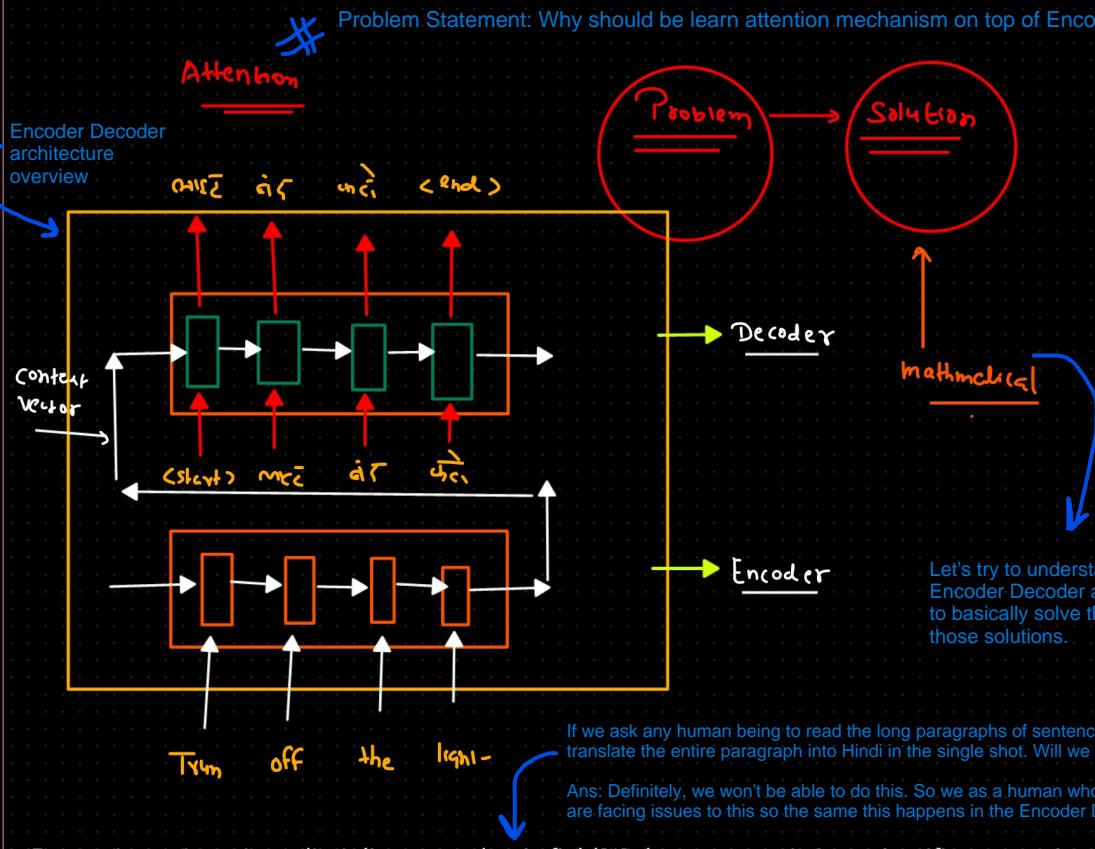


Attention Mechanism

Encoder-decoder



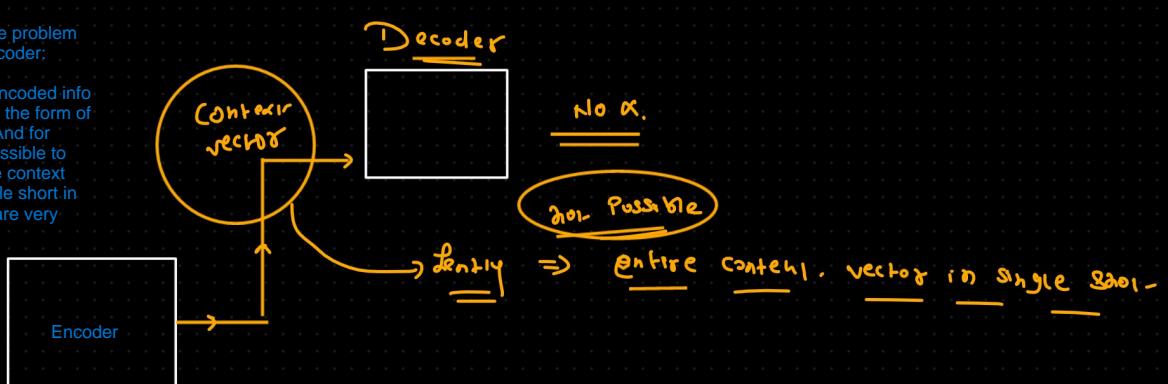
Data science is an interdisciplinary academic field[1] that uses statistics, scientific computing, scientific methods, processes, algorithms and systems to extract or extrapolate knowledge and insights from potentially noisy, structured, or unstructured data.[2]

Data science also integrates domain knowledge from the underlying application domain (e.g., natural sciences, information technology, and medicine).[3] Data science is multifaceted and can be described as a science, a research paradigm, a research method, a discipline, a workflow, and a profession.[4]

Data science is a "concept to unify statistics, data analysis, informatics, and their related methods" to "understand and analyze actual phenomena" with data.[5] It uses techniques and theories drawn from many fields within the context of mathematics, statistics, computer science, information science, and domain knowledge.[6] However, data science is different from computer science and information science. Turing Award winner Jim Gray imagined data science as a "fourth paradigm" of science (empirical, theoretical, computational, and now data-driven) and asserted that "everything about science is changing because of the impact of information technology" and the data deluge.[7][8]

To summarize the problem with Encoder-Decoder:

Encoder sends encoded info to the Decoder in the form of Context Vector. And for Decoder it not possible to decode the entire context vector in the single short in case sentences are very lengthy



Problem statements \Rightarrow we cannot translate longer sentences



Since in Encoder Decoder we were not able to translate/decode the longer sentences in the single short hence, Attention mechanism was introduced.

~~#~~ = Attention \Rightarrow Real time human understanding $\xrightarrow{\text{Focus/blur}}$

Data mining is a particular data analysis technique that focuses on statistical modeling and knowledge discovery for predictive rather than purely descriptive purposes, while business intelligence covers data analysis that relies heavily on aggregation, focusing mainly on business information.[4] In statistical applications, data analysis can be divided into descriptive statistics, exploratory data analysis (EDA), and confirmatory data analysis (CDA).[5] EDA focuses on discovering new features in the data while CDA focuses on confirming or falsifying existing hypotheses.[6][7] Predictive analytics focuses on the application of statistical models for predictive forecasting or classification, while text analytics applies statistical, linguistic, and structural techniques to extract and classify information from textual sources, a species of unstructured data. All of the above are varieties of data analysis.[8]

FOCUS and **BLUR** are 2 main things which we embed in the architecture due to which real time human understanding is achieved. This means when we are reading a part of huge paragraph then that part will be in **FOCUS** and rest entire sentence will be **BLUR**

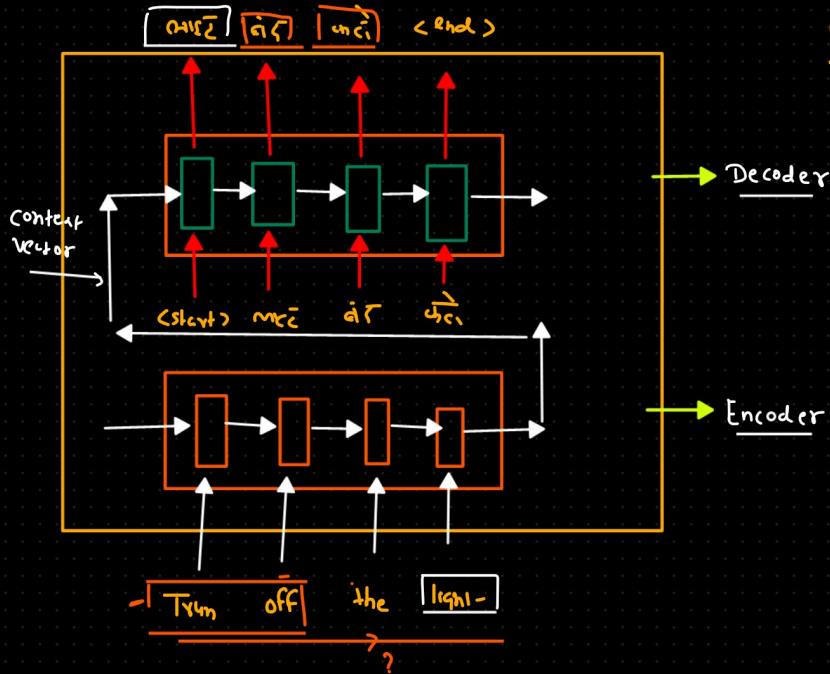
 **Focus** \Rightarrow Whatever part we are reading we are focusing on that.

Bur \Rightarrow Real Part - will be bLgR

↳ Mathematics



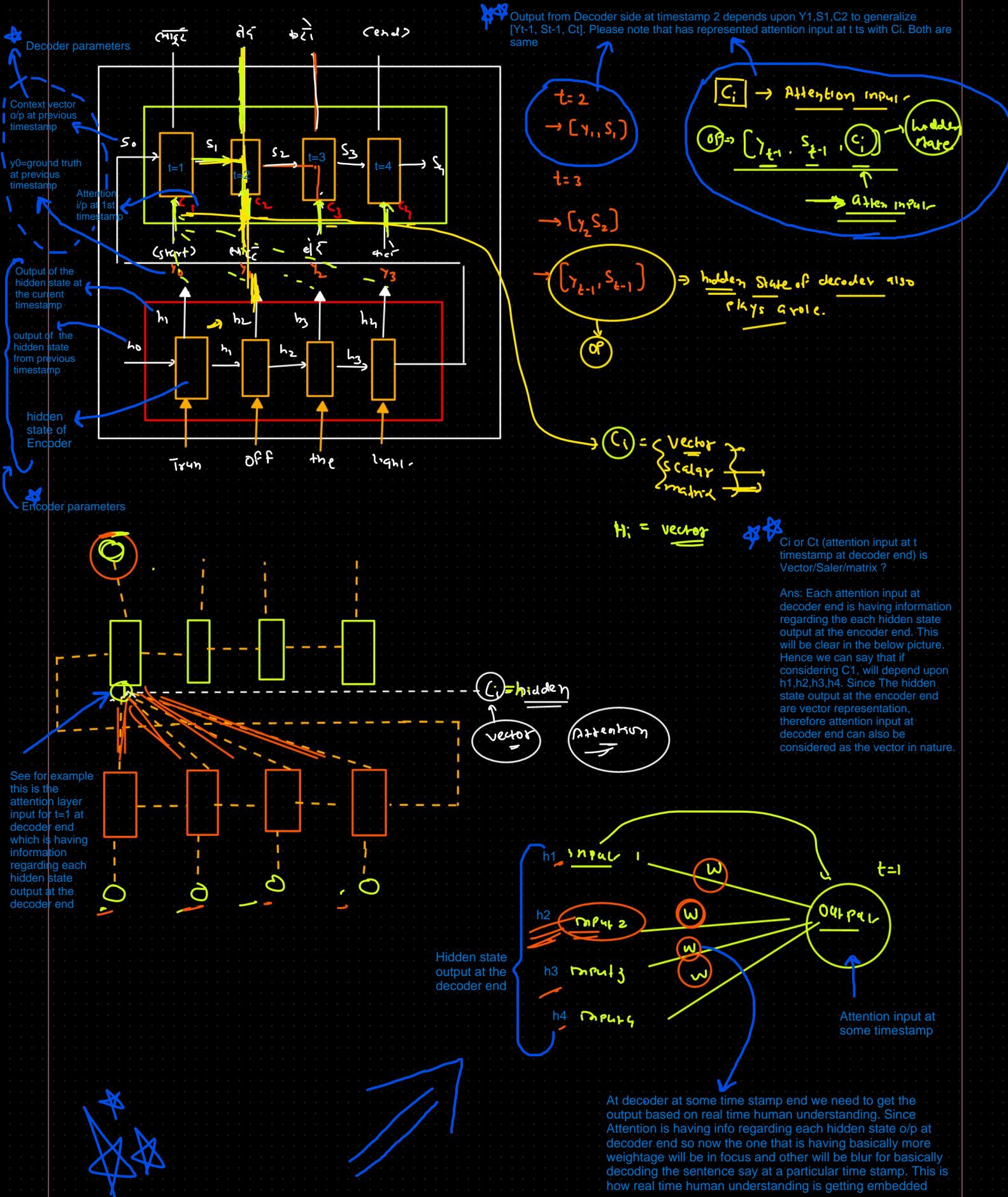
四



Example $\Rightarrow \{ \text{Donor Required entries} \\ \text{seen-} \}$

If we want the Hindi translation of the first two English inputs (ie; Turn off) in the single short then we will FOCUS just on first 2 input at encoder side and 2nd and 3rd output at the decoder side and will BLUR the remaining things

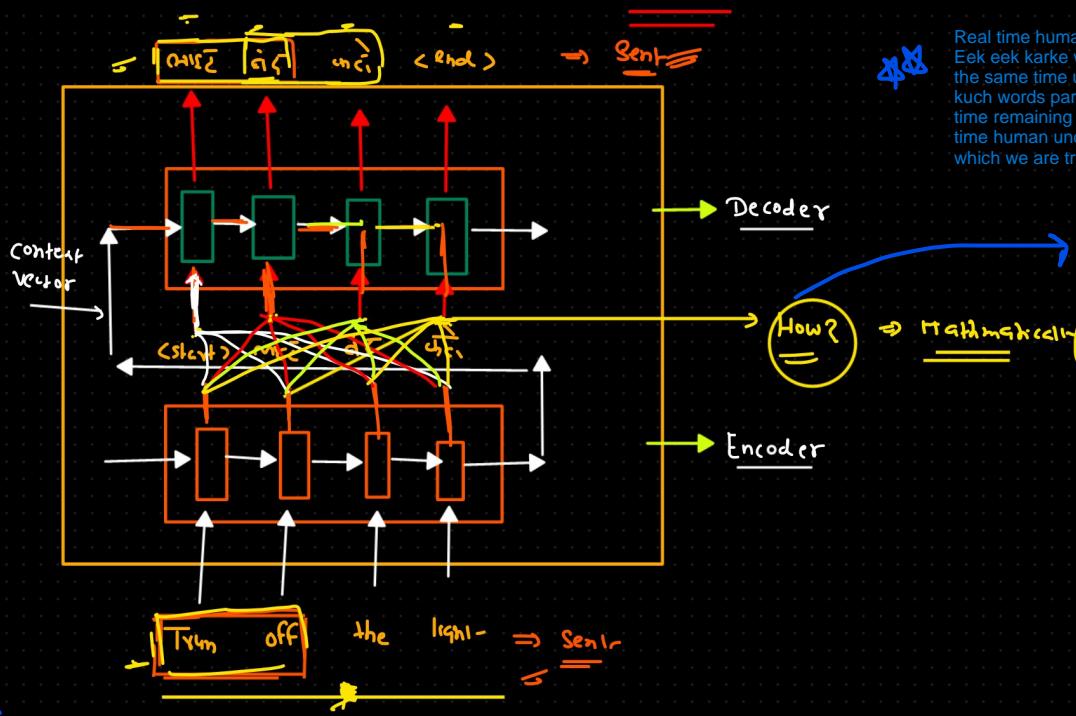
Previous hidden state input and the current hidden state output



Understanding the attention mechanism using neural network:

To understand this first read the below page first, especially the Tea making example-->Good output (Tea) and Bad output, If tea is bad then back propagate and adjust the weights same here also we are feeding the context vector from previous ts followed by ground truth from previous ts followed by attention input after which respective output from decoder will be checked for Good or bad output and accordingly we will be adjusting the weights as discussed above

To summarize attention at a particular timestamp has output from each hidden state but whom it is going to give more weightage to get the respective decoded output at a particular timestamp is decided based on the above explanation.

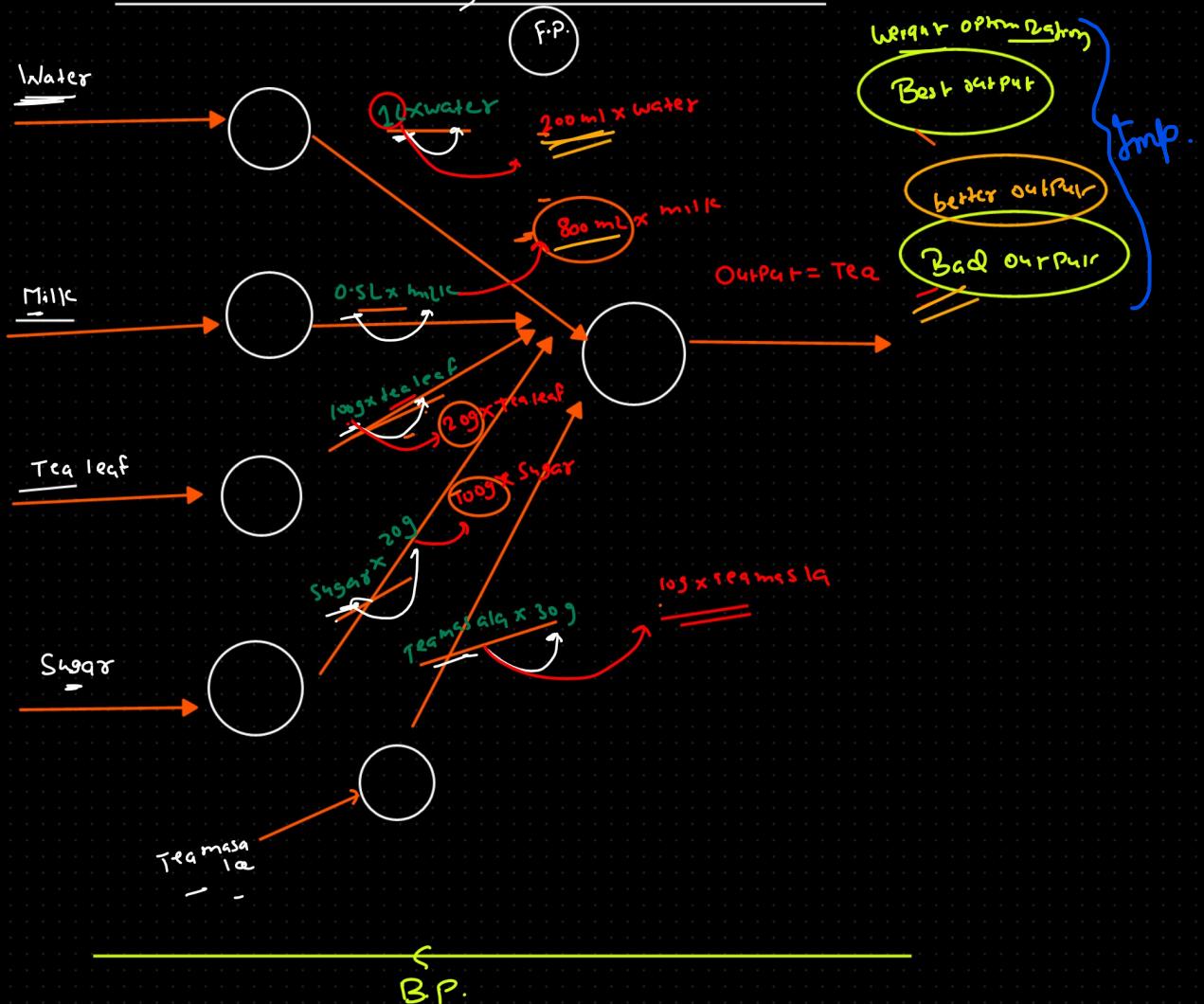


Real time human understanding:
Eek eek karke words padha(FOCUS) and at the same time unka hindi translation kya. Jab kuch words par FOCUS kar rahe then at that time remaining ko BLUR kar diya. This real time human understanding is something which we are trying to embed using attentions

If we want to translate the first 2 words in our sentence ie; "Turn off" then we have to focus of 2nd and 3rd words generated from he decoder side. How we are going to map this relation. that we will be able to do using the connected graph depicted here. But how will this be presented mathematically

WOW!
★ ★
★ ★

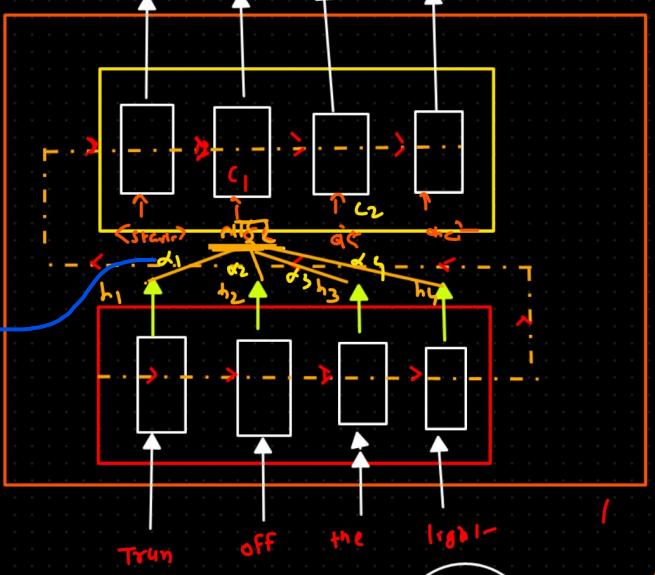
Understanding the Neural Network using Tea making example:



Mathematical Intuition of the Attention:

So, above it is clear that Focus and Blur or real time human understanding is embedded using weights b/w hidden state output from encoder and attention layer input from decoder. So below we will see how to determine this weight





$$c_1 = \alpha_{11} \times h_1 + \alpha_{12} \times h_2 + \alpha_{13} \times h_3 + \alpha_{14} \times h_4$$

$\alpha = \text{Weight}$

$$c_2 = \alpha_{21} \times h_1 + \alpha_{22} \times h_2 + \alpha_{23} \times h_3 + \alpha_{24} \times h_4$$

$$c_i = \sum \alpha_i \times h_i$$

Findout ..

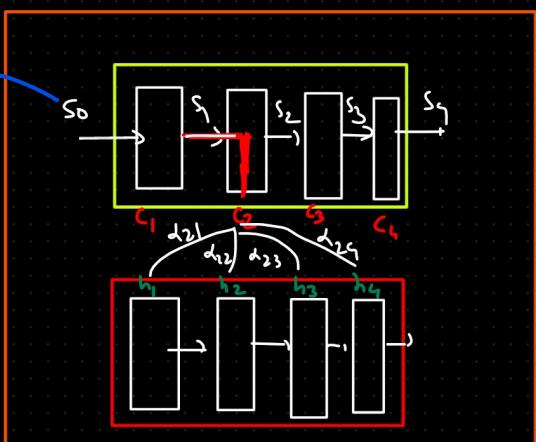
Generalized equation

Here my end goal to find out a weight

Imp.

Please note that S is representing the hidden state input or output at decoder end

Whereas, h represents the hidden state input or output at the encoder end



$$\alpha_{21}$$

h_1

$$\alpha_{21} \rightarrow f(h_1, s_1)$$

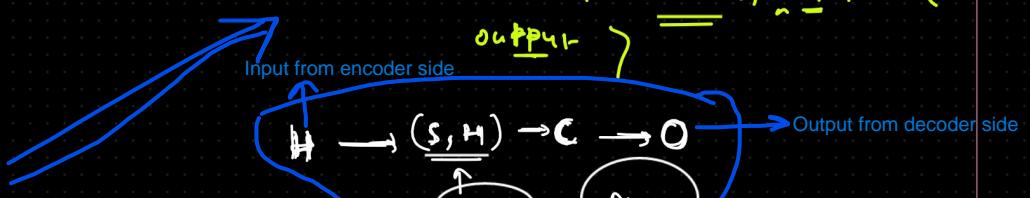
$$\alpha_{1j} \rightarrow f(h_j, s_{i-1})$$

ANN

The intended weights will be depending upon 2 things. Consider $(\alpha)_{21}$ will be depending upon h_1 and s_1

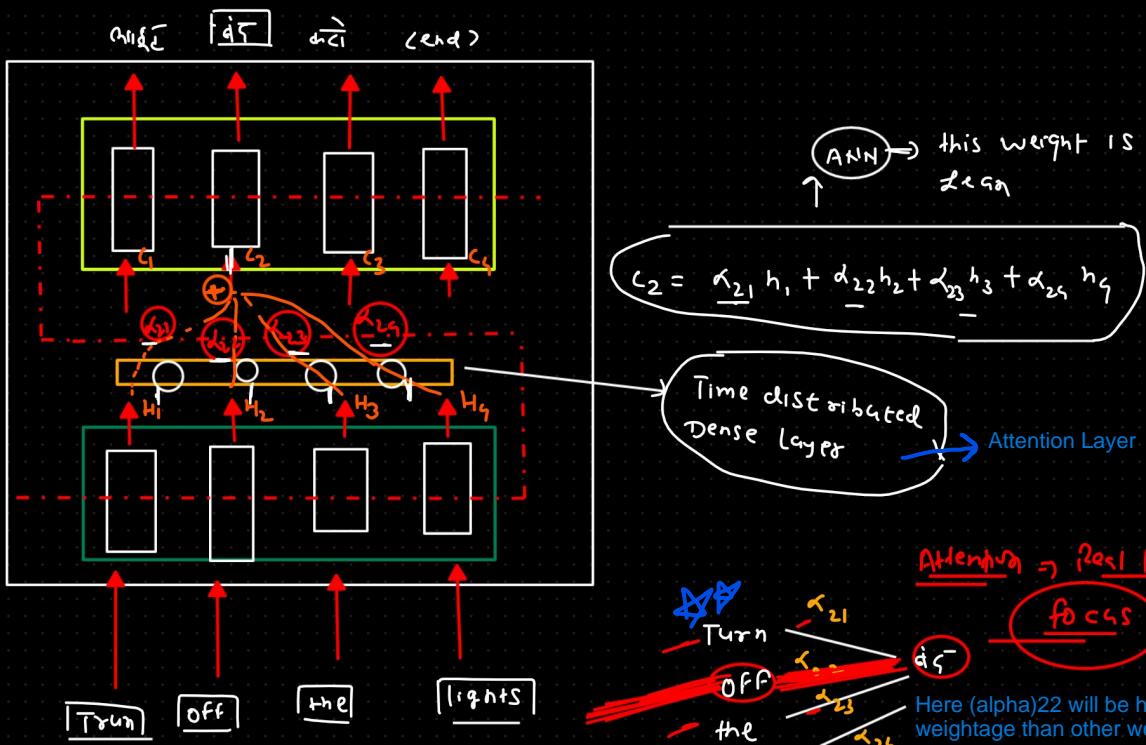
Attention layer \rightarrow ANN \rightarrow to Findout all

a weight between input and output



Imp.

In short to conclude whole discussion until now we can say that Attention is the ANN which is used to determine the weights b/w all inputs(hidden state output at Encoder side) and output(decoder end)



Attention \rightarrow Real time gate

Focus

Here (α_{22}) will be having more weightage than other weights. This is how real time human understanding is embedded where at some instant/timestamp we are able to FOCUS and decode a English word into Hindi by keeping others as BLUR.

Attention layer / Attention mechanism from Research Paper

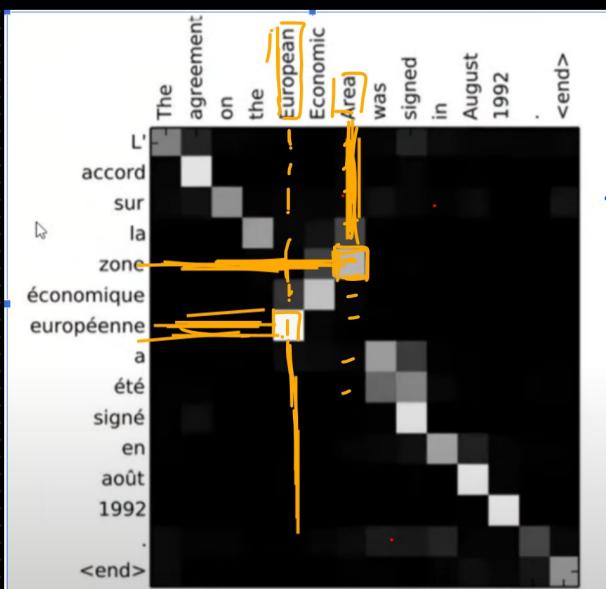
RNN \star Real e) English \rightarrow French

Bi-Directional RNN

In Attention Research paper we are taking the bidirectional RNN. In our previous explanation (English to Hindi) we took simple RNN

Weight
Heatmap

Plotting heatmap that will help to visualize the important weights for a given combination of translation



Relate this heatmap with the 2d weight matrix

$$\begin{matrix} \text{mice} & \text{turn} & \text{on} & \text{the} & \text{lions} \\ \alpha_{11} & \alpha_{12} & \alpha_{13} & \alpha_{14} & \end{matrix}$$

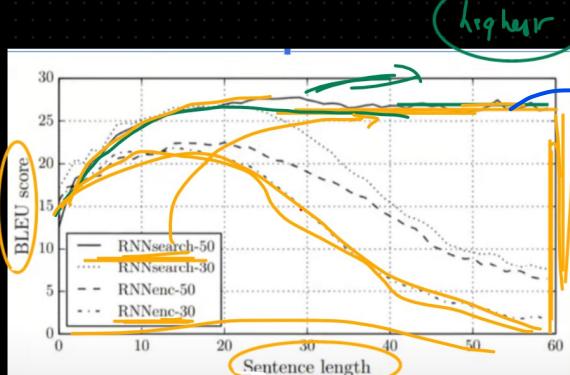
$$\begin{matrix} \text{pig} & \text{shy} & \text{as} & \text{big} & \end{matrix}$$

$$\begin{matrix} \text{red} & \alpha_{31} & \alpha_{32} & \alpha_{33} & \alpha_{34} \\ \alpha_{41} & \alpha_{42} & \alpha_{43} & \alpha_{44} & \end{matrix}$$

Focus
Blur

Real human learning

BLUE score comparison for different models wrt to sentence length as per Research Paper



higher

This curve represents Bidirectional RNN with Attention layer which is having the highest BLEU score

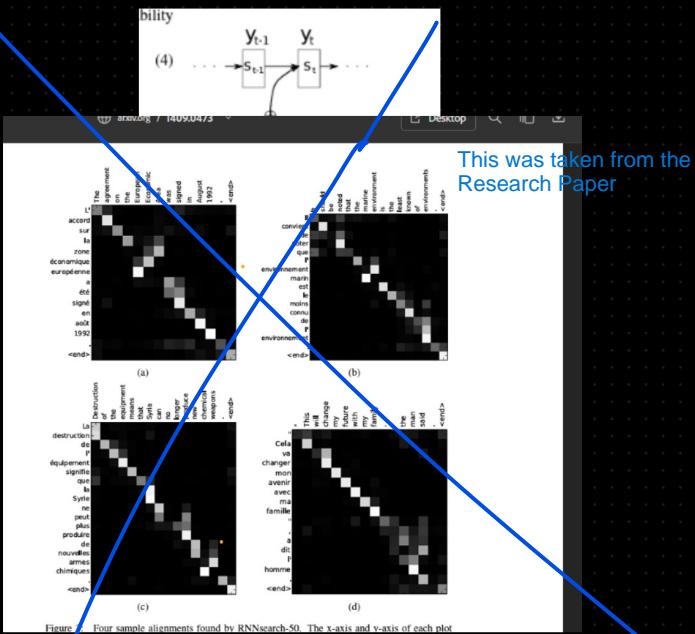


Figure: Four sample alignments found by RNNsearch-50. The x-axis and y-axis of each plot

Transform, Beor, gr-

Projects