

Reza Shokrzad

March 2025

NLP Workshop

Session 1 - NLP Basics & Text mining

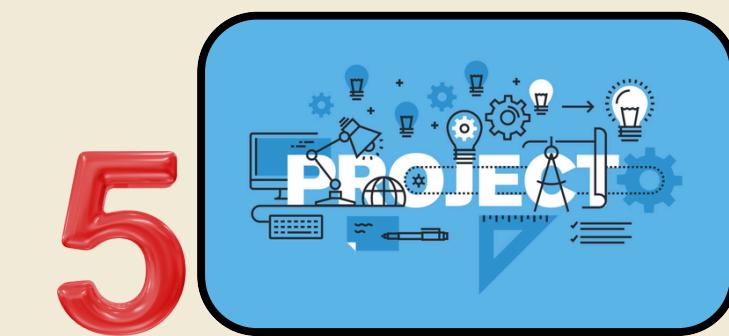
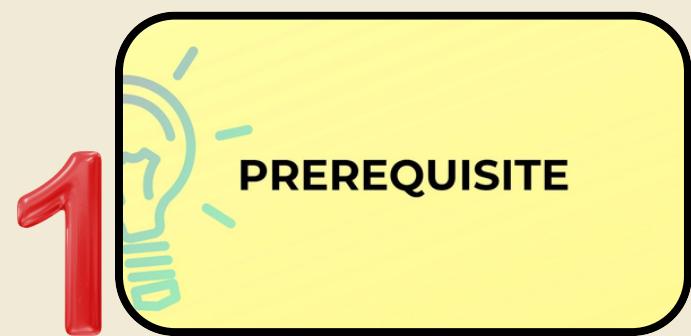
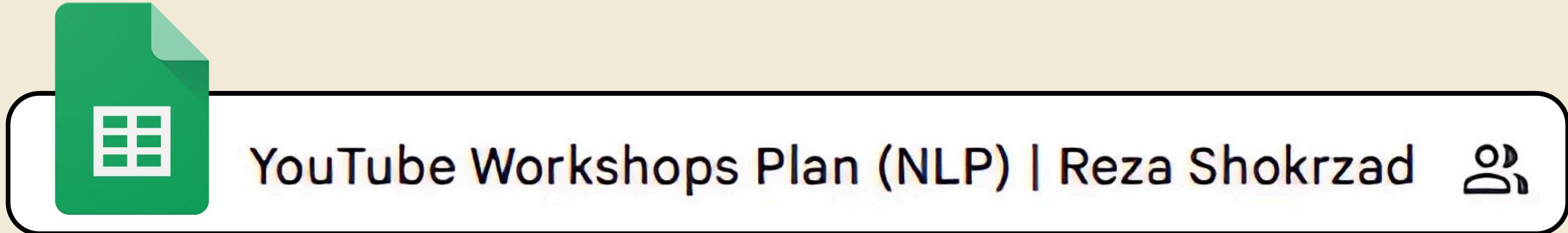


Content

- Course Structure
- Where are we?
- What is NLP?
- NLP Challenges
- Text Data
- Preprocessing
- Regex
- Python Packages



Course Structure



TAs



Mohammad Javad Shamloo · 2nd
Passionate to Learn Everything
Tehran, Tehran Province, Iran · [Contact info](#)

MAPNA Group
Sharif University of Technology



Roqayeh Mohajeri · 1st
Aspiring Data Science | Machine Learning Trainee | Biotech MSc
Gonbad-e Qasr, Semnan Province, Iran · [Contact info](#)

BNUT - Babol Noshirvani University of Technology



Ayda Abdi · 2nd
--
Tehran, Tehran Province, Iran · [Contact info](#)

Freelancer
Khatam University



Hamidreza KAZEMI · 2nd
Enthusiastic about Data Science and Business Intelligence| IUST student
Tehran, Tehran Province, Iran · [Contact info](#)

Iran University of Science and Technology



Afsaneh Shamsaddini · 1st
Artificial Intelligence Engineer /Machine Learning/Python Developer
Kerman, Kerman Province, Iran · [Contact info](#)

#OPENTOWORK
Shahid Bahonar University

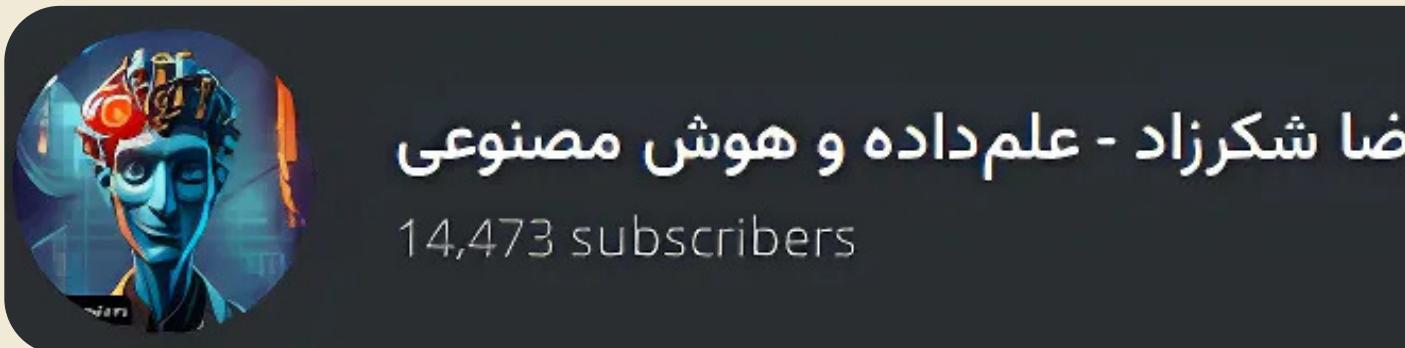


Mohammadmahdi Rouholamini · 2nd
--
Tehran, Tehran Province, Iran · [Contact info](#)

Shahid Beheshti University



Communities



ضا شکرزاد - علم داده و هوش مصنوعی

14,473 subscribers



Reza Shokrzad - Data Science & A

@RezaShokrzad • 3.12K subscribers • 106 videos

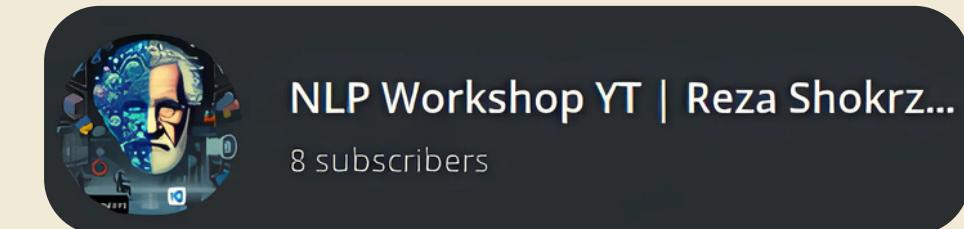
هدف این کتاب آموزش میدهدای (برنامه علم داده است) ...

cafetadris.com/datasci

Manage video

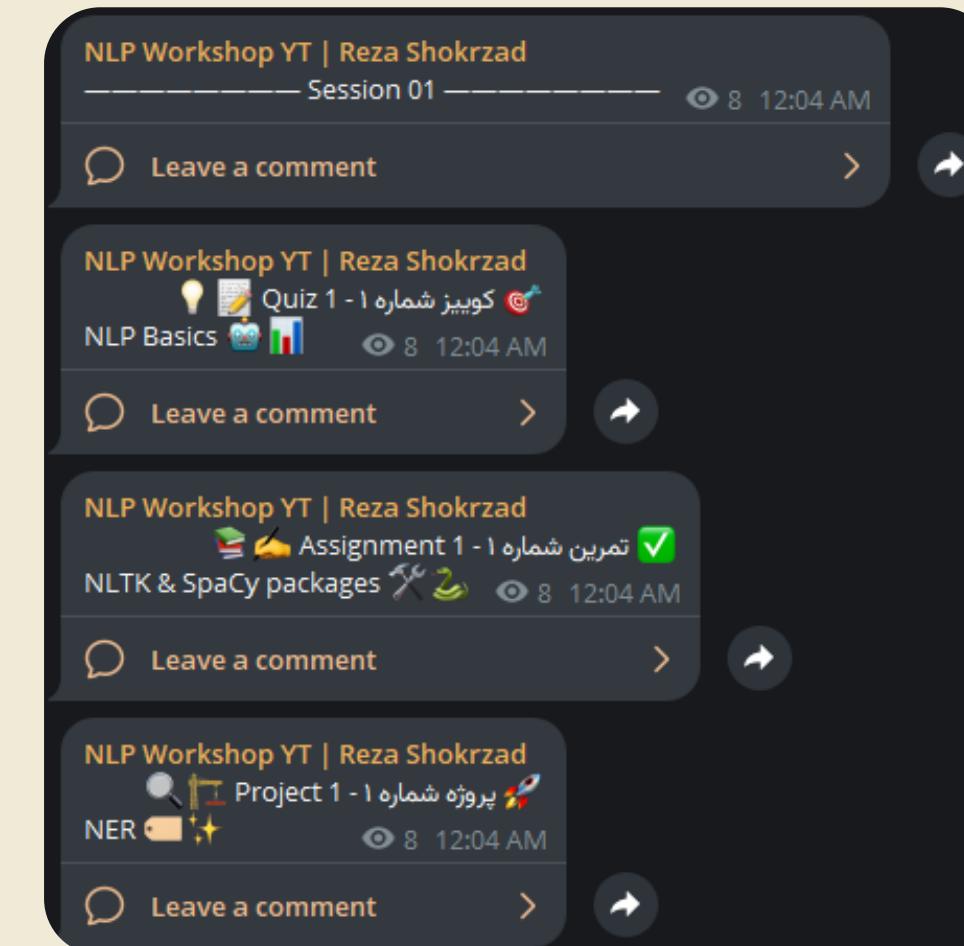


@RezaShokrzad



NLP Workshop YT | Reza Shokrz...

8 subscribers



Where were we?

Can LLMs Generate Novel Research Ideas?
A Large-Scale Human Study with 100+ NLP Researchers

Chenglei Si, Diyi Yang, Tatsunori Hashimoto
Stanford University

{clsi, diyiy, thashim}@stanford.edu

Sep 2024

<https://arxiv.org/abs/2409.04109>



Where are we?

The screenshot shows the bioRxiv homepage with a search bar at the top. Below it, the CSHL logo and the bioRxiv logo with the tagline "THE PREPRINT SERVER FOR BIOLOGY" are displayed. A navigation bar on the right includes "HOME | SUBMIT". The main content area features a "New Results" section and a "Follow this preprint" button with a bell icon. A highlighted preprint is shown with the title "AI mirrors experimental science to uncover a novel mechanism of gene transfer crucial to bacterial evolution". The authors listed are José R Penadés, Juraj Gottweis, Lingchen He, Jonasz B Patkowski, Alexander Shurick, Wei-Hung Weng, Tao Tu, Anil Palepu, Artiom Myaskovsky, Annalisa Pawlosky, Vivek Natarajan, Alan Karthikesalingam, Tiago R D Costa. The DOI is provided as <https://doi.org/10.1101/2025.02.19.639094>.

February 2025

Paper: <https://www.biorxiv.org/content/10.1101/2025.02.19.639094v1.article-info>

News: <https://longportapp.com/en/news/229180857>



Natural Language Processing (NLP)

NLP enables computers to interpret, generate, and interact with human language through computational techniques.



NLP Challenges

- Ambiguity & Context Sensitivity
- Cultural & Linguistic Nuances
- Handling Massive Datasets
- Continuous Innovation for Greater Accuracy & Relevance

Words like "bank" (river vs. financial institution)

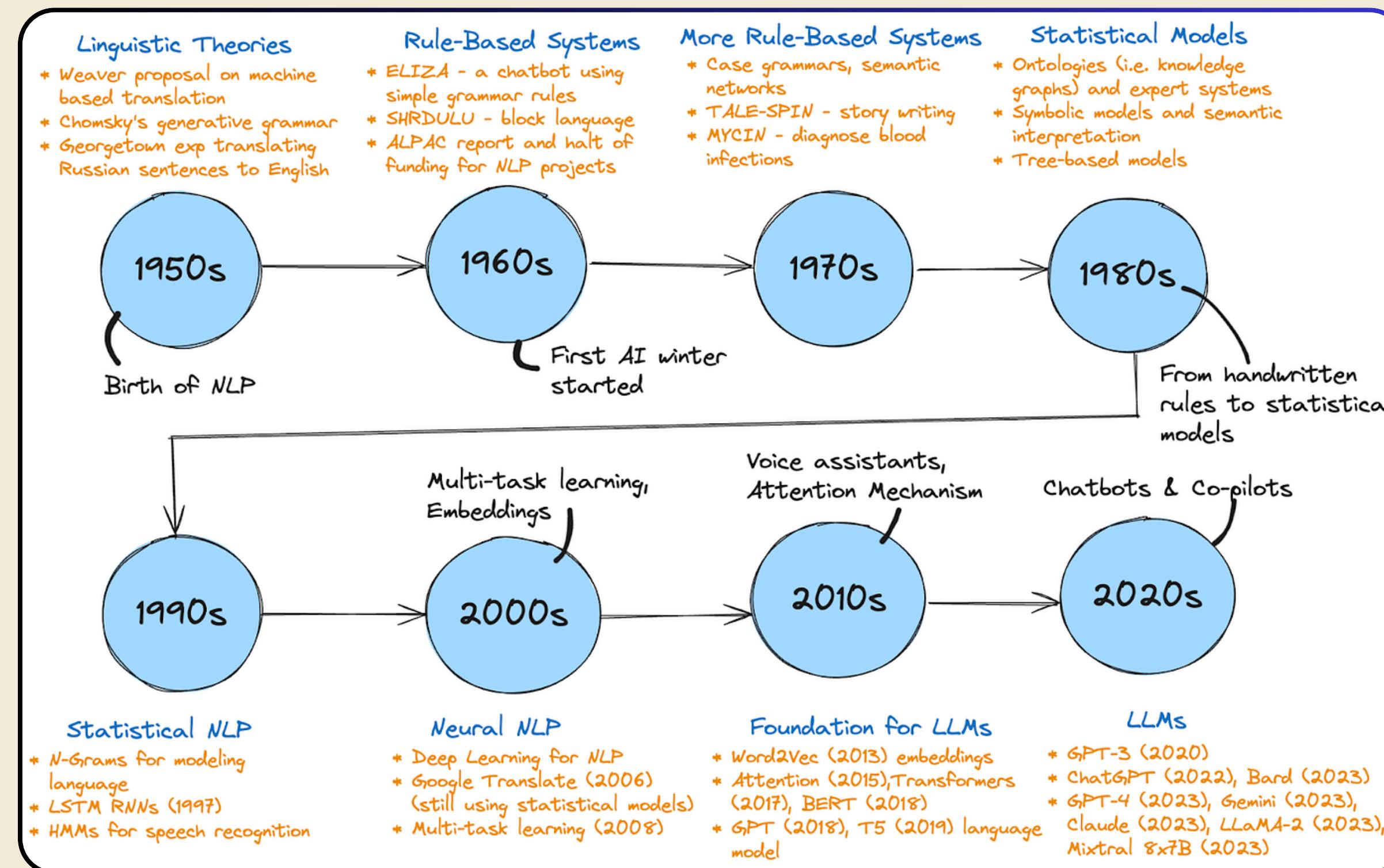
Flying planes can be dangerous

Idioms such as "break a leg"

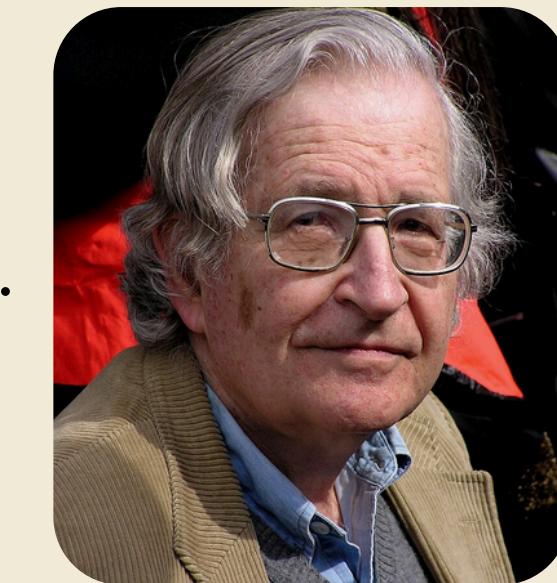
She's on fire



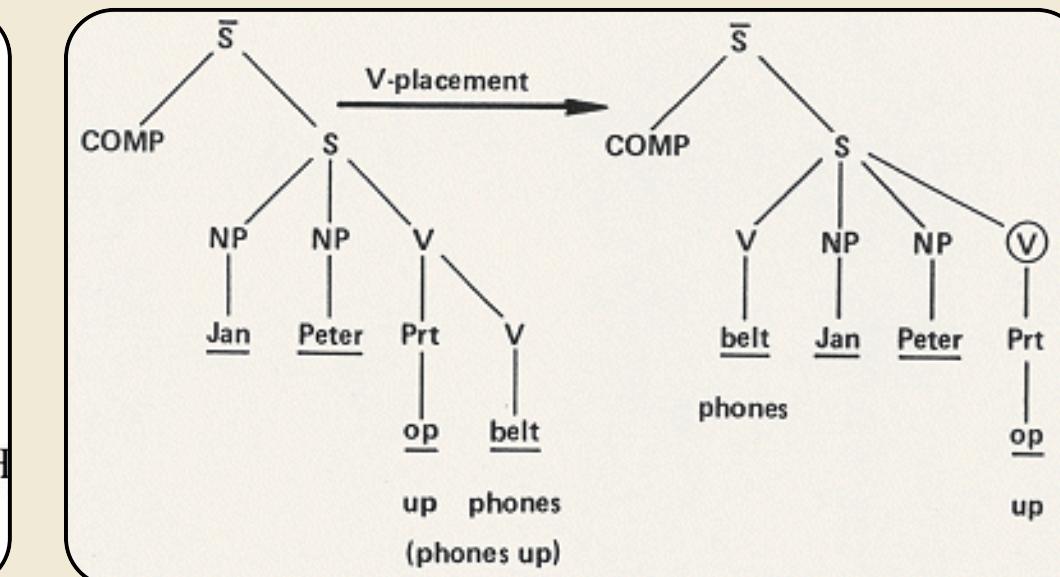
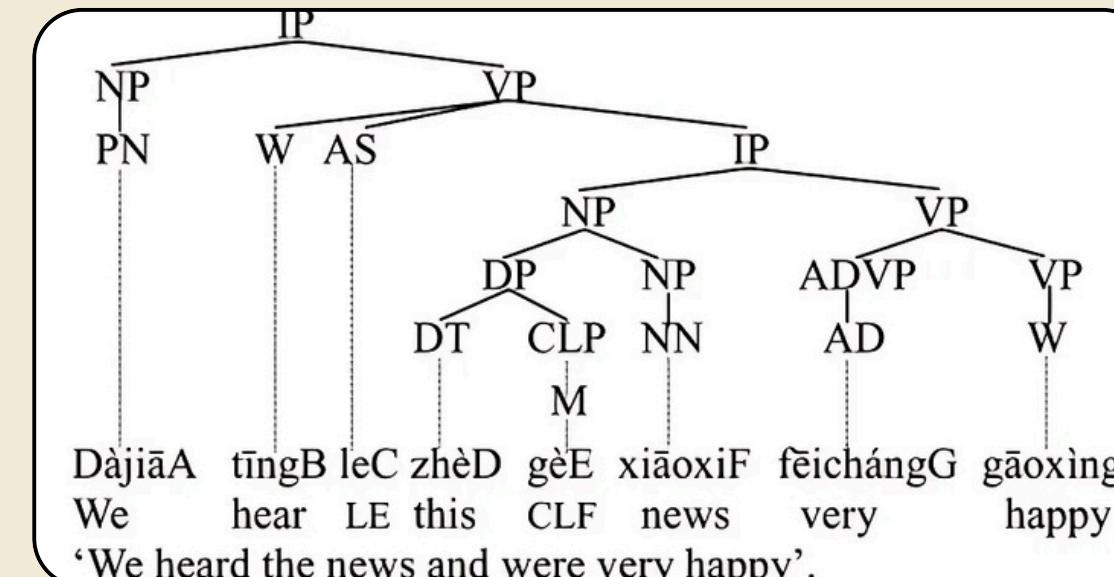
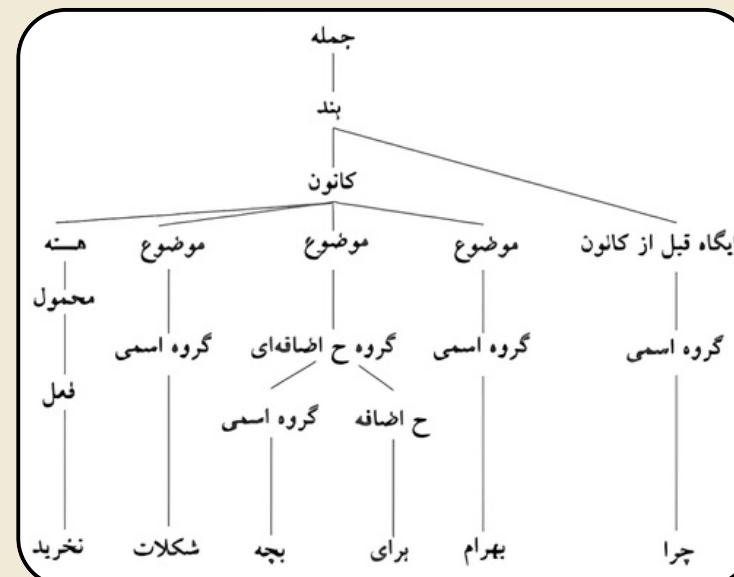
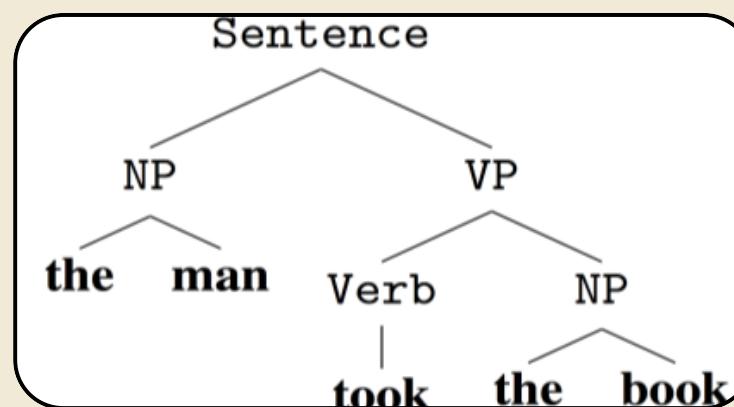
NLP Timeline



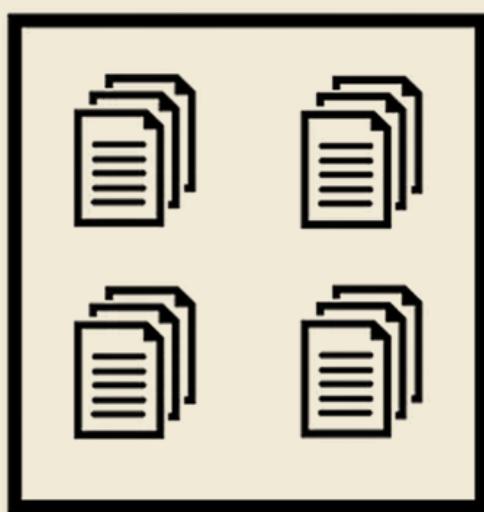
Noam Chomsky Contribution



- Pioneered generative grammar, introducing tree-based syntactic representation
 - Developed phrase-structure rules and transformational approaches.
 - Revolutionized cross-linguistic parsing and modern syntactic theory.



Text Dataset



Corpora



Corpus



Document



Token

- **Brown Corpus:** Foundational American English dataset.
- **Penn Treebank:** Syntactically annotated text collection.
- **British National Corpus:** Comprehensive modern British English.
- **Wikipedia Corpus:** Vast, diverse contemporary text.
- **Google Books Ngram Corpus:** Historical language trends via n-grams.



Text

An unstructured data ...

```
dataset['Reviews'].iloc[2]
```

"Hello e-v-e-r-y-o-n-e!!!@@@!!!!!! 😊 @DONT BUY THIS PHONE at all-first of al
l that say the phone in new , i took it to the lab after 6 month the phone is d
ead dead ,you can save itthay open the phone in the lab and say!!!!the phone is
renew ,and its cheepe commponents.I payed 400\$ for only 6 month ,now i need to
buy new one this LG G4 is dead .not nice 'people say to me dont buy from http://www.amazon.com/gp/aw/d/B01GYUDMFY/ref=ya_aw_od_pi?ie=UTF8&psc=1 at al!!!"

... needs preprocessing



Preprocessing

1. Lowercasing
2. Punctuation & Special Character Removal
3. Tokenization
4. Stop-Word Removal
5. Stemming and Lemmatization

https://en.wikipedia.org/wiki/Text_processing



1. Lowercasing

- Convert text to lowercase.
- Unify case for consistent tokens.
- Reduce vocabulary redundancy.

Pizza

pizza

PIZZA

PIZZA

piZZA



pizza



2. Punctuation & Special Character Removal

- Strip punctuation from text.
- Remove special characters systematically.
- Eliminate non-alphanumeric symbols.
- Reduce noise by cleaning marks.

Wow!!! Visit
<https://example.com> now
- save 30%
 #Deal
@limitedTime

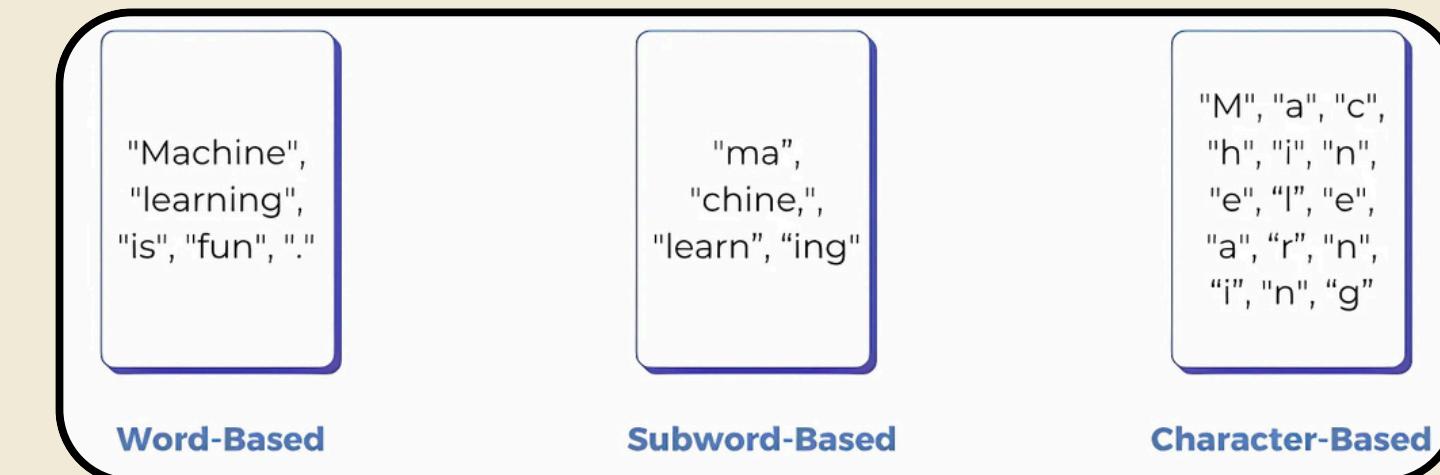


wow visit examplecom
now save 30 deal
limitedtime



3. Tokenization

- Tokenization converts text into manageable units (tokens) for easier processing by NLP models.
- **Types:** Includes word, subword (e.g., Byte-Pair Encoding), and character-level tokenization, each offering different granularity.
- **Impact on Model Performance:** The choice of tokenization affects model size, training time, and accuracy, influencing how models understand text.
- **Handling Special Characters:** Proper tokenization accounts for punctuation, contractions, and special characters to maintain context.
- **Language-Specific Considerations:** Tokenization varies across languages (e.g., word boundaries in English vs. character-based in Chinese).



https://en.wikipedia.org/wiki/Lexical_analysis#Tokenization



4. Stop-Word Removal

- Eliminate common stop words.
- Remove non-informative terms.
- Reduce dataset vocabulary.
- Enhance semantic clarity.

Stopwords

a	it	these
about	its	they
again	itself	this
all	just	those
almost	kg	through
also	km	thus
although	made	to
always	mainly	upon
among	make	use
an	may	used
and	mg	using
another	might	various
any	ml	very
are	mm	was
as	most	we
at	mostly	were

1	
2	و
3	در
4	به
5	از
6	که
7	می
8	این
9	است
10	را
11	با
12	های
13	برای
14	آن
15	یک

https://en.wikipedia.org/wiki/Stop_word

<https://github.com/kharazi/persian-stopwords/tree/master>



5. Stemming and Lemmatization

- Remove affixes (stemming).
- Map words to base forms (lemmatization).
- Reduce vocabulary diversity.
- Enhance text consistency.

Stemming	Lemmatization
adjustable → adjust	was → (to) be
formality → formalit	better → good
formaliti → formal	meeting → meeting
airliner → airlin	

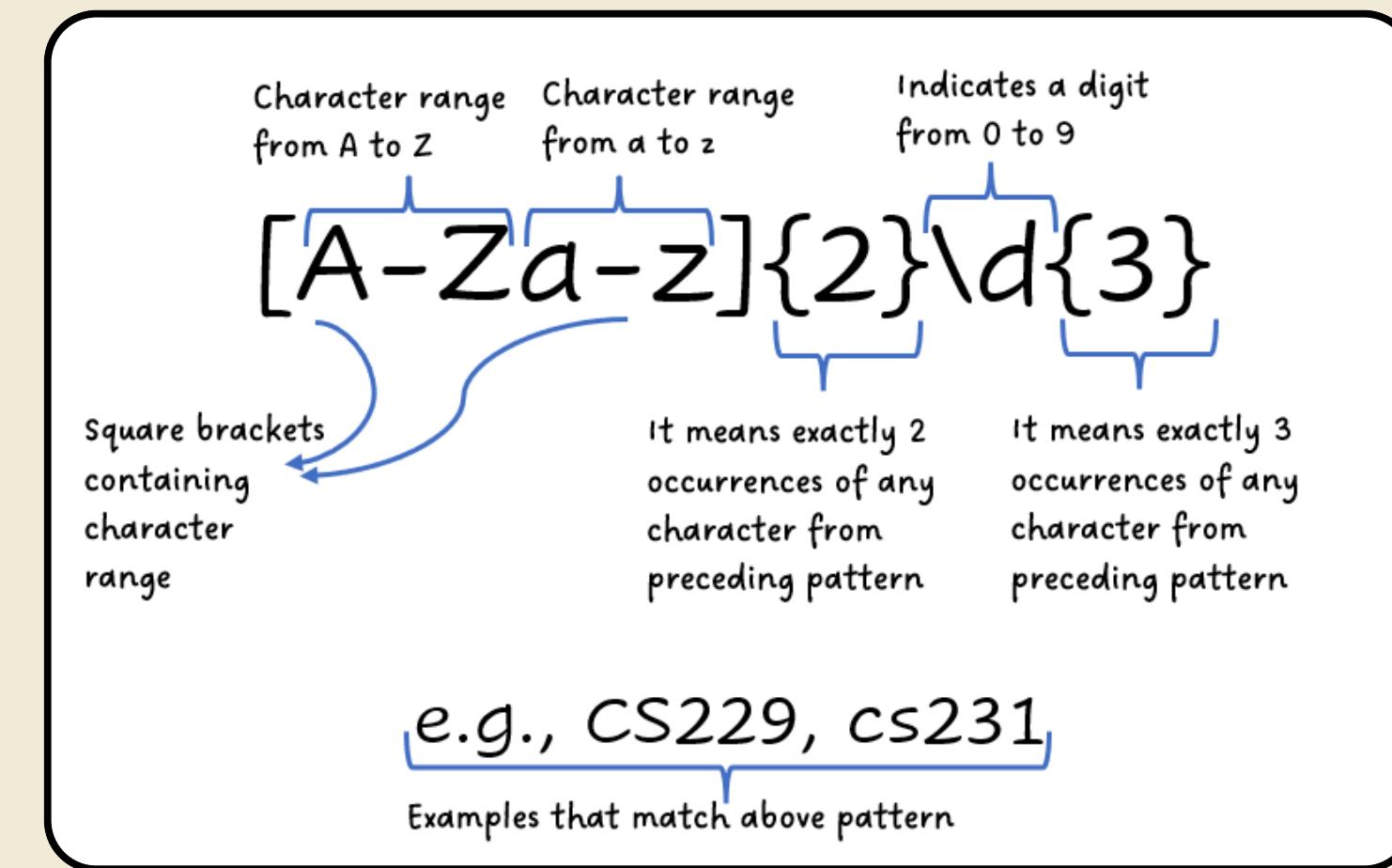
<https://en.wikipedia.org/wiki/Stemming>

<https://en.wikipedia.org/wiki/Lemmatization>

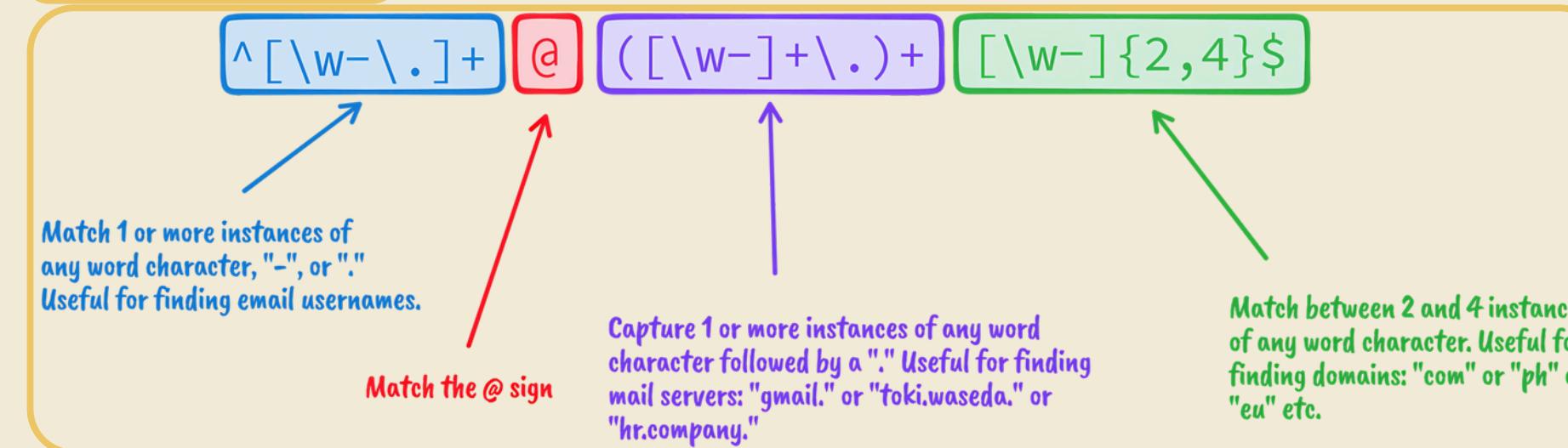


Regular Expression

- Match specific text patterns.
- Extract relevant substrings.
- Replace or remove patterns.
- Filter unwanted text noise.



Email Pattern

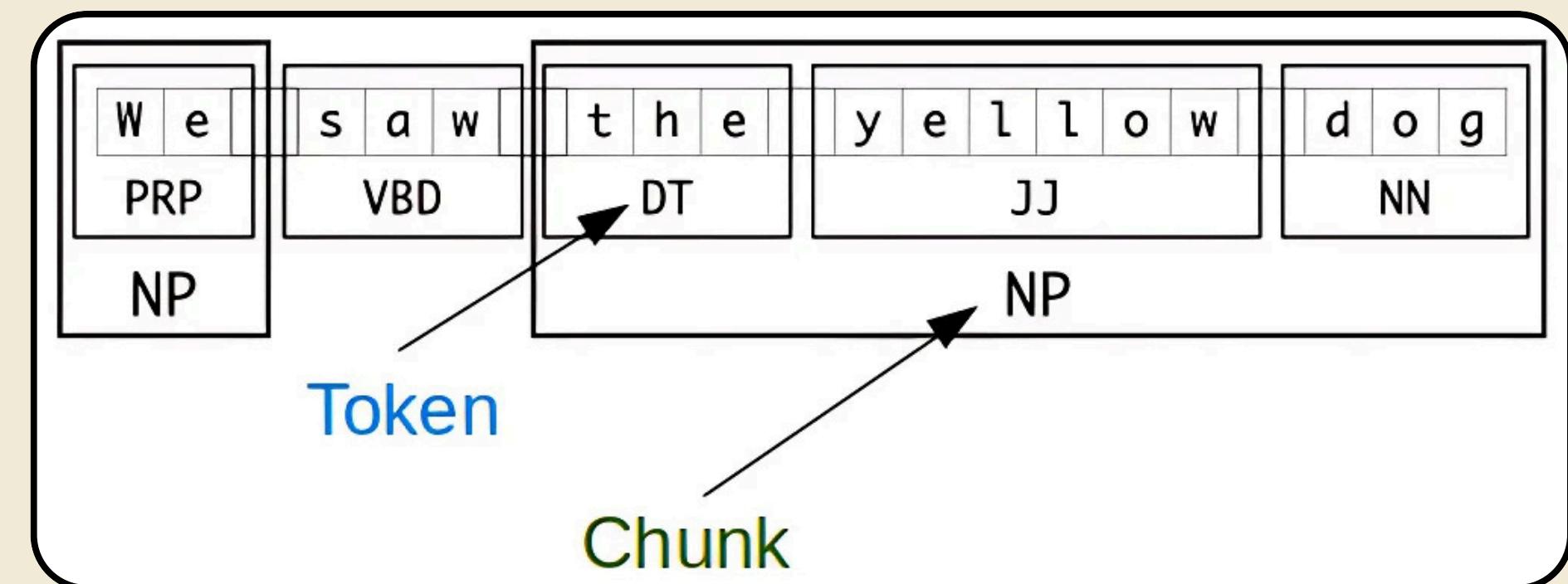


https://en.wikipedia.org/wiki/Regular_expression



Chunking and Parsing

- Group words into meaningful chunks.
- Identify phrase boundaries.
- Analyze grammatical structure.
- Reveal sentence hierarchy.



https://en.wikipedia.org/wiki/Shallow_parsing



POS tagging

- **Noun (N)**- Daniel, London, table, dog, teacher, pen, city, happiness, hope
- **Verb (V)**- go, speak, run, eat, play, live, walk, have, like, are, is
- **Adjective(ADJ)**- big, happy, green, young, fun, crazy, three
- **Adverb(ADV)**- slowly, quietly, very, always, never, too, well, tomorrow
- **Preposition (P)**- at, on, in, from, with, near, between, about, under
- **Conjunction (CON)**- and, or, but, because, so, yet, unless, since, if
- **Pronoun(PRO)**- I, you, we, they, he, she, it, me, us, them, him, her, this
- **Interjection (INT)**- Ouch! Wow! Great! Help! Oh! Hey! Hi!

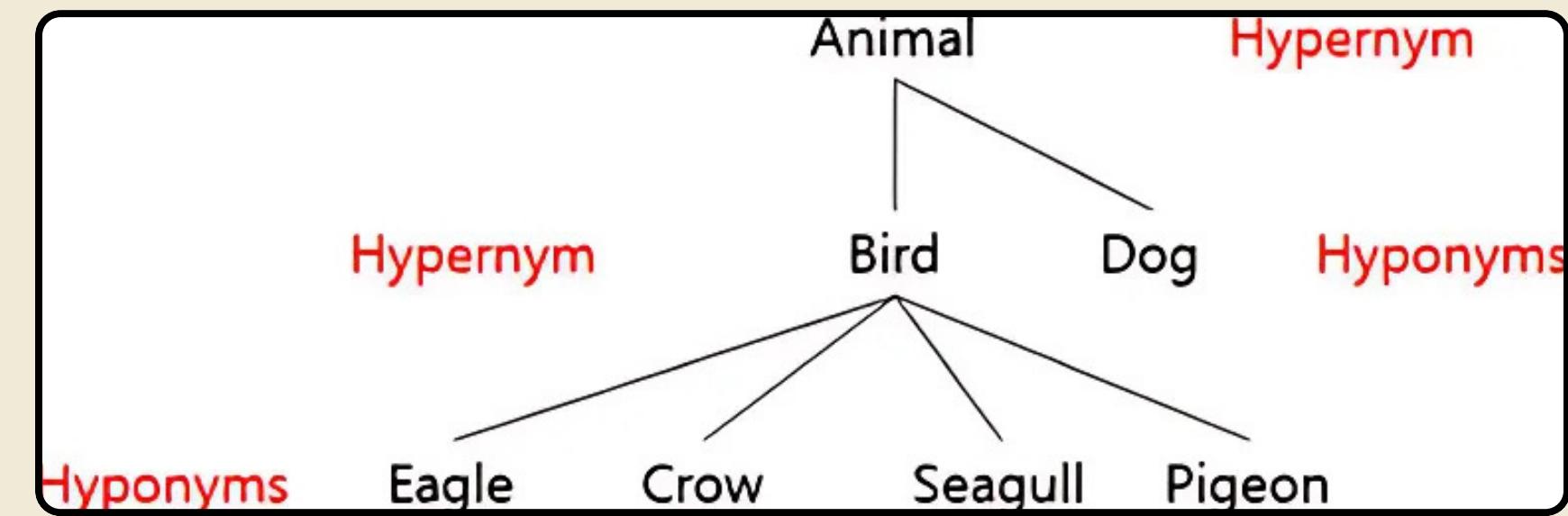
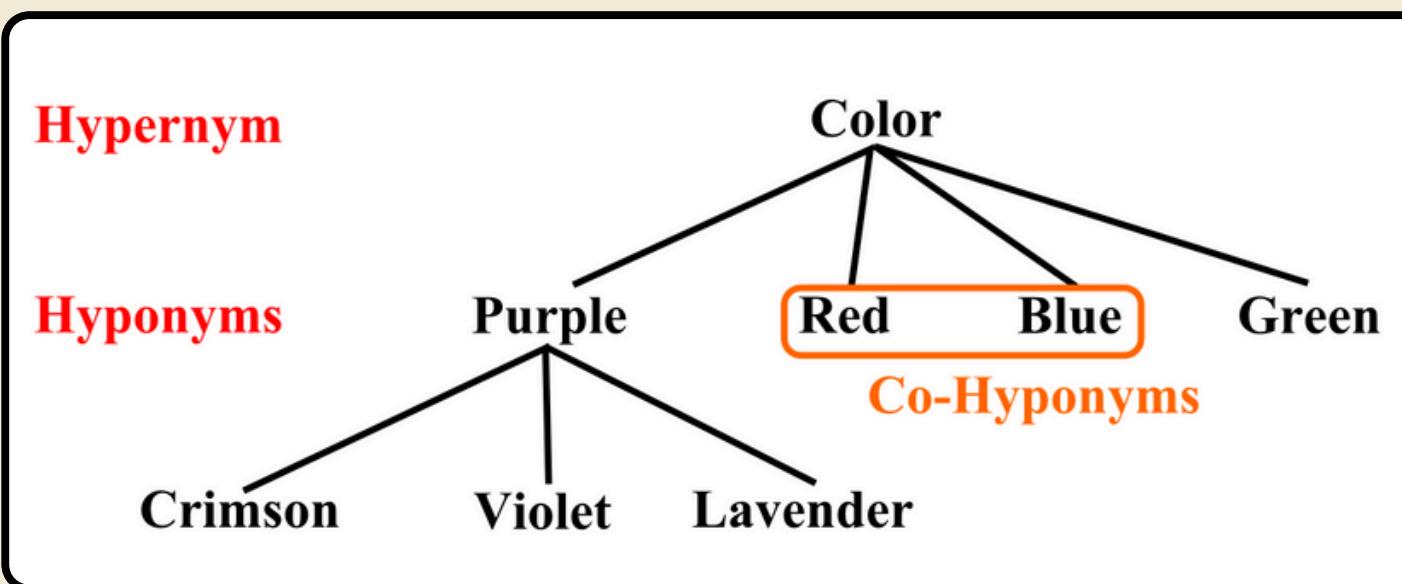
Why POS tagging?

- Clarifies grammatical roles
- Aids in disambiguation
- Enhances downstream NLP tasks



Hypernyms and Hyponyms

- Identify broad category (hypernym).
- Recognize specific instance (hyponym).
- Enhance semantic hierarchy.
- Improve lexical organization.



https://en.wikipedia.org/wiki/Hypernymy_and_hyponymy



NLTK

- Comprehensive NLP toolkit.
- Supports tokenization, tagging, parsing.
- Extensive corpora and lexicons.
- Widely used in education and research.



https://en.wikipedia.org/wiki/Hypernymy_and_hyponymy



spaCy

- Fast, efficient NLP library.
- Robust tokenization, parsing, NER.
- Pretrained language models included.
- Industrial-Strength library.



<https://spacy.io>

https://en.wikipedia.org/wiki/Hypernymy_and_hyponymy



WordNet

- Lexical database for English words.
- Organizes words into synonym sets.
- Maps semantic relationships (e.g., hypernyms, antonyms).
- Supports NLP tasks like similarity analysis.

<https://wordnet.princeton.edu>

<https://en.wikipedia.org/wiki/WordNet>



Jurafsky (YTW) NLP - 1

Semantic σ_{sem}

Syntactic σ_{str} (structure)

Sentiment σ_{sent}

Task) NER : Name Entity Recognition σ_{ner}

Apple is an American Company.

$\xrightarrow{\text{lower}}$ apple is an american Company.

Tokens: a, p, l, e, i, s, ..., l, c, ...

char-level \rightarrow _, -, !, --

Subword-level \rightarrow apple, is, an, america, #n, ...

word level \rightarrow apple, is, an, american, ...

Learning learnable learned learnt

establish
establishment
in establish
in establishment
:
!
assessment
assortment

()

establish
+
ment
+
in

Computation power = hardware

Processor
RAM
Disk

26 (a-z)

a₉ (10) → 2₁₀₀

! ? ≠ --- (...)

Stemming

morphology

run → run
run → run
run → run

active stem
action → act
activity
act

Lemmatization

{ good
better → good { go
best went → go
well
:
:

Regular Expression

(Regex) → Python import re

+ : one or more

a+ : "a", "aa", "aaa", ...

→ asterisk

* : zero or more

a*: "", "a", "aa", ...

\ (Caret) → $\overset{\circ}{N} \overset{\circ}{\sim} \overset{\circ}{\rightarrow} \text{char}(n)$
? → $\overset{\circ}{\sim} \overset{\circ}{\rightarrow} \overset{\circ}{\sim} \overset{\circ}{\rightarrow}$

