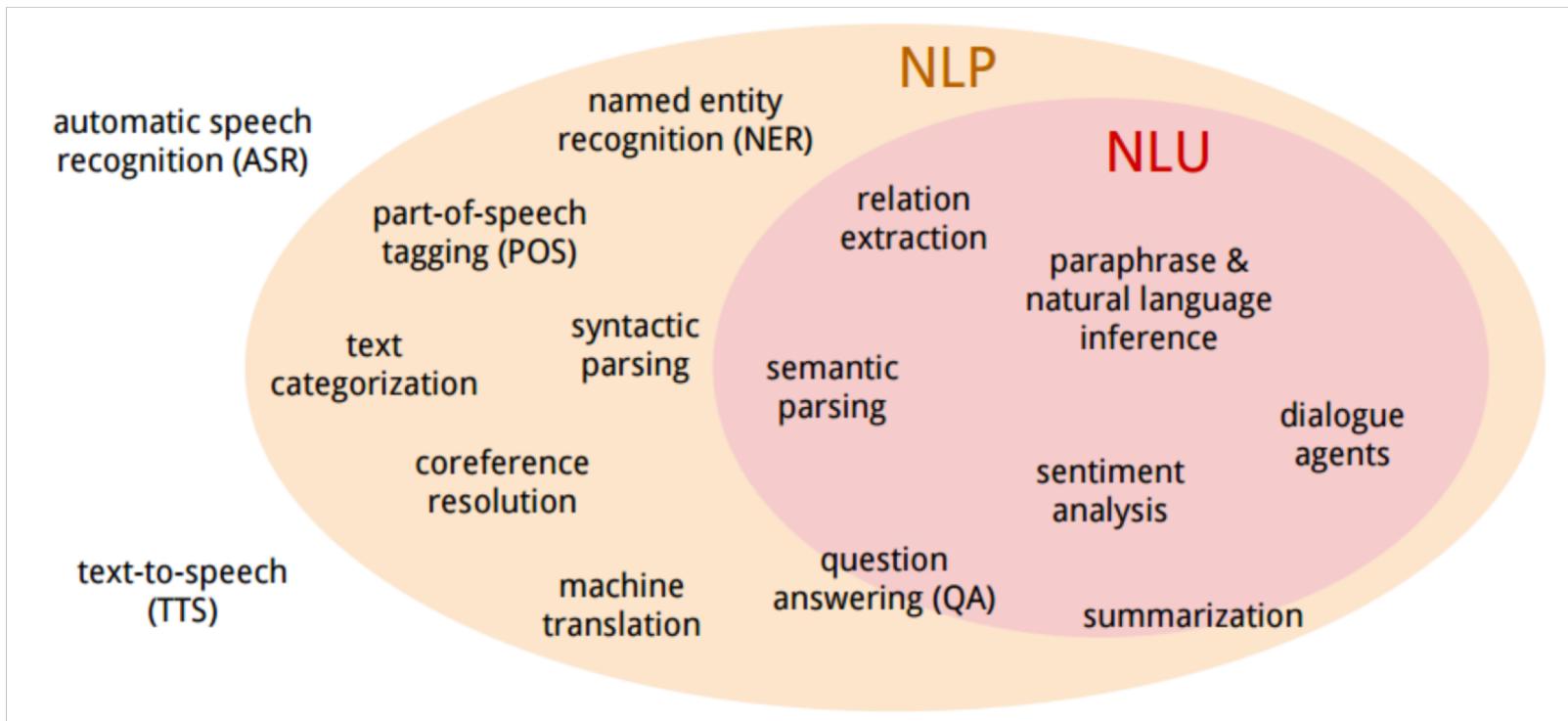


# BASIC CONCEPTS

**Natural Language Processing (NLP)** is a field of computer science and artificial intelligence, concerned with making computers process natural (human) language



**Computational Linguistics (CL)** is the field of using computers to understand language

# OVERVIEW OF THE BASIC NLP PIPELINE

TEXT

Regular Expressions

Tokenization Segmentation

Stemming Lemmatization

Part-of-Speech Tagging

Information Extraction

Named Entity Recognition

Named Entity Disambiguation

Coreference Resolution

Relationship Extraction

Natural Language Understanding

Sentiment Analysis

Semantic Analysis

Knowledge Extraction

Text Summarization

KNOWLEDGE

# TEXT PRE-PROCESSING: REGULAR EXPRESSIONS

A formal language for specifying text strings

- Used as the first model for any text processing
- Used as features in machine learning classifiers

Find all instances of the word “the” in a text

**the** – misses capitalized examples

**[tT]he** – Incorrectly returns *other* or *theology*

**[^a-zA-Z][tT]he[^a-zA-Z]**

False Positives: matching **there, then, other**

False Negatives: not matching **The**

Pattern	Matches	Examples
[A-Z]	An upper case letter	Drenched Blossoms
[a-z]	A lower case letter	my beans were impatient
[0-9]	A single digit	Chapter 1: Down the Rabbit Hole
[^A-Z]	Not an upper case letter	Oyfn pripetchik
[^Ss]	Neither ‘S’ nor ‘s’	I have no exquisite reason”
[^e^]	Neither e nor ^	Look here
a^b	The pattern a carat b	Look up <u>a^b</u> now

# SENTENCE SEGMENTATION

!, ? are relatively unambiguous

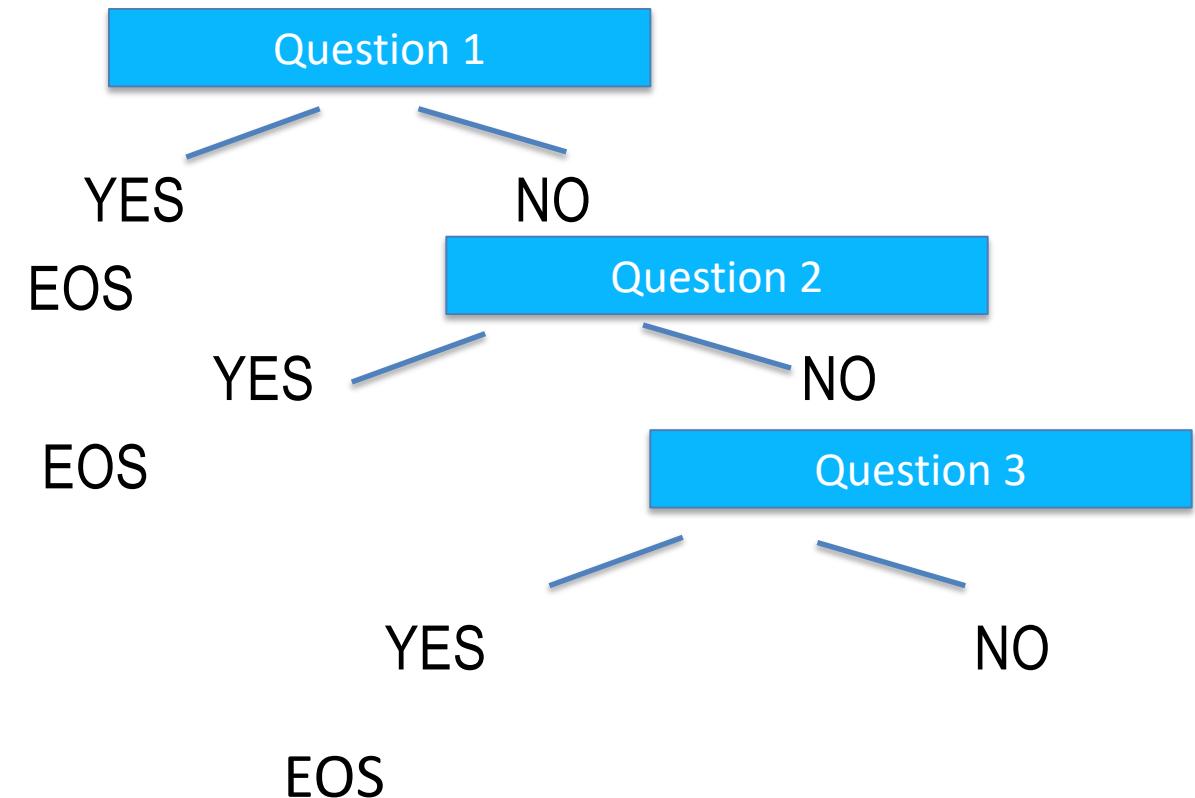
Period “.” is quite ambiguous

- Sentence boundary
- Abbreviations like Inc. or Dr.
- Numbers like .02% or 4.3

Build a binary classifier

- Looks at a “.”
- Decides EndOfSentence/NotEndOfSentence
- Classifiers: hand-written rules, regular expressions, or machine-learning

Decision Trees – Build a Rule for End of Sentence



# SENTENCE SEGMENTATION

!, ? are relatively unambiguous

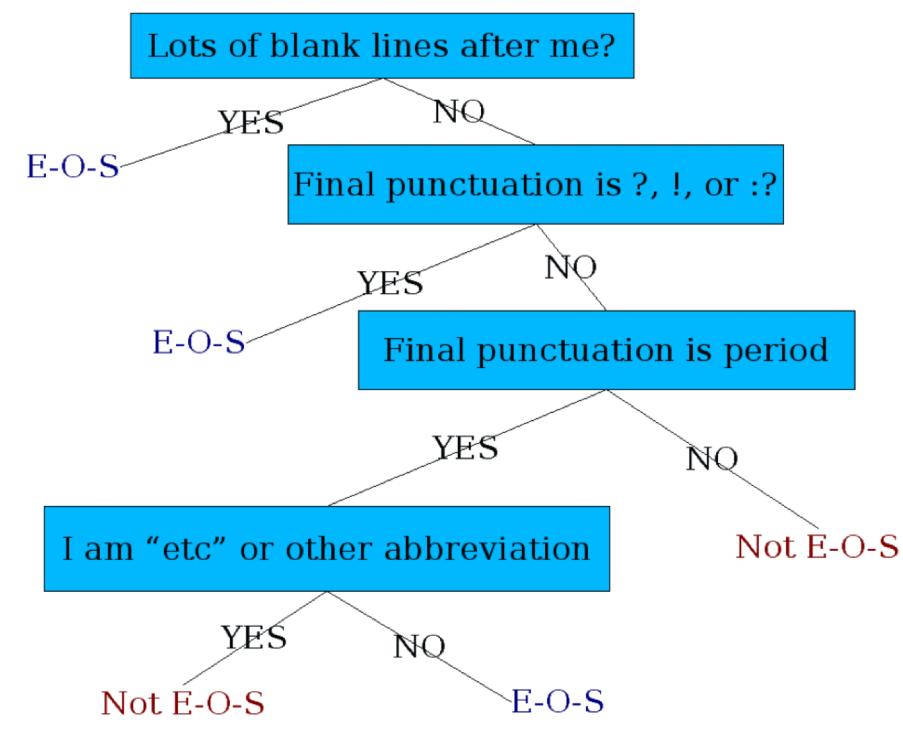
Period “.” is quite ambiguous

- Sentence boundary
- Abbreviations like Inc. or Dr.
- Numbers like .02% or 4.3

Build a binary classifier

- Looks at a “.”
- Decides EndOfSentence/NotEndOfSentence
- Classifiers: hand-written rules, regular expressions, or machine-learning

## Decision Trees



# TOKENIZATION

**Tokenization** is the task of chopping up a character sequence up into *tokens*, i.e. words, *n*-grams, sentences

***n*-gram** is a contiguous sequence of *n* items from a given sequence of text

**Stop words** are the most common words that are of little value in search query

To be, or not to be: To be , or not to be  
Unigrams: to be or not to be  
Bigrams: to be or not to be  
Trigrams: to be or not to be  
Stopwords: to be or not

```
>>> from nltk.tokenize import sent_tokenize, word_tokenize  
  
>>> sent_tokenize(s)  
[ 'Good muffins cost $3.88\nin New York.', 'Please buy me\\ntwo of them.', 'Thanks.' ]  
>>> [word_tokenize(t) for t in sent_tokenize(s)]  
[['Good', 'muffins', 'cost', '$', '3.88', 'in', 'New', 'York', '.'],  
 ['Please', 'buy', 'me', 'two', 'of', 'them', '.'], ['Thanks', '.']]
```

# STEMMING & LEMMATIZATION

- **Stemming** – reducing inflected or derived words to their base form
- **Lemmatization** – grouping together different inflected forms of a word
- Both are linguistic morphology terms with the same goal but different approaches

```
>>> wordnet_lemmatizer.lemmatize('are')
'are'
>>> wordnet_lemmatizer.lemmatize('is')
'is'

>>> wordnet_lemmatizer.lemmatize('is', pos='v')
u'be'
>>> wordnet_lemmatizer.lemmatize('are', pos='v')
u'be'
```

```
>>> from nltk.stem.porter import PorterStemmer
>>> porter_stemmer = PorterStemmer()
>>> porter_stemmer.stem('maximum')
u'maximum'
>>> porter_stemmer.stem('presumably')
u'presum'
>>> porter_stemmer.stem('multiply')
u'multipli'
```

```
>>> from nltk.stem import WordNetLemmatizer
>>> wordnet_lemmatizer = WordNetLemmatizer()
>>> wordnet_lemmatizer.lemmatize('dogs')
u'dog'
>>> wordnet_lemmatizer.lemmatize('churches')
u'church'
>>> wordnet_lemmatizer.lemmatize('aardwolves')
u'aardwolf'
```

# PART-OF-SPEECH (POS) TAGGING

**POS Tagging**, aka grammatical tagging or word-category disambiguation, is the process of marking a word in a text with a particular part of speech, based on both its definition and context

```
>>> import nltk  
>>> text = nltk.word_tokenize("Dive into NLTK: Part-of-speech tagging and POS Tagger")  
>>> text  
['Dive', 'into', 'NLTK', ':', 'Part-of-speech', 'tagging', 'and', 'POS', 'Tagger']  
>>> nltk.pos_tag(text)  
[('Dive', 'JJ'), ('into', 'IN'), ('NLTK', 'NNP'), (:, ':'), ('Part-of-speech', 'JJ'), ('tagging', 'NN'),  
('and', 'CC'), ('POS', 'NNP'), ('Tagger', 'NNP')]
```

## PennTreeBank

Tag	Description	Example	Tag	Description	Example
CC	coordin. conjunction	<i>and, but, or</i>	SYM	symbol	+%, &
CD	cardinal number	<i>one, two</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>	VBN	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, sing.	<i>IBM</i>	\$	dollar sign	\$
NNPS	proper noun, plural	<i>Carolinas</i>	#	pound sign	#
PDT	predeterminer	<i>all, both</i>	"	left quote	' or "
POS	possessive ending	's	"	right quote	' or "
PRP	personal pronoun	<i>I, you, he</i>	(	left parenthesis	[, (, <
PRP\$	possessive pronoun	<i>your, one's</i>	)	right parenthesis	], ), }, >
RB	adverb	<i>quickly, never</i>	,	comma	,
RBR	adverb, comparative	<i>faster</i>	.	sentence-final punc	. ! ?
RBS	adverb, superlative	<i>fastest</i>	:	mid-sentence punc	: ; ... --
RP	particle	<i>up, off</i>			

# CHUNKING – SHALLOW PARSING

Segment and label multi-token sequences  
using regular expressions and POS tags

W e	s a w	t h e	y e l l o w	d o g
PRP	VBD	DT	JJ	NN
B-NP	O	B-NP	I-NP	I-NP

```
>>> sentence = [("the", "DT"), ("little", "JJ"), ("yellow", "JJ"), ①
... ("dog", "NN"), ("barked", "VBD"), ("at", "IN"), ("the", "DT"), ("cat", "NN")]

>>> grammar = "NP: {<DT>?<JJ>*<NN>} " ②

>>> cp = nltk.RegexpParser(grammar) ③
>>> result = cp.parse(sentence) ④
>>> print(result) ⑤
(S
  (NP the/DT little/JJ yellow/JJ dog/NN)
  barked/VBD
  at/IN
  (NP the/DT cat/NN))
>>> result.draw() ⑥
```

