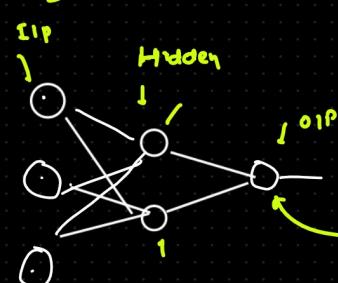


GEN-RI History

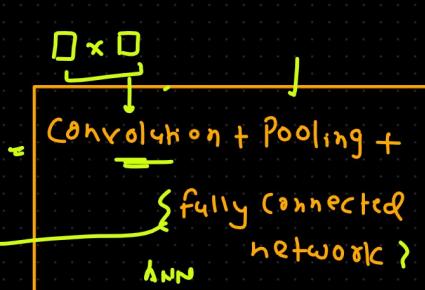
- 1) Language Modelling
- 2) RNN → Transformer
- 3) LLM
- 4) BERT
- 5) GPT
- 6) char-GPT training
- 7) Transfer Learning & Fine tuning NLP

DL Type of NN [NN is a fundamental unit of DL]

1) ANN



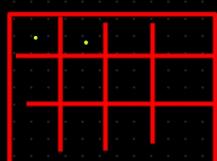
2) CNN



3) RNN

ANN + Feedback
loop

- Structure Data



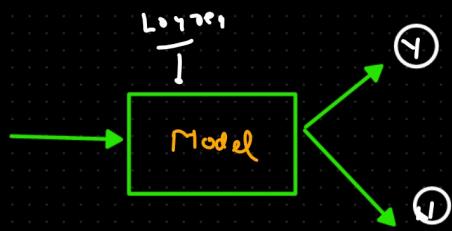
Col - Num
Row - Ord.
cont - Disc.

- Image / videos

Sequence data

Text / Audio

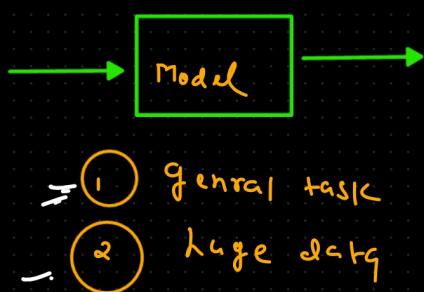
Discriminative Model



- 1 Specific task
- 2 Limited data

v/s

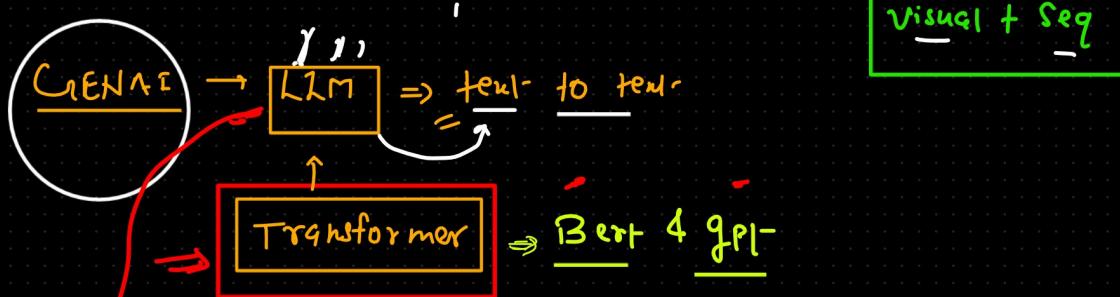
Generative Model



- 1 General task
- 2 Huge data

Generative Model Type

- 1 image to image
- 2 text-to-text } homogenous Model
- { 3 image-to-text
- { 4 text-to-image } heterogeneous Model or (Multi-modal) or (Diffusion model)



Visual + Seq

Vision based task also (CRAFT4V, GeminiPPO, Donut, CLIP, Whisper)

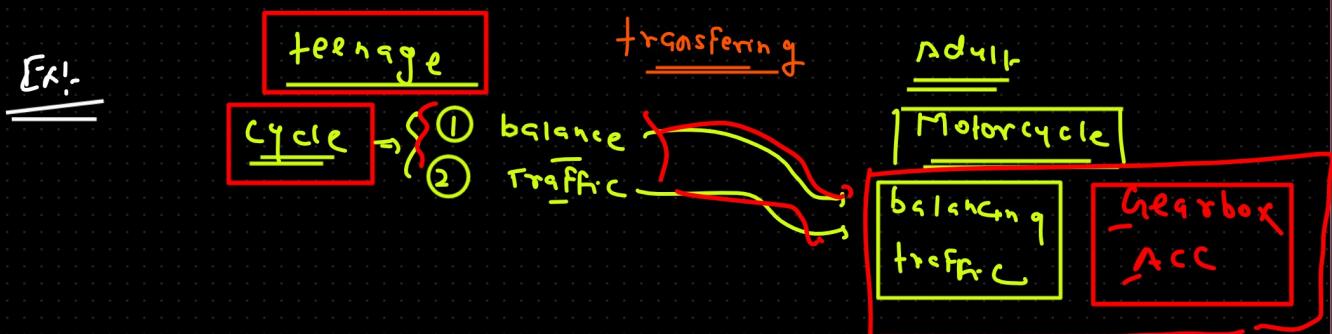


- 1 Efficient Det.
 - 2 VGG
 - 3 Resnet
 - 4 DenseNet
- 1 YOLO
 - 2 SSD
 - 3 fasterRCNN
- 1 Mask R-CNN
 - 2 U-net
 - 3 V-net

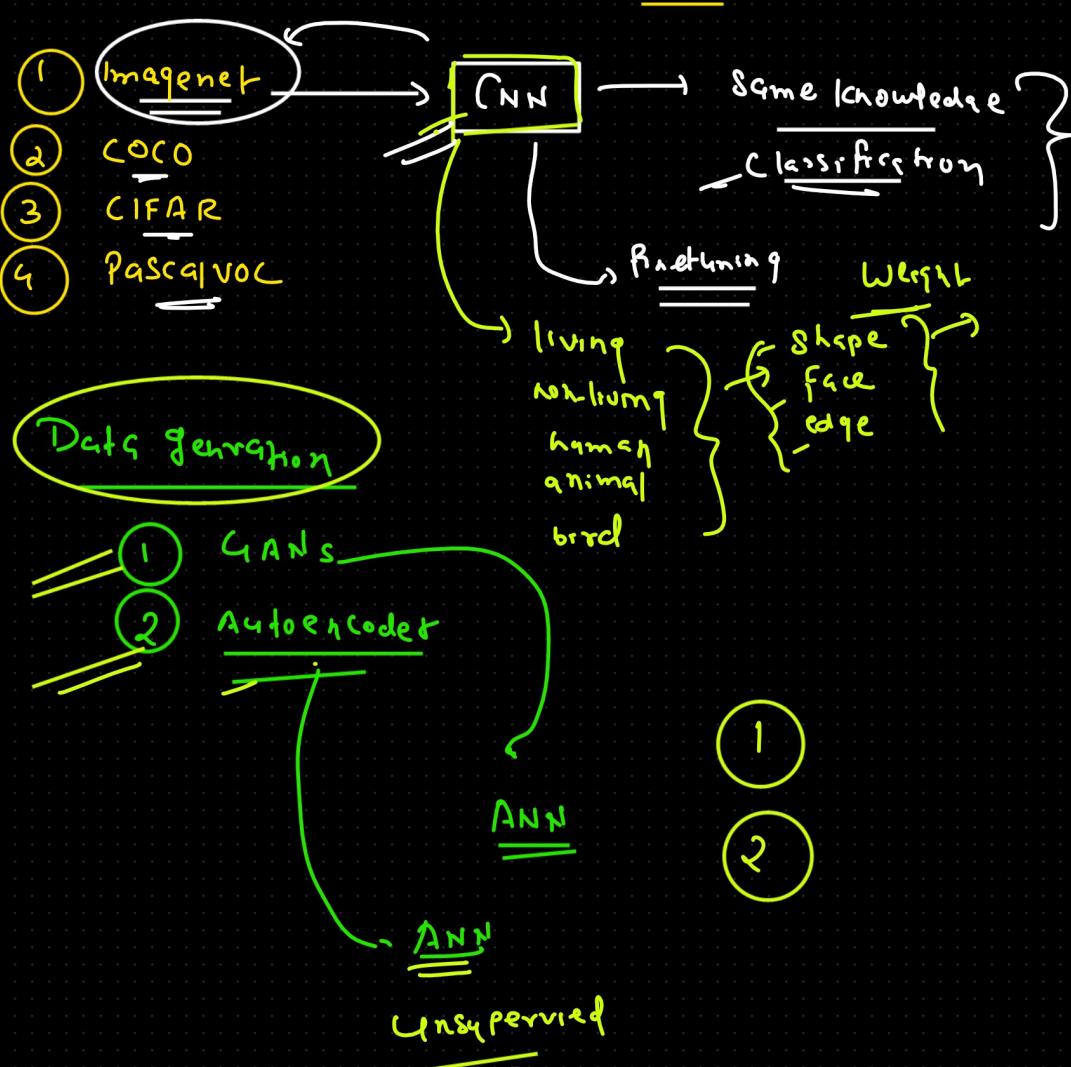
Transfer Learning & Finetuning

Use the Previous Learning for the future task

Modified Previous Learning According to requirement

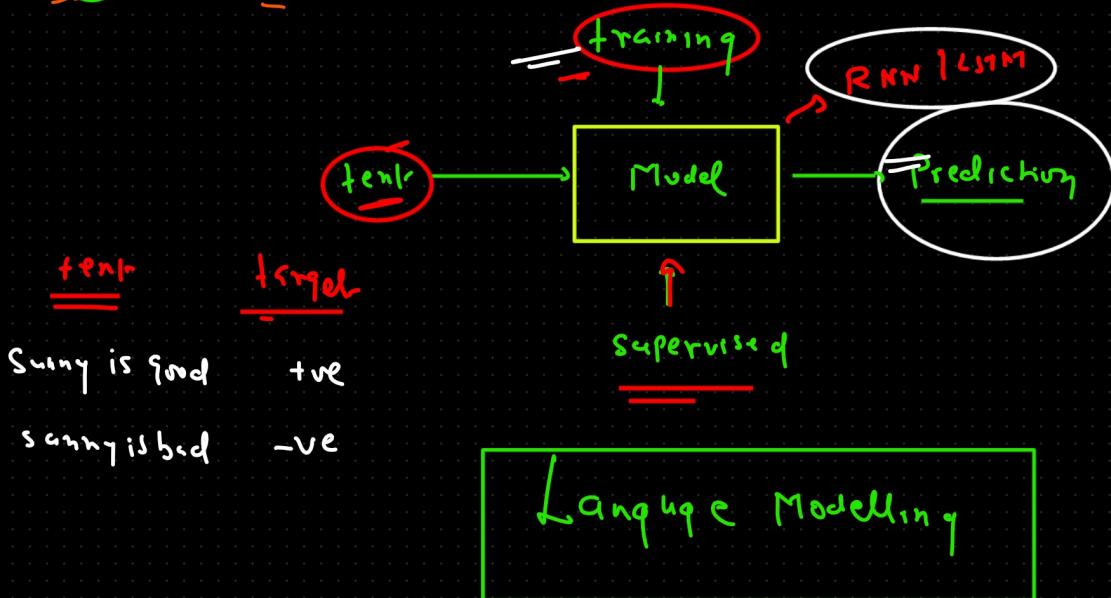
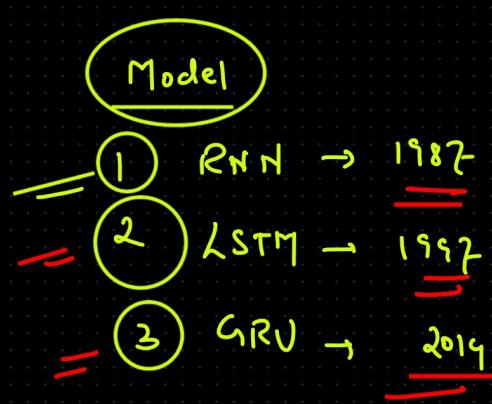


Transfer Learning & Finetuning (CNN) (Computer vision)

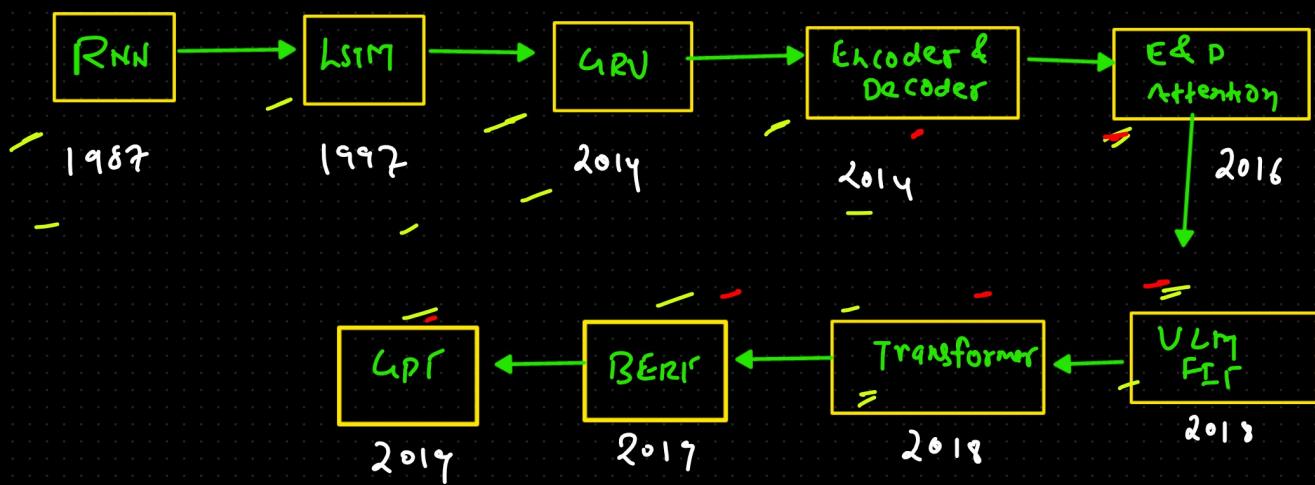


NLP \Rightarrow Sequence data, text data

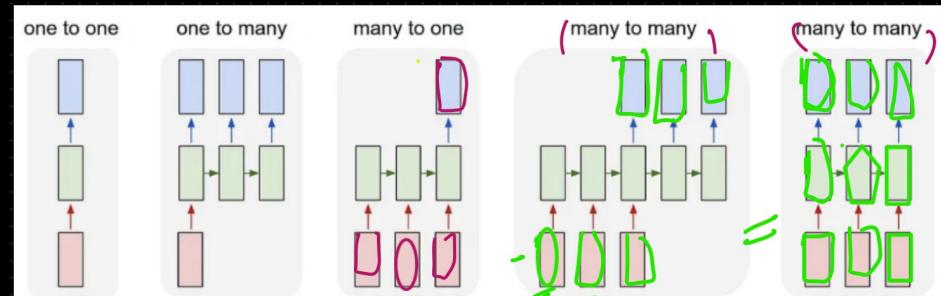
- 1 Summarization
- 2 Generation (word, sentence)
- 3 Translation
- 4 Classification



Complete timeline of the evolution of LLM



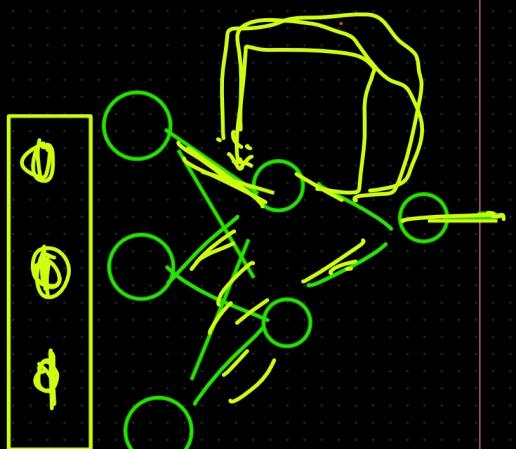
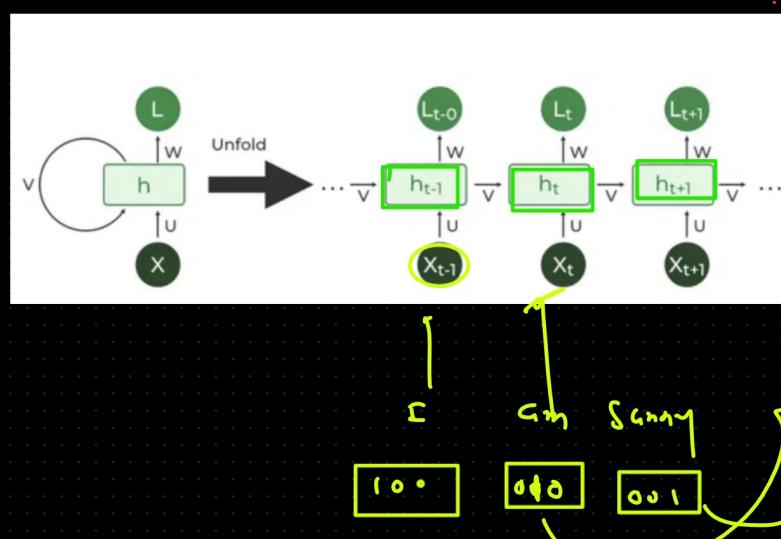
① Mapping (Seq to seq mapping)



↑
image classification
Image caption (CNN + RNN)
Sentiment (RNN, LSTM)

{ Language translation, NER, POS
Same length
Diff length
→ *Korean* *한국*
| *am* *sunny*
↓ *shut, day* and
such off the fc,

② RNN



Working ⇒ time stamp

Segmental Process

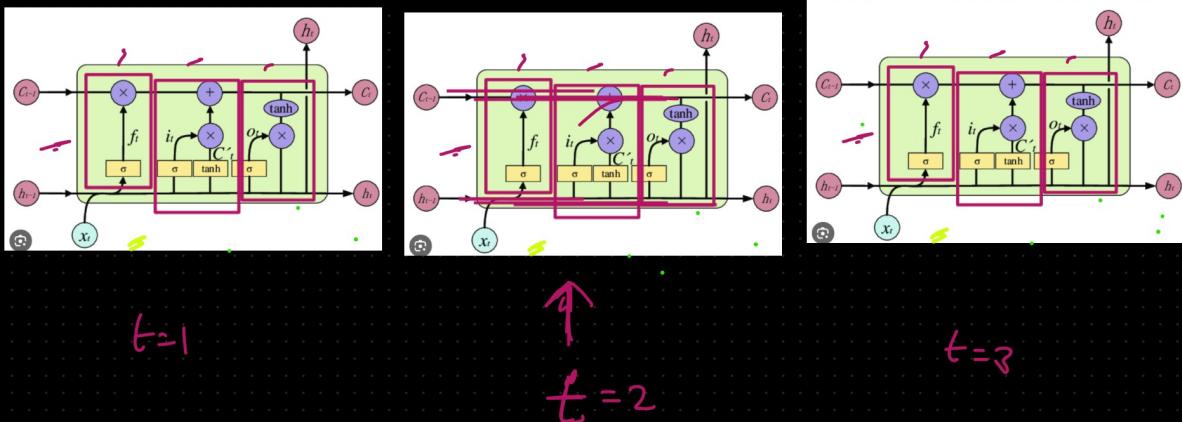
optimization

vanning gradient
exploding ..

Advantage ⇒ Seq

Disadvantage ⇒ long context

3 LSTM (Long short term memory)



little longer sentence

= long - short

Sunny, Sandeep

2000-2010

Sandeep
—> [Developer]

[mentor]

2020
—> [Sunny]
—> [Mentor] → [A2]

Who is the father:

12NN

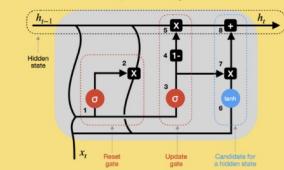
LSTM

4

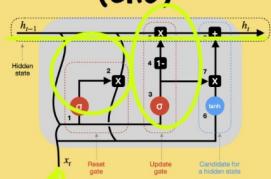
GRU (Gated Recurrent Unit)

 $t=0$

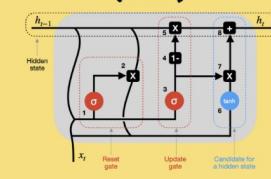
GATED RECURRENT UNIT (GRU)

 $t=1$

GATED RECURRENT UNIT (GRU)

 $t=2$

GATED RECURRENT UNIT (GRU)

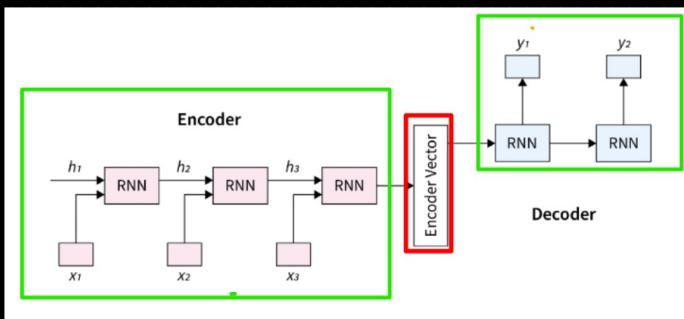


Feature	LSTM	GRU
Architecture Complexity	More complex with three gates (input, forget, output) and a separate memory cell.	Simpler with two gates (update, reset) and a combined hidden state/memory cell.
Memory Handling	Separate memory cell explicitly manages what information to keep or discard.	Combines hidden state and memory cell into a single state vector.
Gating Mechanism	Explicit forget, input, and output gates for controlling information flow.	Single update gate for combining past and new information, and a reset gate.
Parameter Count	Higher due to the additional parameters from the complex architecture.	Lower, as GRUs have fewer parameters, making them computationally more efficient.
Training Convergence	May require more time to converge during training due to the complexity.	Often converges faster during training, making them easier to train in some scenarios.

Choosing between LSTM and GRU depends on the specific task, dataset, and computational resources. LSTMs might be more effective in capturing longer dependencies due to their more intricate structure, but GRUs often require fewer parameters and may be computationally more efficient in certain scenarios. It's common practice to experiment with both architectures to determine which one performs better for a particular use case.

6

Encoder & Decoder



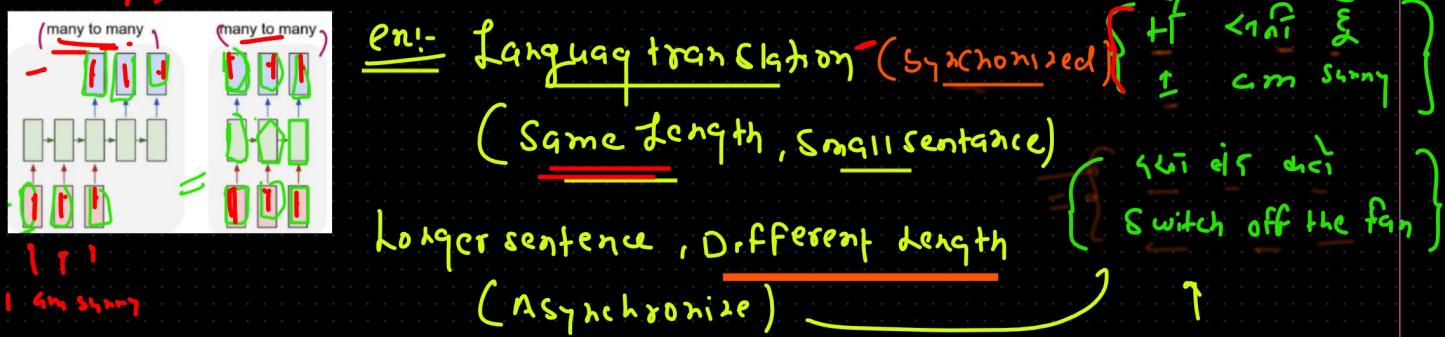
1 Encoder

2 Decoder

3 RNN / LSTM / GRU → LSTM \Rightarrow Bi-Directional

4 Encoder → Content vector → Decoder

Why \rightarrow Seq to seq mapping (many to many) - (Lang. translator)



Research only \rightarrow

30 - 35 words -

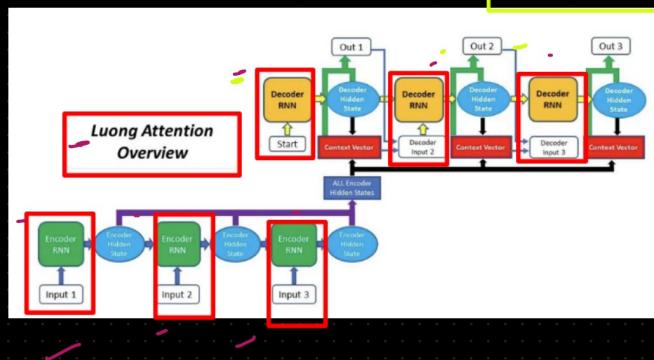
It was fine otherwise giving the issue

Who was using it \rightarrow Google was using it in their earlier days

What was the solution \Rightarrow E&D with Attention

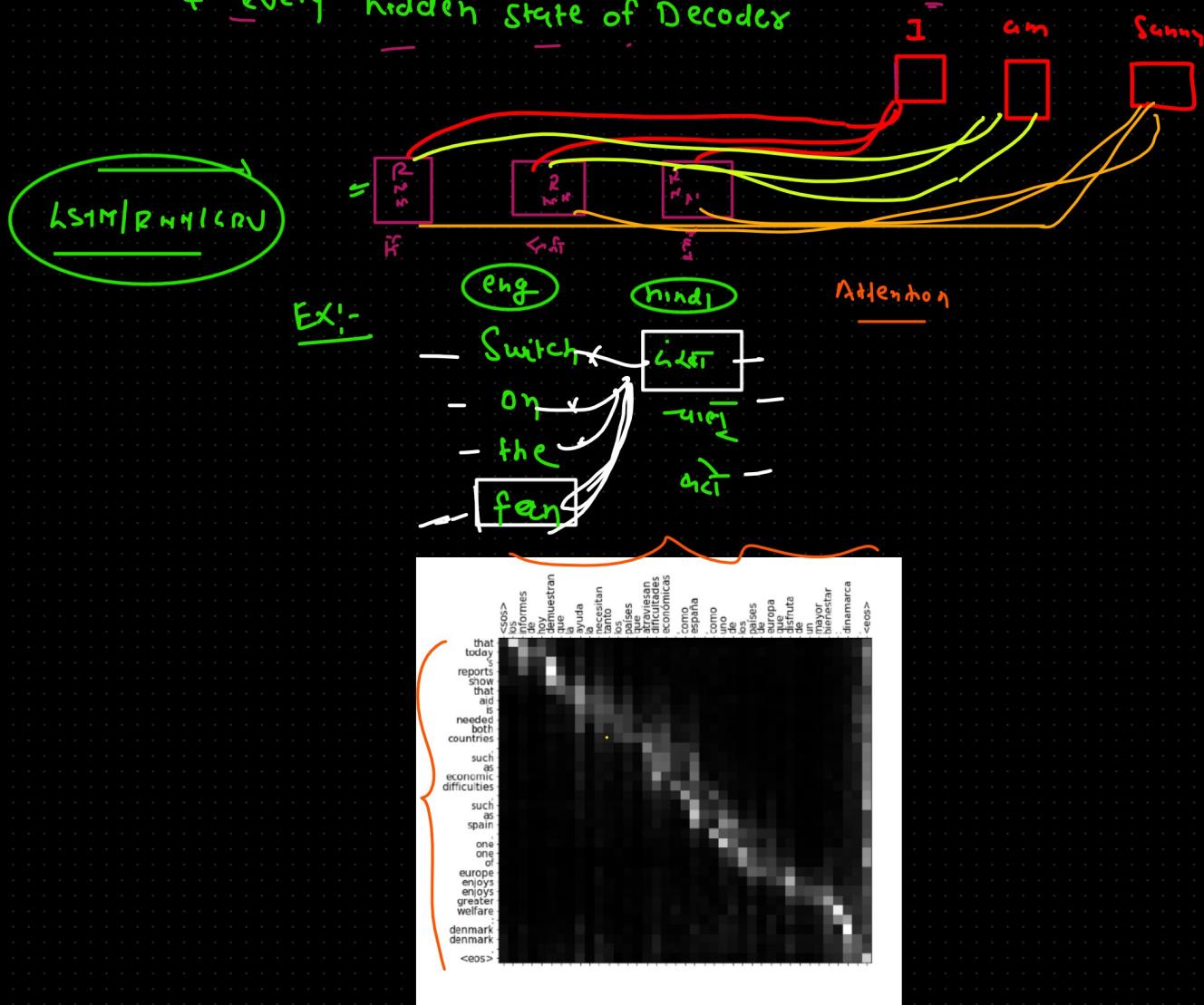
7

Encoder & Decoder with attention



1 Not only single Content vector

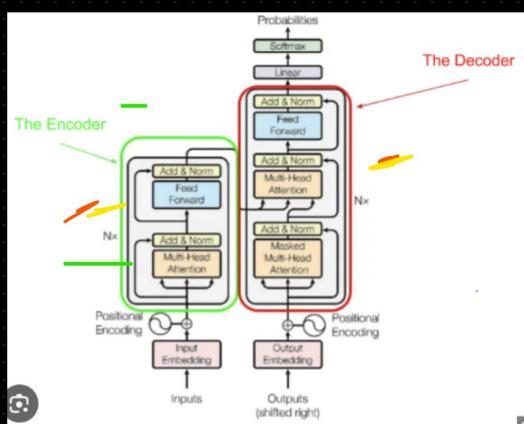
2 Connecting each & every hidden state of encoder with each & every hidden state of Decoder



Intensity of mapping

More white → strong relation Black → No relation
Less white → weak Relation

⑧ Transformer (Attention all you need)



- 1 Remove LSM/RNN
- 2 Self attention → (Multihead attention)
- 3 Positional encoding
- 4 Parallel processing
- 5 Normalization
- 6 Artificial Neural Network
- 7 Skip Connection
- 8 Application beyond NLP

thing required

- 1 Data ⇒ huge data
- 2 Hardware ⇒ GPU, Distributed
- 3 Time

9 ULMFIT =>

- 1 Universal Lang. modelling
- 2 Fine tuning for text classification

Universal Language Model Fine-tuning for Text Classification

Jeremy Howard^{*}
fast.ai
University of San Francisco
^{*}jef@fast.ai

Sebastian Ruder^{*}
Insight Center, NUI Galway
Alylou Ltd., Dublin
sebastian.ruder@tudor.io

Abstract

Inductive transfer learning has greatly impacted computer vision, but existing approaches in NLP still require task-specific modifications and training from scratch. We propose Universal Language Model Fine-tuning (ULMFiT), an effective learning method that can be applied to any task in NLP and introduce techniques that are key for fine-tuning a language model. Our method achieves state-of-the-art performance on six text classification tasks, reducing the error by 18-24% on the majority of datasets. Furthermore, on 100 labeled examples, it matches the performance of training from scratch on 100x more data. We open-source our pretrained models and code.

1 Introduction

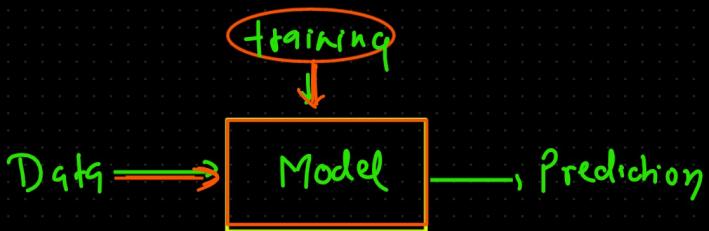
Inductive transfer learning has had a large impact on computer vision (CV). Applied CV models (including object detection, classification, and segmentation) are typically trained from scratch, instead are fine-tuned from models that have been pretrained on ImageNet, MS-COCO, and other datasets (Shafiq Razavian et al., 2014; Long et al., 2015).

arXiv:1801.06146v5 [cs.CL] 23 May 2018

- Question
- 1 Does TL Possible in NLP?
 - 2 ... FT " " "
 - 3 How?

Language modelling =>

- 1 classification
- 2 generation
- 3 summarization
- 4 translation



Supervised Learning
Supervised Language modelling
(teaching to my model)

LLM
Descriptive modelling ✓ Less data, Specific task
Generative modelling ✗ More data, General task

Unsupervised Language modelling

Huge data → Label not required

TASK => Data generation (Next-word Prediction)

- How:
- 1 Statical Pattern (Dist.)
 - 2 Semantic Word
 - 3 Contextual relationship

(word-word)

- 1 So that Model can learn
- 2 based on Previous word can predict next word

I am sunny seville
I am sunny | am sunny seville
sunny ?

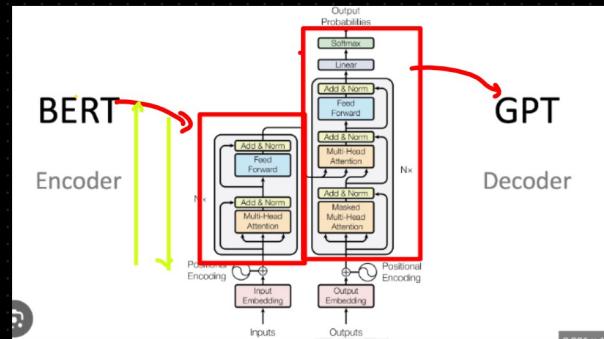
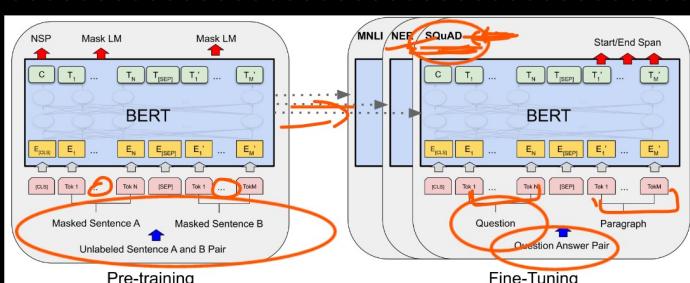
translation

I am sunny | hi can i
= = = = =

Supervised lang. Modelling

Label will be req

10

BERTLanguage Modelhuge amountULMGoogleteching(next word Pre)1 Encoder only model2 MLM3 NWP, NSP4Language Model1 Unsupervised Pretraining2 Supervised fine tuning5Bi-Directional Model =>6Use-Case =>1 Sentiment analysis2 NER3 POS4 question Answering

my name is [REDACTED] scavite.
 am [REDACTED] who is working
 in banglore.

11

GPT \Rightarrow Generative Pre-training

GPT1 :Improving Language Understanding by Generative Pre-Training

Generative Pre-training:

The term "generative" refers to the model's capability to generate coherent and contextually relevant text. "Pre-training" indicates that the model is initially trained on a large corpus of text data in an unsupervised manner before being fine-tuned for specific tasks. During the pretraining phase, the model learns to predict the next word in a sequence, capturing language patterns and semantics.

Key Objectives:

The paper outlines the methodology of training a transformer-based language model on a diverse range of tasks without task-specific labelled data.

The generative pretraining approach allows the model to acquire a broad understanding of language, enabling it to perform well on various downstream tasks.

word word

\rightarrow VLM FT

GPT2: language models are unsupervised multitask learner

In the context of the GPT-2 (Generative Pre-trained Transformer 2) paper, the term "unsupervised multitask learner" refers to the model's ability to perform a variety of natural language processing (NLP) tasks without task-specific supervision during the pretraining phase. GPT-2 is designed as a large-scale language model that is pretrained on a diverse corpus of text data without explicit annotations for specific tasks, given to it at runtime

\Rightarrow VLM \rightarrow huge

GPT3: language model are fewshot learner

In the context of the GPT-3 research paper, the term "few-shot learner" refers to the model's ability to perform a task or answer questions with only a few examples or prompts provided during inference. GPT-3 is known for its remarkable ability to generalise from a small number of examples or demonstrations given to it at runtime.

\Rightarrow very huge

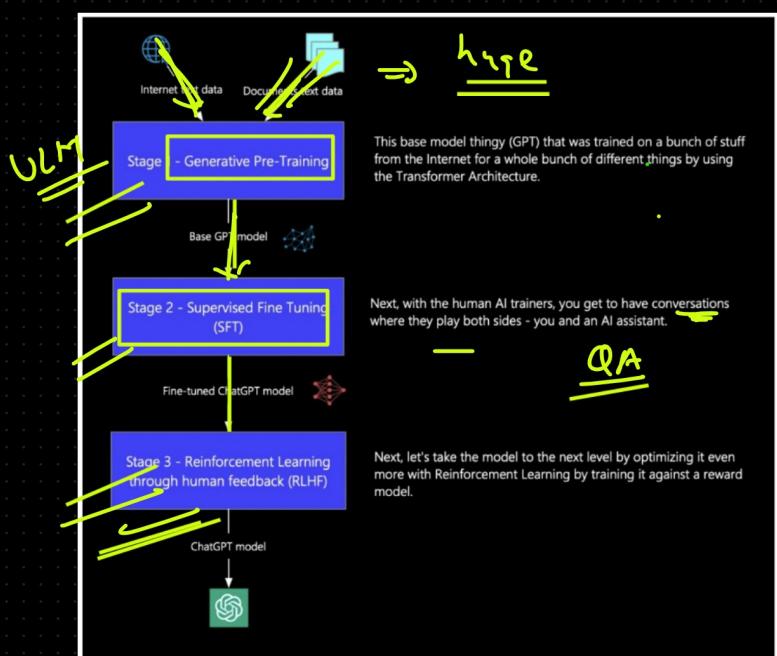
175B

few shot \rightarrow [] \rightarrow

Aspect	GPT-1	GPT-2	GPT-3
<u>Model Size</u>	117 million parameters	1.5 billion parameters	Up to 175 billion parameters
<u>Context Window Size</u>	Limited due to model size	Larger context window	Extremely large context window
<u>Fine-Tuning Capabilities</u>	Limited fine-tuning capabilities	Improved fine-tuning capabilities	Enhanced fine-tuning capabilities
<u>Released Date</u>	June 2018	February 2019	June 2020
<u>Prompt Engineering</u>	Limited influence on model behavior	More influence on model behavior	Substantial influence on model behavior
<u>Downstream Task Performance</u>	Less impressive on specific tasks	Improved performance on diverse tasks	Exceptional performance on wide range
<u>Ethical Considerations</u>	Earlier models raised ethical concerns	Increased awareness and scrutiny	Continued focus on addressing biases
<u>Research Impact</u>	Pioneered large-scale unsupervised learning	Advanced understanding of transfer learning	Broader applications and benchmarks
<u>Use Cases and Applications</u>	Limited by model size and capabilities	Broader applications in various domains	Extensive applications across industries
<u>Public Access</u>	Limited access initially	Increased access with GPT-2	OpenAI API access for GPT-3

12

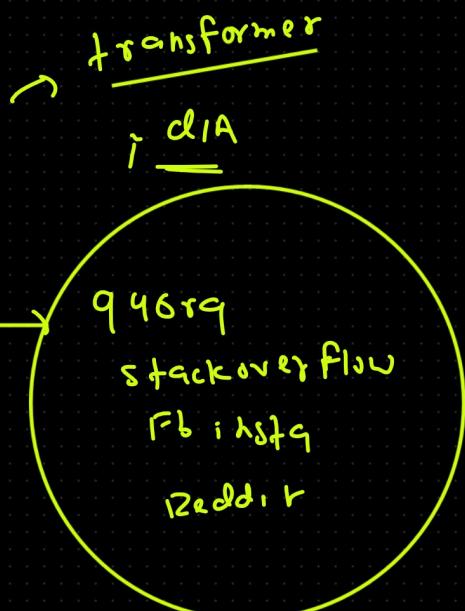
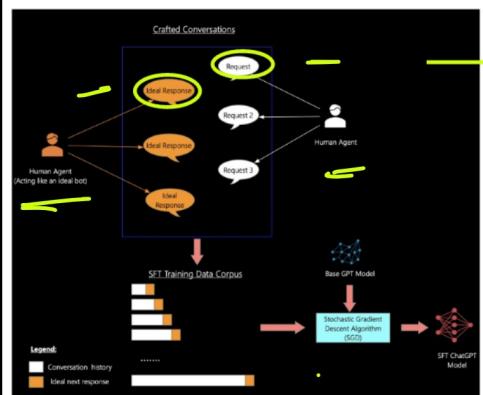
Training of ChatGPT \Rightarrow Application



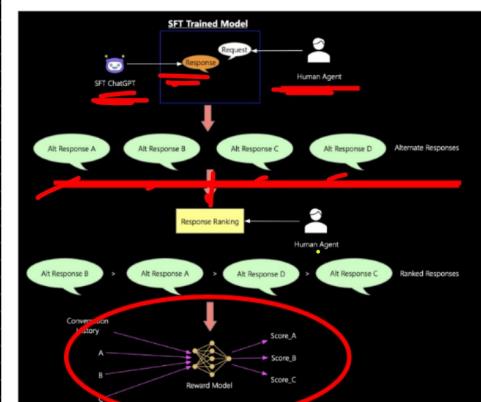
$\frac{\uparrow}{\downarrow}$
GPT-3

3.5, 3.5 turbo, GPT-4, 4 turbo

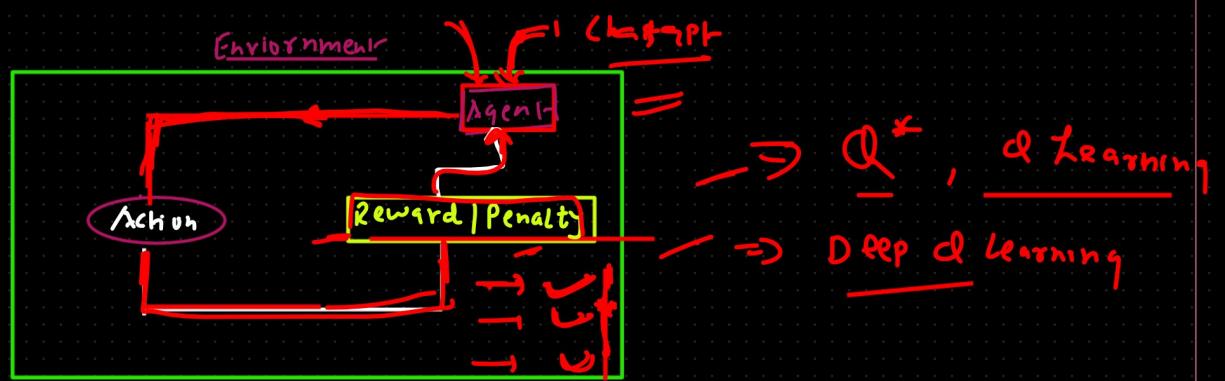
Supervised Fine-Tuning (SFT)



Reinforcement Learning through Human Feedback (RLHF)



RL HF



Imp Points

- 1 Ethics
- 2 Maintaining Memory content
- 3 Differentiate specific training \Rightarrow Conversation
- 4 RLHF (iterative Feedback)
- 5 Bias / error

Conclusions \rightarrow

