# Transformer:

Attention $\rightarrow$ Update the Embedding

Key, Query and Value

| | token 1 | token 2 |
|---------|---------|---------|
| token 1 | 1 | 0.8 |
| token 2 | 0.8 | 1 |

token1 = 0.8 token2 + 1. token1

normalize

$\nearrow$ Softmax

= 0.7 token + 0.3 token

update the numbers

---

Key · token1     Query · token 2

linear transformation

Attention $\longrightarrow$ Multi head Attention

Key, Query

Key 1 Query 1

key 2 Query

key 3 Qy 3

...

$$Score = Softmax \left( \frac{key1 \cdot token1 + Qy1 \cdot token2}{\sqrt{Dimension}} \right)$$

Score * Value 1

$\Downarrow$

Update embedding for each token

Feed Forward lays

$\vdots$

Encoder 2

# Decoder

Masked Multi Head Attention

## Masking the future

- $-inft \Rightarrow$ future token
- $e^{-inf} \rightarrow 0$

---

Decoder output

## One Scalar

number of distinct token neurons
softmax

## Predict one token at a time

---

Decoder $\begin{cases} I/P \rightarrow Vector \text{ (sentence, encoder)} \\ O/P \rightarrow Scalar \text{ (one token)} \end{cases}$

# Decoder → Generative AI

OpenAI | Transformer → 2017 | → Google

| GPT-1 → 2018 |

GPT-1 → 117 million

Input → Book Corpus → 7000 ~ 6GB

task → Predict next token by sending input with some context length

| Hi how are you? I am good |

Context length → 3

Hi how are → you
how are you → ?
are you ? → I
you ? I → am
? I am → good

Unsupervised - Pre-trained
model

## Supervised fine tuning

### ULMFit

Transfer learning is done on top of Pre-trained model

NanoGPT ⟹ Pre-trained Model from Shakespeare data set

tuning SQL Query to write the book SQL statement based on text input

---

### GPT-2    2019    1·5 billion paramtr

Huge Data ⟹ 40 GB

Redit ⟹ post ⟹ > 3 kauma votes + link ↓ scrapped

Webtext

Strategos : Unsupervied Pre Training

Surpaise → translator

## Evaluation

Papers → Glue MMLU, ✓
HellaSwag ✓
ARC ✓

Im_eval

accuracy
Baseline
Mistal 7B → accuracy

Mistal 7B Instruct → accuracy

GPT - 3 → OpenAI

Input → Unknown
Parameter → 1.7B billion ← 1.7 biln
100 X

Strategy :

1. Instruct Model

( Mistral AI 7B Base line )

Trained with Input and output
with Instruction and Response

# Instruction
can you perform Sentiment for
below data

# Content
The Product is good

# Response
positive

Fine tuning

Mistral AI Instruct Model

Fine Tuning → Peft → Lora

Faster
loose the Avoid
older
information

Update new weight
instead of updating
existing Pre trained
model weights

Instruct model ——→ Update it to better model

(Reward model)

↓

[Classification] model

Reinforcement learning

Collect the Data

(I/p) → O/P → good !

I/p     O/P ↛ bad 0

Human
labeller
contractor

I/P | I want to hack my neighbor Internet?

O/P

This is really
bad practice

⑦

O/P

Yes, I can
help you with

⓪

O/P

this is
really ✓
not a
good idea.
but I can
help you. If
you want
to access any
free Internet

⑩

# Reward Model →

I/P → Question & Answer

O/P → Score

Algorithm

PPO → Proximal policy Optimization

Update the model based on reward score