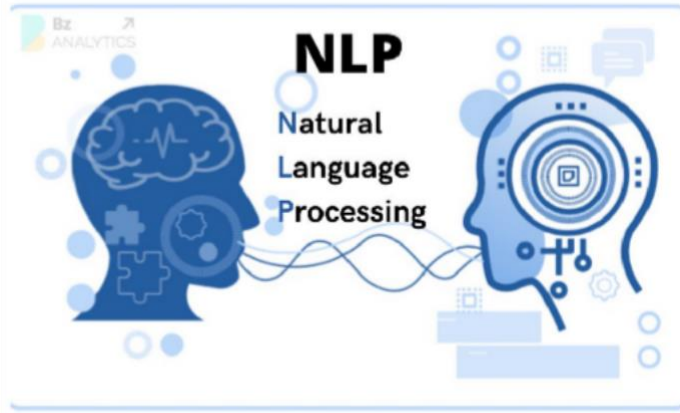
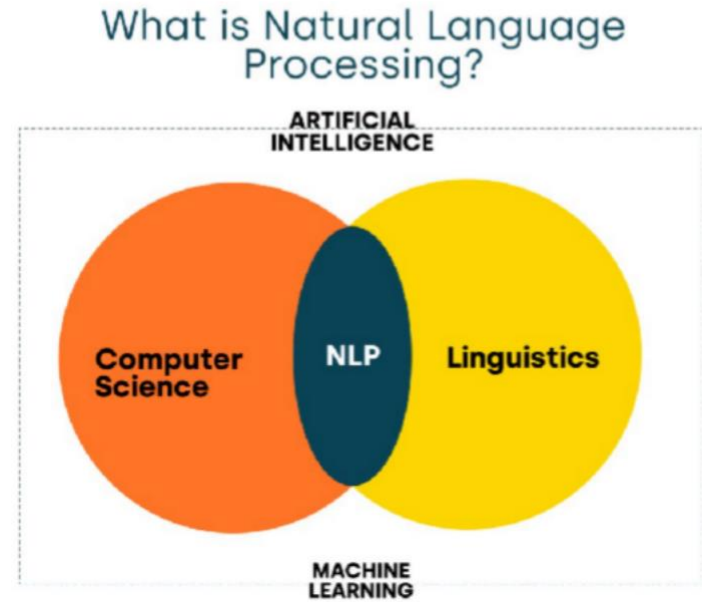

Introduction to NLP

Natural Language Processing (NLP) is a field of artificial intelligence (AI) that focuses on enabling computers to understand, interpret, and generate human language.



NLP bridges the gap between human language (which is complex and ambiguous) and computer language (which is highly structured and precise).



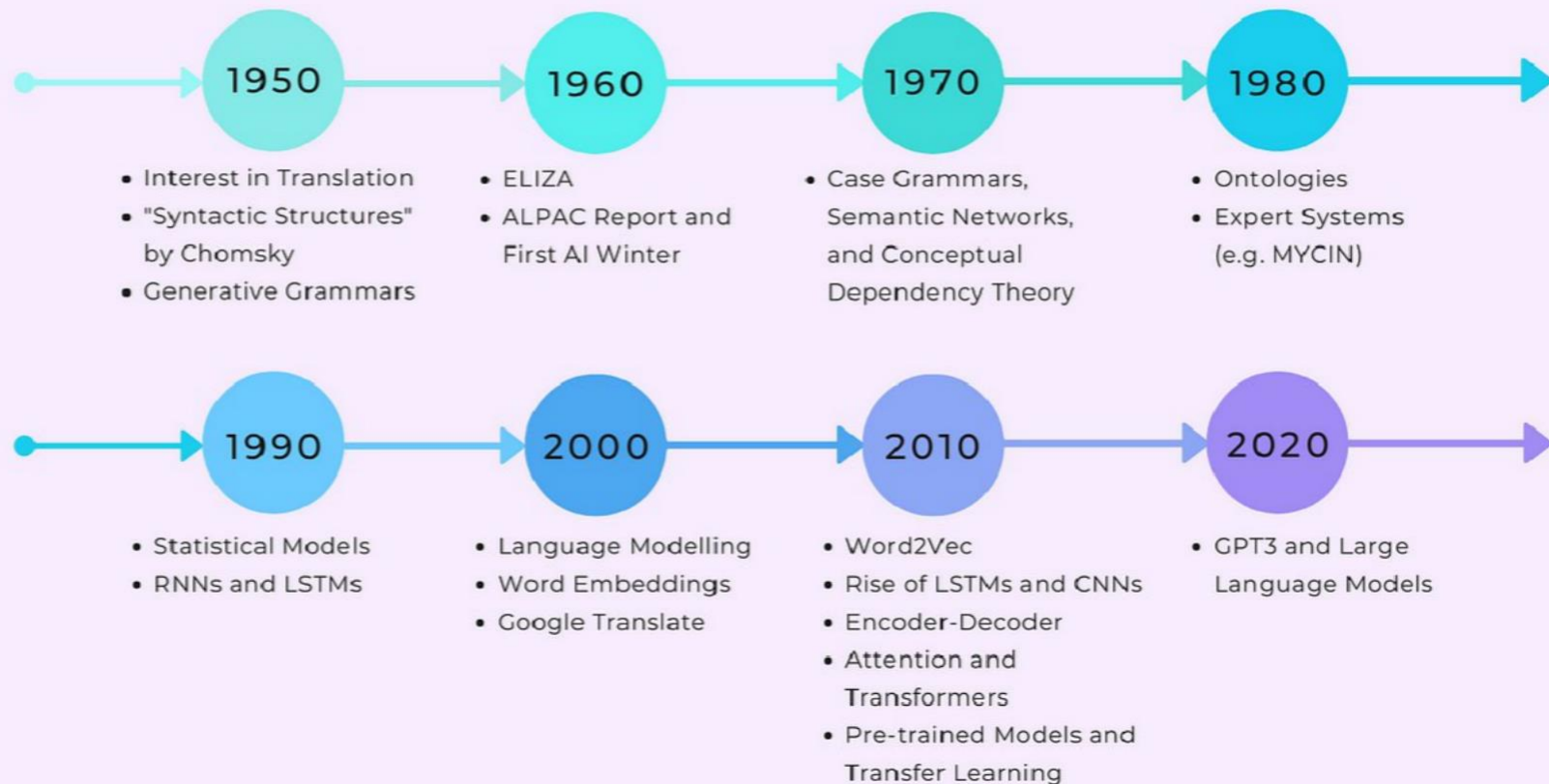
The ultimate goal of NLP is to enable computers to perform a variety of language-related tasks, such as

- ❖ understanding spoken commands,
- ❖ translating languages,
- ❖ summarizing documents, and
- ❖ generating conversational responses.
- ❖ Many more tasks related to language

Applications of Natural Language Processing



A Brief Timeline of NLP



Key Challenges in Natural Language Processing (NLP)

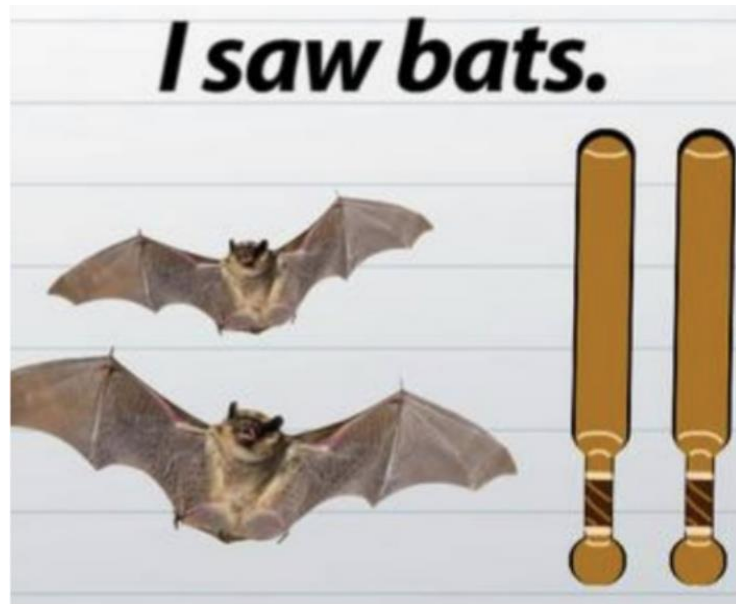


1. Ambiguity in Language

Ambiguity is one of the most fundamental challenges in NLP. Human language is filled with words and structures that can have multiple meanings, depending on the context. There are three main types of ambiguity in NLP:

1.1 Lexical Ambiguity

- **Definition:** Lexical ambiguity occurs when a word has multiple possible meanings.
- **Example:** The word "bank" can refer to a financial institution or the side of a river. In the sentence "He went to the bank," it's unclear if the reference is to a place for financial transactions or the riverbank.



1.2 Syntactic Ambiguity

- **Definition:** Syntactic ambiguity arises when a sentence has multiple possible grammatical structures, leading to different interpretations.
- **Example:** The sentence "I saw the man with the telescope" can mean either "I used a telescope to see the man" or "I saw a man who had a telescope."

1.3 Semantic Ambiguity

- **Definition:** Semantic ambiguity occurs when the overall meaning of a sentence is unclear, even if each word is understood individually.
- **Example:** The sentence "Visiting relatives can be annoying" is ambiguous, as it could mean either that visiting relatives is annoying or that relatives who visit can be annoying.

Language Complexity Issues

Beyond ambiguity, other language complexities such as polysemy, homonymy, and context dependency further complicate NLP tasks. These complexities require a nuanced understanding of language that is challenging to model computationally.

1. Polysemy

- **Definition:** Polysemy refers to a single word having multiple related meanings.
- **Example:** The word “run” has many meanings, such as “to run a race,” “to run a business,” or “a computer program run.” Although these meanings are related, they apply in different contexts.

2. Homonymy

- **Definition:** Homonymy occurs when two words have the same spelling or pronunciation but different, unrelated meanings.
- **Example:** “Bat” can mean either a flying mammal or a piece of sports equipment used in baseball.

3. Context Dependency

- **Definition:** The meaning of many words and phrases depends heavily on the context in which they are used, making context a crucial component in NLP.
- **Example:** The word “cold” can refer to temperature, an illness, or a lack of friendliness. Without understanding the context, it is challenging to determine the correct interpretation.

Additional Complexities in Language

Other complexities, such as dialects, code-switching, sarcasm, and idiomatic expressions, make language even harder for machines to interpret accurately. These elements reflect social, cultural, and situational aspects of language, which are challenging to model computationally.

1. Dialects and Regional Variations

- **Definition:** Dialects are variations of a language spoken in different regions or by different social groups, often involving unique vocabulary, grammar, and pronunciation.
- **Example:** In American English, the term “elevator” is used, while British English uses “lift” for the same concept. Similarly, “subway” in American English refers to underground transit, while in British English, it can mean a pedestrian underpass.

2. Code-Switching

- **Definition:** Code-switching is the practice of mixing languages or dialects within a single conversation, sentence, or even phrase.
- **Example:** In multilingual communities, people may switch between languages, as in Hinglish (Hindi-English) sentences like, “Let’s go to the mall and do some shopping jaar.”

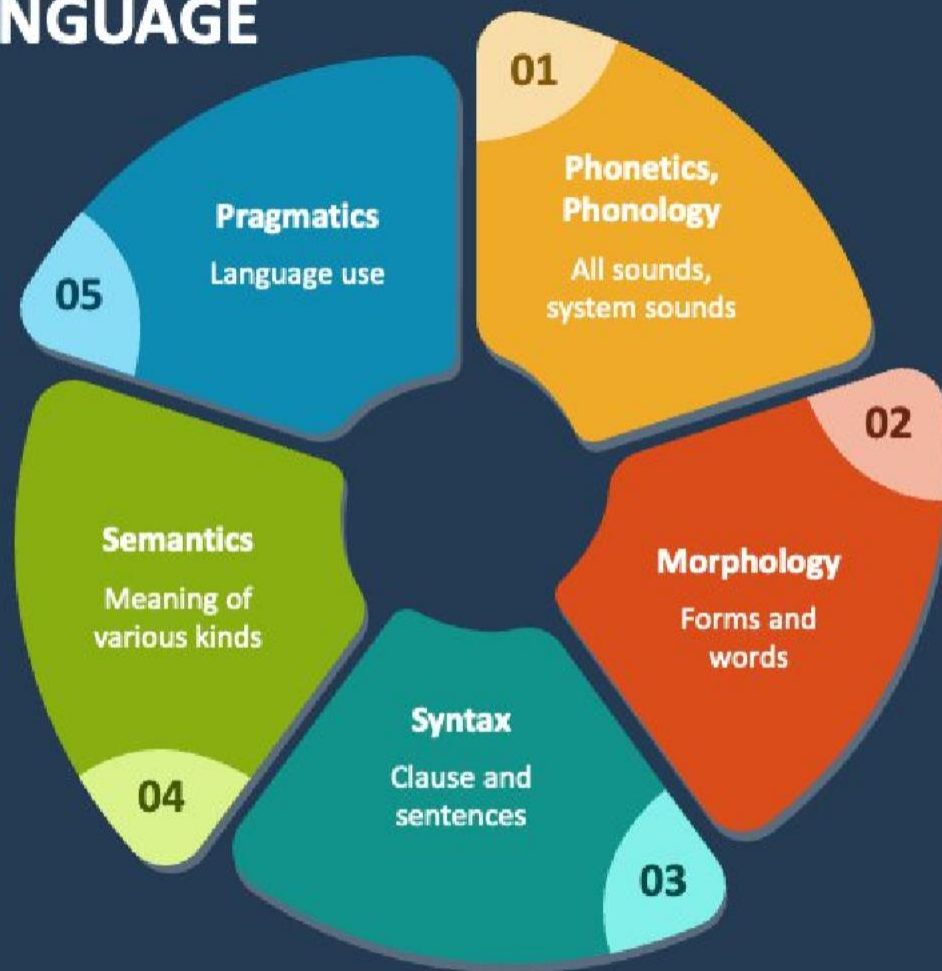
3. Sarcasm and Irony

- **Definition:** Sarcasm and irony involve saying something that means the opposite of what the words convey, often in a humorous or critical way.
- **Example:** "Oh, great! Another traffic jam!" Here, "great" is used sarcastically to mean the opposite.

4. Idiomatic Expressions

- **Definition:** Idiomatic expressions are phrases with meanings that cannot be inferred from the individual words alone.
- **Example:** "Kick the bucket" means "to die," not "to kick" a literal bucket.

LEVELS OF LANGUAGE



Syntax: Grammar and Structure of Sentences

Syntax refers to the rules that govern the structure of sentences in a language. It defines how words are arranged to form meaningful sentences and phrases.

NLP systems need to understand syntax to parse sentences and extract structural information, which is essential for many downstream tasks.

Key Concepts in Syntax

- **Parts of Speech (POS):** Parts of speech are categories of words based on their syntactic roles, such as nouns, verbs, adjectives, adverbs, and prepositions.
- **Phrase Structure:** Phrases are groups of words that function as a single unit within a sentence, such as noun phrases ("the quick brown fox") or verb phrases ("jumps over the lazy dog").
- **Syntactic Parsing:** Parsing is the process of analyzing a sentence to determine its grammatical structure. There are two common types of parsing:
 - **Dependency Parsing**
 - **Constituency Parsing**

Semantics: Meaning of Words and Sentences

Semantics deals with the meaning of words, phrases, and sentences. It is concerned with understanding what the text is about, going beyond just identifying the structure.

It is concerned with understanding what the text is about, going beyond just identifying the structure.

In NLP, semantics is critical for interpreting the content and answering questions about it.

Key Concepts in Semantics

- **Lexical Semantics:** The study of word meanings and relationships between words, such as synonyms (words with similar meanings), antonyms (words with opposite meanings), and polysemy (words with multiple related meanings).
- **Word Sense Disambiguation (WSD):** WSD is the task of determining which sense of a word is used in a given context. This is essential for understanding sentences with ambiguous words.
- **Named Entity Recognition (NER):** NER identifies and categorizes entities like names, locations, dates, and organizations in a text.

Pragmatics: Language in Context (Intent, Social Cues)

Pragmatics involves understanding language in context, which is crucial for interpreting the intent, tone, and social cues behind a statement

Pragmatics goes beyond the literal meaning of words and sentences to consider how language is used in real-world situations.

Key Concepts in Pragmatics

- **Contextual Meaning:** Pragmatics focuses on how meaning changes based on the situation, prior conversation, or shared knowledge.
 - **Example:** "Can you pass the salt?" is literally a question, but pragmatically, it's understood as a polite request.
- **Implicature:** Implicature refers to meanings implied by a speaker but not explicitly stated. The listener must infer these meanings based on context.
 - **Example:** If someone says, "It's getting late," it might imply that it's time to leave without directly saying so.
- **Deixis:** Deictic expressions depend on context to convey meaning. These include words like "this," "that," "here," and "there," which require contextual knowledge to understand.
 - **Example:** In the sentence "Let's meet here tomorrow," "here" and "tomorrow" depend on the speaker's location and the time of the statement.

Core NLP Techniques

1. Text Preprocessing

Text preprocessing transforms raw text into a usable format, ensuring it is consistent and meaningful. Key preprocessing techniques include:

- **Tokenization:** Tokenization is the process of splitting text into smaller units, or tokens, which could be words, subwords, or sentences, depending on the level of tokenization.
 - **Word Tokenization:** Splits text into individual words. For example, "I love NLP!" becomes ["I", "love", "NLP", "!"].
 - **Subword Tokenization:** Breaks down words into smaller parts, particularly useful for languages with complex morphology or for representing rare words.
 - **Sentence Tokenization:** Splits paragraphs or documents into individual sentences. For example, "NLP is interesting. I want to learn it." becomes ["NLP is interesting.", "I want to learn it."]

- **Stemming and Lemmatization:** Both techniques aim to reduce words to their root or base forms, but they approach it differently.
 - **Stemming:** Strips suffixes to obtain a rough base form, often producing non-standard words. For example, “running,” “runner,” and “runs” all become “run.”
 - **Lemmatization:** Uses linguistic knowledge to reduce words to their dictionary root or lemma. For example, “running” becomes “run” and “better” becomes “good.”
 - **Importance:** These techniques ensure consistency in text data, especially for tasks like information retrieval, where words with similar meanings need to be treated the same way.
- **Examples of Tools:**
 - **NLTK (Natural Language Toolkit):** Offers tokenizers, stop-word removal, and stemming/lemmatization functions, making it a powerful tool for text preprocessing in Python.
 - **spaCy:** A popular NLP library that provides efficient and accurate tokenization, POS tagging, dependency parsing, and named entity recognition. spaCy’s tokenizer is widely used for high-performance applications.

Word Embeddings

- Word embeddings are dense vector representations that capture semantic relationships between words by mapping them to a continuous vector space.
- Unlike BoW and TF-IDF, word embeddings preserve context by embedding words with similar meanings closer together in the vector space. Each word is represented by a fixed-length, low-dimensional vector.
- **Importance:** Word embeddings allow NLP models to understand word relationships and semantic similarity, enhancing performance in downstream tasks, such as sentiment analysis, machine translation, and information retrieval.
- Example:
 - a. Word2Vec
 - b. GloVE
 - c. Elmo
 - d. FastText

Introduction to Deep Learning for NLP

Deep learning has transformed NLP by allowing models to learn complex patterns and relationships in data. Deep learning models use neural networks, which can capture sequential and contextual information, making them more suitable for NLP tasks.

1. Recurrent Neural Networks (RNNs)

RNNs are neural networks designed for sequential data, where the output at each step depends on previous steps. They are useful for text data because they can process each word in sequence and retain information about earlier words.

- **Application:** RNNs are used in tasks like language modeling, where predicting the next word depends on previous words in the sequence.

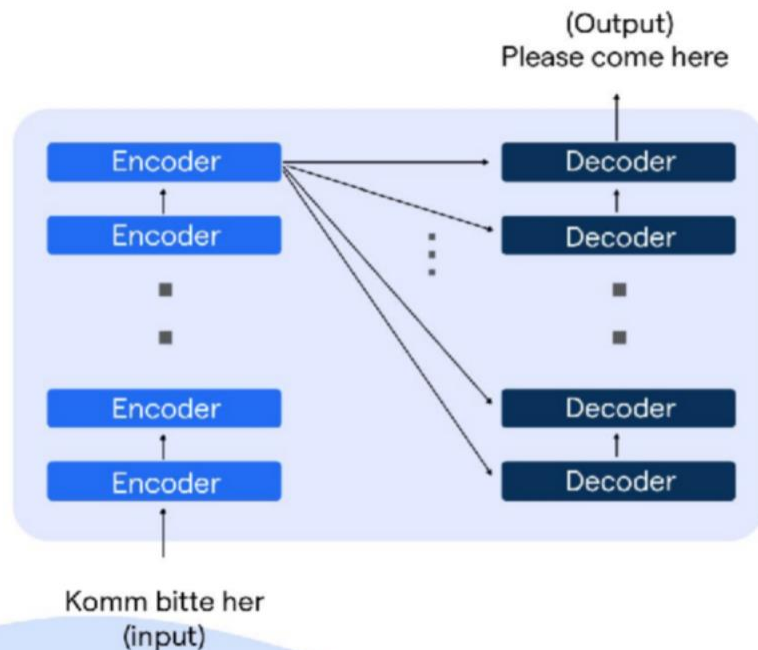
2. Long Short-Term Memory Networks (LSTMs)

LSTMs are a type of RNN designed to overcome the limitations of traditional RNNs. They include memory cells that retain information over longer sequences, allowing the model to remember important information over time.

- **Application:** LSTMs are widely used in sequence tasks like machine translation, text generation, and sentiment analysis, where long-range dependencies are critical.

NLP Transformer

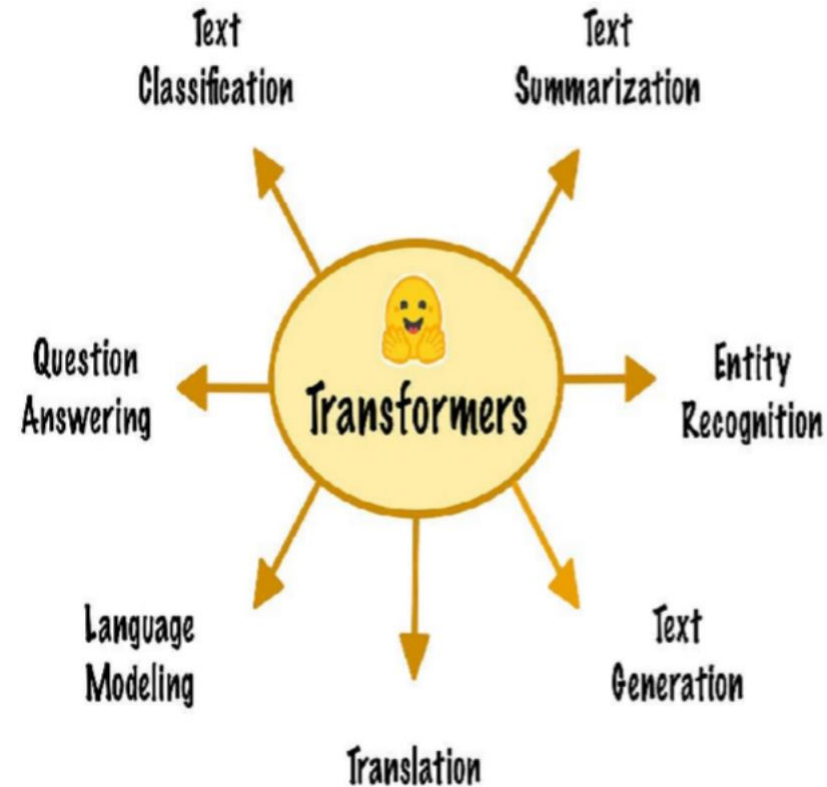
Models: Revolutionizing Language Processing



3. Transformers

Transformers are a type of neural network architecture that replaces RNNs for most NLP tasks. They use a mechanism called **self-attention**, which allows the model to focus on relevant words in a sentence, regardless of their position. Transformers enable parallel processing and improve performance on long sequences.

- **Application:** Transformers are foundational in modern NLP, powering models like BERT and GPT. They have revolutionized NLP by achieving state-of-the-art results in tasks like text classification, question answering, and machine translation.
- **Significance:** Transformers allow models to capture context and dependencies over long distances, making them highly effective for complex language tasks.



Examples of Real-World NLP Applications

NLP is embedded in a wide array of applications, transforming industries and enhancing user experiences. Some notable applications include:

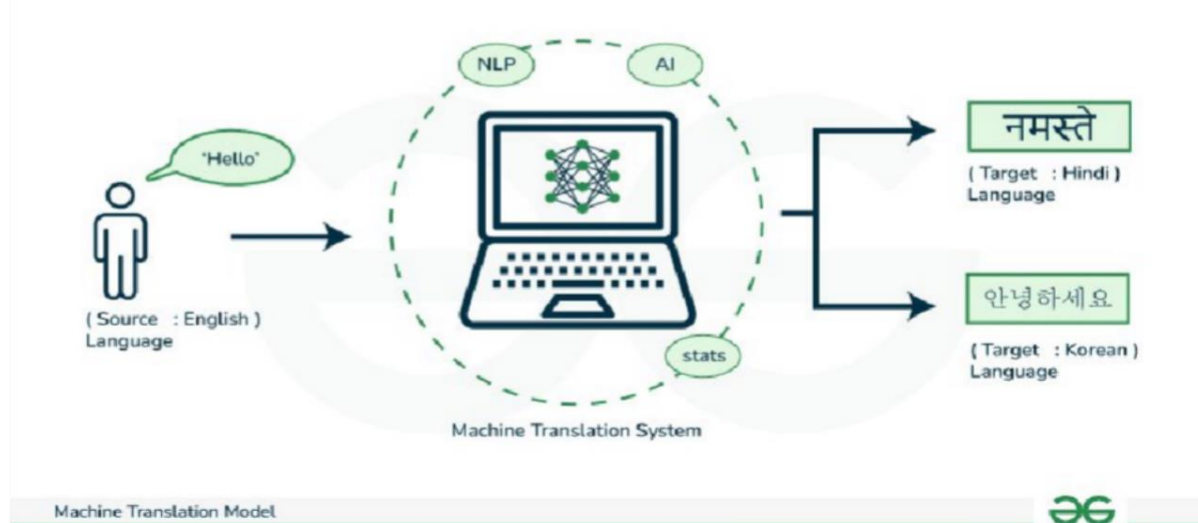
1. Chatbots and Virtual Assistants:

- **Description:** Chatbots and virtual assistants like Siri, Google Assistant, and Alexa rely on NLP to understand user queries and respond appropriately. They are programmed to understand a variety of natural language inputs, providing information, performing tasks, and interacting with users in a conversational manner.
- **Key NLP Techniques Used:**
 - Speech recognition and synthesis for spoken commands.
 - Intent recognition to understand the user's objective.
 - Dialogue management for maintaining coherent conversation flow.



2. Machine Translation:

- **Description:** Machine translation systems, such as Google Translate, enable automatic translation of text or speech from one language to another. These systems have improved dramatically due to advancements in NLP, particularly with the advent of deep learning and transformers.
- **Key NLP Techniques Used:**
 - Neural machine translation (NMT) for improved fluency and accuracy.
 - Sequence-to-sequence (Seq2Seq) models with attention mechanisms to better handle long sentences.
 - Transfer learning to support multilingual models, allowing translation across multiple languages without separate models for each language pair.



3. Sentiment Analysis:

- **Description:** Sentiment analysis is used to determine the emotional tone behind a body of text, such as customer reviews, social media posts, or survey responses. This is widely used by companies to gauge public opinion and customer satisfaction.
- **Key NLP Techniques Used:**
 - Text classification to label text as positive, negative, or neutral.
 - Feature extraction techniques, such as TF-IDF or word embeddings, to represent text data numerically.
 - Advanced sentiment analysis may use BERT or other transformers to capture subtle nuances and context in text.

SENTIMENT ANALYSIS



POSITIVE

"Great service for an affordable price.
We will definitely be booking again."



NEUTRAL

"Just booked two nights
at this hotel."

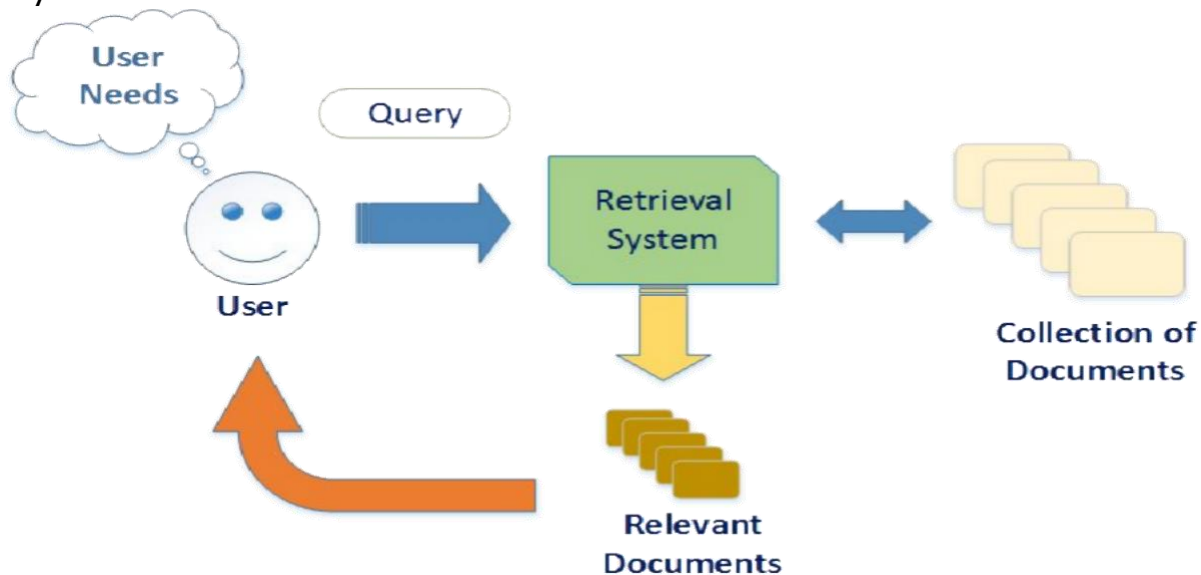


NEGATIVE

"Horrible services. The room
was dirty and unpleasant.
Not worth the money."

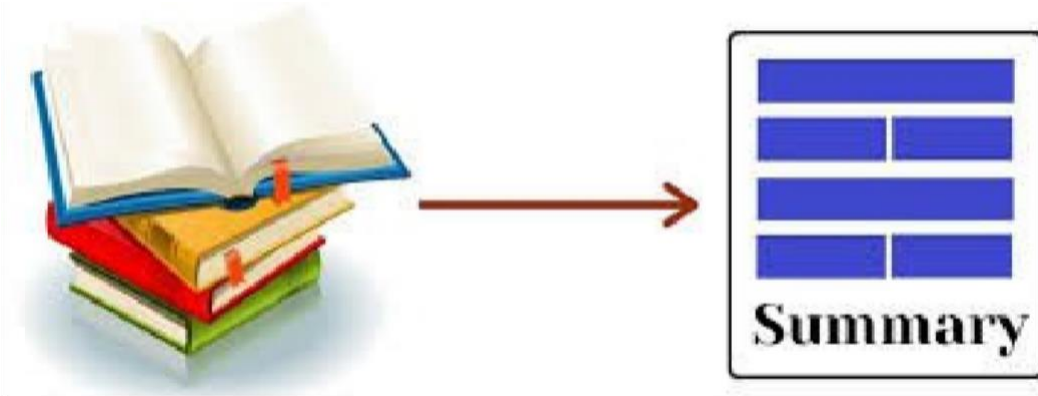
4. Information Retrieval and Search Engines:

- **Description:** Search engines like Google use NLP to understand queries and retrieve relevant information. NLP enables the search engine to go beyond simple keyword matching, understanding user intent and ranking results based on relevance.
- **Key NLP Techniques Used:**
 - Query understanding, which involves parsing and analyzing user queries to interpret intent.
 - Ranking algorithms that use NLP techniques to prioritize search results based on query relevance.
 - Synonym and entity recognition, allowing search engines to return results even if keywords don't match exactly.



5. Text Summarization:

- **Description:** Text summarization automatically creates concise summaries of larger text documents, which is particularly useful in journalism, legal, and academic fields.
- **Key NLP Techniques Used:**
 - Extractive summarization, where key sentences or phrases are extracted from the original text.
 - Abstractive summarization, where the system generates new sentences to capture the meaning of the text, often using Seq2Seq models with attention or transformers.



Thank You