# GPE in NER

in ner tagger, countries are identifies as gpe , what is gpe

In Named Entity Recognition (NER), **GPE** stands for **Geopolitical Entity**. It refers to entities that are associated with countries, cities, states, or other political regions. GPE captures locations that are politically defined, such as:

- Countries (e.g., "India", "USA")

- States or provinces (e.g., "California", "Telangana")

- Cities (e.g., "New York", "Hyderabad")

For example, in a sentence like **"Barack Obama was born in Hawaii."**, Hawaii would be tagged as **GPE** because it refers to a geopolitical region.

Let me know if you'd like further examples or need clarification!

can explain in detail about Tagging in NLP

Sure! Let's dive into **Tagging in NLP (Natural Language Processing)** with a beginner-friendly explanation:

---

## What is Tagging in NLP?

Tagging in NLP involves labeling or assigning metadata (tags) to words or phrases in a given text to understand their meaning, function, or relationship. It helps machines process and analyze text more effectively.

---

## Types of Tagging in NLP

1. **Part-of-Speech (POS) Tagging:**

   - Assigns grammatical categories (like noun, verb, adjective, etc.) to each word in a sentence.

   - Example:

- Sentence: "The dog barks loudly."
- Tags: `The/DET dog/NOUN barks/VERB loudly/ADV`

2. **Named Entity Recognition (NER) Tagging:**
   - Identifies specific entities in text, such as names of people, organizations, locations, dates, etc.
   - Example:
     - Sentence: "Google was founded in California in 1998."
     - Tags: `Google/ORG was founded in California/GPE in 1998/DATE`

3. **Chunk Tagging (Syntactic Parsing):**
   - Groups related words into meaningful chunks, like noun phrases or verb phrases.
   - Example:
     - Sentence: "The quick brown fox jumps over the lazy dog."
     - Chunk Tags: `[The quick brown fox]/NP [jumps]/VP [over the lazy dog]/PP`

4. **Semantic Role Labeling (SRL):**
   - Assigns roles to words or phrases, explaining their role in a sentence (e.g., who did what to whom).
   - Example:
     - Sentence: "John gave Mary a book."
     - Tags: `John/Agent gave/Action Mary/Recipient a book/Theme`

5. **Sentiment Tagging:**
   - Labels the sentiment or emotion in the text (e.g., positive, negative, neutral).
   - Example:
     - Sentence: "I love this movie!"
     - Tag: `Positive`

6. **Dependency Parsing (Dependency Tagging):**
   - Identifies grammatical relationships between words in a sentence.
   - Example:
     - Sentence: "She enjoys reading books."
     - Dependency: "enjoys" (root) → "She" (subject), "reading" (object)

# NER Tagging in Detail

In NER, specific entities are categorized into pre-defined classes. Common NER tags include:

- **PERSON**: Names of people (e.g., "Albert Einstein")

- **ORG**: Organizations (e.g., "Google", "UN")

- **GPE**: Geopolitical Entities (e.g., "India", "Paris")

- **DATE**: Dates (e.g., "January 26, 2025")

- **TIME**: Times (e.g., "5:00 PM")

- **LOC**: Non-political locations (e.g., "Mount Everest")

# Why is Tagging Important?

1. **Text Understanding**: Helps computers interpret language structure and meaning.

2. **Search and Information Extraction**: Improves search engines and systems that extract specific details.

3. **Machine Translation**: Assists in correctly translating words with multiple meanings.

4. **Chatbots and Virtual Assistants**: Helps them understand user queries better.

5. **Sentiment Analysis**: Identifies the tone of reviews or comments.

# Tools for Tagging in NLP

Popular NLP libraries with tagging functionalities:

- **NLTK**: Provides POS and NER tagging.

- **spaCy**: Offers state-of-the-art tagging capabilities.

- **Stanford NLP**: Widely used for POS, NER, and parsing.

- **Hugging Face Transformers**: Used for more advanced tagging models like BERT.

## How Does Tagging Work?

Tagging uses algorithms and models, including:

1. **Rule-Based Models**: Uses predefined linguistic rules.

2. **Statistical Models**: Employs probabilistic models like Hidden Markov Models (HMM) or Conditional Random Fields (CRF).

3. **Deep Learning Models**: Uses neural networks like BiLSTMs or transformers (e.g., BERT) for high accuracy.

---

Let me know if you'd like examples or practice problems to better understand tagging! 😊

explain about Sequential tagger in detail

Sure! Let's dive into **Sequential Tagging** in detail.

---

## What is Sequential Tagging in NLP?

Sequential tagging refers to the process of assigning a tag to each token (word or symbol) in a sequence, such as a sentence or a document, based on its context in that sequence. The goal is to label each token with a specific category or class, taking into account not only the token itself but also its surrounding tokens.

This approach is commonly used for tasks like:

- **Part-of-Speech (POS) Tagging**

- **Named Entity Recognition (NER)**

- **Chunking (Syntactic Parsing)**

- **Semantic Role Labeling (SRL)**

Sequential tagging works because it processes the text in a sequence (word by word), and the tag for each word depends on the words around it.

---

## Why Sequential Tagging?

In many NLP tasks, the correct label for each token is influenced by its neighbors. For instance, the word **"bank"** could be tagged as a **"location"** (GPE) or a **"financial institution"** (ORG), depending on the context provided by surrounding words. So, sequential tagging methods look at the entire sequence of words to make the best prediction.

---

## How Does Sequential Tagging Work?

1. **Input Sequence**: A sequence of words (tokens) is provided as input. For example:

   - **"John lives in New York."**

2. **Tagging Process**: A tagging algorithm processes each word in the sequence and assigns a tag. The algorithm can consider:

   - The word itself.

   - The surrounding context (previous and next words).

   - Part-of-speech patterns, syntax, and semantics.

3. **Output**: Each token in the sequence is assigned a label. For example, in Named Entity Recognition (NER):

   - **"John lives in New York."**

   - Output tags: `John/PERSON lives/VERB in/ADP New York/GPE`.

---

## Types of Sequential Taggers

Sequential taggers can be based on different approaches and models:

**1. Rule-Based Sequential Taggers**

- **Definition**: Uses predefined rules to assign tags based on word patterns, syntax, or other linguistic features.

- **Example**: If a word starts with a capital letter and is followed by a known location name, tag it as a **GPE**.

- **Pros**: Can be very accurate if rules are carefully crafted.

- **Cons**: Doesn't scale well to large datasets and is not flexible enough to adapt to new or unseen data.

## 2. Statistical Models

- **Hidden Markov Models (HMM)**: A probabilistic model that assigns tags to tokens by considering the probability of a tag given the previous tag.

  - It works by defining a set of hidden states (tags) and observing the sequence of words to calculate the most probable tag sequence.

  - **Example**: The tag for **"lives"** would be influenced by the previous tag (e.g., if **"John"** is tagged as a **PERSON**, it's more likely that **"lives"** will be tagged as a **VERB**).

- **Conditional Random Fields (CRF)**: A type of statistical model that also considers the context of neighboring tags and words.

  - CRFs are often used for tasks like POS tagging and NER because they can model dependencies between labels in a sequence.

## 3. Deep Learning Models

- **Bi-directional LSTM (BiLSTM)**: An advanced neural network-based sequential tagger that reads text from both directions (left-to-right and right-to-left). This allows it to capture the full context of each word in a sentence.

  - **Example**: For the sentence "New York is a city", the BiLSTM might tag "New York" as **GPE** because it sees the context both before and after the phrase.

- **Transformers (e.g., BERT)**: A transformer-based architecture that processes the entire input sequence at once, considering the context of each token relative to all others.

  - **Example**: BERT, fine-tuned for NER, can predict whether **"John"** refers to a **PERSON** or another type of entity based on the entire surrounding text.

# Key Concepts in Sequential Tagging Models

## 1. Labeling Sequence Dependency

Sequential tagging models take advantage of the sequential nature of text. For example:

- In **POS tagging**, the label of a word might depend on the tag of the previous word. The word **"dog"** (likely a **NOUN**) could be preceded by a determiner (**the** → **DET**).

- In **NER**, the entity label (like **GPE** for "New York") could be influenced by the surrounding words (e.g., "is" being a verb might give a clue about a place).

## 2. Contextual Information

One of the main strengths of sequential tagging is that the context of the word (from both directions) can help determine its correct tag.

- **Example**: "He went to the bank" (context suggests a financial institution).

- "He sat by the bank" (context suggests the edge of a river).

## 3. Sequence Labeling vs. Classification

- **Sequence labeling**: Every word in the sequence gets a label (as in NER or POS tagging).

- **Classification**: A single label is assigned to the entire sequence (as in document-level sentiment analysis).

---

# Challenges in Sequential Tagging

1. **Ambiguity**: Some words can have multiple meanings (e.g., **"Apple"** could refer to the company or the fruit).

2. **Long-Range Dependencies**: Sometimes, long sentences or larger contexts can make it difficult to capture all dependencies between words.

3. **OOV (Out-of-Vocabulary)**: Words not seen in the training data might be challenging for the model to tag correctly.

---

# Example of Sequential Tagging

Consider the sentence: **"Steve Jobs founded Apple in California."**

**NER Tagging:**

- **Steve Jobs** → PERSON

- **founded** → VERB

- **Apple** → ORG

- **California** → GPE

In a deep learning model like **BiLSTM** or **BERT**, the model would look at the entire sentence and use the surrounding context to decide the tags for each word.

---

## Applications of Sequential Tagging

- **Named Entity Recognition (NER)**: Identifying names of people, places, organizations, dates, etc.

- **Part-of-Speech (POS) Tagging**: Assigning grammatical categories to words in a sentence.

- **Chunking**: Grouping words into syntactic units like noun phrases or verb phrases.

- **Information Extraction**: Automatically extracting structured information from unstructured text.

---

## Tools for Sequential Tagging

- **spaCy**: Provides excellent pre-trained models for sequential tagging tasks like POS tagging and NER.

- **NLTK**: Offers tools for basic tagging like POS tagging.

- **Stanford NLP**: Well-known for its part-of-speech taggers and NER.

- **Hugging Face Transformers**: Provides state-of-the-art models (like BERT) that can be fine-tuned for sequential tagging tasks.

---

# Conclusion

Sequential tagging is a foundational technique in NLP that helps to label tokens in a sequence based on their context. It's widely used for tasks like **POS tagging**, **NER**, and **chunking**. The models range from simple rule-based systems to advanced deep learning models like BiLSTM and BERT, with each offering different strengths in terms of accuracy and context handling.

Let me know if you'd like more details on any specific part!

N-gram tagger in detail

Sure! Let's go over **N-gram Tagging** in detail.

---

# What is an N-gram Tagger?

An **N-gram tagger** is a type of sequential tagger that assigns labels (tags) to words in a sequence (like a sentence) based on the **N-gram model**. An N-gram is a contiguous sequence of **N** items (usually words) from a given text. The idea behind N-gram tagging is to use the previous **N-1 words' tags** (context) to predict the tag of the current word.

For example:

- In a **bigram (2-gram)** model, the tag of the current word depends on the previous word's tag.

- In a **trigram (3-gram)** model, the tag of the current word depends on the previous two words' tags.

---

# How Does an N-gram Tagger Work?

An N-gram tagger works by considering sequences of **N** consecutive words (or tokens) to predict the tag of the current word. This model primarily relies on the assumption that the probability of a word's tag depends on a fixed number of previous words (or their tags).

**Steps in an N-gram Tagging Model:**

1. **Training Phase:**

- The tagger is trained on a labeled corpus (a dataset where each word is already tagged with the correct label, such as Part of Speech (POS) tags).

- The model learns the patterns of N-grams in the corpus. For example, it will learn that a noun is likely to follow a determiner (e.g., "the" → **DET**, "dog" → **NOUN**).

2. **Tagging Phase:**

- For a new sentence, the N-gram tagger applies the learned patterns to predict the tags of each word in the sentence, based on the previous N-1 words (or tags).

- Example:

  - Sentence: **"The dog runs."**

  - In a **bigram tagger**, the tag for **"runs"** depends on the tag for **"dog"** (likely **NOUN**).

---

# N-gram Models

The term **N-gram** refers to the number of words the model considers at a time. Here are common types of N-gram models:

**1. Unigram Model (1-gram)**

- A **unigram tagger** only considers the current word and its tag, without looking at any previous words.

- It's simple but doesn't capture any context.

  - Example: **"dog"** → NOUN

**2. Bigram Model (2-gram)**

- A **bigram tagger** considers the current word's tag and the tag of the previous word.

- This helps capture some local context.

  - Example: For **"The dog"**, the tag for **"dog"** might depend on **"The"** (which is likely a **DET**), so it could predict **"dog"** → NOUN.

**3. Trigram Model (3-gram)**

- A **trigram tagger** looks at the current word and the tags of the two previous words.

- This allows the model to capture a broader context and sequence of dependencies.

- Example: For **"The dog runs"**, the tag for **"runs"** could depend on **"dog"** → NOUN and **"The"** → DET, improving the prediction.

### 4. Higher-order N-grams

- You can extend this idea to **4-grams**, **5-grams**, etc., where you consider even more context. However, the model becomes more complex and requires more data to train effectively.

---

## Mathematics of N-gram Tagging

The key idea in N-gram tagging is **probabilistic tagging**, where the model computes the probability of a given tag based on the context of the previous N-1 tags. The basic formula used is:

$$P(t_n|w_n, w_{n-1}, ..., w_1)$$

Where:

- $t_n$ is the tag for the current word $w_n$.

- $w_n$ is the current word.

- $w_{n-1}, ..., w_1$ are the previous words in the sequence (context).

For an **N-gram** tagger, this probability can be approximated as:

$$P(t_n|t_{n-1}, ..., t_1)$$

Which simplifies as we use the **Markov assumption**, where the probability of the current tag depends only on the previous tag (for a bigram model):

$$P(t_n|t_{n-1})$$

---

## Training an N-gram Tagger

To train an N-gram tagger, you need a **tagged corpus**. The corpus is a collection of sentences where each word is already labeled with its correct tag. For example, in POS tagging, a sentence might look like:

- **"The dog barks."**

  - Tags: **The/DET dog/NOUN barks/VERB**

The training process involves:

1. Extracting all N-grams from the training data.

2. Calculating the probabilities of different tags occurring after specific previous tags (this can be done using **Maximum Likelihood Estimation (MLE)** or smoothing techniques).

3. Building a table of probabilities for each N-gram.

---

## Smoothing in N-gram Taggers

A key challenge in N-gram models is handling **unseen N-grams**—those that never appeared in the training data. This is known as the **zero-frequency problem**. For example, the trigram **"dog barks loudly"** might never have been seen during training.

To address this, **smoothing** techniques are used. Common smoothing techniques include:

1. **Additive Smoothing (Laplace Smoothing):**

   - Adds a small constant (usually 1) to the count of each N-gram to ensure no probability is zero.

   - Example: If "dog" → **NOUN** was never seen in the training data, Laplace smoothing assigns it a non-zero probability.

2. **Good-Turing Smoothing:**

   - Adjusts the probabilities of unseen N-grams based on the frequency of N-grams that appeared once.

3. **Kneser-Ney Smoothing:**

   - A more advanced smoothing technique that improves the handling of rare N-grams.

---

## Example of N-gram Tagging

Consider the sentence **"The cat sleeps."** and let's assume we're using a **bigram tagger** trained with a small corpus.

- **Training data**:

  - "The/DET cat/NOUN sleeps/VERB"

- **Prediction**:

  - Tag "The" → **DET**

  - Tag "cat" → **NOUN** (based on "The" → **DET**)

  - Tag "sleeps" → **VERB** (based on "cat" → **NOUN**)

---

## Limitations of N-gram Tagging

- **Context Limitations**: N-grams are limited by the fixed window of context they consider. For example, a **bigram** tagger only looks at the previous word, missing out on longer dependencies.

- **Data Sparsity**: As the N increases, the number of possible N-grams grows exponentially, leading to issues with sparse data (e.g., not enough occurrences of specific N-grams).

- **Memory Usage**: Higher-order N-grams require more memory to store all the possible sequences.

---

## Applications of N-gram Tagging

- **Part-of-Speech (POS) Tagging**: Assigning grammatical tags (noun, verb, adjective, etc.) to words in a sentence.

- **Named Entity Recognition (NER)**: Identifying and classifying entities such as people, locations, and organizations in text.

- **Text Prediction and Autocompletion**: Predicting the next word in a sequence based on previous words.

---

## Tools for N-gram Tagging

- **NLTK**: A popular library that provides tools for building and training N-gram models.

- **spaCy**: While primarily using deep learning models, it also supports N-gram-based features.

- **Stanford NLP**: Provides a rule-based POS tagger that includes N-gram features.

---

## Conclusion

N-gram tagging is a simple and powerful method for sequential labeling in NLP. It models the relationship between words and their tags based on the previous N-1 words, making it effective for tasks like **POS tagging**, **NER**, and **text generation**. However, it has its limitations in terms of context depth and data sparsity, which is why it is often combined with more advanced techniques like Hidden Markov Models (HMMs) or deep learning-based models.

Let me know if you need more details or examples! 😊

> Regex tagger in nlp in detail, no need info about tools

Sure! Let's dive into **Regex Tagger** in NLP.

---

## What is a Regex Tagger?

A **Regex Tagger** in NLP is a type of tagger that assigns tags to words in a sequence based on **regular expressions (regex)**, which are patterns that match certain sequences of characters in text. It's a **rule-based** approach to tagging, where predefined regular expression patterns are used to match and tag words in a sentence.

The key idea behind the **Regex tagger** is to define patterns that can match specific types of words or sequences, such as **nouns**, **verbs**, **dates**, **email addresses**, etc. Once a pattern is matched, a tag is assigned to the corresponding word or token.

---

# How Does a Regex Tagger Work?

A **Regex Tagger** operates by applying regex patterns to a given input text (usually a sequence of words or tokens) and matching words to these patterns. Once a match is found, a tag (label) is assigned to the word.

**Steps in Regex Tagging:**

1. **Define Patterns**: Regular expressions are defined to match certain word patterns based on their characteristics (e.g., the form of the word, its ending, or special symbols).

   - **Example**:

     - **Regex for Noun**: `r'\b([A-Z][a-z]+)\b'` (matches capitalized words that could be proper nouns).

     - **Regex for Verb**: `r'\b\w+ed\b'` (matches words ending with "ed" that are likely past tense verbs).

2. **Match Patterns**: The input sequence of words is processed, and each word is checked against the defined regex patterns. If a word matches a pattern, the corresponding tag is assigned.

3. **Assign Tags**: Once a match is found for a word, a tag is assigned. This tag can be a part-of-speech (POS) tag (like **NOUN**, **VERB**, **ADJ**) or a custom tag for a specific use case (like **EMAIL**, **PHONE_NUMBER**).

4. **Return Tagged Sequence**: The output is a sequence of words with their respective tags, such as:

   - Input: **"John runs."**

   - Tags: `John/PERSON runs/VERB`

---

# Regex Patterns for Tagging

Regex taggers use various types of patterns to capture the characteristics of words and assign tags. Below are some common patterns used in regex tagging:

1. **Capitalization-based Patterns**: These patterns target capitalized words that could be proper nouns, names, or the start of a sentence.

- Example: `r'\b[A-Z][a-z]*\b'` (matches words starting with a capital letter, e.g., **John**, **New York**).

2. **Word-ending Patterns**: These patterns match words that end with specific suffixes, commonly used for parts of speech (e.g., verbs, adjectives).

   - Example: `r'\b\w+ed\b'` (matches past tense verbs ending with "ed", e.g., **walked**, **jumped**).

3. **Digit-based Patterns**: Used to match numbers or specific digit sequences, often useful for detecting dates or numerical expressions.

   - Example: `r'\b\d+\b'` (matches numbers, e.g., **123**, **2025**).

4. **Apostrophe-based Patterns**: These patterns help in identifying possessive nouns or contractions.

   - Example: `r'\b\w+\'\w+\b'` (matches contractions or possessive nouns, e.g., **I'm**, **John's**).

5. **Specific Word Patterns**: Regex can be used to match specific words or phrases that don't follow a regular pattern but are important for tagging.

   - Example: `r'\b(Mr|Ms|Dr)\.\b'` (matches titles like **Mr.**, **Ms.**, **Dr.**).

6. **Whitespace and Sentence Boundaries**: Regex can be used to identify sentence boundaries or ensure that tags are applied correctly in context.

   - Example: `r'\b[.!?]\b'` (matches sentence-ending punctuation marks).

---

## Advantages of Regex Taggers

1. **Simplicity**: Regex-based taggers are relatively simple to implement and understand, especially for small or rule-based tasks.

2. **Interpretability:** The rules (patterns) are explicit, making it easier to explain why a word received a certain tag.

3. **Efficiency**: Regex taggers are fast because they apply predefined patterns and don't rely on computationally expensive models or large datasets.

4. **Flexibility**: You can design regex patterns for a wide variety of tasks, such as detecting dates, email addresses, or specific word categories.

## Limitations of Regex Taggers

1. **Lack of Context Awareness**: Regex taggers are purely pattern-based and do not consider the broader context of the sentence. For instance, they can't account for word meanings or sentence structure.

   - Example: The word **"bank"** could refer to a financial institution or the edge of a river. A regex tagger may not be able to differentiate this without context.

2. **Complexity with Large Patterns**: As the number of patterns grows, the regex becomes harder to manage and less readable.

3. **Poor Performance with Ambiguity**: Words with multiple meanings may not be handled effectively unless additional patterns or post-processing rules are applied.

   - Example: **"bat"** could be a **NOUN** (the animal) or a **VERB** (as in "to bat a ball"), but regex taggers won't differentiate without context.

4. **Data Dependency**: Regex-based systems require manual rule crafting, which can be time-consuming and may not scale well to large datasets. If new patterns or words appear that don't match existing regex rules, they might not be tagged correctly.

---

## Example of Regex Tagging

Let's walk through an example using regex to tag parts of speech (POS).

**Example 1:**

- Input: **"The quick brown fox jumped over the lazy dog."**

We define the following regex patterns:

- `r'\b[A-Z][a-z]*\b'` : Captures capitalized words (likely proper nouns or sentence starters).

- `r'\b\w+ed\b'` : Captures verbs in past tense (words ending with "ed").

- `r'\b\w+\b'` : Captures any other word (could be a noun, adjective, etc.).

**Matching:**

- "The" → Matches: `\b[A-Z][a-z]*\b` → Tag: **DET** (Determiner)

- "quick" → Matches: `\b\w+\b` → Tag: **ADJ** (Adjective)

- "brown" → Matches: `\b\w+\b` → Tag: **ADJ** (Adjective)

- "fox" → Matches: `\b\w+\b` → Tag: **NOUN** (Noun)

- "jumped" → Matches: `\b\w+ed\b` → Tag: **VERB** (Verb)

- "over" → Matches: `\b\w+\b` → Tag: **ADP** (Adposition/Preposition)

- "the" → Matches: `\b[A-Z][a-z]*\b` → Tag: **DET** (Determiner)

- "lazy" → Matches: `\b\w+\b` → Tag: **ADJ** (Adjective)

- "dog" → Matches: `\b\w+\b` → Tag: **NOUN** (Noun)

**Output:**

- **"The/DET quick/ADJ brown/ADJ fox/NOUN jumped/VERB over/ADP the/DET lazy/ADJ dog/NOUN."**

---

## Applications of Regex Tagging

- **POS Tagging**: Regex can be used for basic POS tagging based on patterns such as capitalization, suffixes, and word structure.

- **Entity Recognition**: Regex is often used in simple **Named Entity Recognition (NER)** tasks, like recognizing email addresses, phone numbers, dates, and URLs.

- **Text Preprocessing**: It's used for text cleaning tasks, like removing punctuation, identifying tokens, or extracting specific patterns from text (e.g., dates or times).

- **Custom Tagging**: Regex can be applied for domain-specific tagging tasks, such as detecting product names, brand names, or even custom lexicons in text.

---

## Conclusion

The **Regex Tagger** is a rule-based approach to tagging in NLP that uses regular expressions to match patterns in text and assign appropriate tags. It is simple, efficient, and interpretable, but it comes with limitations, such as a lack of context sensitivity and difficulty

handling ambiguity. Regex taggers are ideal for simpler, more structured tasks but may not perform well on more complex, nuanced tasks that require deeper linguistic understanding.

Let me know if you need further clarification or examples!

> explain how we wrote that patter why r at start, why is there b , what are \ for

Great questions! Let's break down the **regular expression (regex)** pattern and explain each part in detail.

We'll take this example pattern:

`r'\b[A-Z][a-z]*\b'`

## 1. `r` at the Start

The `r` before the string (e.g., `r"pattern"`) is a way to define a **raw string** in Python.

- **Raw strings** are used to avoid having to escape backslashes (`\`). Normally, in regular strings, backslashes are special characters that escape other characters (e.g., `\n` for newline, `\t` for tab). But in regular expressions, backslashes are frequently used to denote special meanings, so to avoid confusion and the need to escape them, we use raw strings.

- By using the `r`, we tell Python to **treat the string literally**, without interpreting escape sequences.

Example:
- Without `r`: `"\\b[A-Z][a-z]*\\b"` → This string contains literal backslashes.
- With `r`: `r"\b[A-Z][a-z]*\b"` → This string is treated as a raw string, so we don't need to escape the backslashes.

So, when writing regular expressions in Python, we generally use the `r` to make things simpler.

---

## 2. `\b`

The `\b` represents a **word boundary**.

- It matches the position where a word starts or ends. A "word" in regex refers to a sequence of characters that is considered a word by the regex engine (typically, letters, digits, or underscores).

- The word boundary `\b` is **not a character**, but a position or "location" where a word begins or ends.

  - It's useful when we want to match whole words but not part of words.

  - For example, if we use `\bcat\b`, it will match the word "cat" but not "scattered" or "catalog."

In the pattern `r'\b[A-Z][a-z]*\b'`, the `\b` at both ends ensures that the word we're matching is a standalone word, not part of another word.

---

## 3. `[A-Z]`

The `[A-Z]` part is a **character class**.

- It matches **any single character** in the range **A to Z**, i.e., any uppercase English letter.

- The square brackets `[]` indicate a set or range of characters.

In this case, `[A-Z]` matches a **capital letter** at the beginning of the word.

---

## 4. `[a-z]*`

The `[a-z]` part is another **character class**.

- It matches **any single character** in the range **a to z**, i.e., any lowercase English letter.

- The `*` after the character class means **zero or more occurrences** of the preceding character set.

So, `[a-z]*` matches **zero or more lowercase letters** (a to z).

- This part of the pattern allows matching words like **"The"**, **"Dog"**, or **"John"**, but not words like **"tHE"** or **"123"**.

## Putting It All Together

Now, let's break down the complete pattern `r'\b[A-Z][a-z]*\b'` :

1. `r'...'` : This is a **raw string** to handle special characters like backslashes.
2. `\b` : A **word boundary** to ensure we match a standalone word (not part of another word).
3. `[A-Z]` : Matches the first letter, which should be an **uppercase letter** (capital letter).
4. `[a-z]*` : Matches the **remaining part of the word**, which can be any number of **lowercase letters** (or even no letters at all, as `*` allows zero occurrences).
5. `\b` : A **word boundary** again to ensure the word ends after matching.

**In Summary:**

- The pattern `r'\b[A-Z][a-z]*\b'` matches a word that starts with a capital letter followed by zero or more lowercase letters. This is typically used to match **proper nouns** or the first word in a sentence (e.g., **"John"**, **"Paris"**).

## Common Regex Components

Let's go over a few more components of regex to deepen your understanding:

- `[A-Za-z]` : Matches any letter (uppercase or lowercase).
- `\d` : Matches any digit (0-9).
- `\w` : Matches any word character (letters, digits, and underscores).
- `\s` : Matches any whitespace character (spaces, tabs, newlines).
- `+` : Matches **one or more** occurrences of the preceding pattern.
- `?` : Matches **zero or one** occurrence of the preceding pattern.
- `^` : Matches the beginning of a string.
- `$` : Matches the end of a string.

**Example:**

- `r'\d{2,4}'` : Matches a sequence of 2 to 4 digits (e.g., a year, like **2025**).
- `r'\b\w+\b'` : Matches a whole word (a sequence of letters, digits, or underscores).

## Why Use `\` in Regex?

The **backslash** `\` in regular expressions is used as an **escape character**. It changes the meaning of the character that follows it. For example:

- `\b` : Represents a **word boundary**.

- `\d` : Represents a **digit**.

- `\s` : Represents any **whitespace**.

In a regular string (not raw), backslashes are also escape characters, so we use `\\` to represent a single backslash. This is why raw strings are so useful for regular expressions— they allow us to use backslashes naturally without needing to escape them.

## Final Example:

Let's look at a more complete regex example:

`r'\b\d{4}-\d{2}-\d{2}\b'`

- `r'...'` : Raw string.

- `\b` : Word boundary.

- `\d{4}` : Matches exactly **4 digits** (e.g., the year part of a date).

- `-` : Matches the literal dash character ( `-` ).

- `\d{2}` : Matches exactly **2 digits** (e.g., the month and day).

- `\b` : Word boundary.

This pattern would match dates like **"2025-01-26"**.

## Conclusion

To summarize:

- The `r` is used for raw strings, allowing us to use backslashes without escaping them.

- The `\b` indicates word boundaries, ensuring we match entire words.
- `[A-Z]` and `[a-z]` are character classes that match uppercase and lowercase letters, respectively.
- The `*` allows repetition of characters (zero or more occurrences), and `\` is used to escape special characters.

Regular expressions are powerful for pattern matching in text, and they can be customized to match various text structures. By combining these components, you can create complex patterns for a wide range of text-processing tasks.

Let me know if you need further clarification!

explain about Brill tagger in detail

# Brill Tagger in NLP

The **Brill Tagger**, also known as the **Transformation-Based Error-Driven Learning (TBL) Tagger**, is a type of **rule-based** tagger used for **Part-of-Speech (POS)** tagging in Natural Language Processing (NLP). It is a **hybrid tagger** that combines **machine learning** and **rule-based approaches**. The tagger was developed by **Eric Brill** in the early 1990s.

## Overview

The Brill Tagger first assigns initial POS tags to the words in a sentence using a **simpler tagger** (like a probabilistic tagger or a lexicon-based tagger). After that, it applies **transformational rules** to **correct mistakes** and improve the tagging accuracy. These rules are learned through the process of **error-driven learning**.

The main idea behind the Brill Tagger is that **rule-based corrections** can improve the performance of an initial automatic tagging system by focusing on where the system makes errors.

---

## How Does the Brill Tagger Work?

### 1. Initial Tagging

- The first step is to assign **initial POS tags** to the words in the text. This is usually done using a **simpler tagger** (such as a dictionary-based tagger or a probabilistic tagger).

Example:

For the sentence "He runs fast," the initial tagger might assign:

- "He" → PRP (Pronoun)

- "runs" → VBZ (Verb, third-person singular present)

- "fast" → RB (Adverb)

## 2. Rule Creation (Transformation Rules)

- Next, the Brill Tagger uses **error-driven learning** to create **transformation rules**. These rules are applied to the initial tags to correct errors in the POS tagging.

  A transformation rule has the form:

  ```php
  <condition> => <action>
  ```

  - **Condition**: This specifies when the rule should be applied. It checks the tags and possibly the context of surrounding words.

  - **Action**: This is the change that should be made if the condition is met. It modifies the tag or applies a new tag to the word.

  For example, a rule might look like this:

  ```vbnet
  If the word is tagged as "VBZ" and the previous word is tagged as "PRP", change
  the tag of the word to "VB".
  ```

  This rule could fix errors like incorrectly tagging "He runs" as "runs" (instead of the correct infinitive verb form).

## 3. Error-Driven Learning

- The **error-driven learning** mechanism is at the core of the Brill Tagger's functioning. It iteratively **learns** the transformation rules by looking at where the initial tagger makes mistakes and creating rules to correct those mistakes.

- **Learning Process**:

  - The tagger examines the output of the initial tagging.

- It then identifies **mistakes** where the POS tag is incorrect, based on training data (a labeled corpus).

- It applies transformation rules that modify the incorrect tags to correct them.

- The rules are applied in **sequential order**, and the system continues to refine the tagging process by applying corrections based on the learned rules.

**4. Rule Application**

- After the transformation rules are learned, they are applied to the tagged sentence to refine the results.

- Each rule looks for specific conditions (like the tag of the word and its context), and if the conditions are met, it modifies the tag accordingly.

  For example, a word tagged as a verb in a context where it should be a noun might be corrected by a rule.

**5. Iteration**

- The process of error-driven learning and rule application is repeated over multiple iterations. Each iteration refines the accuracy of the POS tagging by fixing more errors.

**6. Final Output**

- After the rules have been applied and refined over multiple iterations, the final output consists of a set of corrected POS tags for the input sentence.

---

## Key Concepts in the Brill Tagger

1. **Transformation Rules**: The Brill Tagger is rule-based, and its strength lies in the transformation rules that modify the initial POS tags. These rules are crafted to capture common errors that occur in the tagging process.

   A transformation rule typically has:

   - A **condition** that defines when the rule applies.

   - An **action** that defines what happens if the condition is true.

   Example rule:

   ```
   vbnet
   ```

```
If a word is tagged as "NN" (noun), and the following word is tagged as "VBZ"
(verb), change the first tag to "NNP" (proper noun).
```

2. **Error-Driven Learning**: The Brill Tagger uses **error-driven learning**, which means that it learns rules by identifying **errors** in the initial tagging and then learning how to correct them.

3. **Contextual Correction**: Unlike a simple rule-based tagger, which might assign a tag based only on the word itself, the Brill Tagger can correct errors based on **context**. It looks at surrounding words and their tags to determine the correct tag.

4. **Iterative Refinement**: The tagging process is **iterative**, meaning that the tagger continuously improves its performance with each round of applying transformation rules, gradually reducing errors.

## Advantages of the Brill Tagger

1. **High Accuracy**: The Brill Tagger typically achieves higher accuracy than simple rule-based or probabilistic taggers because it refines the initial tagging based on real errors.

2. **Flexibility**: The Brill Tagger can be applied to many different languages and can be adapted to different types of linguistic features.

3. **Combining Rule-Based and Statistical Approaches**: It combines the strengths of both rule-based systems (which are interpretable) and statistical approaches (which can generalize well over data), making it more robust.

4. **Improvement Over Initial Tagger**: The use of error-driven learning allows the Brill Tagger to improve the output of simpler taggers and correct mistakes that might otherwise be overlooked.

## Disadvantages of the Brill Tagger

1. **Slow Learning Process**: The process of **learning transformation rules** can be slow, especially with a large corpus, because it needs to check many conditions and iterate over the data.

2. **Requires Initial Tagging**: The Brill Tagger relies on an initial tagging step, which may be inaccurate if the initial tagger is not very good. If the initial tagging is wrong, the Brill Tagger may still make mistakes even after applying rules.

3. **Limited to Rule-Based Modifications**: While the Brill Tagger is effective at improving tagger output, it doesn't have the ability to recognize new patterns in the data without explicit rules. This means it's less adaptable than fully statistical systems like Hidden Markov Models (HMMs).

4. **Large Number of Rules**: In order to achieve high accuracy, the Brill Tagger may require a large number of transformation rules, which can become difficult to manage and may lead to overfitting on the training data.

---

## Example of Brill Tagging

Consider the sentence:

**"She quickly ran home."**

1. **Initial Tagging**: The initial tagger might assign:

   - "She" → PRP (Pronoun)

   - "quickly" → RB (Adverb)

   - "ran" → VBD (Verb, past tense)

   - "home" → NN (Noun)

2. **Transformation Rule**: The Brill Tagger might learn a rule like:

   - **Condition**: If the word is tagged as **NN** and the previous word is tagged as **VBD**, change the tag of the word to **NNP** (proper noun).

   - **Action**: Change "home" from **NN** to **NNP** (if the word "home" refers to a place).

3. **Revised Tags**:

   - "She" → PRP

   - "quickly" → RB

   - "ran" → VBD

   - "home" → NNP (proper noun, referring to a place name)

## Applications of the Brill Tagger

1. **POS Tagging**: The primary application of the Brill Tagger is **POS tagging**, where it helps improve the accuracy of word categorization (verbs, nouns, adjectives, etc.).

2. **Named Entity Recognition (NER)**: Transformation-based methods can also be used in NER tasks to correct entity classifications.

3. **Morphological Tagging**: In languages with rich morphology, the Brill Tagger can help in tagging based on word structure and its context.

---

## Conclusion

The **Brill Tagger** is a powerful hybrid approach for **POS tagging** that combines **initial tagging** with **error-driven learning** and **rule-based transformations**. It is highly effective at improving the accuracy of automatic taggers, but it can be slow to train and requires a good initial tagging system. While it is not as flexible as fully statistical models, its ability to use contextual corrections through transformation rules makes it an effective and interpretable tool for many NLP tasks.

Let me know if you need further clarification on any specific part!

explain about ner tagger in detail

## Named Entity Recognition (NER) Tagger in NLP

### Overview

Named Entity Recognition (NER) is a crucial task in Natural Language Processing (NLP) that involves identifying and classifying named entities in text into predefined categories. These categories might include **names of people, locations, organizations, dates, times, monetary amounts, percentages, and more**.

For example, given the sentence:

- **"Barack Obama was born in Hawaii and served as the president of the United States."**

An NER tagger would identify:

- "Barack Obama" → **PERSON**

- "Hawaii" → **GPE** (Geopolitical Entity)
- "United States" → **GPE**

---

## Key Tasks in NER

1. **Entity Detection**: Recognizing the words or phrases that are entities (e.g., identifying "Barack Obama").

2. **Entity Classification**: Assigning a predefined label or category to each recognized entity (e.g., labeling "Barack Obama" as **PERSON**).

---

## How Does an NER Tagger Work?

An NER tagger processes text and assigns **tags** to words or phrases based on their role in the sentence. For instance:

- "Microsoft" → **ORGANIZATION**
- "New York" → **LOCATION**
- "1995" → **DATE**

**Steps in NER Tagging**

1. **Preprocessing**:

   - **Tokenization**: Breaking text into tokens (words or phrases).

   - **POS Tagging**: Part-of-Speech (POS) tagging can help understand the role of each word (e.g., "Obama" as a proper noun).

2. **Feature Extraction**:

   - **Word Features**: The NER model might consider:

     - Is the word capitalized? (e.g., proper nouns)

     - Does the word appear in a dictionary of names or locations?

     - Is the word numeric? (e.g., years or monetary values)

   - **Contextual Features**: Surrounding words and their tags (e.g., "President Obama").

- **Orthographic Features**: Patterns like capitalization, punctuation, or word length.

3. **Modeling**: NER systems use different approaches to assign entity tags:

   - **Rule-Based Approaches**: Use predefined patterns or regular expressions to identify entities (e.g., recognizing capitalized words or known prefixes like "Mr." or "Dr.").

   - **Machine Learning Approaches**: Train models to learn patterns in data using annotated examples. Common models include:

     - Hidden Markov Models (HMM)

     - Conditional Random Fields (CRF)

     - Support Vector Machines (SVM)

     - Decision Trees

   - **Deep Learning Approaches**: Use neural networks for end-to-end learning. Popular models include:

     - Long Short-Term Memory (LSTM)

     - Bidirectional LSTM (BiLSTM)

     - Transformers (e.g., BERT, RoBERTa)

4. **Named Entity Tagging**:

   - The system assigns a label to each token in the text. Common tagging schemes include:

     - **IO Scheme**:

       - I → Inside an entity.

       - O → Outside any entity.

     - **BIO Scheme**:

       - B → Beginning of an entity.

       - I → Inside an entity.

       - O → Outside any entity.

     - **BILOU Scheme**:

       - B → Beginning of an entity.

       - I → Inside an entity.

       - L → Last word of an entity.

- O → Outside any entity.
- U → Unit (single-word entity).

5. **Postprocessing**:

   - Aggregating multi-token entities (e.g., combining "New" and "York" into "New York").
   - Resolving ambiguities (e.g., deciding whether "Apple" refers to the company or the fruit).

---

# Tagging Categories in NER

Common categories used in NER include:

1. **PERSON**: Names of people (e.g., "Barack Obama").
2. **ORGANIZATION**: Names of companies, agencies, or institutions (e.g., "Microsoft", "United Nations").
3. **LOCATION**: Geographical locations (e.g., "Paris", "Mount Everest").
4. **GPE (Geopolitical Entity)**: Countries, cities, states (e.g., "India", "California").
5. **DATE**: Dates and times (e.g., "January 26, 2025", "5 PM").
6. **MONEY**: Monetary values (e.g., "$100", "10 Euros").
7. **PERCENT**: Percentages (e.g., "50%", "20 percent").
8. **TIME**: Specific times (e.g., "10:30 AM").
9. **EVENT**: Names of events (e.g., "World War II").
10. **PRODUCT**: Product names (e.g., "iPhone", "Coca-Cola").
11. **FACILITY**: Physical structures (e.g., "Eiffel Tower").

The specific tags depend on the NER system and the application.

---

# NER Tagger Approaches

## 1. Rule-Based NER

- Uses manually crafted rules and patterns to identify entities.

- Example:

  - Capitalized words → Proper nouns.

  - Words followed by "Inc.", "Ltd.", etc. → Organizations.

- Advantage: Simple and interpretable.

- Disadvantage: Cannot generalize to unseen data or complex contexts.

## 2. Statistical NER

- Uses **machine learning models** trained on annotated data to predict entity types.

- Features include:

  - POS tags.

  - Capitalization patterns.

  - Context words.

- Popular algorithms:

  - **Hidden Markov Models (HMM).**

  - **Conditional Random Fields (CRF).**

- Advantage: More flexible than rule-based methods.

- Disadvantage: Requires labeled training data.

## 3. Deep Learning-Based NER

- Uses neural networks to model text sequences and learn features automatically.

- Common architectures:

  - **BiLSTM-CRF**: Combines Bidirectional LSTMs with Conditional Random Fields for sequence tagging.

  - **Transformers (e.g., BERT)**: Pretrained language models like BERT provide contextual embeddings for words, enabling state-of-the-art NER performance.

- Advantage: High accuracy and ability to handle complex contexts.

- Disadvantage: Computationally expensive and requires large datasets.

# Challenges in NER

1. **Ambiguity**: Words can belong to multiple categories depending on the context. Example: "Apple" (the fruit) vs. "Apple" (the company).

2. **Variations**: Entities can have multiple forms (e.g., "New York City" vs. "NYC").

3. **Out-of-Vocabulary (OOV) Words**: Handling unseen or rare entities not present in the training data.

4. **Domain Adaptation**: NER models trained on one domain (e.g., news articles) may not perform well in another domain (e.g., medical texts).

5. **Multilingual Support**: Entity recognition in languages with different structures and scripts is challenging.

---

# Applications of NER

1. **Information Extraction**: Extracting relevant entities from unstructured text (e.g., extracting people or organization names from news articles).

2. **Search Engines**: Improving search results by identifying key entities in queries.

3. **Chatbots**: Recognizing user inputs like dates, names, or locations for better responses.

4. **Question Answering**: Identifying entities in a question to provide precise answers.

5. **Medical Text Processing**: Extracting drug names, symptoms, or diseases from clinical records.

6. **Sentiment Analysis**: Understanding sentiment towards specific entities in reviews or social media.

---

# Example of NER Tagging

**Input Sentence:**

**"Amazon launched its new service in Seattle on January 25, 2025."**

**NER Output:**

| Word | Entity Type |
|---|---|
| Amazon | ORGANIZATION |
| launched | O |
| its | O |
| new | O |
| service | O |
| in | O |
| Seattle | GPE |
| on | O |
| January 25, 2025 | DATE |

## Conclusion

NER taggers are essential tools in NLP for identifying and classifying entities in text. They use a combination of rule-based, statistical, and deep learning methods to achieve high accuracy. Despite their challenges, NER taggers play a critical role in applications like information extraction, question answering, and data mining.

If you'd like an example of implementing an NER tagger in Python or further clarification on any part, let me know!

explain affix tagger in detail

## Affix Tagger in NLP

An **Affix Tagger** is a type of Part-of-Speech (POS) tagger in Natural Language Processing (NLP) that uses **prefixes** or **suffixes** of words (collectively known as **affixes**) to predict their POS tags. It is particularly useful for tagging words that follow specific morphological patterns, such as inflected forms or derivational forms, and works well when encountering **unknown words** (words not present in the training vocabulary).

## Key Concepts

1. **Affixes**:
   - **Prefix**: A sequence of characters at the beginning of a word (e.g., "un" in "unable").
   - **Suffix**: A sequence of characters at the end of a word (e.g., "ing" in "running").

2. **Hypothesis**:

   Words with similar affixes often belong to the same part of speech. For example:
   - Words ending in **-ing** are often **verbs** or **gerunds** (e.g., "running", "singing").
   - Words ending in **-ly** are often **adverbs** (e.g., "quickly", "slowly").
   - Words starting with **re-** are often verbs (e.g., "redo", "rewrite").

3. **Rule-Based Tagging**:

   The Affix Tagger applies rules based on observed affixes to predict the POS tag of a word.

---

## How the Affix Tagger Works

### 1. Training Phase

The Affix Tagger is trained on a labeled corpus where each word is already associated with a POS tag. During training, the tagger identifies patterns in the **prefixes** and **suffixes** of words and associates these patterns with the most frequently occurring POS tag for each affix.

For example:

- If the tagger sees that most words ending with **-ing** are tagged as **VERB**, it will create a rule associating the suffix **-ing** with the tag **VERB**.
- Similarly, it might learn that words ending with **-ed** are usually past-tense verbs (**VBD**).

### 2. Tagging Phase

When tagging new text, the Affix Tagger checks if a word matches any of the learned affix rules. If a rule applies, the corresponding POS tag is assigned to the word. If no rule matches, the tagger may fall back to a **default tagger** or other methods to predict the tag.

---

## Example of Affix Tagging

**Training Corpus:**

| Word | POS Tag |
|------|---------|
| running | VERB |
| singing | VERB |
| quickly | ADV |
| slowly | ADV |
| happiness | NOUN |

**Learned Rules:**

- Words ending in **-ing** → VERB

- Words ending in **-ly** → ADV

- Words ending in **-ness** → NOUN

**Test Sentence:**

**"He is running happily towards his goals."**

**Steps:**

1. The Affix Tagger checks each word:

   - **running**: Matches the rule **-ing** → **VERB**, so it is tagged as **VERB**.

   - **happily**: Matches the rule **-ly** → **ADV**, so it is tagged as **ADV**.

   - **towards**: No matching affix, so a fallback tagger (e.g., default tagger) assigns a tag.

   - **goals**: No matching affix, so fallback logic applies.

2. Tagged Sentence: **"He/PRON is/AUX running/VERB happily/ADV towards/ADP his/PRON goals/NOUN."**

---

# Advantages of the Affix Tagger

1. **Handles Unknown Words:**
   The Affix Tagger is effective for tagging **out-of-vocabulary (OOV)** words, especially in morphologically rich languages where affixes often indicate a word's grammatical role.

2. **Simplicity:**
   It is straightforward to implement and relies on simple rules.

3. **Efficient for Morphological Languages**:
   In languages with rich morphology (e.g., Turkish, Finnish, or Hindi), affixes play a critical role in determining the POS of a word.

4. **Can Be Combined with Other Taggers**:
   The Affix Tagger often serves as a fallback tagger or part of a **backoff chain** in conjunction with other taggers like unigram or bigram taggers.

---

## Disadvantages of the Affix Tagger

1. **Limited Context Awareness**:
   The Affix Tagger only considers word affixes and does not account for the context of surrounding words. This makes it less accurate for ambiguous words.

2. **Dependency on Training Data**:
   If the training corpus is not comprehensive, the tagger may fail to learn useful affix rules, leading to incorrect predictions.

3. **Overgeneralization**:
   An affix can sometimes be associated with multiple POS tags. For example:

   - **-ing** can indicate a verb (e.g., "running") or a noun (e.g., "building"). The Affix Tagger might overgeneralize and predict one tag incorrectly.

---

## Applications of Affix Taggers

1. **Morphologically Rich Languages**:
   In languages with a strong reliance on affixes for grammatical structure, an Affix Tagger can play a significant role in POS tagging.

2. **Fallback for Hybrid Tagging Systems**:
   When combined with other tagging techniques, the Affix Tagger can handle unknown or rare words.

3. **Text Processing in Low-Resource Languages**:
   For languages with limited annotated datasets, rule-based approaches like Affix Taggers can provide reasonable performance.

## Implementation in Python

Below is a simple implementation of an Affix Tagger using **NLTK** (Natural Language Toolkit):

```python
import nltk
from nltk.tag import AffixTagger
from nltk.corpus import treebank

# Load the training data (POS-tagged sentences)
train_data = treebank.tagged_sents()[:3000]

# Create an Affix Tagger (using suffixes of length 3)
affix_tagger = AffixTagger(train=train_data, affix_length=-3)

# Test sentence
sentence = ["The", "company", "is", "expanding", "rapidly"]

# POS tagging
tagged_sentence = affix_tagger.tag(sentence)
print(tagged_sentence)
```

**Output:**

```python
[('The', None), ('company', None), ('is', None), ('expanding', 'VBG'), ('rapidly',
'RB')]
```

Here, the Affix Tagger correctly tags "expanding" as **VBG** (gerund/verb) and "rapidly" as **RB** (adverb) based on their suffixes.

## Conclusion

The Affix Tagger is a simple yet effective method for POS tagging, especially for handling morphologically rich languages and unknown words. While it is not as powerful as context-

aware statistical or neural taggers, it can be a valuable component in hybrid systems, providing interpretable rules and fallback capabilities. Its reliance on affixes makes it a practical choice for domains where morphological patterns are critical.