

## ① PROCESAR TEXTO CON TIKKA

QUEREMOS CON:

- TÍTULO (TITLE)
- ⊕ - DESCRIPTION
- X-TIKA: content: Nos podemos quedar con las frases que tengan longitud  $\geq 20$  (para quitar todo lo que no nos interesa)

~~TIKA~~ →

tika-reader.py

tika-processor.py

~~Process-text.py~~

## ② SACAR VECTORES

## ③ APLICAR ALGORITMO CONGLOMERADOS

## ④ ENTIDADES NOMBRADAS

- SPACY → APLICAR A ⊕
- STANFORD
- CoNLL2002 CHUNKER (BIGRAMAS) | ⊕

ner-stanford.py

spacy-exercise2.py

ner-train-chunker.py

## TRATAMIENTO DEL TEXTO

### Ⓐ Stop words

Puntuación

lematizar

Process-text.py

- ### Ⓑ Traducción de textos de inglés a español y español a inglés (funciona mejor un diccionario que otro en un idioma que en otro?)

Textbleb