

spaCy

Présentation par Morgann Sabatier

Inspirée par :

- Formation spaCy.io : **spaCy 101 : Everything you need to know**
- **La documentation spaCy – spacy.io**
- Tutoriel realpython.org : **Natural Language Processing With spaCy in Python**
- **Traitement Automatique du Langage Naturel en français (TAL / NLP)** par Maël Fabien

Objectifs

01.

COMPRENDRE

02.

PRÉSENTER

03.

NUANCER

2015

Matthew Honnibal et Ines Montani



Librairie Python **open-source**
Outils flexibles



Licence MIT – logiciel libre



65 modèles, 17 langues



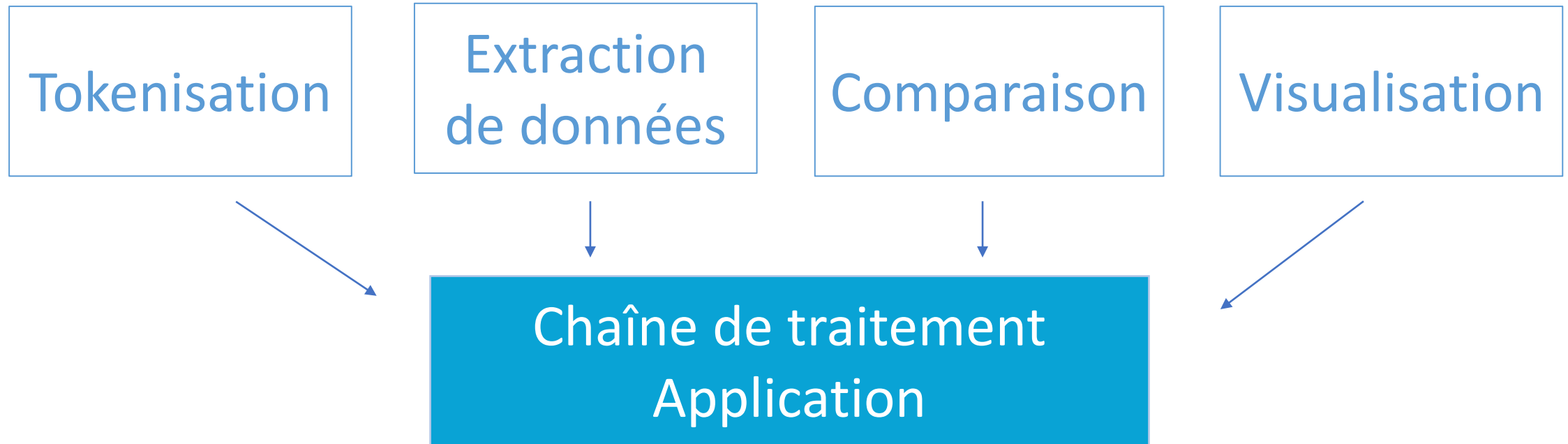
Communauté active

spaCy,
c'est quoi ?

Pour quoi faire ?

NLTK → Enseignement et recherche

spaCy → production et application



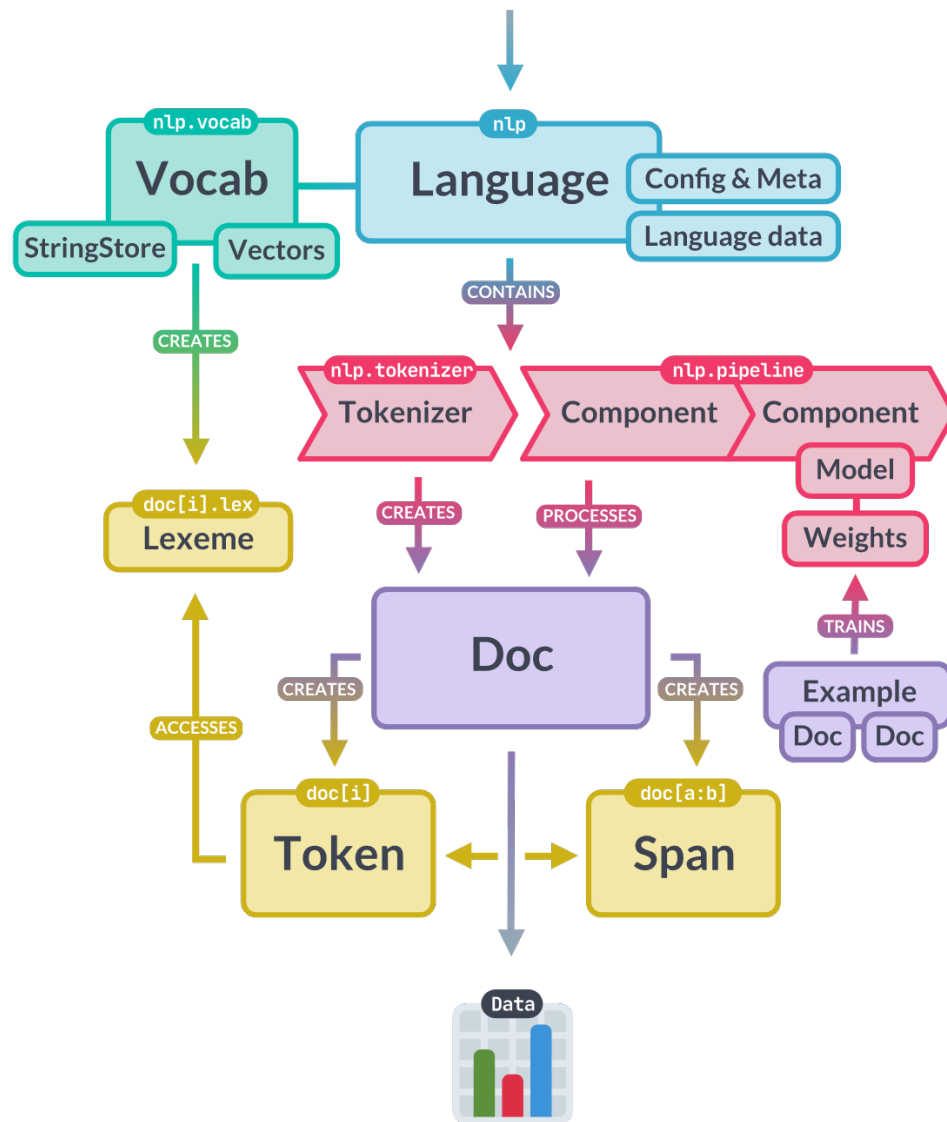
Comment ça marche ?

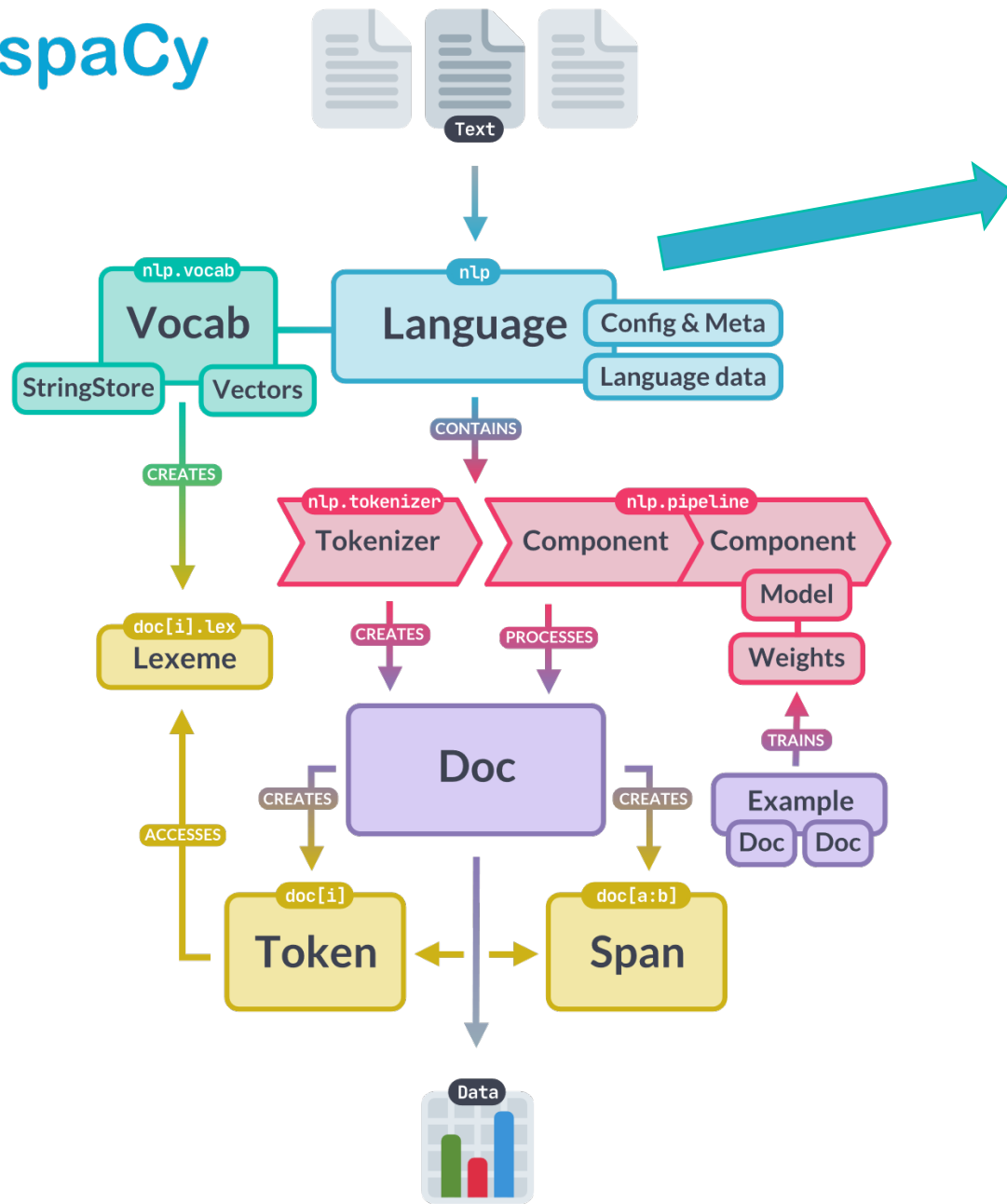
- Architecture
- Pipeline
- Format de données
- Modèles



`text = "Morgann explains how spaCy works"`

- Chaines de caractères
- Documents
- Corpus

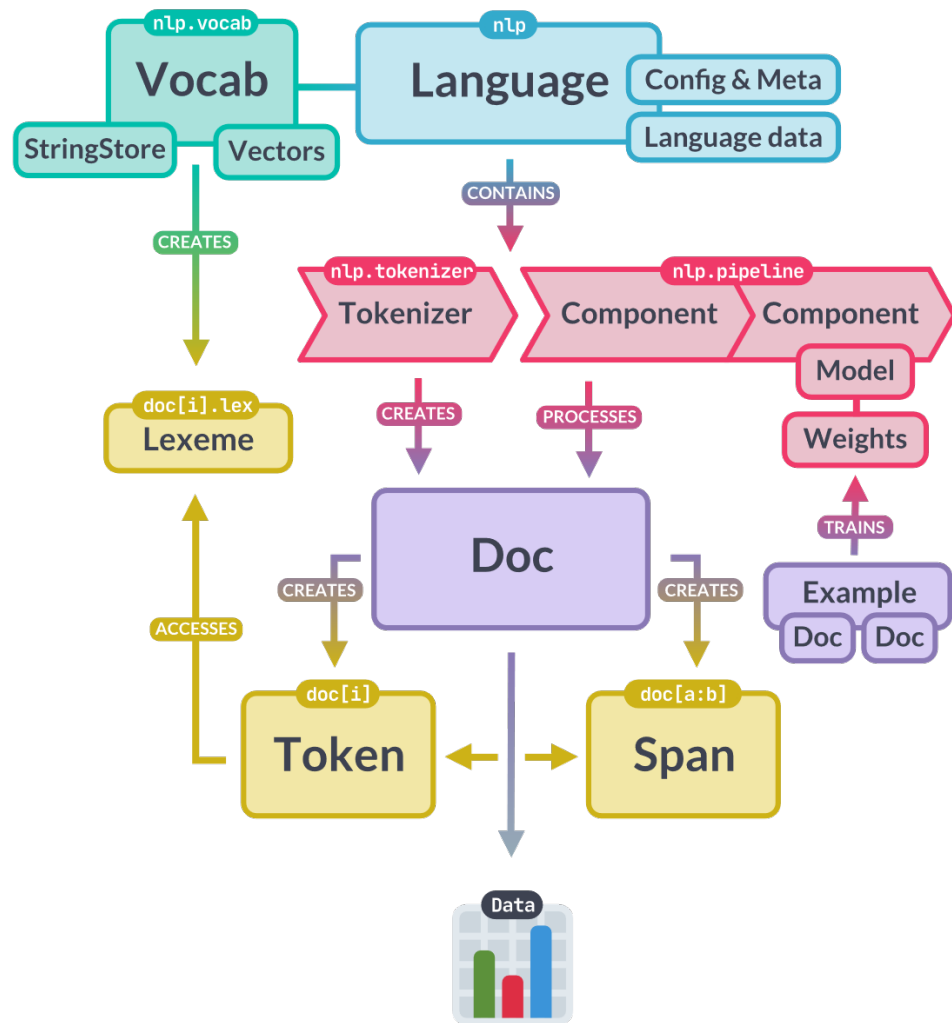




```
text = "Morgann explains how spaCy works"
```

```
nlp = spacy.load("en_core_web_sm")
```

- Vocabulaire du modèle
 - Vecteurs + strings
- Lexèmes

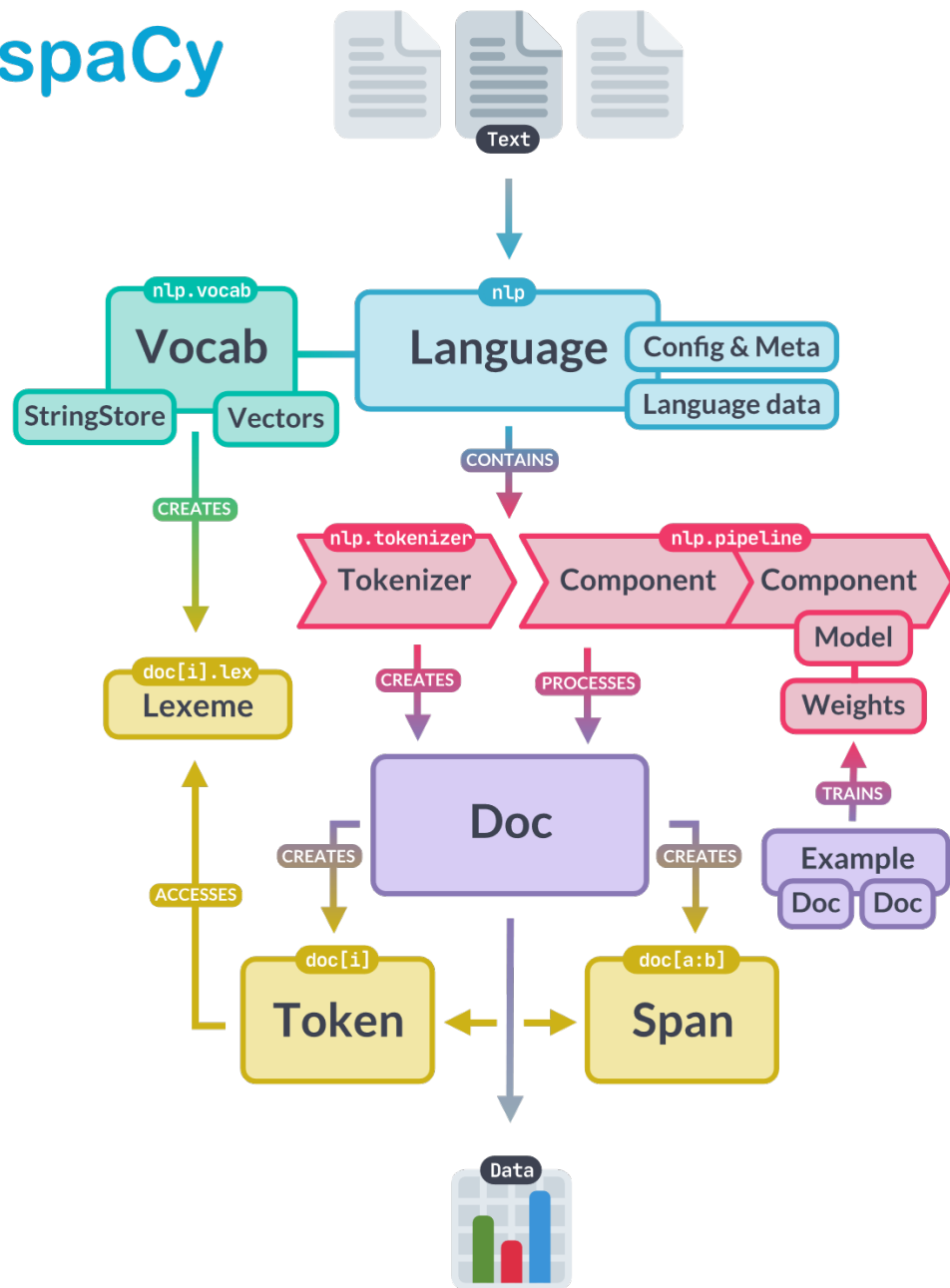


```
doc = nlp(text)
```

ANNOTATION

```
Doc[0].dep = nsubj          #Sujet
```

```
[ (ent.text, ent.label_) for ent in doc.ents]
>>> [ ('Morgann', 'PERSON') ]
```



```
text = "Morgann explains how spaCy works"
```

```
nlp = spacy.load("en_core_web_sm")
```

```
doc = nlp(text)
```



TOKENISATION



ANNOTATION

```
Doc[0] = Morgann #Token
```

```
Doc[0].pos_ = PROPN #Nom propre
```

```
Doc[0].dep_ = nsubj #Sujet
```

```
[(ent.text, ent.label_) for ent in doc.ents]
>>> [('Morgann', 'PERSON')]
```



ENRICHISSEMENT

```
spacy_ent = Span(doc, 3, 4, label="LIB")
doc.ents = list(doc.ents)+[Spacy_ent]
>>> [('Morgann', 'PERSON'), ('spaCy', 'LIB')]
```

Fonctionnalités

Quelles fonctionnalités ?

- Tokenisation
- Token attributes
- Entités nommées
- Matcher

Tokenisation

spacy.tokens.doc.Doc VS list NLTK

NLTK	<code>['d', '', 'enseignants'], [l'encouragement]</code> <code>['O.N.U', '.']</code> <code>['#', 'France']</code>
spaCy	<code>['d', 'enseignants'], [l', 'encouragement']</code> <code>['O.N.U.']</code> <code>['#', 'France']</code>

Et pour le mandarin ?

NLTK	<code>['我想自我介绍一下']</code>
spaCy	<code>['我', '想', '自我', '介绍', '一下']</code>

Et les langues agglutinantes ?

Token attributes

- **Extraire** des informations

.lemma_	.tag_	.morph	Spacy.explain(token.tag_)
we	PRP	Number=Plur Person=1 PronType=Prs	Pronoun, personal

- **Interroger** un type

Token.text	token.is_alpha	Token.is_punct	Token.shape_	Token.is_stop
Name	True	False	Xxxx	True

Stop Words

NLTK, mots vides :

`il', 'ils', 'je', 'la',
'le', 'les', 'leur',
'lui', 'ma', 'mais'`

spaCy, des mots **vraiment** vides ?

`'différents', 'via',
'nombreuses', 'néanmoins',
'certains', 'toujours'`

Limite des modèles

Entités nommées

Reconnaissance d'entités nommées + visualisation

Performant en anglais

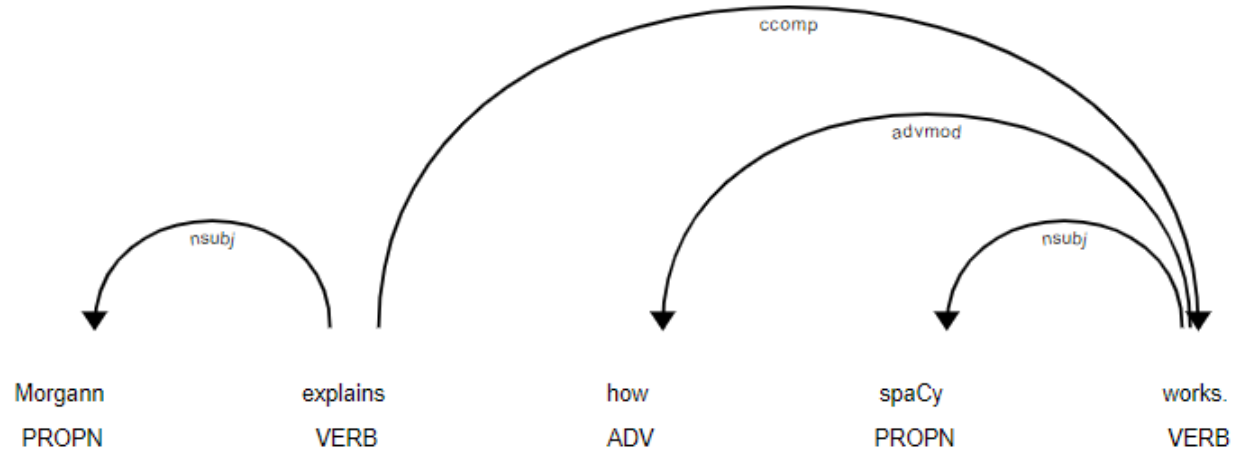
Marie Curie PERSON, née Maria Sklodowska PERSON, was born in Warsaw GPE on November 7, 1867 DATE

Résultats **dépendants** de l'annotation corpus modèle

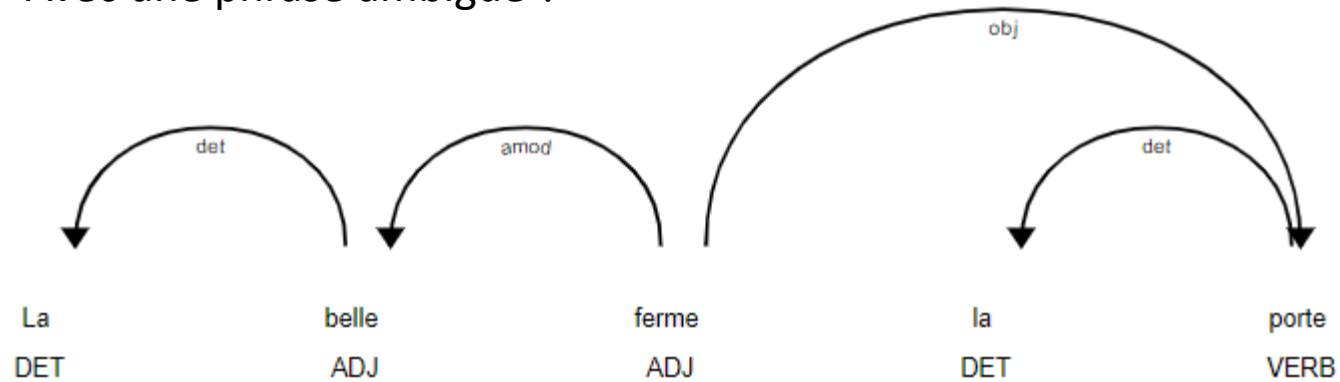
Maria Sklodowska PER naît le 7 novembre 1867 à Varsovie LOC,

Гарри Кíмович Каспáров PER (фамилия при рождении Вайнштéйн; род. 13 апреля 1963, Баку LOC, Азербайджанская ССР LOC),

Dependencies

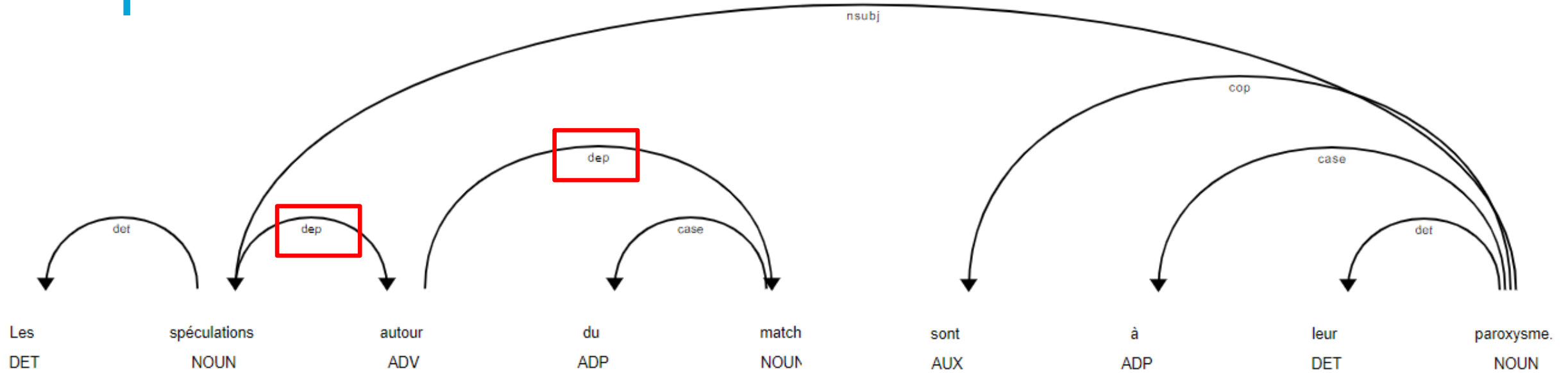


Avec une phrase ambiguë ?

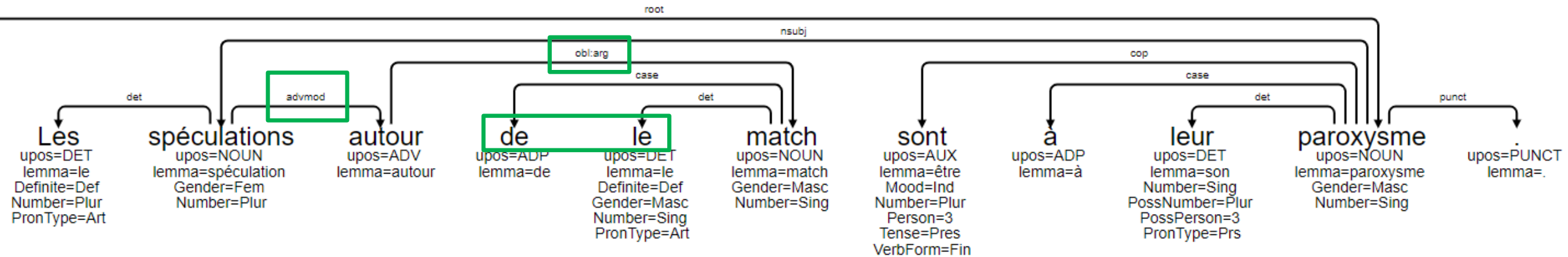


```
displacy.render(doc_fr, style="dep")  
spacy.explain("dep") >>> 'adjectival modifier'
```

Dependencies

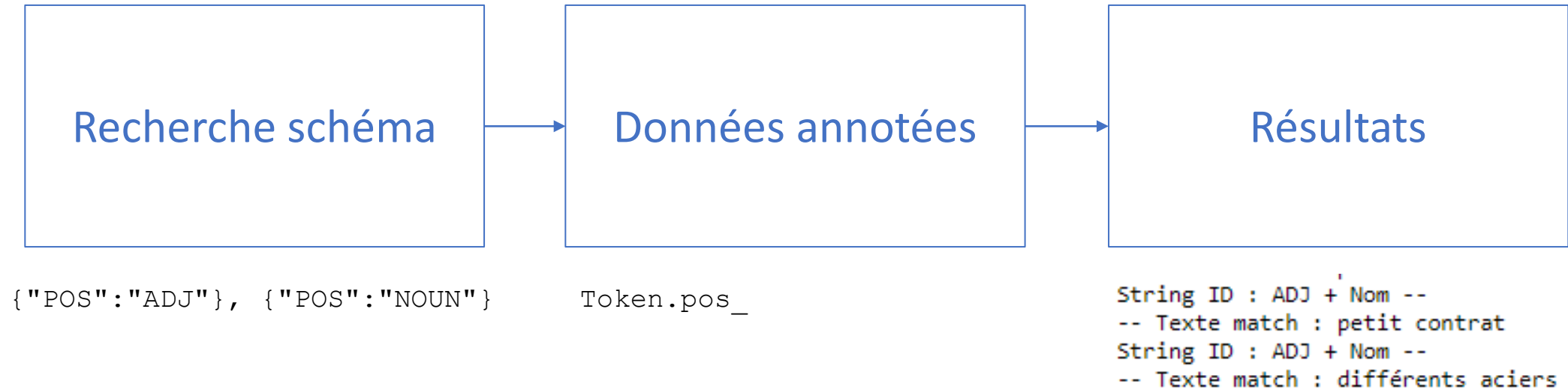


spaCy visualisation displacy



GREW Match

Matcher



The background is a solid blue color with various abstract geometric elements. In the top-left corner, there is a grid of small white dots. In the top-right corner, there is a circular pattern of small white dots. In the bottom-left corner, there is a large white circle. In the bottom-right corner, there is a grid of small white dots. There are also several white plus signs and small white dots scattered across the background.

Démo

Bilan

Avantages

- + Simple d'accès
- + Communauté active
- + Performant
- + Customisable
- + Open source et licence libre

Inconvénients

- Modèles hors anglais peu performants
- Pas de langues peu dotées
- Puissance des outils limités
- Stopwords
- Ambiguïté

The background is a solid blue color with various abstract geometric elements. In the top-left corner, there is a grid of small white dots. In the top-right corner, there is a circular pattern of small white dots. In the bottom-left corner, there is a large white circle. In the bottom-right corner, there is a grid of small white dots. There are also several white lines and dots scattered across the background. The text "Merci de votre attention" is centered in the middle of the image in a large, white, sans-serif font.

Merci de votre
attention

Je mets à votre disposition

- La présentation en format PDF
- Un notebook avec des exemples simples et en différentes langues de nombreux outils de spaCy
- Le notebook de la démo présentée

Quelques ressources :

- <http://datacamp-community-prod.s3.amazonaws.com/29aa28bf-570a-4965-8f54-d6a541ae4e06> - Cheatsheet fonctions
- <https://betterprogramming.pub/extract-keywords-using-spacy-in-python-4a8415478fbf> - Trouver Keywords pipeline
- <https://towardsdatascience.com/named-entity-recognition-ner-using-spacy-nlp-part-4-28da2ece57c6> - NER