



06/05/2021

L'emprunt à l'anglais dans le journal L'Est Républicain (1999-2011)

Projet du Dictionnaire et Néologisme



Yuyan QIAN



M1 Langue et Informatique, Sorbonne université



Plan



1. Introduction



2. Matériel



3. Méthodologies



4. Résultats



5. Limites



5. Perspectives



Introduction

Pourquoi ?

Marketing : « Les chercheurs sont parvenus à des mesures de bien-être beaucoup plus prédictives en matière de fidélité »

L'évolution de l'emprunt à l'anglais dans les journaux français :

- plus fréquent ou au contraire en diminution ?
- le pourcentage d'emprunts linguistiques par rapport à celui de la néologie en français ?



L'EST
RÉPUBLICAIN

Nancy

ER La start-up nancéienne Abby lève 1,2 million d'euros

L'entreprise fondée à Nancy par quatre étudiants du CESI est spécialisée dans l'accompagnement des ...



Matériel



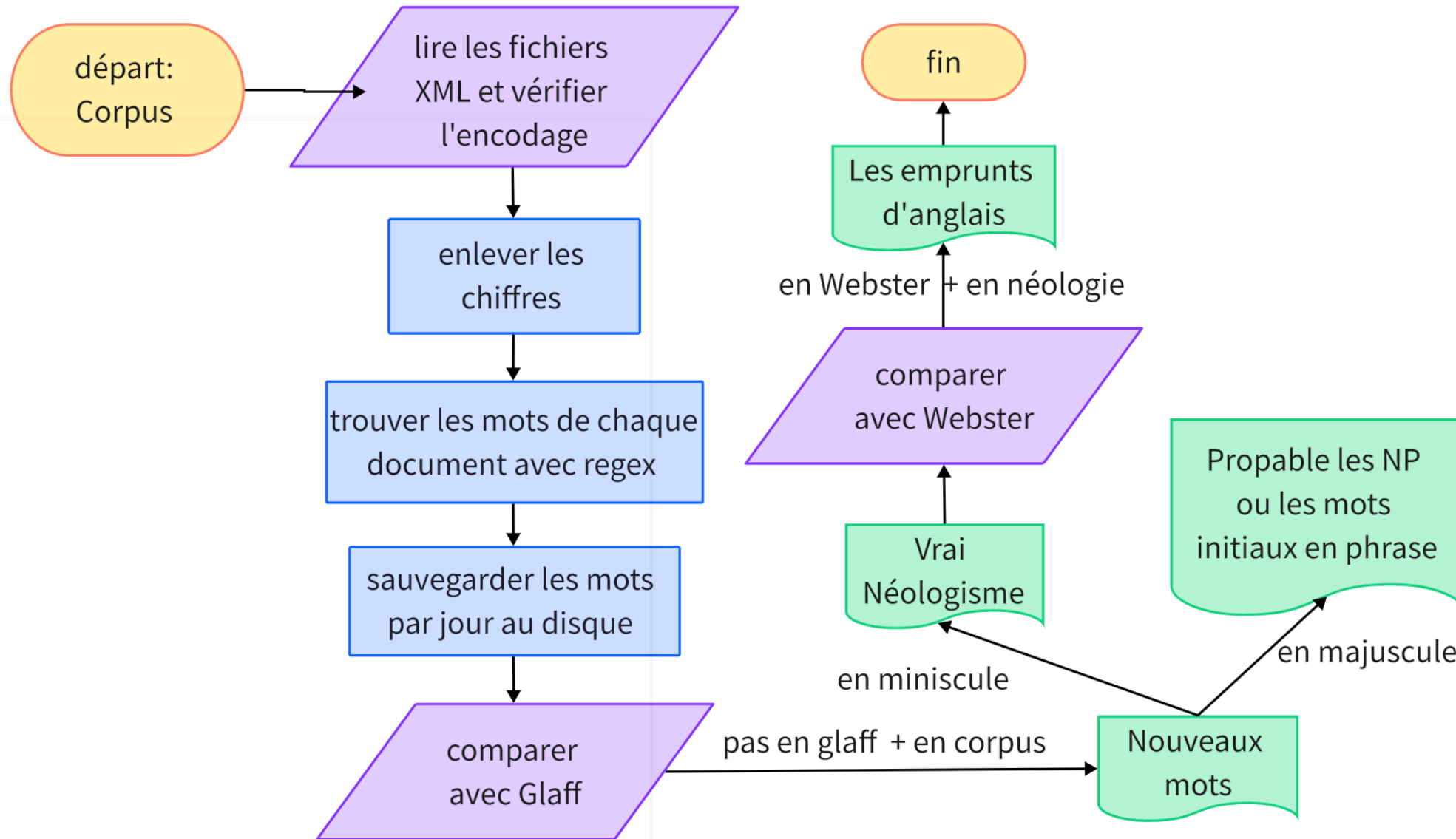
Corpus	Collection	Le premier jour	Le dernier jour	Taille
Annee1999	132 jours	1999-05-17	1999-09-30	249,4 MB
Annee2002	169 jours	2002-01-02	2002-12-27	318,5 MB
Annee2003	54 jours	2003-01-02	2003-02-24	101,3 MB
Annee2006	357 jours	2006-02-01	2006-12-31	1,89 GB
Annee2008	361 jours	2008-02-01	2008-12-31	1,72 GB
Annee2010	170 jours	2010-02-01	2010-05-31	672,2 MB
Annee2011	13 jours	2011-02-01	2011-02-14	46,7 MB

Dictionnaires

Corpus journalistique issu de l'Est Républicain

1. Glaff (1, 406, 857 entrées) - un lexique du français construit à partir du Wiktionnaire, la branche francophone de Wiktionary, partagé dans la licence CC BY-SA 3.0.
2. Webster (86, 036 entrées) - un texte brut dans le fichier dictionary.txt. Son contenu est soumis aux termes de la licence du Projet Gutenberg.

Ex: "BAHAISM": "The religious tenets or practices of the Bahais."





Traitements



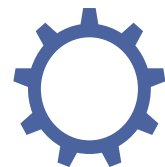
Vérification d'encodage

```
import chardet
for path_file in glob.glob("Annee2008/*"):
    with open(path_file, 'rb') as f:
        result = chardet.detect(f.read())
    if result['encoding'] != 'utf-8':
        print(path_file)
        #print(result['encoding'])
```

Par exemple : certains alias en byte

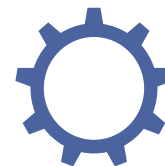


Lecture avec BeautifulSoup



Prétraitement

2. Lang (Sui) à 21'' ; 3. Taramarcas (Sui) à 30''
 (Sui) à 29'' ; 3. Brandau (All) à 37''... 8. Colin ('



Tokenisation

```
re.findall("[A-ZÀ-ÿ][A-ZÀ-ÿ0-9.]+" |
           "[cCdDjJlLmnNsSt]" |
           "[qQ]u'|\w+[\w'-]+\w+" |
           "\w\w+)", texte)
```

Par exemple : BW's, contre-productive,
publicité-marketing



Sauvegarde



Réécriture en json

```
try:
    with open(path, 'r') as f:
        old_dict = json.load(f)
except:
    old_dict = {}
new_dict = {**old_dict, **dico}
```



Glaff et Webster

Différence et intersection

Mise à jour de vocabulaire néo



Vérification

Division les tokens en deux groupes :
capitalisé ou pas



Connexion des documents

Actualisation des néologismes pour
chaque année



Visualisation

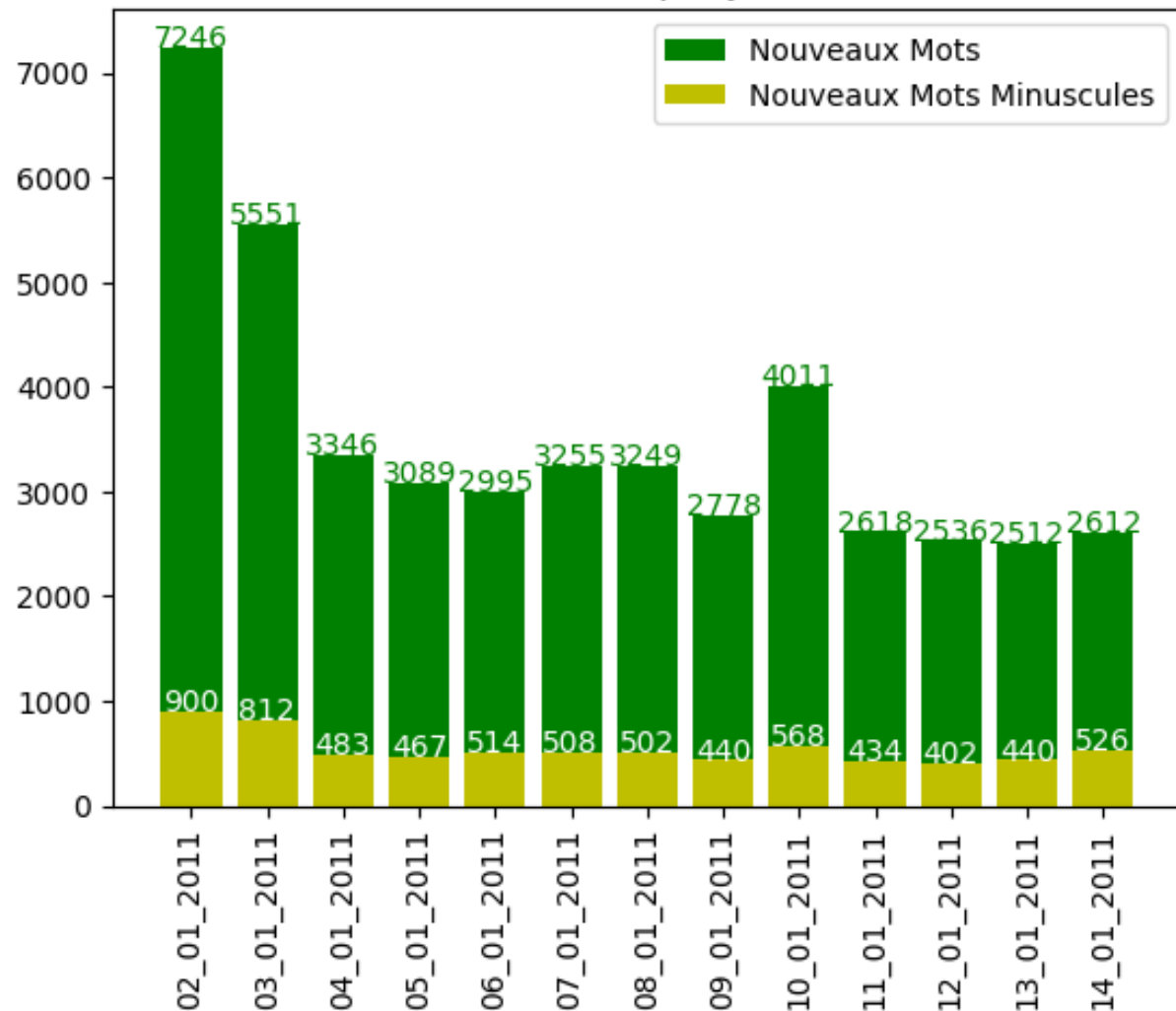
Matplotlib - histogramme et courbe

Wordcloud - nuage de mots

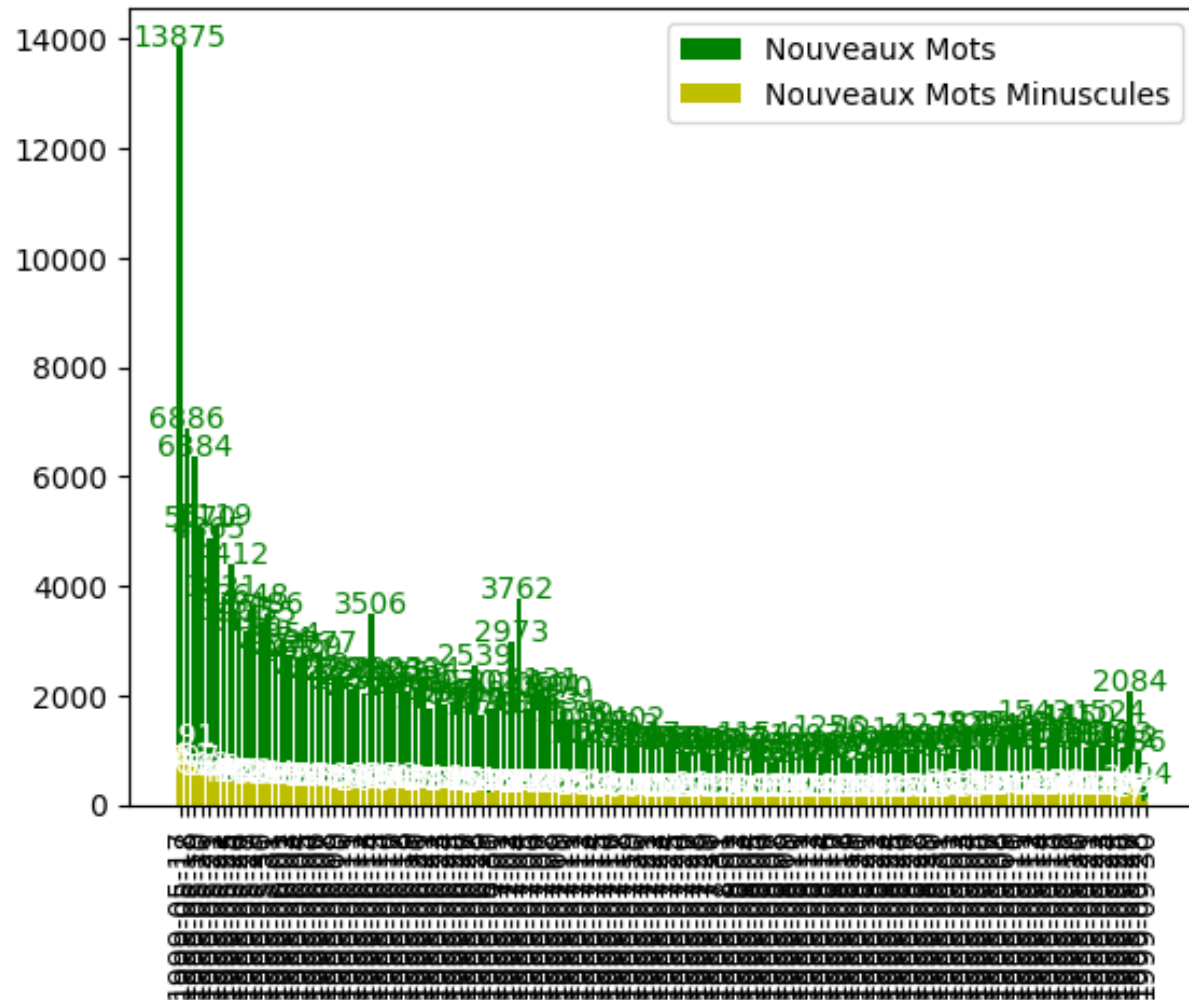


Résultats

NB de nouveaux Mots par jour (minuscules)



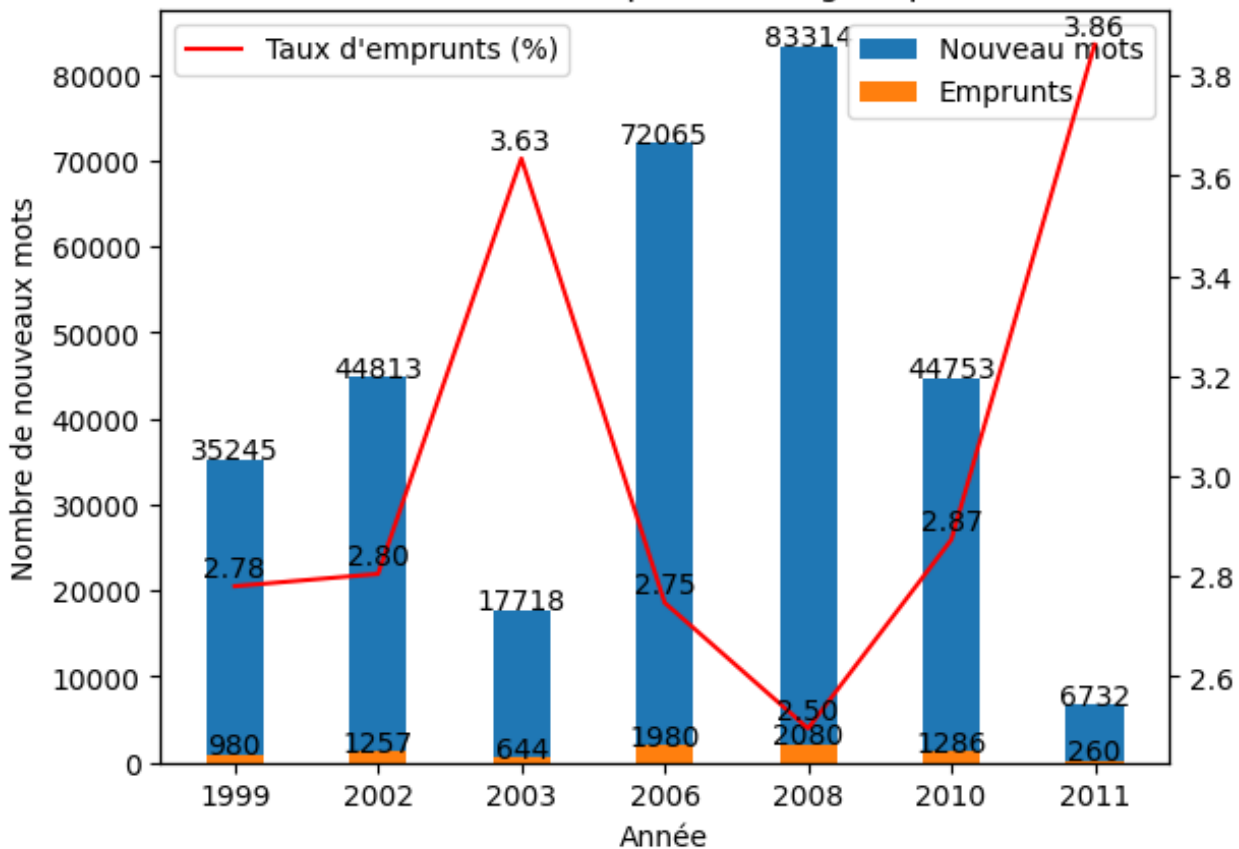
NB de nouveaux Mots par jour (minuscules)



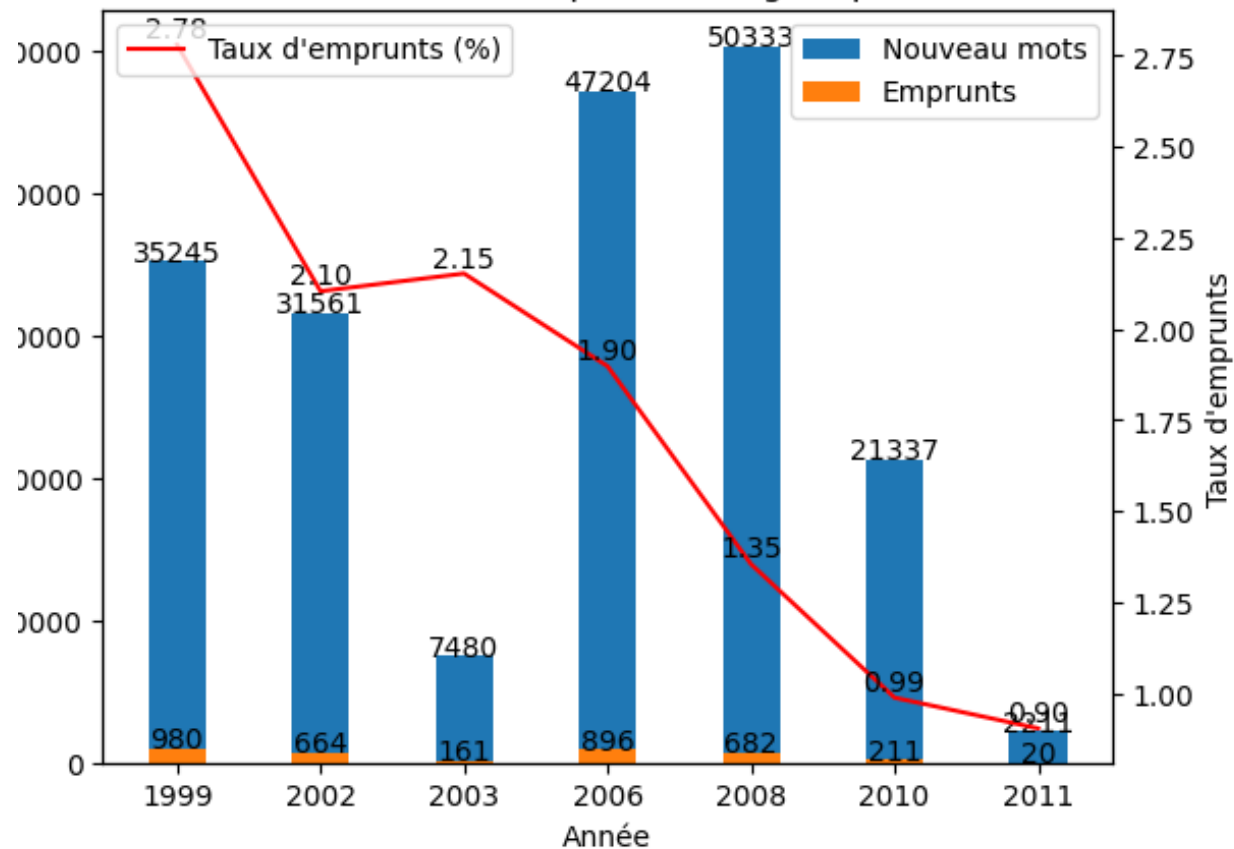


Résultats

Nouveaux mots et emprunts d'anglais par année

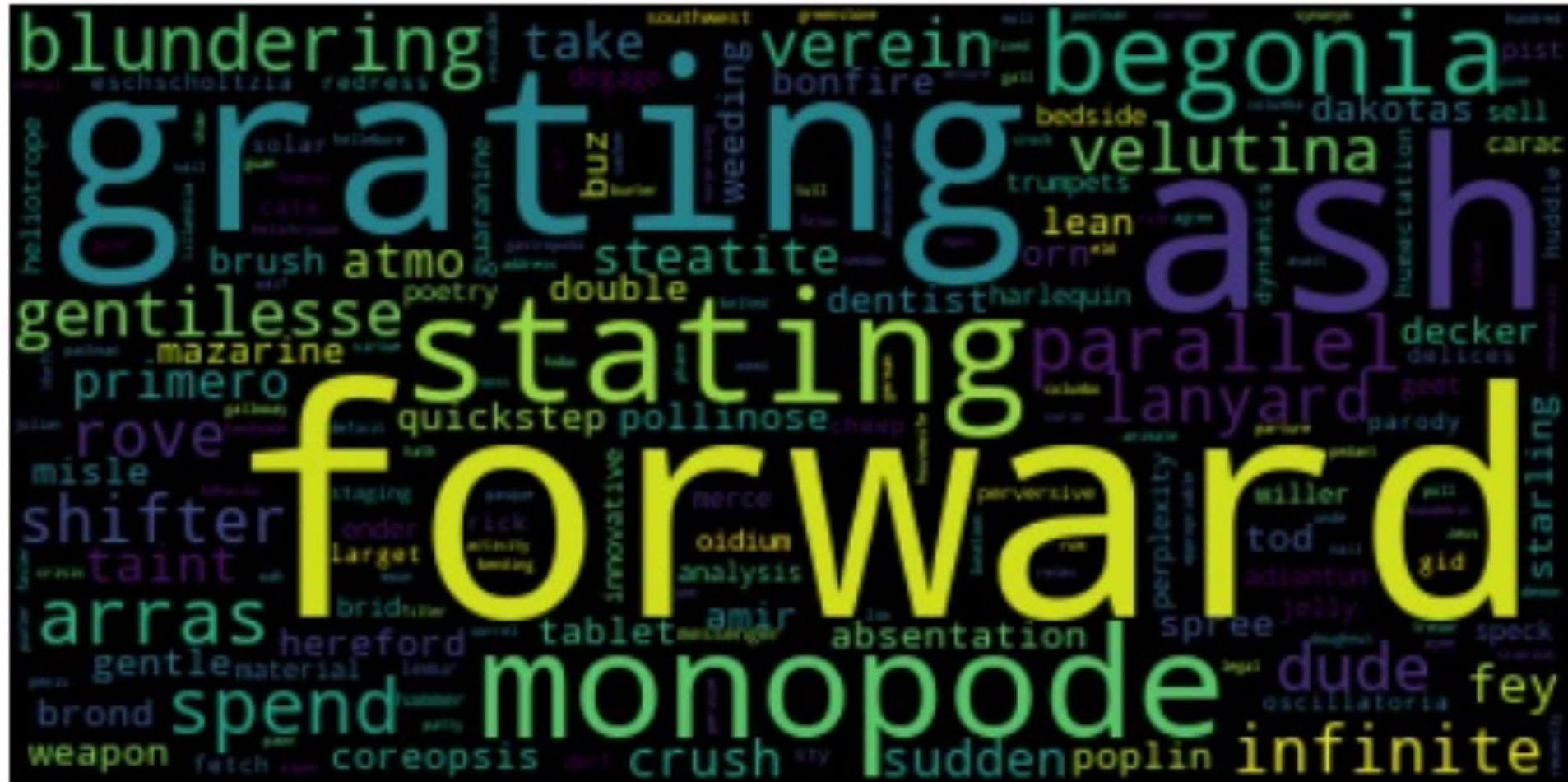


Nouveaux mots et emprunts d'anglais par année





Résultats





Limites

01

Les données du corpus

02

La mélange des vrais néologismes et des NP

03

Les manques d'entrées dans deux dictionnaires

04

La taille déséquilibre des données annuelles

```
<div type="article">
  <div type="texte">
    <div>
      <div>
        <head type="h4">Le Dakar, c'est parti
        <p>Après les vérificationstechniques,
      </div>
    </div>
  </div>
</div>
```

EX: L'histoire de Tom-Tom

EX: Glaff ne contient pas les mots-composés

La fréquence sert à réduire l'impact de la quantité de données sur les résultats.





Références

[1] ATILF and CLLE. Corpus journalistique issu de l'est républicain, 2020. ORTOLANG (Open Resources and TOols for LANGuage) – www.ortolang.fr.

[2] Denis Jamet and Adeline Terry. Les néologismes anglais issus de l'emprunt : l'étude diachronique. ELAD-SILDA, 1, 2018.



Merci !

