

Decoding Language: A Deep Dive into Natural Language Processing

Data Science Workshop – Breakout Session

Indian Institute of Technology, Bhilai



Our Expedition



Part 1: The Foundation

What, Why, and How to Clean Text



Part 2: Giving Words Meaning

Embeddings, Classic ML, and Intro to Neural Nets



Part 3: The Modern Era

Transformers, LLMs, and the Future



Why Should a Student Care About NLP?

It's the engine behind the modern Indian digital economy, powering everyday interactions and complex analytical tasks.



E-commerce Insights

Analyzing Flipkart/Amazon reviews in English and Hinglish to gauge customer sentiment and product reception.



FinTech Automation

Your bank app automatically categorizing UPI spends (Food, Travel, Bills) for better financial management.



Social Media Analytics

Gauging sentiment for CSK vs MI during an IPL match from thousands of tweets in real-time.



Customer Support Chatbots

Chatbots on sites like IRCTC or Swiggy efficiently handling your queries and providing instant solutions.



The Core Challenge: Ambiguity

Human language is inherently complex and full of double meanings, making it a significant hurdle for machines to understand.

Example 1: "I saw a man on a hill with a telescope." Who has the telescope? Is it you, or the man on the hill?

Example 2 (Indian Context): "The bank of the river Ganga is sacred." Here, 'bank' refers to a river bank, not a financial institution like SBI.

Example 3 (Hinglish): "Achha, you are coming?" 'Achha' can convey various meanings like 'okay', 'really?', 'I see', depending on the intonation and context.



The NLP Pipeline: A "Safai Abhiyan" for Text

A step-by-step process to convert messy, raw text into clean, structured data that a machine can understand, much like a meticulous cleaning drive.

1. Normalization

Making text uniform and removing "noise" like case differences and punctuation.

2. Tokenization

Breaking down a continuous string of text into individual components or "tokens" (words, sentences).

3. Stopword Removal

Eliminating high-frequency words that carry little semantic weight (e.g., "the", "is", "a").

4. Stemming/Lemmatization

Reducing words to their root form for consistent analysis (e.g., "running" to "run", "better" to "good").

5. POS Tagging

Assigning a grammatical category (e.g., noun, verb, adjective) to each token to understand sentence structure.



Step 1: Normalization

The crucial first step in NLP is to make text uniform and remove "noise" that can confuse machine learning models.

Goal: Make Text Uniform

- Convert to lowercase: THE and the become the same.
- Remove punctuation: amazing!!! becomes amazing.
- Remove URLs, hashtags, and user mentions.

Before:

WOW! The finale of MERAZ'24 at
@iit_bhilai was AMAZING!!
#TechFest

After:

wow the finale of meraz at iit bhilai
was amazing techfest



Step 2: Tokenization

Tokenization is the process of breaking down a string of text into individual components or "tokens," which are the basic units for further processing.

Goal: Break Down Text

- **Sentence Tokenization:** Breaking a paragraph into individual sentences.
- **Word Tokenization:** Breaking a sentence into individual words.

Example:

- Sentence: MS Dhoni plays cricket.
- Tokens: ['MS', 'Dhoni', 'plays', 'cricket', '.']



Step 3: Stopword Removal

Stopword removal focuses on eliminating high-frequency words that carry little semantic weight, allowing the model to concentrate on the more meaningful terms.

Goal: Remove Semantic Noise

- These words are like the "filler" in a sentence.
- **English Stopwords:** a, an, the, is, are, was, in, on, at
- **Hindi/Hinglish Stopwords:** hai, hai, ka, ki, ke, ne, par, mein, ek

Example:

the match is in the stadium

Becomes:

match stadium



THE MATCH IS IN STADIUM



Step 4: Stemming vs. Lemmatization

Goal: To reduce a word to its root form.

Stemming

A crude, rule-based process of chopping off suffixes. It's fast but can be inaccurate, sometimes creating non-existent words.

Lemmatization

A smarter, dictionary-based process to get the actual root word (lemma). It's slower but more accurate, using linguistic knowledge.



Stemming Examples

studies → studi

running → run



Lemmatization Examples

studies → study

better → good

चले (chale) → चलना (chalna)



Step 5: Part-of-Speech (POS) Tagging

POS Tagging assigns a grammatical category to each token, providing crucial structural understanding to the machine, similar to how we analyze grammar.

Goal: Assign Grammatical Category

It helps the machine understand the structure of a sentence and the role each word plays.

Example:

Virat (NNP) calmly (RB) hit (VBD) the (DT) ball (NN)

Common Tags:

- **NNP**: Proper Noun, singular (e.g., Virat)
- **RB**: Adverb (e.g., calmly)
- **VBD**: Verb, past tense (e.g., hit)
- **DT**: Determiner (e.g., the)
- **NN**: Noun, singular (e.g., ball)



The Problem: How Does a Computer Understand "Biryani"? "Biryani"?

Computers don't understand human language; they only process numbers. The challenge in NLP is converting words like "Biryani," "Samosa," or "Dhokla" into numerical representations. Crucially, these numbers must reflect semantic meaning, ensuring "Biryani" is mathematically closer to "Pulao" than to "Laptop" in vector space.



The Old Way: One-Hot Encoding

One-Hot Encoding assigns a unique numerical vector to each word in the vocabulary. For a dictionary of 10,000 words, each word becomes a vector of 10,000 zeros with a single '1' at its specific index.

Example:

Vocabulary: [a, ..., biryani (idx 500), ..., laptop (idx 3000), ..., the] Biryani: [0, 0, ..., 1, ..., 0, ..., 0] (at index 500) Laptop: [0, 0, ..., 0, ..., 1, ..., 0] (at index 3000)

Huge & Sparse Vectors

Inefficiently uses memory, especially with large vocabularies, as vectors are mostly filled with zeros.

No Semantic Similarity

Words are treated as completely unrelated (orthogonal), failing to capture vital linguistic relationships. For example, "biryani" is as different from "pulao" as it is from "laptop" in this representation.



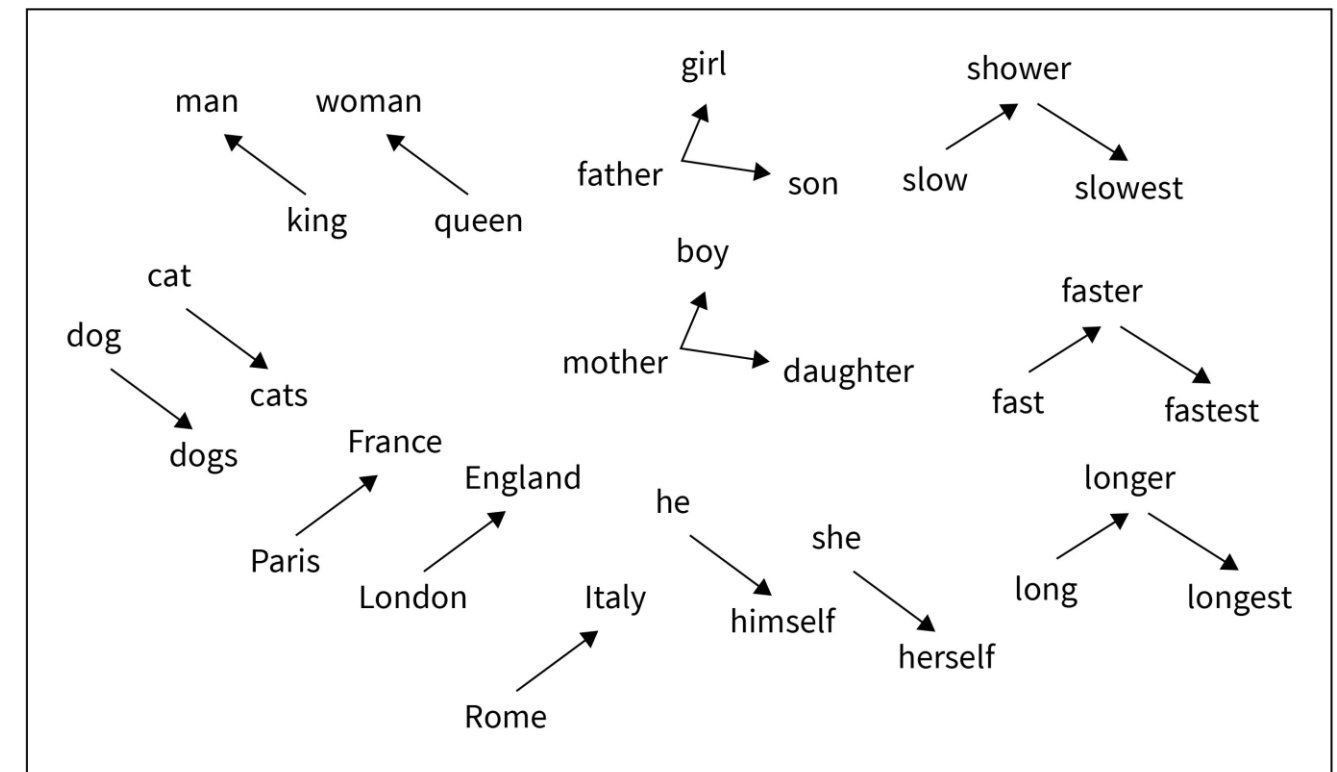
The Breakthrough: Word Embeddings (e.g., Word2Vec)

Compact Semantic Representation

Words are transformed into short, dense numerical vectors (e.g., 300 numbers), drastically reducing sparsity from one-hot encoding.

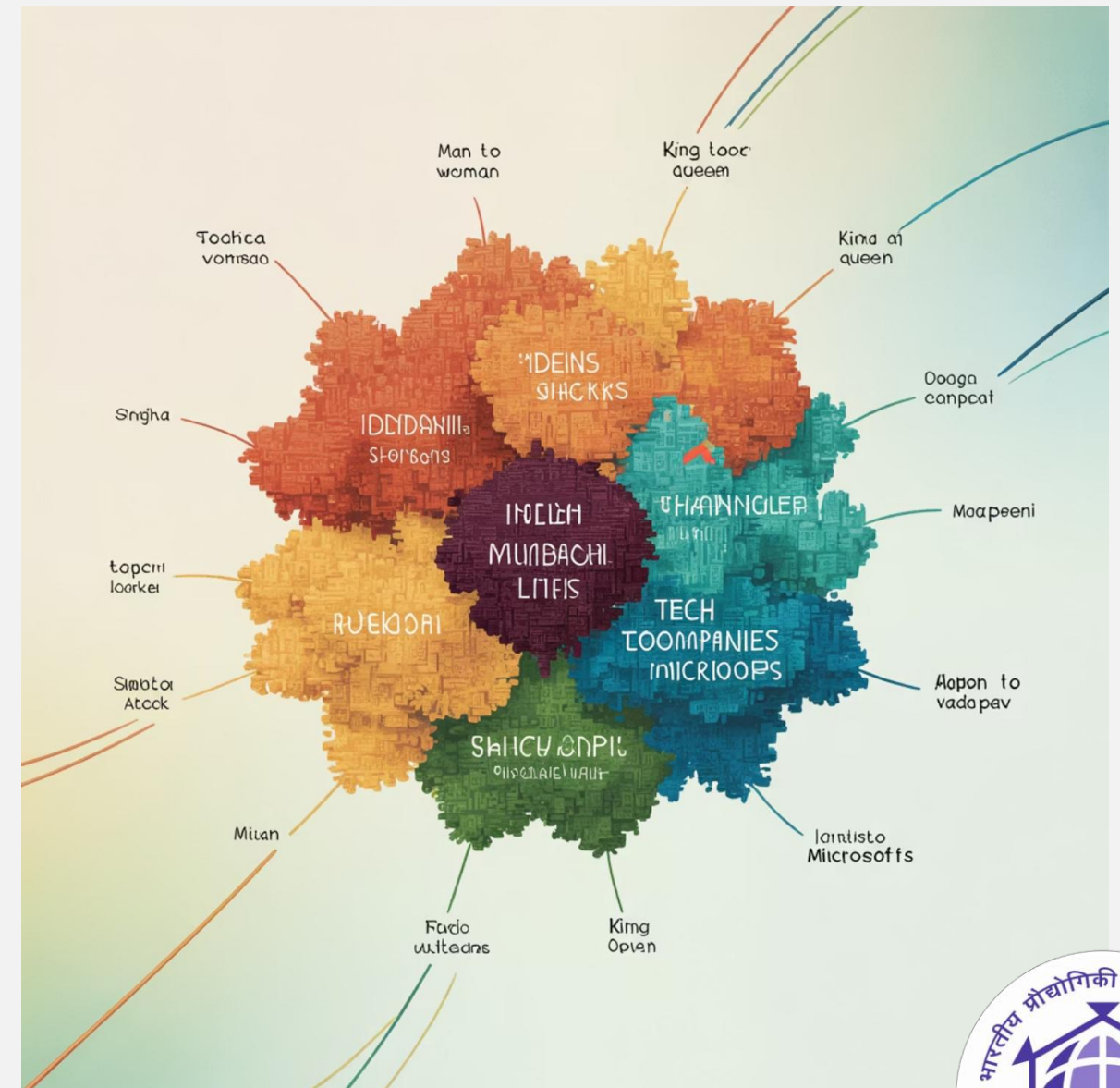
Vectors are learned from huge text datasets. Critically, words appearing in similar contexts will have similar vector representations.

"A word is characterized by the company it keeps." This principle ensures semantic closeness between related words like "Biryani" and "Pulao."



Exploring the "Meaning Space"

Word embeddings map words into a high-dimensional "meaning space." In this space, vector proximity reflects semantic similarity. Words like **Delhi**, **Mumbai**, and **Kolkata** form distinct clusters. Furthermore, vector directions capture relationships: the transformation from **King** to **Queen** parallels **Man** to **Woman**, revealing profound linguistic analogies.



The Famous Analogy: Arithmetic on Words

The structure of the meaning space in word embeddings allows us to perform "conceptual arithmetic," revealing nuanced relationships between words.

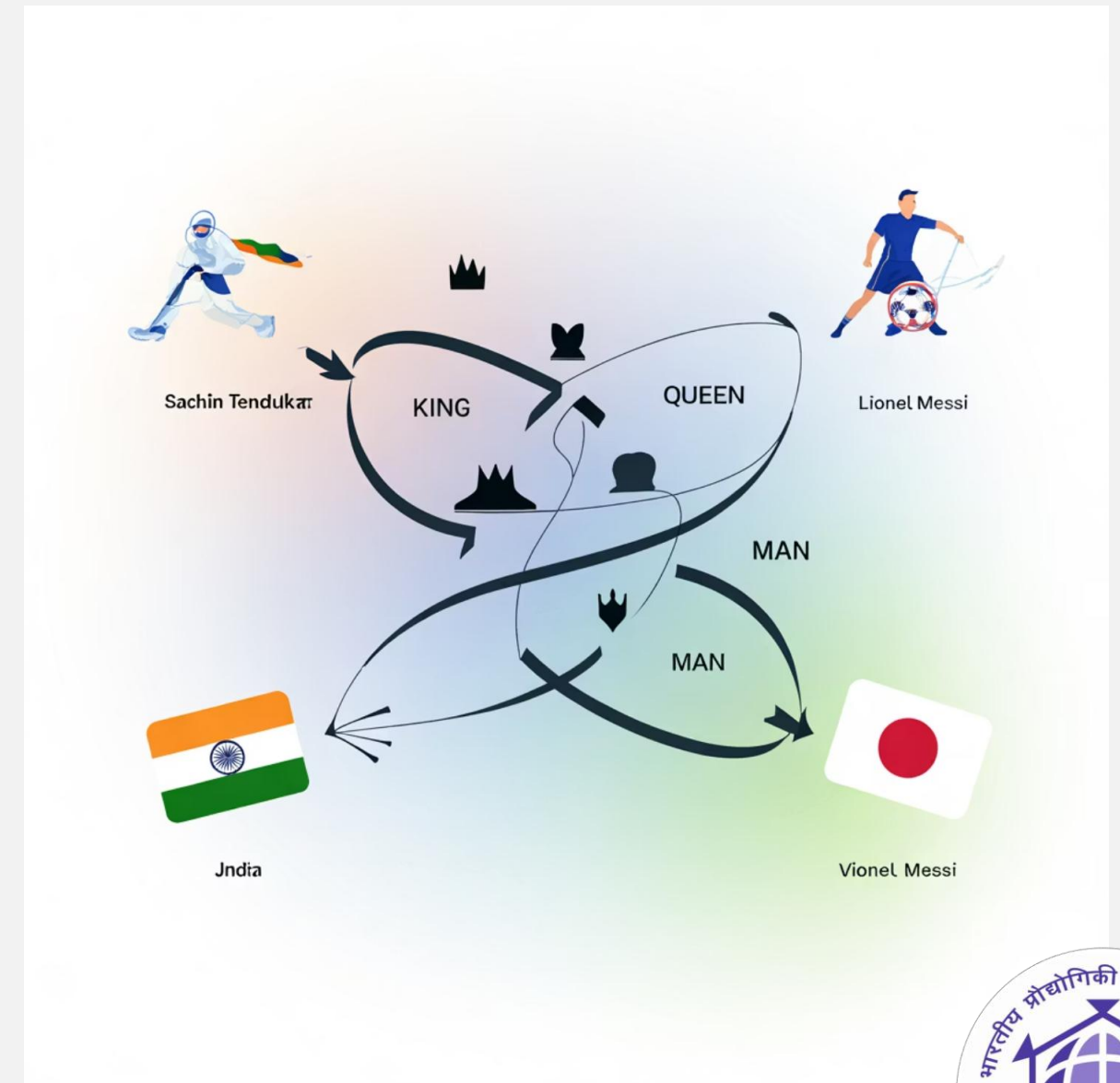
Classic Example

$\text{vector('King')} - \text{vector('Man')} + \text{vector('Woman')} \approx \text{vector('Queen')}$

Indian Context Examples

$\text{vector('Delhi')} - \text{vector('India')} + \text{vector('Japan')} \approx \text{vector('Tokyo')}$

$\text{vector('Sachin Tendulkar')} - \text{vector('Cricket')} + \text{vector('Football')} \approx \text{vector('Lionel Messi')}$



Classic NLP: Bag-of-Words (BoW)

The Bag-of-Words (BoW) model is a fundamental approach for document classification, like categorizing movie reviews.

It works by:

- Building a vocabulary of all unique words.
- Representing each document as a vector of word frequencies.
- Ignoring grammar and word order, treating text as a "bag" of words.

For instance, a review with "paisa", "vasool", "superhit" strongly suggests a positive sentiment.



A Smarter Way to Count: TF-IDF

The Bag-of-Words model struggles with common words, which often dominate counts but lack semantic value.

TF-IDF (Term Frequency-Inverse Document Frequency) addresses this by:

- **Term Frequency (TF):** Measures how often a word appears in a specific document.
- **Inverse Document Frequency (IDF):** Down-weights words that are common across many documents.

This method highlights unique, informative terms like "paisa vasool" while reducing the impact of widespread words like "movie."



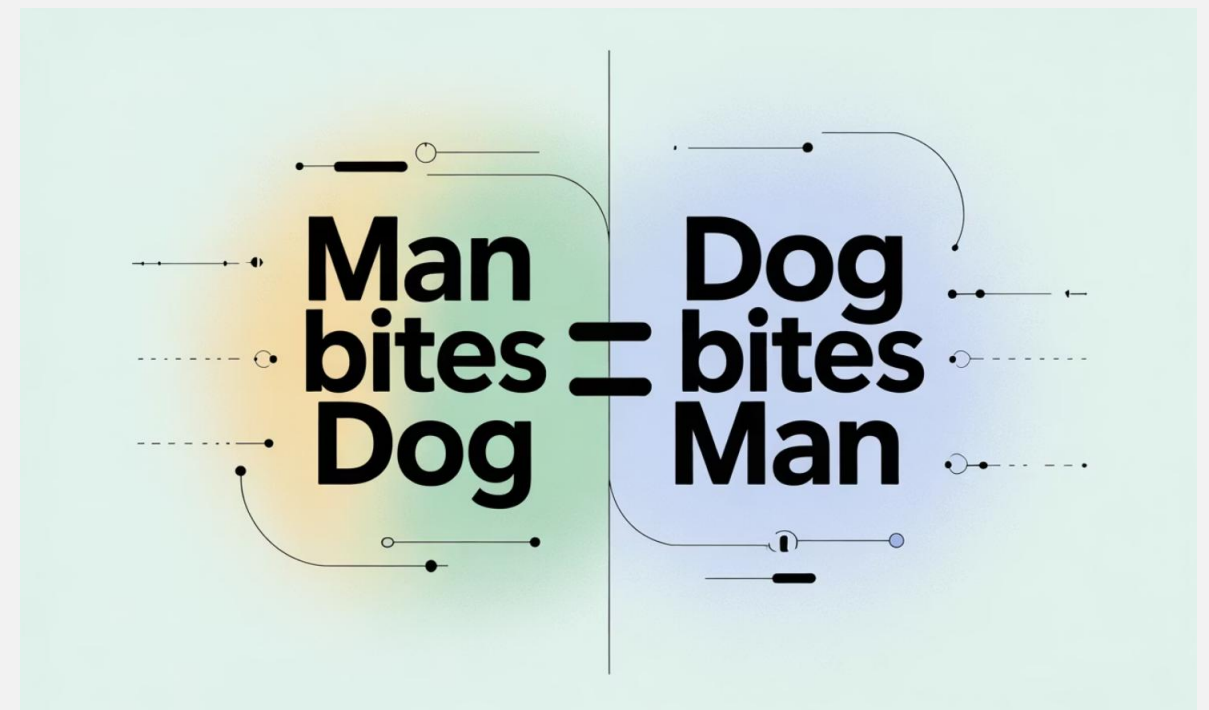
TF-IDF: Balancing common vs. impactful words.

The Big Limitation: Word Order Matters!

Bag-of-Words (BoW) and TF-IDF completely ignore the sequence of words. To these models, "Man bites dog" and "Dog bites man" are treated as identical word collections, losing crucial meaning.

This lack of sequential understanding also makes detecting nuances like sarcasm impossible. For example, "The movie was not good, it was brilliant" is misinterpreted.

We need advanced models that can process and understand sentences sequentially, just like human cognition.



BoW and TF-IDF struggle with context and sequence!



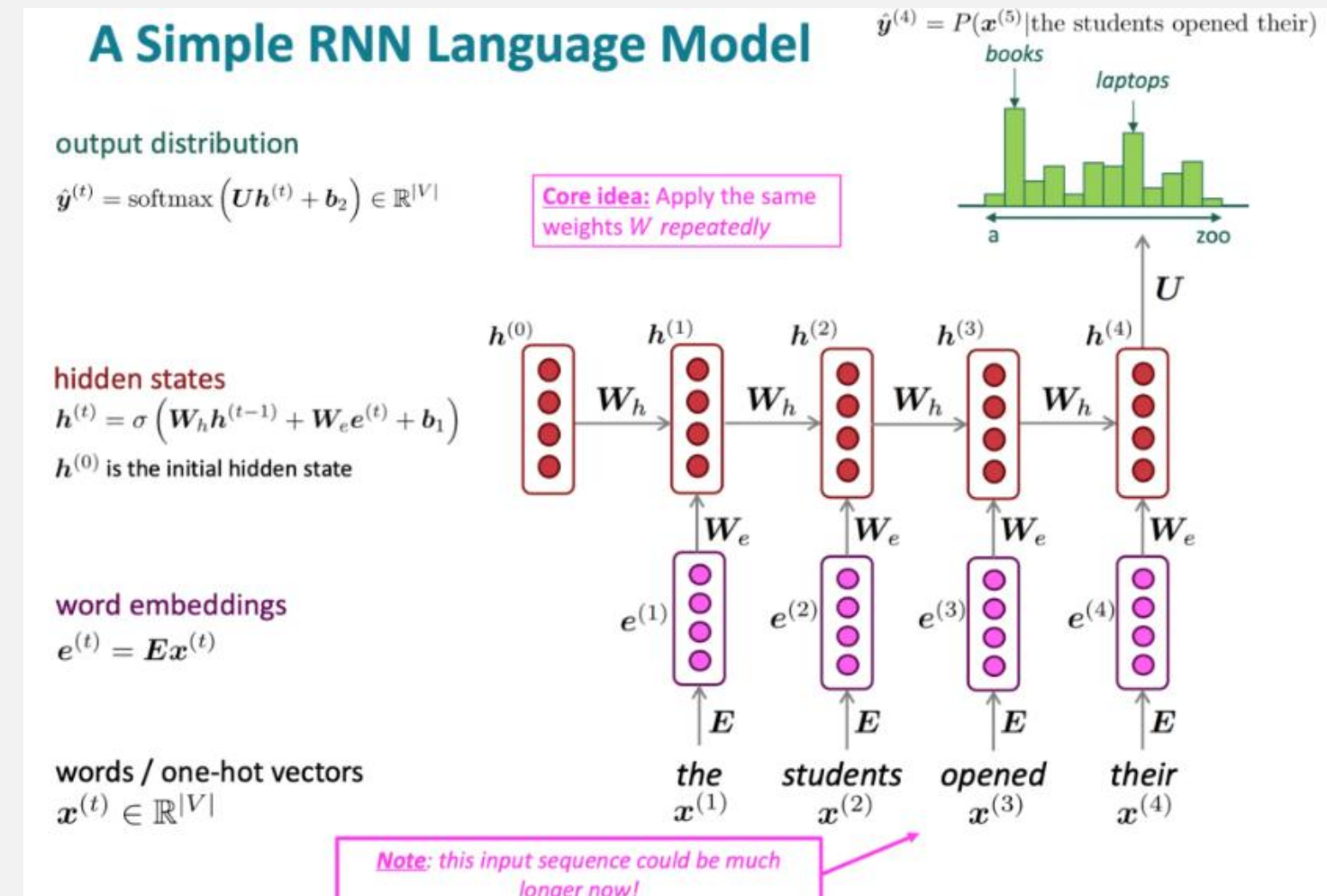
Introducing Sequential Models: Recurrent Neural Networks

Networks (RNNs)

Addressing the **word order limitation**, Recurrent Neural Networks excel at processing sequential data.

Unlike previous models, RNNs maintain a "memory" or **hidden state** that captures context as they read text word-by-word.

- Each step takes the current word's embedding.
- It combines this with the previous step's hidden state.
- Resulting in an updated memory for the next word, enabling deep contextual understanding.



How an RNN "Remembers"

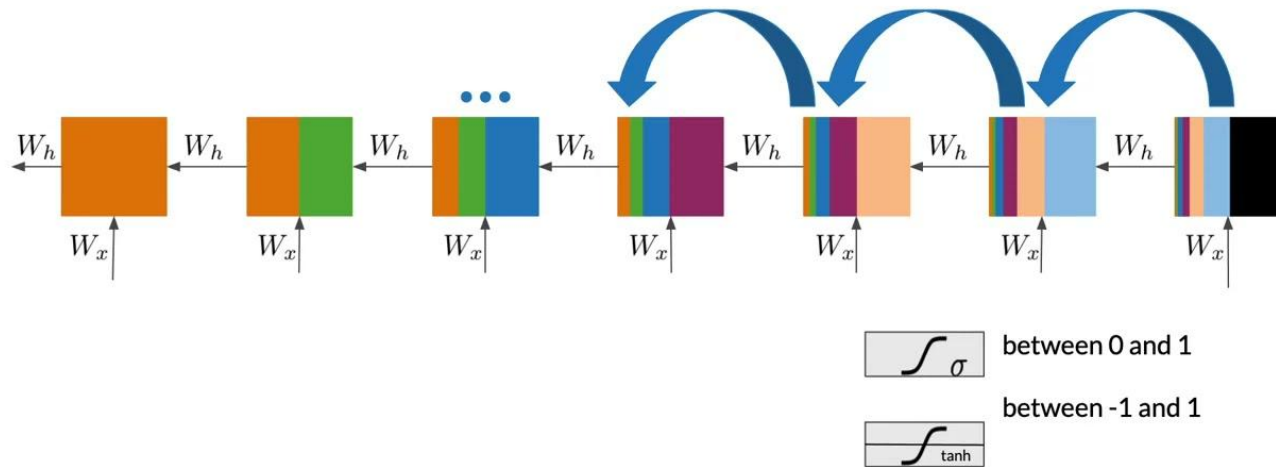
Imagine an RNN reading a sentence, building context word by word:

- **The:** RNN starts with an empty, clean memory slate.
- **capital:** Memory updates, now it registers we're talking about a capital city.
- **of:** Memory continues to refine the context.
- **India:** Memory strongly focuses: "capital of India."
- **is:** Memory continues to integrate the sentence structure.
- **New:** Based on its rich, contextual memory state, the RNN can confidently predict the next word: **Delhi**.

This dynamic memory allows RNNs to understand sequence and context, making predictions more accurate.



Backpropagation through time



The RNN Memory Problem: Vanishing Gradients

Standard Recurrent Neural Networks (RNNs) face a significant challenge: a short-term memory problem. In processing long sentences, the influence of earlier words can diminish, or "vanish," by the time the model reaches the end.

Consider the sentence: "I grew up in a small village in Rajasthan... many years later, when I moved to the US, I really missed the daal baati churma my grandmother used to make." By the time the RNN sees "daal baati churma", it may have forgotten the crucial context of "Rajasthan". This phenomenon is known as the vanishing gradient problem, where the signal from early words gets weaker as it travels through the sequence.

The Solution: LSTMs (Long Short-Term Memory)

Long Short-Term Memory (LSTMs) networks are an advanced type of RNN specifically designed to combat the vanishing gradient problem. Their key innovation lies in a "cell state"—acting like a conveyor belt for information—and a sophisticated system of "gates" that control the flow of information.

Forget Gate

Decides what old information to discard from memory, preventing irrelevant details from cluttering the cell state.

Input Gate

Determines what new information from the current input should be stored in the memory cell, filtering for relevance.

Output Gate

Regulates what part of the memory cell's content will be used for the current prediction and passed to the next hidden state.

This intricate gating mechanism allows LSTMs to remember relevant information for extended periods, enabling them to process longer sequences effectively.

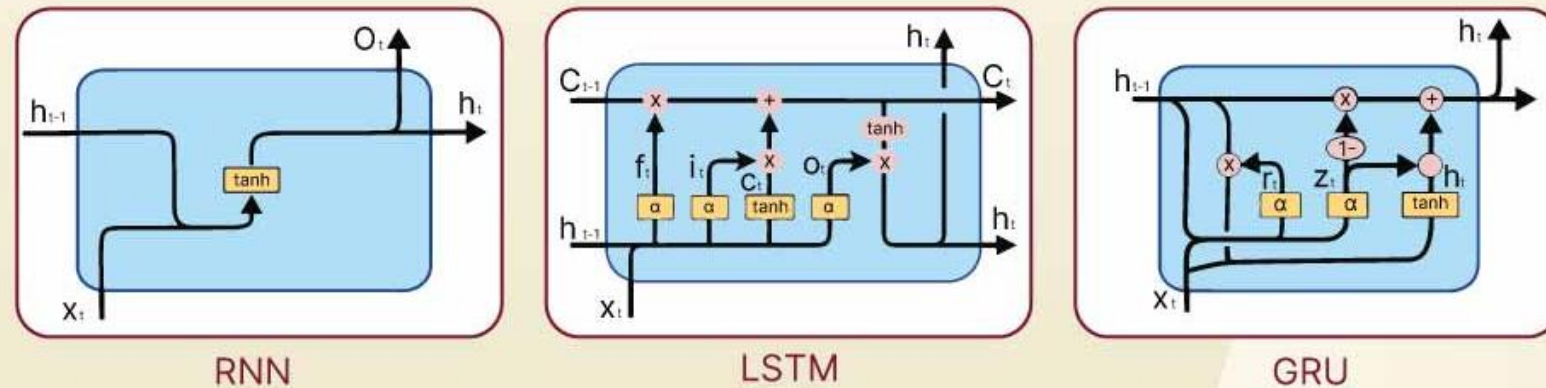


The Next Leap: Introducing Transformers

In 2017, the groundbreaking paper "Attention Is All You Need" introduced the Transformer architecture, fundamentally changing the landscape of NLP. The core issue with traditional RNNs and LSTMs was their sequential nature; they had to process text word-by-word, which was inherently slow and prevented parallel computation.

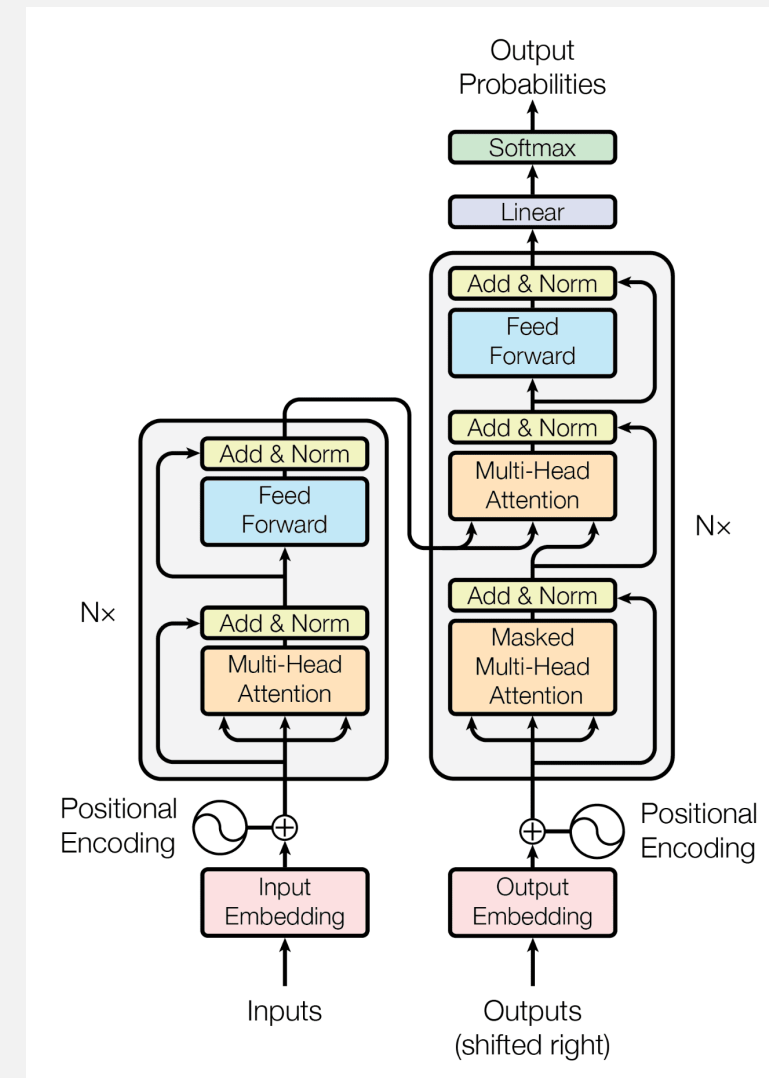
RNN/LSTM Processing

RNNs and LSTMs process information sequentially, one word at a time, leading to slow training and inference, especially for long sequences.



Transformer Processing

Transformers abandon recurrence, processing all words in a sentence simultaneously using the self-attention mechanism, enabling fast parallel computation.



The Secret Sauce: Self-Attention

Attention Mechanism

The fundamental goal of the self-attention mechanism is to determine, for each word in a sentence, which other words are most important for understanding its context. The model learns to calculate "attention scores" between every pair of words.

A higher attention score signifies that the model should "pay more attention" to that particular word when processing the current word. This enables the model to grasp complex dependencies and relationships within a sentence, regardless of how far apart the words are from each other.

This mechanism is crucial for understanding nuances and ambiguities in language, allowing Transformers to build highly contextualized representations of words.

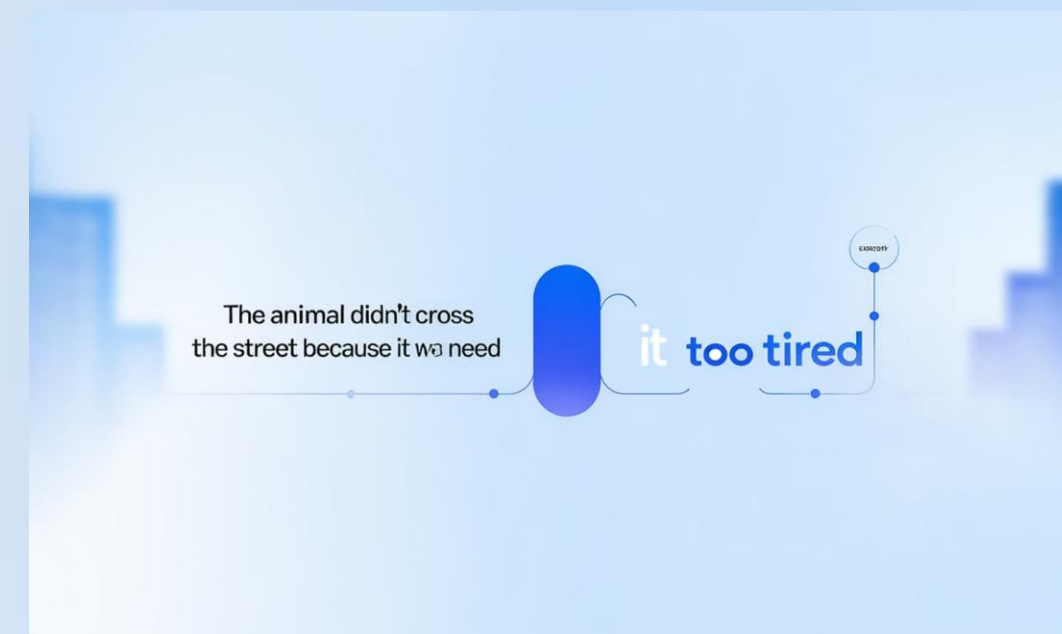


Visualizing Attention in Action

To illustrate the power of self-attention, consider the classic example sentence: "The animal didn't cross the street because it was too tired." A common ambiguity arises with the pronoun "it." Does "it" refer to the "animal" or the "street"?

Through its training, the attention mechanism learns to assign a very high attention score between the word "it" and the word "animal." This strong connection allows the model to correctly infer that "it" refers to the animal, not the street, based on the context of being "too tired."

This dynamic capability allows Transformer models to build complex and accurate relationships between words for every sentence they process, effectively resolving linguistic ambiguities that challenge traditional models.



The Transformer Family: BERT & GPT

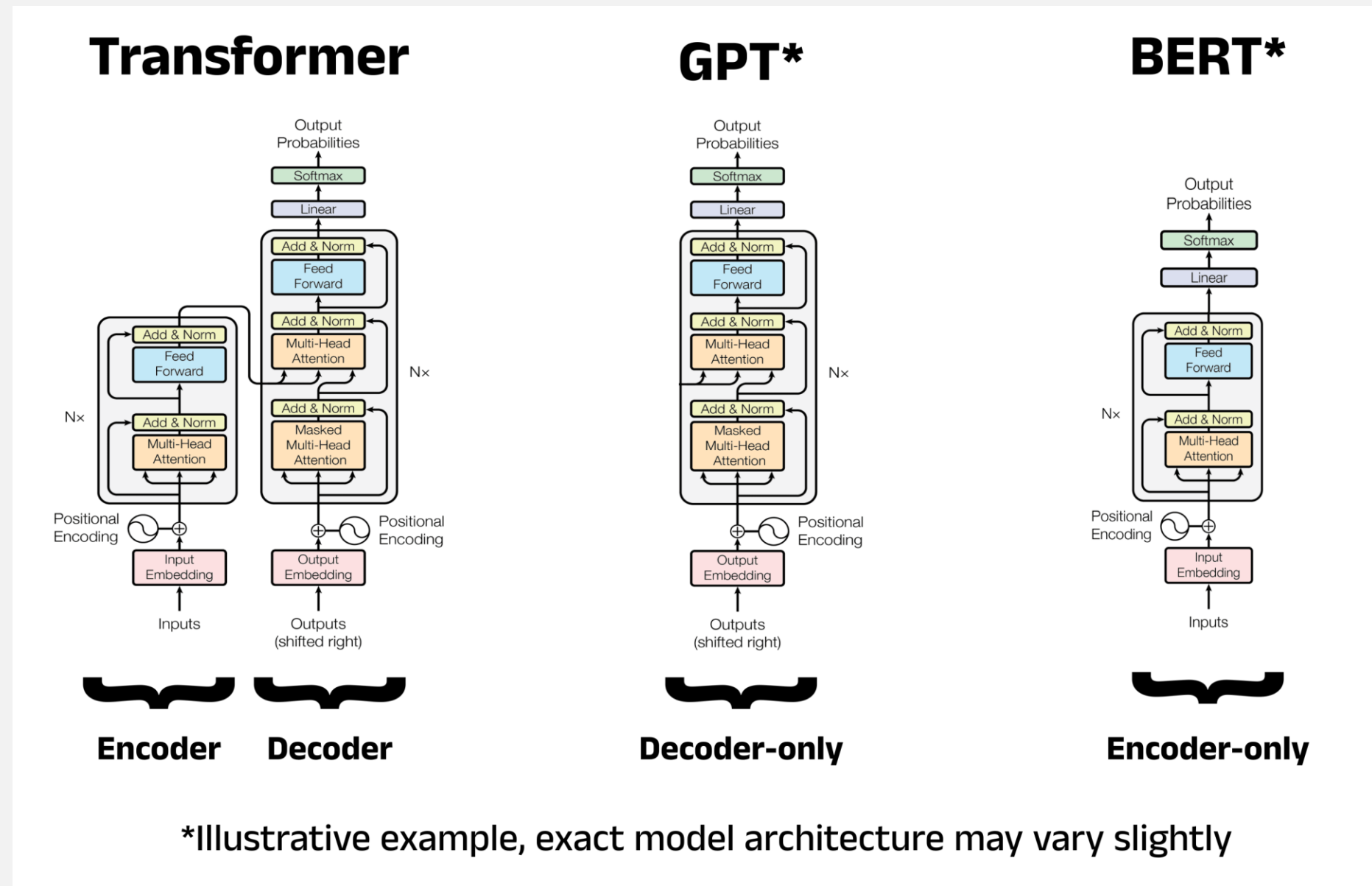
The groundbreaking Transformer architecture has given rise to two primary types of powerful language models, each optimized for different tasks based on their processing direction and design.

BERT (Bidirectional Encoder Representations from Transformers)

- Reads the entire sentence at once, allowing for a bidirectional understanding of context.
- **Job:** To deeply understand text from both left and right contexts.
- **Use Cases:** Google Search (improving relevance), sentiment analysis, named entity recognition, and question answering.

GPT (Generative Pre-trained Transformer)

- Processes text from left-to-right, predicting the next word in a sequence (autoregressive).
- **Job:** To generate human-like text by predicting the most probable next word.
- **Use Cases:** ChatGPT, content creation, summarization, creative writing, and code generation.



The Age of Giants: Large Language Models (LLMs)

Large Language Models (LLMs) are essentially Transformer models scaled to an unprecedented size. They are trained on a vast portion of the internet's text data, including encyclopedias, books, articles, and websites.

These models possess hundreds of billions, and sometimes trillions, of "parameters" – the internal configurations that the model learns and adjusts during its extensive training process. This sheer scale leads to what are known as "emergent abilities."

At this massive scale, LLMs begin to exhibit surprising capabilities they were not explicitly trained for, such as complex reasoning, accurate language translation across many languages, and even writing functional code. These abilities are a direct result of the immense data and computational power invested in their creation.



The Art of Asking: Prompt Engineering

To unlock the full potential of Large Language Models (LLMs) and obtain the most accurate and useful results, it is crucial to formulate your queries effectively. This skill is known as **Prompt Engineering**.

Bad Prompt: Generic and Vague

Write about IIT Bhilai.



IIT Bhilai, established in 2016, is a premier institute of technology in India. Located in Chhattisgarh, it emphasizes innovation, interdisciplinary research, and holistic education. With modern infrastructure and a vibrant academic environment, it nurtures future leaders in science, engineering, and technology through cutting-edge programs and industry collaborations.

Good Prompt: Specific and Contextualized

Role: Act as a prospective student's guide.

Task: Write a 4-paragraph overview of IIT Bhilai.

Context: Focus on its unique academic programs, vibrant campus life, and the vision for it to become a leading tech institute. Use an encouraging and informative tone.



Indian Institute of Technology Bhilai, established in 2016, is one of India's newest IITs, rapidly emerging as a hub of innovation and academic excellence. Located in the serene city of Bhilai, Chhattisgarh, the institute offers cutting-edge programs in engineering, sciences, and liberal arts, with a strong emphasis on interdisciplinary learning. Its modern campus fosters creativity, collaboration, and research-driven education. Students enjoy a vibrant campus life enriched by cultural fests, tech events, and student-led initiatives. With a vision to become a global leader in technology and innovation, IIT Bhilai is shaping the next generation of thinkers and problem-solvers.



Summary of Our Entire Journey

Raw, Ambiguous Text

We began with unstructured and often unclear textual data.

Cleaned and Structured

Processed text using the fundamental NLP pipeline for clarity.

Word Embeddings

Gave words semantic meaning through dense vector representations.

RNNs/LSTMs

Explored sequential processing with memory capabilities for context.

Transformers & Attention

Discovered parallel processing and contextual understanding with attention mechanisms.

LLMs & Prompt Engineering

Concluded with massive Transformers controlled by precise prompting.

From the initial ambiguity of raw text to the sophisticated generation capabilities of LLMs, we've covered the evolutionary milestones in neural networks for language processing.





Thank You & What's Next?

Let's continue to build amazing things with the power of language.

Workshop Materials

Access all code notebooks and resources on our [GitHub Repository](#).

Recommended Things to Explore

- **Hugging Face Transformers:** For state-of-the-art Large Language Models.
- **Course :** Stanford: CS224N: Natural Language Processing with Deep Learning