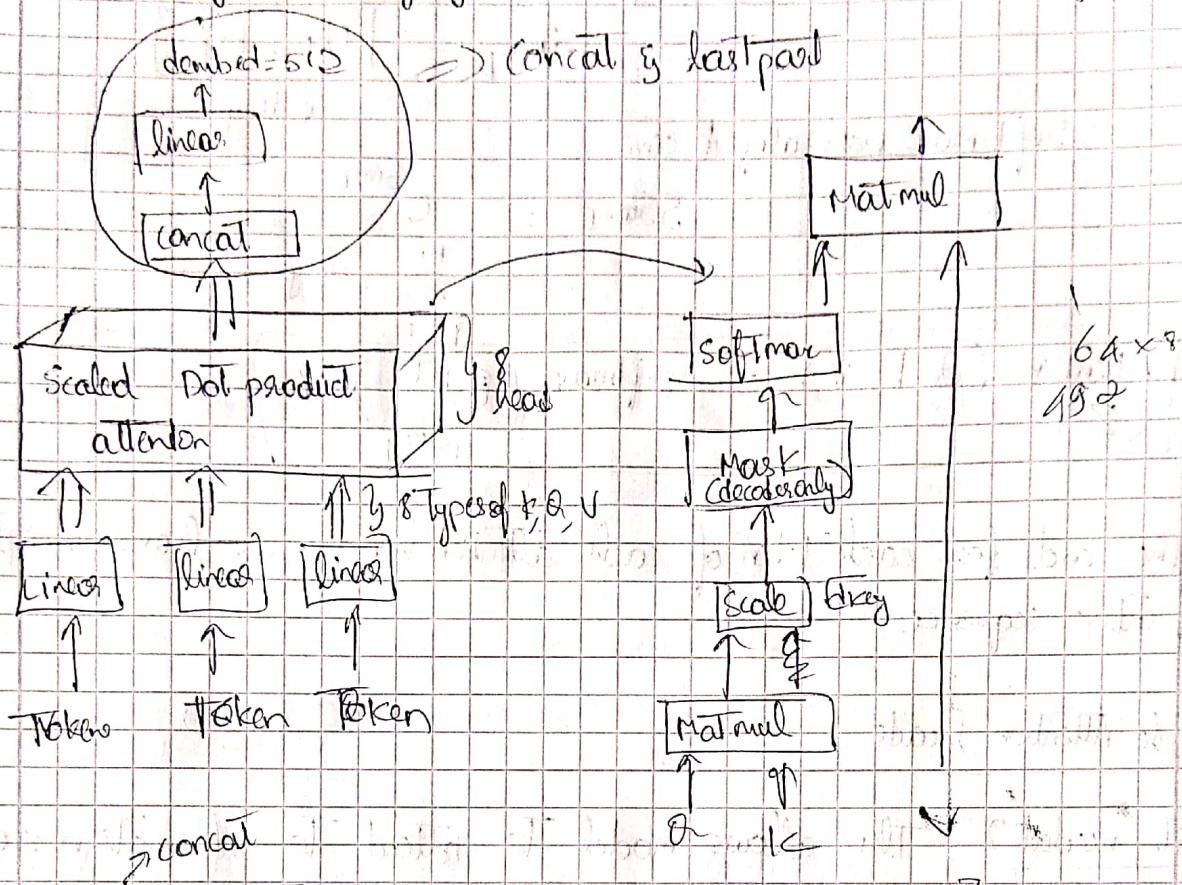
But is this sufficient?

For example:- if we ask the question "Who saw the golden deer?" we will get the answer as "She". Then we ask who is "she"? model has no clue, because it's not tracking gender. We humans tend to know it intuitively, but for machines we need to provide multiple perspectives of the same sequence

→ In Transformer model we have 8 attention heads, hence for each token we have 8 perspectives hence 8 different vectors of weighted value per token

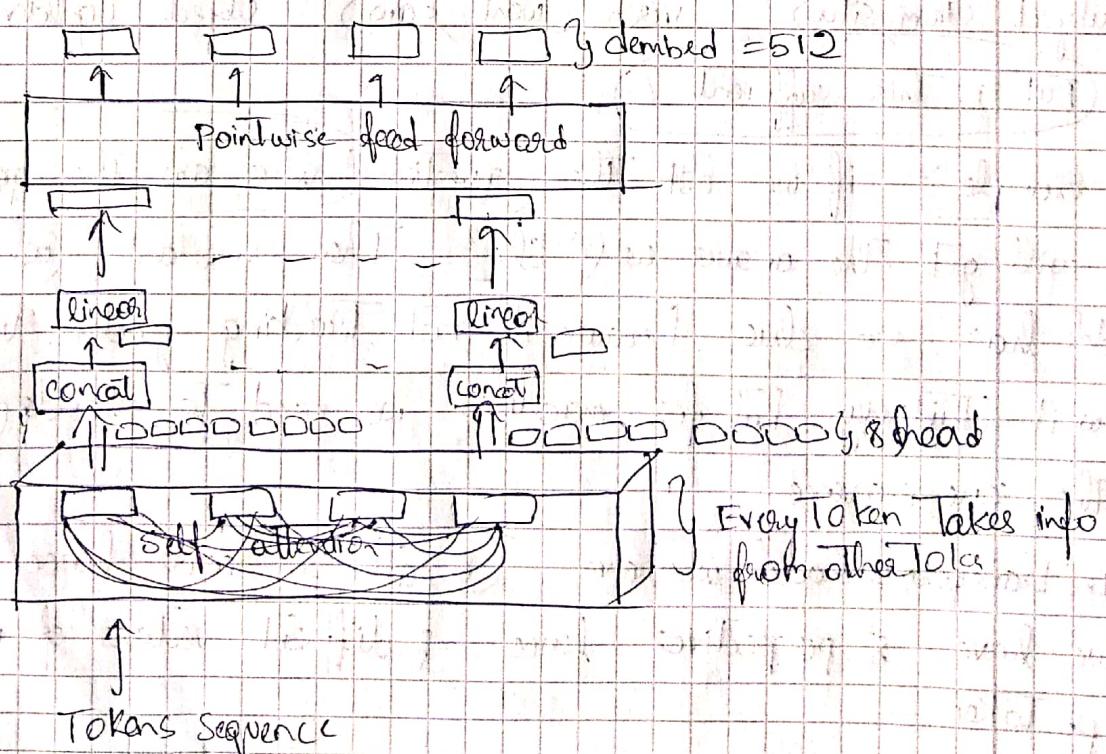
→ At the end we just concat all these 8 representations of each token to one vector per token. Then bring it back to embedding representation

$d_{\text{embed}} = 512$ by multiplying this with a learnable weight matrix W_0



Example: $Z_{\text{from}} = [Z_{1\text{from}} + Z_{2\text{from}} + \dots + Z_{8\text{from}}]$

Pointwise linear feed forward :- This same network is applied to all weighted value tokens in parallel. We also use a residual connection with layer normalization and dropout. Hence this operation is also done fully. ReLU activation is used



Decoder Side :- At decoder side we have both intra(self) attention and inter attention (encoder-decoder). Both are again multihead attentions.

→ As usual the decoder input is modulated by positional encoder and fed to decoder self attention layer.

→ MASKED decoder Self attention :-

Consider decoder predicting next word in below steps

- a) <st>
- b) <st> Je
- c) <st> Jesuis
- d) <st> Je suis Student
- e) <st> Je suis student <eos>

→ We can see that during prediction we don't know the next word hence we don't know all future words.

→ Hence it's better to mask out and not consider all future states during our attention calculation, This is done by using $-\infty$ as during softmax

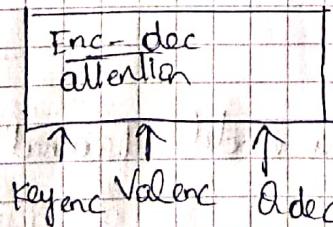
$$\frac{e^{-\infty}}{\sum e^x} = 0 \text{ hence zero attention}$$

a) [<st>, $-\infty$, $-\infty$, $-\infty$, $-\infty$] b) [<st>, Je, $-\infty$, $-\infty$, $-\infty$]

→ Apart from masking which forces to look into only previous decoder states we see that the attention calculation is same as that of encoder. This is like forming a question based only on previous experience.

→ Encoder-Decoder Intra Attention :- We receive key, values from last layer of encoder and use decoders query to lookup for answers on next token

→ This is like a question-answer type of scenario where the question



is asked by attention models decoder query and answer is taken from encoders states regarding next token.

→ Here masking is not present as we are actually looking for future token

→ At the end we use a linear + softmax layer to get next token and calculate loss & propagate gradients.

Training: We actually train the model by teacher forcing.
 For example: "Je suis étudiant" the first Token "Je" is predicted as "chound" which is wrong. We won't be feeding it back to decoder model for next cycle. Instead no matter what we will feed correct Token "Je" as token for next cycle.

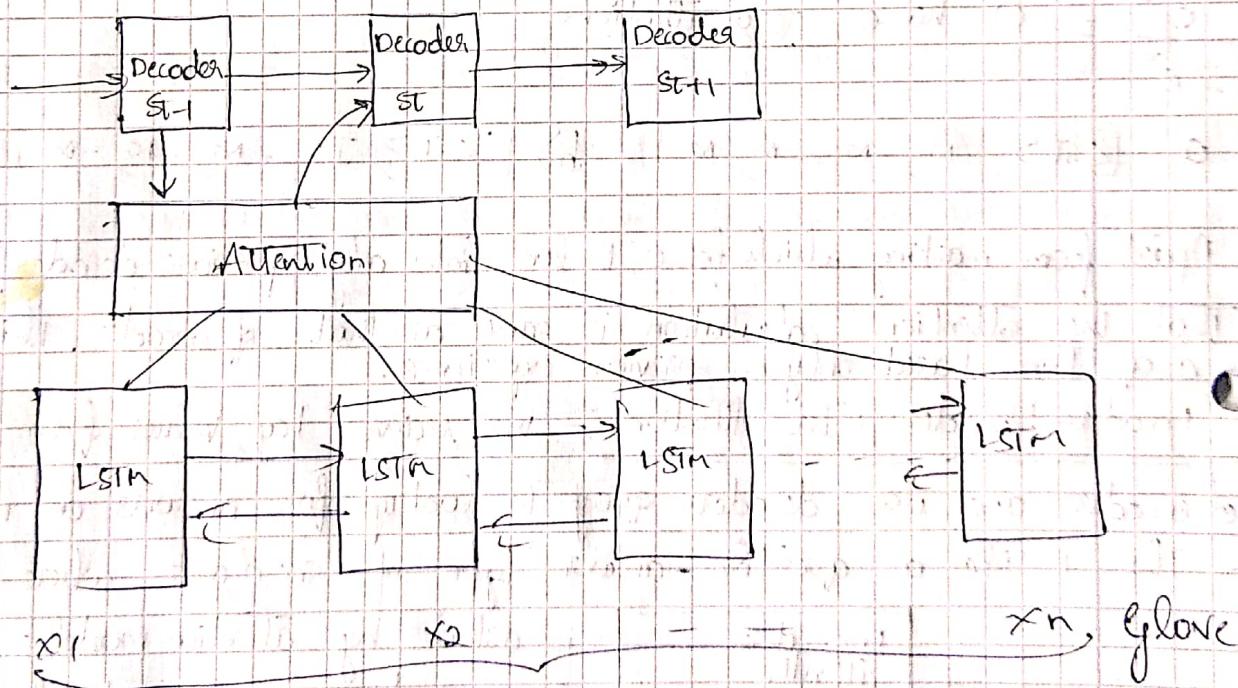
→ Also decoder model is shifted to right by one place hence we stop predicting as soon as we find <eos> Token

→ ~~3rd~~

Language Models

Contextual Word Vector (cove)

→ Contextual word embedding learnt from sequence - sequence machine translation model. Represents tokens in perspective of input sequence.



→ IT is nothing but a machine translation model trained on English-German or some other translation task.

→ Encoders are bidirectional whilst decoders can only look into past. We have additive attention to facilitate translation

→ Here the theory is during translation encoder stack has to learn about

Tokens syntactic (pos tag, sub-Vb+Obj) and semantic (gender, location, person) is hence after Training we detach the 'decoder' and keep only encoder.

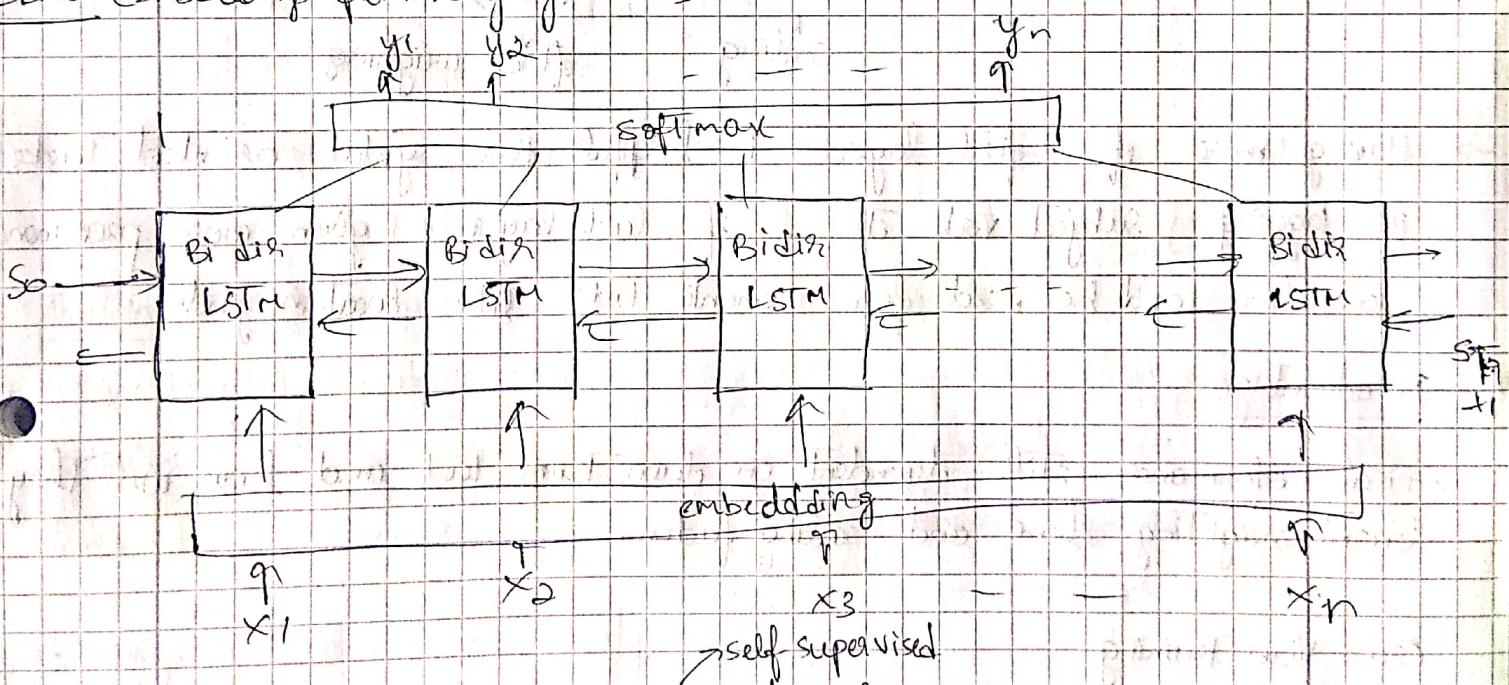
Usage during Runtime

- glove vector and core vectors are both concatenated as one provides skip gram perspective and another one contextual
- This is piped into Task specific model for further prediction

Disadvantages

- i) Pretraining phase := Supervised learning, hence corresponding datasets must be present for Training
- ii) Downstream prediction := Prediction quality directly depends on Task specific model, hence it's a major constraint

ELMO (Embeddings from Language Model)



- ELMO can be pre-trained in unsupervised fashion.
- The training task is to predict next word in the sequence given history and future.
- Hence we have bidirectional LSTM encoders which takes i/p from embedding layer
- Furthermore we stack LSTM on top of each other and to get

Predicting current word based on past and future

$$P(x_i) = \prod_{i=1}^n P(x_i | x_1, x_2, \dots, x_{i-1}) \rightarrow \text{past}$$
$$\prod_{i=i}^n P(x_i | x_1, x_2, \dots, x_n) \rightarrow \text{future}$$

ELMO representations

- After Training we detach the output layer and have $L+1$ total layers with us (LSTM + embeddings).
- According To Task we must learn to combine representation at each layer in weighted fashion. We also use a scaling factor γ to scale representations to ranges required by the task. This weighing scheme and scaling factor must be done by Task specific model.

$$\text{Final embedding } v_i = \gamma \sum_{l=0}^{L-1} \underbrace{\text{Task}_l}_{\text{scaling}} \underbrace{\text{Softmax}_l}_{\text{softmax weighting}}$$

- Having more of first layer is helpful for syntax oriented Tasks like pos Tagging, Subject Verb etc. iff last layer is given more prominence then elmo could be used for semantic tasks like location, gender

Disadvantages

- final performance still dependent on downstream Task and how well they learnt weighting scheme and scaling factor.

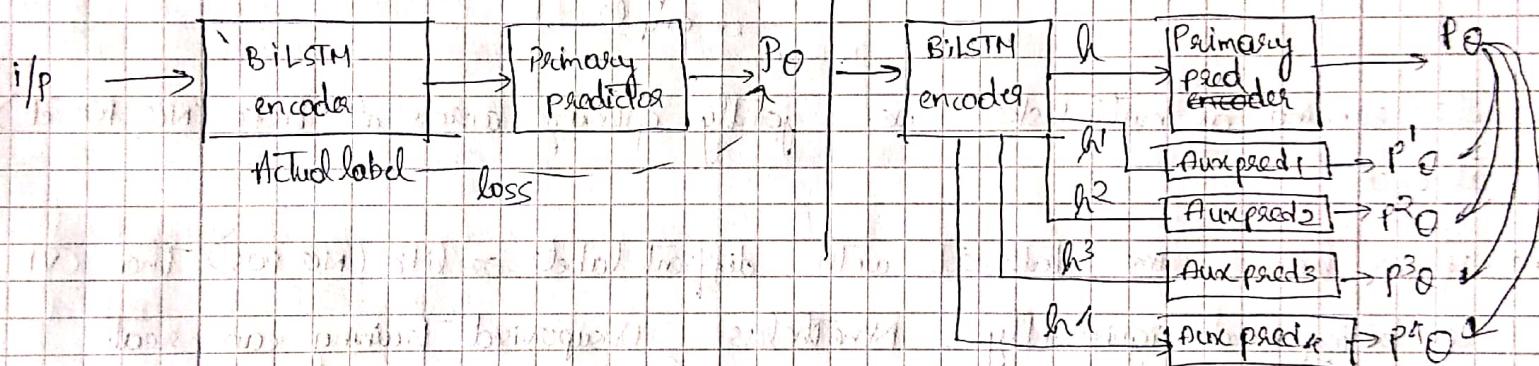
Cross View Training

- One of the best as it tries multiple new methodologies at once.
- We have a huge corpus, most of it is lab. unlabelled. We also have lots of small subsets of data annotated differently (like ner, pos, sentiment). We can combine all of the data in a big CVI module which at the end provides good generalizable embeddings over these multi tasks.

- Kannada NLP / NLP for languages, shortlisting resources
- During learning phase, model acts both as a both teacher that makes unsupervised main prediction and student that are trained on those predictions (auxiliary). For supervised, learning is usual with loss b/w actual & predicted propagated back.

unsupervised learning
 → During practical setting, if "auxiliary" has same i/p as main then prediction is redundant, Hence we give masked input To auxiliary modules and try to match them with main prediction.

Type of learning in CRF



Training method

- The prediction modules are really simple neural nets with softmax.
- We alternate the training b/w batch with labels (supervised) and batch with no labels (unsupervised)

Example of Training:

- First we randomly select a task of supervised learning
- We do the pass the i/p through Bidirectional LSTM encoder and make prediction. Loss is calculated as usually propagated
- Batch 1 → supervised Task 1 (Ner recog)
- Batch 2 → unsupervised Task
- Batch 3 → supervised Task 2 (POS Tag)
- Batch 4 → unsupervised Task
- Batch 5 → supervised Task 3 (dependency parsing)
- Batch 6 → unsupervised

$$\text{Loss}_{\text{sup}} = \frac{1}{\text{Num of sup examples}} \sum \text{CrossEntropy}(\text{Predicted label}, \text{Actual label})$$

Unsupervised Training

- This is the highlight of CNT which is main factor for improving quality of embeddings.
- For the batch we have K auxiliary modules which predict or are made to predict output as close to main prediction and auxiliary modules don't have access to complete encoder states, hence they learn the context of the word and its neighbouring tokens as well.
- We get loss by comparing auxiliary modules off with main module through distance function

$$\text{Loss}_{\text{unsup}} = \frac{1}{n} \sum_{k=1}^K \text{Distance Similarity} (\text{Main module}, \text{ } k^{\text{th}} \text{ auxiliary module})$$

| num of unsup examples |

- While running on multi tasks, we jointly optimize across all tasks (NER, POS, etc) at once
- If we have same data set with different labels like (NER, POS) then CNT can be improved significantly. Nonetheless, Unsupervised Training can create this effect from unlabelled data hence improves efficiency & run time.

Example = Sequential Tagger (tagging POS, NER, Locations, general purpose)

Example : He is visiting Washington, Tomorrow by Train. → Full sequence

Main (h^1) → Past context + current token

Aux 1 $\vdash (h^1)$ He is visiting, washington

Aux 2 $\vdash (h^2)$ He is visiting _____ → past context

Aux 3 $\vdash (h^3)$ _____ washington, Tomorrow by Train. → current Token, future context

Aux 4 $\vdash (h^4)$ _____ Tomorrow by Train. → Only future context

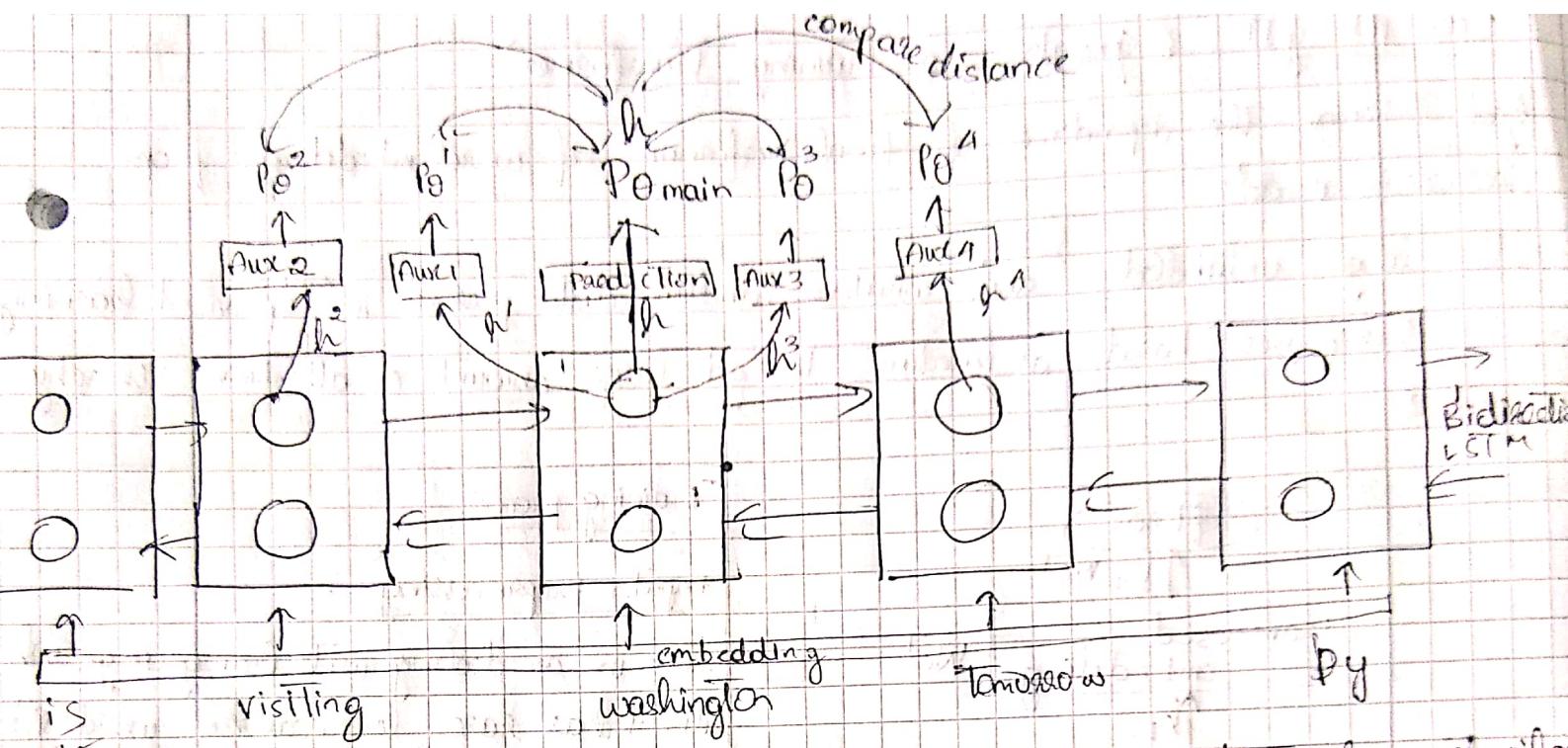
Loss unsup calculation \star \star → blanked out Rep.

Pred main = $[0.01, 0.05, 0.5, 0.25, 0.05, 0.05, 0.09]$

Pred aux1 = $[0.01, 0.1, 0.4, 0.1, 0.1, 0.1, 0.19]$

Pred aux2 = $[0.1, 0.1, 0.3, 0.3, 0.2, 0, 0]$

Pred aux3 = $[0.01, 0.05, 0.4, 0.2, 0.05, 0.19, 0.1]$



* We try to find distance between predicted softmax distro from auxiliary modules to main module softmax. By doing so system learns that "visiting" and "tomorrow" may be tagged to location.

* The encoders are taken from first layer of stack as top layers usually have seen all states and we cannot derive h_1, h_2, h_3, h_4 if top layers already seen all states.

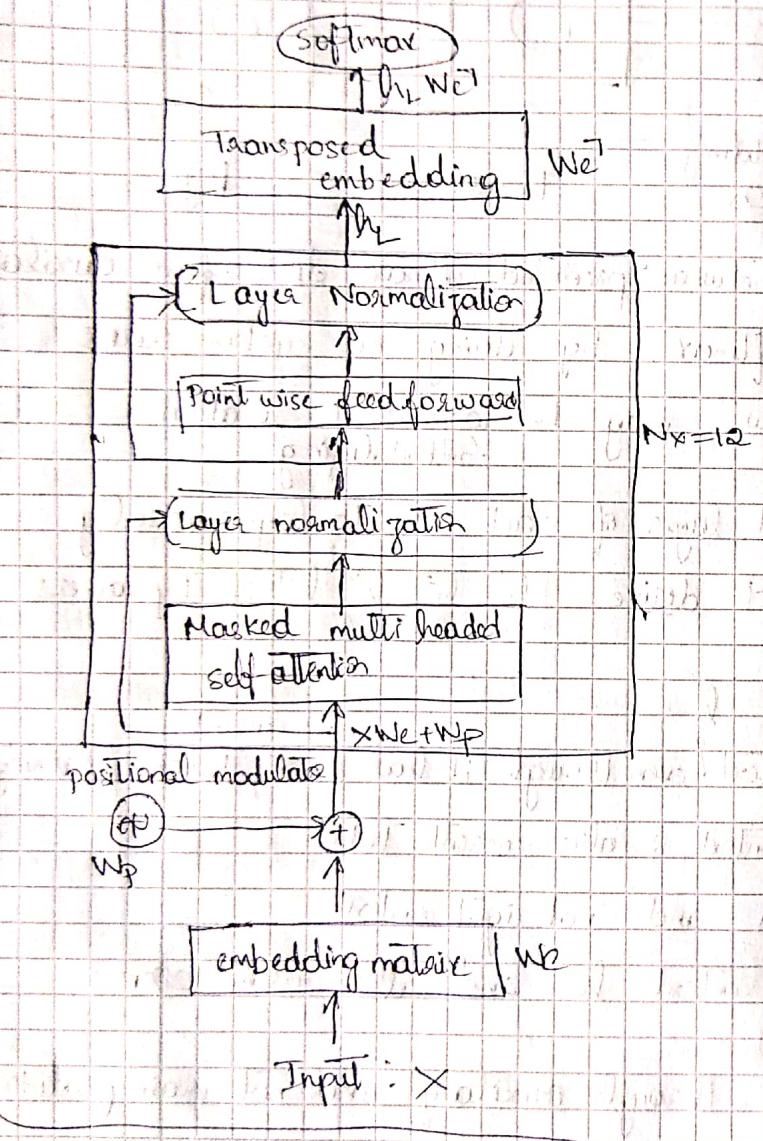
* The $h_1 \rightarrow$ predict w/o context even though it has seen Token "wash".
 $h_2 \rightarrow$ No access to recent context & also current Token
 $h_3 \rightarrow$ Access to current Token and not right context
 $h_4 \rightarrow$ No access to recent context "h" and also current token

* * CNT loss is back propped only through auxiliary and not main predict

Open AI GPT (Generative pre-training Transformer)

Goal: Remove the dependence of final performance of word embeddings from last specific models

- Use large unlabelled and whatever labelled data present in ^{semi} supervised learning
- Use Transformer based architecture to get large amount of attention thereby context



Output

- Finally The last layer's output is taken and we calculate next probable word through softmax.

Unsupervised training (Predicting next word in sequence)

- We predict next word of sequence given its k' predecessor/context window

$$\text{Loss unsup} = - \sum \log p(x_i | x_{i-k'}, \dots, x_{i-1})$$

Training process

Inputs preprocessing:

Since we are dealing with language which almost has rare, less occurring words, we cannot have every such word in our vocabulary.

- Use BPE To divide word to different segments and get most frequent sub-token that already present in vocabulary.

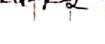
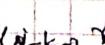
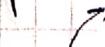
- We also add positional embedding so that Transformer decoder is position aware.

Contextual Model

- We use decoder part of Transformer with masked self attention as we are predicting next word in sequence.

- There are 12 such blocks which are stacked upon each other.

(for books This will be 10^6 but for Tweets around 10^3)

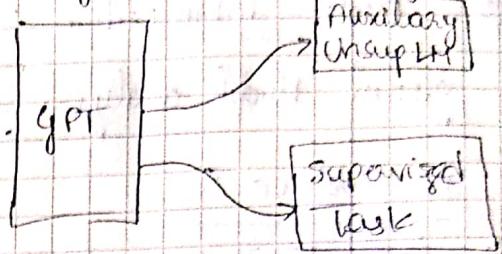


Supervised fine-tuning

→ Depending on final task we either choose last layers last embedding OR all embeddings.

→ Let's say our supervised task is sentiment analysis.

Then there is only one o/p per sequence, hence choose last layers last ~~/first~~ embedding



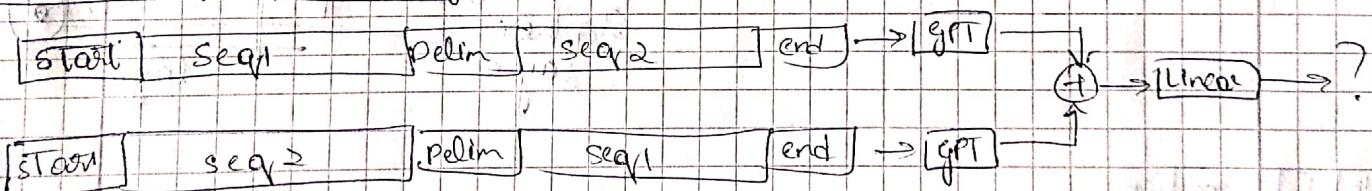
$$L_{\text{sup}} = \sum_{(x,y) \in D} \log P(y|x_1 - x_n) \quad L_{\text{aux}} = - \sum_i \log p(x_i|x_{i-1} - x_{i+1})$$

$$\text{Total loss.} = L_{\text{sup}} + \lambda L_{\text{aux}}$$

- Having auxiliary unsupervised training has lot of advantages as it leads to faster convergence and generalizability. However this is true for larger datasets. For smaller datasets, better to reduce Auxiliary Task by setting $\lambda \approx 0$.

→ We can use GPT for other purpose like Question-answer, finding duplicates, NLP solving by bringing data to common format as taken by GPT.

Consider Duplicate solving

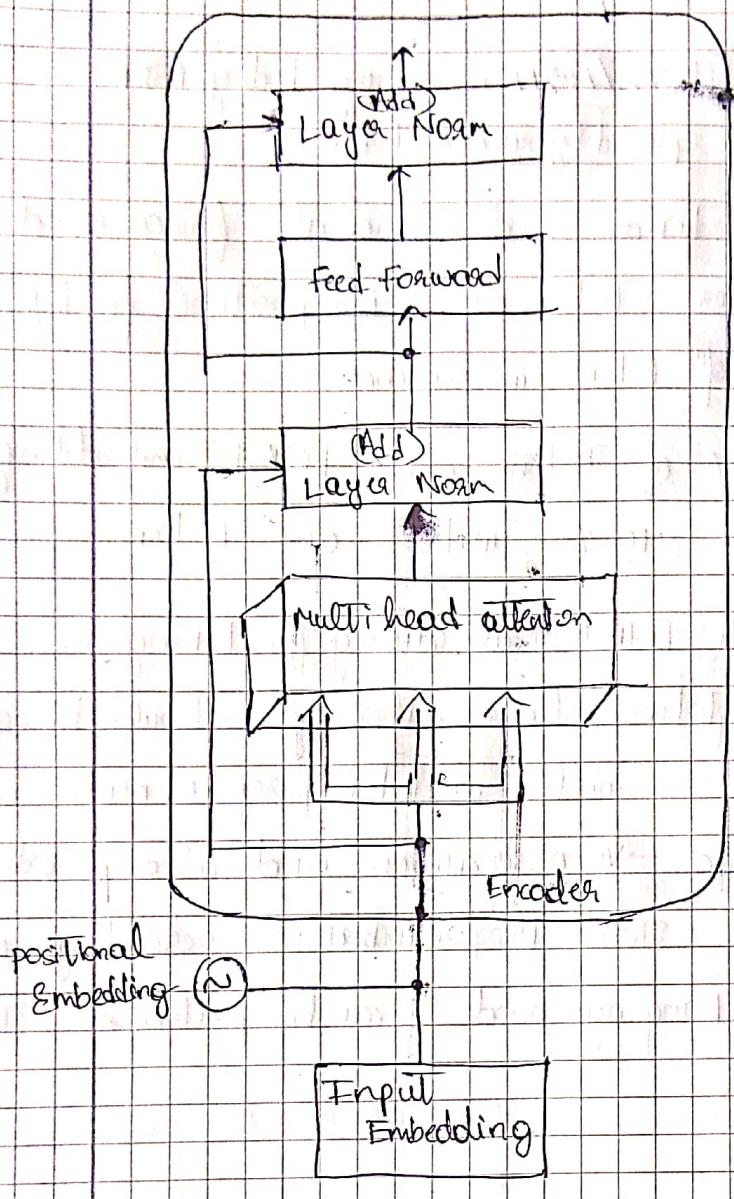


→ Since we don't have order here, we try both and specify '\$' as delimiter. Which Transformer has to learn by itself and not specified explicitly.

→ We can also remove as many layers from Transformer stack and put up our own during Supervised fine-tuning. But generally it's better to not remove any layer as it leads to bad performance.

→ Since unsupervised is trained on variety of tasks like books, Tweets having extreme long & short range attention we see that model is better.

BERT (Bidirectional Encoder model for Transformers)



Principle : Use encoder of Transformer
To generate contextual info and embeddings

NLP Pretraining :-

- a) Masked Language Model (MLM)
- b) Next Sentence Prediction (NSP)

Advantages w.r.t GPT

- Bidirectional context consideration during generation process
- Pretrained using NSP provides additional boost on learning sentence relations
- Learns [SEP], [CLS], [A/B] tokens during pretraining itself unlike gpt.
- Trained on additional Wikipedia corpus

Advantages of Bidirectional Training and need for NLM

- In gpt Version 1 we were predicting nextword given previous words, We were not considering future context. Hence may not work in some cases as below.
- This is especially helpful for languages where verbs may appear at end
Ex:- German \Rightarrow Ich möchte buch schreiben \Rightarrow I would like to write a book
- This future context is also helpful in some question - answer / similar sentence prediction scenarios which gpt v1 cannot fully cater to: Hence bidirectional
- Problem with the bidirectional stuff is we would already know word after first layer as its embedding would be used for self attention calculation of other tokens. Hence learning doesn't happen. Need a novel training approach.

BERT Pretraining procedures

a) Masked Language model (MLM)

Ex: [CLS] The sky is [MASK] and The ~~Saturn~~ ^{is} shining Today [SEP]
masked replaced by random Token

→ To avoid seeing the token in bidirectional encoder and also force encoder to learn context at each token rather than at specific position first / last/middle we randomly mask 15% of tokens in sequence.

→ Masking process : Replace 80% of tokens with [MASK] embedding 10% with random token and 10% with same token

→ We don't replace all 100% with [MASK] as during fine tuning we most likely won't use [MASK] tokens, hence other types should be represented

→ 10% with random token and 10% with same token replacement helps model to learn what may cause → ve association and also provide slight bias to our pretraining. Since the number or percentage of random replacements are less, language model can be parallelized with relative ease.

b) Next Sentence Prediction (NSP) :-

→ Helps to contextualize relation b/w two sentences and long range thinking.

→ The embeddings of two sentences are learnt

Ex:- [CLS] Weather is good [SEP] Sky is clear [SEP]

Token embedding $E_{[CLS]}$ $E_{[\text{Weather}]}$ - - - $E_{[\text{SEP}]}$ - - - $E_{[\text{Sky}]}$ $E_{[\text{clear}]} [SEP]$

segment embedding EA EA - - - EA EB - - - EB EB

Position embedding E_1 E_2 - - - E_5 E_6 - - - E_g

Input :- Tokens are further divided to Bipe pair encodings to cover for infrequent words or new words. This also helps to standardize input vocabulary and Training.

Fine Tuning :- Fine Tuning using BERT is relatively simple

Single Sequences :- As similar to other models, enclosed with [CLS] & [SEP]
embeddings [CLS] SA [SEP]

Pair Sequences :- Concatenate Two sentences and Train them in usual manner, we don't need to encode them separately as we have learnt segment embeddings in our pretraining Ex :- [CLS] SA [SEP] SB [SEP]

Output :- If it's an entailment / sentiment prediction task then we can make use of [CLS] Token representation, pipe it through a simple Linear layer and fine tune it for our output

→ Complete representations too can be feed to output layers for tasks such as Pos Tagging

Language Models

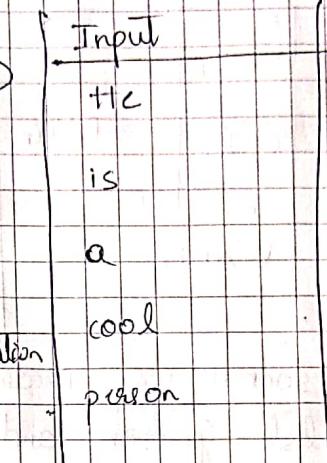
Skipgram model :- Given a word predict its context

Ex:- He is a cool person \Rightarrow

→ Much more successful in predicting infrequent words than CBOW

→ More training time than CBOW

→ More helpful for embedding calculation due to context



Output Context

(is, a)

(he, a, cool)

(he, is, cool, person)

(is, a, person)

(a, cool)

CBOW Model :- Given contexts predict a word

→ Much success in predicting frequent words

→ less training time and hidden embeddings are averaged out to get smooth dist. to

→ More helpful for "Fill in the blanks" type of prediction (autocomplete)

Category

CBOW

CBOW

Skipgram

Training Time

Comparatively less

More

Prediction Task

Fill in the blanks / message
autocomplete / (synonyms)

Contextual rich embedding and
useful for finding similar words
(neighboring concepts)

Failures

Cannot predict infrequent words

(cannot do autocomplete tasks)

Ex:- Today was

and difficult to use for prediction
of words.

a) beautiful

b) good \rightarrow most frequent

c) delightful \rightarrow rare

Ex:-

? ?

delightful

In CBOW, we use "good" or "beautiful"

than "delightful" even though it was correct

This is due to averaging of contexts

Problems with Word2Vec models

- i) Often we have a vocabulary of millions of words, hence million inputs and million outputs. This causes weight matrix to blow up and slow down computation & Training.

Solutions :- Negative Sampling, Subsampling, Hierarchical softmax, Phrases, etc

A) Subsampling :- waste

Ex:- The quick brown fox jumped over the lazy dog.

→ Lots of stopwords/prepositions do exist in each sentence which don't add much to context

→ But cannot remove it all together as we need to predict these words too

→ Use a frequency based approach and decide the inclusion of word during training based on frequency. Sample these words using sampling factor

→ Lesser the sampling factor lesser are chances of word to be picked

$$\text{Prob of picking the word } P(w_i) = \left(\sqrt{\frac{z(w_i)}{\text{sampling factor}}} + 1 \right) * \frac{\text{sampling factor}}{z(w_i)}$$

where $z(w_i)$ = $\frac{\text{frequency of occurrence of word } w_i \text{ in corpus}}{\text{sampling factor}}$

⇒ "Quick" would appear less than "The" hence Quick would mostly be used during Training and "The" would most likely be not

Let say $z(\text{"The"}) = 300$

Total words = 1000

$$P(\text{"The"}) = \left(\sqrt{\frac{300/1000}{10^{-3}}} + 1 \right) \left(\frac{10^{-3}}{300/1000} \right) = \frac{1.07 * 10^{-6}}{300} = 3.39 * 10^{-9}$$

$z(\text{"Quick"}) = 10$

Total words = 1000

sampling factor = 10^{-3}

$$P(\text{"Quick"}) = \left(\sqrt{\frac{10/1000}{10^{-3}}} + 1 \right) \left(\frac{10^{-3}}{10/1000} \right) = 31.8 * 10^{-5} = 3.17 * 10^{-3}$$

$P(\text{"Quick"}) \gg P(\text{"The"})$

2) Hierarchical softmax

→ Put the output words in a tree structure so that we need to compute only the $O(\log N)$ softmax probs rather than $O(N)$ direct.

$$P(\text{for}) = P(\text{left}, 1) * P(\text{right}, 2) \Rightarrow \text{node num} = 3 \Rightarrow \text{"quick fox"}$$

→ For input "fox" we only go to "quick" path by computing prob of taking branch

* Applicable only during Training as we know where to go and what is the target node. During prediction, we still need to compute softmax of all words

3) Negative sampling :- For each target word there may be 10-20 target words based on window size, rest all are negative non-required outputs

input word \Rightarrow 1 million output \Rightarrow ggg, ggo \Rightarrow no \Rightarrow Sample \Rightarrow Take only 100/1000 random example

→ Instead of complete random selection, we can attach some probability of picking a -ve word based on Term frequency

$$\text{Prob of picking -ve sample (Wordi)} = \frac{\text{freq (Wordi)}}{\sum_{j=0}^n \text{freq (Wordj)}}$$

→ More frequent words like "The, in, but, a, an" etc would be picked a little less. and infrequent -ve samples more