

# AGODA

ANALYSE SÉMANTIQUE ET GRAPHES RELATIONNELS POUR L'OUVERTURE ET L'ÉTUDE DES DÉBATS À L'ASSEMBLÉE NATIONALE

---

Marie Puren (Epitech)  
Pierre Vernus (LARHRA)

Charles Paul Renouard  
*Les Présidents Deschanel et Waldeck-Rousseau à la Chambre  
des Députés*  
Musée d'Orsay  
Sans date (entre 1898 et 1902 ?)





Source : BnF, <https://gallica.bnf.fr/html/und/droit-economie/iiie-republique-0?mode=desktop>

## OBJECTIFS



- Donner accès aux débats parlementaires
- Faciliter la recherche dans le corpus
- Offrir de nouveaux modes de visualisation des documents



# OBJECTIFS

Projet en mode « preuve de concept » sur une partie réduite du corpus (1889-1893)



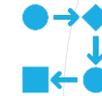
## Exploration

Créer une plateforme de consultation des **débats parlementaires à la Chambre des députés** (1881-1940), retranscrits dans le *Journal officiel de la République française. Débats parlementaires. Chambre des députés : compte rendu in-extenso*, et disponibles sur Gallica



## Annotation

Produire des données textuelles structurées et sémantiquement enrichies à partir de ces débats numérisés



## Réutilisation

Contribuer à la conception d'un *workflow* adapté à l'analyse de gros corpus de documents historiques (initié par les travaux réalisés au cours de l'ANR TIME US 2018-2021)



## COMPOSITION DE L'EQUIPE

Pierre Vernus – LARHRA

Marie Puren – Epitech (MNSHS)

Nicolas Bourgeois – Epitech (MNSHS)

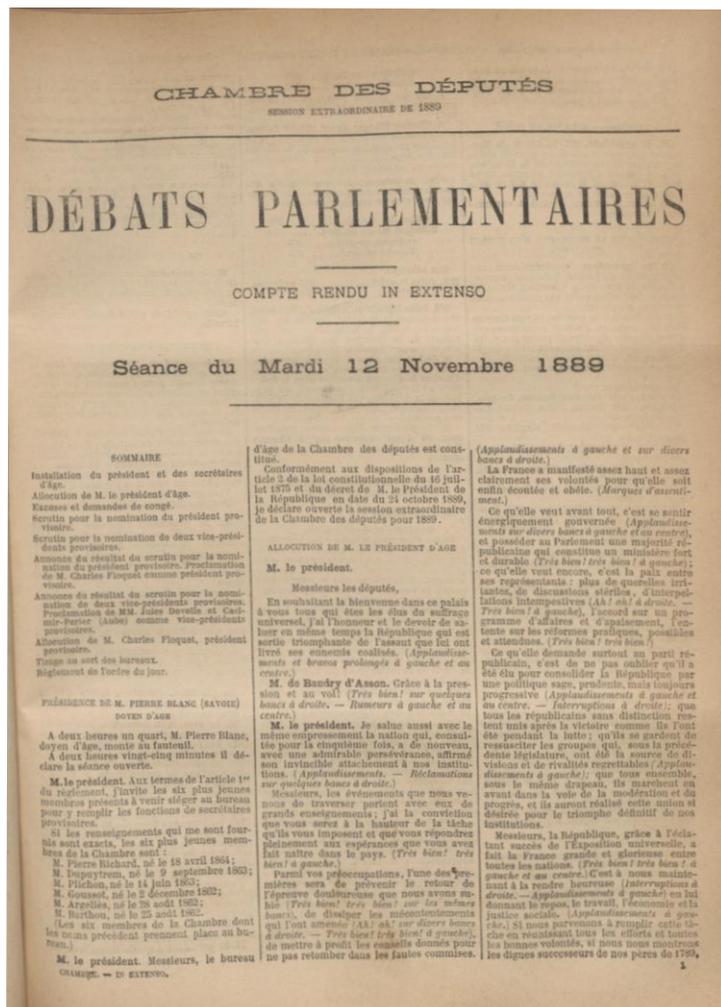
Eric de La Clergerie – Inria (ALMAAnaCH)

**LARHRA**  
UMR 5190

LABORATOIRE DE RECHERCHE  
HISTORIQUE **RHÔNE-ALPES**

*Inria*  
inventeurs du monde numérique

{ EPITECH. }



<https://gallica.bnf.fr/ark:/12148/cb328020951/date>

# LES DÉBATS PARLEMENTAIRES

- ❑ Disponibles sur Gallica : images + textes OCRisés
- ❑ Source intéressante pour :
  - ❑ Histoire
  - ❑ Science politique
  - ❑ Sociologie
  - ❑ Linguistique ...
- ❑ Mise en ligne du Hansard (retranscription des débats parlementaires des gouvernements de type Westminster) a favorisé de nouvelles recherches



**GÉNÉRAL BOULANGER**  
Né à Rennes, le 29, Avril 1837.  
Décédé à Bruxelles, le 30, Septembre 1891.

Général Boulanger - KU Leuven, Belgium -  
Public Domain.

<https://www.europeana.eu/fr/item/2024903/photography> ProvidedCHO KU Leuven 9983157830101488



Bibliothèque nationale de France  
<https://gallica.bnf.fr/ark:/12148/bpt6k1206804x/f1>



Wikimedia Commons

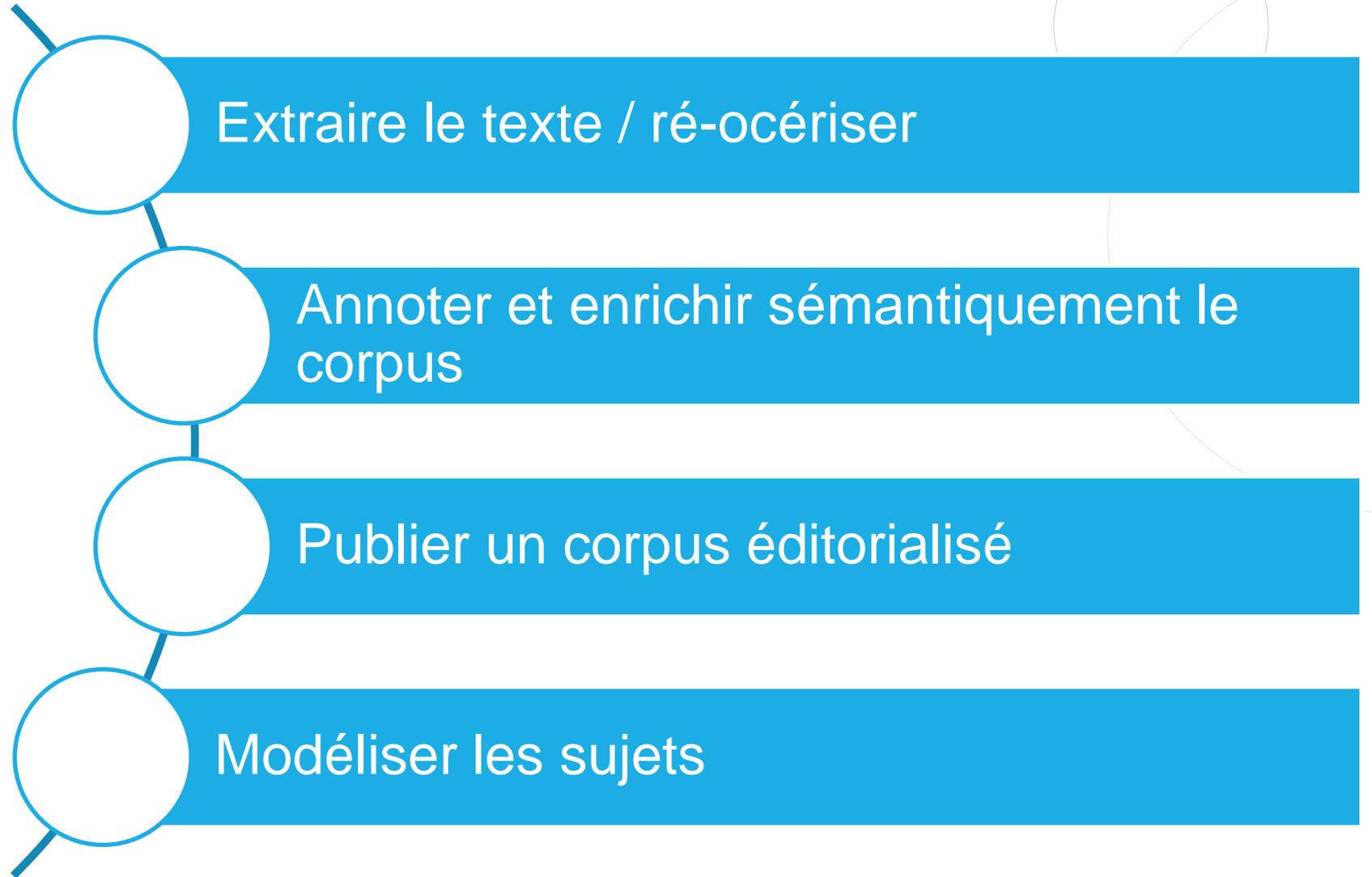
<https://commons.wikimedia.org/wiki/File:ActionPanama.JPG>

# VÈME LÉGISLATURE (1889-1893)

- Renouvellement partiel du personnel politique (boulangisme et scandale de Panama)
- Premières manifestations du Ralliement des catholiques à la République
- Tournant de la politique douanière (lois Méline)
- Essor du socialisme et du syndicalisme (Fourmies)
- Premiers attentats anarchistes
- ...



# WORKFLOW





Icone conçue par Freepik

## EXTRAIRE LE TEXTE / RÉ-OCÉRISER

- ❑ Récupération des textes océrisés via API Document de Gallica
  - ❑ Erreurs dues à la courbure de la page au niveau de la reliure
- ❑ Deux options pour ré-océrer :
  - ❑ Utiliser une solution propriétaire : Abby Fine Reader (excellentes performances sur les corpus de presse)
  - ❑ Entraîner un modèle avec une solution libre (Tesseract, OCR4all, Kraken)
- ❑ Post-corrections :
  - ❑ A la main
  - ❑ Emploi d'expressions régulières et librairie Python pypellchecker

# ANNOTER ET ENRICHIR

## Annotation des fichiers transcrits en XML-TEI (Text Encoding Initiative)

### Choix des annotations

- Sections très formalisées
- Entités nommées (Personnes)
- Entités et segments liés aux mouvements, partis et idées politiques
- Perspective « données liées » : identification des annotations avec des URI fournies par data.bnf et Wikidata

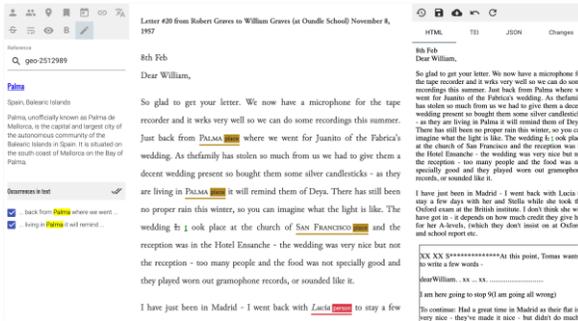
### Création d'un schéma TEI

- ODD (One Document Does it all)
- Adapté aux débats parlementaires
- Tenant compte des annotations choisies

### Phase d'annotation

- Scripts Python (par ex. script LSE-OD2M : <https://github.com/TimeUs-ANR/LSE-OD2M>)
- Annotations d'abord réalisées à la main = données d'entraînement pour l'automatisation de la tâche par l'analyseur syntaxique FrMG
- Annotation des entités nommées avec FastText embeddings et le modèle linguistique contextuel neuronal français CamemBERT dans une architecture LSTM-CRF

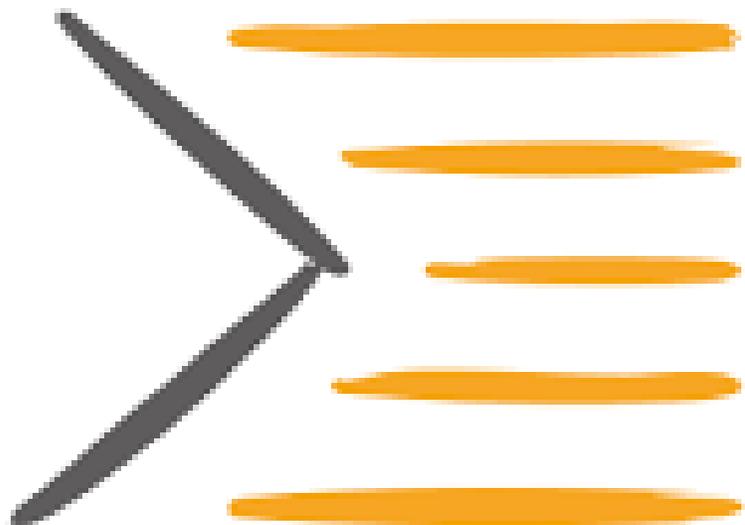
# UN CORPUS ÉDITORIALISÉ



Corpus encodé en TEI et stocké dans  
une base de données XML eXist-db

TEI Publisher : transformer les données  
stockées dans la base de données en  
pages Web HTML pour leur publication

Accès à un corpus éditorialisé, intégré dans  
un contexte d'application avec navigation,  
recherche plein-texte et affichage des fac-  
similés



## PUBLIER UN CORPUS ÉDITORIALISÉ

---

- Technologie *Web Components* → possibilité de développer de nouvelles fonctionnalités
- Accès aux annotations linguistiques
- Affichage des entités nommées
- Exploration avec la modélisation de sujets

12 JANVIER

Session ordinaire de 1893 11

DE CONGÉ

ac-Adaras et de  
t de ne pouvoir  
r et demandent

oyées à la com-

IDENT DU SÉNAT

u de M. le pré-  
unication sui-

1 janvier 1893.

nt,

ns ses séances  
élection de son  
ve composé de

demôle, Challe-

et que les **espérances** les plus **généreuses**,  
les plus **viriles pensées** viennent de ceux  
qui ne vieillissent pas. On reste jeune quand  
on s'oublie soi-même pour ne songer qu'à  
son pays. (*Applaudissements.*)

Non, ce ne sont pas des **défaillances** indi-  
viduelles qui pourront atteindre la **Répu-  
blique**. Le **suffrage universel** a moins de  
**passion**, plus de bon sens et d'**équité** que  
les **mœurs** politiques. (*Très bien! très  
bien!*) Il se fait aux **mœurs** de la **liberté**, et  
il sait qu'à d'autres époques le silence et  
l'**impunité** étaient acquis aux **fautes** que la  
**République** veut **dévoiler** et saura **punir**.  
(*Vifs applaudissements à gauche et au  
centre.*)

C'est en vain qu'on tente de se faire une  
**arme** contre les **institutions** de la **rigueur**  
que les **pouvoirs publics** et la **justice** ap-  
porteront dans la **répression**; c'est en vain  
qu'on espère que le **suffrage universel** ne  
se montrera pas assez **éclairé** pour démêler  
ceux qui ont failli aux **lois** de l'**honneur** et  
ceux que la **calomnie** cherche à mettre en  
cause. (*Très bien! très bien!*)

Exemple de détection automatique de topics  
dans un rapport de l'assemblée (CC-BY-SA  
Nicolas Bourgeois)

## MODÉLISER LES SUJETS

- ❑ Modélisation de sujets ou *topic modeling* : une méthode d'apprentissage non supervisée qui permet de découvrir les structures sémantiques latentes d'un corpus de textes, sans faire appel à des ressources sémantiques et lexicales.
- ❑ Particulièrement bien adapté pour l'étude de grands corpus, comme les corpus de presse par exemple
- ❑ Réalisé avec le Natural Language Toolkit (NLTK) et la librairie Python Gensim



Source : BnF <https://gallica.bnf.fr/ark:/12148/btv1b530893305>

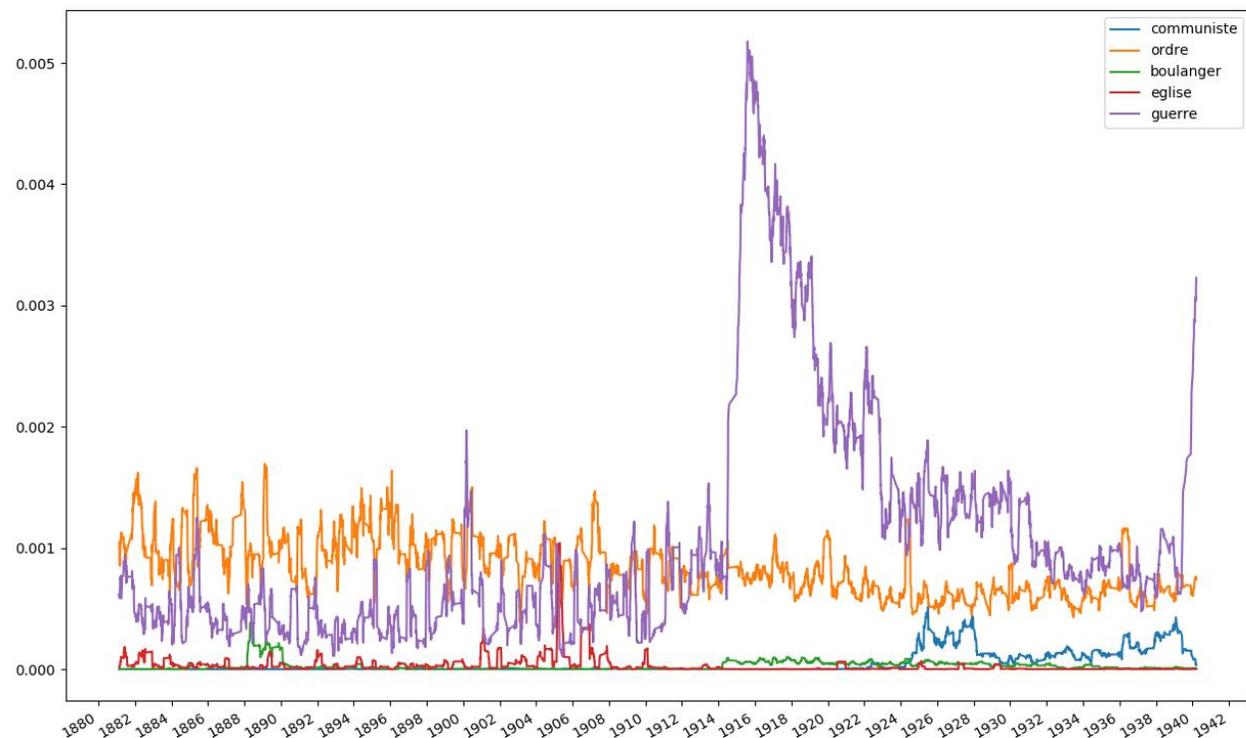
## RÉSULTATS ATTENDUS



- ✓ Une plateforme de consultation
- ✓ Le code source, les données annotées en XML-TEI accessibles
  - ✓ Une documentation pour utiliser la plateforme
- ✓ De nouvelles analyses (articles scientifiques, conférences...)



# PREMIERS TESTS : FRÉQUENCE DES MOTS



Evolution de la fréquence de mots  
choisis dans les débats  
parlementaires (1880-1942)

(CC-BY-SA Nicolas Bourgeois)



# MERCI



[marie.puren@epitech.eu](mailto:marie.puren@epitech.eu) [pierre.vernus@msh-lse.fr](mailto:pierre.vernus@msh-lse.fr)

Eugène Atget  
*Palais Bourbons – Chambre des Députés*  
Bibliothèque Nationale de France  
Entre 1900 et 1927 (d'après négatif 1900)