

# **Natural Language Processing**

## **Course code: CSE3015**

### **Module 2**

### **Parts of Speech Tagging**

**Prepared by**  
**Dr. Venkata Rami Reddy Ch**  
**SCOPE**

# Syllabus

- Parts of Speech Tagging and Named Entities
  - Tagging in NLP,
  - Sequential tagger,
  - N-gram tagger,
  - Regex tagger,
  - Brill tagger,
  - NER tagger;
- Machine learning taggers-MEC, HMM, CRF,

# Part-of-Speech (POS) Tagging

- Part-of-speech (POS) tagging is a process in NLP where each word in a text is labeled with its corresponding part of speech.
- Assigning a part-of-speech to each word in a text.
- This can include nouns, verbs, adjectives, and other grammatical categories.
- It helps algorithms understand the grammatical structure and meaning of a text.
- POS tagging is useful for a variety of NLP tasks, such as information extraction, named entity recognition, and machine translation.



# Part-of-Speech (POS) Tags

- 1.Noun:** A noun is the name of a person, place, thing, or idea.
- 2.Pronoun:** A pronoun is a word used in place of a noun.
- 3.Verb:** A verb expresses action or being.
- 4.Adjective:** An adjective modifies or describes a noun or pronoun.
- 5.Adverb:** An adverb modifies or describes a verb, an adjective, or another adverb.
- 6.Preposition:** A preposition is a word placed before a noun or pronoun to form a phrase modifying another word in the sentence.
- 7.Conjunction:** A conjunction joins words, phrases, or clauses.
- 8.Interjection:** An interjection is a word used to express emotion.
- 9.Determiner or Article:** A grammatical marker of definiteness (the) or indefiniteness (a, an).

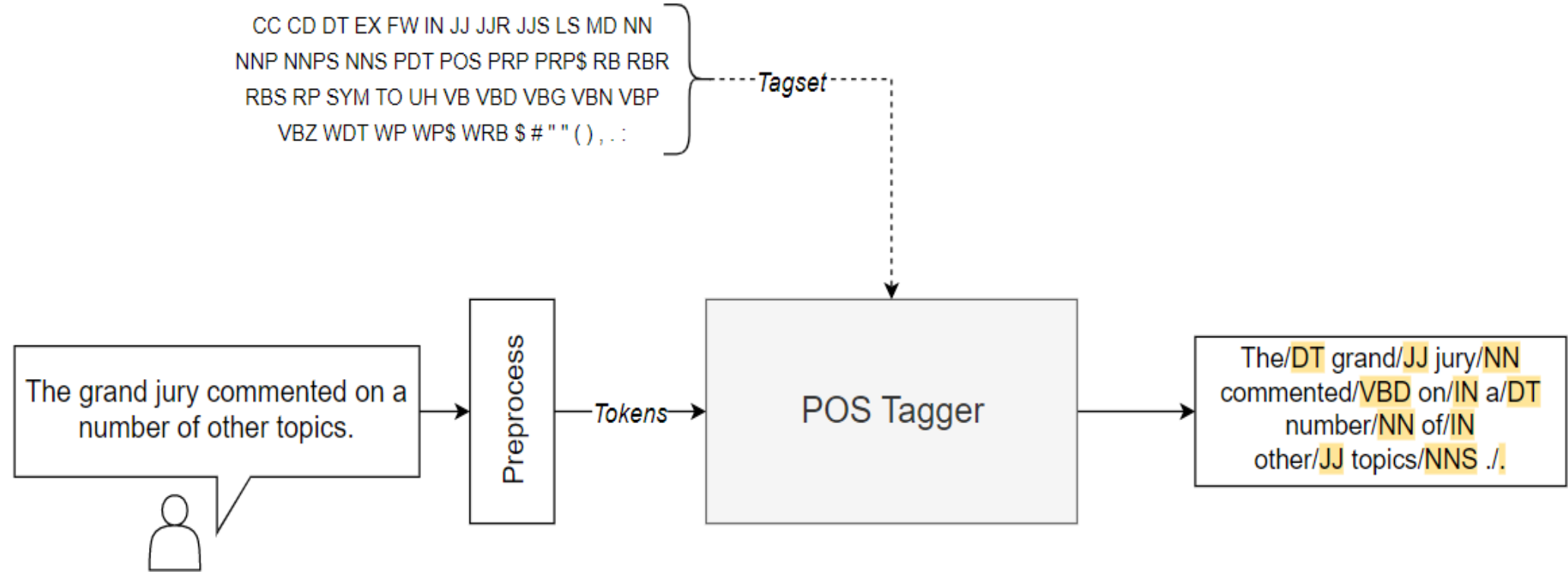
Let's take an example,

Text: "The cat sat on the mat."

POS tags:

- The: determiner
- cat: noun
- sat: verb
- on: preposition
- the: determiner
- mat: noun

# POS Tagging architecture



# POS Tagset

- Tagset is the collection of tags from which the tagger finds appropriate tags and attaches to the word
- Different POS tag sets are used depending on the language and the application.

Tagset	Number of Tags
Penn Treebank	36+ (includes punctuation)
Universal Dependencies (UD)	17
Brown Corpus	87
CLAWS	70-100
Stanford POS Tagset	47
Simplified Tagset (UD)	12

# Penn Treebank POS Tagset

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	<i>+, %, &amp;</i>
CD	cardinal number	<i>one, two, three</i>	TO	“to”	<i>to</i>
DT	determiner	<i>a, the</i>	UH	interjection	<i>ah, oops</i>
EX	existential ‘there’	<i>there</i>	VB	verb, base form	<i>eat</i>
FW	foreign word	<i>mea culpa</i>	VBD	verb, past tense	<i>ate</i>
IN	preposition/sub-conj	<i>of, in, by</i>	VBG	verb, gerund	<i>eating</i>
JJ	adjective	<i>yellow</i>	VCN	verb, past participle	<i>eaten</i>
JJR	adj., comparative	<i>bigger</i>	VBP	verb, non-3sg pres	<i>eat</i>
JJS	adj., superlative	<i>wildest</i>	VBZ	verb, 3sg pres	<i>eats</i>
LS	list item marker	<i>1, 2, One</i>	WDT	wh-determiner	<i>which, that</i>
MD	modal	<i>can, should</i>	WP	wh-pronoun	<i>what, who</i>
NN	noun, sing. or mass	<i>llama</i>	WP\$	possessive wh-	<i>whose</i>
NNS	noun, plural	<i>llamas</i>	WRB	wh-adverb	<i>how, where</i>
NNP	proper noun, singular	<i>IBM</i>	\$	dollar sign	<i>\$</i>
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	<i>#</i>
PDT	predeterminer	<i>all, both</i>	“	left quote	<i>‘ or “</i>
POS	possessive ending	<i>’s</i>	”	right quote	<i>’ or ”</i>
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	<i>[, (, {, &lt;</i>
PRP\$	possessive pronoun	<i>your, one’s</i>	)	right parenthesis	<i>], ), }, &gt;</i>
RB	adverb	<i>quickly, never</i>	,	comma	<i>,</i>
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	<i>. ! ?</i>
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	<i>: ; ... --</i>
RP	particle	<i>up, off</i>			

# lexicons

- **lexicons** are structured collections of words or phrases that include additional information, such as part of speech, meanings, synonyms, or domain-specific attributes.
- A dictionary that contains words and their possible tags.
- The **WordNet Lexicon** is a widely used lexical database in NLP that groups words into **nouns, verbs, adjectives, and adverbs**.
- For example, the word "run" might have tags such as *verb* or *noun*.



# POS Tagger

- A part-of-speech tagger, or **POS-tagger**, processes a sequence of words, and attaches a part of speech tag to each word.
- It is a program that carries out POS Tagging
- Taggers utilize a various types of data: lexicons, dictionaries, rules, etc. for POS tagging.

```
import nltk
```

```
nltk.download('averaged_perceptron_tagger')
```

```
text = "The quick brown fox jumps over the lazy dog."
```

```
tokens = nltk.word_tokenize(text)
```

```
# Perform POS tagging
```

```
pos_tags = nltk.pos_tag(tokens)
```

```
print("Tokenized Words with POS Tags:")
```

```
for word, tag in pos_tags:
```

```
    print(f"{word}: {tag}")
```

Tokenized Words with POS Tags:

The: DT

quick: JJ

brown: NN

fox: NN

jumps: VBZ

over: IN

the: DT

lazy: JJ

dog: NN

..: .

- The **Averaged Perceptron Tagger** is the default tagger used in pos\_tag().
- **Averaged Perceptron Tagger** assigns tags based on learned features from a large annotated corpus, such as the **Penn Treebank**

# Approaches to POS Tagging

- **Rule-based Approach**

- Uses set of rules to tag input sentences

- e.g. RegExp Tagger

- **Statistical approaches(Machine Learning Based)**

- Use training corpus to assign a tag to every token in given text.

- e.g N-gram tagger, HMM(Hidden Markov Model),  
CRF(conditional random field)

- **Transformation based(Hybrid)**

- Rules + machine learning( 1 gram tagger)

- e.g Brill Tagger

# Regex tagger

- The Regex tagger assigns tags to words based on matching patterns specified using regular expressions.
- You can specify both a regular expression and an associated tag for identifying pos of each word and assign associated tag to that word.
- For instance, we might guess that any word ending in *ed* is the past participle of a verb, and any word ending with 's is a possessive noun. We can express these as a list of regular expressions

<b>Determiners:</b>	<code>\b(the a an)\b</code>
<b>Adjectives ending in 'able' :</b>	<code>\b\w+able\b'</code>
<b>Past tense verbs:</b>	<code>.*ed\$</code>
<b>Adverbs:</b>	<code>.*ly\$'</code>
<b>Pronouns:</b>	<code>\b(I my he him his she her we you it they)\b</code>
<b>Prepositions:</b>	<code>\b(on in at by with about into to)\b</code>

`\b`: Word boundary to ensure you're matching whole words.

**Write a Python script** to tag parts of speech in the given sentence. Define patterns for pronouns, conjunctions, prepositions, determiners, adjectives, verbs, adverbs, and nouns. Then, print the tagged words.

```
import nltk
from nltk.tag import RegexpTagger
```

```
# Define patterns for the regular expression tagger
```

```
patterns = [
    (r'^\d+$', 'CD'),          # Cardinal numbers
    (r'\b(I|me|my|he|him|his|she|her|we|you|it|they)\b', 'PRP'), # pronouns
    (r'\b(on|in|at|by|with|about|into|to)\b', 'IN'),
    (r'\b(and|or|but|also)\b', 'CC'),
    (r'\b(The|the|A|a|An|an)\b', 'DT'), # Determiners
    (r'\b\w+able\b', 'JJ'),      # Adjectives ending in 'able'
    (r'.*ing$', 'VBG'),         # Gerunds
    (r'.*ed$', 'VBD'),          # Past tense verbs
    (r'.*es$', 'VBZ'),          # 3rd person singular verbs
    (r'.*ly$', 'RB'),           # Adverbs
    (r'.*', 'NN')               # Default: Noun
]
```

```
# Create a RegexpTagger using the defined patterns
```

```
regexp_tagger = RegexpTagger(patterns)
```

```
sentence = "John is running quickly and he catch the 9am train."
```

```
words = nltk.word_tokenize(sentence)
```

```
[('John', 'NN'), ('is', 'NN'), ('running', 'VBG'),
 ('quickly', 'RB'), ('and', 'CC'), ('he', 'PRP'),
 ('catch', 'NN'), ('the', 'DT'), ('9am', 'NN'), ('train',
 'NN'), ('.', 'NN')]
```

```
# Tag the words using the regular expression tagger
```

```
tagged_words = regexp_tagger.tag(words)
```

```
print(tagged_words)
```

# Unigram Tagger

- A **Unigram Tagger** is a type of POS tagger that assigns tags based on individual words.
- In a unigram model, the tag for a word is determined independently, without considering the POS tags of the previous words.

## How it Works:

### Training:

- The Unigram Tagger is trained on a collection of tagged words (word-tag pairs).
- e.g. the [Brown Corpus](#) for English. In such corpora each word is associated with its PoS.

*[('The', 'DET'), ('Fulton', 'NOUN'), ('County', 'NOUN'), ('Grand', 'ADJ'), ('Jury', 'NOUN'), ('said', 'VERB'), ('Friday', 'NOUN'), ('an', 'DET'), ('investigation', 'NOUN'), ('place', 'NOUN'), (',', ',')]*

- It learns the most frequent tag for each word in the training data.
- The result of the training is a table of two columns, the first column is a word and the second the most-frequent PoS of this word

Word	Most Frequent Tag
control	noun
run	verb
love	verb
red	adjective

# Unigram Tagger

## Tagging:

- During tagging, for a given word in a sentence, the tagger looks up the most likely tag associated with that word from the training data.
- If the word was seen during training, it assigns the most common tag associated with that word.
- If the word wasn't seen during training, it assigns a default tag (usually NN).

## Cons:

- It doesn't consider the context or POS tags of previous words. Consequently, a word is always tagged with the same POS, independent of its context.

## Implementation steps:

- Import necessary modules.
- Load the `brown corpus` and divide the data into training and testing data
- Train the `UnigramTagger` on training data
- Tag a sentence using the `tag()` method of `UnigramTagger`.

# Unigram Tagger

```
import nltk
from nltk.corpus import brown
from nltk.tag import UnigramTagger

# Download the required NLTK resources
nltk.download('brown')
nltk.download('punkt')

# Load the brown dataset for training
train_data = brown.tagged_sents()[ :3000] # Use first 3000 sentences for training
test_data = treebank.tagged_sents()[3000:] # Use remaining sentences for testing

# Create a UnigramTagger (1-gram model)
unigram_tagger = UnigramTagger(train_data)

# Tag a sentence using the trained UnigramTagger
sentence = "The dog sat on the mat".split()
tagged_sentence = unigram_tagger.tag(sentence)

print(tagged_sentence)

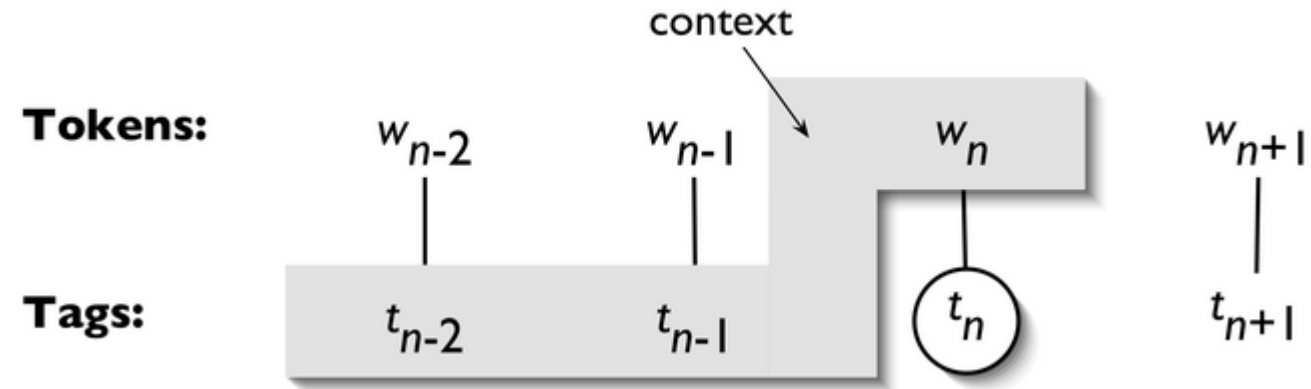
print(unigram_tagger.evaluate(test_data))
```

`[('The', 'AT'), ('dog', 'NN'), ('sat', 'VBD'), ('on', 'IN'), ('the', 'AT'), ('mat', 'NN')]`

`0.7641347348147013`

# N-gram tagger

- **N-gram tagger** in NLP uses the POS tags of  $N-1$  previous words to predict the tag for a current word.
- **N-gram-Tagger** assigns the PoS-tag of the current word by taking into account the current word itself and the PoS-tag of the  $N-1$  preceding words.



## N-gram Models:

- **Unigram Tagger (1-gram)**: It predicts the tag for the current word based on the word alone.
- **Bigram Tagger (2-gram)**: It predicts the tag for the current word based on the previous word and its tag.
- **Trigram Tagger (3-gram)**: It predicts the tag for the current word based on the previous two words and their tags.



# N-gram(bigram) tagger training

The goal during training is to compute the probabilities needed for tagging. This involves:

## a. Building the bigram model

- A **bigram** refers to a pair of consecutive tags,  $(t_{i-1}, t_i)$ , where:
  - $t_{i-1}$  is the tag of the previous word.
  - $t_i$  is the tag of the current word.
- Training data consists of sentences where words are already tagged (e.g., The/DT cat/NN sleeps/VBZ ).

From the training data:

1. **Tag Transition Probabilities:** Compute the probability of one tag following another:

$$P(t_i | t_{i-1}) = \frac{\text{Count}(t_{i-1}, t_i)}{\text{Count}(t_{i-1})}$$

This models how likely a tag  $t_i$  is given the previous tag  $t_{i-1}$ .

2. **Emission Probabilities:** Compute the probability of a word  $w_i$  being associated with a tag  $t_i$ :

$$P(w_i | t_i) = \frac{\text{Count}(t_i, w_i)}{\text{Count}(t_i)}$$

# N-gram tagger Tagging

- The tagger calculates the most probable tag sequence using the learned probabilities.
- For example:
  - In a **bigram model**, the tag  $t_i$  for a word  $w_i$  is chosen based on:

$$\operatorname{argmax}_{t_i} P(t_i | t_{i-1}) \cdot P(w_i | t_i)$$

- In a **trigram model**, it's based on:

$$\operatorname{argmax}_{t_i} P(t_i | t_{i-1}, t_{i-2}) \cdot P(w_i | t_i)$$

## Tag Assignment:

- The tagger assigns the tag with the highest probability to each word in the sentence.

# bi-gram tagger Tagging Example

## Training Data:

"The/DT cat/NN sat/VBD on/IN the/DT mat/NN

1. Compute Emission Probabilities:  $P(\text{word}|\text{tag})$  Count occurrences of each word given a tag.

Example:

- $P(\text{cat}|\text{NN}) = \frac{\text{Count}(\text{cat}, \text{NN})}{\text{Count}(\text{NN})} = \frac{1}{2}$
- $P(\text{sat}|\text{VBD}) = \frac{1}{1} = 1.0$

2. Compute Transition Probabilities:  $P(\text{current tag}|\text{previous tag})$  Count transitions between tags.

Example:

- $P(\text{NN}|\text{DT}) = \frac{\text{Count}(\text{DT} \rightarrow \text{NN})}{\text{Count}(\text{DT})} = \frac{2}{2} = 1.0$
- $P(\text{VBD}|\text{NN}) = \frac{\text{Count}(\text{NN} \rightarrow \text{VBD})}{\text{Count}(\text{NN})} = \frac{1}{2} = 0.5$

## Tagging Phase:

$$P(\text{tag}_t|\text{word}_t, \text{tag}_{t-1}) = P(\text{word}_t|\text{tag}_t) \times P(\text{tag}_t|\text{tag}_{t-1})$$

Ex: the cat is running

For word cat:

$$P(\text{cat}|\text{NN}) \times P(\text{NN}|\text{DT}) = 0.5 \times 1.0 \\ = 0.5$$

$$P(\text{cat}|\text{VBD}) \times P(\text{VBD}|\text{DT}) = 0 \times 0.5 \\ = 0$$

NN tag for cat having highest probabilities so assign tag NN to cat

```
import nltk
from nltk.corpus import brown
from nltk.tag import UnigramTagger, BigramTagger, TrigramTagger

nltk.download('brown')
nltk.download('punkt')

train_data = brown.tagged_sents()[ :3000] # First 3000 sentences for training
test_data = brown.tagged_sents()[3000:] # Remaining sentences for testing

# Create and train a Unigram Tagger
unigram_tagger = UnigramTagger(train_data)

# Create and train a Bigram Tagger
bigram_tagger = BigramTagger(train_data, backoff=unigram_tagger)

# Create and train a Trigram Tagger
trigram_tagger = TrigramTagger(train_data, backoff=bigram_tagger)

# Test the trigram tagger on a new sentence
sentence = "The dog sat on the mat".split()
tagged_sentence = trigram_tagger.tag(sentence)

print(tagged_sentence)

# Print Performnces
print(unigram_tagger.evaluate(test_data))
print(bigram_tagger.evaluate(test_data))
print(trigram_tagger.evaluate(test_data))
```

[('The', 'AT'), ('dog', 'NN'), ('sat', 'VBD'),  
('on', 'IN'), ('the', 'AT'), ('mat', 'NN')]

0.7641347348147013  
0.7734657846307614  
0.772435283624883

# Examples

<S> I like to play with it </S>

<S> You like to play </S>

Total count of unigrams (without <S> and </S> tags) :10

Total count of bigrams (including <S> and </S> tags): 12

<S> I, I like, like to, to play, play with, with it, it </S>

<S> You, You like, like to, to play, play </S>

**Unigram probabilities for:  $P(\text{like}) =$**

$$P(\text{word}) = \frac{\text{Count}(\text{word})}{\text{Total count of unigrams}}$$

$P(\text{like})$ :

- Count of "like" = 2
- Total unigrams = 10

$$P(\text{like}) = \frac{2}{10} = 0.2$$

# Examples

<S> I like to play with it </S>

<S> You like to play </S>

**Conditional probabilities:** P(like | it)

---

Conditional probability:

$$P(w_2|w_1) = \frac{\text{Count}(w_1 \rightarrow w_2)}{\text{Count}(w_1)}$$

- $P(\text{like}|\text{it}) = \frac{\text{Count}(\text{it} \rightarrow \text{like})}{\text{Count}(\text{it})} = \frac{0}{1} = 0.$

**Bigram probability for:** P(you like)

$$P(w_2|w_1) = \frac{\text{Count}(w_1 \rightarrow w_2)}{\text{Count}(w_1)}$$

- $P(\text{you like}) = \frac{\text{Count}(\text{you} \rightarrow \text{like})}{\text{Count}(\text{you})} = \frac{1}{1} = 1.0.$

**Trigram probability for:** P(you like it)

Trigram probability:

$$P(w_3|w_1, w_2) = \frac{\text{Count}(w_1 \rightarrow w_2 \rightarrow w_3)}{\text{Count}(w_1 \rightarrow w_2)}$$

P(you like it):

- Count of ("You like it") = 0 (No trigram "You like it")
- Count of ("You like") = 1

$$P(\text{you like it}) = \frac{0}{1} = 0$$

# N-gram

## Pros of N-gram Tagger

- Captures Local Context
- Improves Predictions
- Easy to Implement
- Flexibility in N-gram Size

## Cons of N-gram Tagger:

### **Requires Large Amount of Training Data:**

- N-gram models, especially for higher values of N, require large amounts of labeled data to accurately estimate the probabilities of sequences.

### **Memory Intensive:**

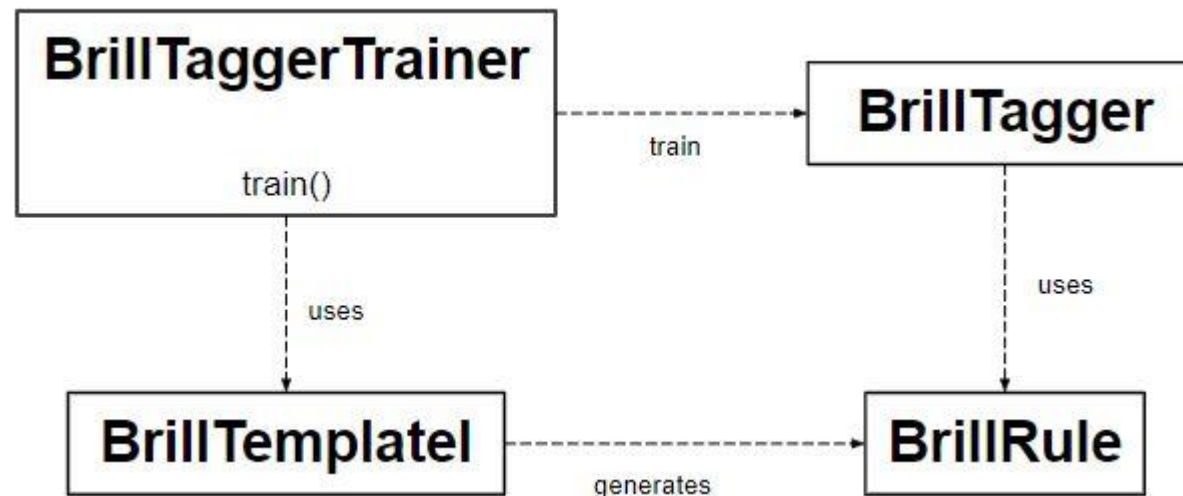
- Higher-order N-grams (trigrams, 4-grams, etc.) can become memory-intensive because they require storing a large number of tag combinations

### **Lack of Flexibility:**

- N-gram taggers are not able to capture global syntactic dependencies or hierarchical structures in sentences.

# Brill tagger

- The **Brill Tagger** is a **transformation-based tagger** introduced by Eric Brill.
- It uses an initial tagger (unigram tagger) and then applies transformation rules iteratively to correct errors made by the initial tagger.
- Moreover, it uses a series of rules to correct the results of an initial tagger.
- These rules it follows are scored based.  
This score is equal to the no. of errors they corrected minus the no. of new errors they produced.





# Training the Brill Tagger

## a. Initial Tagging

- Apply an initial tagger, which can be:
  - A simple default tagger or statistical tagger (e.g., unigram or bigram tagger) to assign tags initially.

## b. Learning Transformation Rules

### 1. Define a Template:

1. Specify templates for rules, such as:
  1. Change the tag of a word if the previous word is a specific tag.
  2. Change the tag of a word if the next word is a specific tag.

**Ex:** Change NN to VB if the previous word is a noun

### 2. Error Analysis:

- Compare the initial tagging results with the ground truth from the training data to identify tagging errors.

### 3. Rule Generation, Evaluation & Selection

- Generate potential rules from the template to correct the errors.
- Score each rule by no. of errors they corrected minus the no. of new errors they produced.
- Select the best rule based on its score and add it to the list of transformation rules.

### 4. Iterative Process:

1. Apply the selected rule to the training data and repeat the process until:
  1. A stopping condition is met (e.g., no significant improvement, maximum number of rules reached).

# Testing the Brill Tagger

## **a. Apply Initial Tagging**

- Use the same initial tagger as in the training phase to assign initial tags to the testing set.

## **b. Apply Learned Rules**

- Apply the learned transformation rules, in the same order as learned, to refine the initial tags.

## **c. Evaluate Performance**

- Compare the resulting tags with the ground truth annotations in the test set.

```
import nltk
from nltk.tag import brill, brill_trainer, UnigramTagger
from nltk.corpus import brown
```

```
nltk.download('brown')
```

```
train_data = brown.tagged_sents()[ :3000]
test_data = brown.tagged_sents()[3000:]
```

```
# Step 2: Define the Initial Tagger
initial_tagger = UnigramTagger(train_data)
```

```
# Step 3: Define Brill Tagger Templates
templates = brill.fntbl37()
```

```
# Step 4: Train the Brill Tagger
trainer = brill_trainer.BrillTaggerTrainer(initial_tagger, templates)
brill_tagger = trainer.train(train_data)
```

```
sentence = "The dog sat on the mat."
tokenized_sentence = nltk.word_tokenize(sentence)
tagged_sentence = brill_tagger.tag(tokenized_sentence)
```

```
print(tagged_sentence)
```

```
# Evaluate Brill Tagger
accuracy = brill_tagger.evaluate(test_data)
print(f"Accuracy of the Brill Tagger: {accuracy:.2f}")
```

```
[('The', 'AT'), ('dog', 'NN'), ('sat', 'VBD'),
 ('on', 'IN'), ('the', 'AT'), ('mat',
 'NN'), ('.', '.')]

```

0.79

# NER Tagger

- Named Entity Recognition (NER) is a NLP technique to find and classify entities from textual data into predefined categories called named entities.

## Types of Named Entities:

**Person:** Names of individuals (e.g., *Elon Musk, Marie Curie*).

**Organization:** Names of organizations (e.g., *Google, United Nations*).

**Location:** Geographical entities (e.g., *Paris, Mount Everest*).

**Date/Time:** Temporal expressions (e.g., *January 1, 2025, 10:30 AM*).

**Monetary values:** Financial amounts (e.g., *\$10,000*).

**Percentages:** Percentage figures (e.g., *50%*).

**Miscellaneous:** Other domain-specific categories (e.g., product names, scientific terms).

- A Named Entity Recognition (NER) tagger is a tool used in NLP to identify and classify entities within a text.
- Named entity recognition is important because it enables organizations to extract valuable information from unstructured text data.
- For example, an NER system could be used to extract the names of all the companies mentioned in a set of news articles, along with their stock prices, market capitalization, and other relevant attributes.

# NER using Spacy library

```
import spacy
```

```
# Load the pre-trained SpaCy model
nlp = spacy.load("en_core_web_sm")
```

```
# Define the text
text = """
```

```
Elon Musk, the CEO of SpaceX, announced that the company will launch a new mission to Mars
in 2025. The announcement was made in Paris on January 20, 2025, and the mission will cost
around $10,000,000. The United Nations has also shown interest in the project, especially
due to its potential to address climate change. Approximately 50% of the project's funding
will come from private investors.
```

```
"""
```

```
# Process the text
doc = nlp(text)
```

```
#print entities
print(doc.ents)
```

```
# Iterate over the entities detected in the text
for ent in doc.ents:
    print(f"Text: {ent.text}, Label: {ent.label_}")
```

(Elon Musk, SpaceX, Mars, 2025, Paris, January  
20, 2025, around \$10,000,000, The United Nation  
Approximately 50%)

Elon Musk: PERSON  
SpaceX: GPE  
Mars: LOC  
2025: DATE  
Paris: GPE  
January 20, 2025: DATE  
around \$10,000,000: MONEY  
The United Nations: ORG  
Approximately 50%: PERCENT

# Applications of NER

## **Information Extraction:**

- Extracting structured information from unstructured text (e.g., extracting dates, names, and locations from news articles).

## **Search Engine Optimization:**

- Enhancing search relevance by identifying key entities in user queries.

## **Document Categorization:**

- Automatically categorizing documents based on detected entities.

## **Question Answering Systems:**

- Understanding user queries and identifying key entities.

## **Customer Support:**

- Recognizing entities like product names or customer details in support tickets.