# ATTENTION IS ALL YOU NEED PAPER EXPLANATION AND IMPLEMENTATION FROM SCRATCH

## TABLE OF CONTENTS

RNN



I like eating apple but I like using moble phones from the apple company more.

$$\{ apple : [ \; . \; - \; . \; - \; . \; ]$$

$$orange : [ \qquad ] \}$$

128

## Disadvantage

1. No context awareness
2. Sequential processing
   X
   ↳ Parallel processing
3. Long context failure

## Advantages

1. Generate context rich vectors by using Attention
2. Parallel processing
3. It can handle long contexts.

# TRANSFORMERS

-) Attention is All you Need - Google team
                                                        -2017.

→ Encoder
  Decoder
  ) Enc & Dec

→ Architecture Deep Dive.

Introduction → Architecture → Positional
                    Deep Dive         encoding
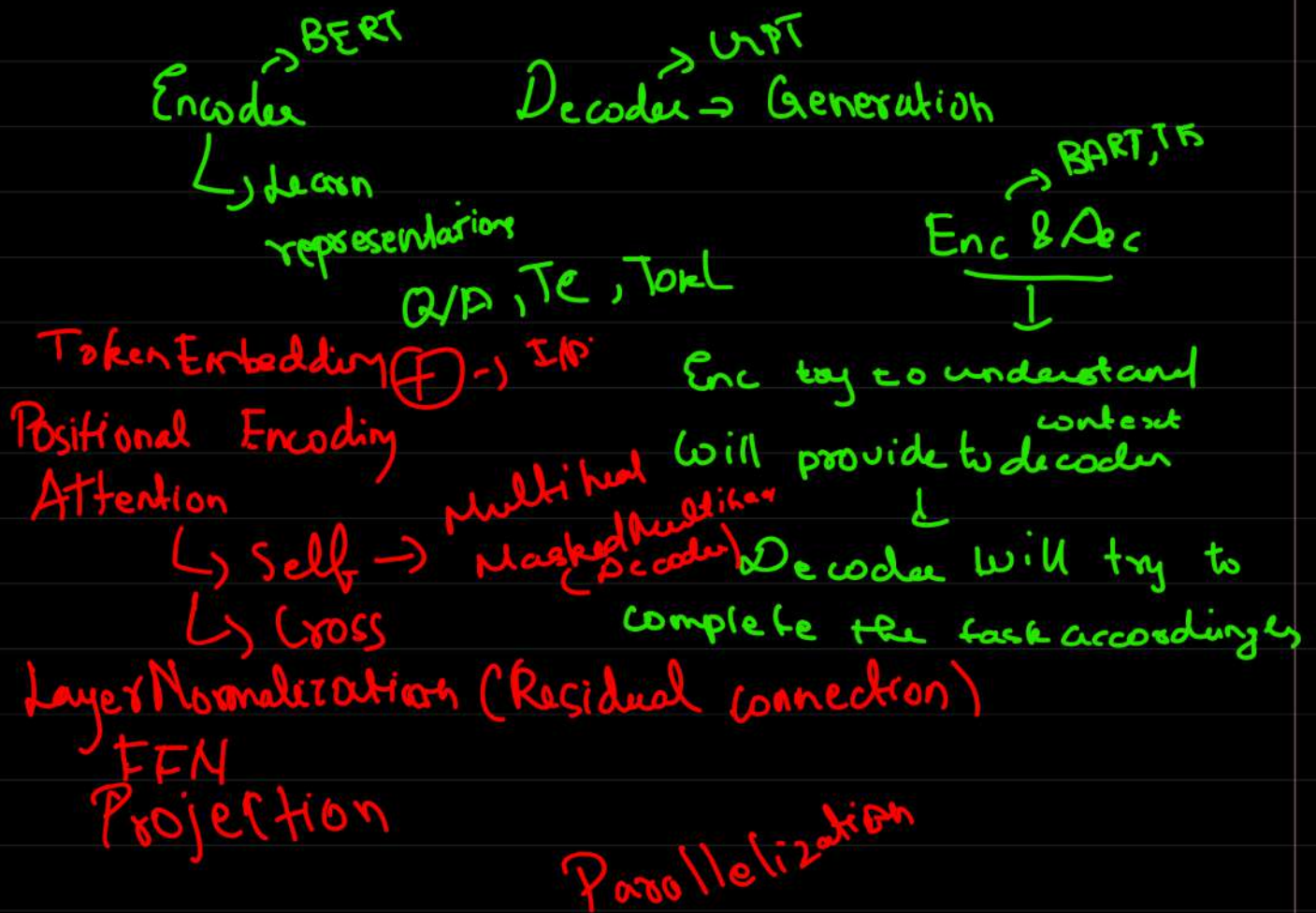                                          ⊥
                                      Attention
                                          ⊥
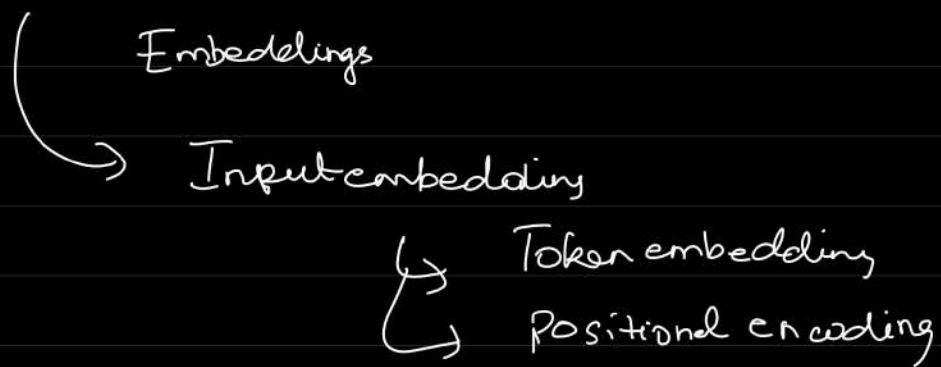Decoder ← Encoder ← FFN ← Layer Normalization
  ⊥
Transformer) Dataset building -) Tokenizer
                                          ⊥
Decoding ← Attention ⟵ Training
strategies   Visualization    ∩ & Inference
              How tokens are processed

Encoder $\longrightarrow$ BERT
   $\searrow$ Learn
   representations

Decoder $\rightarrow$ GPT
   $\rightarrow$ Generation

$\rightarrow$ BART, T5
Enc & Dec
$\downarrow$

Q/A, Te, Toul

Token Embedding $\oplus$ $\rightarrow$ I/P

Positional Encoding
Attention
  $\hookrightarrow$ Self $\rightarrow$ Multi head  Will provide to decoder
  $\hookrightarrow$ Cross   Masked Multihead (Decoder)

Enc try to understand
     context
Will provide to decoder
    $\downarrow$
Decoder will try to
complete the task accordingly

Layer Normalization (Residual connection)
  FFN
Projection

Parallelization

# POSITIONAL ENCODING

Embeddings

↳ Input embedding

↳ Token embedding
Positional encoding

Token embedding: - Tokenizers.

$dmodel = 512$

The apple is tasty → bs→1, seq-len→4,

[The, apple, is, tasty]
    ↓      ↓      ↓    ↓
    1     10     8    7

{ 1: [                              ] 512 → vectors / embeddings

}

[ [                    ],  512
  [                    ],
  [                    ],
  [                    ] ]     → Token
                                embedding

0   1   2      3        4    5      6    7
(I) had two   apples  (1 orange )apple and
another 1 (red) apple    where the (second one
was slightly      fastier it comes from Apple Inc.
        apple → 6
        apples 11

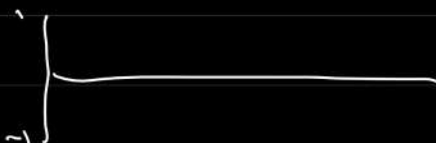# Sinusoidal

0 , 1, 2, 3 _____       ⟨512⟩

→ Scaled
    in nature.

cosine, sin

Encodes relative
position word.

sin( )
  even
cos( ).

$$PE(POS, 2i) \rightarrow Sin\left(POS / 10000^{2i/dmodel}\right)$$
         0                              ↑
                                   frequency, h

i=0

$$PE(POS, 2i+1) \rightarrow cos\left(POS / 10000^{2i+1/dmodel}\right)$$
                    1

## Why Sinusoidal

* Unique encoding of position of words
** Encode relative to other position
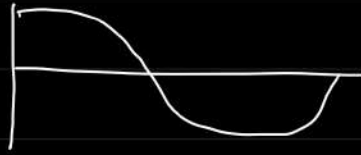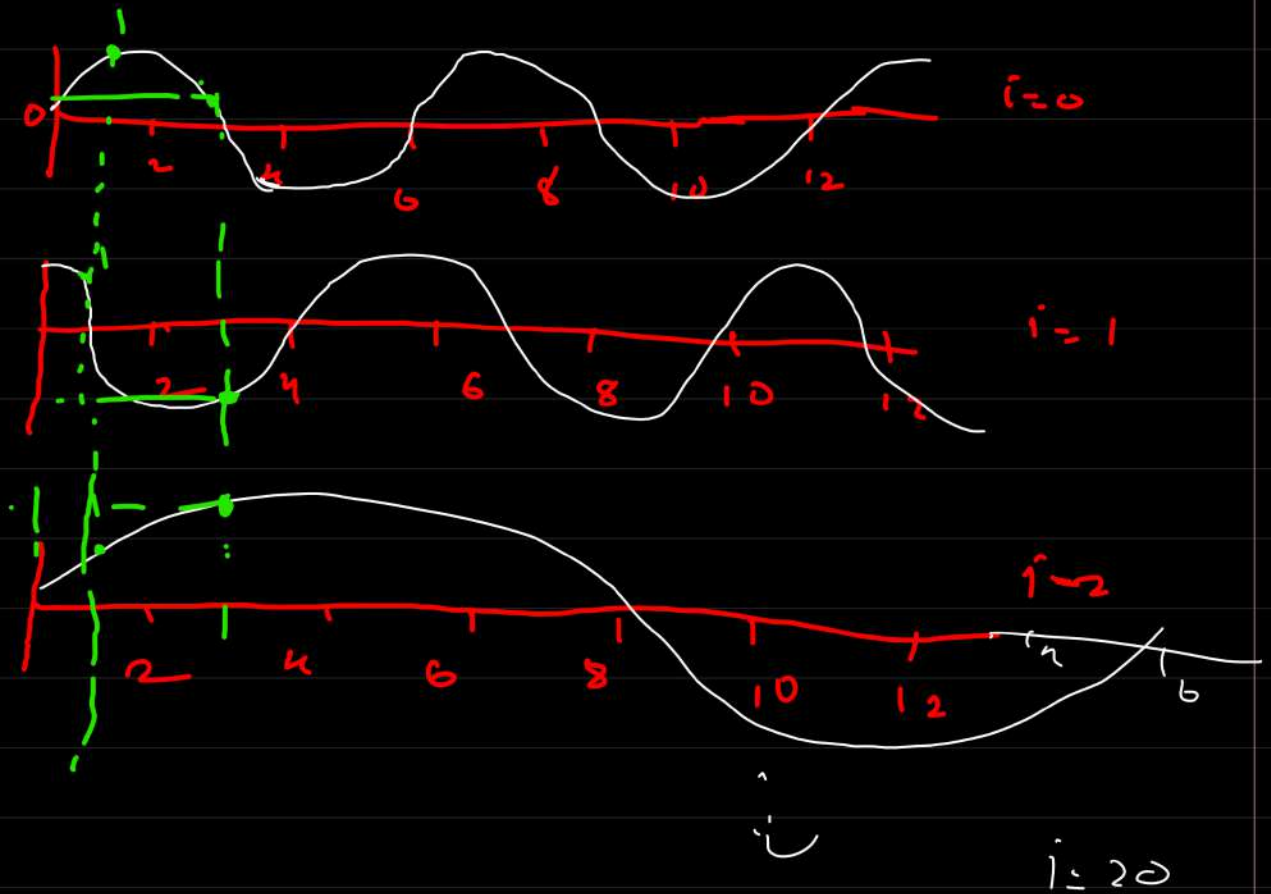* Range → Normalized [-1, 1]

## Why sin & cosin

sin & cosin

sin & cosin

512 → embed
        i=0

{          }

P1
P0  P2
i3

|————|————|————|————|————|————→ position
0    1    2    3    4  5 6

Sin

$\sin(0) = 0$

cos

$\cos(0) = 1$

seq → 4
①

# DIMENSIONAL DEPENDENCY:.

i=0

0
1
2
4
6
8
10
12

i=1

2
4
6
8
10
12

i=2

2
4
6
8
10
12
6

i = 20

Formula :-

$$PE(pos, 2i) \rightarrow \sin(pos/10000^{2i/dmodel})$$
                                              $\underbrace{}_{n}$

$$PE(pos, 2i+1) \rightarrow \cos(\quad\quad\quad^{2i+1}\quad)$$

The value is tasty $\rightarrow$ POS [0, 3]

$$d_{mat} = 4.$$

$$i \quad\quad\quad n \rightarrow 100$$

| | 0 | 0 | 1 | 1 |
|---|---|---|---|---|
| pos 0 The | $\sin(0)$ $= 0$ | $\cos(0)$ $= 1$ | $\sin(0)$ $= 0$ | $\cos(0)$ $= 1$ |
| 1 water | $\sin(1/1)$ $= 0.84$ | $\cos(1/1)$ $= 0.54$ | $\sin(1/10)$ $0.10$ | $\cos(1/10)$ $\simeq 1.0$ |
| 2 is | | | | |
| 3 tasty | | | | |

POS = 0   i = 1

$$(0/ \text{- - -}) \rightarrow 0$$

POS = 1   i = 0 $\rightarrow$ i = 1

$$i = 0, pos = 1$$
$$\sin(1/1)$$
$$\cos(1/1)$$

$$\sin\left(1/100^{\frac{2\textcircled{1}4}{4}}\right) \Rightarrow \sin\left(1/100^{1/2}\right)$$
$$= \sin(1/10)$$

$$\cos\left(1/100^{\frac{2\textcircled{1}+1}{4}}\right) \rightarrow \cos\left(1/10\right) = \frac{a}{10}$$

512→ [seqlen, embed dim]

[          ]

I/P→ Token Embed + Pos Encoding.

512→ [seqlen, embed dim]

# ATTENTION

<span style="color:green">* Generation of context rich vectors
* Parallelization</span>

My name is Vasanth → Scaled Dot Product Attention

| | My | name | is | Vasanth |
|---|---|---|---|---|
| My | * | * | - | + |
| name | * | * | - | * |
| is | - | + | * | * |
| Vasanth | * | * | - | * |

Increase seq length ↑ Use ↑

→ Generating context rich vectors

↳ Encoder → Decoder (Generate)

Query, key & Value.

## <span style="color:red">Query :-</span>
<span style="color:red">Represents what you want</span>

## <span style="color:green">Key :-</span>
<span style="color:green">The location of the answer you want</span>

**Value:** -

The answer.

Eg.

$V, K$   $V_i$ $^k$ $V$ $^k$ $V$ $^k$   $V$ $^k$ $^k$ $V$ $^k$   $V$ $^k$

Context $V, K$ (Vasanth) is a youtuber who has bad handwriting $\rightarrow$ Encoder.

Question: Who has a bad handwriting?

$\rightarrow$ Encoder

$\rightarrow$ Query   $\rightarrow Q$

Decoder $\rightarrow$ Question.. Who has a bad handwriting?

$Q \rightarrow$   $Q^* K (Vasanth) \rightarrow \boxed{0.8}$ $\Big\}^{1512}$

$Q^* K (is) \rightarrow 0.01$   $1$

$\vdots$   $0.05$

$Q^* K(handwriting) \rightarrow \cdot$

$$Attention = softmax\left(\frac{QK^T}{\sqrt{d_{model}}}\right) * V$$

Vasanth.

$0.1$   $0.01$   $0.05$   $0.02$   $0.02$

$softmax\left(\frac{(Q \cdot K^T)}{\sqrt{d_{model}}}\right)$   $\overline{Q \cdot K^T}$   $\overline{Q \cdot K}$   $\overline{Q \cdot K}$   $\overline{Q \cdot K}$

$Q$   $\rightarrow$ 1 head.

# Multihead Attention



Concat

→ VASANTH

Why we divide $\sqrt{d_{model}}$ Variance ↑↑?
↳ To reduce variance

## MASKED MULTIHEAD ATTENTION

↳ Autoregressive

Vasanth has bad handwriting

$\rightarrow$ (inf)

[ [ 1 , 0 , 0 , 0 ], → Teacher forcing

[ 1, 1, 0 0 ]

[ 1, 1, 1, 0 ]

[ 1, 1, 1, 1 ] ]

LOOK AHEAD MASK

PADDING MASK       $0 \to 1, 8, 9, 17, 170$
                   $170, 170$

Self Attention $\to$ Gross Attention $\to$ Decoder Architecture

## Flow

Input → Query $\quad$ x $\quad$ W $_q$

$\quad$↳ Key $\quad$ x $\quad$ W $_k$

$\quad$↳ Value $\quad$ b $\quad$ W $_v$

↓

Split among heads

Scaled dotproduct attention

At each head: $\quad$ Softmax$\left(\dfrac{QK^T}{\sqrt{d_{model}}}\right) * V$

↓

Concat o/p of each head

↓

Project/learn the concatenation

↓

ATTENTION OUTPUT

( Context rich vectors)

Add & Layer Norm → Paper → ①

ResNet → ②

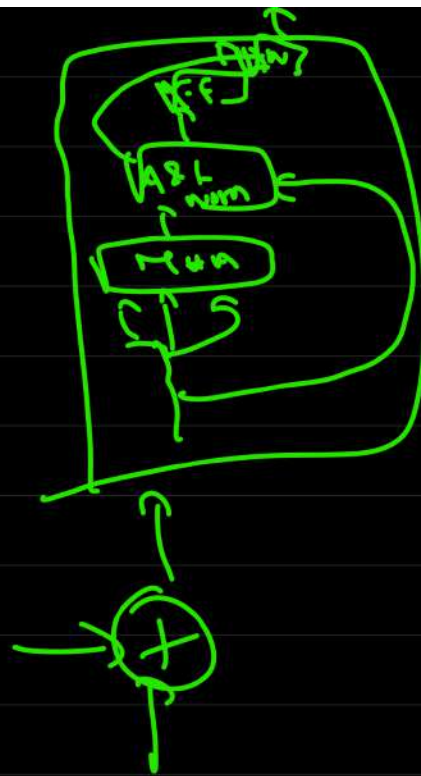VGG 19, 16 → Im$^a_n$ge Net → CNN.

Problem:
increased

On an extended training, loss

Soln: ResNet → Residual Network → Skip connection

if accuray >= layer is consider

else
skip layer

# Add & Layer Norm

## Add :-

# Layer Normalization

**Formula:** $X \Rightarrow$ Layer Norm $[Sublayer(x) + Output]$

**What?.**

Normalization is a process of scaling the values close to zero & around

$$(-3, 3) \Rightarrow [0, 1]$$

**Why?.**

* Stable training
* Reach the global minima (optimum value) faster.

↳ Training becomes faster.

# Layer Norm:-

$$X_i^1 = b [W^T x_i + bi] \rightarrow \text{Forward O/p}$$

$$y = \gamma_i \left[ \frac{X_i^1 - \mu_i}{\sigma_{i} + \epsilon} \right] + \beta_i$$

Model will cale.

$\rightarrow \left( \begin{array}{c} W^T x \\ +b \end{array} \middle| \rightarrow Ad() \right) \rightarrow X_i^1$

$\rightarrow$ Standardication

$\gamma_i, \beta_i \rightarrow$ learnable

Result

All the values in this layer
will be come [-3 to 3]
mean = 0, std → 1

token → 2        xcomb → 3

$$\begin{bmatrix} 0.1 & 0.2 & 0.3 \\ 0.4 & 0.2 & 0.3 \end{bmatrix}$$

$$\mu_{11} = \frac{1}{3}[0.1 + 0.2 + 0.3]$$
$$= 0.2$$

$$\mu_{22} = \frac{1}{3}[0.9] = 0.3$$

$$\sigma_{11} = \sqrt{\frac{1}{3}\left[[0.1-0.2]^2 + (0.2-0.2)^2 + [0.3-0.2]^2\right]}$$

$$= \sqrt{\frac{1}{3}[0.02]} = \sqrt{0.066}$$

$$= 0.08$$

$$\sigma_{21} = 0.08$$

$$\text{mean } \mu = \begin{bmatrix} \mu_{11} \\ \mu_{21} \end{bmatrix} = \begin{bmatrix} 0.2 \\ 0.3 \end{bmatrix}$$

$$\text{std } \sigma = \begin{bmatrix} \sigma_{11} \\ \sigma_{21} \end{bmatrix} = \begin{bmatrix} 0.08 \\ 0.08 \end{bmatrix}$$

$$Y = \begin{bmatrix} \dfrac{0.1-0.2}{0.08} & \dfrac{0.2-0.2}{0.08} & \dfrac{0.3-0.2}{0.08} \\[2mm] \dfrac{0.4-0.3}{0.08} & \dfrac{0.2-0.3}{0.08} & \dfrac{0.3-0.3}{0.08} \end{bmatrix}$$

$\sigma\sigma_1 = \gamma Y + \beta$

$$Y = \begin{bmatrix} -1.2248 & 0 & 1.2248 \\ 1.2248 & -1.2248 & 0 \end{bmatrix}$$