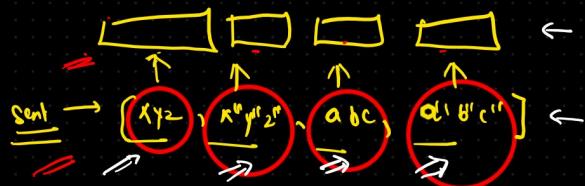
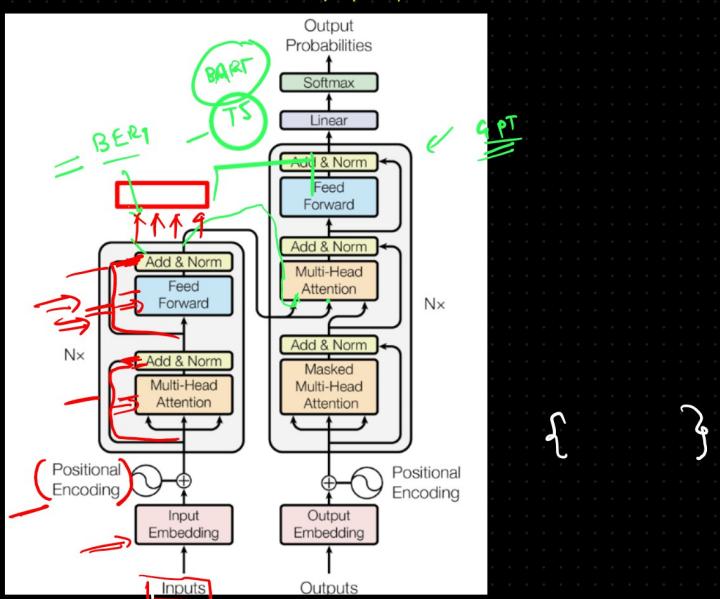


Transformer code → Pytorch → file corrupted  
20 min recording  
short  
Sunday X, Wednesday X  
quick revision

Transformer → BERT  
→ GPT  
→ Mamo  
→ Mistral } Language text  
= 99.5 (ε-6)  
→ CLIP → OpenAI  
→ DALLE → OpenAI  
→ Vision Transformer

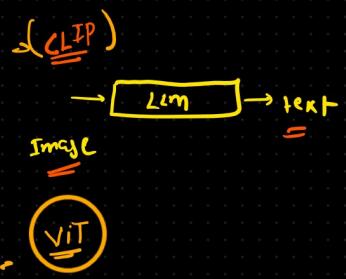
(Diffusion) Generating high resolution image

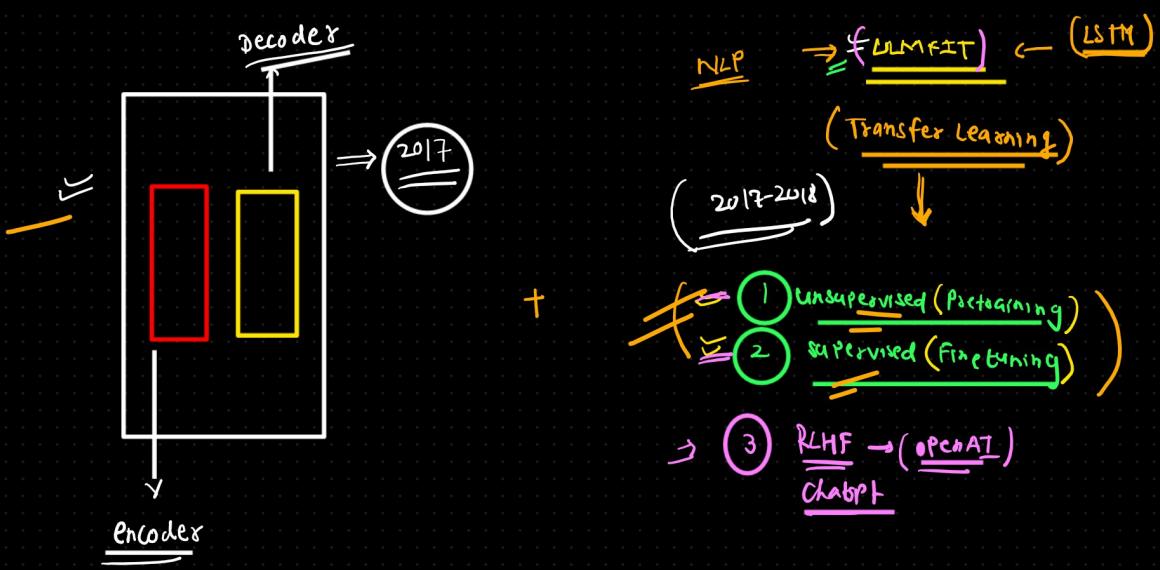
API ⇒ huggingface → openAI → GPT Gemini → Qwen CMstudio }



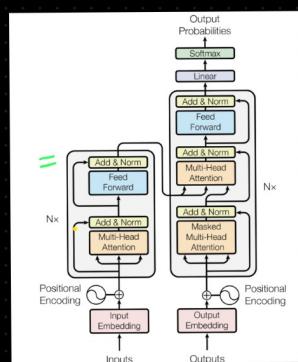
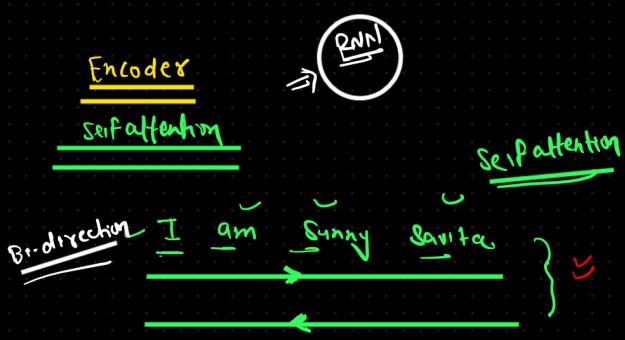
$\times$   
 $(\text{ViT} + \text{BERT})$

$\hookrightarrow (\text{GPT}, \text{LLaMA}, \text{Mistral})$





## Bidirectional encoder Representation from transformer

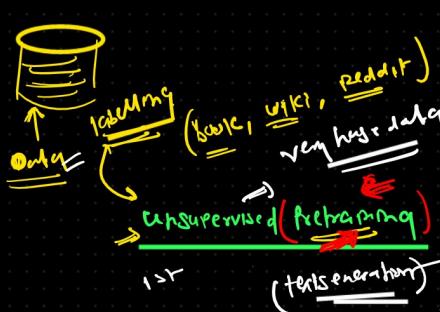
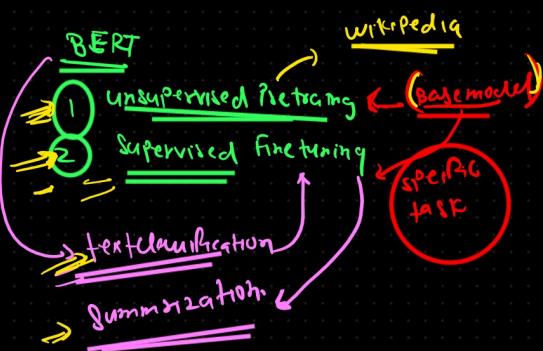
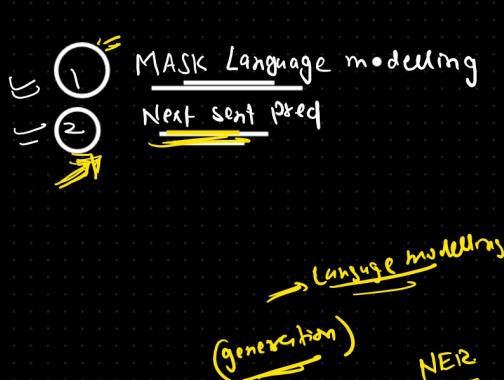


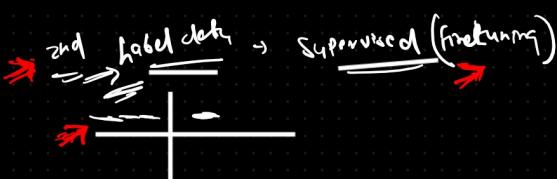
## Positional encoding

## Self attention

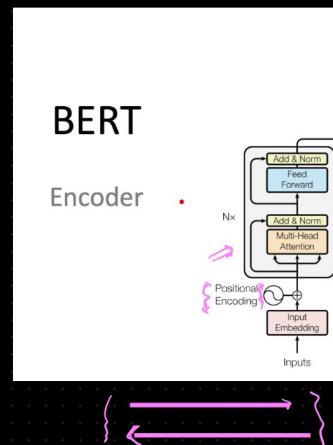
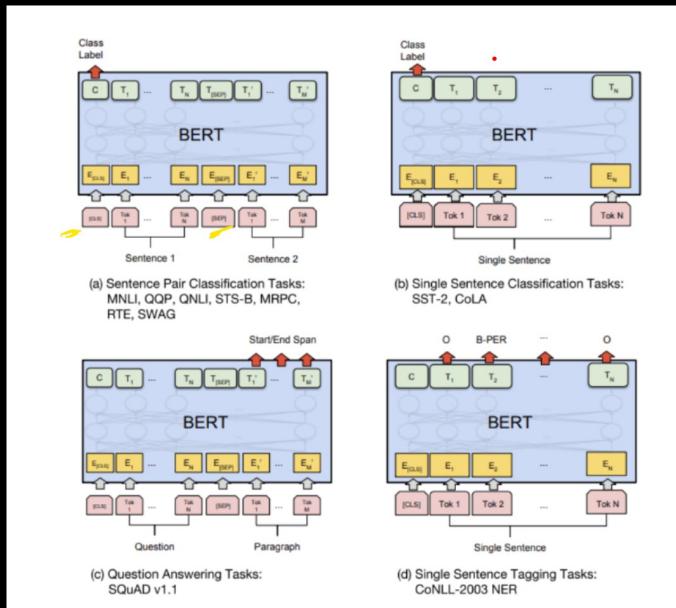
Layer norm, feed forward

BERT  $\Rightarrow$  12 Layer  $\Rightarrow$  hidden layer size  $\Rightarrow$  768 } FF  
 $\underline{(110M)}$       24 Layer  $\Rightarrow$  hidden layer size  $\Rightarrow$  1024  
 $\uparrow$   
 $\underline{340M}$



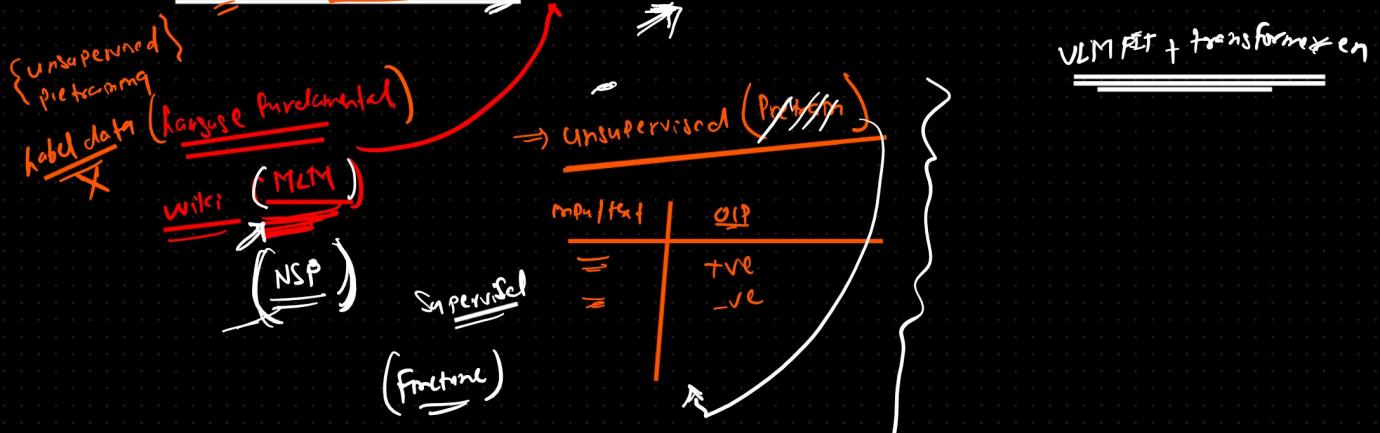


BERT  $\Rightarrow$  (MLM, NSP)  
GPT  $\Rightarrow$  lang modeling  
(auto regressive)



Transformer en  $\rightarrow$  BERT  $\rightarrow$  Bi

MASK lang. Random word 1st  
training [CLS] am [MASK] who [MASK] A [sep] [MASK] [sep] You [MASK] are [MASK]

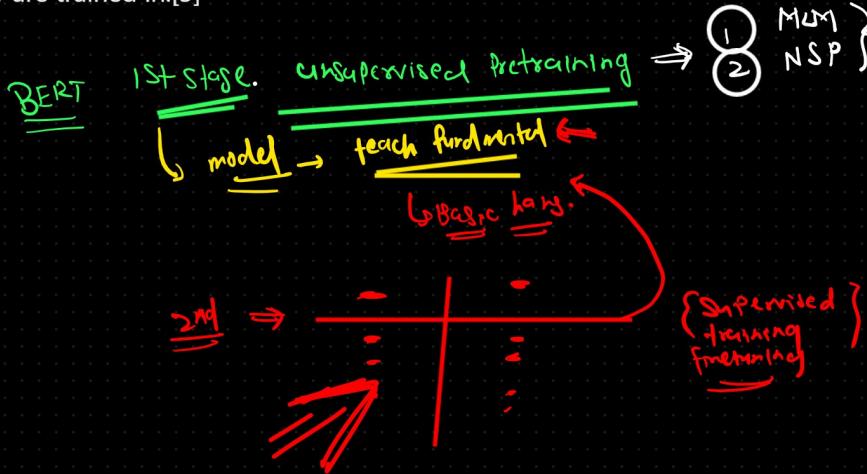


VLM fit + transformer en

text  
labeling

~~sent 1~~ A large language model (LLM) is a type of machine learning model designed for natural language processing tasks such as language generation. LLMs are language models with many parameters, and are trained with self-supervised learning on a vast amount of text.

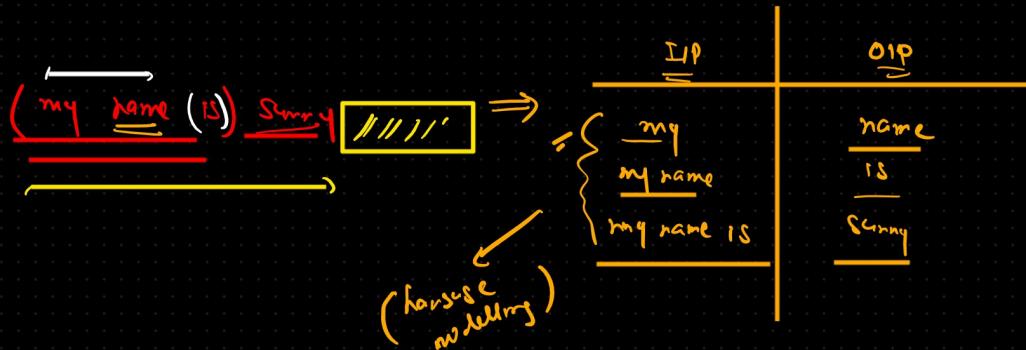
The largest and most capable LLMs are generative pretrained transformers (GPTs). Modern models can be fine-tuned for specific tasks or guided by prompt engineering.[1] These models acquire predictive power regarding syntax, semantics, and ontologies[2] inherent in human language corpora, but they also inherit inaccuracies and biases present in the data they are trained in.[3]



GPT  $\Rightarrow$

A large language model (LLM) is a type of machine learning model designed for natural language processing tasks such as language generation. LLMs are language models with many parameters, and are trained with self-supervised learning on a vast amount of text.

- 1 Unsupervised pretraining  $\Rightarrow$  Autoregressive modelling (dans modelling)  
2 Supervised finetuned



1 Bi-directional Attention  $\Rightarrow$  strong contextual understanding  
Pretrain on massive Wiki data  $\hookrightarrow$  M(MLM, NLP)

$\hookrightarrow$  Sentiment analysis, NER, Summarization

Encoder based architecture.

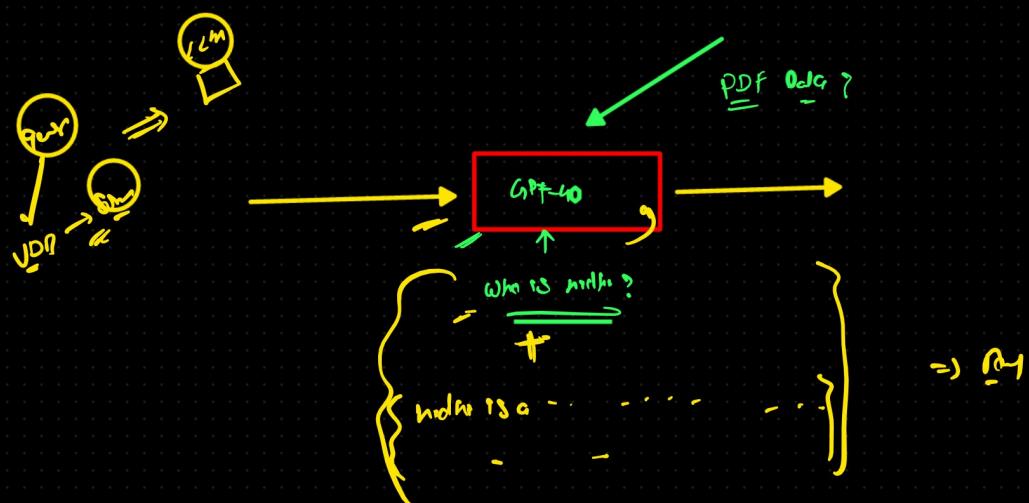
2 GPT  $\Rightarrow$  Unidirectional attention

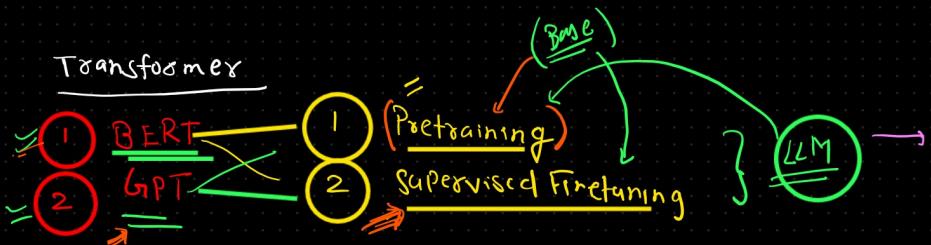
Autoregressive modelling  $\Rightarrow$  (Pretrained)  
(next word pred)

decoder based archi  $\hookleftarrow$  transformer reddit

massive (book, wiki, social platform)

$\Leftarrow$  text generation  $\Rightarrow$  Opt 3.5  $\hookrightarrow$  Chatspt

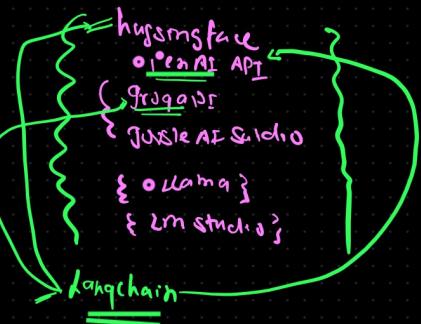




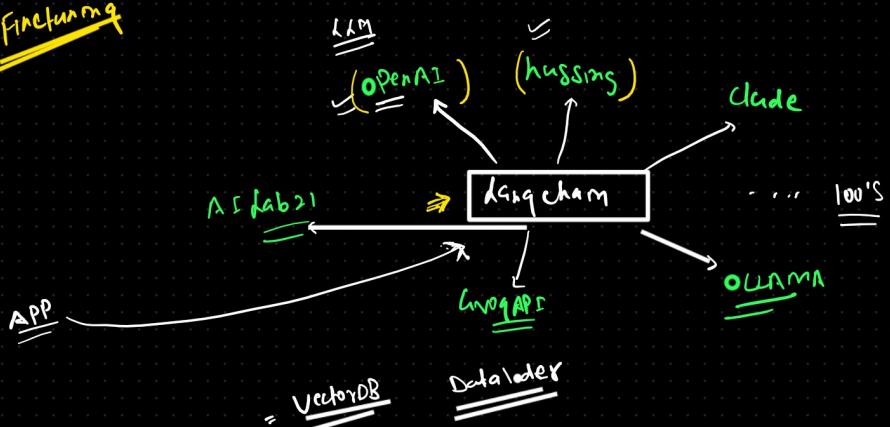
student  
foundation  
 $(1-10\text{nm})$   
 $\frac{11-12}{11-12} \left( \begin{matrix} \text{PCM} \\ \text{PCB} \end{matrix} \right)$

{ 3 Llama  
4 Mistral }

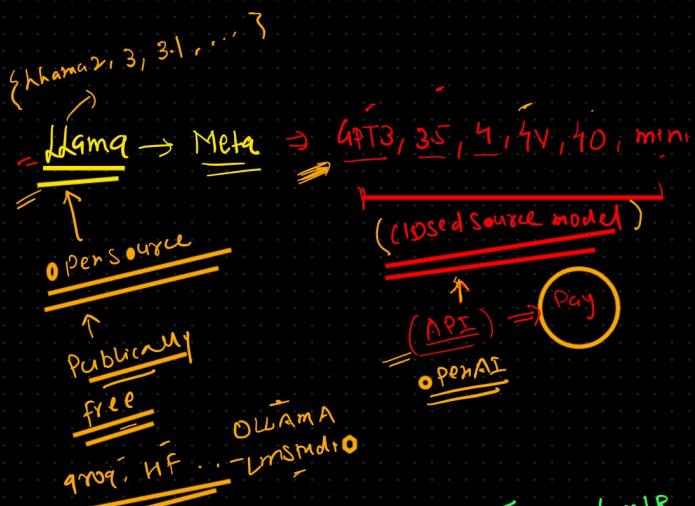
Multimodality  
{ 5 CLIP  
6 DALLE  
7 VISION tools.  
8 Diffusion }



Finetuning

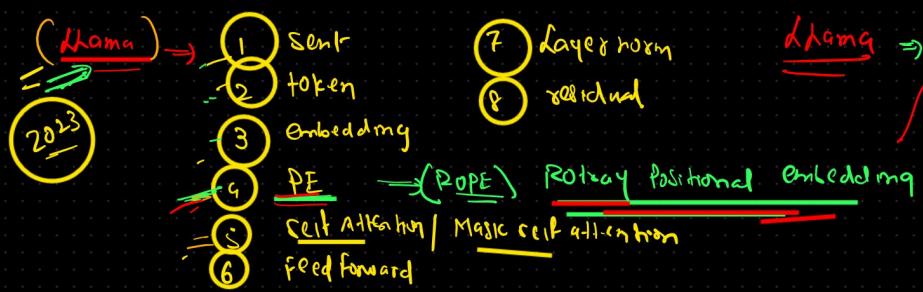


PAY



Web ← trainable

Transformer ⇒ decoder-based archi



llama ⇒ larger context

↑  
culture layer shift

## Supervised



↳ dialog.

Mistral → French  
→ mistral AI  
→ 7B, 70B

- ↳ decoder based model
- ↳ efficient computational perf

Benchmark

reduce redundancy

state sharing

{ long }  
cent.

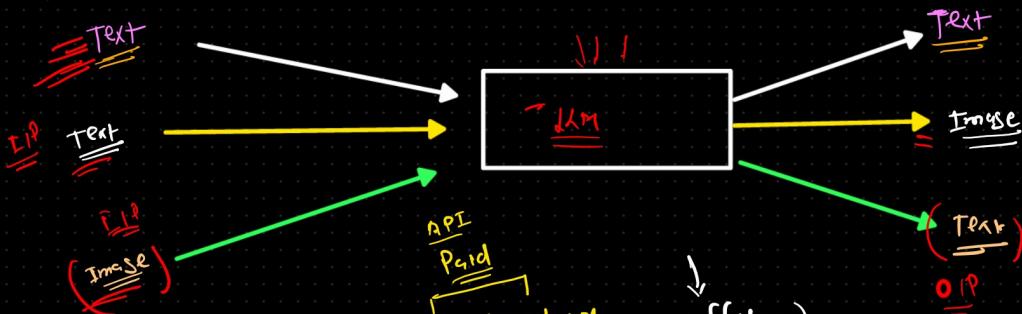
- (1) (sequential)  
query attention
- (2) sliding window attention

FFNN

Masked Self Attention



{ Khamla }  
CPT  
Mistral



API  
Paid  
Dance / Midjourney

Diffusion

Image

Image

Text

Birds are flying

Image  
probability

(Nesting,  
pending)

= (NN + LSTM)

= CNN + BERT

= 2014-2015  
2017-2018, 19

OpenAI  
CLIP  
ViT

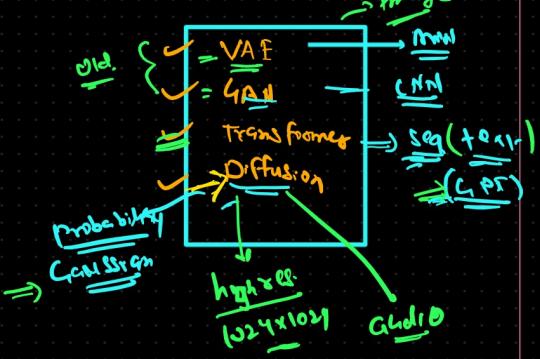
OCR

ChatGPT

Image  
(MM)

Text

Birds are flying



CLIP  $\rightarrow$  Transformer  $\Rightarrow$  (Encoder)

dual encoder



One person standing

first encoder  $\rightarrow$  Image

second encoder  $\rightarrow$  text

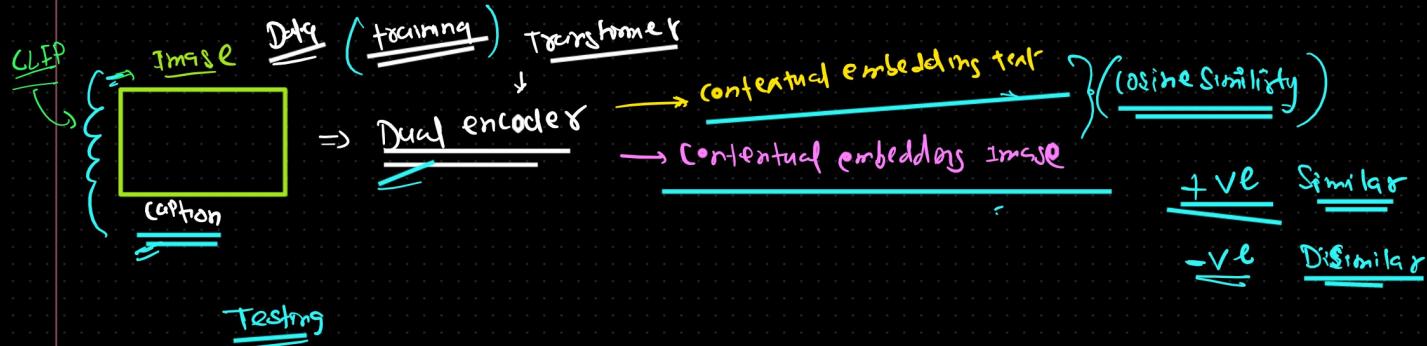
black white

$\downarrow$

  $\Rightarrow$  Pixel  $\Rightarrow$  Numbers (0-255)

Flatten Image  $\rightarrow$  (1D array)  $\Rightarrow$  Transformer encoder  $\Rightarrow$  contentual embedding

Text  $\Rightarrow$  embedding  $\Rightarrow$  Transformer encoder  $\Rightarrow$  contentual embedding



1 hundreds of millions Image-text pairs

2 Dual encoder setup

Image encoder  
Text encoder

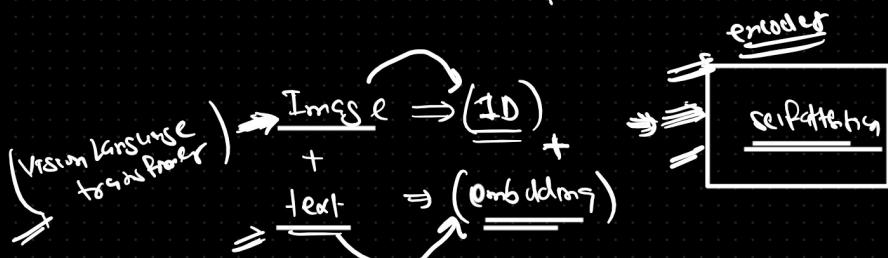
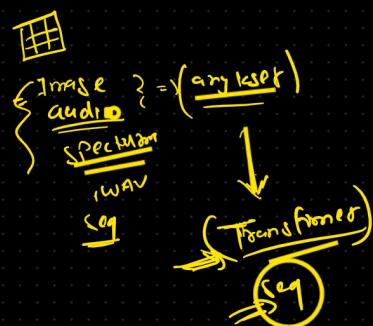
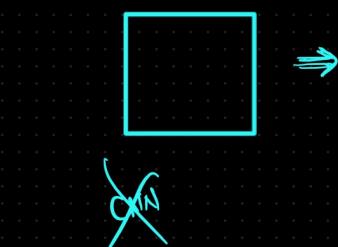
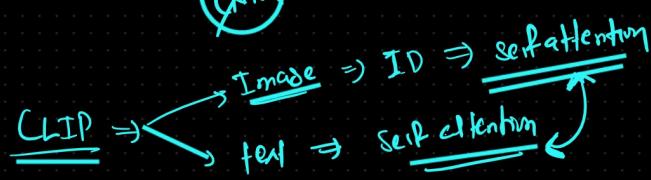
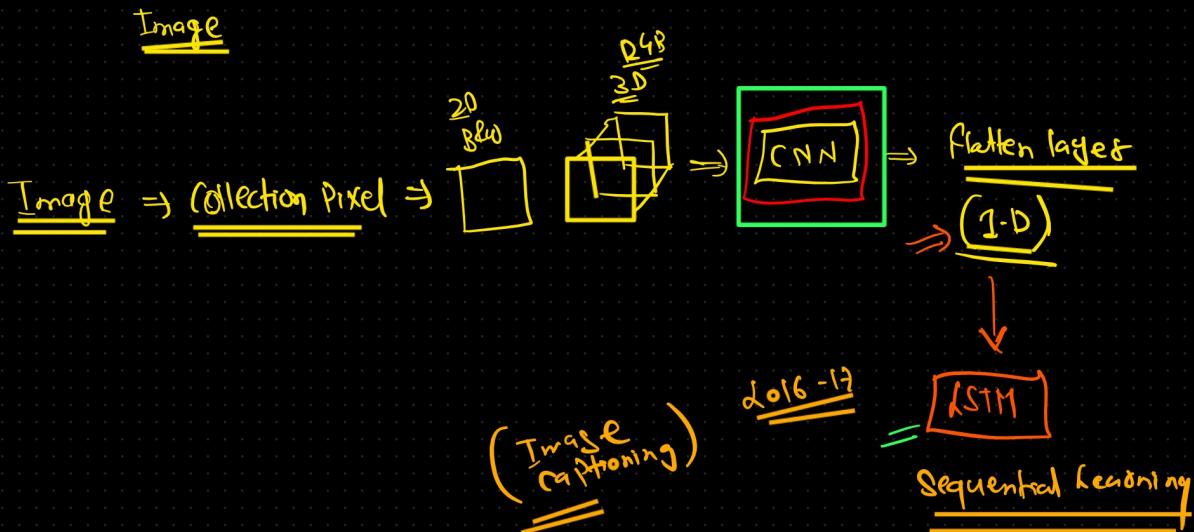
3 Similarity

Cosine Similarity

4 Zero-shot capabilities

  $\rightarrow$  Description

## Vision and lang transformet (ViLT)



Dall-E

Text → Image generation

≡ GAN  
≡ dVAE  
⇒ ? denoising

Transformers

DALLE

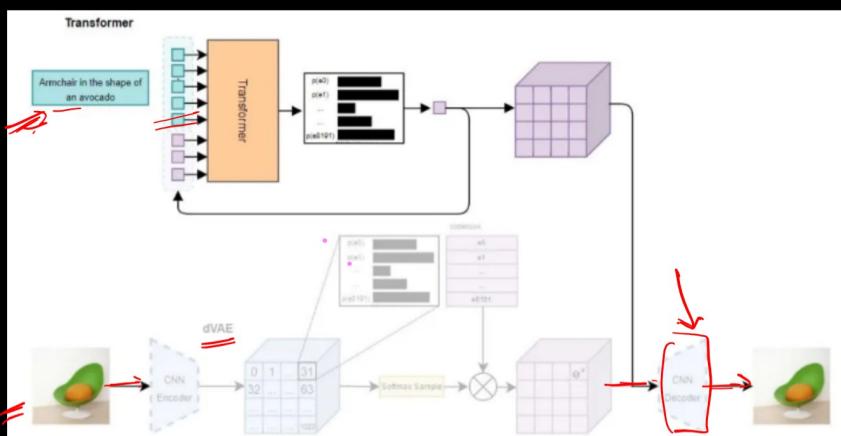
=  
Patch image  
Pixel  
GPT

128×128  
512×512

DALLE-2

→ Diffusion  
high resolution → HD ⇒ 1024×1024

Two Stage Process



training

VAE (CNN)

2021

Dalle 1 ⇒ GPT like transformer arch.

Dalle 2 ⇒ CLIP + Diffusion

2022

similarly

training

VIT

Image

CLIP

text 2 image

text

VAE

text & image encoder

relationship

B

Image Gen

- autoencoder ⇒ GPT ⇒ Image Patching

- Diffusion ⇒ nd, HD

Dalle → CLIP + Diffusion  
= Transformers

Langchain = (AI)  
Storage {String, Ching} → Unified code base → compact  
↓  
Code → LLM

VectorDB = CloudDB, Inmemory, OnDisk

API :> HF, OpenAI, Code

Document => WQb  
HTML, JSON, XML

huggingface  
VectorDB =