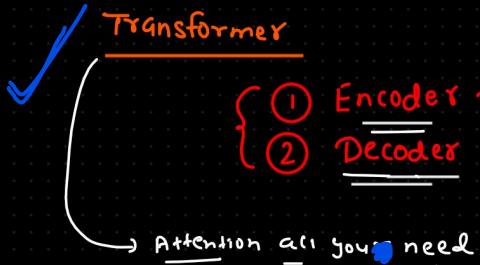
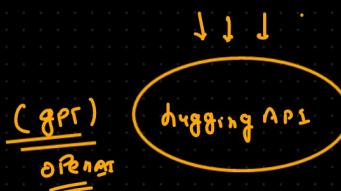


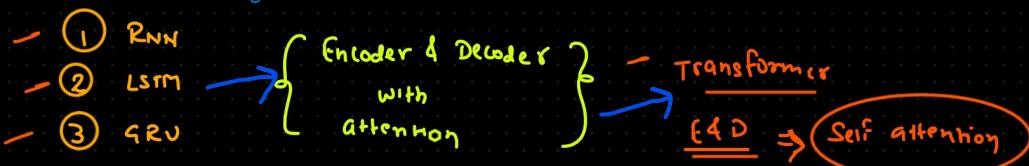
## NLP

- ① RNN
  - ② LSTM
  - ③ GRU
  - ④ encoder & Decoder
  - ⑤ Attention
  - ⑥ transformer
  - ⑦ GPT & BERT (Language modeling) ULMFiT
- = Project { ① - BERT  
② Google Pegasus }

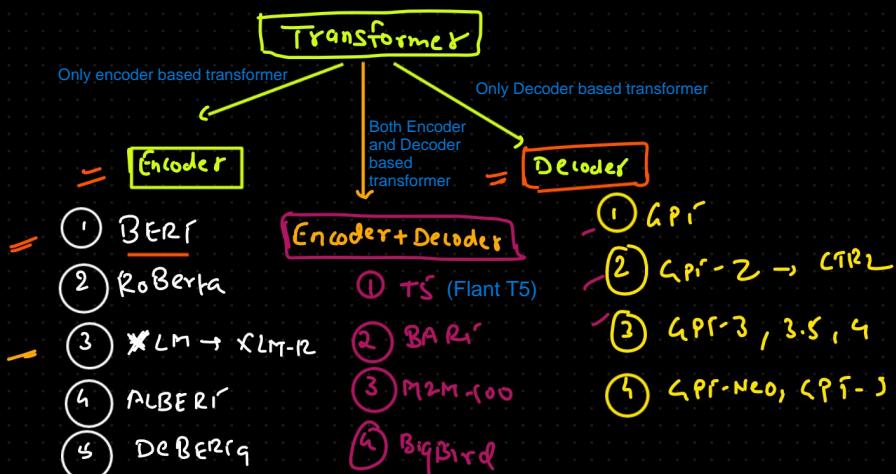
Embedding  
text-cleaning & rear reposition



Progression of NLP in chronological order:



If in the interview we are asked about how are we going to identify which transformer we would be requiring. Then tell them based on divisions mentioned alongside?

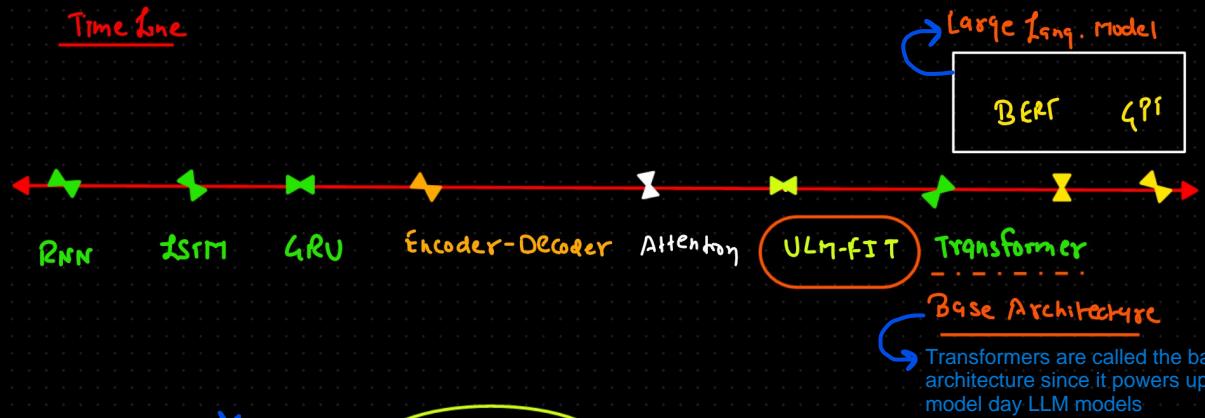


There are many LLM models. But the ones mentioned here are the milestones which even people are using in industries in Production

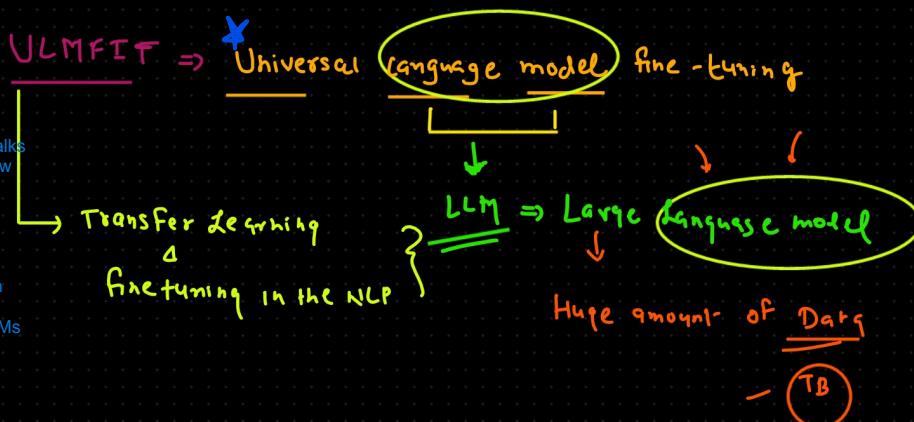
## 6 ELECTRA

## 5 Megatron (Megatron)

Chronological Progression in NLP domain



This research paper talks about how can we apply Transfer Learning and Fine Tuning in NLP by using LLMs

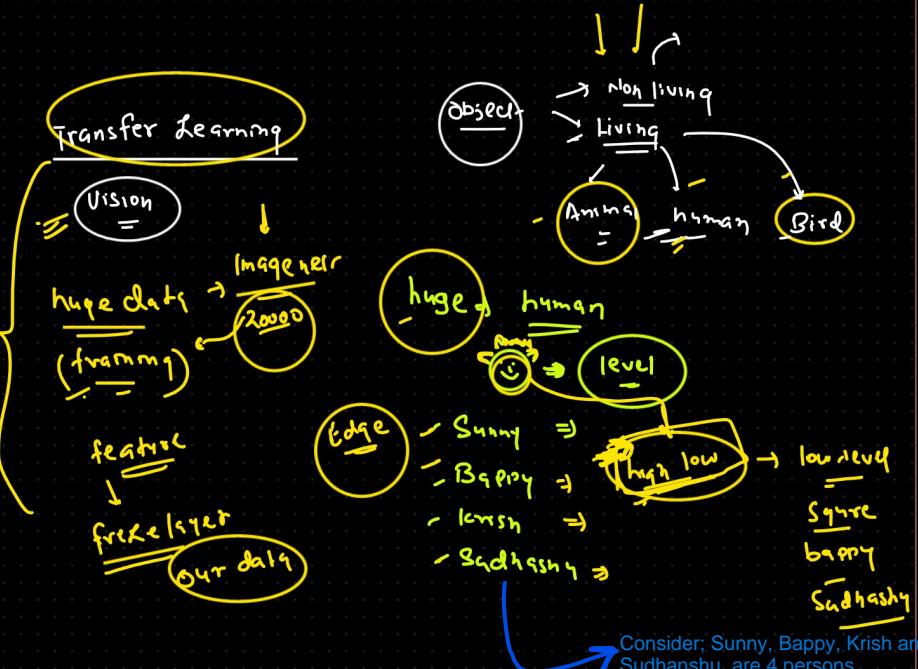


### Computer Vision

Transfer Learning in Computer Vision domain:

- 1 Image Classification
- 2 obj. Detection
- 3 Segmentation
- 4 OCR
- 5 tracking

We were able to apply Transfer learning on top of Different Tasks in the Computer Vision domain (like: Image Classification, Obj. Detection, Segmentation, OCR, Obj tracking). The main analogy behind this was that all the objects (Living or non living) have same high level features up to a point. So we can use the existing models to that point by freezing the most layers and later further fine-tuning on top of low level features on top of which model was not pre-trained.



Consider; Sunny, Bappy, Krish and Sudhanshu are 4 persons.

So here high level feature can be overall edges detection as a human which will be same for all the human entity. But when we wish to establish individual distinction then we have to consider low level features like face type (Square, round, oval etc.), nose type etc. In other words here we can use a common model to identify human entity and train the existing model on top of pictures of Sunny, Bappy, Krish and Sudhanshu to make a model capable to identify these persons differently. Here we were able achieve this using a single model since there very less diversity or high level features are common mostly.



Applying Transfer Learning in the NLP domain:

On the other hand NLP use cases have **diversity**. For example, if we consider Sentiment analysis and Text Generation then we will find that each of these are very different from each other and a **common model** cannot be used for the same. Hence, initially it was not possible to use any pre-trained NLP models for a particular use case and apply Finetuning and transfer Learning on top of it due to Diversity factor present in the NLP use cases.

Luckily, later a research paper was published coining the term **ULM-FIT** that gave an approach to apply finetuning and transfer learning.

NLP → **Diversity**

- 1 Sentiment Analysis
- 2 NER
- 3 Text Generation

**Common**

**finetuning**

**transfer learning**

**ULM-FIT** ⇒ **finetuning**  
**transfer learning**

Language modelling  $\Rightarrow$  Understanding the Language

Q: How ULM-FIT was able to apply transfer learning in NLP use cases?  
Ans: Using Language modeling which means getting the understanding of the language.

Big corpora  $\Rightarrow$  Next word Prediction

BERT (NSP, MLM)

Model of Language

Q: How the understanding of the language is obtained?  
Ans: Using next word Prediction approach. For Ex; Understanding of the Language using BERT LLM is obtained using NSP(Next sentence prediction) and MLM(Masked language modeling)

Next word prediction means for eg giving this much data to language based model

And then giving our model the power to predict the next word

Tesla, Inc. (/tɛslə/ TESS-lə or /'tɛzla/ TEZ-lə[a]) is an American multinational automotive and clean energy company headquartered in Austin, Texas, which designs and manufactures electric vehicles (cars and trucks), stationary battery energy storage devices from home to grid-scale, solar panels and solar shingles, and related products and services. Its subsidiary Tesla Energy develops and is a major installer of photovoltaic systems in the United States and is one of the largest global suppliers of battery energy storage systems with 6.5 gigawatt-hours (GWh) installed in 2022.

Tesla is one of the world's most valuable companies and, as of 2023, is the world's most valuable automaker. In 2022, the company led the battery electric vehicle market, with 18% share.

Tesla was incorporated in July 2003 by Martin Eberhard and Marc Tarpenning as Tesla Motors. The company's name is a tribute to inventor and electrical engineer Nikola Tesla. In February 2004, via a \$6.5 million investment, Elon Musk became the company's largest shareholder. He became CEO in 2008. Tesla's announced mission is to create products which help "accelerate the world's transition to sustainable energy."

The company began production of its first car model, the Roadster sports car, in 2008. This was followed by the Model S sedan in 2012, the Model X SUV in 2015, the Model 3 sedan in 2017, the Model Y crossover in 2020, the Tesla Semi truck in 2022 and the Cybertruck light-duty pickup truck in 2023.[7] The Model 3 is the all-time bestselling plug-in electric car worldwide, and in June 2021 became the first electric car to sell 1 million units globally.[8]

Google, MEG, MS operate

BERT

GPT

XLm

Chatgpt

Google  
3.5  
Gemini

RHF

huge amnt data

Language modeling + fine tuning

+

Chatgpt

Google  
3.5  
Gemini

RHF

UNSUPERVISED finetuning

BERT  
base

Understanding  
the language

{huge  $\rightarrow$  Pred. the near  
word}

Specific task

- 1 sentiment analysis
- 2 QA
- 3 Summarization

Supervised  
finetuning

Q: What is Language modeling?  
Ans: Language modelling is the way through which the language based models can be made to predict the next word by training it on top of huge data (supervised learning) which enables the language based models(LLMs) to get the understanding of the language by capturing complex patterns within it.

Unsupervised Finetuning means directly passing the data for inference or prediction to the model and using model's pre trained weights to generate the output

Transformer (encoder/Decoder)

- 1 huge Data
- 2 Unsupervised finetuning
- 3 Supervised finetuning
- 4 RHF

$\Rightarrow$  Chatgpt, Bard, Gemini

Supervised finetuning means updating some of the pre trained weights of LLMs by finetuning/training the existing model using dataset specific to the use case that we are trying to solve like: Sentiment analysis, question answering, Summarization etc.

ULMFIT → Transfer Learning & fine tuning 2018 → Published year

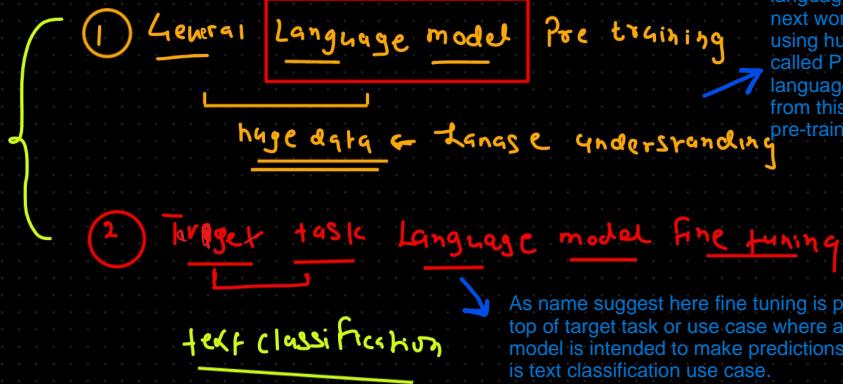
text classification → Here transfer learning and Fine tuning is done on top of text classification use case

① Model Architecture

- ① RNN  
② LSTM

Instead of using transformers this used the traditional RNN and LSTM to build the language model which can be enabled for Transfer learning and Fine tuning

② Fine tuning Methodology



In this language modeling is done meaning enabling language model to predict next word by training it using huge data. Also called Pre-training and the language model obtained from this is called pre-trained model

As name suggest here fine tuning is performed on top of target task or use case where a pre-trained model is intended to make predictions. In ULMFIT it is text classification use case.

Please note that:

- Through General Language model Pre Training we are enabling the language model for Transfer Learning.
- Whereas, through Target task language model fine tuning we are enabling the language model for Fine Tuning.

Similarly, decoding the BERT research paper:

BERT ⇒ ① Model

Unlike ULMFIT BERT is build using Transformer(Encoder Only Transformer)

Transformer → Encoder

BERT  
BASE

12 Encoder

BERT  
Large

24 Encoder

These are the two variants of BERT model

② Fine tuning methodology

Language modeling is used for Pre-training or building the Pre-trained BERT model.

① Language modeling ⇒ Pre-training ⇒ {Predicting the next word in sentence)  
↓  
Understanding of the language

② target - task ⇒ Sentiment classification

According to the research paper Target task language model fine tuning is done on top of Sentiment analysis use case.



Some common terms related to the BERT research paper:

1 BERT  $\Rightarrow$

Bi-Directional Encoder Representation from transformer



Both Direction

2 Prese

2 Transformer  $\Rightarrow$  Self attention, skip connection, Norm, NN  $\Rightarrow$  Encoder

3 Pre-training  $\Rightarrow$  init Phase  $\rightarrow$  Large Corpus  $\rightarrow$  Understanding [next word pred]

4 Fine-tuning  $\Rightarrow$  Task Specific Data  $\rightarrow$  train

5 MLM  $\Rightarrow$  Mask Language modeling

6 NSP  $\Rightarrow$  Next sent Pred

7 Contextual embedding  $\Rightarrow$  position Embed

8 CLS Token  $\Rightarrow$  init token  $\Rightarrow$  classification

9 SEP token  $\Rightarrow$  Segmentation

If time permits read about the ULMFIT, GPT, BERT research papers from Jay Alammar's blog



Difference b/w BERT and GPT language models:

Feature	BERT	GPT
Training obj.	MLM (Masked Language modelling) Predict the masked based on context  Masked word here represents the next word which we are intend to predict	Autoregressive Language Model: Predicting the next word based on the preceding context.
Architecture	Transformer $\Rightarrow$ Encoder with Bi-Directional context understanding	Transformer $\rightarrow$ Decoder Unidirectional context generation.
Context understanding	Bi-directional context	Unidirectional context left-to-right $\Rightarrow$ next word
Usecase	text classification NER QA  Powerful { Huge amr }	Text generation Story completion Code generation
Pre-training		huge variety <sup>and</sup> of Data $\Rightarrow$ 175B

## Fine-tuning

{ Commonly fine-tuned  
for specific task }

it can be fine-tune



## Token masking

Uses mask lang. model

- masking of the token

Here masking of tokens (Input) is done and the same is used to draw prediction based on the context.



## Model Size

BERT BASE

BERT Large

Does not perform token masking  
(Auto-regressive)

Doesn't require masking instead it is a autoregressive model meaning it can generate the next word by using the context of the preceding words.

GPT-1  
GPT-2  
GPT-3, 3.5, 3.5+4.860  
GPT-4  
GPT-4V

Mathematical  $\Rightarrow$  transformer

encoder

decoder

10% of 100  
10  $\Rightarrow$  BASIC  
Predict

X Understanding the mask language modelling concept used in the BERT:

I am sunny I am a ~~DATA~~ who love to work with ML, DL,  
NLP etc. I know the ~~Big data~~ I know how to ~~design~~

Datq Pipeline.

Predict

So, this masking is applied during pre-training the model.

During pre-training 10-15% of a corpus or single sentence is masked(Hidden). Then model is enabled to perform the prediction of the masked words by using the context captured from the preceding and succeeding words.



Ques: Say for example your manager asked you to implement question and answering use case. Also, you are provided with very less or limited resources. So which variant of BERT will you use or how are you going to approach this problem statement?

Ans: Well there are different variants of BERT like: ALBERT(A Lite BERT), RoBERTa, ELECTRA, DistilBERT, TinyBERT. As we know BERT language based model can be used for solving the mentioned usecase - Question and Answering. So, firstly we are going to read all the research papers describing the different variants of BERT. Then we are going to see which one consists of less parameters also we are going to see the Evaluation or Benchmark evaluation which will basically help us in understanding that how well the chosen variant will work on top our use case.

One can also google about the different BERT variants, read the different blogs and then take a reasonable decision.