

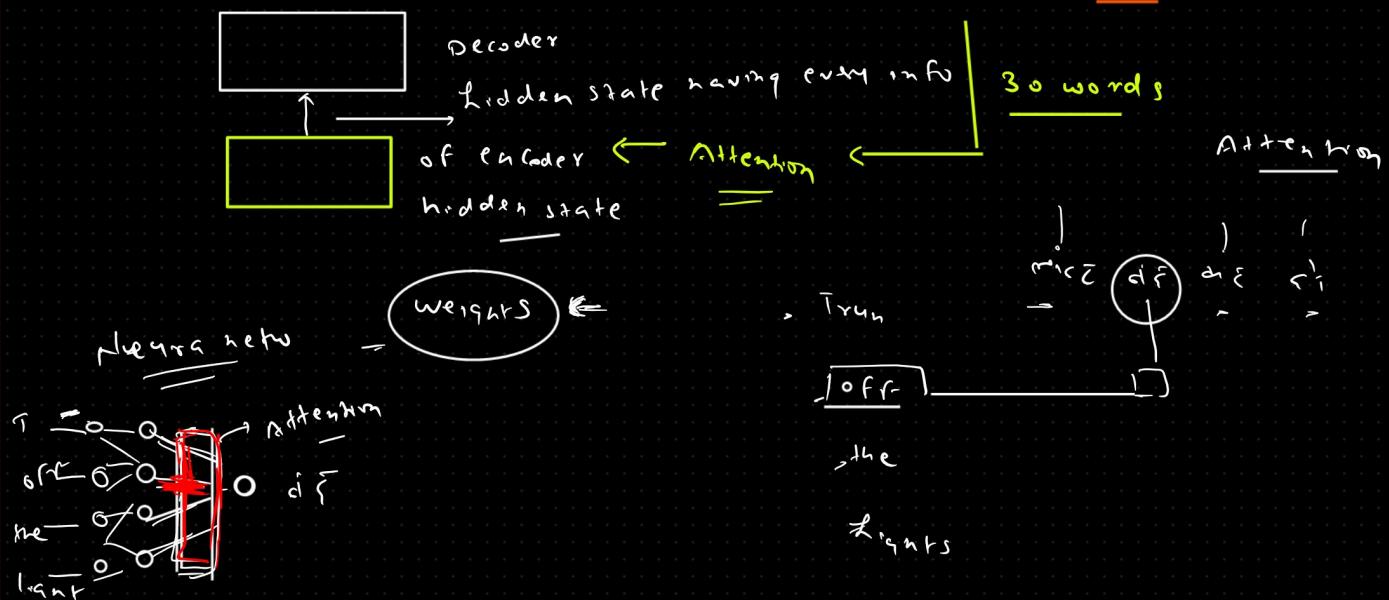
Google brain => attention is all you need

$$\{R_{NN}, L_{SM}, G_{RN}\}$$

Preregressing the sequence

RNN \Rightarrow LSTM \Rightarrow GPT \Rightarrow {Sequence to sequence mapping}

2014



Transformer => RNN, LSTM, GRU

The diagram illustrates the architecture of a transformer layer. It starts with an input labeled "embedding". This is followed by a red circle containing the word "Sentence". An arrow points from "embedding" to the red circle. From the red circle, an arrow points to the right, leading to two stacked boxes labeled "Self Attention". Finally, another arrow points from "Self Attention" to the right, leading to two stacked boxes labeled "Linear".

—

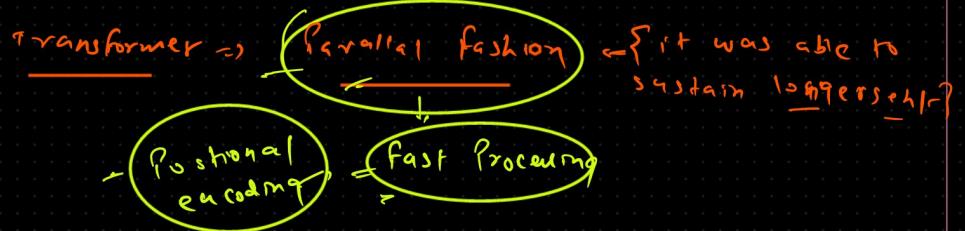
1084 109

feedforward

{ multi headed }

sequential \Rightarrow RNN, LSTM, GRU

timestamp



Datetime \Rightarrow sequential \Rightarrow process \Rightarrow transformer

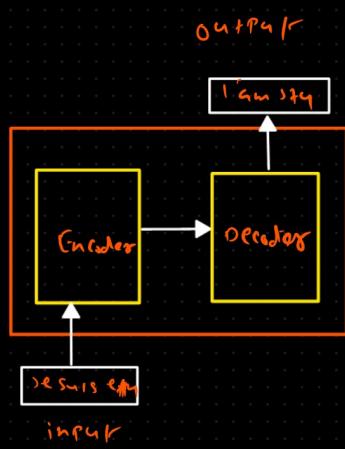
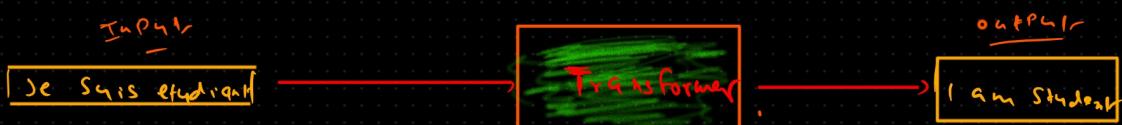
Encoder Decoder \Rightarrow this was the original

encoder decoder

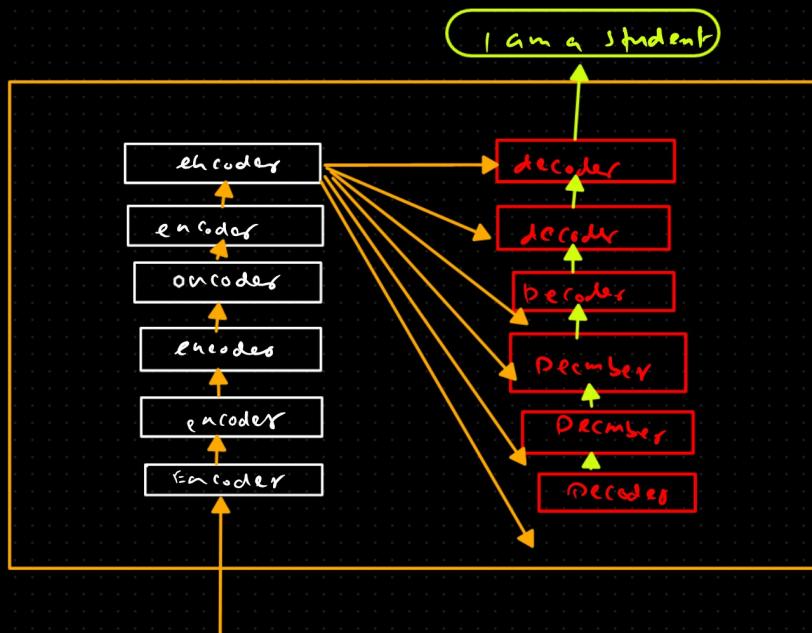
Encoder & Decoder \downarrow
Stack

- ① Architecture
 - ② Attention
 - ③ Multi-head attention
 - ④ Embedding
 - ⑤ Positional encoding
 - ⑥ Softmax
 - ⑦ {training, evaluation, decoding}
 - ⑧ GPU
- {Some interviews \Rightarrow train \rightarrow transform}

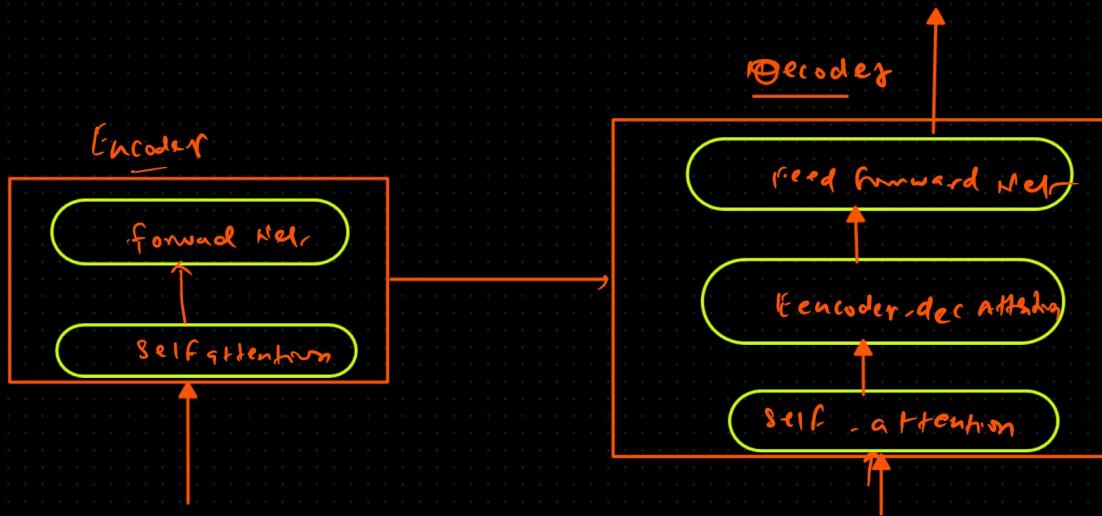
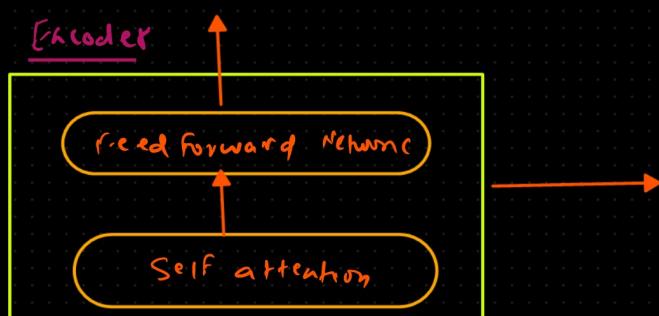
Machine translation



Research Paper

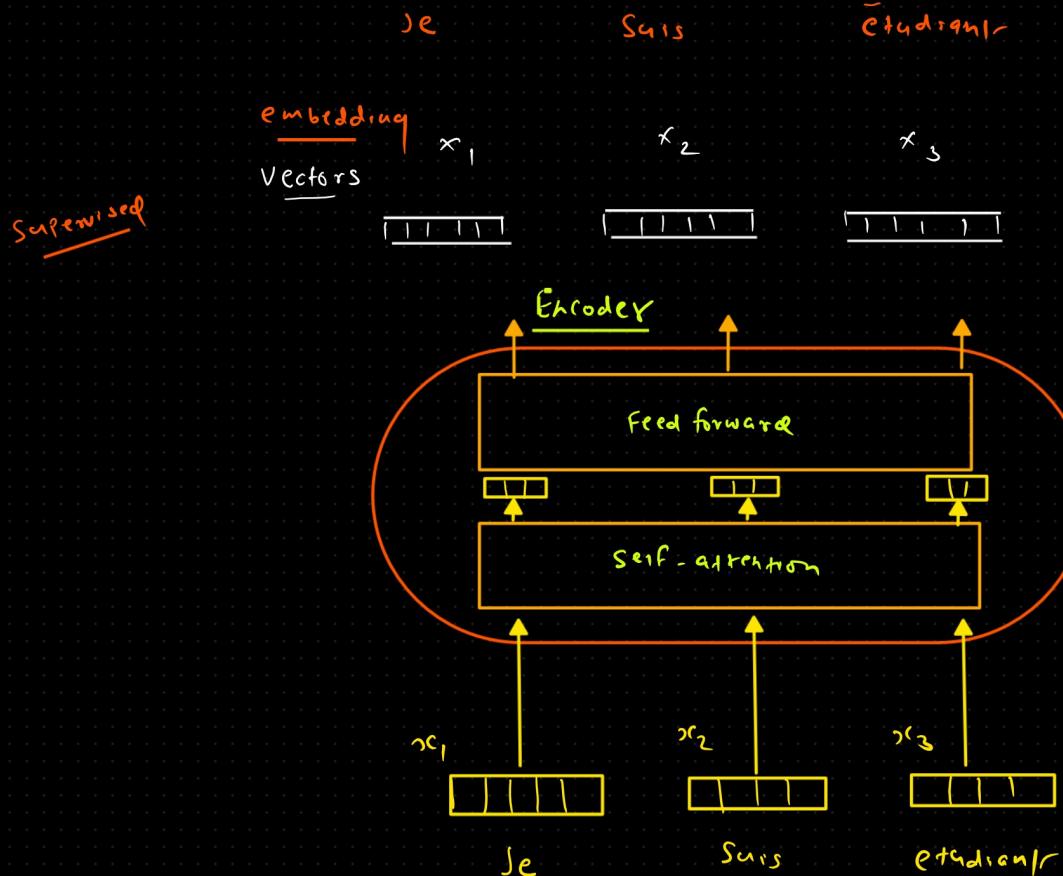


Je suis étudiant

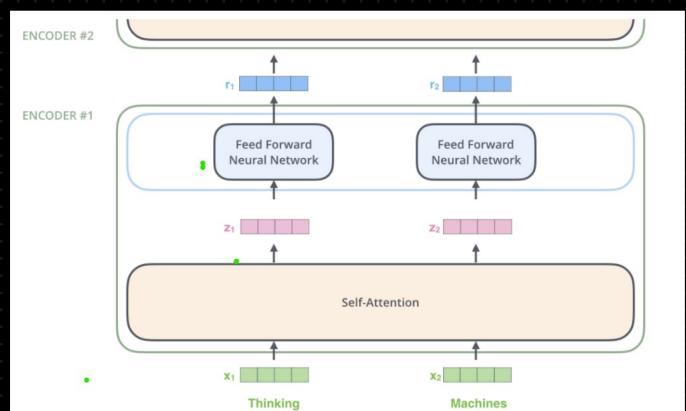


Sentence

french sentence



Self attention



Attention

one language

translation

another language

turn off the lights

မြန်မာ အင်ဂျင်ဘ်

Connection

Input —————> Output

decoder hidden state

turn

မြန်မာ အင်ဂျင်ဘ်

Attention matrix

off

မြန်မာ အင်ဂျင်ဘ်

encoder hidden state

the

မြန်မာ အင်ဂျင်ဘ်

light

မြန်မာ အင်ဂျင်ဘ်

Self-attention

Retence =) The animal didn't cross the street because it was too tired

Attention

Sequence

(RNN LSTM)

Processing

$t = 1 \times 2 \times 3 \dots$

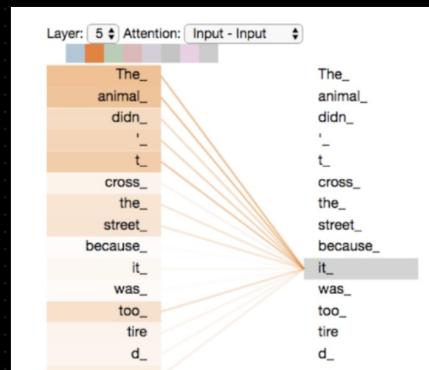
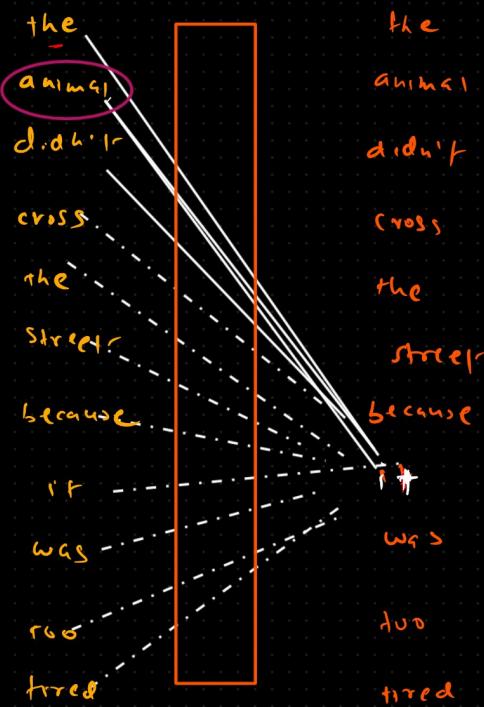
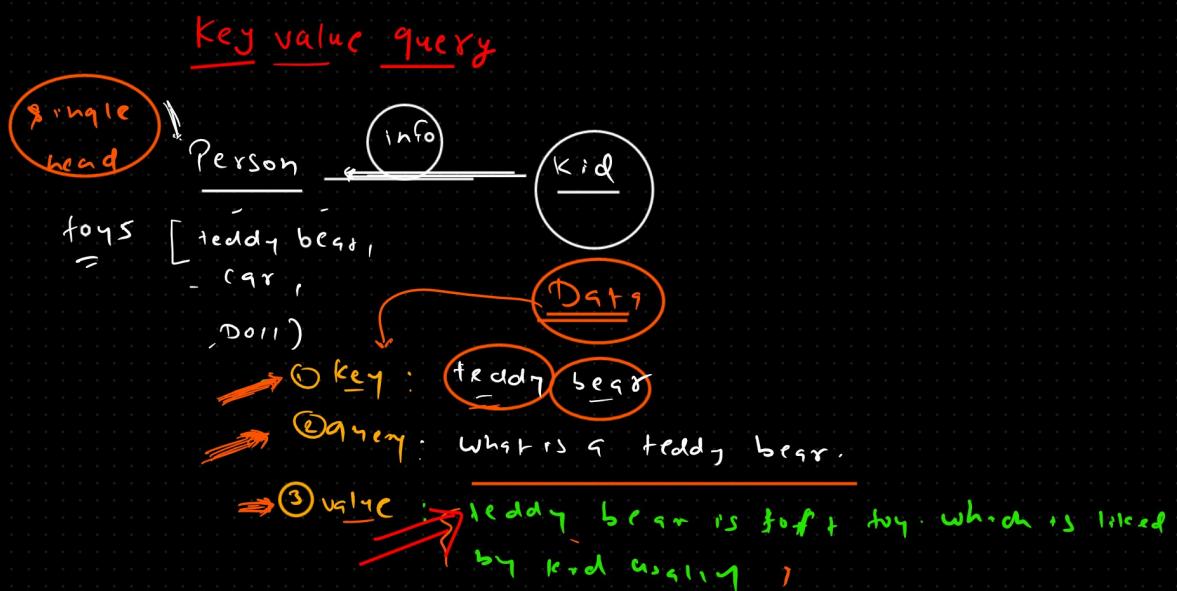
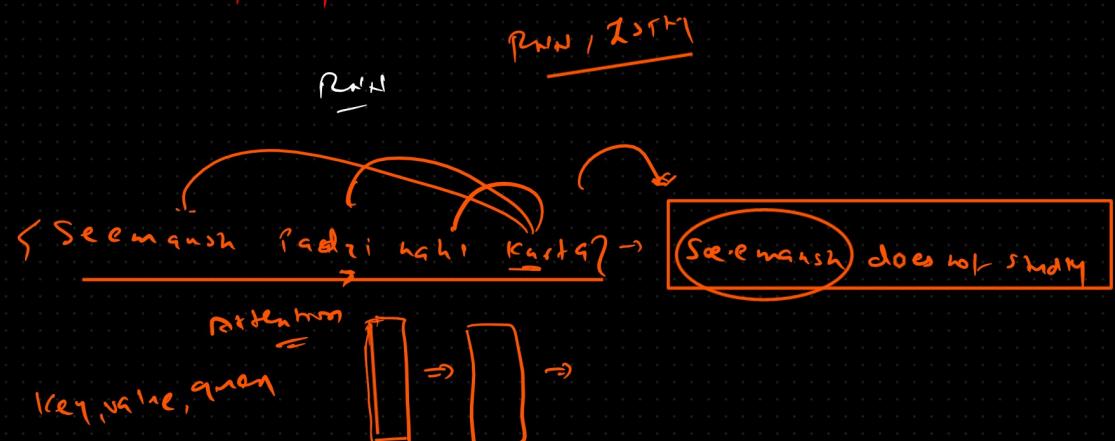


Figure 4: Two attention heads also in Layer 5 of 6, apparently involved in n-gram resolution. Top: Full attention for beam 5. Bottom: Biased attention from just the word 'its' for attention heads 5 and 6. Note that the attentions are very sharp for this word.



$\{ \text{key}, \text{query}, \text{value} \}$ → reversing the context
Self attention → decoding the sequence
 Himanshu got go-i. in CBSE Exam because he is very
hardworking boy



Single Person } Attention
one head

Multiple Person } multihead
multi-head

Attention all you need

