

# **Modélisation de l'encodage automatique des débats parlementaires : transformation JSON vers XML TEI**

## **1- Fonctionnement général de l'encodage automatique**

Il est nécessaire au préalable de :

- Créer un dossier json\_data pour stocker l'ensemble des fichiers JSON annotés avec les étiquettes, corrigés et dont les boxes sont dans le bon ordre.
- Créer un dossier xml\_data pour stocker le schéma de l'encodage (agoda\_schema.rng) et le fichier corpus (ex : FR\_3R\_5L.xml) rédigé à la main sans les inclusions (xi:include). Les fichiers composants XML balisés en TEI seront inclus dans ce dossier et les inclusions seront ajoutées automatiquement dans le fichier corpus.

Une fois cela effectué, il est possible de faire tourner les scripts. Les fichiers composants XML sont créés, les données sont balisées en XML TEI. Les inclusions des fichiers composants dans le fichier XML corpus sont effectuées.

Pour plus de lisibilité des fichiers XML, il est possible de les ouvrir manuellement et de les indenter à l'aide d'un éditeur XML.

Les données obtenues sont encodées en UTF-8.

## **2- Note sur cette présente modélisation**

La modélisation est pensée pour :

- l'ensemble des CR de la 5e législature de la IIIe république (1889-1893).
- excepté les cas particuliers de la première séance de chaque année (CR à traiter par la suite selon ses spécificités et à inclure dans cette présente modélisation).

Certaines modifications et améliorations seront à effectuer sur les scripts lorsqu'il s'agira de traiter d'autres corpus de débats (évolutions des étiquettes).

Cette modélisation utilise des étiquettes pour appliquer les balises TEI de façon automatique. Ces dernières sont pensées en amont et incluses manuellement (pour l'instant) dans les fichiers JSON, au niveau de la clef "comment". Toutes les boxes contenant la clef "text\_ocr" avec les données textuelles doivent être annotées à l'aide de ces étiquettes.

### **3- Tableaux descriptifs de la modélisation des scripts**

Les tableaux présentés ci-dessous détaillent, script par script, les méthodes utilisées pour traiter l'ensemble des étiquettes que nous utilisons pour baliser le texte :

- chaque étiquette est associée à un résultat XML TEI attendu,
- certaines étiquettes peuvent avoir un traitement particulier qui est notifié,
- certaines étiquettes nécessitent d'être complétées (prise en compte des attributs XML par exemple).

#### **❖ *Script balisage formel***

	Étiquettes	Résultat XML TEI	Méthode et incertitudes	À ajouter
<b>add_seg</b> <b>(paragraphes)</b>	seg	<seg>text</seg>	<b>Attention</b> : traitement particulier pour les seg couplés avec quote-beginning et quote-end car chevauchement de balises. Pour ces cas, le seg est traité dans la	<b>Ajouter</b> : @xml:id

			fonction add_quote.	
	seg-beginning	<seg>text		
	seg-end	text</seg>		
<b>add_signed</b> <b>(signature)</b>	signed	<signed> <seg>text</seg> </signed>		
<b>add_page_number</b> <b>(changements de pages)</b>	body	<pb n=""/>	<p>Première page de la séance : à récupérer dans le texte.</p> <p>La récupération de la page permet de supprimer le restant de la ligne qui ne nous sert pas.</p>	
	body1	<pb n="1"/>	<b>Attention</b> : traitement différent à opérer pour la première page de l'année, le numéro de page est présent en bas de page.	<b>À inclure plus tard car traitement spécifique</b>
	page-number	<pb n="(zwt +1)"/>	Numéro de la page à obtenir grâce à l'incréméntation effectuée à partir de la première page.	
	page-number-ref	<ref target="#"(zwt +1)"/>		

			La récupération de la page permet de supprimer le restant de la ligne qui ne nous sert pas.	
--	--	--	---	--

❖ *Script balisage sémantique*

	Étiquettes	Résultat XML TEI	Méthode et incertitudes	À ajouter
<b>add_utterance</b> (énoncés)	u	<u><seg>Texte</seg></u> >		<b>Ajouter :</b> @xml:id  <b>Ajouter plus tard :</b> @who @ana
	u-beginning	<u><seg>Texte</seg>		
	u-end	<seg>Texte</seg></u>		
<b>add_comment</b> (commentaires généraux)	comment	<note type="comment"> <seg>Texte</seg> </note>		<b>Ajouter :</b> @xml:id  <b>Ajouter plus tard :</b> @corresp pour les notes des rectifications
	comment-beginning	<note type="comment"> <seg>Texte</seg>		
	comment-end	<seg>Texte</seg> </note>		
	result	<note type="result">		<b>Ajouter :</b> @xml:id

		<seg>Texte</seg> </note>		
	opening	<note type="opening"> <seg>Texte</seg> </note>		
	closing	<note type="closing"> <seg>Texte</seg> </note>		
<b>add_incident</b>  <b>(commentaires atmosphères)</b>	incident	<incident><desc>( Texte)</desc></incident>	Utilisation d'un point de repère pour placer les balises sur le texte en question → parenthèses  <b>Attention :</b> certains incident ne sont pas annotés car le point de repère est manquant dans l'OCR.	
	incident-beginning	<incident><desc>		
	incident-end	</desc></incident>		
<b>add_quote</b>  <b>(citations)</b>	quote	<quote>texte</quote>	Utilisation d'un point de repère pour placer les balises sur le texte en question → guillemets  <b>Attention :</b> certaines quote ne sont pas annotées car le point de repère est manquant dans l'OCR.	

	quote-beginning	<quote><seg>texte	<b>Attention :</b> traitement particulier pour les seg couplés avec quote-beginning et quote-end car chevauchement de balises. Pour ces cas, le seg est traité dans la fonction add_quote.	
	quote-end	texte</seg></quote>		

❖ *Script balisage logique*

	Étiquettes	Résultat XML TEI	Méthode et incertitudes	À ajouter
<b>add_structure</b> (structure générale)	body	<text> <body> texte	Faut-il inclure le texte ? Oui car besoin de récupérer numéro de page	<b>Ajouter :</b> @ana sur text
	body1	<text> <body> <pb n="1"/>		<b>À inclure plus tard car traitement spécifique.</b>
	text	text	<b>Attention :</b> pas besoin	<b>À vérifier sur séances</b>

		<div> <div> </div> </div>	d'inclure le texte avant les balises, car le texte de la box en question n'est pas à inclure (pied de page contenant l'information sur l'imprimerie).	<b>sans parties complémentaires ni annexes</b>
	back	<div> <div> <div> </div> </div> </div>		
	text-back	<div> <div> <div> </div> </div> </div>	<b>Attention :</b> pas besoin d'inclure le texte avant les balises, car le texte de la box en question n'est pas à inclure (pied de page contenant l'information sur l'imprimerie).	
<b>add_division (divisions)</b>	part	<div> <div> <div> </div> </div> </div>		<b>Ajouter plus tard : @corresp</b>
	part1	<div> <div> <div> </div> </div> </div>		
	agenda	<div> <div> <div> </div> </div> </div>		

	appendices	</div> <div type="appendices"> <head>text</head>		
	part1-appendices	<div type="appendices"> <head>text</head>		<b>À vérifier sur séances sans parties complémentaires avant</b>
	erratum	</div> <div type="erratum"> <head> <label>text</label>		<b>À vérifier sur séance avec erratum</b>
	part1-erratum	<div type="erratum"> <head> <label>text</label>		<b>À vérifier sur séances avec erratum</b>
	lists	</div> <div type="lists"> <head> <label>text</label>		<b>À vérifier sur séances avec listes</b>
	part1-lists	<div type="lists"> <head> <label>text</label>		<b>À vérifier sur séances avec listes</b>
	offices	</div> <div type="offices"> <head>text</head>		<b>À vérifier sur séances avec nomination des bureaux</b>



	part1-offices	<div type="offices"> <head>texte</head>		<b>À vérifier sur séances avec nomination des bureaux</b>
	sitting	<div type="sitting">	<b>Attention :</b> obligation de traiter sitting et contents en même temps car impossible d'appliquer 2 traitements différents sur une même boxe.	
	other-sitting	</div> <div type="other-sitting"> <head>texte</head>		<b>À vérifier sur les CR contenant plusieurs séances</b>
	contents	<div type="contents"> <head>texte</head> <list>	<b>Attention :</b> obligation de traiter sitting et contents en même temps car impossible d'appliquer 2 traitements différents sur une même boxe.	
	voting1	<div> <div type="voting"> <head> <label>texte</label>		<b>Ajouter :</b> 1ere div → @xml:id  2e div → @xml:id  <b>Ajouter plus tard :</b> 2e div → @corresp

	voting	</div> <div type="voting"> <head> <label>text</label>		<b>Ajouter : @xml:id</b>  <b>Ajouter plus tard : @corresp</b>
	div-end	</div>		
	rectification	</div> </div> <div type="rectification"> <head>text</head>		<b>Ajouter plus tard : @corresp</b>
	petition	</div> <div type="petition"> <head> <label>text</label>		<b>À vérifier sur séance avec petition</b>
	part1-petition	<div type="petition"> <head> <label>text</label>		<b>À vérifier sur séance avec petition</b>
<b>add_structural_comment</b>  <b>(commentaires structurels)</b>	note	<note>  </note>		<b>Ajouter : @xml:id</b>
	voterslist-beginning	<note type="voterslist"> <desc>text</desc>		
	note-beginning	<note type="numbersannounced">text		

	note-end	</note>		
<b>add_item</b> <b>(item)</b>	item	<item>text</item>		<b>Ajouter plus tard :</b> @xml:id
	item-list	<item>text</item> </list>		
<b>add_title</b> <b>(titres et titres divisés)</b>	head	<head>text</head>		
	desc	<desc>text</desc>		
	note-head	<note>text</note> </head>		<b>Ajouter :</b> @xml:id
<b>add_table</b> <b>(tableau)</b>	table	<table> <row> <cell>Tout le texte du tableau</cell> </row> </table>		<b>À ajouter par la suite :</b> changer le balisage en fonction de l'analyse des tableaux à réaliser (garder <table> ou mettre <measure>)

### ❖ *Script compilation*

	Méthode	Problèmes
<b>compilation</b>	Appel de l'ensemble des fonctions permettant de baliser le contenu du CR.	

❖ *Script nettoyage*

	Méthode	Problèmes
<p><b>clean_xml</b></p> <p><b>(sauts de ligne)</b></p>	<p>\n : à remplacer par un espace            -\n : ne pas rajouter d'espace</p> <p><b>À prendre en compte :</b></p> <ul style="list-style-type: none"> <li>- Cas des mots divisés entre 2 boxes :</li> </ul> <p>→ Si c dû à un changement de colonne alors il faut recoller le mot.            → Si c dû à un changement de page, il faut garder le tiret car balise &lt;pb&gt; sera inclus dans le mot.            → À appliquer directement sur les fichiers XML.</p> <p><b>Attention :</b></p> <ul style="list-style-type: none"> <li>- Cas de l'enchaînement de deux boxes :</li> </ul> <p>→ Ajouter un espace, traitement à réaliser directement dans la boucle permettant d'écrire le texte dans le fichier XML.</p>	<p><b>Défaut de la méthode :</b> quand le mot est réellement divisé en deux, et qu'il est sur deux lignes, il est recollé (exemple grand-père).</p>
<p><b>delete</b></p>	<p>Utilisation de l'étiquette "useless" afin de supprimer toutes les boxes que nous ne</p>	

(suppression élément inutile)	<p>souhaitons pas intégrer dans l'élément "text" (l'en-tête par exemple).</p> <p>Certaines boxes sont supprimées via le traitement de d'autres étiquettes directement (bandeau supérieur traité avec la gestion des changements de page, pied de page avec la gestion des étiquettes "text" et "text-back".</p>	
-------------------------------	---	--

### ❖ *Script métadonnées*

	Étiquette	Résultat XML TEI	Méthode et incertitudes	À ajouter
<b>build_teiheader</b>  (construction à la main du header complété par les variables issues de var_metadata)		<code>&lt;date when="2022-07-25"/&gt;</code>	Utilisation de datetime :  <code>date.today().strftime("%Y-%m-%d")</code>	<b>Ajouter</b> : une gestion automatique des éléments <code>&lt;tagsDecl&gt;</code> et <code>&lt;extent&gt;</code>
<b>var_metadata</b>  (variables contenant le texte utile pour les métadonnées issus	date-pub	<code>&lt;date&gt;27 novembre 1889&lt;/date&gt;</code> <code>&lt;publicationStmt&gt;</code>		<b>Ajouter</b> : <code>@when</code>

de data)	meeting-session	<meeting n="E1" ana="#parla.lower #parla.session">Sessio n extraordinaire de 1889</meeting>	<b>Attention :</b> informations traitées dans une même condition car présentes sur la même ligne.	
	meeting-legislature	<meeting n="5L" ana="#parla.lower #parla.legislature">5e législature</meeting>		
	meeting-sitting	<meeting n="10" ana="#parla.lower #parla.sitting">10e séance</meeting>		
	date-sitting	<date>mardi 26 novembre<date>		

❖ *Script main*

	Méthode	Problèmes
<b>Chemin relatif</b>	Mettre un chemin relatif pour aller chercher les : - json (path_to_json) - xml (path_to_xml)	
<b>Variables</b>	Déclaration des variables utiles pour la création de l'élément racine :	

	<ul style="list-style-type: none"> <li>- beginning_elements</li> <li>- end_elements</li> </ul>	
<b>Boucle générale</b>	<ul style="list-style-type: none"> <li>- Ouverture des fichiers JSON et .load()</li> <li>- Gestion des changements de page</li> <li>- Appel des métadonnées (fonction build_teiheader)</li> <li>- Appel des données à inclure dans l'élément "text" (fonction compilation)</li> <li>- Sous-boucle permettant de créer un fichier XML pour chaque séance, que l'on nomme selon le nom du fichier json, dans lequel on écrit l'ensemble des éléments et pour lequel on nettoie les espace</li> </ul>	
<b>Gestion des inclusions dans le fichier corpus</b>	<p>Le fichier corpus doit être situé dans le dossier xml_data.</p> <p>Utilisation de la librairie lxml pour parser le fichier et inclusion automatique des éléments xi:include.</p>	<b>Défaut :</b> si l'on supprime un fichier XML composant, alors la suppression de l'inclusion doit se faire à la main.
<b>Amélioration des scripts : éléments à traiter</b>		
<b>Tests unitaires</b>		
<b>Vérifier la validité de l'XML</b>	<b>Solution provisoire :</b> la validation du	<b>Solution à long terme :</b> reprendre le

	<p>schema se fait manuellement pour l'instant. Le schema est précisé dans chaque en-tête des fichier XML, il faut ouvrir un éditeur XML et voir si la validation fonctionne.</p> <p><b>Solution à long terme</b> : utilisation de la librairie lxml pour valider le schéma RNG.</p>	code présent à la fin du fichier main et le modifier, car il ne fonctionne pas.
<b>Gestion différents OS</b>	<p>Utiliser le module pathlib :</p> <ul style="list-style-type: none"> <li>- <a href="https://docs.python.org/3/library/pathlib.html">https://docs.python.org/3/library/pathlib.html</a></li> <li>- <a href="https://www.digitalocean.com/community/tutorials/how-to-use-the-pathlib-module-to-manipulate-filesystem-paths-in-python-3-fr">https://www.digitalocean.com/community/tutorials/how-to-use-the-pathlib-module-to-manipulate-filesystem-paths-in-python-3-fr</a></li> </ul>	
<b>Vérification étiquettes</b>	Script à intégrer dans la chaîne de traitement.	
<b>Gestion de certaines métadonnées du fichier corpus</b>	<p>Mettre en place une construction automatique pour :</p> <ul style="list-style-type: none"> <li>- l'élément &lt;tagsDecl&gt;</li> <li>- l'élément &lt;particDesc&gt;</li> <li>- l'élément &lt;standOff&gt;.</li> </ul>	