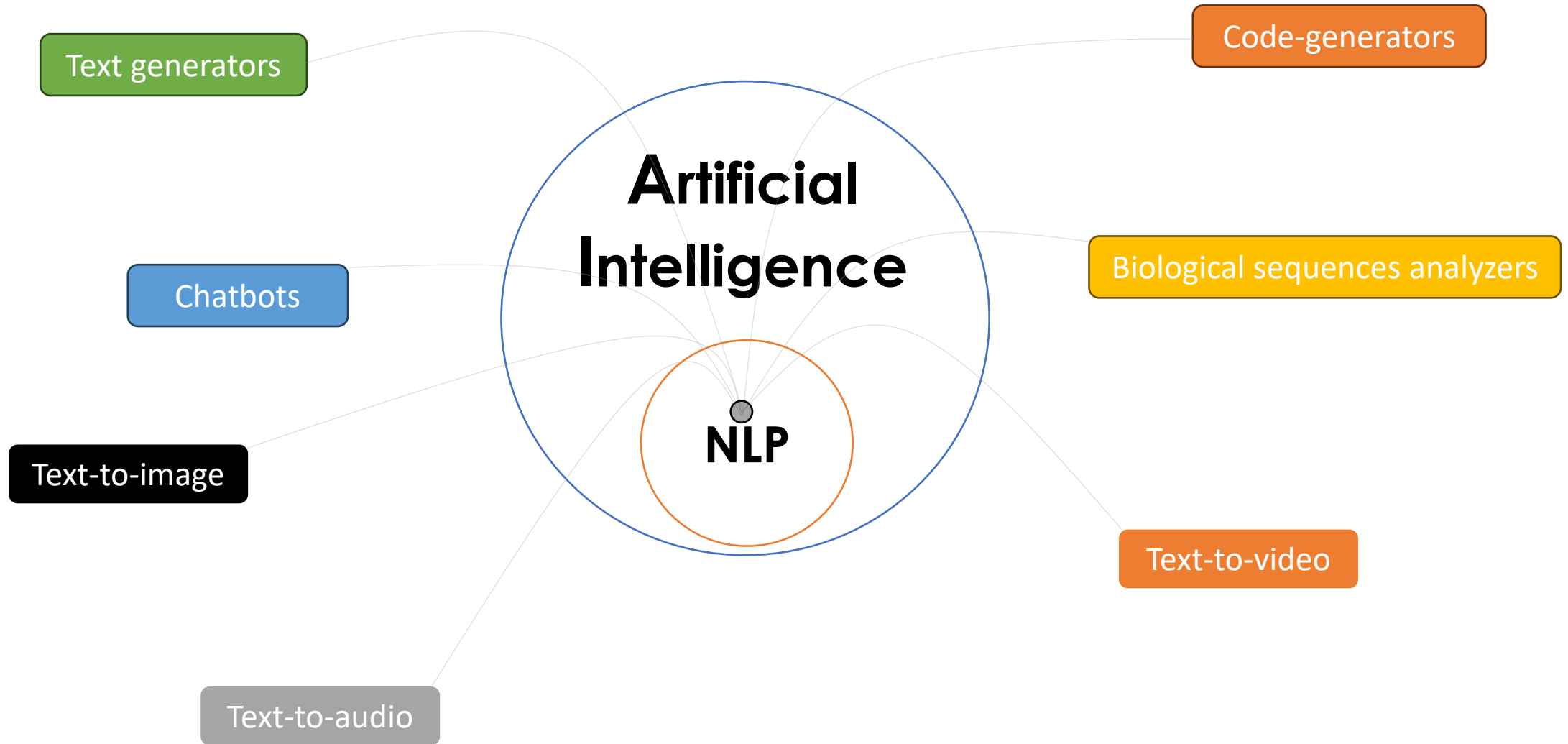


PLAN

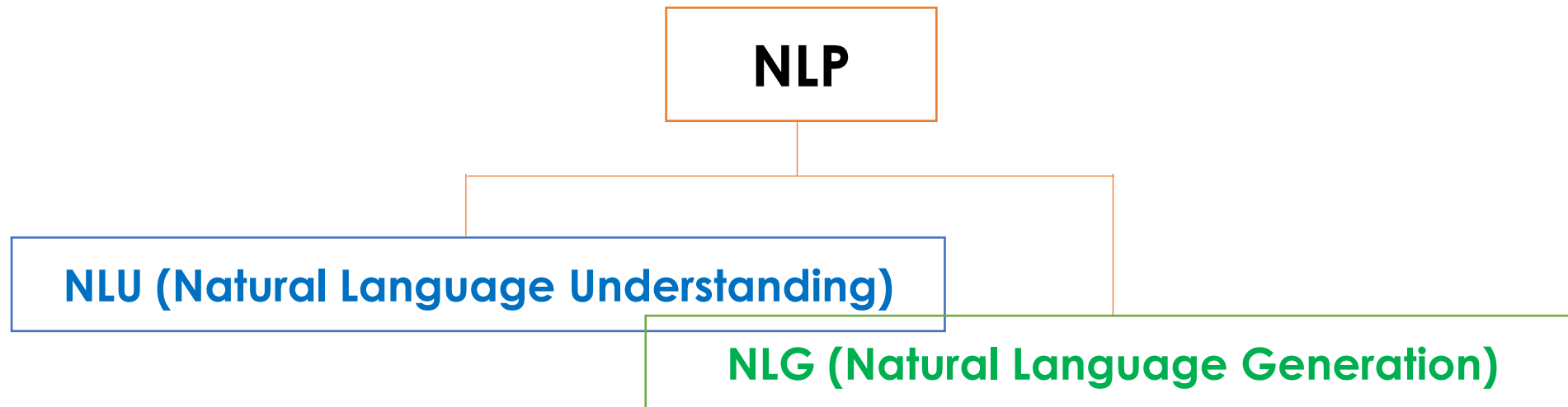
- What is NLP?
- NLP Applications
- How does NLP work?
- NLP techniques
- Libraries and Frameworks for NLP
- Python examples

NLP: A fast-growing research field in AI

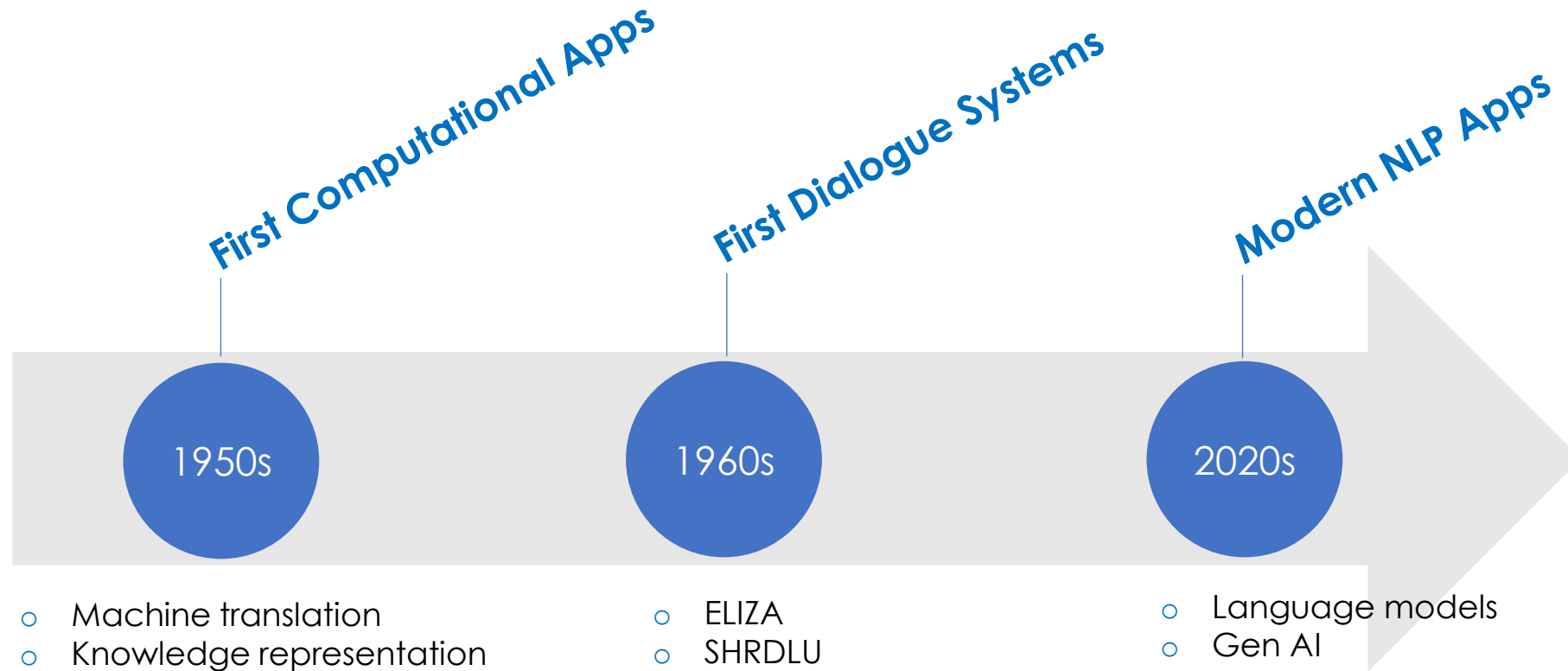


What is NLP ?

How to program computers to **analyze** the meanings of input text and **generate** meaningful, expressive output.

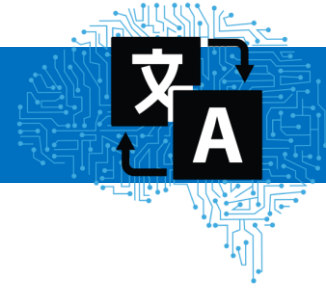


The early days of NLP

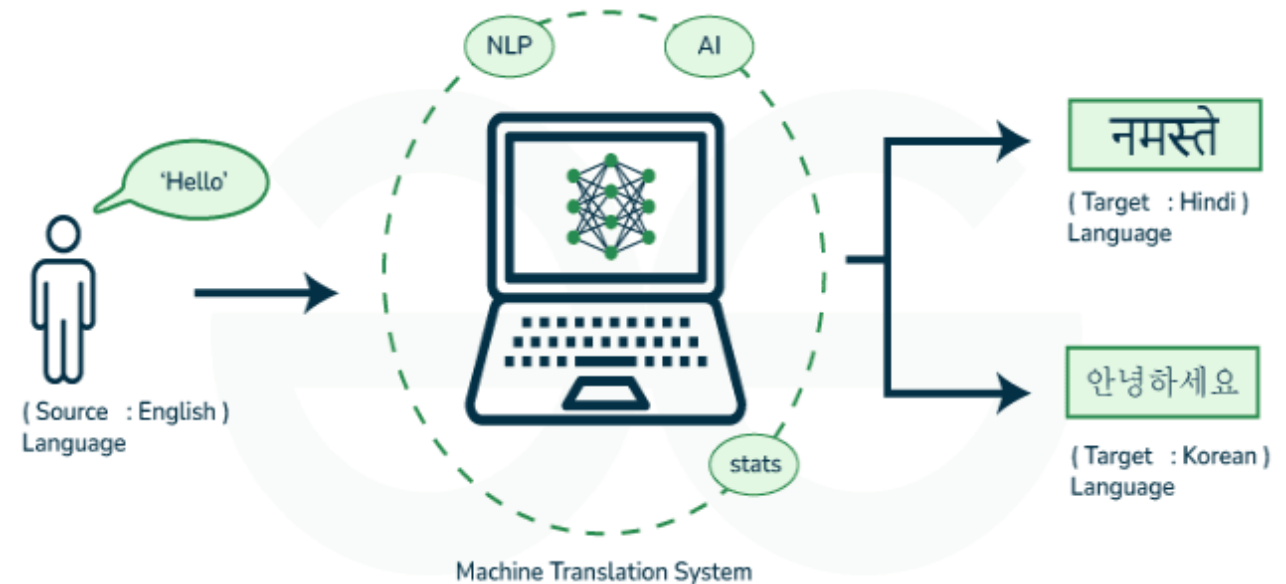


Main NLP Applications

1. Machine Translation

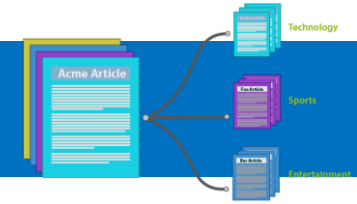


- **Rule-based**
(Grammar rules and dictionaries)
- **Statistical**
(Examine extensive human translations)
- **Neural**
(Training on Source-Target language dataset)
- **Hybrid**
(Use of multiple machine translation models)



Main NLP Applications

2. Text Classification



- **Document classification**
(Document categorization: Techno, Sport, Art,...)
- **Sentiment analysis**
(Classifying emotional quality)
- **Toxicity classification**
(Detecting threats, insults, hatred towards entities)
- **Spam detection**
(Classify emails as either spam or not)
- **Hadith authentication**
(Verify originality of Prophetic Hadiths)
- **Misinformation and Fake news detection,...**



Main NLP Applications

3. Named Entity Recognition



Extract entities in a piece of text into predefined categories such as personal **names**, **organizations**, **locations**, and **quantities**.

Andrew Yan-Tak Ng **PERSON** (**Chinese NORP** : 吳恩達; born **1976 DATE**) is a **British NORP** -born **American NORP** computer scientist and technology entrepreneur focusing on machine learning and **AI GPE** .
Ng was a co-founder and head of **Google Brain ORG** and was the former chief scientist at **Baidu ORG** ,
building the company's **Artificial Intelligence Group ORG** into a team of **several thousand CARDINAL** people.

Main NLP Applications

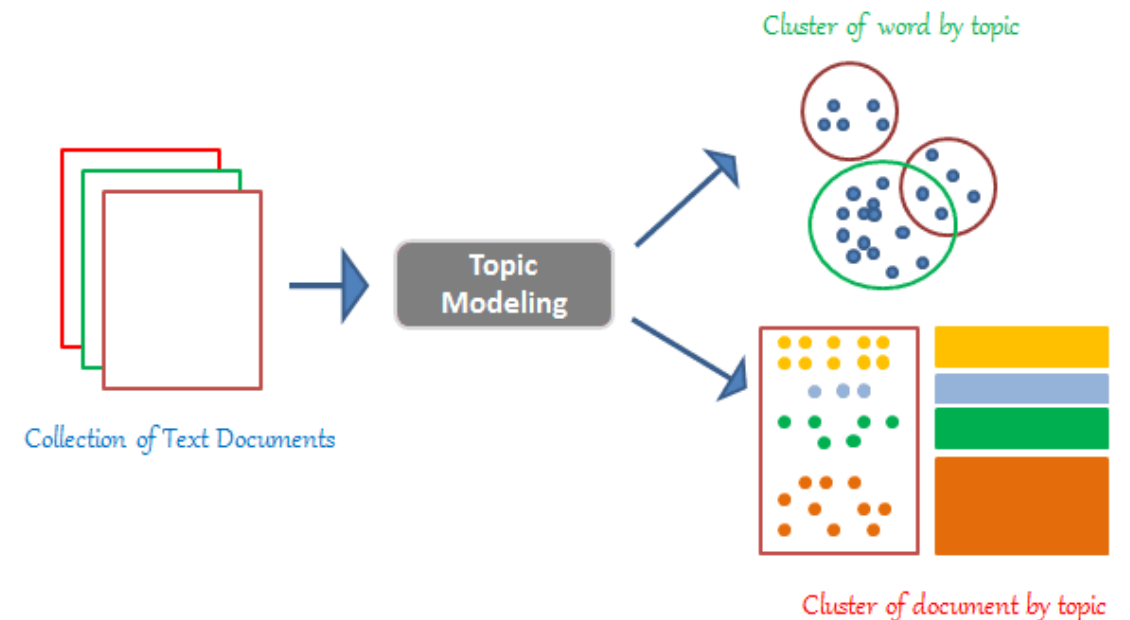
4. Topic Modeling



Unsupervised text mining task that takes a corpus of documents and discovers abstract topics within that corpus.

Techniques:

- Latent Semantic Analysis (LSA)
- Latent Dirichlet Allocation (LDA)
- LDA2Vec
- BERTopic



Main NLP Applications

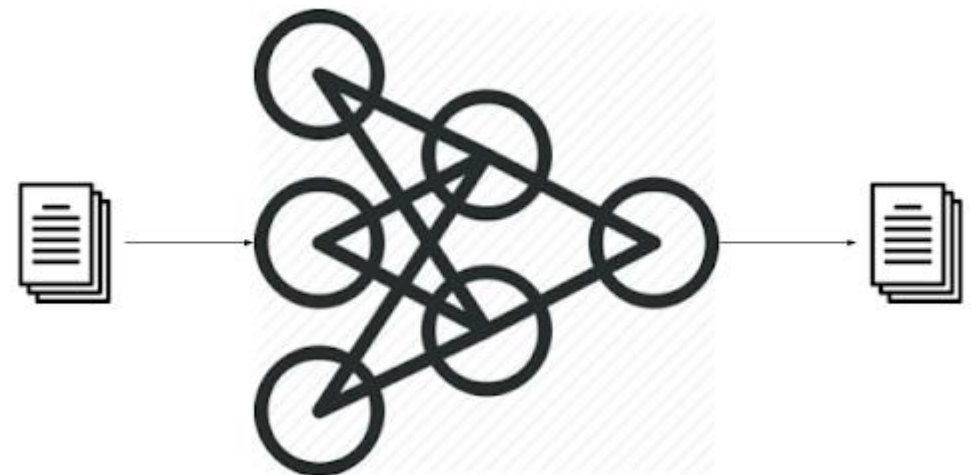
5. Text Generation



Automatically produces text that is similar to human-written text (such as: Tweets, Blogs, Essays, Computer code,..): LSTM-RNN, BERT, BARD, ChatGPT,...

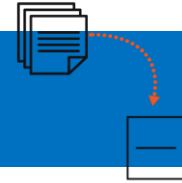
Variations:

- **Autocomplete:** predicts what word comes next
- **Chatbots:** automate one side of a conversation
 - Questions & Answers database
 - Conversation generation



Main NLP Applications

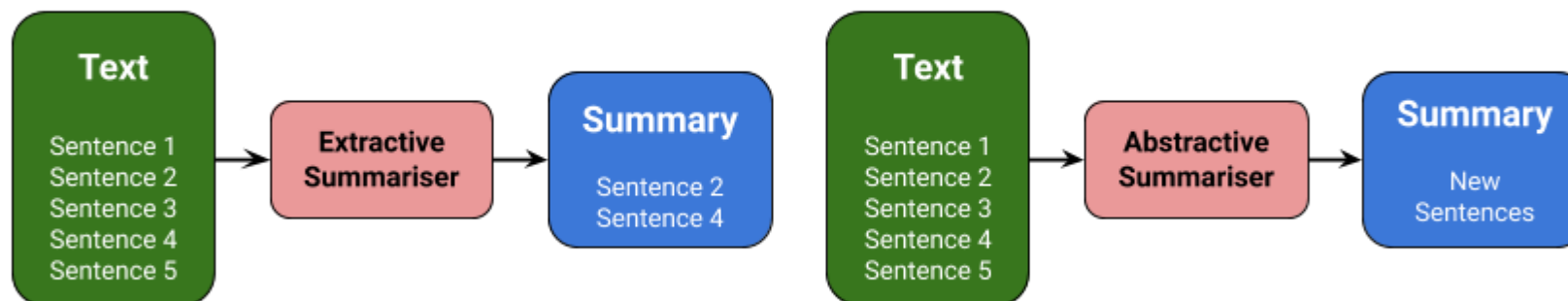
6. Text Summarization



Shortening text to highlight the most relevant information

Variations:

- **Extraction:** extracting the most important sentences from a long text and combining these to form a summary
- **Abstraction:** writing the abstract that includes words and sentences that are not present in the original text

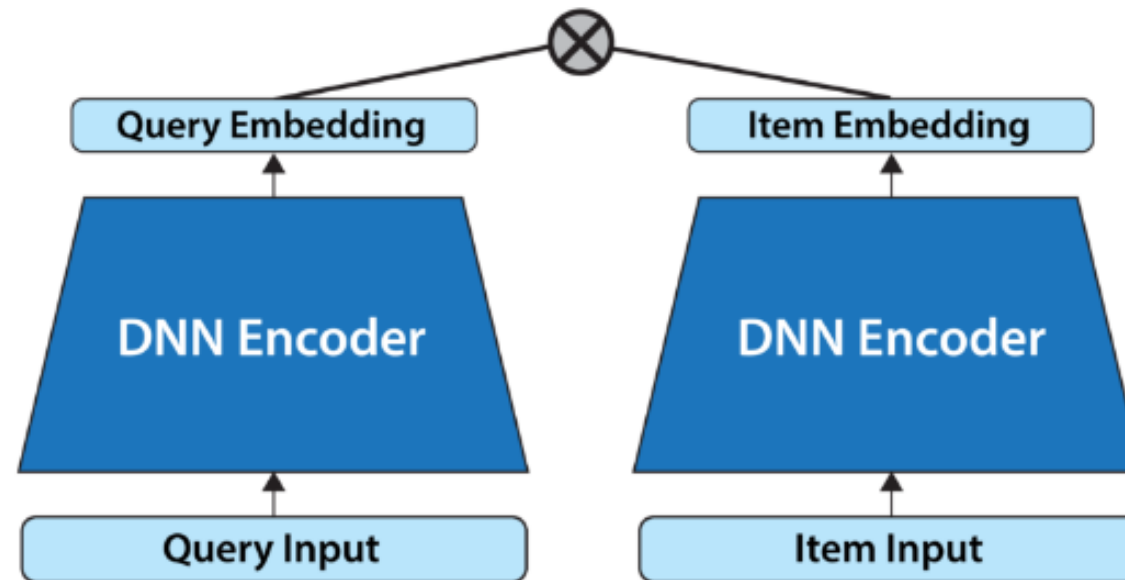


Main NLP Applications

7. Information Retrieval



Finds (indexing and matching) the documents that are most relevant to a query.



- **Indexing:** using a vector space
- **Matching:** using similarity score

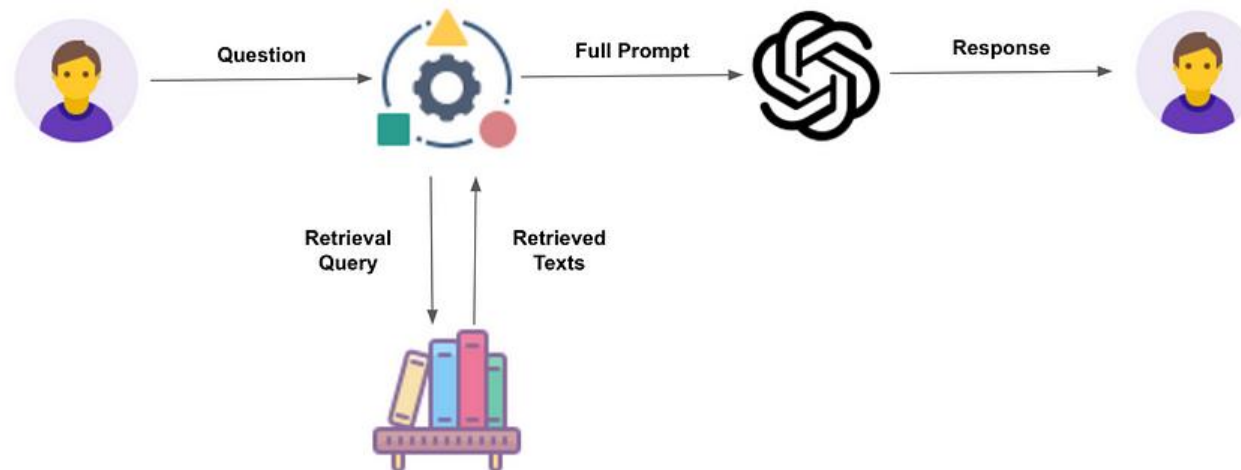
Main NLP Applications

8. Question/Answering



Answering questions asked by humans in a natural language

- **Multiple choice:** question problem is composed of a question and a set of possible answers
- **Open-domain:** the model provides answers to questions in natural language without any options provided



Main NLP Applications

9. Other NLP apps

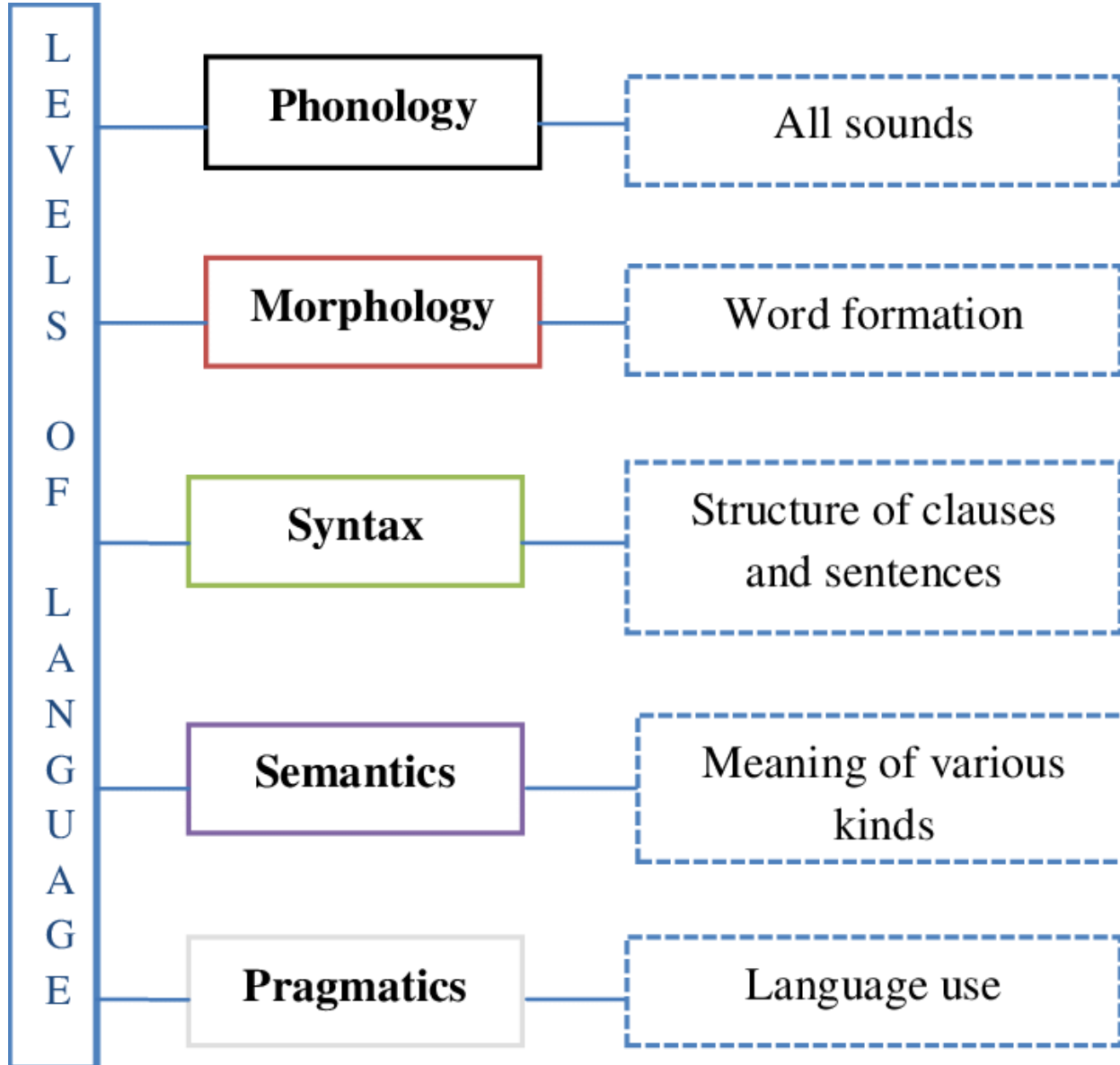
- **Grammatical error correction:** encode grammatical rules to correct the grammar within text.
- **Part-of-Speech Tagging:** classifying words in a text according to their grammatical categories (such as noun, verb, and adjective).
- **Language modeling:** building models that predict the probability of a sequence of words.
- **Speech recognition:** transform spoken language into a machine-readable format.

Main NLP Applications

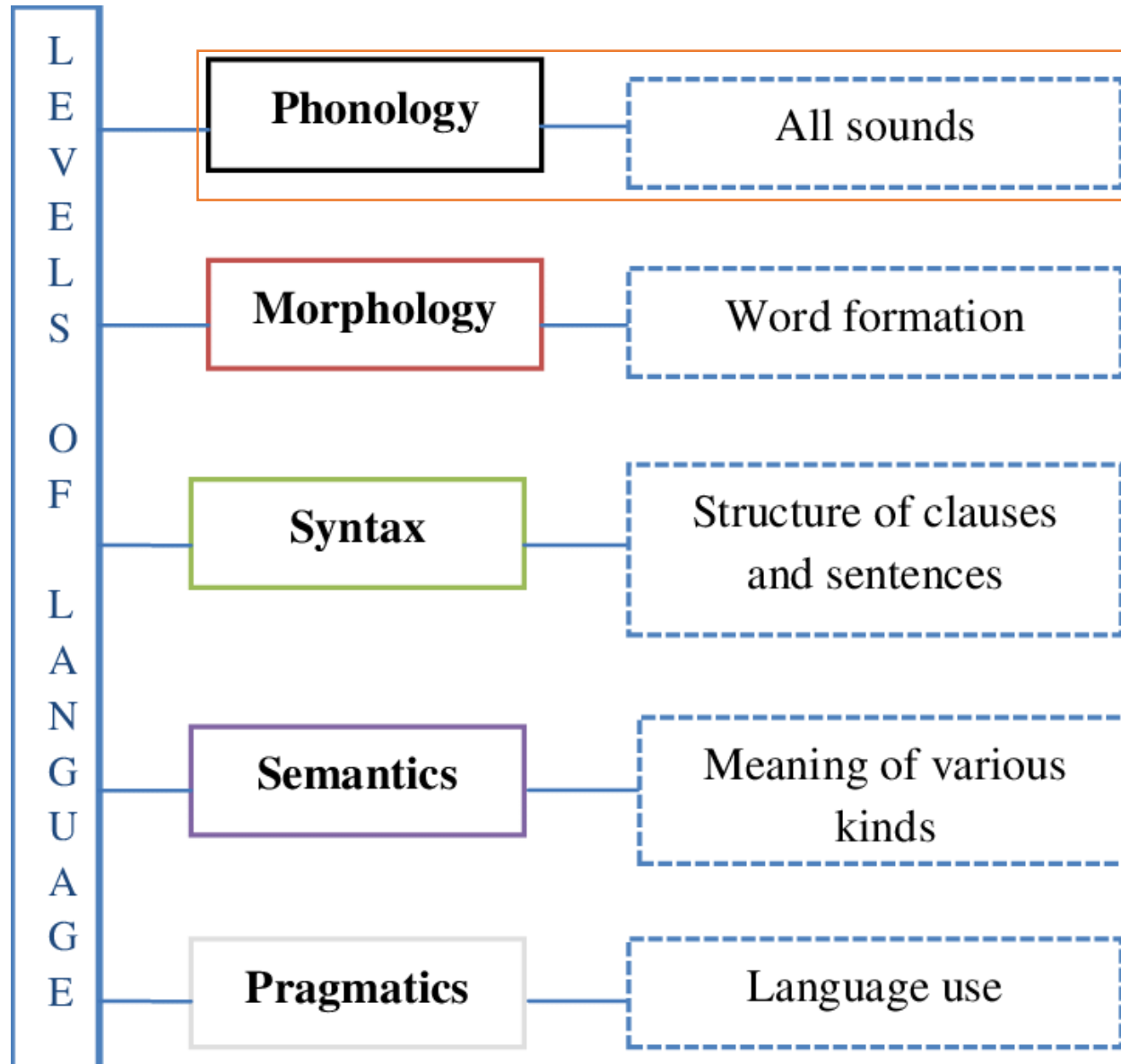
9. Other NLP apps

- **Grammatical error correction:** encode grammatical rules to correct the grammar within text.
- **Part-of-Speech Tagging:** classifying words in a text according to their grammatical categories (such as noun, verb, and adjective).
- **Language modeling:** building models that predict the probability of a sequence of words.
- **Speech recognition:** transform spoken language into a machine-readable format.

NLP Processing levels



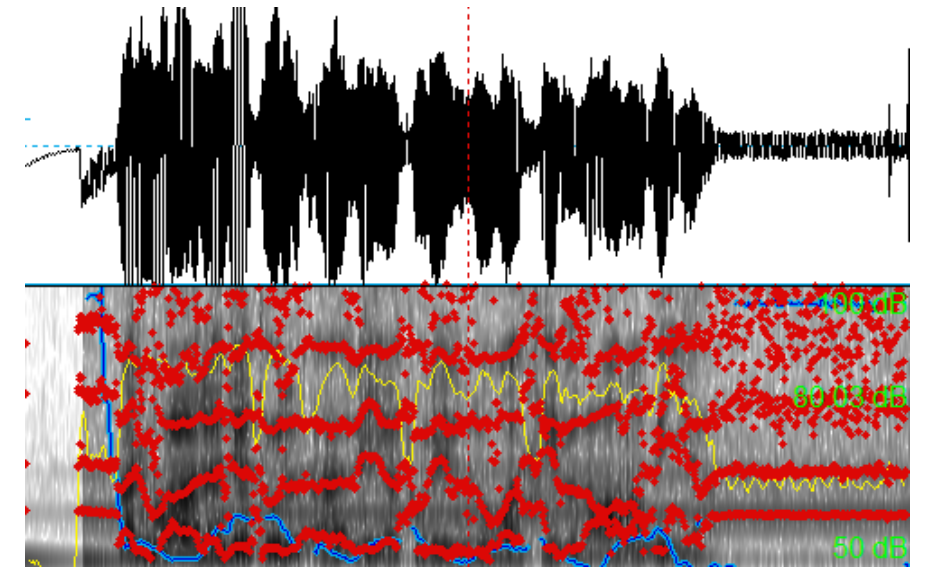
NLP Processing levels



- Phoneme detection
- Prosody identification
- Transitions

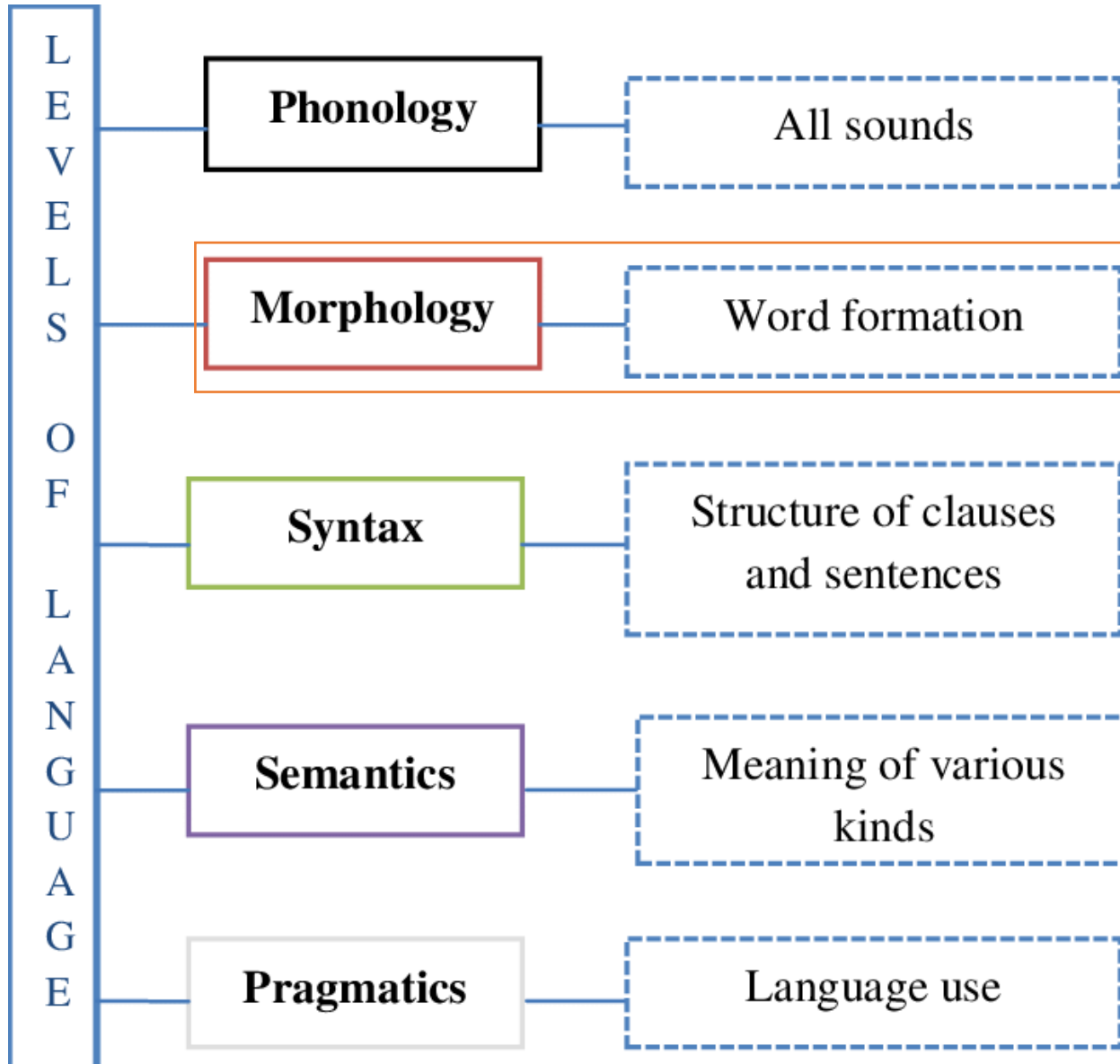


العربية اللغة الآلية المعالجة



al-mu'ālağatū al-'ālīatū liluğatī al-'ārabīa

NLP Processing levels



- Morphemes
- Tokens

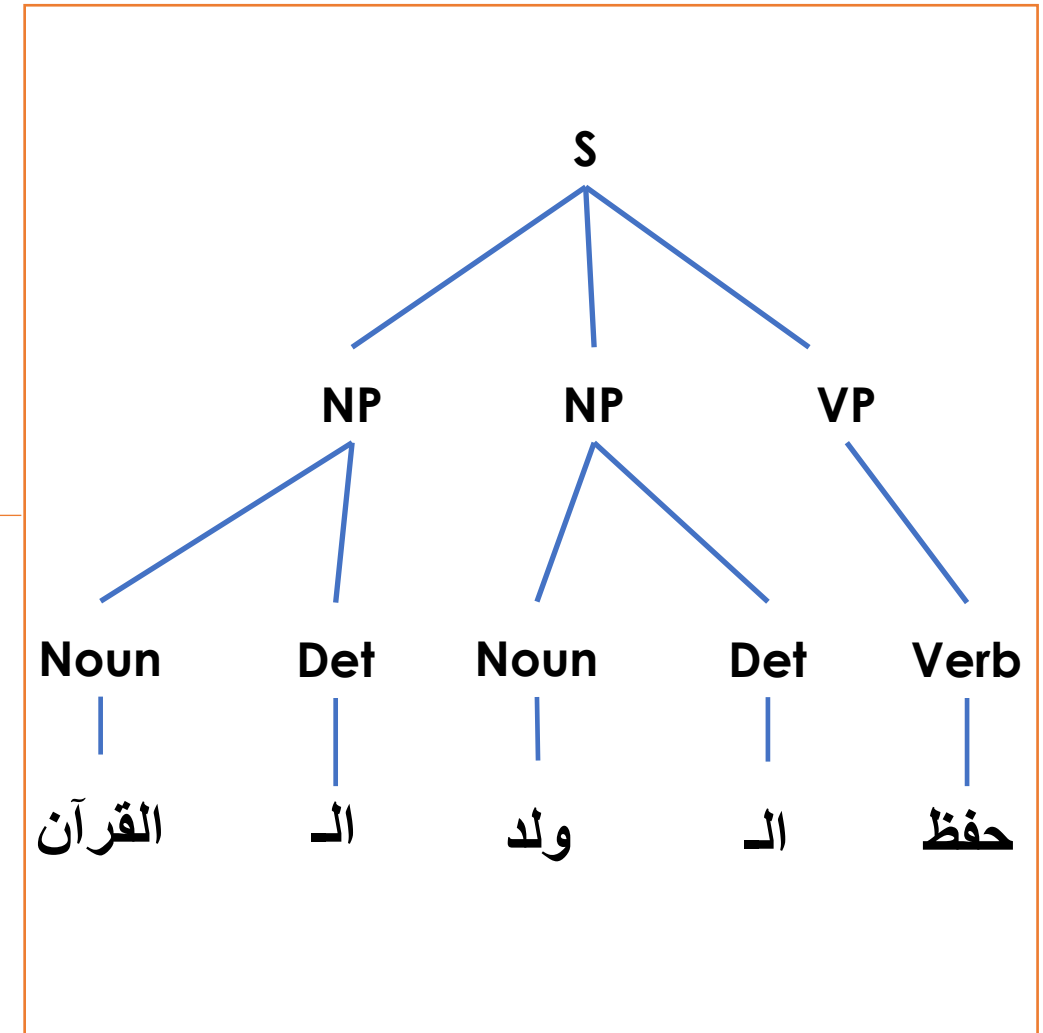
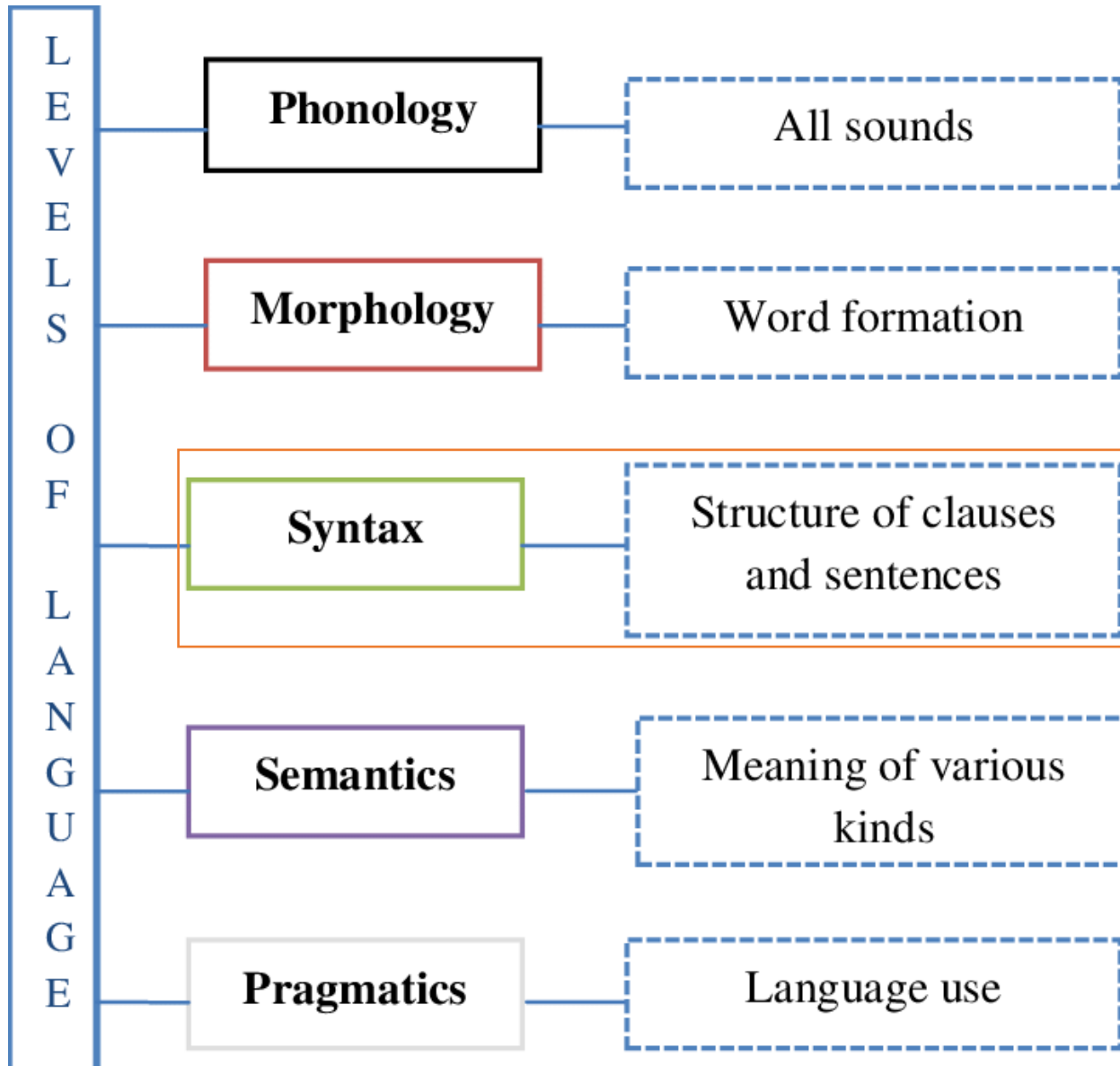
فأسقيناكموه

ف أ س ق ي ن ا ك م و ه

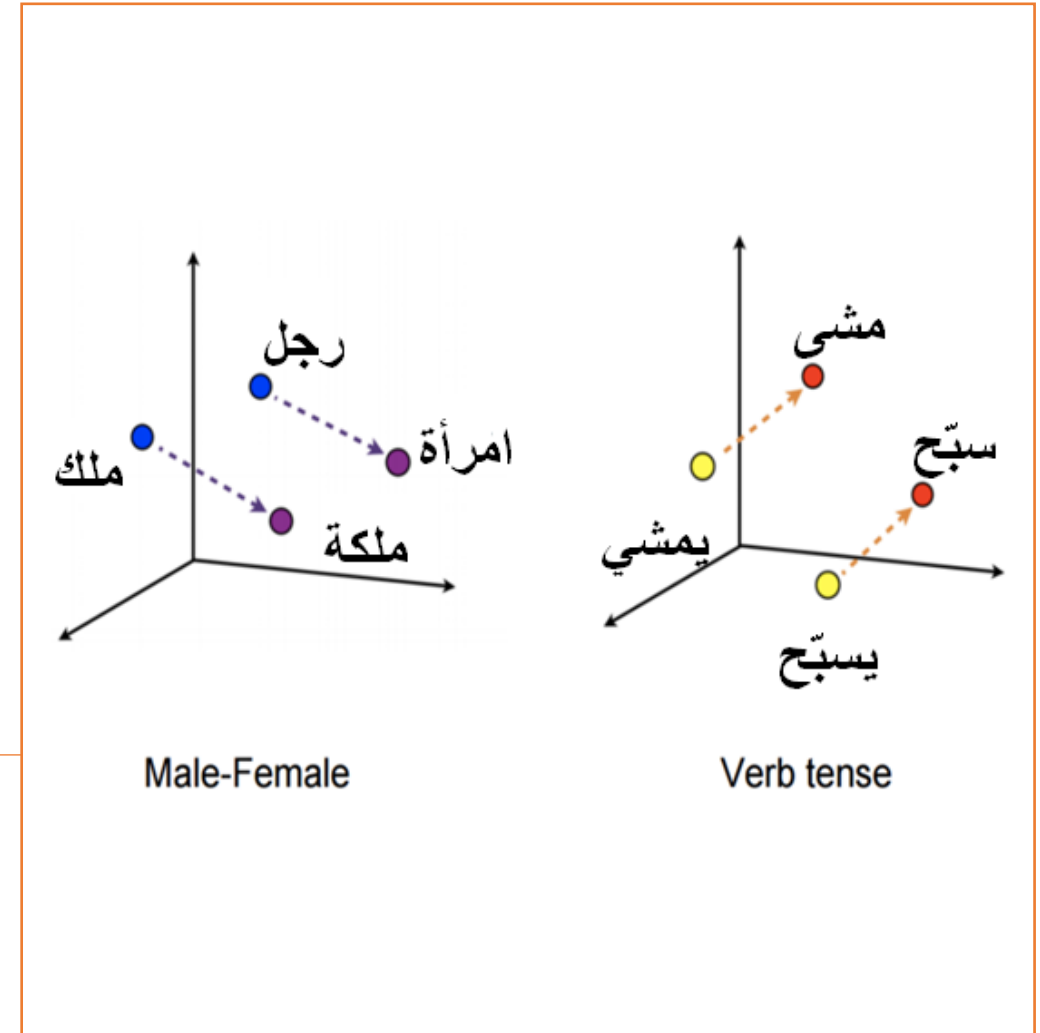
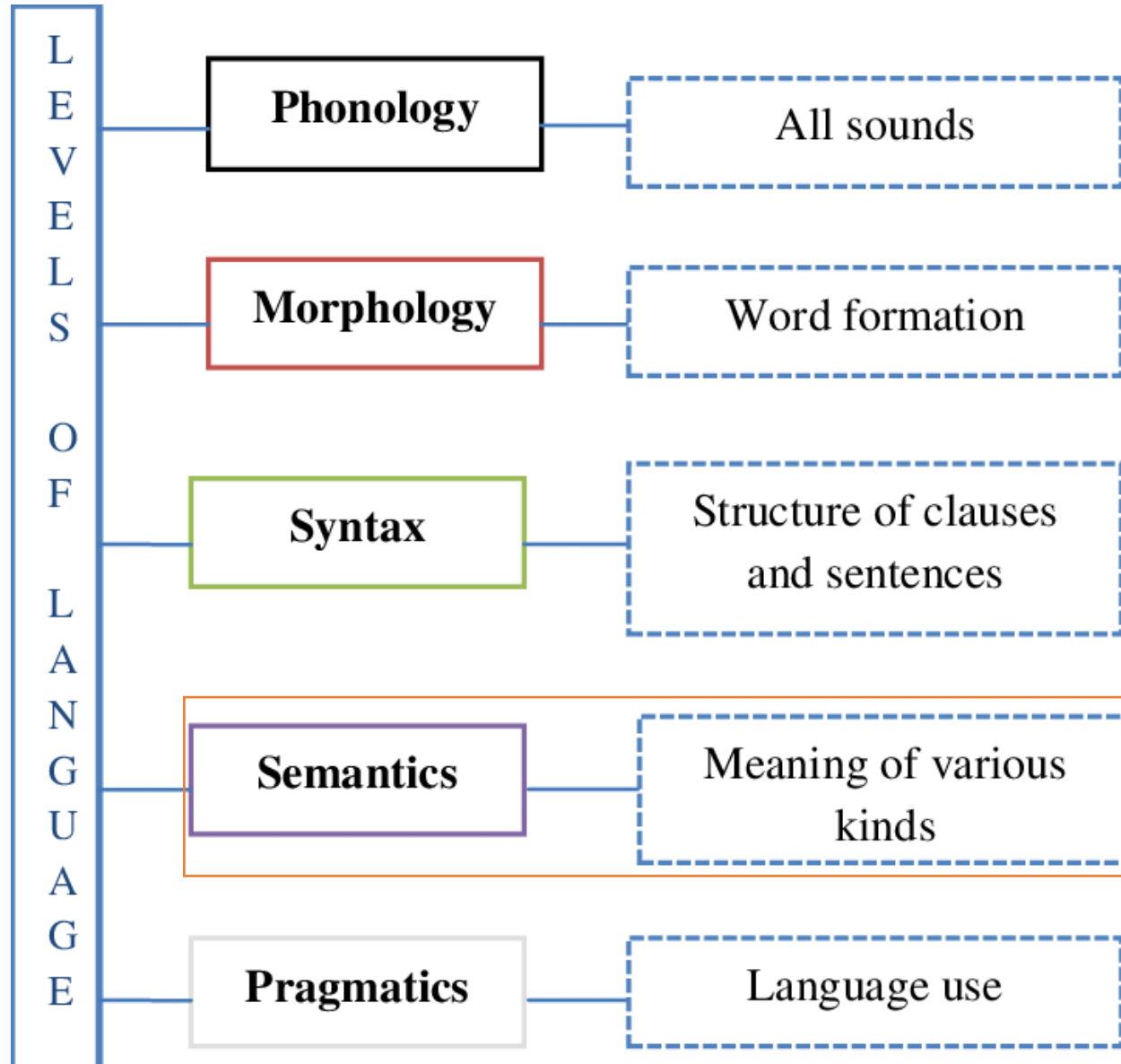
يذهب محمد إلى المسجد كل يوم

يذهب، محمد، إلى، المسجد، كل، يوم

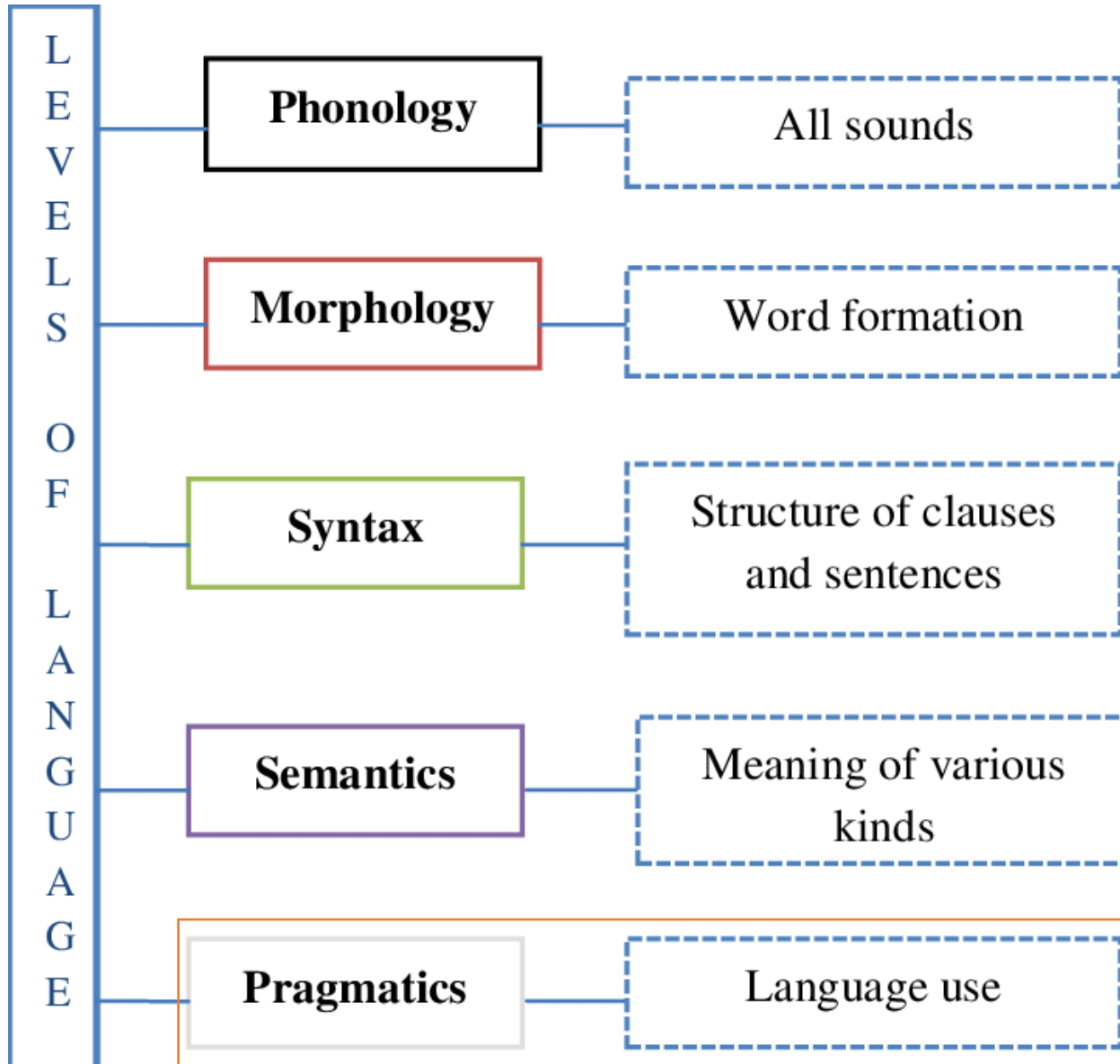
NLP Processing levels



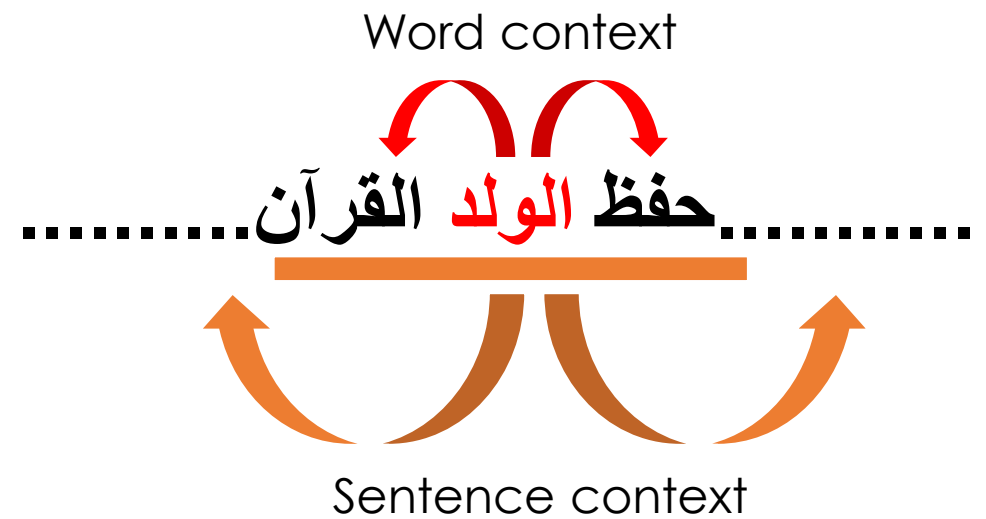
NLP Processing levels



NLP Processing levels



- o Meaning in the context



How does NLP work?

Data preprocessing



Feature extraction



Fed into NLP architecture

How does NLP work?

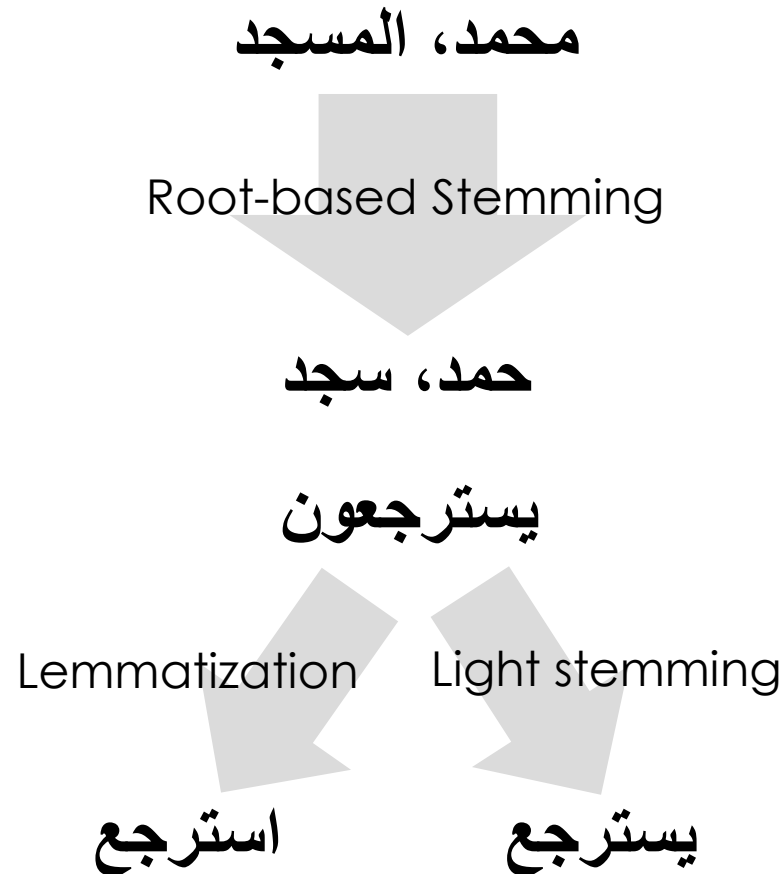
1. Data preprocessing

- **Stemming and Lemmatization:** converting words to their base forms.
- **Sentence segmentation:** breaks a large piece of text into meaningful sentence units.
- **Stop word removal:** remove words that don't add much information to the text.
- **Tokenization:** splits text into individual words.

How does NLP work?

1. Data preprocessing

- Stemming and Lemmatization:

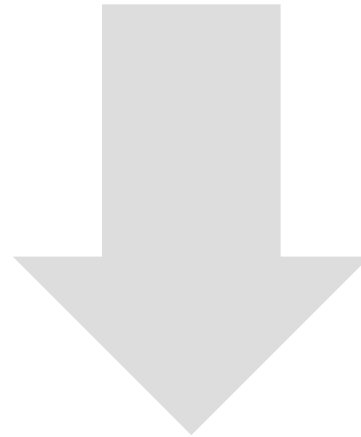


How does NLP work?

1. Data preprocessing

- Stop word removal:

يذهب محمد إلى المسجد كل يوم



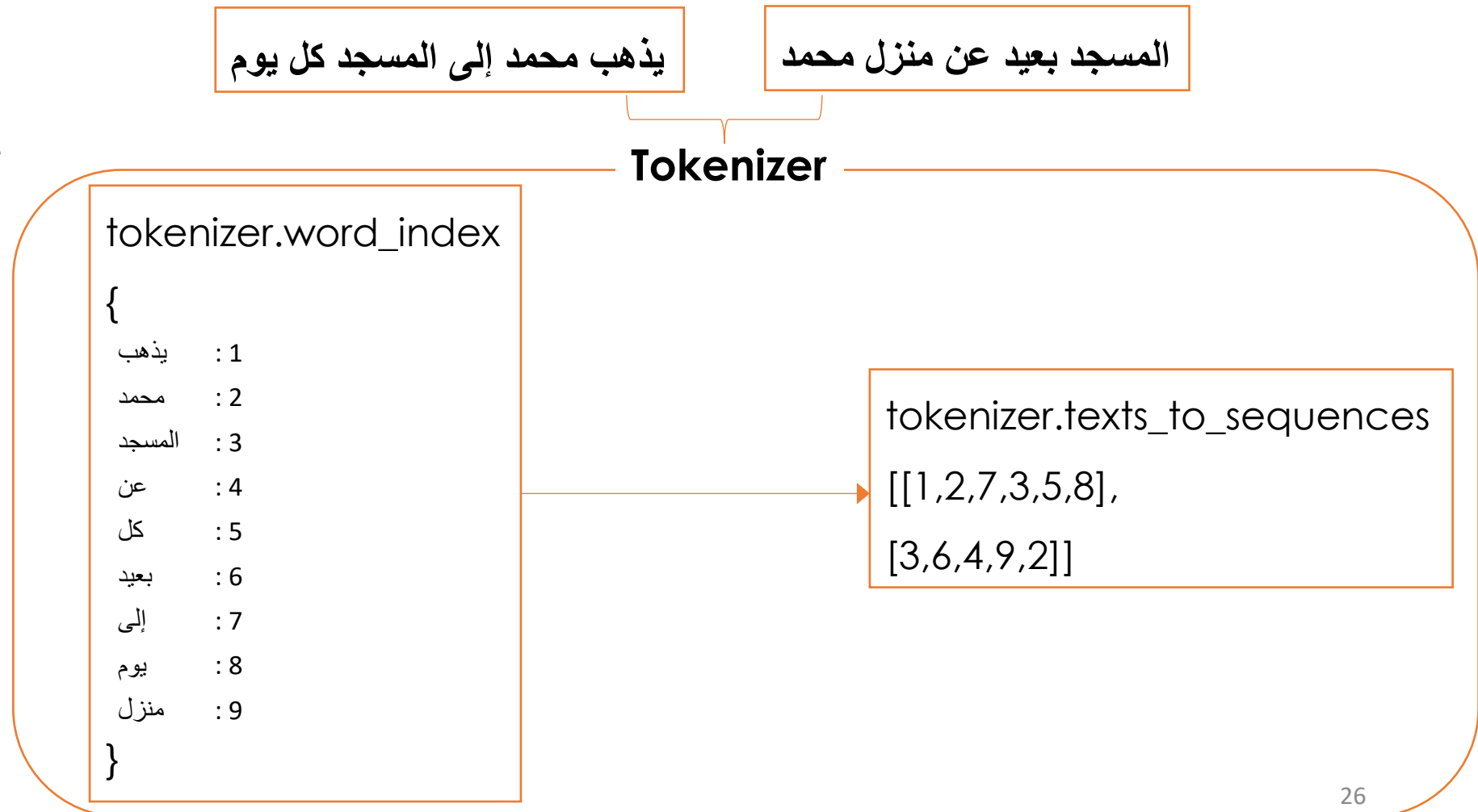
يذهب، محمد، المسجد، يوم

How does NLP work?

1. Data preprocessing

○ Tokenization:

- Word_piece
- Sentence_piece



How does NLP work?

2. Feature extraction

- **Bag-of-Words**
- **One-Hot-Encoding**
- **N-Grams**
- **TF-IDF**
- **Word Embeddings**
 - **Word2Vec (CBoW, Skip-Gram)**
 - **GLoVE**

How does NLP work?

2. Feature extraction

- Bag-of-Words (BoW)

يذهب محمد إلى المسجد كل يوم، كل يوم

المسجد بعيد عن منزل محمد

Vectorizer

word_index

```
{
  يذهب : 1
  محمد : 2
  المسجد : 3
  عن : 4
  كل : 5
  بعيد : 6
  إلى : 7
  يوم : 8
  منزل : 9
}
```

	يذهب	محمد	المسجد	عن	كل	بعيد	إلى	يوم	منزل
S1	1	1	1	0	2	0	1	2	0
S2	0	1	1	1	0	1	0	0	1

How does NLP work?

2. Feature extraction

- One-Hot-Encoding

يذهب محمد إلى المسجد كل يوم

المسجد بعيد عن منزل محمد

Vectorizer

word_index

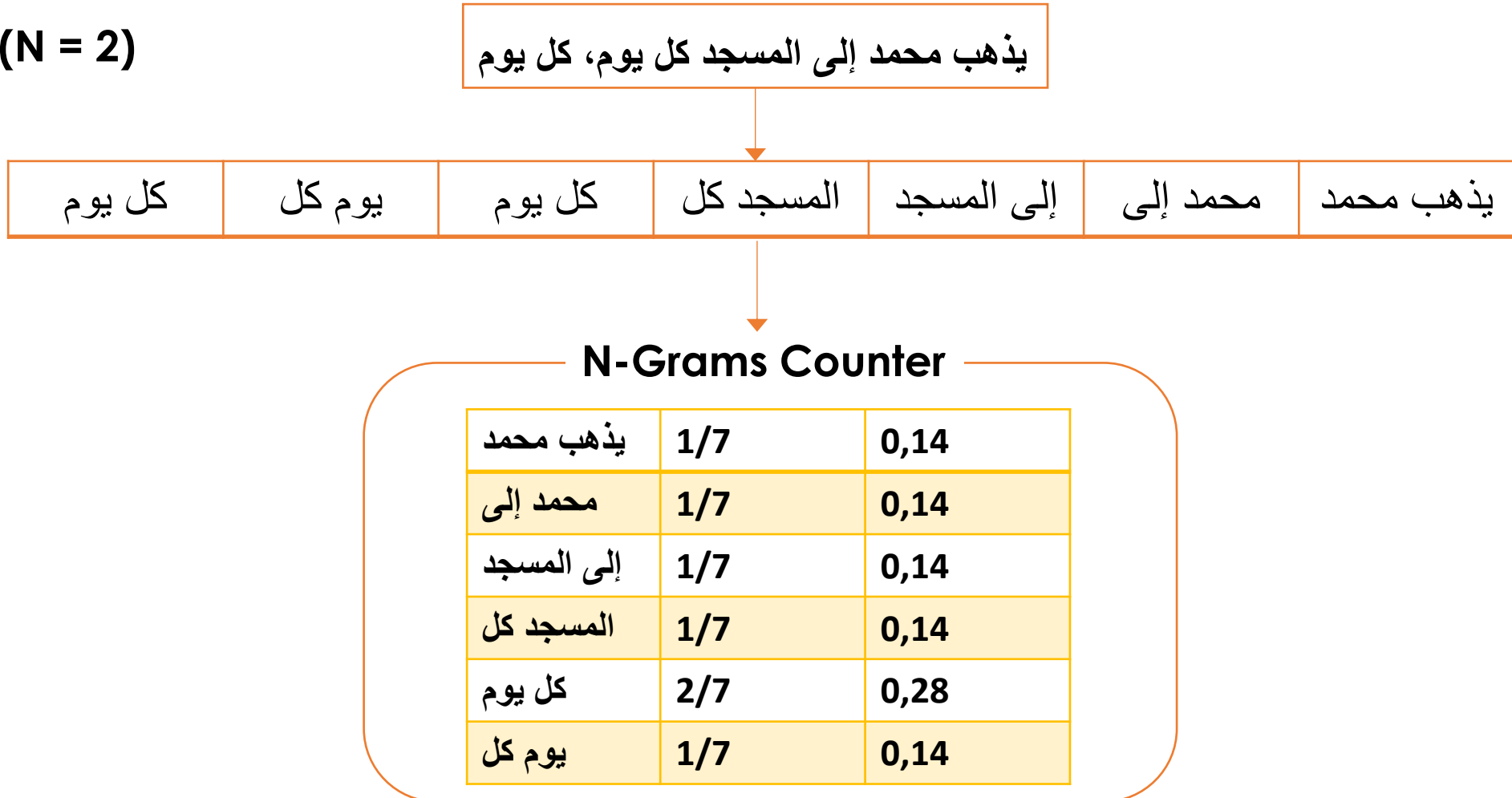
```
{
  يذهب : 1
  محمد : 2
  المسجد : 3
  عن : 4
  كل : 5
  بعيد : 6
  إلى : 7
  يوم : 8
  منزل : 9
}
```

	يذهب	محمد	المسجد	عن	كل	بعيد	إلى	يوم	منزل
يذهب	1	0	0	0	0	0	0	0	0
محمد	0	1	0	0	0	0	0	0	0
المسجد	0	0	1	0	0	0	0	0	0
عن	0	0	0	1	0	0	0	0	0
كل	0	0	0	0	1	0	0	0	0
بعيد	0	0	0	0	0	1	0	0	0
إلى	0	0	0	0	0	0	1	0	0
يوم	0	0	0	0	0	0	0	1	0
منزل	0	0	0	0	0	0	0	0	1

How does NLP work?

2. Feature extraction

- N-Grams (N = 2)



How does NLP work?

2. Feature extraction

○ TF-IDF:

- Weights each word by its importance
- TF (Term Frequency) = $\text{Number of occurrences of the word in document} / \text{Number of words in document}$
- IDF (Inverse Document Frequency) = $\log(\text{number of documents in the corpus} / \text{number of documents that include the word})$

D1	هذا أمر جيد وممتاز
D2	هذا أمر سيء للغاية

TF-IDF Vectorizer

TF

	هذا	أمر	جيد	سيء	و	للغاية	ممتاز
D1	1/5	1/5	1/5	0	1/5	0	1/5
D2	1/4	1/4	0	1/4	0	1/4	0

IDF

هذا	أمر	جيد	سيء	و	للغاية	ممتاز
$\log(2/2)$	$\log(2/2)$	$\log(2/1)$	$\log(2/1)$	$\log(2/1)$	$\log(2/1)$	$\log(2/1)$

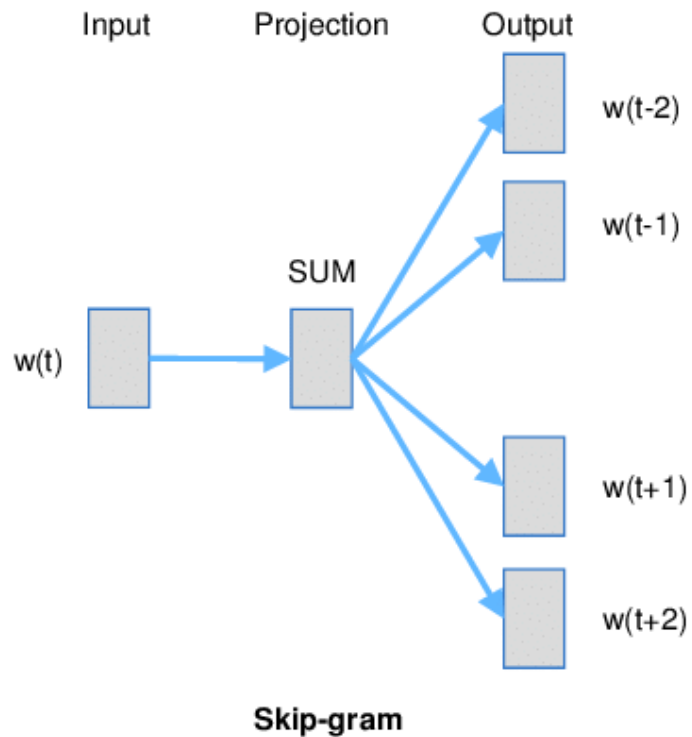
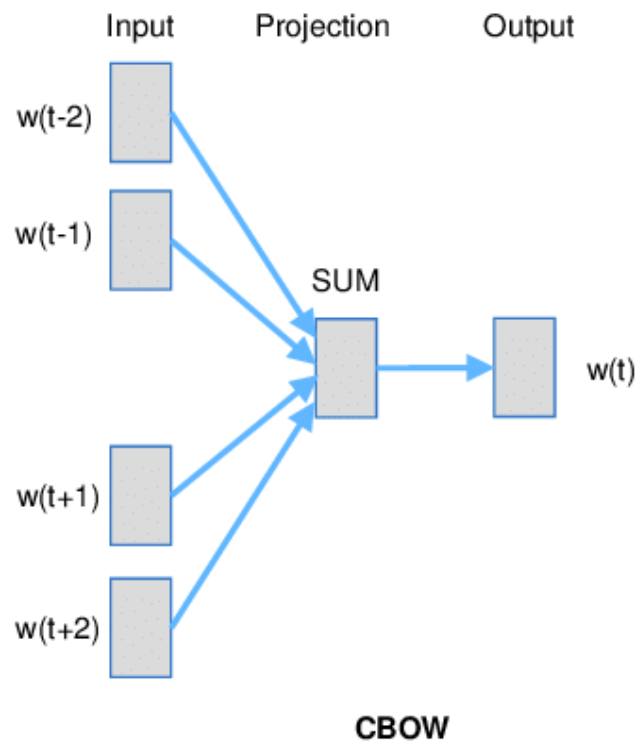
TF-IDF Features

	هذا	أمر	جيد	سيء	و	للغاية	ممتاز
D1	0	0	0,060	0	0,060	0	0,060
D2	0	0	0	0,075	0	0,075	0

How does NLP work?

2. Feature extraction

- Word embeddings (Word2Vec):



يذهب محمد إلى المسجد كل يوم



Word2Vec

يذهب	[0.2, 0.3, -0.1, 0.5, ...]
محمد	[0.1, -0.4, 0.6, -0.2, ...]
إلى	[0.3, -0.2, 0.4, 0.1, ...]
المسجد	[-0.5, 0.2, 0.3, -0.1, ...]
كل	[0.4, 0.1, -0.3, 0.2, ...]
يوم	[0.6, -0.3, 0.2, 0.4, ...]

NLP techniques

- **Traditional machine learning techniques**

- Logistic regression
- Naive Bayes
- Decision trees
- LDA/LSA
- Hidden Markov Models (HMM)

- **Deep learning techniques**

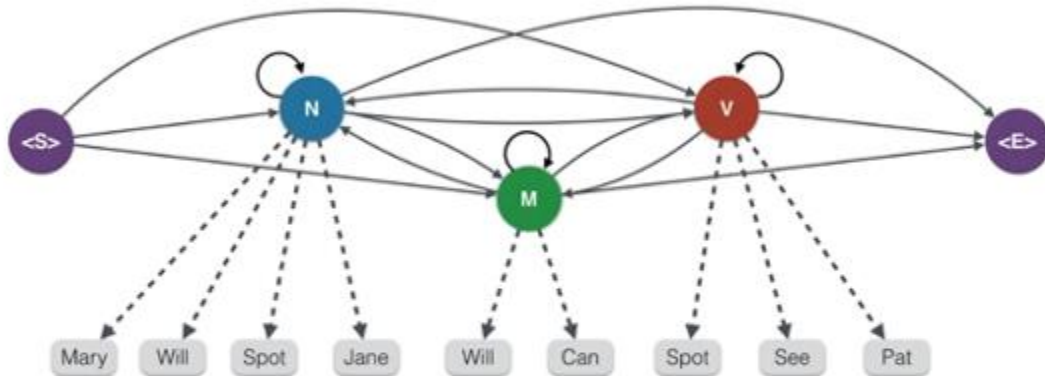
- Convolutional Neural Networks (CNN)
- Recurrent Neural Networks (RNN)
- Autoencoders
- Seq2Seq models
- Transformers

NLP techniques

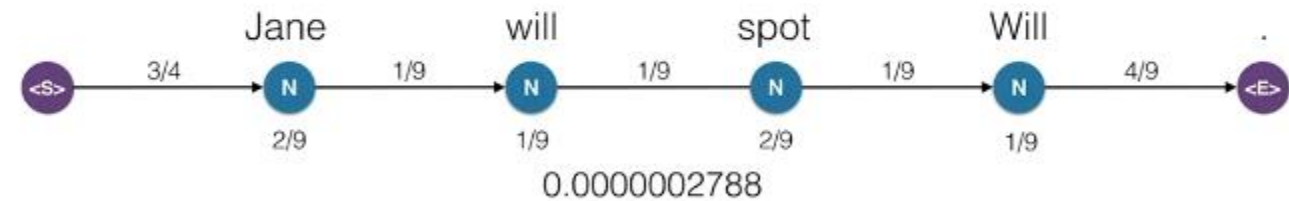
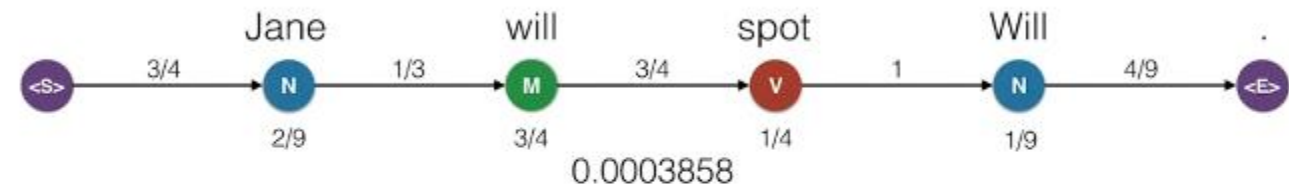
POS-Tagging with HMM

S = Jane will spot Will

What will be the most likely assignment for each word?



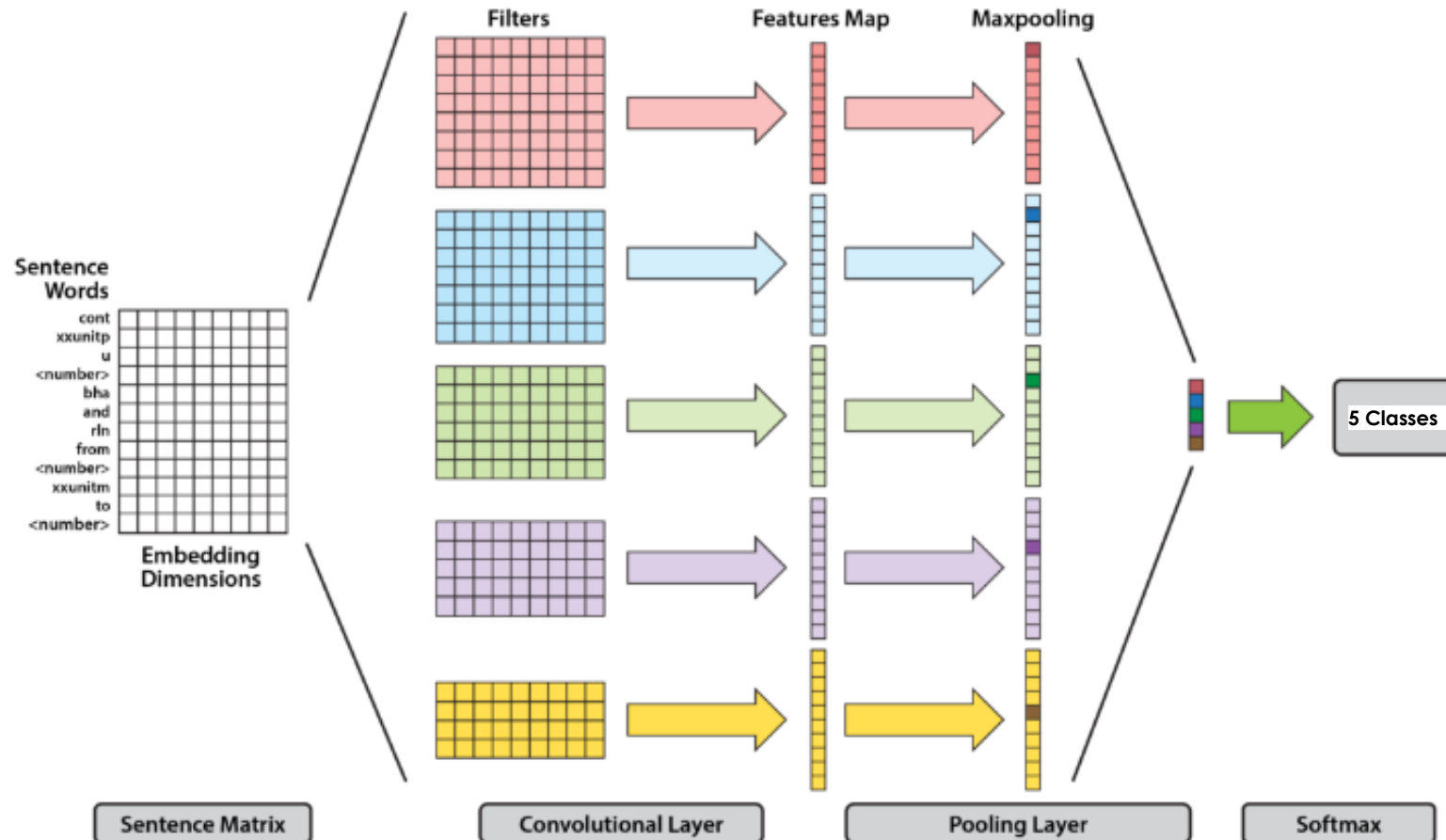
High probability sequence



NLP techniques

CNN-Based text classification

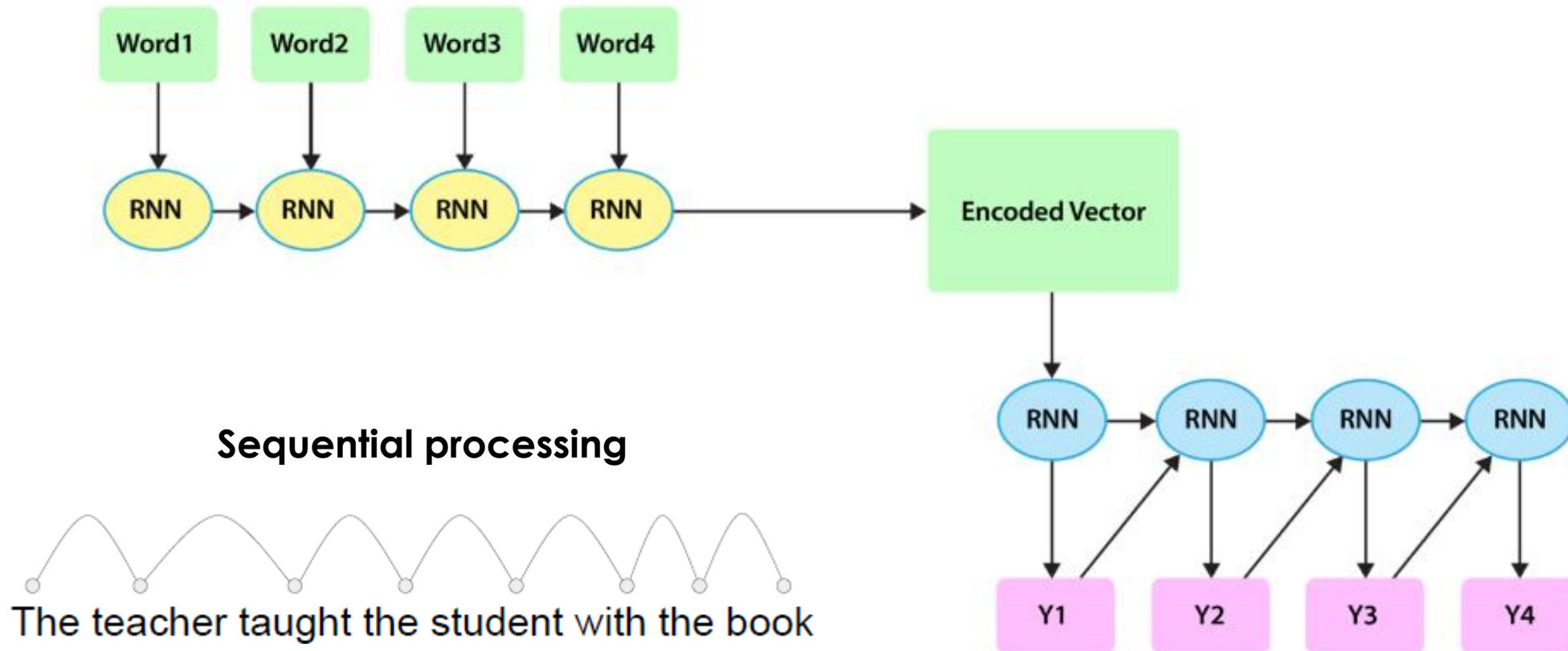
Given a sentence, a CNN uses convolutional layers to refine representations of input words, before combining them to render a classification



NLP techniques

RNN-Based Seq2Seq model for Machine translation

Given a sentence, a RNN encodes the sequence and then iteratively generates a translation

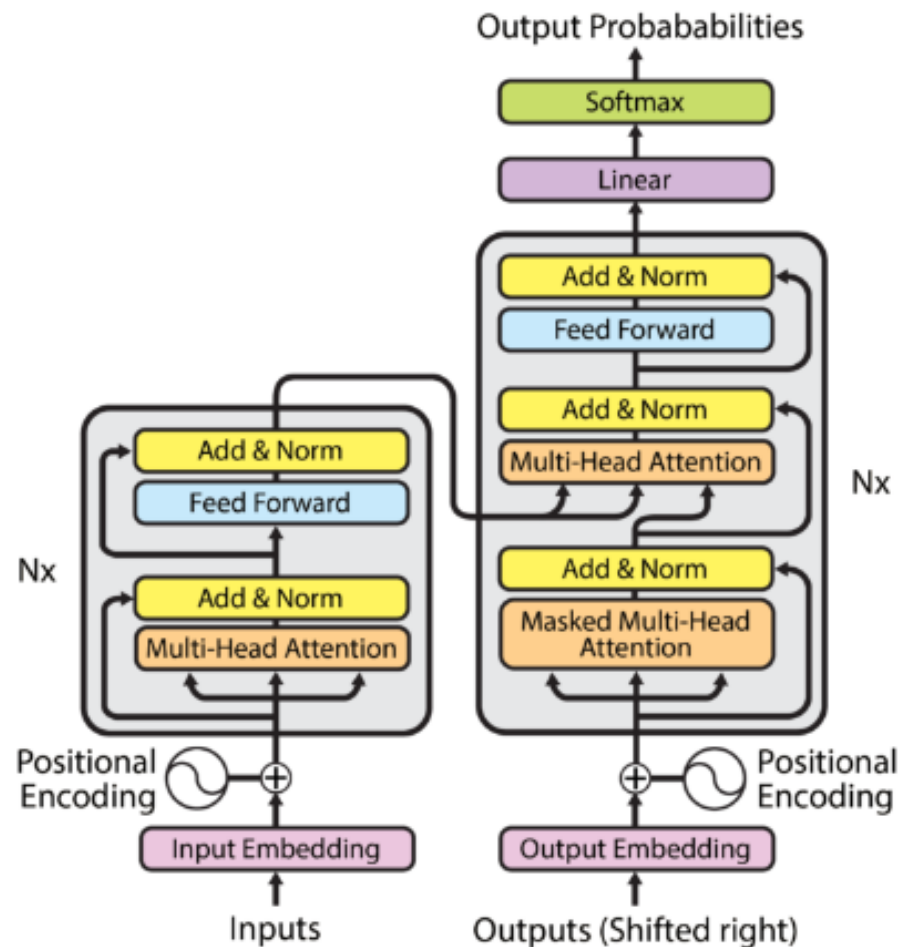
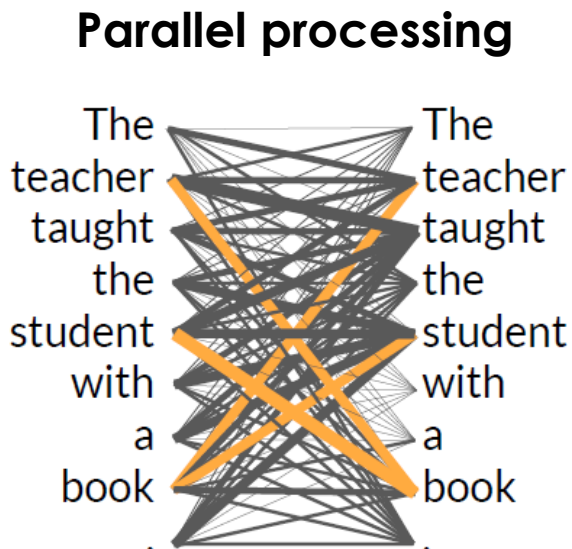


NLP techniques

Transformer architecture

Relies entirely on a self-attention mechanism to draw global dependencies between input and output, It is at the core of new language models:

- **Autoencoder (Encoder only):** BERT, ROBERTA
- **Autoregressive (Decoder only):** GPT, BLOOM
- **Seq2Seq (Encoder-Decoder):** T5, BART



NLP: Programming languages, libraries and Frameworks

- **Python**

- Natural Language Toolkit (NLTK)
- scikit-learn (Traditional machine learning algorithms)
- spaCy
- Deep learning libraries (keras, Tensorflow, PyTorch)
- Gensim
- Hugging Face: open-source models and implementations

- **R**

- TidyText
- Weka
- Word2Vec, SpaCyR, Tensorflow, PyTorch

- **JavaScript, Java, Julia**