Course 2: Tokenization

What is tokenization?

Turning text...

```
I love playing soccer!
```

...into tokens

```
['I', 'love', 'play', 'ing', 'soccer', '!']
```

Course 2: Tokenization 2

Historical Notions

Tokenization Origins

The concept comes from linguistics

" non-empty contiguous sequence of graphemes or phonemes in a document

 \approx

split on blanks

"

Tokenization Origins

```
old_tokenize("I love playing soccer!") = ['I', 'love', 'playing', 'soccer!']
```

- Different from word-forms !
 - damélo → da/mé/lo (=give/me/it)

Course 2: Tokenization 5

Tokenization Origins

Natural language is split into...

- Sentences, utterances, documents... (*macroscopical*) that are split into...
 - Tokens, word-forms... (*microscopical*)
- → Used for linguistic tasks (POS tagging, syntax parsing,...)

6

Tokenization & ML

Machine Learning relies on tokenization:

- Gives better performance
- Fixed-size vocabulary often required

Tokenization & ML

→ Need for a good downsampling mechanism

8