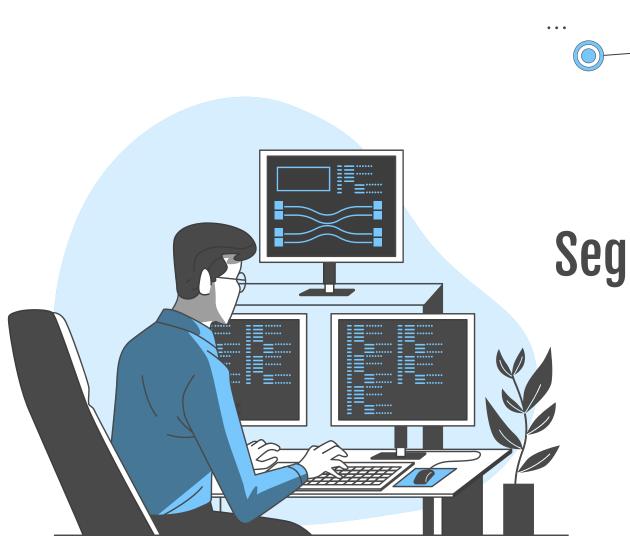


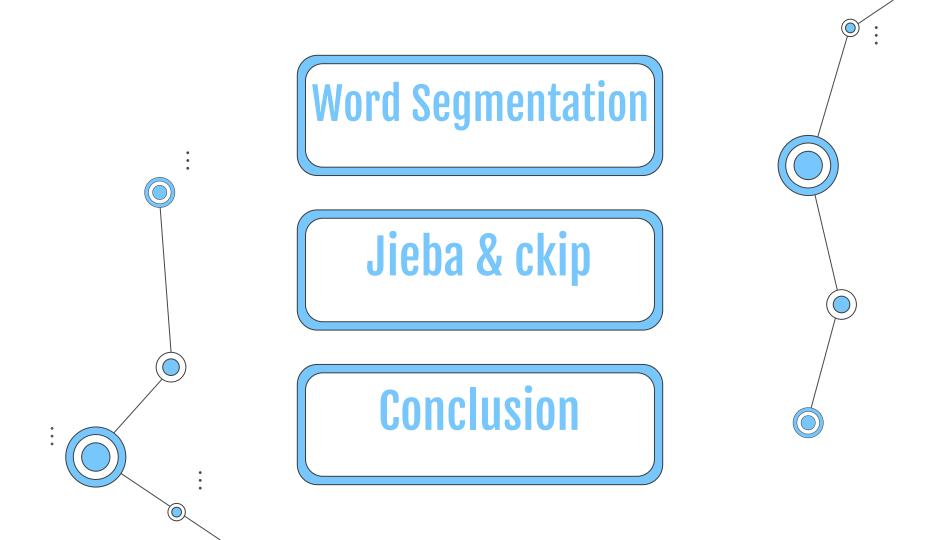
Join at slido.com #1540124

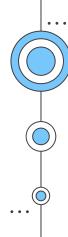
① Start presenting to display the joining instructions on this slide.





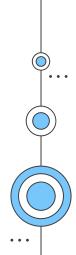
2023/03/09





O1 Quick Introduction

Why do we need word segmentation and POS tagging?





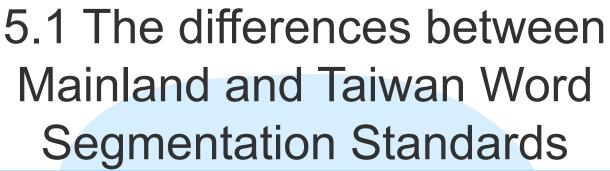
4.1 Word Segmentation Standard

Definition: A segmentation unit is the smallest string of character(s) that has both an independent meaning and a fixed grammatical category. Bound morpheme?

Reduplication?

Basic principles:

- A string whose meaning cannot be derived by the sum of its components should be treated as a segmentation unit. [Combination principle]
- The string whose grammatical category cannot be derived. by the sum of the grammatical categories of its components should be treated as a segmentation unit. [Combination principle]



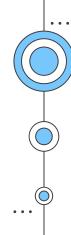
Segmentation units did not equal to words and the target is not the linguistic word, but a processing unit for information processing of Chinese texts.

- The Word Segmentation Standard of Contemporary Chinese Language for Information Processing
 - A word is the smallest element that may be uttered in isolation.
 - A segmentation unit is the smallest element that may
 be adopted in Chinese information processing and still
 have a semantic or syntax function. It includes word and
 phrases in the standard.



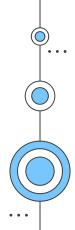
Join at slido.com #1540124

① Start presenting to display the joining instructions on this slide.



How to subcategorize these words?

Behavioral Characteristics!!!





Tokenization



POS tagging



Dependency parsing



NER

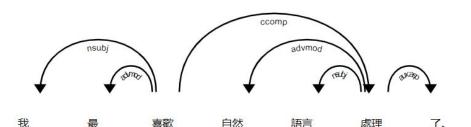
Chinese NLP pipeline in spaCy

下課 NOUN noun NT temporal noun 我 PRON pronoun PN pronoun 要 VERB verb VV other verb 知道 VERB verb VC 是 (copula) 實雅 PROPN proper noun NR proper noun 屈臣氏 PROPN proper noun NR proper noun 他們 PRON pronoun PN pronoun

> token.pos_: coarse-grained pos

token.tag_: fine-grained pos

```
('我', 'PRON', 'PN', 'nsubj', 喜歡)
('最', 'ADV', 'AD', 'advmod', 喜歡)
('喜歡', 'VERB', 'VV', 'ROOT', 喜歡)
('自然', 'ADV', 'AD', 'advmod', 處理)
('語言', 'NOUN', 'NN', 'nsubj', 處理)
('處理', 'VERB', 'VV', 'ccomp', 喜歡)
('了', 'PART', 'SP', 'aux:asp', 處理)
('。', 'PUNCT', 'PU', 'punct', 喜歡)
('所以', 'ADV', 'AD', 'advmod', 課)
('秒選', 'PROPN', 'NR', 'name', 謝)
('謝', 'PROPN', 'NR', 'dep', 課)
('舒凱', 'PROPN', 'NR', 'compound:nn', 老師)
('老師', 'NOUN', 'NN', 'nmod:assmod', 課)
('的', 'PART', 'DEG', 'case', 老師)
('課', 'NOUN', 'NN', 'ROOT', 課)
('!', 'PUNCT', 'PU', 'punct', 課)
```



ADV

NOUN

VFRB

PART

PRON

ADV

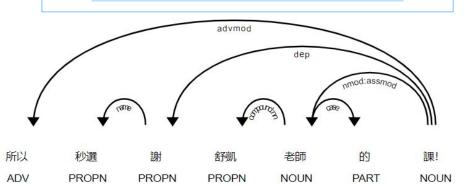
VFRB

Chinese NLP pipeline in spaCy



Dependency parsing

token 對這個 head 有這樣的依存關係





To what extent, can we create new verbs?

① Start presenting to display the poll results on this slide.



To what extent, can we create new prepositions?

① Start presenting to display the poll results on this slide.



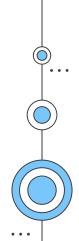
To what extent, can we create new nouns?

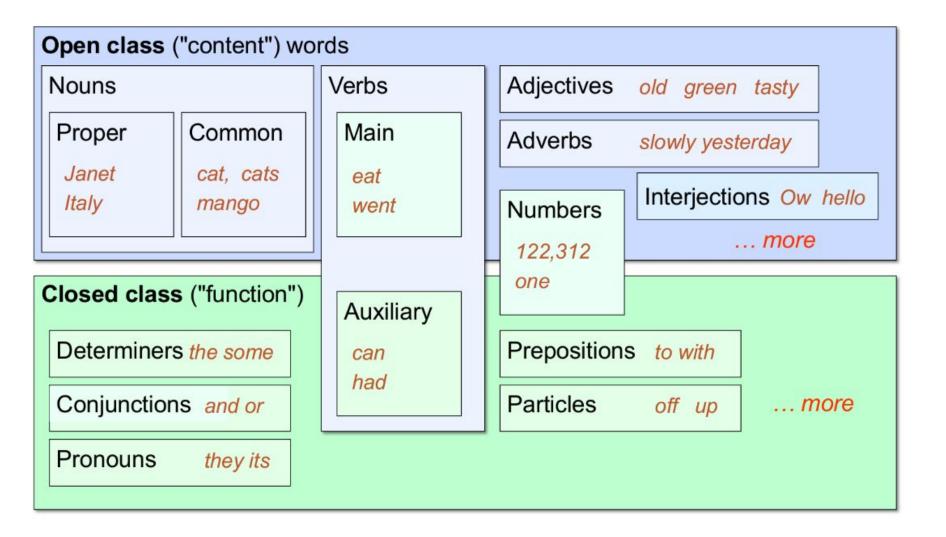
① Start presenting to display the poll results on this slide.



Open class vs. Closed class

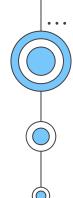
- Open class words
 - Usually content words: nouns, verbs, adjectives, adverbs ...
 - New words like iPhone or to fax
- Closed class words
 - Usually function words with grammatical function: determiners, pronouns, prepositions ...
 - Relatively fixed membership



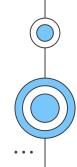


"Universal Dependencies" Tagset

	Tag	Description	Example
6	ADJ	Adjective: noun modifiers describing properties	red, young, awesome
Open Class	ADV	Adverb: verb modifiers of time, place, manner	very, slowly, home, yesterday
D D	NOUN	words for persons, places, things, etc.	algorithm, cat, mango, beauty
Sen	VERB	words for actions and processes	draw, provide, go
Ō	PROPN	Proper noun: name of a person, organization, place, etc	Regina, IBM, Colorado
J.	INTJ	Interjection: exclamation, greeting, yes/no response, etc.	oh, um, yes, hello
	ADP	Adposition (Preposition/Postposition): marks a noun's	in, on, by under
S		spacial, temporal, or other relation	
Words	AUX	Auxiliary: helping verb marking tense, aspect, mood, etc.,	can, may, should, are
≥	CCONJ	Coordinating Conjunction: joins two phrases/clauses	and, or, but
Closed Class	DET	Determiner: marks noun phrase properties	a, an, the, this
\Box	NUM	Numeral	one, two, first, second
sed	PART	Particle: a preposition-like form used together with a verb	up, down, on, off, in, out, at, by
65	PRON	Pronoun: a shorthand for referring to an entity or event	she, who, I, others
	SCONJ	Subordinating Conjunction: joins a main clause with a	that, which
		subordinate clause such as a sentential complement	
ii.	PUNCT	Punctuation	; , ()
Other	SYM	Symbols like \$ or emoji	\$, %
0	X	Other	asdf, qwfg

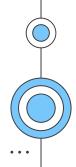


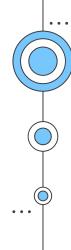
- Can be useful for other NLP tasks
 - Improve syntactic/dependency parsing
 - Help ML in reordering of constituents
 - Sentiment analysis may want to distinguish adjectives or other POS
 - Text-to-speech needs to contrast between homonyms/allophones etc. (e.g. 兒的生活好痛苦一點也沒有糧食多病少掙了很多錢)





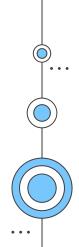
- And also useful for linguistic or language-analytic computational tasks
 - Control for POS when studying linguistic change (creation of new words, or meaning shift)
 - Control for POS in measuring meaning similarity or difference

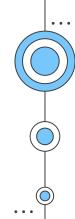




O2 Jieba & ckip

Python modules for word segmentation and POS tagging





- If you haven't installed the modules
 - To install jiebapip install jieba
 - To install ckip-transformers

pip install -U ckip-transformers





To load module

```
# dependencies
import jieba
import jieba.posseg as pseg
import logging
jieba.setLogLevel(logging.INFO)
import ckip transformers
from ckip_transformers.nlp import CkipWordSegmenter, CkipPosTagger
```



Jieba



POS tagging

```
text_jb = jieba.cut(text)

print(text)
print('\t'.join(text_jb))
```

自然

我最喜歡自然語言處理了我 最 喜歡



Jieba



POS tagging

```
words = pseg.cut(text)

print(text)
print('\t'.join(text_jb))
print('\t'.join([f'{word}/{tag}' for word, tag in words]))
```

我最喜歡自然語言處理了



ckip

```
02
```

POS tagging

```
texts = text.strip().splitlines()
ws_texts = ws_driver(texts)
pos texts = pos driver(ws texts)
```

```
# with GPU
ws_driver = CkipWordSegmenter(device=0)
pos_driver = CkipPosTagger(device=0)

# no GPU
# ws_driver = CkipWordSegmenter(device=-1)
# pos_driver = CkipPosTagger(device=-1)
```

```
for sentence, sentence_ws in zip(texts, ws_texts):
    print(sentence)
    print('\t'.join(sentence_ws))
```

```
我最喜歡自然語言處理了
我 - 喜歡 自然 語言 處理 う
```



ckip



POS tagging

```
for sentence, sentence_ws, sentence_pos in zip(texts, ws_texts, pos_texts):
    print(sentence)
    print('\t'.join(sentence_ws))
    print('\t'.join([f'{word}/{tag}' for word, tag in zip(sentence_ws, sentence_pos)]))
```

我最喜歡自然語言處理了

```
    我
    最
    喜歡
    自然
    語言
    處理
    了

    我/Nh
    最/Dfa
    喜歡/VK
    自然/Na
    語言/Na
    處理/VC
    了/Di
```



Reference

Speech and Language Processing

