

# EXTRAIRE DU TEXTE DES IMPRIMÉS ANCIENS : QUELS DÉFIS, QUELLES PROMESSES ?

Illustration par le projet AGODA

---

Marie Puren<sup>1 2</sup> (et Aurélien Pellet et Pierre Vernus)

8 février 2024 - Séminaire de l'axe de recherche en histoire numérique du LARHRA

<sup>1</sup>LRE, EPITA <sup>2</sup>CJM, Ecole nationale des chartes

# LES COMPTES-RENDUS DES DÉBATS PARLEMENTAIRES FRANÇAIS

574

ASSEMBLÉE NATIONALE - 1<sup>re</sup> SÉANCE DU 12 JUIN 2023

qui devra rendre des comptes au Parlement, et pas le contraire) (Approuvé par le Sénat par le groupe SOC et par plusieurs députés du groupe LR-FRANÇES).

Jean-Jérôme Ghall, sur l'Assemblée nationale, vous l'avez fait et en expliquant le début de la présidence de votre référent des retraites. Il est de votre devoir, puisque que le Gouvernement a des positions sur le positionnement de l'Assemblée nationale pour que ce débat n'est pas fin.

**M. Jérôme Ghall.** Eh oui ! C'est bien sûr pour le président !

**M. Jean-François Cornu.** C'est fait et impromptu ! Des projets ?

**Mme Valérie Rabault.** Cette décision est un très grand précédent. Pour la première fois dans l'histoire de la V<sup>e</sup> République, des amendements reprenant une disposition pour dont fin déclarée recevable par l'Assemblée nationale – une première faite par vous l'année de 2015, une seconde fois par le président de la commission des finances le 30 mai – n'ont pu être ni en discussion dans l'hémicycle, ni vos collègues qui ont voté en ce moment faire entendre, le répètez que la décision est très claire, elle a même été faite par Eric Woerth le 25 février 2022, lorsqu'il était président de la commission des finances. (Approuvé par plusieurs députés du groupe SOC et LR-FRANÇES) Je vous en donne lecture : « L'Assemblée a donc proposé inviolable une explication personnelle de ne pas poursuivre autrement les initiatives parlementaires... En l'espèce, le « doit proposé », c'est la proposition de loi. Comme je l'ai expliqué depuis 1958 sans aucune exception, mes chers collègues, vous pouvez présenter les comptes rendus publiés au Journal officiel, sous le titre des amendements acceptés réceptifs.

**M. Erwan Billaut.** Il n'y a pas plus sans argument !

**Mme Valérie Rabault.** Malheureusement à mon collègue qui demandait de cette pratique pendant longtemps depuis 1958, je prendrai un exemple très récent, celui de la proposition de loi n° 1092 visant à assurer au regard de la Sécurité sociale, les mandataires dans notre hémicycle le 9 février dernier. Cette proposition de loi est en ce moment en débat public avec un article 1<sup>er</sup> modifié, qui reconnaît le bénéfice du repos à long terme de notre système de retraite de retraite. Très simplement, l'auteur de la proposition a déposé ses amendements pour modifier la volonté législative de créer un repos à long terme pour nos collègues. (Approuvé par plusieurs députés du groupe SOC et par plusieurs députés du groupe LR-FRANÇES) Concrètement à l'usage courant depuis 1958 à l'Assemblée nationale, soit amendement n° 1 est jugé recevable financièrement.

**M. Jérôme Ghall.** Eh oui !

**Mme Valérie Rabault.** Il a donc pu être discuté dans l'hémicycle et la censure du vote (Réponse occasionnelle) l'absence de cette lecture de l'article 49 n'a donc été faite aucun député des groupes parlementaires qui ont soutenu le Gouvernement, aucun d'entre vous, mes chers collègues, et ce non sans l'absence de la loi de règlement de l'hémicycle national, que vous pourrez à votre retour dans l'hémicycle, de soutenir la recevabilité financière d'un amendement. Vous ne l'avez pas fait (M. Benjamin Lecomte applaudit) Car exemple, qui n'est pas un pari d'argent, découvrez que la décision d'acceptabilité prise la semaine

dernière est contraire à tous les usages de notre assemblée. (Approuvé par le Sénat par le groupe SOC et par plusieurs députés du groupe LR-FRANÇES et LR-FRANÇES)

**M. Jean-François Cornu.** C'est fait !

**Mme Valérie Rabault.** En procédant ainsi, vous avez ouvert une voie dangereuse : celle de l'arbitraire, celle qui absorbe tout ce droit que nous avons consacré et consacré au sein de fil des ans depuis la Révolution française. Qui, malade la Première ministre, si s'agit bien d'arbitraire, jusqu'à ce point à l'instant, tous les députés, quelle que soit leur appartenance politique, étaient réunis à l'Assemblée, conformément à un usage consacré depuis 1958. Depuis ce 8 juin, c'est un sans-dieu l'instant, c'est le régime de la recevabilité partielle, sans voix de mort. (Approuvé par plusieurs députés du groupe SOC, LR-FRANÇES et LR-FRANÇES) Unanimes, aucun député de l'Assemblée nationale ne pourra perdre le vote qui sera obtenu aux amendements de stabilisation d'articles de loi d'un projet parlementaire de loi.

**M. Jérôme Ghall.** C'est bien sûr le contraire !

**Mme Valérie Rabault.** Malheureusement la Première ministre, en faisant pression sur le président de l'Assemblée nationale, vous avez fait dans l'article de la Constitution : l'article 48, article 5, qui garantit le droit d'initiative des groupes d'opposition, et l'article 51-3 qui confie aux parlementaires – et à moi – la responsabilité de « déterminer les sujets des groupes parlementaires contrôlés, au sein de chaque assemblée.

C'est ce que j'ai justement eu avec cette motion de censure, c'est la possibilité de mettre un terme au caractère du Gouvernement sur une Assemblée nationale. (Approuvé par plusieurs députés du groupe SOC, LR-FRANÇES et LR-FRANÇES) – Mme Bernadette Deshayes et M. Stéphane Flipo applaudit) – Cela veut dire que nous sommes réunis et vous nous entendrez pour débattre ce scandale, et surtout pour y mettre un terme. C'est une question de responsabilité de chacun politiquement ; c'est une question de responsabilité...

**M. Benjamin Lecomte.** Entièrement.

**Mme Valérie Rabault.** «... que chaque et chaque d'entre nous, sans voix de mort Française et des Français que nous représentons. (Mme Cécile Clément et M. Benjamin Lecomte applaudit) C'est une question de responsabilité que nous sommes réunis ; c'est une question de responsabilité de l'hémicycle et de la République française que nous sommes réunis dans les députations. Car qui a voté sans avoir été entendu ne réussira à expliquer à nos électeurs et électeurs et ce n'est ni la justice ni le premier pas sur une voie de la réforme des retraites) (Approuvé par plusieurs députés du groupe SOC et LR-FRANÇES)

**Mme Chloé Guichard.** La France à qui ?

**M. Stéphane Mullard.** C'est fait !

**Mme Valérie Rabault.** Qui d'entre nous peut occuper la duplicité gouvernementale qui consiste d'un côté à déclarer la recevabilité d'un amendement de la Première ministre – que l'on sait il se agit la dimension parlementaire sans la dernière étape de l'acte, à corriger le vote d'une proposition de loi qui nous amène parfois de dévoter d'incertitude au sein de la réforme des retraites) (Émotionnel par plusieurs députés du groupe SOC)

**M. Stéphane Mullard.** Rendez-vous Valérie Rabault !

- Naissance spontanée des comptes-rendus des débats en 1789
- Publication de discussions entre parlementaires jusqu'à nos jours
- A partir de 1881 : diffusion des débats à la Chambre des députés et au Sénat via des publications dédiées
- Forme très similaire à celles des débats publiés aujourd'hui

FIGURE – Séance parlementaire du 12 juin 2023

# LES DÉBATS PARLEMENTAIRES DURANT LA III<sup>E</sup> RÉPUBLIQUE

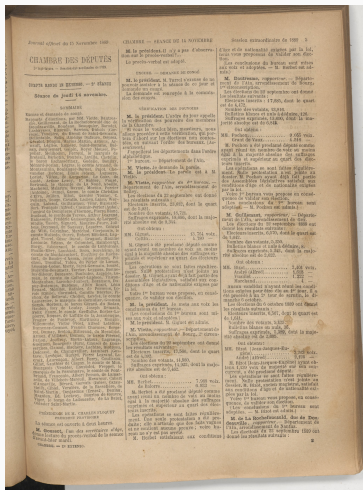


FIGURE – Séance parlementaire du 14 novembre 1889

- **Projet AGODA : Analyse sémantique et Graphes relationnels pour l'Ouverture et l'étude des Débats à l'Assemblée nationale**
- **Débats à la Chambre des députés (chambre basse du parlement) transcrits :**
  - entre 1876 et 1880, dans les **Annales du Sénat et de la Chambre des députés**
  - à partir de 1881, dans le **Journal officiel de la République française. Débats parlementaires (1881-1940)**
- Disponible en ligne via **Gallica** (bibliothèque numérique de la Bibliothèque nationale de France)

- Permettent d'observer l'évolution des questions politiques et sociales traitées à l'assemblée
- « Corpus remarquable » pour l'histoire politique et intellectuelle
- Sources précieuses pour l'histoire sociale, économique, culturelle et l'histoire du droit

# OBJECTIFS

- Créer une plateforme de consultation
- Produire des données textuelles structurées et sémantiquement enrichies à partir de ces débats numérisés
- Contribuer à la conception d'un workflow adapté à la préparation, à la publication et à l'analyse de grands corpus de documents historiques

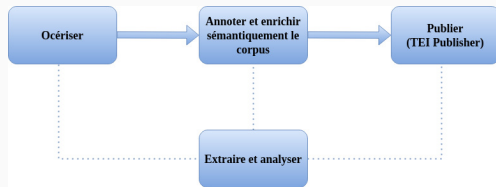
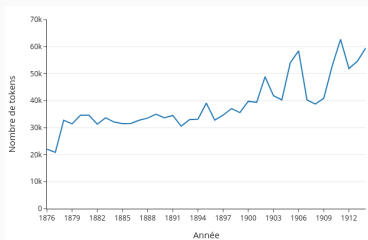


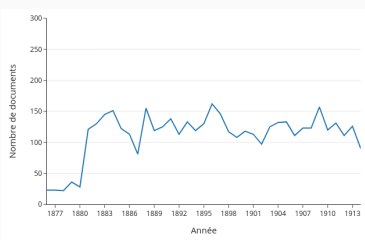
FIGURE – Les étapes de la chaîne de traitement

- Projet « preuve de concept »
- Ne pas traiter l'ensemble des débats parlementaires depuis 1789 mais travailler sur des sous-corpus test
  - Débats à la Chambre des députés
  - Plateforme de publication : période 1889-1893
  - Analyses en TAL menées sur un corpus allant de 1876 à 1914

Un corpus « massif ». Exemple : 15 législatures entre 1881 et 1940 à la Chambre des députés / 10-12.000 images par législature



(a) Médiane du nombre de tokens par années (1876-1914)



(b) Évolution du nombre de documents par années 1876-1914

**FIGURE** – Un corpus qui ne cesse de grandir au cours du temps

# OCÉRISER LES DÉBATS

---



# LE CAS DES DÉBATS PARLEMENTAIRES

- Récupération des textes océrisés via l'API Document de Gallica  
=> qualité inégale de l'OCR
- Erreurs dues à :
  - qualité du document : tâches et surimpression
  - la courbure de la page au niveau de la reliure

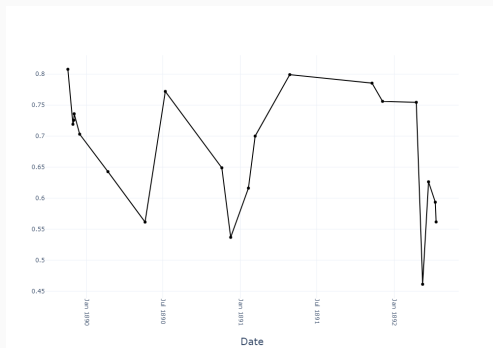


FIGURE – Evaluation de la qualité de l'OCR fourni par Gallica

# EFFET DE LA COURBURE SUR LES RÉSULTATS DE L'OCR

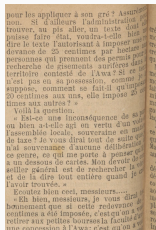


FIGURE – 20  
octobre  
1890  
(p.1718)

pour les appliquer à son gré ? Assure.

non. Si d'ailleurs l'administration f, trouver, au pis aller, un texte (10111 en puisse faire état, voudra-t-elle dire le texte l'autorisant à imposer un devance de 25 centimes par hectar , personnes qui prennent des permis Pour l recherche de gisements aurifères a teSI territoire contesté de l'AwA ? Si ce tes

n'est pas en sa possession, comas le suppose, comment se fait-il qu'ilWjs3j 20 centimes aux uns, elle impose 20 cet times aux autres ? » , Voilà la question. é

« Est-ce une inconséquence de sa te ou bien a-t-elle agi en vertu d'un v l'assemblée locale, souveraine en OIqallC de taxe? Je vous dirai tout de suite II n'ai souverance d'aucune délibératiIV ce genre, ce qui me porte à penser a un dessous de cartes. Mon devoir seiller général est de rechercher la cJ1" et de la dire tout entière quand Je l'avoir trouvée. »

Eoutez bien ceci, messieurs. l «Eh bien, messieurs, je vous dirà de bonnement que si cette redevance J centimes a été imposée, c'est qu'on a, retirer aux petites bourses la faculté" une concession à l'AwA; c'est qu'on a rj;3l

FIGURE – Résultat de l'OCR

Décision de ré-océreriser le corpus

Améliorer la qualité de l'image avec une méthode de « dewarping »  
=> résultats peu probants

- Gérer la courbure des pages avec le dewarping?
- Utiliser des outils plus avancés?



(a) Image d'origine



(b) Image « dewarpée »

FIGURE – Dewarping : pas adapté à nos documents

# NETTOYAGE DE L'IMAGE



(a) Image d'origine



(b) Image nettoyée

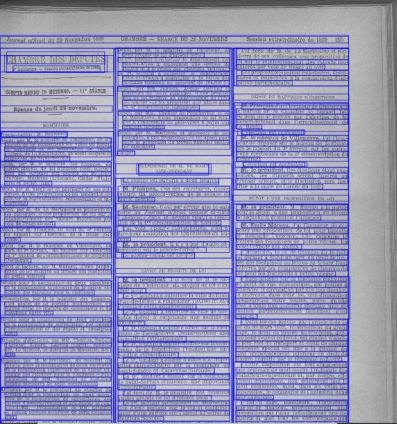
FIGURE - Exemple d'image nettoyée

# ANALYSE DE LA MISE EN PAGE

Directory

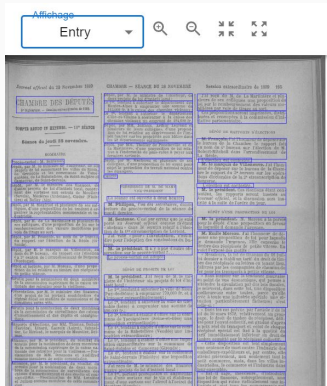
Directory Viewer    zz\_jo\_18891128.pdf    1    EXPORT    SAVE

Affichage    🔍    🔍    ↶ ↷    ✎    ✎    Edit inlr



The image shows a PDF viewer interface with a blue header bar. The header contains the text 'Directory Viewer', a dropdown menu showing 'zz\_jo\_18891128.pdf', a page indicator '1', and two buttons: 'EXPORT' and 'SAVE'. Below the header is a toolbar with a dropdown menu labeled 'Affichage', two magnifying glass icons, four arrow icons for navigation, two pencil icons, and a toggle switch labeled 'Edit inlr'. The main content area displays a document page with a blue header and a grid of text. The text is mostly illegible but appears to be a technical or administrative document. The document content is mostly illegible but appears to be a technical or administrative document.

# ANALYSE DE LA MISE EN PAGE



(a) Zones "entrées"



(b) Zones "Titre II"

FIGURE – Distinction entre différentes "zones"

The screenshot displays a web-based interface for document viewing and analysis. At the top, a blue header bar contains the text "Directory Viewer" and the filename "zz\_jo\_18891128.pdf". Navigation icons for back, forward, and search are present, along with "EXPORT" and "SAVE" buttons. Below the header, a control bar includes a search input field with the text "Affichage Entry, Title II", a magnifying glass icon, and a "Edit inline" toggle switch. The main content area is split into two panels. The left panel shows a thumbnail of a document page with a grid of text blocks. The right panel shows a zoomed-in view of a text snippet: "M. François PER . J PER 'ai l'honneur de déposer sur le bureau de la Chambre le rapport fait au nom du 4e bureau sur l'élection de M. Robert-Mitchell dans l' arrondissement de Réole LOC". The words "PER" and "LOC" are highlighted in orange and pink respectively, indicating Named Entity Recognition. Above this text are two buttons labeled "OCR" and "NER". Below the text is a "Comment" field containing the text "u-beginning seg".

FIGURE – Démonstration de l'outil sur une page de débat parlementaire

1. Évaluation, analyse comparative
2. Outils prêts pour la production, vraiment en libre accès
3. D'autres méthodes que le NER pour l'extraction de données structurées (entryseg + structured output + reading order + ...)
4. Projet Mezanno (financé dans le cadre du plan quadriennal de la BnF) : DH sans informaticiens ?

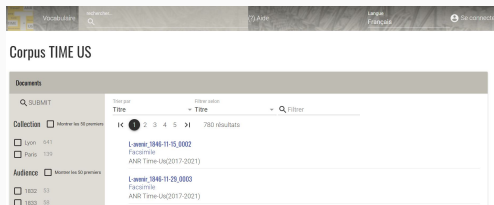


- Page web sans installation ni serveur
- Interface de constitution de corpus à partir de ressources IIF publiques
  - Gallica
  - Images sur Nakala, Zenodo...
- Interface de création et modification d'annotations
  - Adaptée à la transcription
  - Stockage local ou serveur d'annotation collaboratif
- Possibilité d'appeler des API distantes standardisées
  - Layout detection
  - OCR
- Facilité d'export des données
  - CSV...

**Réutiliser tous les outils existants, et maximiser  
l'indépendance chercheurs/ses en SHS**

# ANNOTER LES DÉBATS EN XML-TEI

---



(a) Interface de recherche



(b) Affichage d'un document

FIGURE – Corpus du projet ANR TIME-US publié avec TEI Publisher

# CORPS D'UN DÉBAT ENCODÉ EN XML-TEI

```
<div type="part" corresp="#rapportelections">
  <head>SUITE DE LA VERIFICATION DES POUVOIRS</head>
  <u who="#pers_ID" xml:id="CR_1889-11-26_u23" ana="#chair">
    <seg xml:id="CR_1889-11-26_u23.1"><persName ref="#pers_ID">M. le <roleName ref="#pers_ID">président</roleName></persName>. L'ordre du jour
appelle la suite de la vérification des pouvoirs.</seg>
    <seg xml:id="CR_1889-11-26_u23.2"><persName ref="#pers_ID">M. Reybert</persName> a la parole pour donner lecture d'un rapport sur une
élection non contestée.</seg>
  </u>
  <u who="#pers_ID" xml:id="CR_1889-11-26_u24" ana="#rapporteur">
    <!-- Lecture d'un rapport -->
    <quote>
      <seg xml:id="CR_1889-11-26_u24.1"><persName ref="#pers_ID">M. Reybert, <roleName ref="#pers_ID">rapporteur</roleName></persName>.
--<placeName ref="#lieu_ID">Département de la Corrèze, arrondissement de Tulle, <num>1</num> circonscription</placeName>.</seg>
      <seg xml:id="CR_1889-11-26_u24.2">Les élections du <date when="1889-09-22">22 septembre</date> ont donné les résultats suivants :</seg>
      <seg xml:id="CR_1889-11-26_u24.3">Electeurs inscrits, <num>17,950</num>, dont le quart est de <num>4,263</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.4">Nombre des votants, <num>12,322</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.5">Bulletins blancs et nuls, à déduire, <num>133</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.6">Suffrages exprimés, <num>12,189</num>.</seg>
      <seg xml:id="CR_1889-11-26_u24.7">Ont obtenu : MM. <persName ref="#pers_ID">Borie, <roleName ref="#pers_ID">député
sortant</roleName></persName>... <num>7,508</num> voix</seg>
      <seg xml:id="CR_1889-11-26_u24.8"><persName ref="#pers_ID">Vachal, <roleName ref="#pers_ID">ancien député</roleName></persName>...
<num>4,748</num></seg>
      <seg xml:id="CR_1889-11-26_u24.9"><persName ref="#pers_ID">M. Borie</persName> a été proclamé député comme ayant réuni un nombre de voix
au moins égal à la majorité absolue des suffrages</seg>
    </quote>
  </u>
  <floatingText><body><div><p n="176"/></div></body></floatingText>
  <!-- [...] -->
</div>
```

FIGURE – Encodage - Séance parlementaire du 26 novembre 1889 (extrait)

# ANNEXES D'UN DÉBAT ENCODÉES EN XML-TEI

Annexes au procès-verbal de la séance du mardi 26 novembre 1889.

---

SCRUTIN

Sur les conclusions du 7<sup>e</sup> bureau tendant à l'annulation des opérations électorales de la 1<sup>re</sup> circonscription de l'arrondissement de Lorient (Morbihan).

Nombre des votants.....	506
Majorité absolue.....	253
Pour l'adoption.....	330
Contre.....	176

La Chambre des députés a adopté.

---

NOT VOTÉ POUR :

MM. Abeille, Aïnos (Emanuel), Armez, Arribat, Audifred, Aynard (Eduard).

Balle (Marial), Bergy, Barodet, Barthou, Bataillon, Bédit (Adrien), Benjard, Beauquier, Bérard, Berger (Georges) (Séclé), Bertrand, Bézine, Bisarville, Bisol, Bissonard-Bert, Boute (Pierre), Boulay-d'Aupiais, Boudigoy-Sibour, Bony-Casténon, Bourglione, Bouchet (Vierge), Bouchonnet, Boudetille, Bouge, Boulanger-Bennet, Boullay, Bourgeois (Jules), Bourgeois (Léon) (Marie), Boulière de Bouchéger, Boyer-Lapierre, Buvard, Brand, Breton, Briens, Brisson (Henri), Broseau (Emile), Brognon, Buvard, Bully, Boudou, Davignier.

---

Rectifications aux scrutins de la séance du 25 novembre 1889.

M. Michau (Nord), porté comme s'étant abstenu dans le scrutin sur l'urgence de la proposition de M. Maxime Lecomte, déclare avoir voté pour ».

(a) Source numérisée

```
<!-- ANNEXES -->
<back>
<head>Annexes au procès-verbal de la séance du <date when="1889-11-26">mardi 26 novembre 1889</date>.</head>

<div xml:id="vot18891126">
<!-- VOTE 1 -->
<div xml:id="vot18891126_vot1" type="voting" corresp="mdiscussion7ebureau">
<head>
<label>SCRUTIN</label>
<note>seg<!-- Sur les conclusions du <num>7</num> bureau tendant à l'annulation des opérations électorales de la
<placeName ref="#lieu_ID"><num>1</num> circonscription de l'arrondissement de Lorient (Morbihan)</placeName>.</seg--></note>
</head>
<!-- Détail du vote -->
<desc>
<measure type="nbvotants" quantity="506">Nombre des votants <num>506</num></measure>
<measure type="maj" quantity="254">Majorité absolue <num>254</num></measure>
<measure type="ayes" quantity="330">Pour l'adoption <num>330</num></measure>
<measure type="noes" quantity="176">Contre <num>176</num></measure>
</desc>
<note type="result">seg<!-- La <orgName ref="#org_ID">Chambre des députés</orgName> a adopté.</seg--></note>
<floatingText>body<div pb="192"></div></body></floatingText>
<!-- Liste des votants -->
<note type="voterslist">
<desc>Ont voté pour :</desc>
<seg>MM.<persName ref="#pers_ID">Abeille</persName>. <persName ref="#pers_ID">Arène (Emanuel)</persName>. <persName
ref="#pers_ID">Armez</persName>. <persName ref="#pers_ID">Arribat</persName>. <persName ref="#pers_ID">Audifred</persName>. <persName ref="#pers_ID">Aynard (Eduard)</persName>.</seg>
<!-- [...] -->
</note>
</div>
<!-- RECTIFICATIONS -->
<div corresp="vot18891125" type="rectification">
<head>Rectifications aux scrutins de la séance du <date>25 novembre 1889</date>.</head>
<note corresp="vot18891125_vot1">seg<!-- <persName ref="#pers_ID">M. Michau</persName> <placeName
ref="#lieu_ID">(Nord)</placeName>, porté comme s'étant abstenu dans le scrutin sur l'urgence de la proposition de <persName
ref="#pers_ID">M. Maxime Lecomte</persName>, déclare avoir voté pour ».</seg--></note>
<!-- [...] -->
</div>
</div>
</back>
```

(b) Modèle d'encodage

FIGURE – Séance parlementaire du 26 novembre 1889 - votes, liste des votants, rectifications (extrait annexes)

# STRUCTURE GÉNÉRALE DES FICHIERS XML TEI

```
<teiCorpus xmlns="http://www.tei-c.org/ns/1.0" xsl:id="FR_3R_5L" xsl:lang="fr">
  <!-- Métadonnées du corpus (non développées) -->
  <teiHeader>
    <fileDesc>
      <titleStnt>
        <title></title>
      </titleStnt>
      <publicationStnt>
        <sp></sp>
      </publicationStnt>
      <sourceDesc>
        <sp></sp>
      </sourceDesc>
    </fileDesc>
  </teiHeader>
  <!-- Données liées et informations contextuelles -->
  <standoff>
    <listPerson>
      <person></person>
    </listPerson>
    <listOrg>
      <org></org>
    </listOrg>
    <listPlace>
      <place></place>
    </listPlace>
  </standoff>
  <!-- Stockage du composant correspondant à la séance du 26 novembre 1889 -->
  <xli:include xmlns:xli="http://www.w3.org/2001/XInclude" href="FR_3R_5L_1889-11-26.xml"/>
  <!-- Stockage des autres composants du corpus de façon identique -->
  <!-- ... -->
</teiCorpus>
```

(a) Structure générale d'un fichier corpus

```
<TEI xmlns="http://www.tei-c.org/ns/1.0" xsl:id="FR_3R_5L_1889-11-26" xsl:lang="fr">
  <!-- Métadonnées du composant (non développées) -->
  <teiHeader>
    <fileDesc>
      <titleStnt>
        <title></title>
      </titleStnt>
      <publicationStnt>
        <sp></sp>
      </publicationStnt>
      <sourceDesc>
        <sp></sp>
      </sourceDesc>
    </fileDesc>
  </teiHeader>
  <text>
    <!-- Transcription du compte rendu -->
    <body>
      <div></div>
    </body>
    <!-- Annexes du compte rendu -->
    <back></back>
  </text>
</TEI>
```

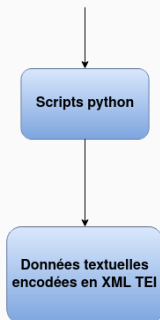
(b) Structure générale d'un fichier composant

FIGURE – Encodages - Structure générale des fichiers corpus et composant

# APPLIQUER L'ENCODAGE : AUTOMATISATION

```
▼ 19:
activities:      []
addresses:      []
▼ box:
  0:             220.12998972250773
  1:             202.1299897225077
  2:             561.8527187504491
  3:             00.4418437486525
checked:        true
comment:        "seg"
id:             276
▼ ner_xml:
  origin:        "<PER>M. Borie</PER> a «ACT>déjà fait partie des Assemblées\u2029législatives et satisfait aux conditions d«ACT>«ACT>âge\u2029et de
  parent:        "computer"
  text_ocr:      "M. Borie a déjà fait partie des Assemblées\nlégislatives et satisfait aux conditions d'âge\net de nationalité exigées par la loi."
  type:          "ENTRY"
```

Données textuelles au format JSON



- Données textuelles dans les fichiers JSON = valeurs des clés « text\_ocr »
- Clés contenues dans des objets qui correspondent à la segmentation au niveau paragraphe
- Appliquer des règles de transformation sur les valeurs des clés « text\_ocr »



# A LA RECHERCHE DE « FEATURES »

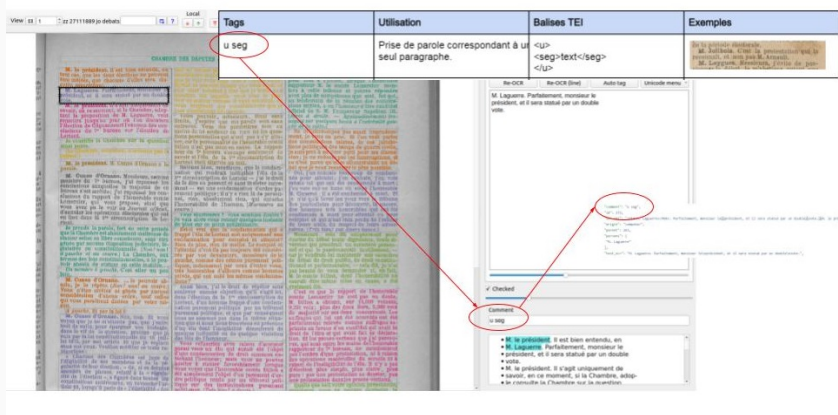


FIGURE – Du guide d'annotation à l'export JSON

```
{
  "activities": [],
  "addresses": [],
  "box": {
    "x1": 57.616701158597,
    "y1": 1710.0,
    "x2": 58.3883298849495,
    "y2": 58.07766597698992
  },
  "checked": true,
  "comment": "u seg",
  "id": 305,
  "key": {
    "x": 0,
    "y": 1735
  },
  "text_xml": "<PER-M, Paul D roul de</PER>. Je demand  la pa-rol .",
  "origin": "computer",
  "parent": 269,
  "persons": [
    "M. Paul D roul de"
  ],
  "text_ocr": "M. Paul D roul de. Je demand  la pa-rol .",
  "type": "ENTRY"
}
```

(a) Fichier JSON avec une prise de parole

```
def add_utterance(data):
    """
    Ajout de l' l ment TEI "u" pour chaque box  tiquet e "u" ou "u-beginning" et "u-end"
    :param data: dictionnaire contenant l'ensemble des donn es issues des JSON
    """
    for i in range(len(data)):
        if "comment" in data[i]:
            if re.search(r"^(u|u-)?$", data[i]["comment"]):
                data[i]["text_ocr"] = "".join(["<u>", data[i]["text_ocr"], "</u>"])
            elif re.search(r"u-beginning", data[i]["comment"]):
                data[i]["text_ocr"] = "".join(["<u>", data[i]["text_ocr"]])
            elif re.search(r"u-end", data[i]["comment"]):
                data[i]["text_ocr"] = "".join([data[i]["text_ocr"], "</u>"])
            else:
                pass
    return data
```

(b) Fonction add\_utterance

```
<u><seg>M. Paul D roul de. Je demand  la parol .</seg></u>
```

(c) R sultat

FIGURE – Exemple d'application de la fonction add\_utterance

1. Récupération des données JSON stockées dans une variable « data » sous forme d'un dictionnaire.
2. Gestion de l'encodage des changements de page.
3. Intégration des scripts contenant les règles de transformation dans la chaîne de traitement. Permettent d'encoder les données contenues dans la variable « data ».
4. Création de l'élément <teiHeader> contenant les métadonnées :
  - Rédaction d'une partie du <teiHeader> à la main contenant les métadonnées « fixes »
  - Construction de règles permettant de rechercher les métadonnées propres à chaque compte rendu, présentes dans la variable « data », par la suite intégrées au <teiHeader> préalablement établi
5. Création des fichiers XML valides.
6. Nettoyage des fichiers XML (gestion des césures).

1. Génération du fichier corpus.
2. Intégration des métadonnées.
3. Gestion des `<xi :include>`.

1. Annotation manuelle très consommatrice de temps et non exempte d'erreurs
  - Ajout de manière semi-automatique des étiquettes
  - Permettre aux utilisateurs de réutiliser les étiquettes et de créer leurs propres étiquettes (vocabulaire contrôlé?)
2. Utiliser ces scripts pour proposer un export XML-TEI (outil SoDUCo)
  - Ajouter les entités nommées (<persName>, <orgName>, etc.)

## SpanCategorizer

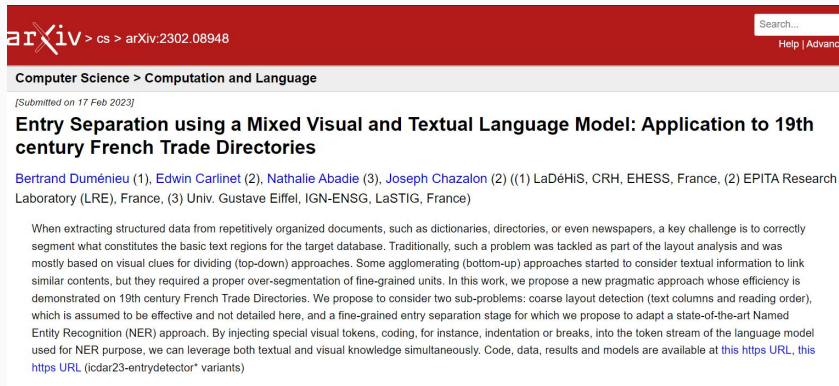
[</> SOURCE](#)**CLASS, EXPERIMENTAL****V3.1** ?**STRING NAME:** `spancat`**TRAINABLE:** 

Pipeline component for labeling potentially overlapping spans of text

---

A span categorizer consists of two parts: a suggester function that proposes candidate spans, which may or may not overlap, and a labeler model that predicts zero or more labels for each candidate.

FIGURE – SpanCategorizer : <https://spacy.io/api/spancategorizer>



The screenshot shows the arXiv interface. At the top left is the arXiv logo and navigation links: 'cs > arXiv:2302.08948'. At the top right is a search bar and 'Help | Advanced Search' link. Below the navigation is a breadcrumb trail: 'Computer Science > Computation and Language'. The main title of the paper is 'Entry Separation using a Mixed Visual and Textual Language Model: Application to 19th century French Trade Directories'. Below the title is the author list: 'Bertrand Duméniou (1), Edwin Carlinet (2), Nathalie Abadie (3), Joseph Chazalon (2) ((1) LaDéHIS, CRH, EHESS, France, (2) EPITA Research Laboratory (LRE), France, (3) Univ. Gustave Eiffel, IGN-ENSG, LaSTIG, France)'. The abstract text follows, starting with 'When extracting structured data from repetitively organized documents, such as dictionaries, directories, or even newspapers, a key challenge is to correctly segment what constitutes the basic text regions for the target database. Traditionally, such a problem was tackled as part of the layout analysis and was mostly based on visual clues for dividing (top-down) approaches. Some agglomerating (bottom-up) approaches started to consider textual information to link similar contents, but they required a proper over-segmentation of fine-grained units. In this work, we propose a new pragmatic approach whose efficiency is demonstrated on 19th century French Trade Directories. We propose to consider two sub-problems: coarse layout detection (text columns and reading order), which is assumed to be effective and not detailed here, and a fine-grained entry separation stage for which we propose to adapt a state-of-the-art Named Entity Recognition (NER) approach. By injecting special visual tokens, coding, for instance, indentation or breaks, into the token stream of the language model used for NER purpose, we can leverage both textual and visual knowledge simultaneously. Code, data, results and models are available at [this https URL](#), [this https URL](#) (icdar23-entrydetector\* variants)'. The URL 'https://arxiv.org/abs/2302.08948' is highlighted in orange in the original image.

FIGURE – Modèle de langue basé sur les Transformeurs (BERT) :

<https://arxiv.org/abs/2302.08948>

# TOPIC MODELING ET WORD EMBEDDING

---



- Corpus massif : long à traiter, impossible à lire dans le détail
- Difficile de repérer les grandes tendances

« Lecture distante » (Franco Moretti) comme condition de la connaissance : permet plus facilement de repérer des régularités ou des anomalies qu'une « lecture proche » n'aurait pas détectées.

- Méthodes d'analyse permettant de « lire à distance » le corpus
  - « Topic modeling » ou modélisation de sujets
  - « Word embedding » ou plongement de mots ou plongement sémantique
- Appliquées au corpus océrisé extrait de Gallica (sans post-correction)

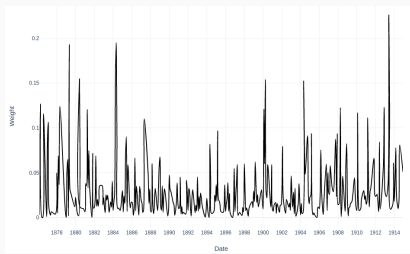
## QUELLE MÉTHODE ?

- Écueil de la lecture distante : imposer des catégories de lecture inappropriées
- Choix d'une méthode inductive : modélisation de sujets qui fait émerger les sujets à partir des textes eux-mêmes, sans intervention humaine

Méthode bien adaptée à l'étude de grands corpus historiques

- Particulièrement bien adaptée aux corpus de type « presse » : grand volume, sériels, traversés par de nombreux sujets différents évoluant dans le temps
- Identifier les sujets dans les débats :
  - Nouveau point d'entrée
  - Nouveau « mode de lecture »
  - Meilleure compréhension de l'évolution des idées/ sujets/ débats au cours du temps

- Méthode d'apprentissage non supervisée permettant de découvrir des sujets abstraits, dits « topics », issus d'un corpus à l'aide d'un modèle probabiliste LDA (Latent Dirichlet Allocation)
- Sujets = champs sémantiques préexistants à l'écriture du texte
- Textes construits à partir de ces champs sémantiques
- Principe de LDA : inverser le processus de génération textuelle, en partant du texte pour extraire les sujets



(a) Variation du poids des sujets liés à l'armée : moyenne par mois

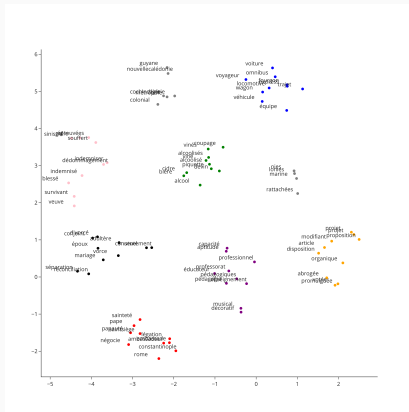
Topic 9	Topic 8	Topic 48	Topic 43
armée	cau	retraite	école
guerre	ville	caisse	primaire
homme	paris	pension	supérieur
militaire	canal	assurance	lycée

(b) 4 sujets ou « topics » parmi les 50 définis avec LDA

FIGURE – Résultats obtenus avec LDA

- Fournit une représentation des mots dans un espace vectoriel en fonction de leur similarité sémantique
- Repose sur l'hypothèse qu'un mot est caractérisé par son contexte, autrement dit par les mots qui l'entourent
  - Si les mots partagent des contextes similaires, ils partagent alors aussi des significations similaires

# LES PLONGEMENTS DE MOTS : WORD2VEC ET TOP2VEC



(a) Projection t-SNE des centroïdes des vecteurs (word2vec)

Cluster 55	Cluster 68	Cluster 70
victimes	divorce	enveloppes
inondations	epoux	timbres
secourir	mariage	poste
eprouves	conjugal	postale
orages	divorces	timbre
sinistres	adultere	recepisses
grele	conjugale	postes
secours	remarier	postaux
venir	separation	telegraphes
infortunes	indissolubilite	colis
ravages	conjoints	fixe
miseres	mutuel	recouvrements
catastrophe	separations	graphes
evenements	mari	postales
repartition	mariages	taxe
incendies	femme	decide
soulager	conjoint	soit

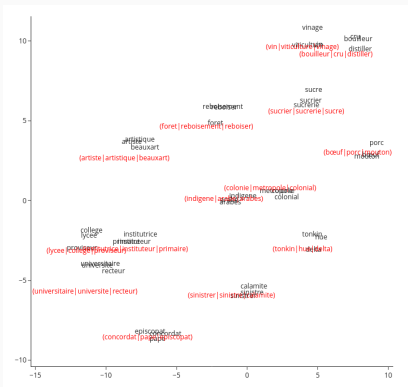
(b) 3 clusters parmi les 113 obtenus avec top2vec : tempêtes (55), divorce (68) et poste (70)

FIGURE – Résultats obtenus avec word2vec et top2vec



Topic 37	Topic 40	Topic 42	Topic 54
institutrice	concordat	artiste	lycée
instituteur	pape	beaux-arts	enseignement
paris	primaire	musée	universitaire
canal	élémentaire	louvre	internat
enseignement	catholique	décoratif	bachelier

(a) 4 sujets obtenus avec Top2vec



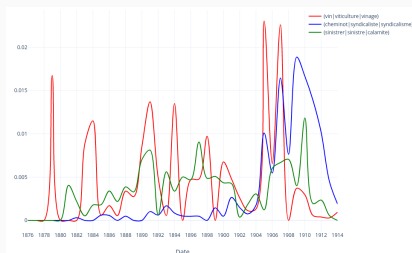
(b) Projection des mots par sujets avec la méthode t-SNE

FIGURE – Résultats obtenus avec Top2vec

Méthode « basée » sur Word2vec et Doc2vec. Nombre de sujets obtenus : 375.

Topic 34	Topic 45	Topic 55
Vin	Cheminot	Sinistré
Viticulture	Syndicaliste	Calamité
Alcoolisation	Confédération	Cyclone
Coupage	Militant	Gelée

(a) 3 sujets et leurs 4 mots les plus représentatifs



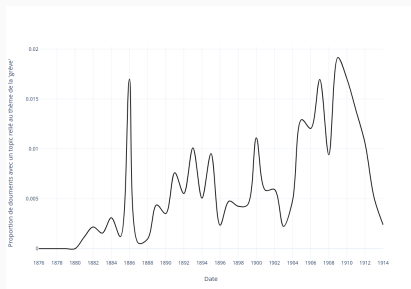
(b) Proportion des documents associés à chacun des trois sujets, agrégation par années

FIGURE – Evolution de sujets au cours du temps

Mais à la fin les ouvriers eurent connaissance du rôle que jouait M. Watrin et qui consistait encore une fois à obliger les chefs de la mine à réduire les salaires convenus; ils apprirent en outre que M. Watrin avait imaginé de réduire, au bout du mois, le salaire que l'ouvrier avait gagné, et cela à l'insu de l'ouvrier. Je m'explique: l'ouvrier croyait recevoir, d'après la quantité de travail qu'il avait faite, une certaine somme; M. Watrin se permettait de la réduire parfois de moitié sans prévenir les intéressés et sans explication.

J'arrive à la coopérative, c'est-à-dire à l'installation des magasins...

M. Laur. C'est là sa faute! Si M. Watrin ne s'était pas tant occupé de l'ouvrier, il ne serait pas mort.



(a) Extrait d'un document ayant une forte similarité avec le sujet de la grève pendant l'année 1886

(b) Proportion des documents associés aux topics sémantiquement similaires au mot "grève", agrégation par années

FIGURE – Evolution de sujets liés au champ sémantique de la grève dans le temps

Topic 1	Topic 2	Topic 3	Topic 4
de	de	suffrages	voté
la	discussion	exprimés	déclare
le	du	de	porté
et	loi	inscrits	ayant

FIGURE – Les 4 sujets obtenus avec CamemBERT



(a) Huit sujets obtenus avec BERTopic, et les cinq mots les plus représentatifs



(b) Projection des huit sujets avec t-SNE

FIGURE – 2ème itération avec paraphrase-multilingual-MiniLM-L12-v2

- Associer LDA et Top2vec :
  - LDA donne une vision plus générale du corpus;
  - Top2vec permet d'approfondir certains sujets.
- BERTopic : résultats assez décevants :
  - Travailler sur un corpus plus « propre » (ré-océréisé);
  - Utiliser un autre modèle de langue développé pour le français comme FlauBERT.

INTÉGRER CES RÉSULTATS AUX FICHIERS TEI

---

# ANNOTER LES DÉBATS AVEC LEURS SUJETS

```
<body>
  <!-- [...] -->
  <u>
    <!-- [...] -->
    <!-- "some of the war material in Madagascar" -->
    <w xml:id="ps1895022_116">
      une</w>
    <w xml:id="ps1895022_117">
      partie</w>
    <w xml:id="ps1895022_118">
      du</w>
    <w xml:id="ps1895022_119">
      matériel</w>
    <w xml:id="ps1895022_120">
      de</w>
    <w xml:id="ps1895022_121">
      guerre</w>
    <w xml:id="ps1895022_122">
      à</w>
    <w xml:id="ps1895022_123">
      Madagascar</w>.
    <!-- [...] -->
  </u>
  <!-- [...] -->
</body>
```

(a) Annotations dans le corps de texte

```
<standOff>
  <spanGrp type="topic">
    <span target="#ps1895022_119">
      army</span>
    <span target="#ps1895022_123">
      colonization</span>
  </spanGrp>
</standOff>
```

(b) Utilisation de l'élément  
<standOff>

**FIGURE** – Traitement des sujets comme une annotation linguistique dans le fichier XML-TEI



- Travailler sur des textes moins fautifs
- Extraire les sujets
  - LDA et Top2vec : complémentaires
  - Nouveau modèle de langue : FlauBERT ou spécialement développé pour le langage politique XIXe-XXe siècle
- Annoter automatiquement les textes avec leurs sujets ? Quelle stratégie ?

- Travailler sur des textes moins fautifs : nouvel OCR
- Extraire les sujets
  - LDA et Top2vec : complémentaires
  - Nouveau modèle de langue : FlauBERT ou spécialement développé pour le langage politique XIXe-XXe siècle
- Annoter automatiquement les textes avec leurs sujets ? Quelle stratégie ?

MERCI POUR VOTRE ATTENTION !



Marie Puren : `marie.puren@epita.fr`