

## Zero-shot classification:-

texts not labelled

specify labels

X FT

return Probability scores of labels

## text-generation:-

provide a prompt

auto-complete it.

can control no. of o/p generated & max.length.

## Mask filling:-

## Text generation:-

NER:-

QA:-

## Summarization:-

## History of Transformers:-

June 2017.

GPT: first PT Xmore.

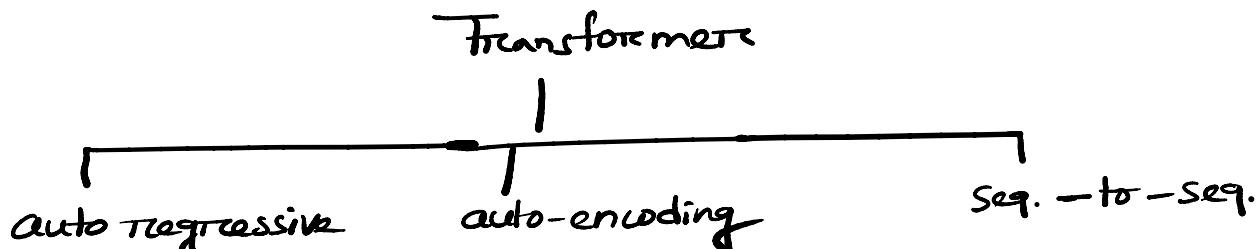
BERT:-

GPT-2:-

DistilBERT:-

BART, T5:-

GPT-3:- Zero-shot learning.



auto-regressive

auto-encoding

Seq. - to - Seq.

Transformers are LMs:-

trained on large amount of raw texts.

Self-supervised w/o label data

Statistical understanding of lang.

✗ useful for specific tasks.

Pre-training

Transfer learning  
FT in supervised way

Fine-tuning.

Causal LM:-

I/P:- past, present I/P

O/P:- predict next word.

Masked LM:-

I/P यहाँ कोई token masked करेंगा इसका क्या O/P होगा?

Transformers are big models:-

→ outliers (DistilBERT) ↑  $\uparrow$  performance मात्राएँ

→  $\uparrow$  model size &  $\uparrow$  amount of data.

Costly in time & compute.

environment impact.

footprint of finding best hyperparameters ↑.

Scratch से train करावे cost ↑

Soln:- sharing

Transfer learning:-

PT training from scratch.

Large amount of data

Large amount of data

Several weeks.

FT after PT

perform + training dataset specific to your task.

Reasons for not training directly for final task:-

Sim betw PT & FT dataset.

take adv. of K acquired during PT

PT data ↑

FT data ↓

amount of time & resource →

environmental impact →

General Architecture:- of Transformer:-

Two blocks:-

Encoder:- receive i/p

build a rep. (features) of i/p

acquire K from i/p.

Decoder:- receive encoder's rep. + other i/p = target seq.  
Optimized for generating o/p.

Encoder only models:- good for task req. understanding of i/p.

Sentence classification

NER

Decoder n n :- generative tasks such as text generation.

En-Decoder / Seq-to-seq. models:- generative task req. i/p.  
translation / summarization

Att<sup>n</sup> layers:- pay specific att<sup>n</sup> to certain words

(Translation note: many certain n depends on certain other words.  
... are some important.)

Translation starts after certain  $n$  depends on certain other words.  
Subject/noun  $\Rightarrow$  ~~is~~ dependent.  
meaning of a word affected by context.

Original Architecture:-  
en. i/p: Source lang. + full sentence  
de. " : target lang. + only gen. words so far. + en. rep.

Decoder's  $Att^h$  layers:-

1st:  $Att^h$  to previous decoder's o/p.

2nd: - " " encoder's o/p.

helpful ~~isn't~~ context ~~comes~~ comes as full sentence access ~~not~~ ~~not~~ !

Arch - Skeleton of model.

def. of each layer & op. within the model.

checkpoints - weights loaded in a given archi.

Encoders Arch.: -

only use encoder. of Xmerc.

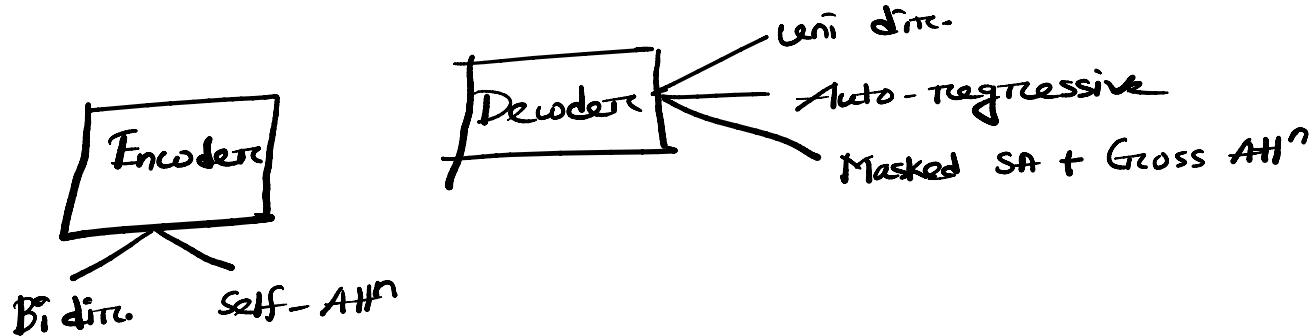
$Att^h$  layers access all words in initial sentence.

bidim.  $Att^h$

auto-encoding models.

PT by MLM (Masked Lang. Modeling)

best suited for tasks req. understanding of full sentence.



Bi-direc.      self-Att<sup>n</sup>

Decoders models:-

only use decoder of Xmer  
Att<sup>n</sup> access only before/left words.

Auto-regressive

generate seq. of texts by modeling the conditional probabilities.  
of previous tokens.

CTRL

GPT - 1, 2, 3

Xmer-XL

Seq-to-Seq / En. - Decoder Models:-

Autoregressive manner  
previous gp of de. acts as ip

Bias & limitations:-

PT on lots of data

best + worst data.