# Recurrent Neural N/w's
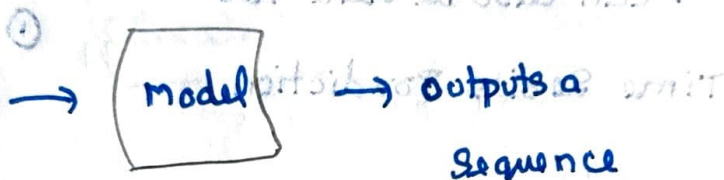
⤷ feed forward neural N/w's
rolled out over time

(Deal with Sequencedata
where the i/p have some
defined ordering)

## Several types of Architecture
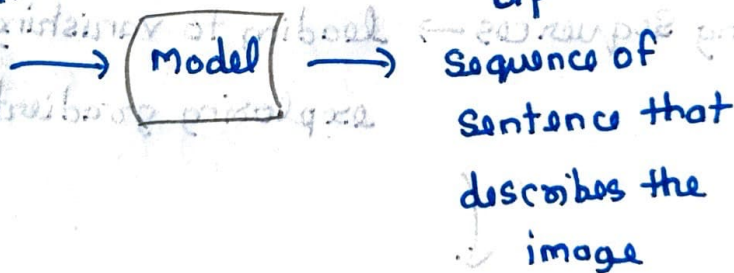
### ① vector - Sequence models

(Take in a fixed
size of vector
as i/p)

Model → outputs a Sequence

Eg → Image Captioning

Vector Representation
of an image
(i/p)

→ Model → o/p Sequence of Sentence that describes the image

### ② Sequence - Vector models

Take in a Sequence as i/p → o/p fixed length vector

The main character
Sucked

→ $\begin{bmatrix} 0.11 \\ 0.89 \end{bmatrix}$ → good ⎱ The person
bad ⎰ taught about the movie

③ Sequence - Sequence models

(Language Translation)

Takes in a
Sequence as           →        o/p's another
i/p                             Sequence

(Sentence in          →        (Sentence in
Spanish)                        English)

RNN's → can also be used for
    Time Series Prediction

Disadvantages

① Slow to train

② Long Sequences → leading to vanishing/
                   exploding gradients

To overcome this problem they
used LSTM N/w to introduce
long memory to the N/w

(This allows us to retain memory
for longer sequences)

we need the i/p of
the previous state in
order to perform any operation
on the current state)

① Normal RNN's
        ↓
      Slow

② LSTM's
        ↓
   Even slower

(i/p data has to
be passed
sequentially/
Serially one
after the
other

        ↓

This doesn't
make use of
GPU's which
can compute
llely)

(Adv of Transfoomers)

⌡

① Sequences can be passed in ||el

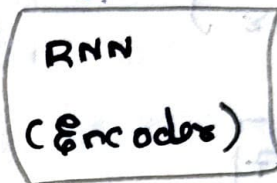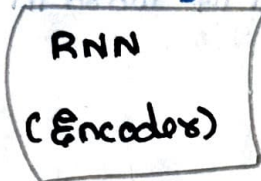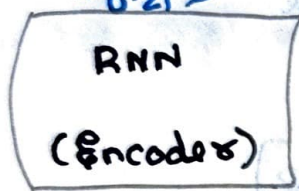$$\begin{bmatrix} 0.44 \\ 0.0 \\ 0.08 \end{bmatrix}$$

Using RNN

we pass the words Sequentially
(under a time Sequence)

$$\begin{bmatrix} 0.76 \\ 0.23 \\ 0.11 \\ 0.21 \end{bmatrix}$$
$$\begin{bmatrix} 0.32 \\ 0.88 \\ 0.41 \\ 0.14 \end{bmatrix}$$
$$\begin{bmatrix} 0.32 \\ 0.19 \\ 0.91 \\ 0.01 \end{bmatrix}$$

| RNN (Encoder) | → | RNN (Encoder) | → | RNN (Encoder) |

↑ The          ↑ red          ↑ dog

- - - - - - - - - - - - - - - - - - - - → TIME

The →
red →  | Transfoomer (Encoder) |  → → →   (||el vector
dog →                                      o/p's)

(input Embeddings)

mapping →
every words
to a space

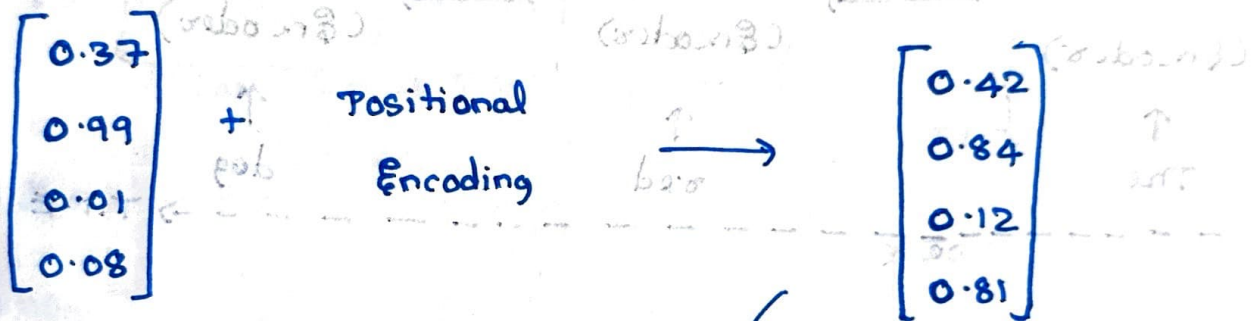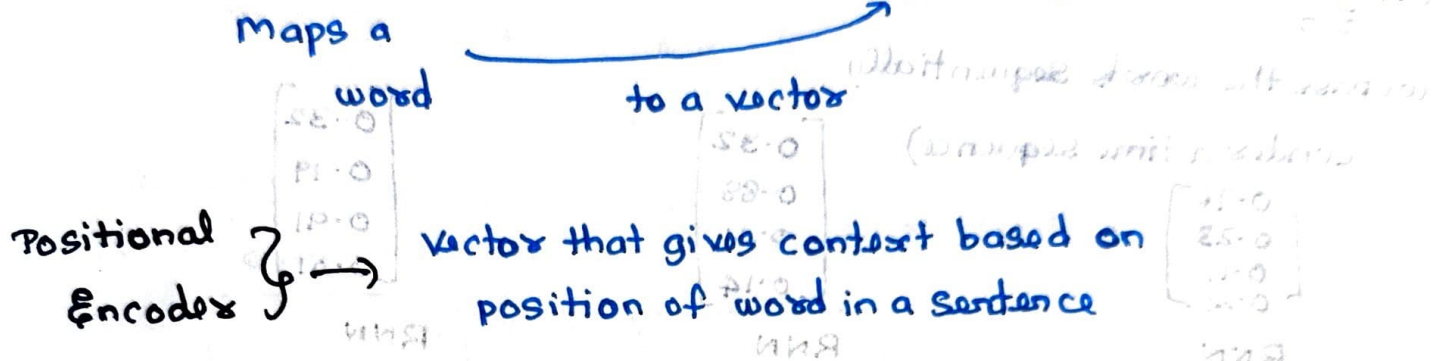| Sun          dog  horse |
|              cat        |
| red                     |

Embedding Space
(we can also
② pre train this
space to save
time)

(where words with similar
meaning are positioned
close)

Dog $\longrightarrow$ Embedding Space $\longrightarrow$ $\begin{bmatrix} 0.37 \\ 0.99 \\ 0.01 \\ 0.08 \end{bmatrix}$

Maps a word $\underrightarrow{\hspace{2cm}}$ to a vector

Positional Encoder $\Big\} \longrightarrow$ Vector that gives context based on position of word in a sentence

$\begin{bmatrix} 0.37 \\ 0.99 \\ 0.01 \\ 0.08 \end{bmatrix}$ + Positional Encoding $\longrightarrow$ $\begin{bmatrix} 0.42 \\ 0.84 \\ 0.12 \\ 0.81 \end{bmatrix}$

(word vector with positional info)

Embedding of Dog (with context info)

we pass this to the Encoder block

① it goes thro'

Ⓐ multi headed attention

Ⓑ feed forward layer

(to tranform the Attention vectors into a form that is digestible by the Encoder block)

# Attention

↳ This tells us on what part of the i/p we must focus on ?

(how relevant is the $i^{th}$ word in a sentence w.r.t other words in the Sentence)

> how relevant is the $i^{th}$ word in a sentence w.r.t other words in a Sentence

→ represented in the → $i^{th}$ **Attention Vector**

(captures the textual relationship b/w words in a Sentence)

The → ~~The~~ big red dog  $[0.71 \quad 0.04 \quad 0.07 \quad 0.18]^T$

big → The ~~big~~ red dog  $[0.01 \quad 0.84 \quad 0.02 \quad 0.13]^T$

red → The big ~~red~~ dog  $[0.09 \quad 0.05 \quad 0.62 \quad 0.24]^T$

dog → The big red ~~dog~~  $[0.03 \quad 0.03 \quad 0.03 \quad 0.91]^T$

→ 8 Soch attention Vectors (and then we compute the weighted average)

here we get the vector

# Decoder

we get the french word

Le → Convert this into vector

(input embedding)

→ Positional Embedding → Pass it to the Decoder block

> Each of the attention nets are independent of each other because of which we can pass them ||ely

**Decoder** ──→ generates vector for every word in a sentence

( Self attention is done for the french words)

Le ──→ ~~Le~~ gros chien rouge

gros ──→ Le ~~gros~~ chien rouge

chien ──→ Le gros ~~chien~~ rouge

rouge ──→ Le gros chien ~~rouge~~

(Attention Vectors)

$$\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0.1 \\ 0.9 \\ 0 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0.05 \\ 0.40 \\ 0.55 \\ 0 \end{bmatrix} \quad \begin{bmatrix} 0.16 \\ 0.09 \\ 0.15 \\ 0.66 \end{bmatrix}$$

Le      gros      chien      rouge

(Another Attention block)

**Encoder-Decoder Attention**

$$\begin{bmatrix} 0.71 \\ 0.04 \\ 0.07 \\ 0.18 \end{bmatrix} \quad \begin{bmatrix} 0.01 \\ 0.84 \\ 0.02 \\ 0.13 \end{bmatrix} \quad \begin{bmatrix} 0.09 \\ 0.05 \\ 0.62 \\ 0.24 \end{bmatrix} \quad \begin{bmatrix} 0.03 \\ 0.03 \\ 0.03 \\ 0.91 \end{bmatrix}$$

The      big      red      dog

(Vectors from Encoder)

(how related each word vector is similar with respected to each other)

O/p of the block is the Attention vectors for every word in English & French Sentence

( Each vector represents a relationship with other vectors in both the languages

( Then each attention Vector is passed
to a feed forward Neural N/w in
order to make it more digestable by
the next layer

① Linear ——————→ This has another
                      feed forward layer
② Softmax ....)
                      (used to expand the dim
        ↓              to the no. of words in
Converts it into           french language)
a prob dist
     ↓
final word is the one
with the highest prob


masked
   Attention ——→   (while generating the next french word
   block              we can use all the words from the
                      english Sentence but only the previous
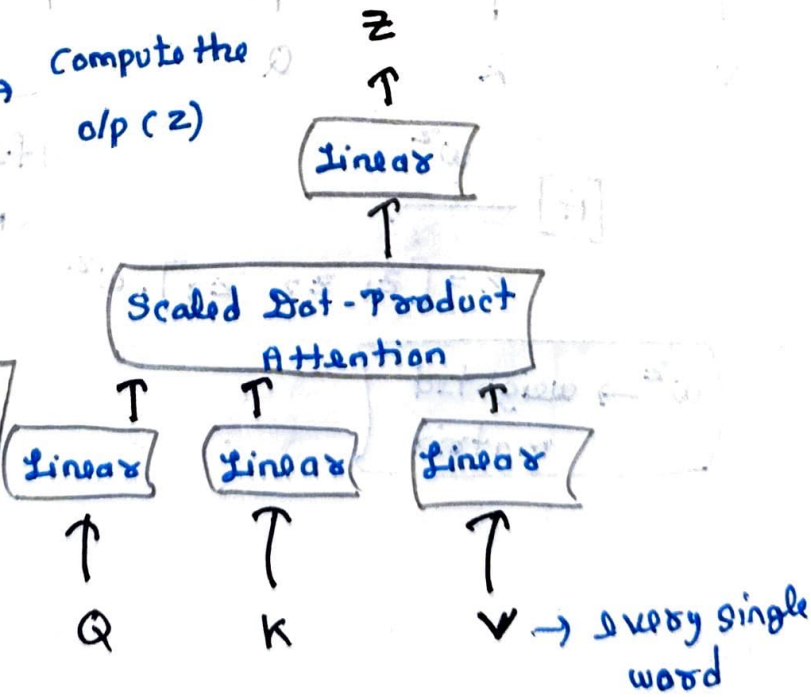                      word from the french Sentence)

we have  Q⎫
         K ⎬ vectors for    ——→  Compute the            Z
         V ⎭ every single         o/p (z)                ↑
              word                                    ⎧Linear⎫
                                                         ↑
   (Single headed                               ┌─────────────────┐
      Attention) ——————→                        │ Scaled Dot-Product│
                                                │    Attention     │
┌──────────────────────────────┐               └─────────────────┘
│                   Q·K^T        │                ↑      ↑       ↑
│ Z = Softmax (───────────────)·V│            ⎧Linear⎫⎧Linear⎫⎧Linear⎫
│                √dimension of    │               ↑      ↑       ↑
│                 vector Q, K or v│               ∫      ∫       ∫
↓ └──────────────────────────────┘               ↑      ↑       ↑
Compute                                           Q      K       V ——→ every single
attention                                                              word
vectors for
every word

# for multi headed Attention

### (we have multiple weight matrices)

$z_1$    $z_2$    $z_3$

[A]

we have multiple attention vectors ($z$) for every word

↑

Linear

↑

Concat

↑

Scaled Dot-Product Attention

↑        ↑        ↑

Linear    Linear    Linear

↑        ↑        ↑

$\omega^v$    $\omega^k$    $\omega^Q$

V        k        Q

$\omega^z$

[A] ———→

$$z = [z_1 \ z_2 \ z_3] \cdot \omega^z$$

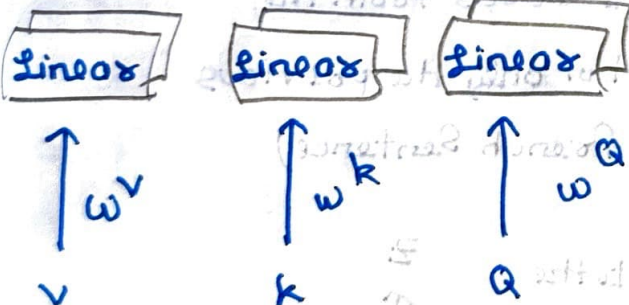$\omega^z$ → weighted matrix

Same

z formula

After each and every layer we perform some kind of normalisation

( batch normalisation

Layer normalisation)

↳ Read about this

feed forward Neural N/w

makes sure it is only one ($z$) per word attention vector