

Large Language Models (LLMs) concept

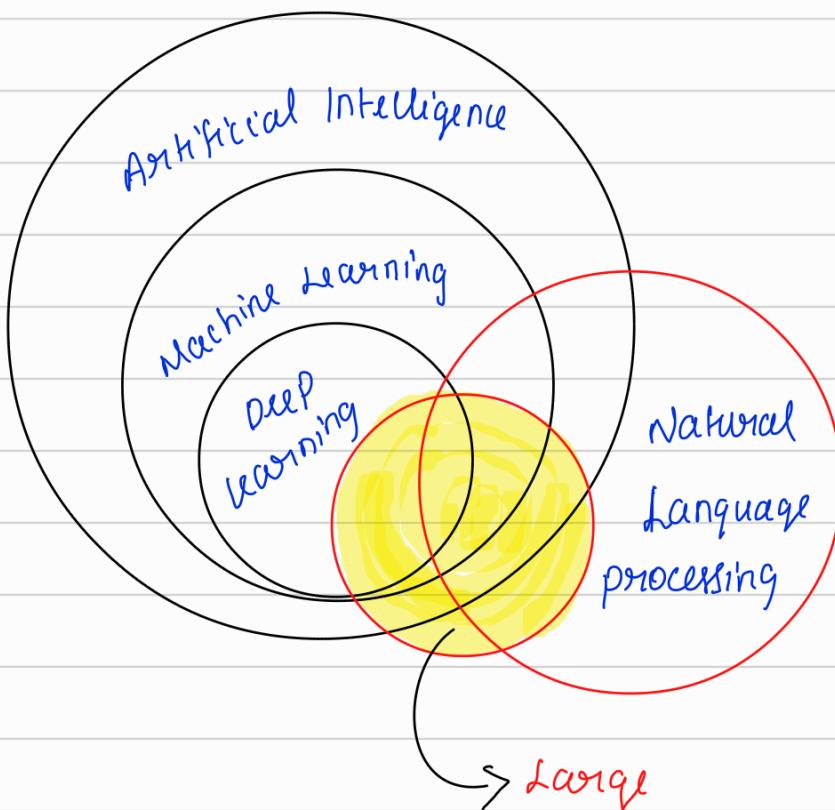
Introduction to LLM

AI performed well in data driven tasks.

(sentiment analysis and fraud detection)

Still, lacked human-like interaction.

So, Large language Model came.



Definition

- * Large (Training data & compute power)
- * Language (Human-like text)
- * Model (Learn complex patterns using text data)

Applications

- Sentiment analysis
- Identifying themes
- Translating text / speech
- Generating code.
- Next word prediction.

Real world application

Business Opportunities



Business constantly seek new and innovative ways to improve their products and services using LLM.

Benefits.

- Automate tasks
- Improve efficiency
- Create revenue stream
- Enable new capabilities.

Transforming finance industry

Investment
outlook

Annual
report

news
articles

Social media
posts



processing a unstructured data
(data that lacks definition and is presented free form)

LLM can analyse such data.



challenges in Healthcare

Analyzing Health records is important for giving personalized recommendations to provide quality healthcare.

Analyze patient record to offer personalized recommendation.

Medical record

→ Must adhere to privacy laws
(data)

Health checkup

(imaging report)

→ LLM

personalized treatment recommendation.

Education

- personalized coaching & feedback.
- Interactive learning experience

- AI-powered tutor
 - (Ask questions)
 - (Receive Guidance)
 - (Discuss Ideas)

- personalizing education (text generation)

Multimodal

- Many types of processing or generation (text & audio & image)
- visual question answering is one such multimodal application (Answers to questions about visual content)

- Object identification & relationship
- Scene description
- Responds with additional images

challenges of language Modelling

- ① sequence matters
- ② context modeling.
- ③ Long range dependency (Recognize and connect distant words in a sentence)
 - challenging for traditional language models.

Single task learning (Traditional model)



Time and resources expensive
less flexible compared to modern LLMs.

Multitask learning (Modern LLM)

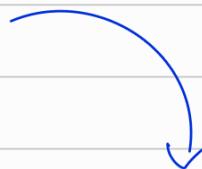
- Improved performance on each individual tasks.
- might impact accuracy and efficiency.
- less training data needed because data is shared.

Novelty of LLMs (unique capabilities)

→ NLP techniques are the foundations of LLMs.

→ Linguistic subtleties.

(irony, humor, intent, sarcasm...)



overcome

data's

unstructured

nature.

How do LLMs understand:

→ Trained on vast amount of data.

→ Large size of LLMs (parameters)

→ parameters represent the patterns and rules

→ more parameters → complex patterns

→ generates sophisticated and accurate responses.

Emergence of new capabilities

* Emergent abilities

(only present in large scale model)

→ Music

→ Poetry

→ code generation

→ Medical diagnosis & treatment

* Scale,

Volume of training data,

Number of Model parameter,

plans

Building Block of LLMs

→ Text pre-processing

→ Text representation

→ pre-training

→ fine-tuning

→ Advanced fine-tuning

Generalized overview of NLP

① Text pre-processing

* tokenization

* Lemmatization

* Stop word removal.

can be done in different order.

Splits text into individual words or token

Stop words do not add meaning & eliminated

* Group slightly different words with similar meanings

* Reduces words to their base form.

* Mapped to root word.

2. Text Representation.

(Text data into numerical form)

* Bag of words

* word embeddings.

Text into a matrix of word counts.

'0' represents the absence of a word.

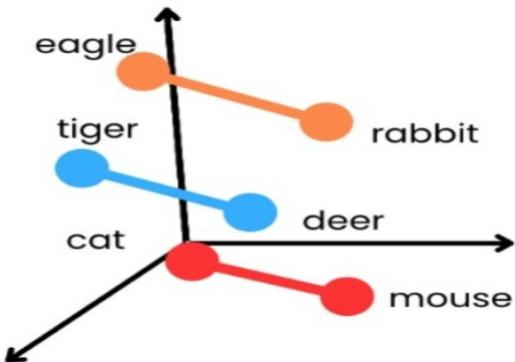
Limitation

→ does not capture the order / context

→ does not capture the semantics between the words.

→ capture the semantic meanings as numbers.

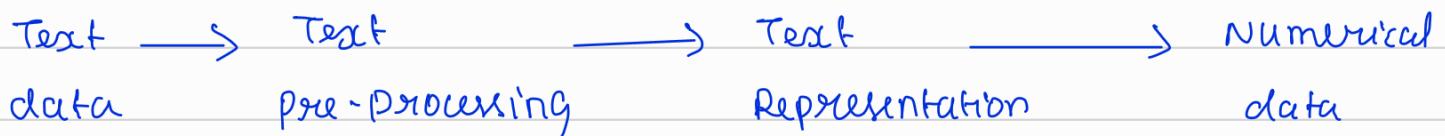
Predator - prey relation



Machine Readable form

start with text
pre-processing

convert pre processed
text into numerical
form.



③ pre-training.

Not everyone needs to train an LLM from scratch

- Pre-training



School education

- Fine-tuning



University specialization

Largeness challenge

- * fine-tuning can help
- * powerful computers
- computing power.

- memory
- processing power
- infrastructure
- Expensive.
- LLM, 10K of GPU & GPU

- * Efficient model training methods
- * Large amounts of training data.

Efficient model

- training.
 - Training time is huge
 - may take weeks or even months
 - Efficient model training = faster training time.
 - 355 years of processing time on a single GPU.

Data availability

- need of high-quality of data
- to learn the complexities & subtleties of language.
- few 100's GB of text data.
- Massive amount of data.

Overcoming the challenges

- **Fine-tuning**
 - Addresses some of these challenges
 - Adapts a pre-trained model
- **Pre-trained model**
 - Learned from general-purpose datasets
 - Not optimized for specific-tasks
 - Can be fine-tuned for a specific problem



Fine-tuning.

- * Compute (1-2 CPU & GPU)
- * Training time
(Hours to day)
- * Data (~1 GBs)

pre-training.

- * Compute
(1000s of CPU & GPU)
- * Training time
(Weeks to Month)
- * Data (~100 of GBs)

Learning techniques

fine-tuning, training a pre-trained model for a specific task.

Transfer learning

→ Learn from one task and transfer to related task.

N-shot learning

→ zero shot, few shot and multi-shot learning.

Building Blocks

- Data preparation workflow
- Fine tuning
- N-shot Learning techniques

Generative pre-training

- * Trained using generative pre-training
 - Input data of next token
 - Trained to predict the token with the dataset.
- * Types:
 - Next word prediction
 - Masked Language modelling.

Next word prediction

- * Supervised Learning technique
(Model trained on input-output pairs)
- * Predicts next word and generates coherent text and captures the dependencies between words.
- * Training data (pairs of input & output examples)
- * More examples = better prediction.

Training data for next word prediction

Input	Output
The quick brown	fox
The quick brown fox	jumps
The quick brown fox jumps	over
The quick brown fox jumps over	the
The quick brown fox jumps over the	lazy
The quick brown fox jumps over the lazy	dog
The quick brown fox jumps over the lazy dog.	

Masked Language Modelling

- Hides a selective word
- Trained model predicts the masked word.

Original Text: "The quick brown fox jumps over the lazy dog."

Masked Text: "The quick [MASK] fox jumps over the lazy dog."

- Objective, predict the missing data
- Based on learnings from training data.

Transformers

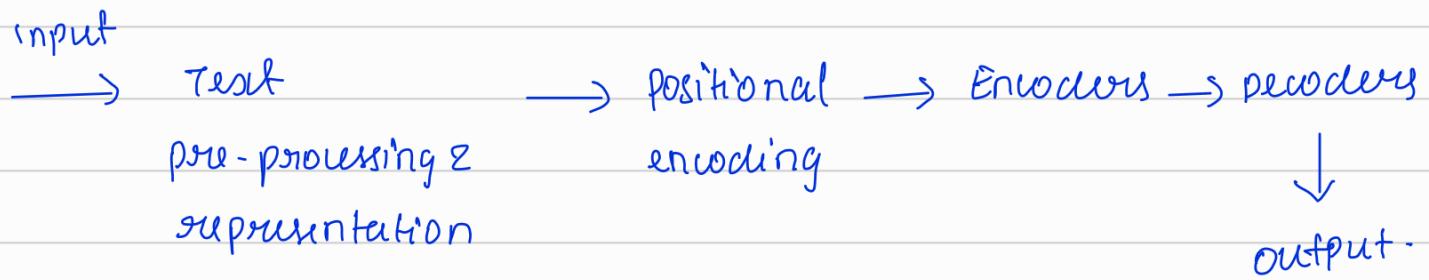
* Attention is all you need.

(Revolutionized language modelling)

* Transformer architecture.

→ Relationship between words

→ Components, pre-processing, positional encoding,
Encoders & decoders.



Positional encoding

- information on the position of each word.
- understand distant words.

Encoder

- Attention Mechanism, directs attention to specific words and relationship
- Neural network, process specific features.

Decoder

- includes attention and neural network.
- generates the output

Transformers & long-range dependencies

- initial challenges, long range dependency.
- Attention, focus on different parts of the input.

processes multiple parts simultaneously.

- limitation of traditional language model.
(sequential - one word at a time)
- Transformer,
 - * processes multiple parts simultaneously.
 - * faster processing.

Attention Mechanism

- understand complex structures
- focus on important words

Self attention

- weights the importance of each word.
- captures long range dependencies

Multi-head attention

- Next level of self attention
- splits inputs into multiple heads with each head focusing on different aspects.

Attention in party

Self-attention

- Focus on each person's words
- Evaluate and compare their relevance
- Weigh each speaker's input
- Combines for a comprehensive understanding

Multi-head attention

- Split attention into "multiple" channels
- Focus on different aspects of conversation
- Speaker's emotions, primary topic, and related side-topics
- Process each aspect and merge

Multi-head attention advantages

example

- "The boy went to the store to buy some groceries, and he found a discount on his favorite cereal."
- Attention: "boy," "store," "groceries," and "discount"
- Self-attention: "boy" and "he" → same person
- Multi-head attention: multiple channels
 - Character ("boy")
 - Action ("went to the store," "found a discount")
 - Things involved ("groceries," "cereal")

multi-head attention is like having multiple self attention mechanisms working simultaneously

Advanced fine-tuning

Reinforcement Learning through Human feedback

- * pre-training
- * Fine-tuning
- * Reinforcement Learning through Human feedback (RLHF)

pre-training.

- large amounts of text data
- next word prediction
- Masked Language Modelling.

Fine-tuning

- N-shot training
- small labeled dataset for related task.

Why RLHF?

→ General purpose training data lacks quality
(Noise, Errors, inconsistencies & reduced accuracy)

Pretraining



learns underlying language patterns

Doesn't capture context complexity

Fine tuning



Quality labelled data improves performance.

RLHF !



Human Feedback.

Simplifying RLHF

- Model output reviewed by human
- updates model based on the feedback.

Step 1:

- Receives a prompt
- Generates multiple responses.

Step 2:

- Human expert checks their responses
- Ranks the responses based on quality.
 - * Accuracy
 - * Relevance
 - * Coherence.

↑ Enter
human Experts.

Step 3:

- Learn from expert's ranking
- To align its responses in future with their preferences.

And it goes on!

- * Continues to generate responses
- * Receives expert's ranking
- * Adjusts the learning.

[RLHF techniques to enhance fine-tuning through human feedback]



Data concerns and consideration

Data considerations

- * Data volume & compute power.
- * Data quality
- * Labelling
- * Bias
- * Privacy.

Accurate data =

better learning = improved

response quality =

increased trust.

→ LLM need a lot of data, similar to a child learning to talk (570GB, ~1.3M books)

→ Extensive computing power, think of the energy consumption

→ can cost millions of dollars?

Labelled data

- correct data label. (accurate learning, generalize patterns, accurate response)
- Labor-intensive (assigning correct label to each article)
- Incorrect labels impact model performances.
- Addresses errors: Identify → analyse → iterate.

Data bias

- influenced by social stereotypes
- lack of diversity in training data
- discrimination & unfair outcomes.

→ spot and deal with the biased data

* Evaluate data imbalance

* promote diversity

* Bias mitigation techniques: more diverse examples.

Data privacy

→ compliance with data protection & privacy regulation

→ sensitive / personally identifiable information (PII)

→ Get permission

→ Privacy is a concern

* Training on data without permission can lead to a breach

* Legal, financial and reputation harm

Ethical concern

* Transparency risk. → challenging to understand the output

* Accountability risk

* Information hazards.



Disseminating

harmful information

→ Harmful content generation,

→ Misinformation spread,

→ Malicious use,

→ Toxicity

Responsibility of LLM's actions.

Environmental concerns

- Ecological footprints of LLMs
- substantial energy resources to train
- Impact through carbon emissions.

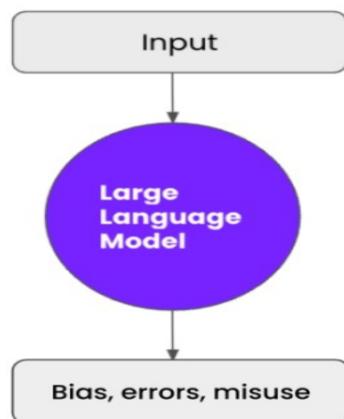
Cooling requires electricity too!

- Produce considerable heat that needs cooling
- Imagine thousands of laptops overheating
 - Require complex cooling systems
 - Adds to environmental impact
- Balance the cost and benefits
 - Use renewable energy
 - Energy-efficient tech



Transparency risk

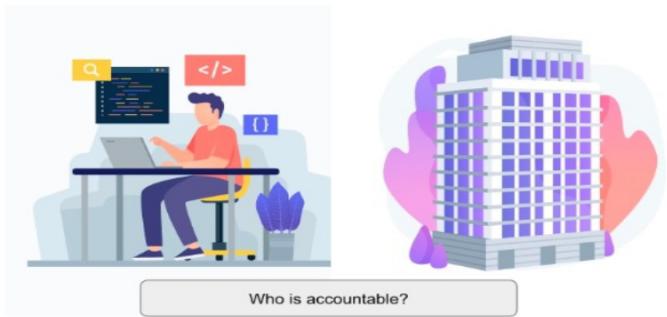
- Challenging to understand the output
- Difficult to identify issues
 - Bias
 - Errors
 - Misuse
- Black box
- Example: reasoning behind predicting disease outcomes



LLM must explain for
making informed treatment
decisions

Accountability risk

- Responsibility of LLMs' actions
- Who is responsible?
 - Incorrect and harmful advice
 - Model developer or the company?
- Game without rules
 - No transparency
 - No accountability



where are LLMs heading?

① Model explainability.

- How do they arrive at their outputs.
- Builds trust and transparency
- Identify and correct the biases or errors.

② Efficiency.

- computational efficiency.
 - (High quality output with less compute)
- Faster & efficient
 - (Model compression & optimization)
- Benefits.
 - (Better storage, lower energy use)
- Accessibility & sustainability
 - (promotes green AI,
Reduces operating costs)

③ unsupervised bias handling.

- Biased data (discrimination)

→ unsupervised bias handling.

* Bias detection and mitigation techniques automatically

* No need of explicit human labelled data.

* identifies & reduces by analysing patterns

→ challenges

* subtler, difficult to detect

* Might introduce new biases.

4.

Enhanced creativity

- Creativity in text-based and visual art forms
- Artistic content: learned patterns, not emotional understanding
- Lack human-like comprehension of art or emotions
- Demonstrate human-like emotional behavior
- Future: emotion inference

