

Paper: ~~Transformer~~ - Unseen

- before appearance of transformers, Rnn, LSTM, gru were used a lot.

1

◆ let's dive into Rnn

- what is Rnn: Rnn is a kind of neural network which is useful for dealing with sequence data.

↳ sequence data: they all have order

↳ like text data → you are great boy

↳ if we destroy order between them, they don't

make sense

* Rnn's work with text data most of time

- ↳ some application of Rnn \Rightarrow
- 1) machine translation
 - 2) sentiment analysis
 - 3) language model
 - 4) text topic

* all of tasks above work with text data which have order

↳ the movie was ~~awful~~ so Funny \rightarrow positive

* Hello, are you okay? \rightarrow translation.

↳ language model \rightarrow predicting next word \star the most important task in nlp

↳ language model is basic, simple like image classifier.

in cnn

- let's dive into know more about language model.

↳ predicting next word in sentence

↳ backbone of nlp tasks

↳ gives a probability of next word based on previous word

- example:

you are great

boy
girl
desk
actor
Football Player

- probability of all words in vocabulary based on previous

- so all of them have to keep somewhere to know about next word

$$P(x^{t+1} | x^t, x^{t-1}, x^{t-2}, \dots, x^0)$$

next word
or
last word

↳ big difference between Rnn with other kinds of neural network:

- neural networks like CNN, ANN are not able to consider time

- there is no difference between input data

- they are not able to understand order between data

you are good boy

- their input must be fixed → what does it mean?

CNN → $48 \times 48 \times 1$ → input data must be same size

ANN → (24,)

- using Padding for text data? → sequences are not same size
increasing number of input data → Parameters *

For CNN, and, it does not matter which word comes first

2

you great are boy \rightsquigarrow don't make sense
are boy your great \rightsquigarrow

because they have order \rightarrow RNN comes to tackle problem.

- they can process words with order through time \rightarrow in order

first: you goes to network

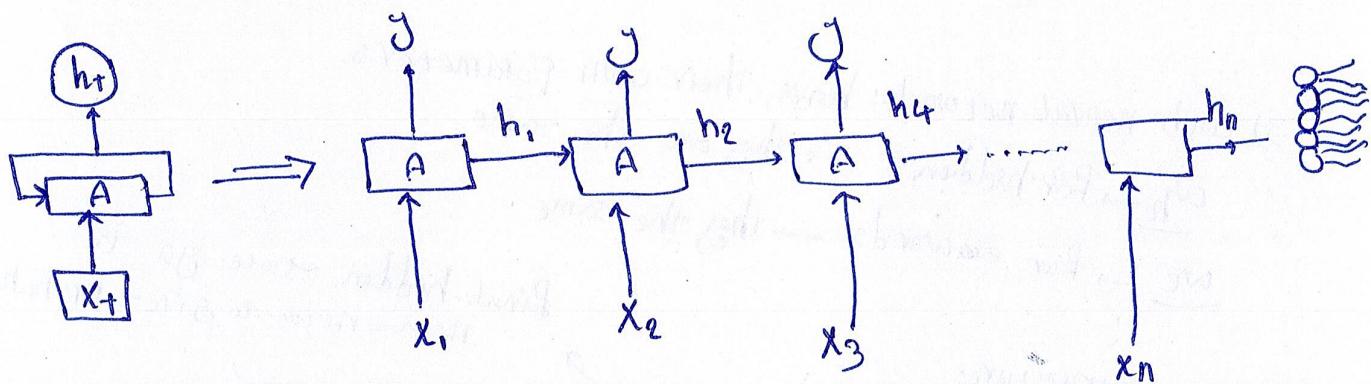
second: are goes to network

third: great goes to —

- For processing each word, we need previous to make sense

for next word prediction

- lets see structure of RNN:



more general:

- *) all of units have two input, haven't first one
- **) we insert a random vector to first one \rightarrow to make them similar

*) some tips:

1) words are not able to enter neural network directly:

\hookrightarrow they go to somewhere called embedding factory

to make a array of them.

you \rightarrow embedding \rightarrow $\begin{bmatrix} 0 \\ 0 \\ \vdots \end{bmatrix}$

are \rightarrow $\sim \sim \sim$

great \rightarrow $\sim \sim \sim$

2) what is hidden-state: it is a context vector which stores

information from each timestep to predict next word

- combined with word at time t , $h_{t+1} \rightarrow$ to predict x_{t+1}

$h_t \rightarrow$ ~~also~~ carries information from previous part

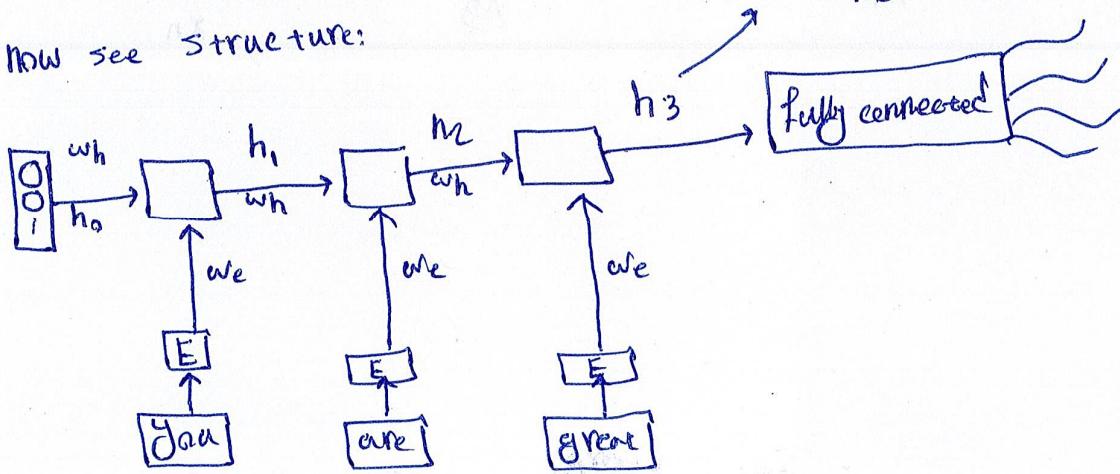
3) each neural networks have their own parameters.

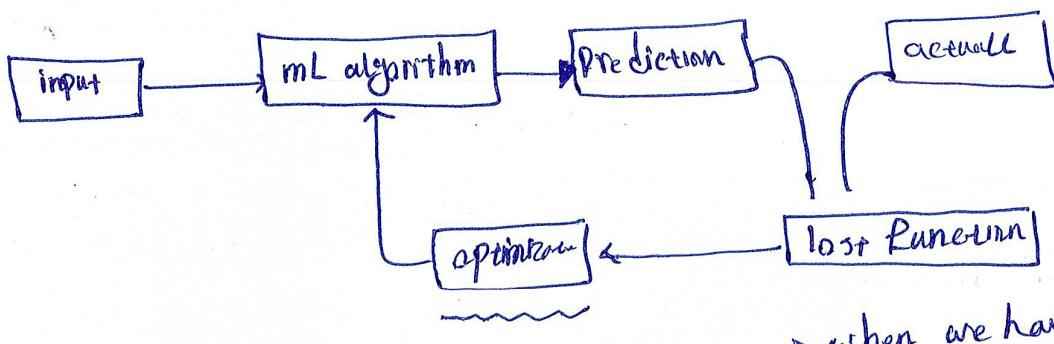
w_h \rightarrow for hidden's \rightarrow they are the same.

w_e \rightarrow for exwords \rightarrow they the same.

Final hidden state go to
neural network to give probability

- now see structure:





* there is big problem in optimization \Rightarrow when we have long input word

For the first time all weights are random. coarse

- we must use backpropagation to update our weights. $\rightarrow w_h$

- after some step \rightarrow gradient will disappear \rightarrow we can't update

our weights \rightarrow training will stop \rightarrow target data which comes

first

* Vanishing gradient descent

* exploding gradient decent

LSTM comes to address problem of vanishing gradients.

4

* LSTM is special version of RNN → solves problem through concept of gates

* language model → predicting next word is so challenging when it comes to backpropagation → vanishing gradient

* LSTM has 3 kinds of gate

- forget gate
- input gate
- output gate.

- LSTM's are capable of learning long-term dependencies.

- LSTM introduced 1997

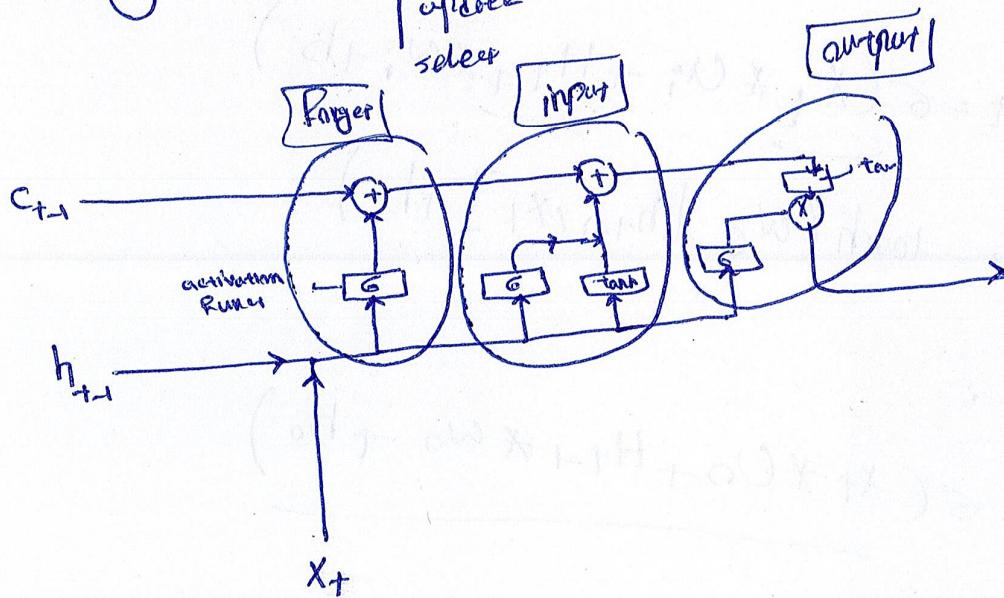
- they work well on wide variety of sequential problems.

* remembering the context of the sentence for long period of time is their default behaviour. → they have memory

* memory is used to

- delete
- insert
- update
- select

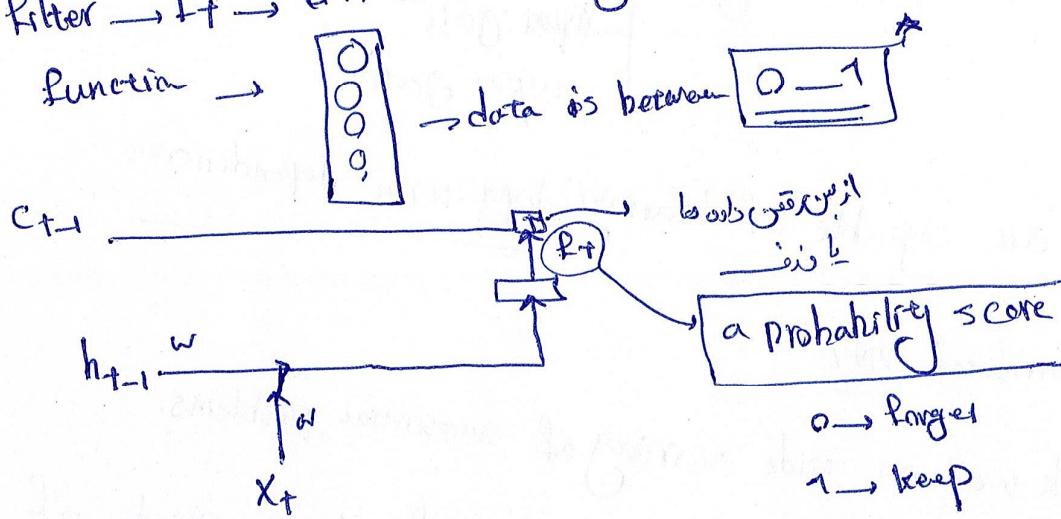
data



Forget gate

example of text prediction:

- Prince was Feeling uncomfortable in the party while chiray enjoyed the party.
- there some important word's here → subject, feeling / adjective...
- in Forget gate all ~~relevant~~ irrelevant information is removed via filter
- Filter → P_f → it is created by multiplying x_t and h_t → Passing sigmoid



$$\text{Forget gate} \Rightarrow P_f = \sigma(x_t * w_f + h_{t-1} * w_f + b_f)$$

$$C_t = P_f * C_{t-1} + i_t * \hat{C}_t$$

$$\text{input gate: } i_t = \sigma(x_t * w_i + h_{t-1} * w_i + b_i)$$

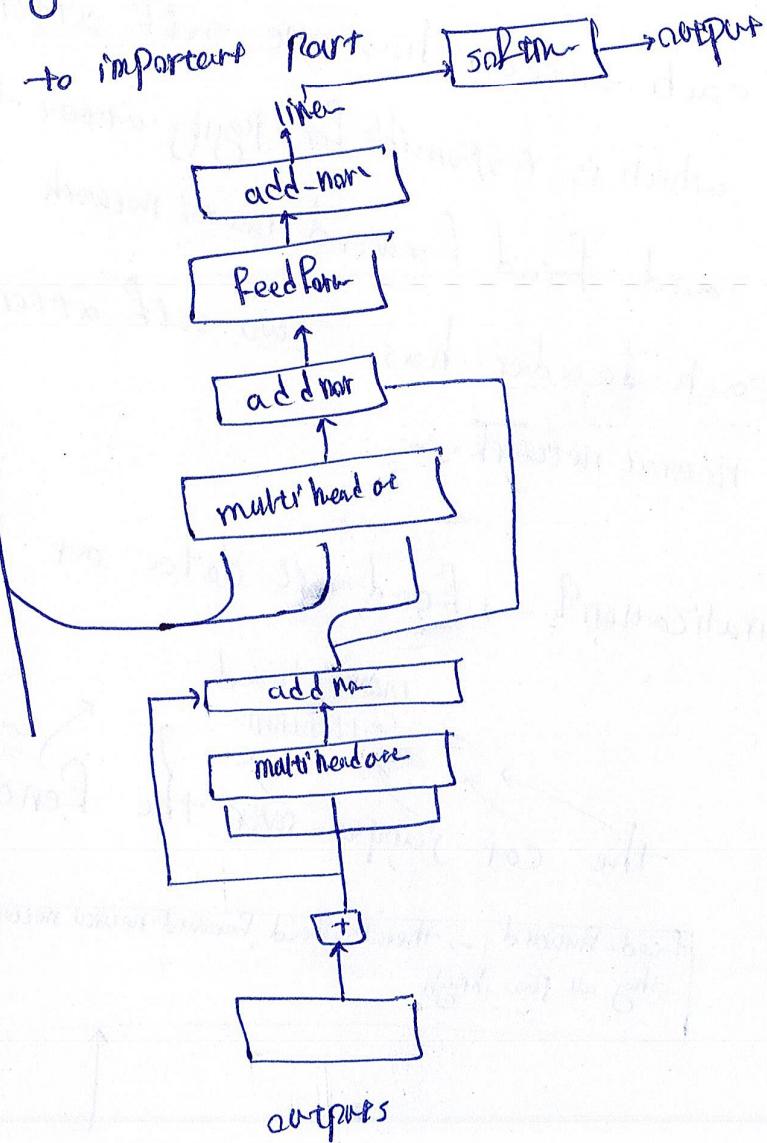
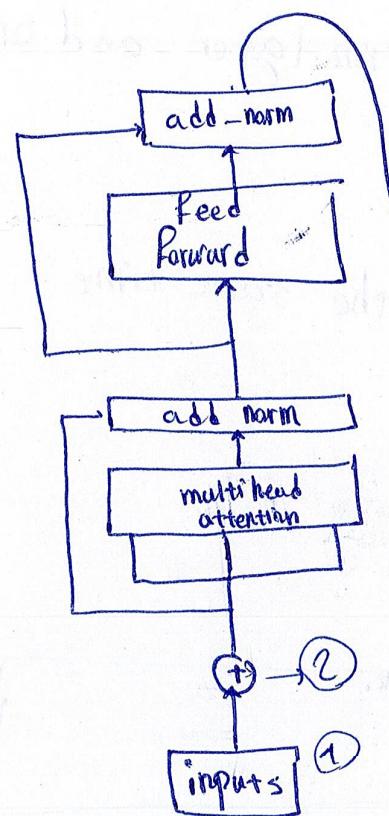
$$C_t = \tanh(w_c [h_{t-1}, x_t] + b_c)$$

output-gate:

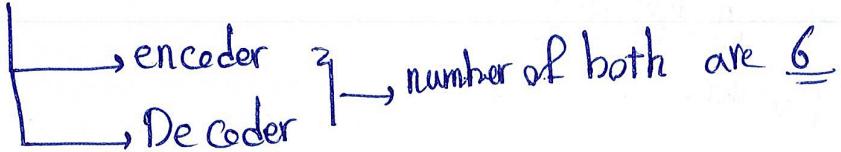
$$o_t = \sigma(x_t * w_o + h_{t-1} * w_o + b_o)$$

Transformers:

- Transformers structure only rely on attention mechanism.
- they do all work parallel.
- they don't use any Rnn, Lstm, — → their problem.
- they more faster, smarter and quicker than → ~~process~~ parallel way.
- * what is attention? in general, attention is the ability of a model to pay more attention to important part



Transformers



- right hand side → Decoder

- left → encoder

* each encoder has one self attention layer
which is responsible for paying attention to itself → input text

and feed forward neural network

* each decoder has two self attention layer and one
neural network.

* parallelization? → Feed all data at the same time

the cat jumped over the fence

Feed-Forward → there's Feed Forward neural network.
they all pass through

multi-head attention

the cat jumped over the fence

there is self attention layer → tries to find relation
and communication between words

cat

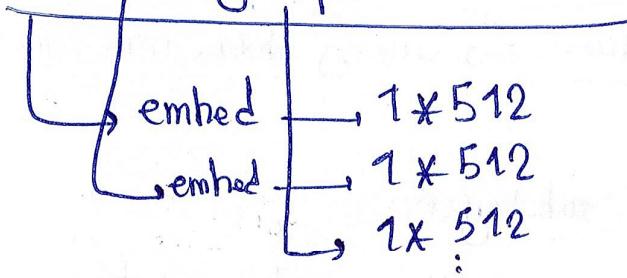
the ~ jumped

over

the cat jumped over the fence

- we have six encoders. → each of them, they are different
- * all input either encoder or decoder ~~or~~ are embedded

* The cat jumped over fence



what is embedding → this is a way for presenting words

* each word represents by 512 numbers.

↳ this is hyperparameter

2) Positional encoding:

- since the model do not any Rnn-LSTM, how does it able to understand order of the data?

- by adding Position encoding → let it know → where words come

* output and linear →

output of decoder can be transformed into something to be understood.

- Normalization layers:

- *) they are between Feed Forward Layer and Self Attention Layer
- * what do they do?
 - *) responsible for \rightarrow normalize output of P sublayer

\rightarrow an improvement over batch normalization \rightarrow collecting data into same region

- skip connection:

- there are some arrows around sub-layers.

- they carry ~~out~~ information doesn't go to ~~either~~ either
self or Feed Forward.

- It helps model not forget things.

helps data go through model.

*) novel idea

- multi-head attention
- Positional encoding.

Lets Dive into them

There are only 2 kinds of multi-head attention:

1) multi-head attention.

2) masked multi-head attention.

Difference:

- in multi-head attention all the words compared with all the other words
- in masked multi-head attention only the words that are coming before a word compared to that word.

Let's go to attention-layer

1) there is scaled dot product

* How does attention calculate?

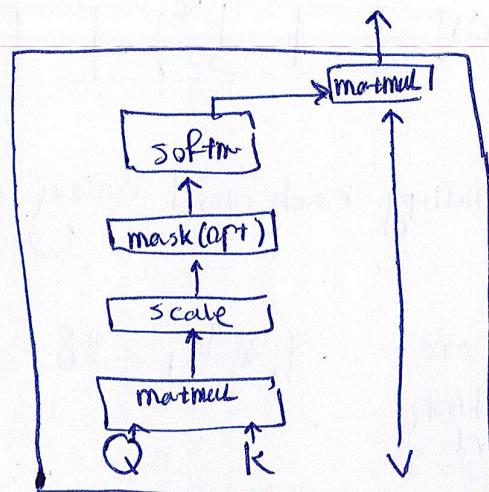
input that was Run

emb 1 512 1 512 1 512

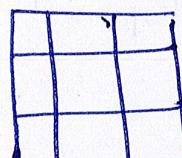
Q 1 64 1 64 1 64

K 1 1 1 1 1 1

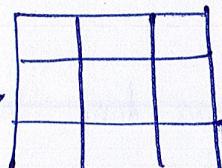
V 1 1 1 1 1 1



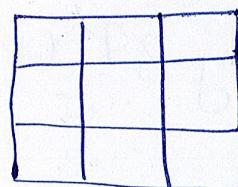
Query



Key



Value



1) multiply this embedding vector to these matrix

2) values of matrix are random, during training they will train to be learned

Let's go to calculation:

attention will give a score → it means how much ~~more less~~ ~~more less~~
between each words.

- Dot product \Rightarrow show similarity \rightarrow similarities

that

was

fun

q_1 [] [] []

q_1 [] [] []

q_2 [] [] []

Q

key k_1 [] [] []

k_2 [] [] []

k_3 [] [] []

value v_1 [] [] []

v_2 [] [] []

v_3 [] [] []

multiply each word query with others word keys

Score

$$q_1 * k_1 = 98$$

$$q_1 * k_2 = 46$$

$$q_1 * k_3 = 78$$

for first
word

second

$$q_2 * k_1 = ?$$

$$q_2 * k_2 =$$

$$q_2 * k_3 =$$

third

→ all of them are done together.

when we have all the score → we divide them by
8 → why 8? roots of 64 → is length of
key, Query, value length

$$\text{Score} \quad q_1 * k_1 = 98 \quad q_1 * k_2 = 46 \quad q_1 * k_3 = 10 \quad 0$$

Divide
by 8

$$12,25$$

$$5,25$$

$$9,5$$

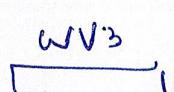
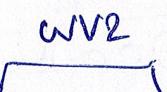
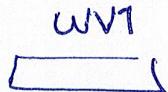
because it's root of keys values ; Queue Dimension

$$\text{softmax } a_{1,93}$$

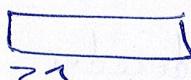
$$a_{1,01}$$

$$a_{1,06}$$

weighted $wv1$

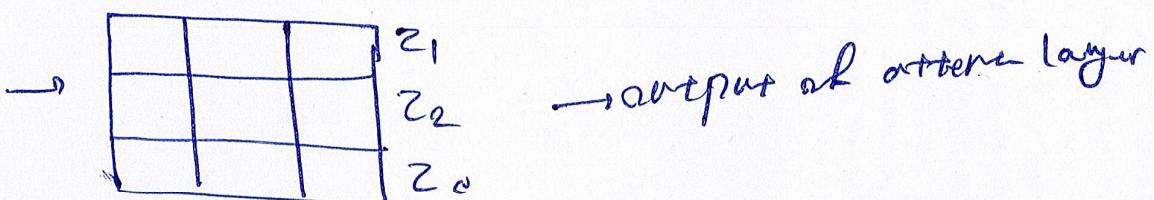


sum



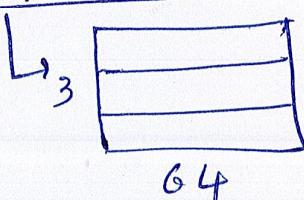
→ Weighted query -

new query -

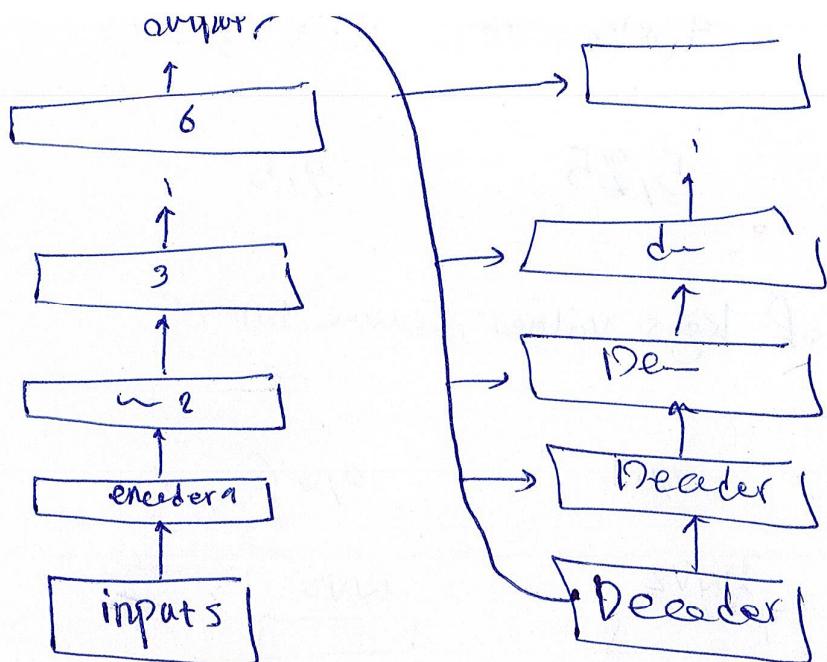


* multi-head atten does 8 times

* output of each multi-head attention



* we have eight output → connect them to make original dimension.



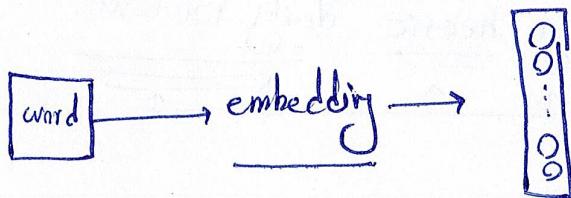
output goes to multi-head attention of decoder

another views

I

- Embedding Layer takes a word, gives a vector of number

- they are const



* علاوه بر سیم کلیسی سیندرلیسی،
اور جملہ کلیسی را دھننا ماند
لماں کی درج.

word \rightarrow embedding \rightarrow Vector بیر تماقی این کا تابع ہے

- سیندریں حیوان خطرناک است.
- سیندر اصلی دلیلی.

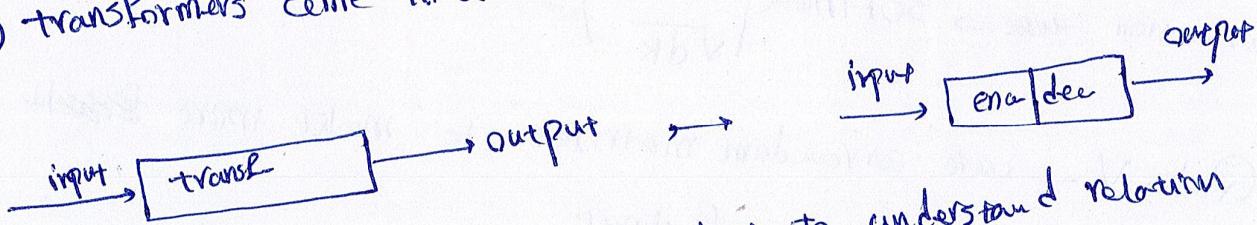
- ملک قبول کرنے ایسا یا کام کا ملک کیا ہے؟

- سیندر جملہ کامل مکار است.

- یعنی نتائج بہ فتح و مخفون و بار جد نہ اور.

- انگلیوں کا کارڈنال سیم دیگر عوامل در مقابلہ
نہ است.

* transformers come to address.



in encoder \rightarrow we use self attention mechanism to understand relations
between data/text \rightarrow this because word like book has different

meaning in each sentence.

- by using / calculating attention score.

you are great boy \rightarrow t. good job displayed

- transformers apply to machine translation task.

* I am really hungry, because of my hectic busy routine.


- all of them go through embedding layer

then by using self attention \rightarrow (explained) \rightarrow calculating scene
by dot product

\rightarrow scaled Dot product \rightarrow find relationship and build up

a embedding layer.

- so, encoder builds up powerful embedding vector
including context.

- attention score \Rightarrow $\text{softmax} \left(\frac{\mathbf{Q}\mathbf{k}^T}{\sqrt{dk}} \right) \mathbf{v}$

- $\mathbf{Q}, \mathbf{k}, \mathbf{v}$ \rightarrow are random matrixes to make more flexible
flexible and reduce size of input

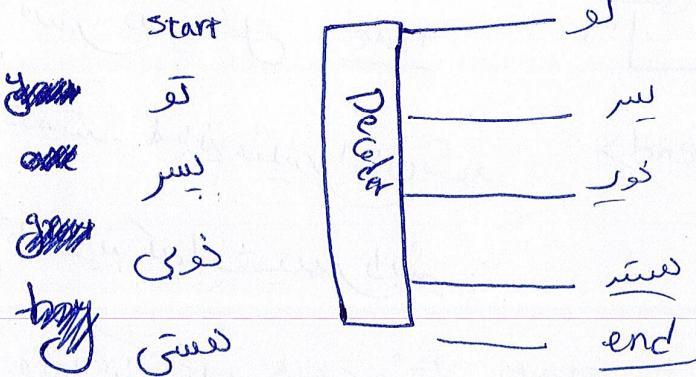
- $\sqrt{dk} \rightarrow$ is root of matrixes dimen $\rightarrow \underline{8}$

- transformers apply on machine translation
- it has difference between training / testing process

- training Part:

input of decoder \Rightarrow a shifted output, adding start

token.



$\langle \text{start} \rangle \xrightarrow{\text{Dec}} \langle \text{end} \rangle$

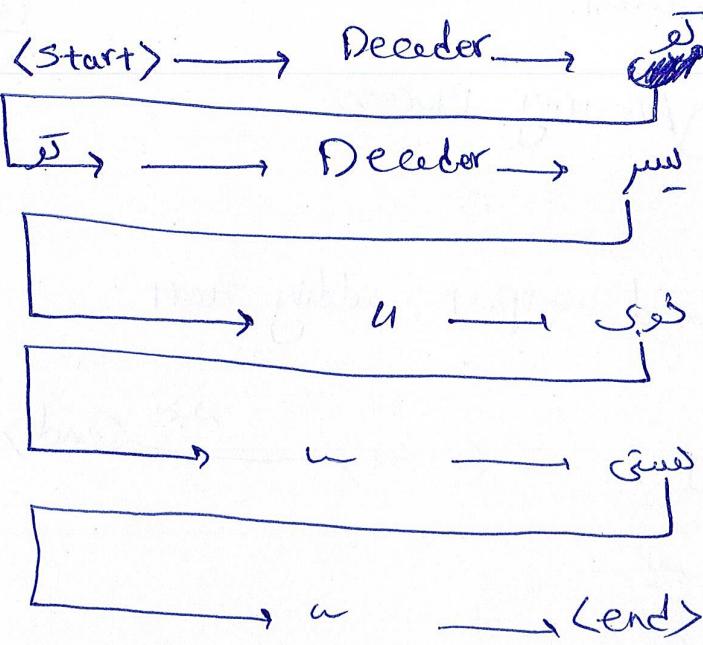
- in decoder, we have generally data

- decoder has two input because of two attention mechanism
input dec

1) is responsible for finding relation between output
2) encoder output with decoder input.

a) mask attention step cheating structure

we prepared all information to produce next word. well done.



test نتایج در داده دور ریاضی کنند
برای این تعلیم داده است. اکنون
بهم توان درست توجه داشتیم.
- در فرآیند train می‌باشد
- ورودی کامل true
لسته. با این فنیز اعمال شود
که تواند کمترین تأخیر را داشته باشد.

- کجا باید کلر کوچک مفهوم start و end را داشته باشد

$$\text{softmax} \left(\frac{QK^T}{\sqrt{dk}} + \frac{\text{multi_head}}{h} \right) V$$

کجا باید کلر کوچک مفهوم start و end را داشته باشد

- softmax multi-head attention \rightarrow all entities are available
masked \rightarrow some of them are available.