# Encoder And Decoder
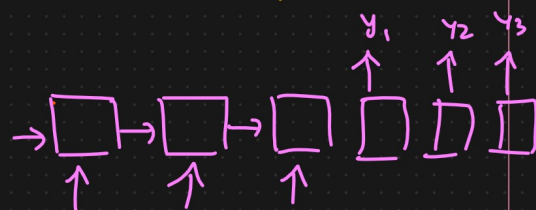
① Simple RNN → Vanishing Gradient Problem

② LSTM RNN →

③ GRU RNN → ⎫⎬⎭ Long Short Term Memory.

④ Bidirectional RNN ←

① Many to Many RNN



## Encoder And Decoder ⎬

Eg: One Language To Other

English → French
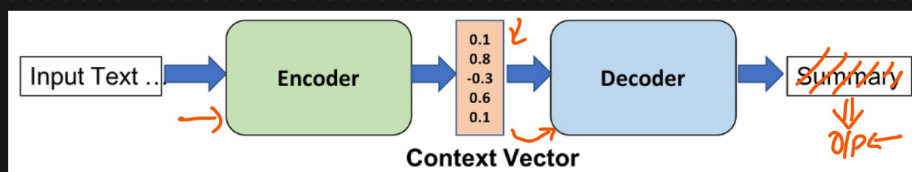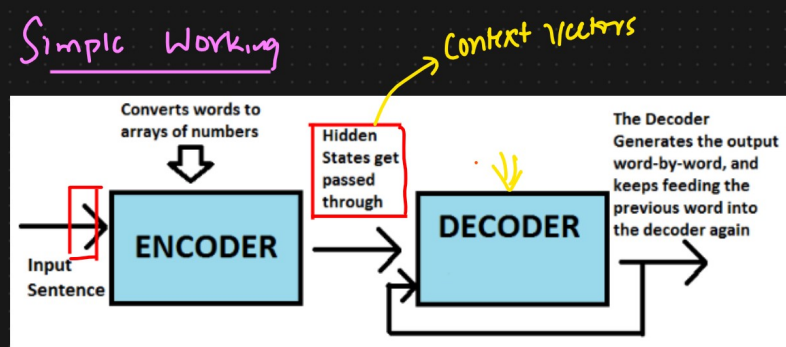
Eg: Linked Chat → Hi, How are you?

Sequences I/p          O/p Sequence Of Words

## Simple Working

Context Vectors



Converts words to arrays of numbers

Hidden States get passed through

**ENCODER**

Input Sentence

**DECODER**

The Decoder Generates the output word-by-word, and keeps feeding the previous word into the decoder again



Input Text ..    Encoder    | 0.1 |    Decoder    Summary
                            | 0.8 |
                            | -0.3 |              O/p
                            | 0.6 |
                            | 0.1 |
                    **Context Vector**

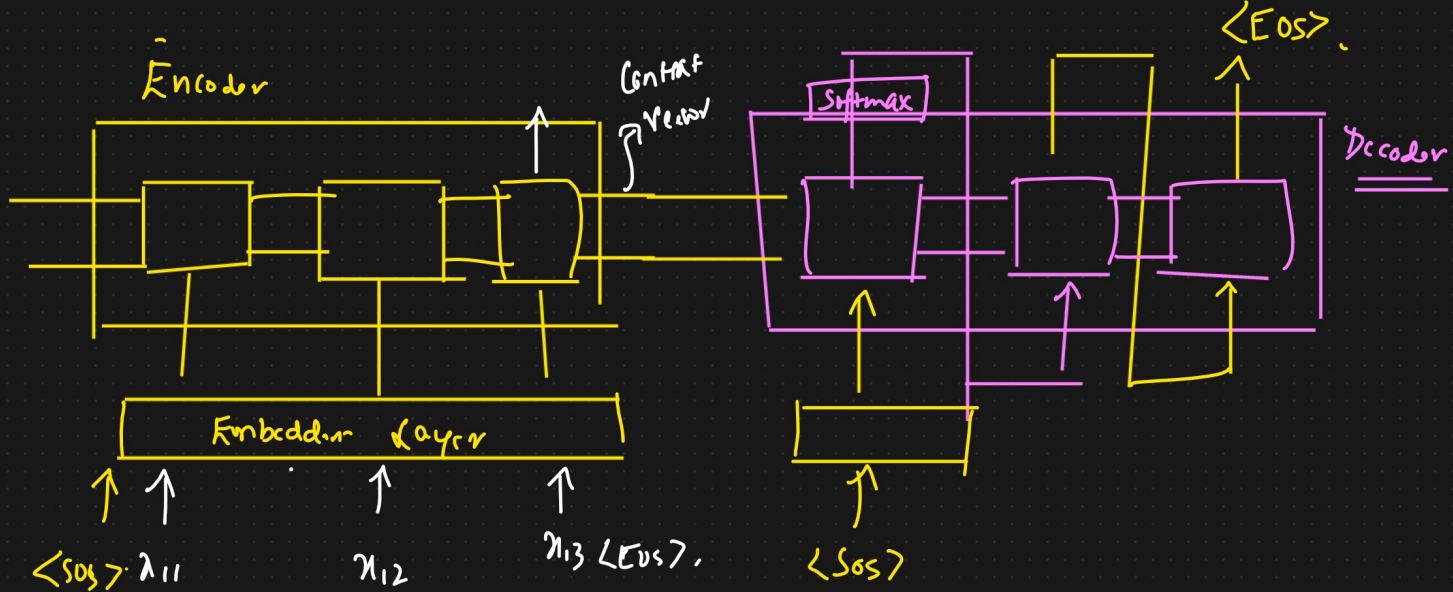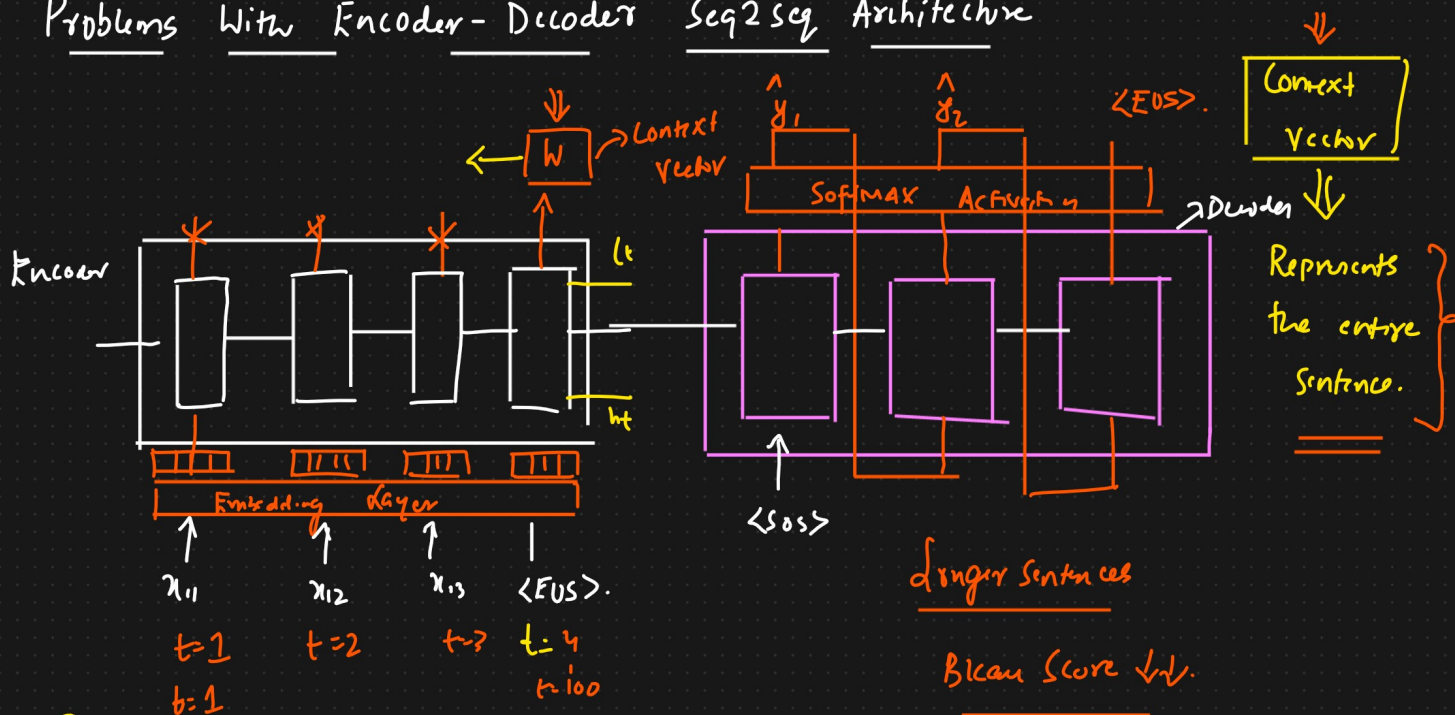① Encoder ⟹ I/p ⟹ Context Vector ⟸ Vectors ⎫
                                                ⎬
② Decoder ⟹ ↙ ⟹ O/p                            ⎭

Useage

① Language Translation

② Text Generation

③ Text Suggestion.

# RNN → Vanishing Gradient Problem

① forget gate
② I/p gate & candidate i/p
③ O/P

**Sequence-to-Sequence (seq2seq)**
**Encoder-Decoder Neural Network**



Dataset

Rue O/P    Y truth
True O/P    { Tracing }

English    French

Ŷ

$[y \, \hat{y}]_{L(x)}$
<EOS>

<SOS> Thank you <EOS>. Gracias → <SOS> GRACIAS <EOS>. → GRACIOUS    <EOS>.

Encoder

Context Vectors

Long Term

$\hat{y}$ [0.1, 0.6, 0.3]
[0.1, 0.3, 0.6]  → $\hat{y}$ [- - -]

Softmax    Softmax

$C_t$    $L_t$    Decoder

$h_6$    $h_t$

Short Term Memory

OHE  ← |1|0|0|  |0|1|0|0|  |0.6|1|0|  |0|0|0|1|

|1|0|0|  |0|1|0|

<SOS>  Thank  You  <EOS>. Memory    <SOS>    GRACIOUS

t=1    t=2    t=3    t=4 ←    t=1    t=2    t=3

Context Of this Sentence.

Loss $(y-\hat{y})^2$

Optimizer

$\hat{y}$ [- - -]  [- - -]  [- - -]
$y$  [ | | | ]  [ | | | ]  [ | | | ]
    [ - - - ]  [ - - - ]  [ - - - ]

## Encoder / Decoder Diagram

Encoder

Context Vector

Softmax

Decoder

$<EOS>$.

Embedding Layer

$<SOS>$ $x_{11}$    $x_{12}$    $x_{13}$ $<EOS>$.    $<SOS>$

---

# Problems With Encoder-Decoder Seq2Seq Architecture

Context Vector

$W$ → Context Vector

$\hat{y}_1$    $\hat{y}_2$    $<EOS>$.

Softmax Activation

Encoder

$l_t$

$h_t$

Embedding Layer

Decoder →

$<SOS>$

$x_{11}$    $x_{12}$    $x_{13}$    $<EOS>$.

$t=1$    $t=2$    $t=3$    $t=4$

$b=1$    $t=100$

Context Vector

Represents the entire Sentence.

Longer Sentences

Blean Score ↓↓.

Researchers : Sentences of varying length
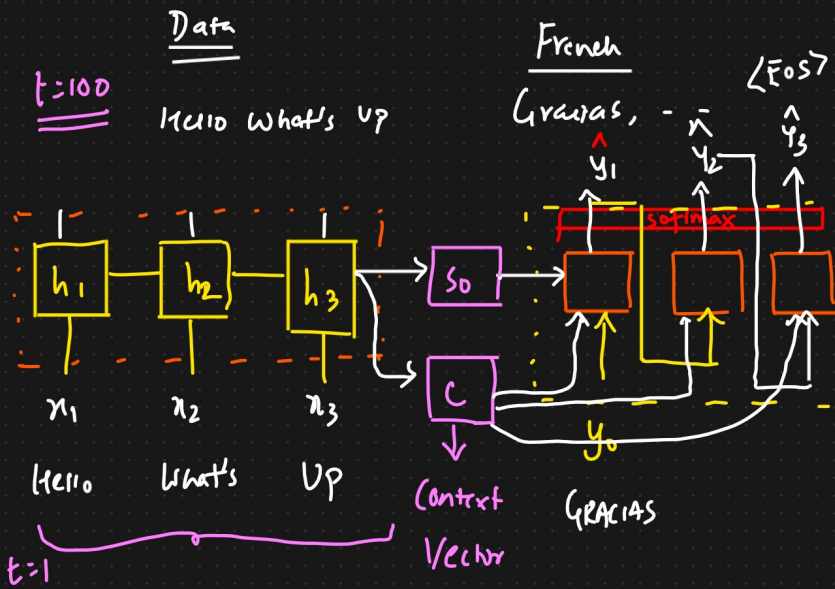
⇒ Seq to Seq Data

Bleu Score

30  40  50  60

Sentence Length

# ✳ Attention Mechanism → Seq2Seq Network

Longer paragraph → { Context Vector }

+

{ Context }

## Attention Mechanism | Seq2Seq Networks

Data

$t = 100$

Hello What's up

$h_1$   $h_2$   $h_3$   →   $S_0$

$x_1$   $x_2$   $x_3$

Hello   What's   Up

$t = 1$

Context Vector

French

Gracias, - - -

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$

⟨EOS⟩

softmax

$y_0$

GRACIAS

### Encoder Decoder Architecture

### Attention Mechanism

https://erdem.pl/2021/05/introduction-to-attention-mechanism

$h_1$   $h_2$   $h_3$   →   $S_0$

$\overleftarrow{h_1}$   $\overleftarrow{h_2}$   $\overleftarrow{h_3}$

Hello $t=3$
$t=1$

What's  $t=2$
$t=2$

Up  $t=1$
$t=3$

---

## 3 LEARNING TO ALIGN AND TRANSLATE

In this section, we propose a novel architecture for neural machine translation. The new architecture consists of a bidirectional RNN as an encoder (Sec. 3.2) and a decoder that emulates searching through a source sentence during decoding a translation (Sec. 3.1).

### 3.1 DECODER: GENERAL DESCRIPTION

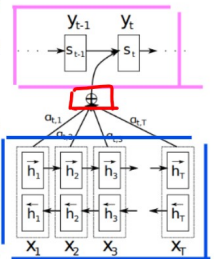In a new model architecture, we define each conditional probability in Eq. (2) as:

$$p(y_i|y_1, \ldots, y_{i-1}, \mathbf{x}) = g(y_{i-1}, s_i, c_i), \qquad (4)$$

where $s_i$ is an RNN hidden state for time $i$, computed by

$$s_i = f(s_{i-1}, y_{i-1}, c_i).$$

It should be noted that unlike the existing encoder–decoder approach (see Eq. (2)), here the probability is conditioned on a distinct context vector $c_i$ for each target word $y_i$.

The context vector $c_i$ depends on a sequence of *annotations* $(h_1, \cdots, h_{T_x})$ to which an encoder maps the input sentence. Each annotation $h_i$ contains information about the whole input sequence with a strong focus on the parts surrounding the $i$-th word of the input sequence. We explain in detail how the annotations are computed in the next section.

---

The context vector $c_i$ is, then, computed as a weighted sum of these annotations $h_i$:

$$\left\{ c_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j. \right\} \qquad (5$$

The weight $\alpha_{ij}$ of each annotation $h_j$ is computed by

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

where

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

$$[a_{11}, a_{12}, a_{13}]$$

Attention Weights

Compute Context Vector

$$C_t = \sum_{i=1}^{t} a_{t,i} h_i$$

⊗ ⊗ ⊗ ⊕

| $a_{1,1}$ | $a_{1,2}$ | $a_{1,3}$ |

SOFTMAX

{Feed Forward Neural Nlw}.

ANN

| $e_{1,1}$ | $e_{1,2}$ | $e_{1,3}$ |

Aligment Scores

$h_1$   $h_2$   $h_3$

$\hat{y}_1$   $\hat{y}_2$   $\hat{y}_3$

Softmax

$\overrightarrow{h_1}$   $\overrightarrow{h_2}$   $\overrightarrow{h_3}$

$S_0$   $S_1$   $S_2$   $S_3$

$\overleftarrow{h_1}$   $\overleftarrow{h_2}$   $\overleftarrow{h_3}$

$C_1$

$y_0$   $y_1$   $y_2$

Encoder

$C_2$

Hello t=3    What's t=2    Up t=1
t=1          t=2           t=3

$C_3$

Hello , What's Up

attention mechanism