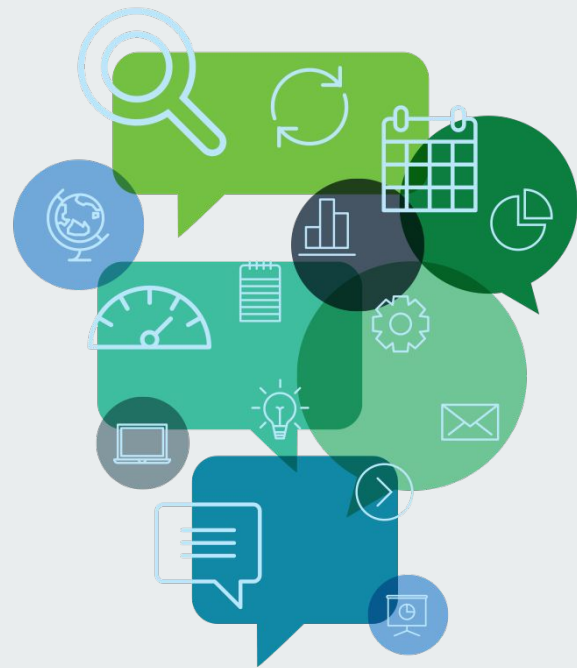


# Deep Dive 1

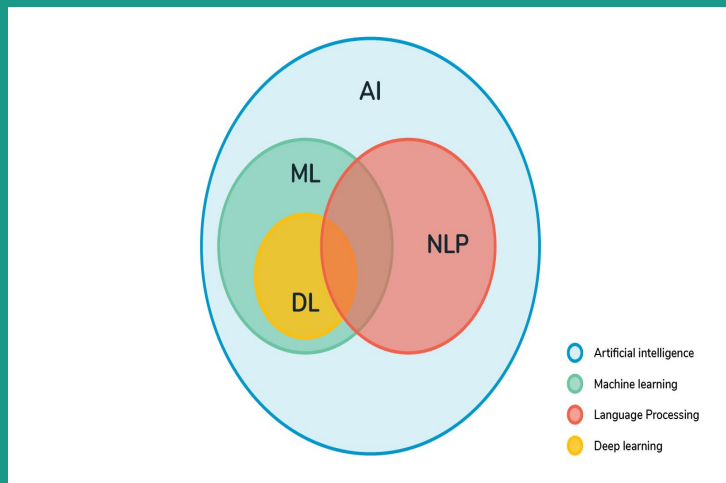
*Feature Engineering -  
POS Tagging, Chunking, Entity  
Parsing, Phrase Detection, N-Grams  
Implement with NLTK, Spacy*



# Agenda

1. Recap Week 1
2. **Feature Engineering** - Syntactic vs Semantic
3. Part of Speech (POS) Tagging
4. Shallow Parsing or Chunking
5. **Feature Engineering** - Entity Parsing
6. Named Entity Recognition
7. N-Grams
8. Implement with Google Colab
9. Questions?
10. Wrap-up and Next Steps

# Recap - Week 1



What?



Where?

# Recap - Week 1

---



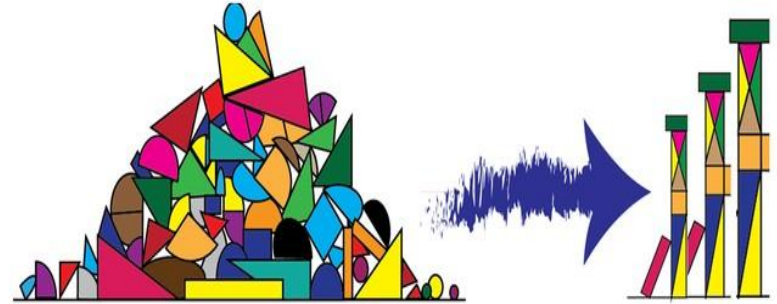
Challenges?



Must do  
Pre-processing

1

# What is Feature Engineering?



- ## One-Hot Word Representations

Source: <https://www.kdnuggets.com/2019/10/introduction-natural-language-processing.html>




# What are features?



## Boston Housing Dataset

1. average number of rooms
2. per capita crime rate by town
3. proportion of non-retail business acres per town
4. index of accessibility to radial highways

# Understanding Syntax & Structure



dog the over he  
lazy jumping is the fox  
and is quick brown

- Syntax and structure are co-dependent
- A set of specific rules, conventions, and principles govern the way words are combined
- English language constituents include: words, phrases, clauses, and sentences
- Unordered words don't convey much information
- Syntactic analysis (syntax) and semantic analysis (semantic) - primary techniques to understand natural language



# Syntactic Vs. Semantic Analysis



- Syntax is the grammatical structure of the text, semantics is the meaning conveyed
- Sentence that is syntactically correct, may not always be semantically correct
- Syntactic analysis basically assigns a semantic structure to text
- Semantic analysis is the process of understanding the meaning and interpretation of words, signs and sentence structure
- Ex, “cats flow supremely” is grammatically valid (subject – verb – adverb) but it doesn't make any sense

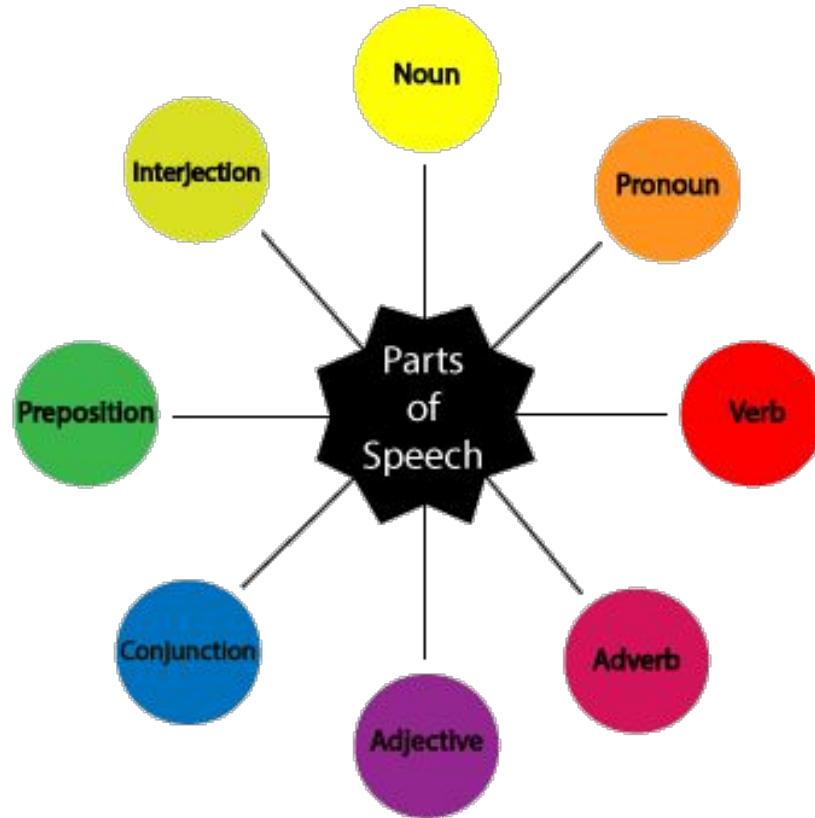
# 2

## Feature Engineering For text

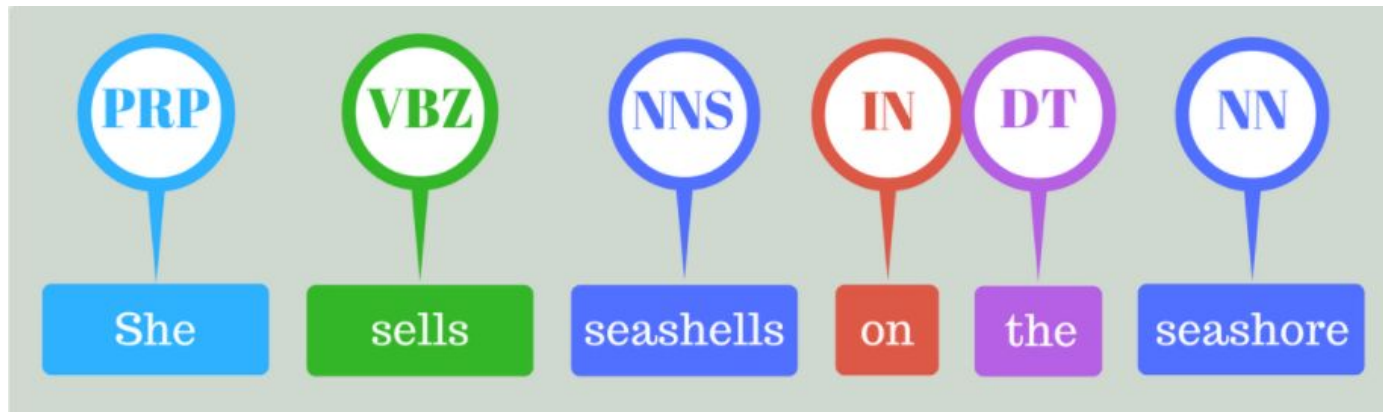


# Techniques to Understand Text

## Part of Speech (POS) Tagging

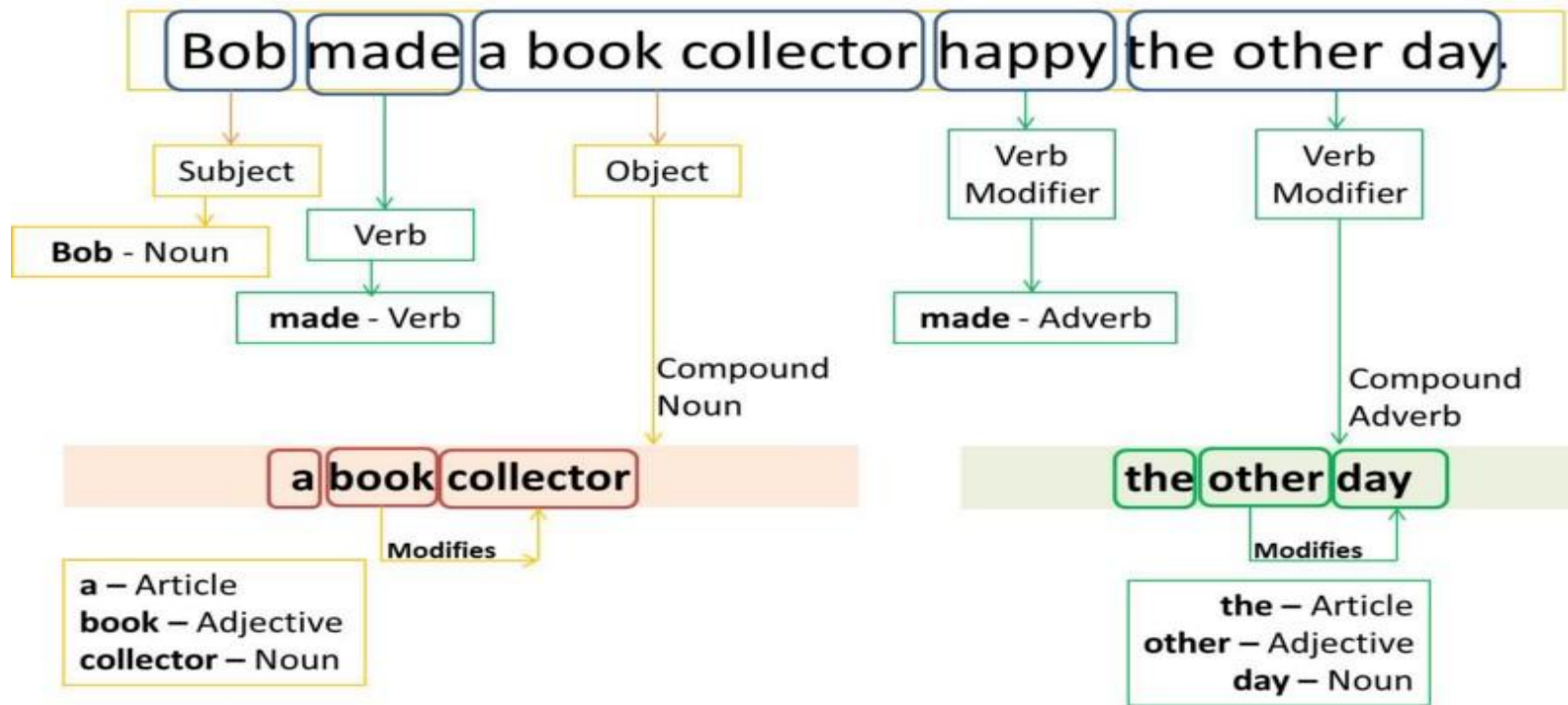


# Part-of-Speech Tagging

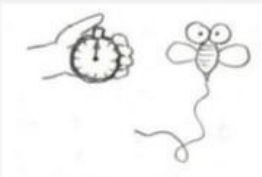
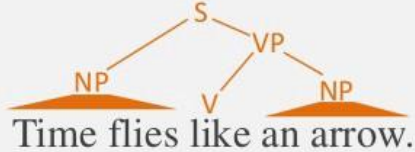
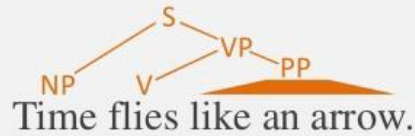


Also called grammatical tagging or word-category disambiguation, is the process of marking up a word in a text (corpus) as corresponding to a particular part of speech, based on both its definition and its context

# Part-of-Speech Tagging



# Why POS tagging is needed?



- To solve word sense disambiguation
- Information Retrieval
- Text to Speech: object(N) vs. object(V)  
E.g. Time (N) vs. Time(V)
- Machine Translation

Refer to Penn Treebank Project: [https://www.ling.upenn.edu/courses/Fall\\_2003/ling001/penn\\_treebank\\_pos.html](https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html)

# Methods for POS tagging

Rule-Based POS tagging –  
e.g. ENGTWOL [ Voutilainen, 1995 ]

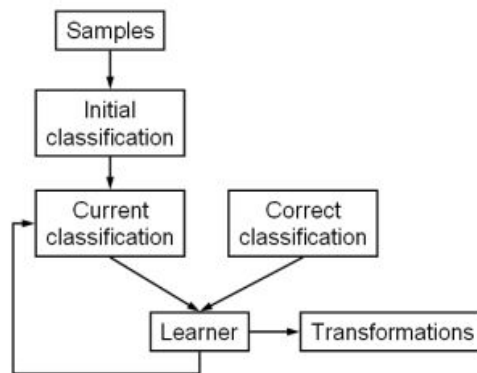
1. Use contextual information to assign tags to unknown
2. Disambiguation is done by analyzing the linguistic features of the word
3. Manual, time consuming, not scalable

Example of a rule:

If an ambiguous/unknown word X is preceded by a determiner and followed by a noun, tag it as an adjective.

Transformation-based tagging –  
e.g., Brill's tagger [ Brill, 1995 ]

1. Transformation-based Error-driven Learning (TEL)
2. Tagger is based on transformations or rules, and learns by detecting errors.



Stochastic (Probabilistic) tagging – e.g.,  
TNT [ Brants, 2000 ]

1. Based on probability of certain tag occurring
2. Necessitates a training corpus
3. Brown Corpus - 1M words
4. Hidden Markov Model (HMM) - uses both tag sequence probabilities and word frequency measurements



# POS Tagging Challenges



POS tags are not generic. Problem is **Ambiguity** in English language

A **single word** can have different tag in different sentences based on **different contexts**

Eg 1:

She saw a bear (Bear → Noun)

Your efforts will bear fruit (Bear → Verb)

Eg 2:

The trash can is hard to find (Can → Noun)

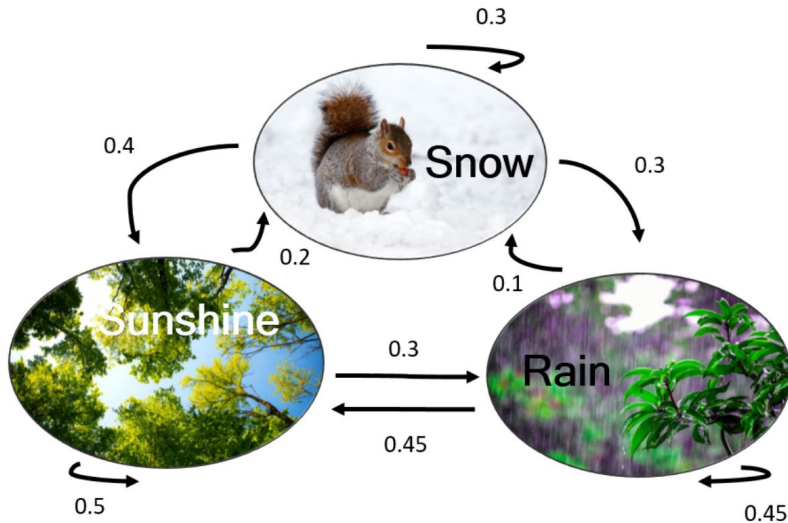
I can do better (Can → Modal Verb)

Retrieving POS tags have 2 part components:

1. Individual words with statistical preferences for their POS
2. Context has an important effect on the POS for a word

# Mathematical concept for POS tagging

## Markov Models



Basic Formula:

$$\begin{bmatrix} \text{NEXT} & \text{STATE} \end{bmatrix} = \begin{bmatrix} \text{MATRIX OF} \\ \text{TRANSITION} \\ \text{PROBABILITIES} \end{bmatrix} \begin{bmatrix} \text{CURRENT STATE} \end{bmatrix}$$

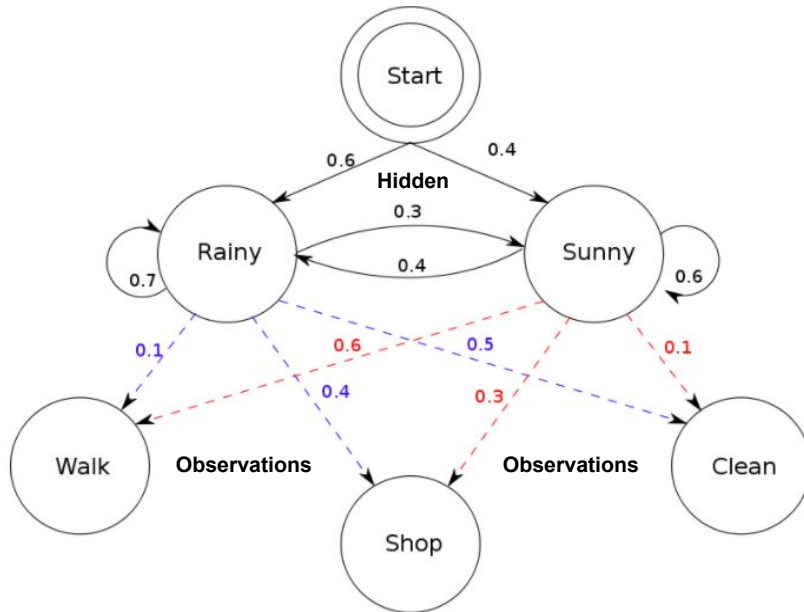
States: SUNSHINE, SNOW & RAIN

State transition probability: Decimal numbers  
(State1 → State2)

E.g.: There is 0.2 probability of SNOW tomorrow if today it is SUNSHINE.

# Mathematical concept for POS tagging

## Hidden Markov Models



**Hidden States:** SUNNY, RAIN

**Observation States:** WALK, SHOP, SUNNY

**Transition Prob:** Prob of 1 hidden state to another

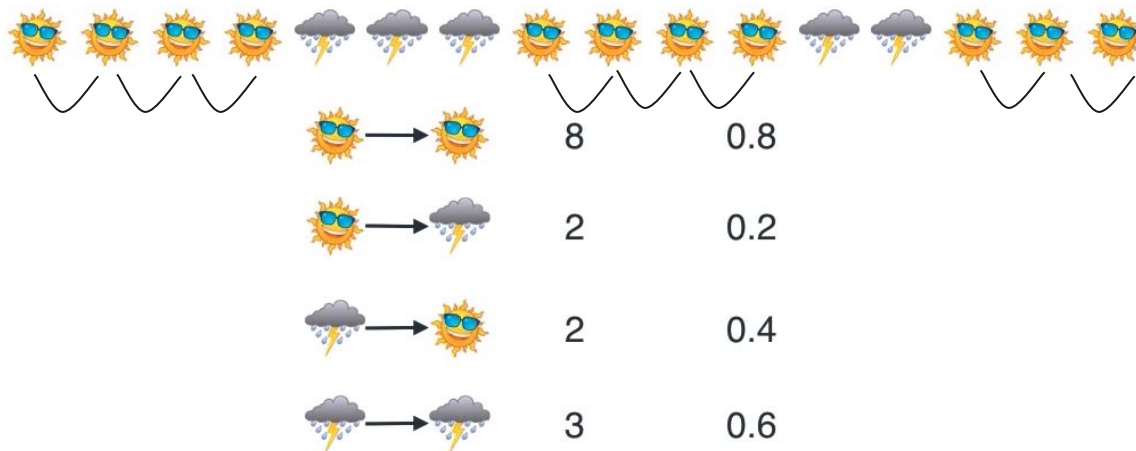
**Emission Prob:** Obs are emitted from Hidden States

**Logic:**

**Predict** sequence of states not directly observable,  
given another sequence of states that are observable  
and hidden states have some **dependence** on the  
observable states

# How are we getting the probability?

From past data observations:



# POS tagging and HMM - Related how?

**Observable States:** words in a sentence

**Hidden States:** POS Tags

**Estimating POS tags:** using HMM

Matrix A: Matrix contains the tag **transition** probabilities  
 $P(t_i | t_{i-1})$

Set of possible Tags: Calculate  $A[\text{Verb}][\text{Noun}]$ :

$$P(\text{Noun} | \text{Verb}) = \text{Count}(\text{Noun \& Verb}) / \text{Count}(\text{Verb})$$

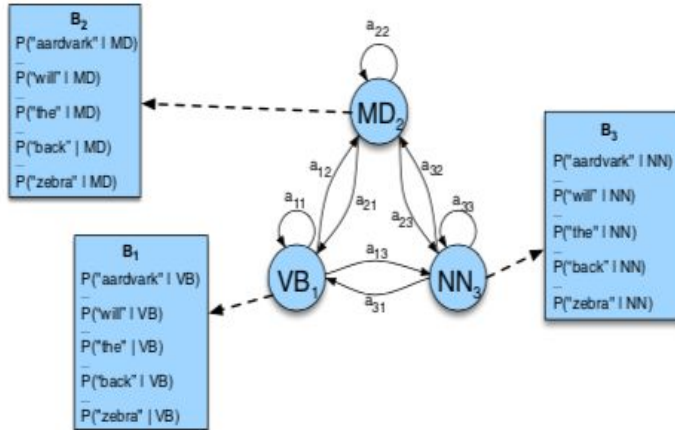
Matrix B = **emission** probabilities,  $P(w_i | t_i)$

Sequence of observation (words in the sentence):

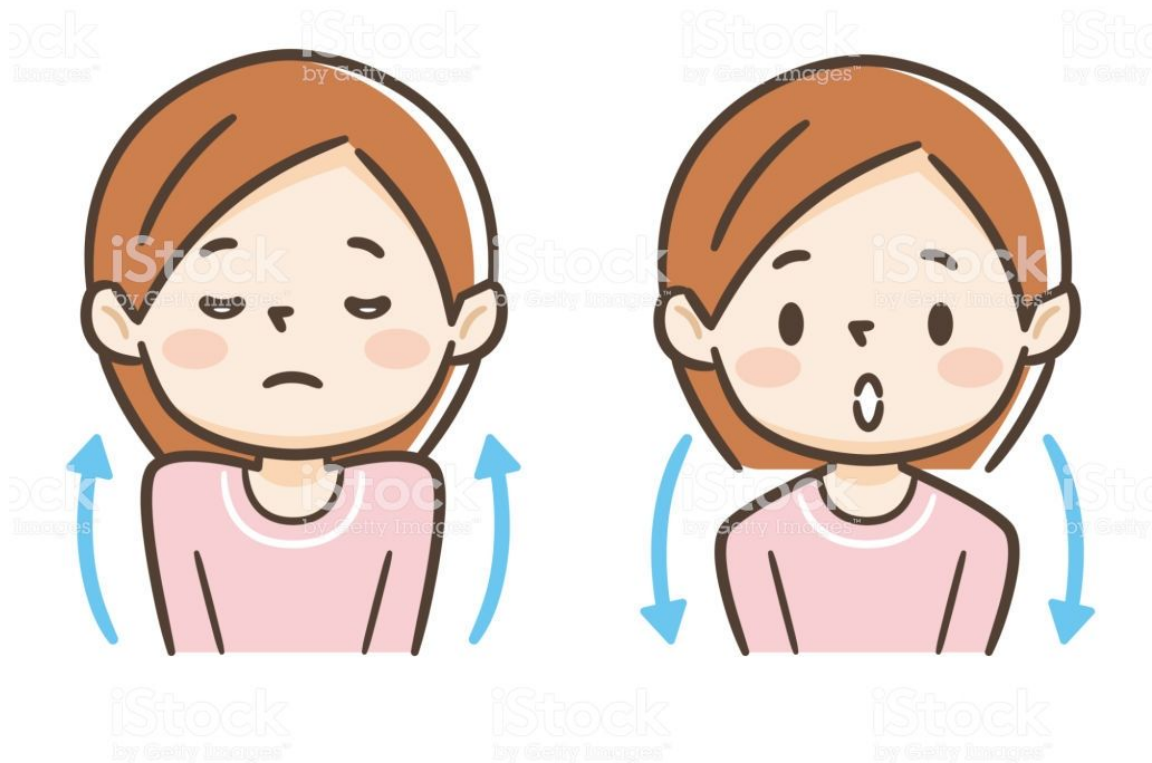
Given a tag (Verb), it's associated with a given word (Running)

The emission probability  $B[\text{Verb}][\text{Running}]$ :

$$P(\text{Running} | \text{Verb}) = \text{Count}(\text{Running \& Verb}) / \text{Count}(\text{Verb})$$



Let's take few deep breaths now!



# Techniques to Understand Text

## Shallow Parsing or Chunking

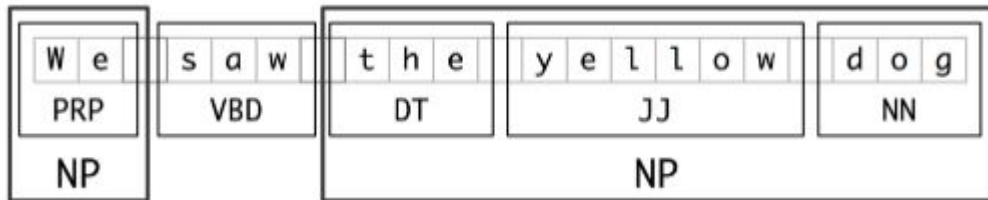
Interpreting Language is Hard!

I saw a girl with a telescope





# What is Shallow Parsing?



- Shallow parsing or chunking is a process dividing a text into syntactically related group
- Divide the whole text into non-overlapping contiguous subsets of tokens
- Segments and labels multi-token sequence
- Crucial for information extraction from text to create sub-components such as Locations, Person Names

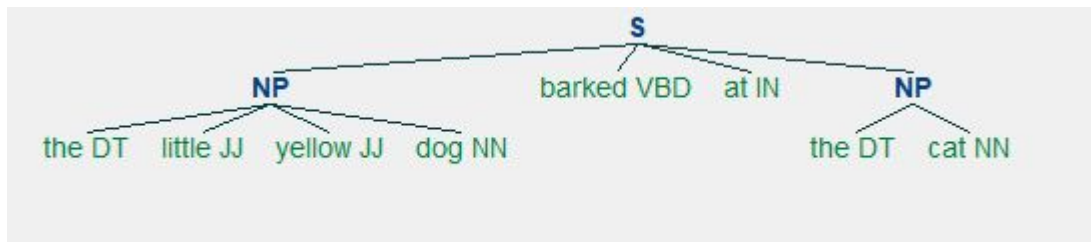
# POS-tagger vs Chunking?

E.g.: Sentence = "the little yellow dog barked at the cat"

POS tagged:

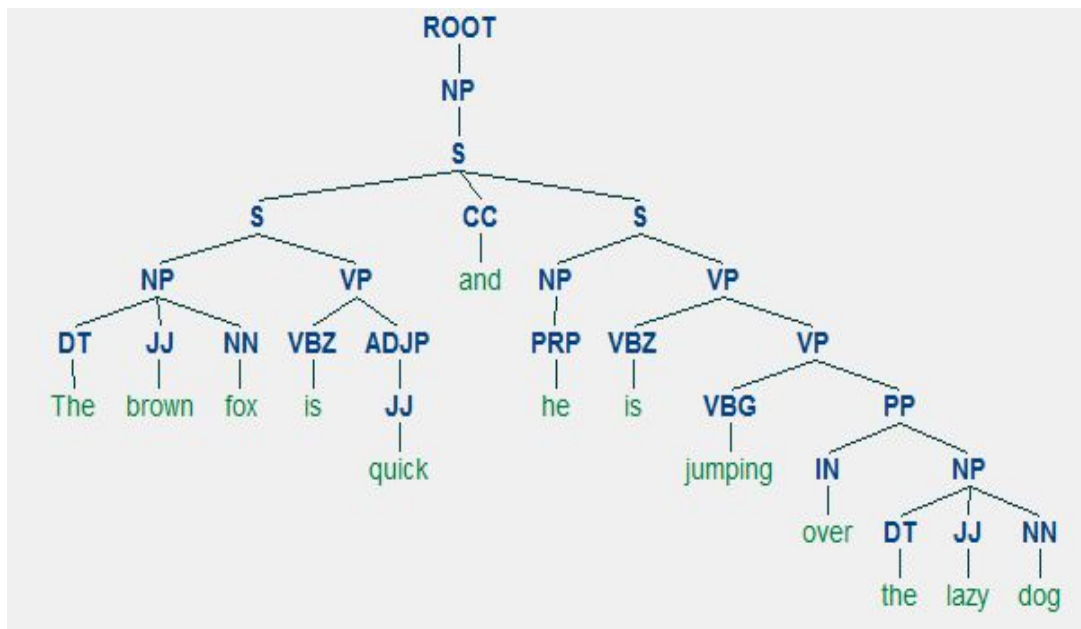
[('the', 'DT'), ('little', 'JJ'),  
('yellow', 'JJ'), ('dog', 'NN'),  
('barked', 'VBD'), ('at', 'IN'),  
('the', 'DT'), ('cat', 'NN')]

Chunked Tree:



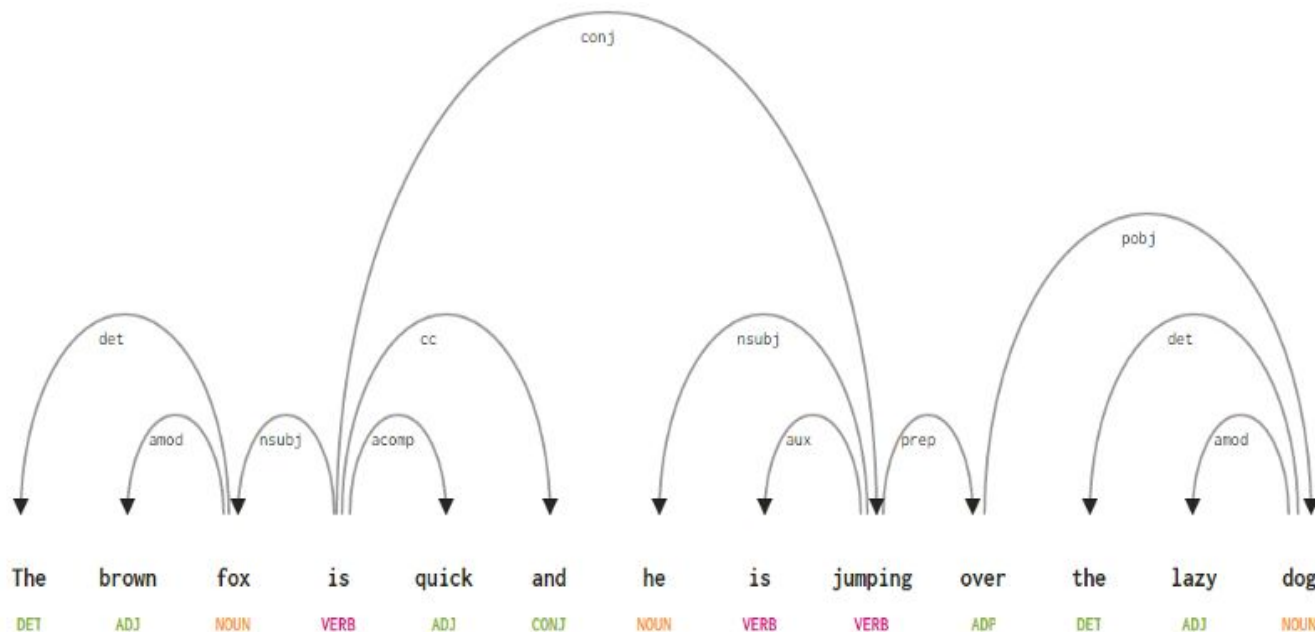
# Other methods of Parsing

## Constituency Parsing

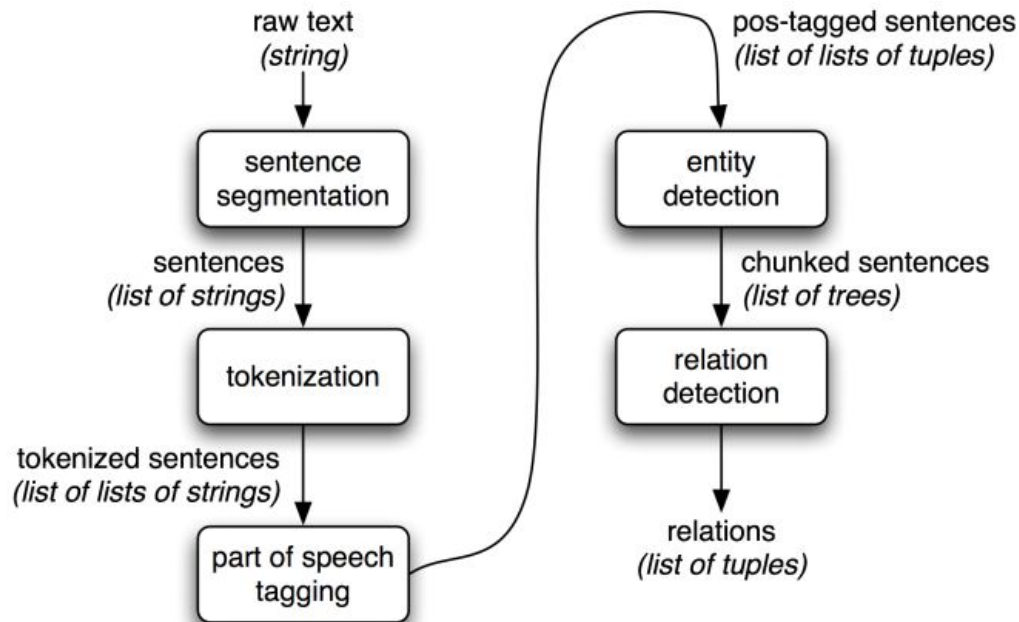


# Other methods of Parsing

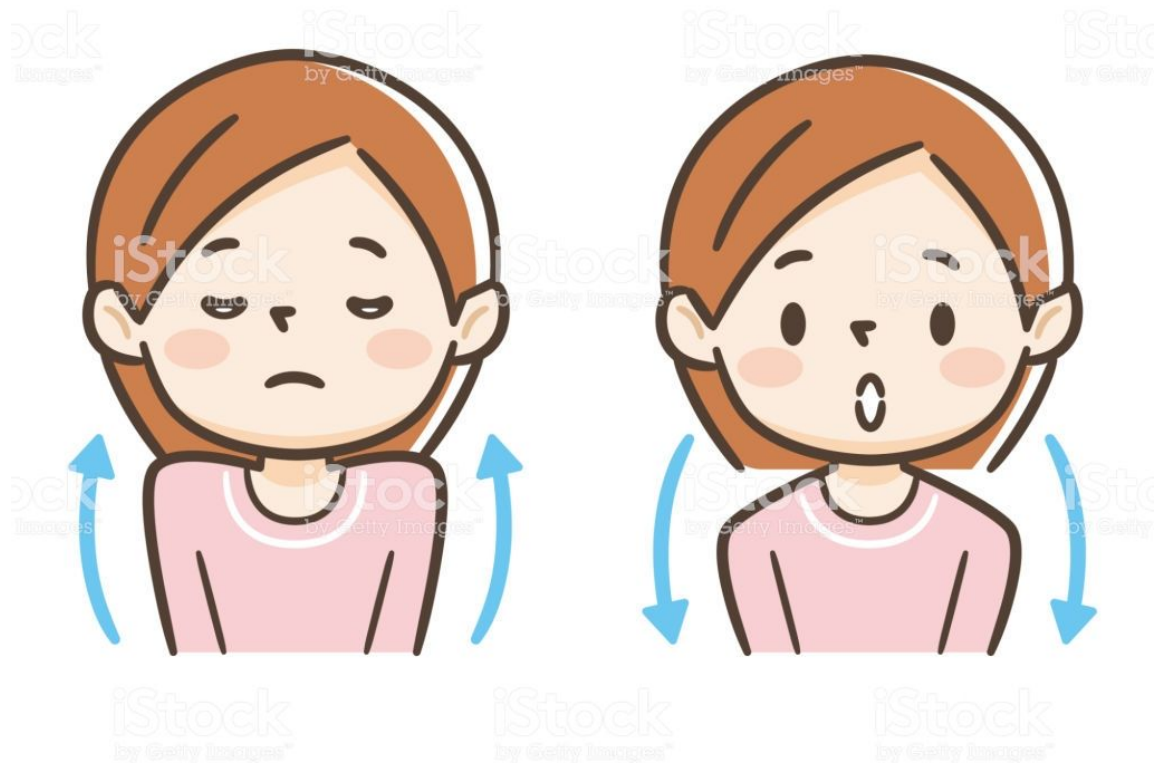
## Dependency Parsing



# Understanding Entity Parsing



Let's take few deep breaths now!



# Techniques to Understand Text

## Named Entity Recognition (NER)

Automatically find  
names of people, places,  
and organizations in text  
across many languages.




# What is NER?

F.B.I. Agent Peter Strzok PERSON, Who Criticized Trump PERSON in Texts, Is Fired GPE - The New York Times ORG SectionsSEARCHSkip to contentSkip to site indexPoliticsSubscribeLog InSubscribeLog InToday's PaperAdvertisementSupported ORG byF.B.I. Agent Peter Strzok PERSON, Who Criticized Trump PERSON in Texts, Is FiredImagePeter Strzok, a top F.B.I. GPE counterintelligence agent who was taken off the special counsel investigation after his disparaging texts about President Trump PERSON were uncovered, was fired. CreditT.J. Kirkpatrick PERSON for The New York TimesBy Adam Goldman ORG and Michael S. SchmidtAug PERSON. 13 CARDINAL, 2018WASHINGTON CARDINAL — Peter Strzok PERSON, the F.B.I. GPE senior counterintelligence agent who disparaged President Trump PERSON in inflammatory text messages and helped oversee the Hillary Clinton PERSON email and Russia GPE investigations, has been fired for violating bureau policies, Mr. Strzok PERSON's lawyer said Monday DATE. Mr. Trump and his allies seized on the texts — exchanged during the 2016 DATE campaign with a former F.B.I. GPE lawyer, Lisa Page — in PERSON assailing the Russia GPE investigation as an illegitimate "witch hunt." Mr. Strzok PERSON, who rose over 20 years DATE at the F.B.I. GPE to become one of its most experienced counterintelligence agents, was a key figure in the early months DATE of the inquiry. Along with writing the texts, Mr. Strzok PERSON was accused of sending a highly sensitive search warrant to his personal email account. The F.B.I. GPE had been under immense political pressure by Mr. Trump PERSON to dismiss Mr. Strzok PERSON, who was removed last summer DATE from the staff of the special counsel, Robert S. Mueller III PERSON. The president has repeatedly denounced Mr. Strzok PERSON in posts on Twitter EVENT, and on Monday DATE expressed satisfaction that he had been sacked. Mr. Trump's ORG victory traces back to June DATE, when Mr. Strzok PERSON's conduct was laid out in a wide-ranging inspector general's report on how the F.B.I. GPE handled the investigation of Hillary Clinton's PERSON emails in the run-up to the 2016 DATE election. The report was critical of Mr. Strzok PERSON's conduct in sending the

- Terms that represent specific entities that are more informative and have a unique context
- Represent real-world objects like people, places, organizations, and so on, which are often denoted by proper names
- Used in information extraction to identify and segment the named entities in predefined classes

# Steps of NER

1. Detect a Named Entity
2. Extract the Entity
3. Categorize the Entity



## Entity Extractor

Extract Entities from text using Named Entity Recognition (NER). NER labels sequences of words in a text which are the names of things, such as person and company names. This implementation labels 3 classes: PERSON, ORGANIZATION and LOCATION.

PERSON
LOCATION
ORGANIZATION

### Test with your own text

SpaceX is an aerospace manufacturer and space transport services company headquartered in California. It was founded in 2002 by entrepreneur and investor Elon Musk with the goal of reducing space transportation costs and enabling the colonization of Mars.

Extract Text

[LIST](#)
[JSON](#)

TAG	VALUE
COMPANY	SpaceX
PERSON	Elon Musk
LOCATION	Mars
LOCATION	California

# Methods for NER extraction

## Lexicon approach

1. relies on a knowledge base called ontology
2. contains all terms related to a particular topic, grouped in different categories
3. system looks for matches with named entities

### Cons:

1. doesn't work to extract new words not in lexicon

### Example:

lexicon of cities, states, and countries to recognize locations in data.

## Rule-based systems

1. series of grammatical rules hand-crafted by computational linguists
2. can get results of high precision but low recall

### Cons:

1. Defining rules takes time
2. Domain specific

### Example:

Build a model to extract "legal terms", you need to manually tag tokens of legal methods, cases, process

## Machine learning-based systems

1. build an entity extractor
2. feed the model with a large volume of annotated training data

### 1. Cons:

Need tagged and clean training data

### Example:

Build a model to extract "legal terms", you need to manually tag tokens of legal methods, cases, process

# Industry Use-Case of NER?



## Problem

Reduce time for processing customer queries and tickets

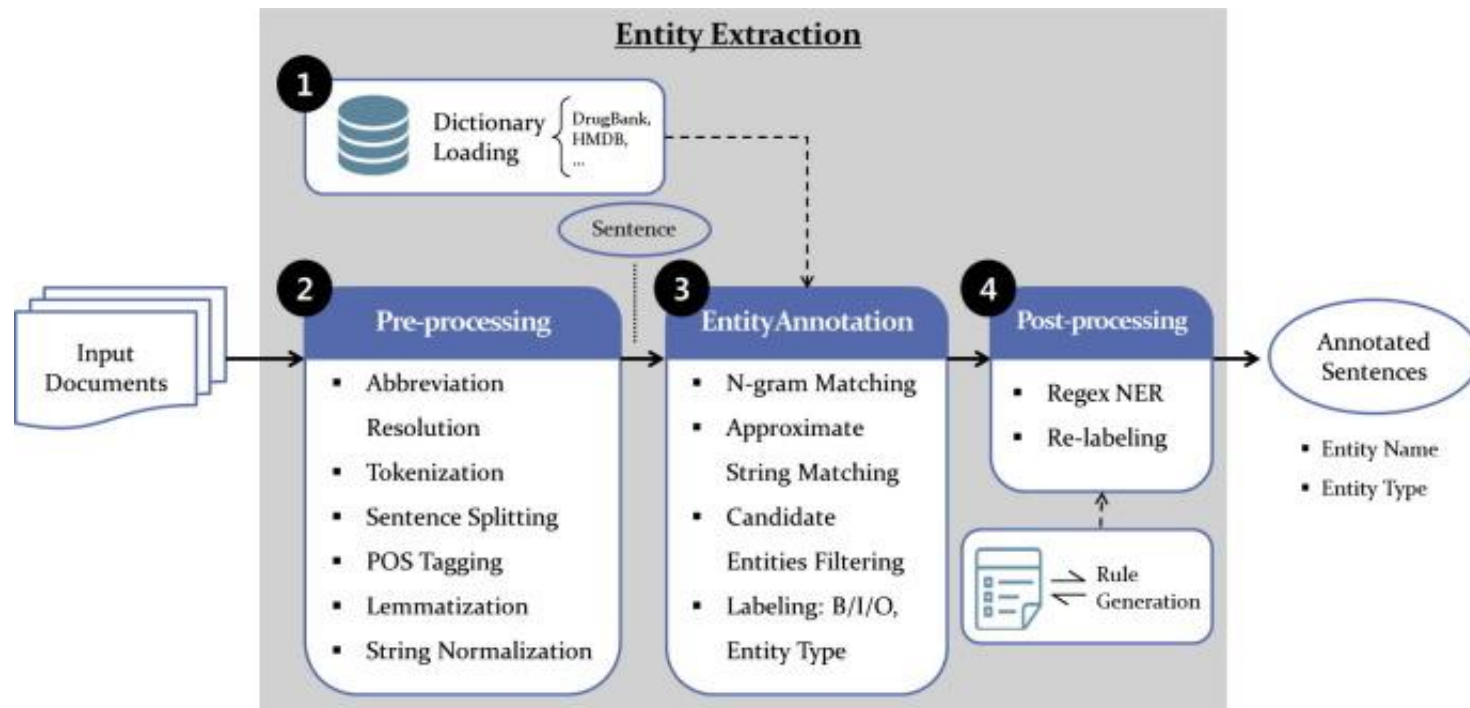
## Goal

Increase customer satisfaction by reducing queue time and solving problems faster

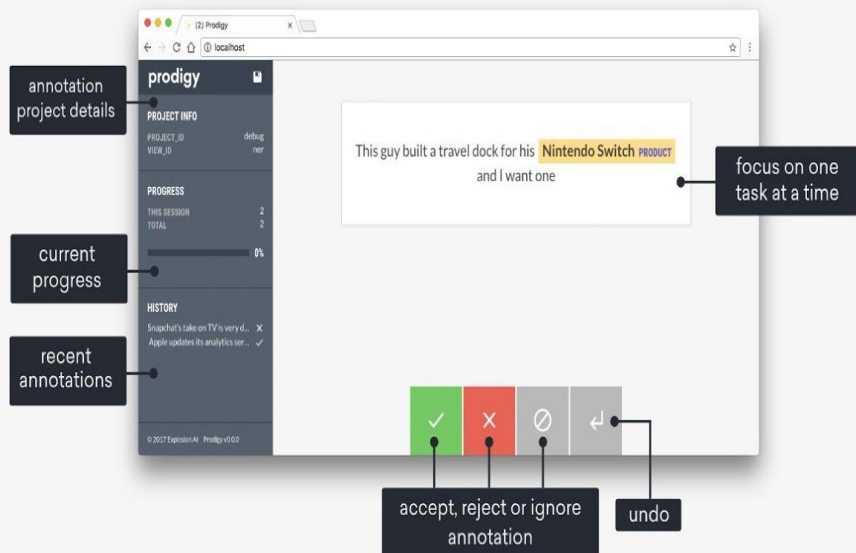
## Constraints

Limited number of agents to respond to tickets, minimize burden

# Training a NER model



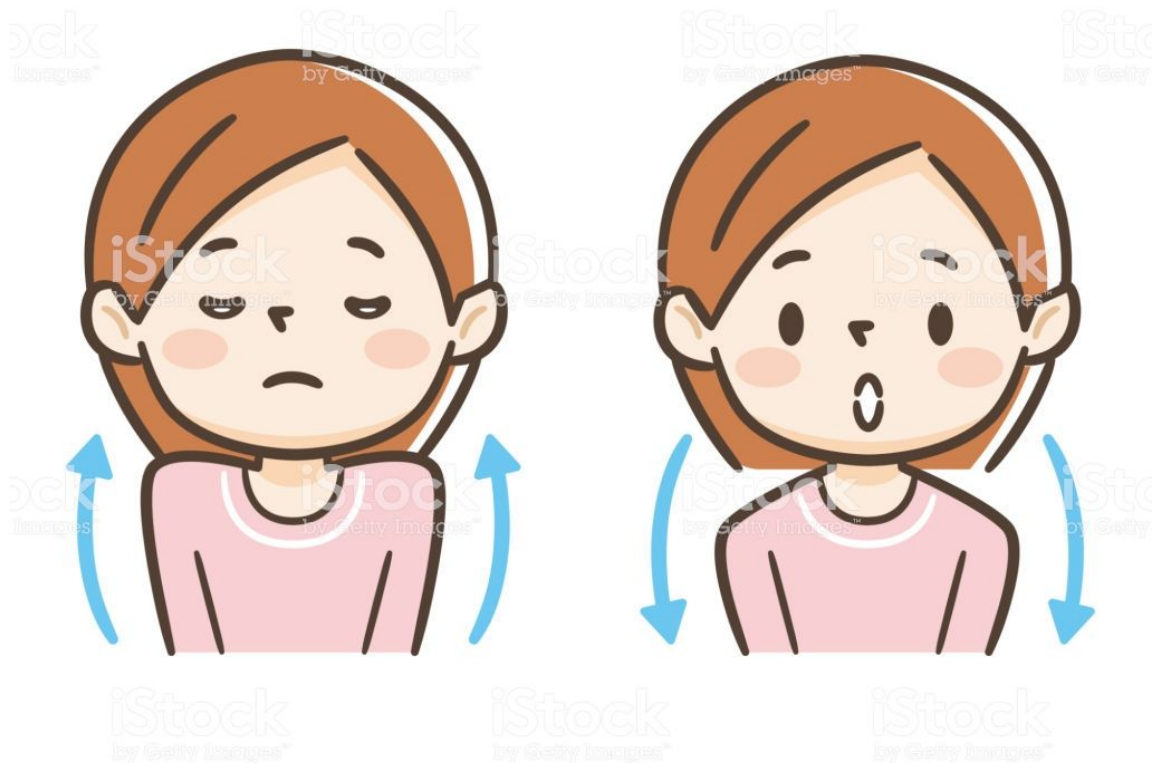
# Annotation for NER



```
$ prodigy ner.manual ner_news_headlines blank:en
./news_headlines.jsonl --label PERSON,ORG,PRODUCT,LOCATION

✨ Starting the web server at http://localhost:8080 ...
Open the app in your browser and start annotating!
```

Let's take few deep breaths now!





# Techniques to Understand Text

N-Grams



# What are n-grams?

This is Big Data AI Book

*Uni-Gram*

This	Is	Big	Data	AI	Book
------	----	-----	------	----	------

*Bi-Gram*

This is	Is Big	Big Data	Data AI	AI Book
---------	--------	----------	---------	---------

*Tri-Gram*

This is Big	Is Big Data	Big Data AI	Data AI Book
-------------	-------------	-------------	--------------

- N-gram model is a type of Language Model (LM), which is about finding the probability distribution over word sequences
- N-gram means a sequence of N words
- N-grams cut out the noise from the data in your analyses
- Identify themes quickly
- NLP applications including speech recognition, machine translation and predictive text input

# Mathematical concept of n-grams

Probability of a word  $w$ , given some history,  $h = P(w|h)$

Eq:  $P(\text{the}|\text{today the sky is so clear that})$

Here,

$w = \text{the}$

$h = \text{today the sky is so clear that}$

**Approach 1:** relative frequency count

Step1: Take a large corpus

Step2: count the number of times **\*\*today the sky is so clear that\*\*** appears

Step3: count the number of times it is followed by **\*\*the\*\***

Eq:

$$P(\text{the}|\text{today the sky is so clear that}) = \frac{C(\text{today the sky is so clear that the})}{C(\text{today the sky is so clear that})}$$

Basically, we need to answer:

Out of the times you saw the history  $h$ , how many times did the word  $w$  follow it

# Mathematical concept of n-grams



## Cons of Approach 1:

1. If we have a large corpus, Approach 1 needs to go over entire corpus
2. Not feasible for scaling as well as time performance
3. To decompose the probability function into smaller chunks leads to the usage of **chain rule**

Instead of computing probability using the entire corpus, chain rule using N-grams would approximate it by just a **few historical words**

# Mathematical concept of n-grams

Probability of a word  $w$ , given some history,  $h = P(w|h)$       **Approach 2: Bigram Model**

Eq:  $P(\text{the}|\text{today the sky is so clear that})$

Here,

$w_n$  = the

$h$  = today the sky is so clear that

$w_{n-1}$  = that

Process: Approximates the probability of a word given all the previous words by using only the conditional probability of one preceding word

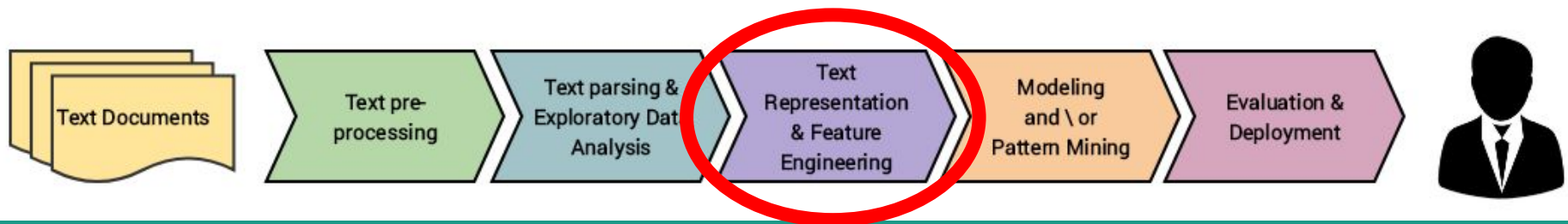
$$P(w_n | w_1^{n-1}) \approx P(w_n | w_{n-1})$$

Assumption that the **probability of a word** depends only **on the previous word** = Markov assumption

Def: Markov models are the class of probabilistic models that assume that we can **predict the probability of some future unit** without looking too far in the past.

# 4

# NLP Workflow



5

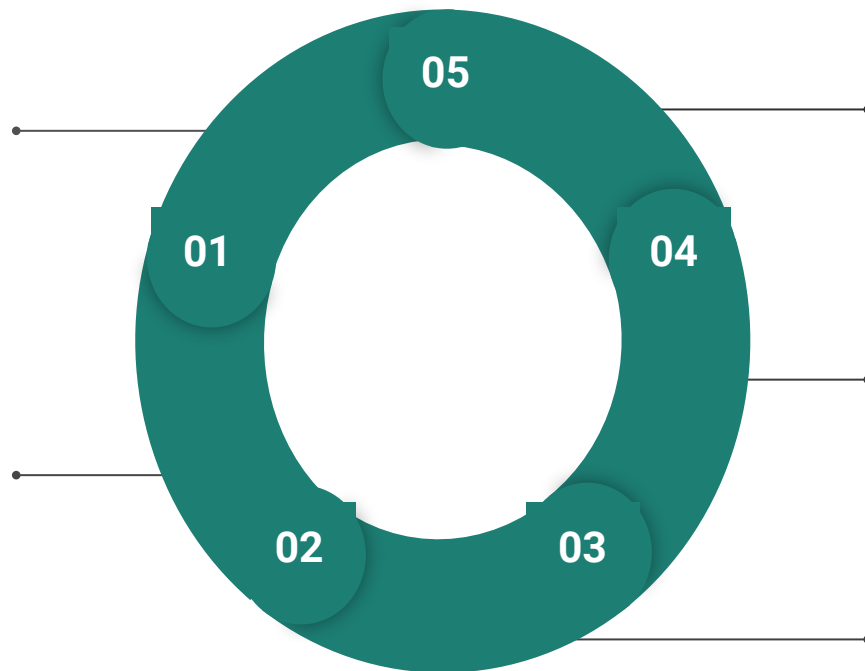
# Theory Wrap-up & Next Steps

# Recap



Syntax vs Semantics

Part of Speech (POS) Tagging



Shallow Parsing or Chunking

Named Entity Recognition

N-Grams



# 6

# Google Colab Project

<https://bit.ly/introtonlp-week2-notebook>

# Homework

## #1

### Additional Resources

- [\(Youtube\) Hidden Markov Model by Luis Serrano](#)
- [\(FreeCodeCamp\) Part of Speech Tag and Hidden Markov Model](#)
- [\(AnalyticsVidhya\) Dependency](#)
- [\(TowardsDataScience\) N-grams](#)

See you  
next week!

## Questions?

Join us on [Slack](#) and post your questions  
to the #help-me channel