

Encoder en TEI

Simon Gabay

Université de Neuchâtel

12 novembre 2018



Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Edition numérique

Une édition numérique

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Qu'est-ce qu'une édition numérique ?

- À une époque où tout le monde travaille sur un ordinateur, est-ce qu'utiliser le medium numérique, à un stade ou un autre (transcription, annotation, index. . .), est suffisant pour parler d'édition numérique ?
- Est-ce que diffuser une édition sur un medium numérique (pdf, clef USB, site internet. . .) est suffisant pour parler d'édition numérique ?
- (Cette question nous renvoie plus généralement au flou qui entoure la définition des humanités numériques.)

Une édition électronique savante

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Qu'est-ce qu'une édition électronique savante ?

- Une édition qui utilise les possibilités de l'outil informatique pour proposer. . .
- . . . une représentation critique ou enrichie sémantiquement d'une source. . .
- . . . des fonctionnalités utiles à leur compréhension et leur analyse.

La philologie numérique

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

L'édition numérique savante s'appuie sur la philologie numérique

- L'ecdotique (l'édition d'un texte) n'est qu'une partie du travail
- La philologie (numérique) s'intéresse aussi à la création des données (transcription) ou leur exploitation (exploitation) en plus de leur diffusion.

Vocabulaire de la philologie

Le travail de philologie ecdotique doit suivre quelques grandes étapes, dont chacune connaît des solutions et des procédures informatiques spécifiques.

- Transcription
- Collation
- Analyse paléographique et codicologique
- Analyse linguistique
- Annotation
- Indexation
- Publication
- (Exploitation)
- Archivage

Une communauté

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

L'ecodtique est un art complexe et ancien, dont des rudiments sont nécessaires pour éditer un texte, numériquement ou pas. L'utilisation d'outils numériques complexifie encore le travail.

Faire de l'édition numérique :

- C'est travailler en équipe
- C'est avoir recours à une communauté de recherche, qui partage des pratiques et une philosophie (cf. Humanistica, Cahier, EADH, DARIAH, Communauté TEI...)

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

XML-TEI

Le XML

- XML signifie *Extensible Markup Language*
- C'est un **langage** de balisage (vs langage de programmation, de définition de données ou de requête)
- Comme tout langage, il est régi par des règles

Règles principales

Le langage de balisage fonctionne de manière simple

```
<élément attribut="valeur">donnée</élément>
```

- ▶ Un `<élément>` est entre chevrons
- ▶ Une `<balise>` doit être fermé `</balise>`
- ▶ Une `<balise1>` ne doit `<balise2>` pas être croisé `</balise2>` avec un autre `</balise1>`
- ▶ Une `<balise/>` peut être auto-fermante
- ▶ Un `<élément>` peut porter un `@attribut` (noté avec un `@`)
- ▶ l'`@attribut` a une "valeur" (entre guillemets)

Du texte à la base de données

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Du texte à la base de données

- Une donnée est enfermée entre deux balises. Pour nous il s'agit de textes, de paragraphes, de phrases. . .
- Les données sont "emboîtées" les unes dans les autres : un document contient des paragraphes, qui contiennent des phrases, qui contiennent des mots. . .
- On transforme ainsi le texte en base de données

Une structure arborescente

Encoder en
TEI

Simon Gabay

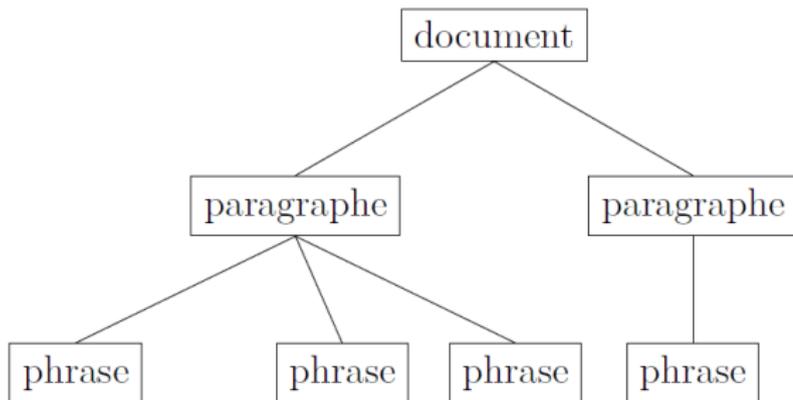
Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion



Document

On emploie a priori les italiques pour les termes empruntés à d'autres langues. On emploie les petites capitales pour les noms propres, comme Léopold Delisle. On emploie en revanche généralement le gras pour des raisons coupables. On retourne à la ligne pour un nouveau paragraphe.

XML comme langage structuré (par des balises)

En pratique, on obtient un document de ce type :

```
1 <document>
2 <paragraphe>
3   <phrase>
4     On emploie <locutionÉtrangère>a priori</locutionÉtrangère> les
5     italiques pour les termes empruntés à d'autres langues.
6   </phrase>
7   <phrase>
8     On emploie les petites capitales pour les noms propres,
9     comme <nom>Léopold Delisle</nom> ou <nom>Jules Quicherat</nom>.
10  </phrase>
11  <phrase>
12    On emploie en revanche généralement le gras pour des raisons
13    coupables.
14  </phrase>
15 </paragraphe>
16 <paragraphe>
17   <phrase>
18     On retourne à la ligne pour un nouveau paragraphe.
19   </phrase>
20 </paragraphe>
21 </document>
```

Une question fondamentale

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Une question fondamentale

- Nous avons ici utilisé `<paragraphe>` ou `<phrase>`, mais nous aurions pu choisir d'autres noms.
- Si nous étions italiens, nous aurions choisi `<paragrafo>` et `<frase>`
- Mais alors les documents sont encodés différemment: comment choisir des noms pour les `<éléments>` et les `@attributs` communs à tous?

La TEI

- ▶ TEI pour *Text Encoding Initiative*
- ▶ Elle est créée en 1987 (donc avant internet)
- ▶ La TEI est pilotée par un consortium qui maintient et développe des recommandations pour l'encodage des textes
- ▶ Ces recommandations sont en constantes évolutions
- ▶ Elles sont disponibles en ligne
<http://www.tei-c.org/guidelines/>

Entre vocabulaire et langage

- ▶ Il existe d'autres vocabulaires XML que la TEI, comme l'EAD pour les archivistes. . .
- ▶ . . . ou Dublin Core (DC) pour les bibliothécaires
- ▶ Ces vocabulaires peuvent d'ailleurs être exprimés avec d'autres langages (RDF-DC).
- ▶ Pour cette raison, on parle de XML-TEI, (ainsi il a existé un SGML-TEI).

Trois particularités

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Trois particularités de la TEI

- ▶ Le vocabulaire est en anglais : on utilise une balise `<w>` (*word*) pour un `<w>mot</w>`
- ▶ Ce vocabulaire est limité: on ne peut (presque) pas inventer de nouvelles balises
- ▶ Elle propose autant que possible un encodage sémantique (à l'inverse de LaTeX, utilisé pour faire ces diapos)

Sémantique et procédural

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

On emploie *a priori* les italiques pour les locutions et termes empruntés à d'autres langues.

Procédural

On emploie `<italique>a priori</italique>` les italiques pour les locutions et termes empruntés à d'autres langues.

semantique

On emploie `<locutionEtrangère>a priori</locutionEtrangère>` les italiques...

semantique II

On emploie `<latin>a priori</latin>` les italiques...

La solution en TEI

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

En XML-TEI, on encoderait le document ainsi :

La solution en TEI

On emploie `<foreign xml:lang="la">a priori</foreign>` les italiques...

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

**Avant toute
chose**

Pourquoi la
TEI ?

Conclusion

Avant toute chose

La première étape est la **Modélisation**

Définition

Opération par laquelle on établit le modèle d'un système complexe, afin d'étudier plus commodément et de mesurer les effets sur ce système des variations de tel ou tel de ses éléments composants.^a

a. J. Giraud, P. Pamart, J. Riverain, *Les nouveaux mots « dans le vent »*, Paris, France, 1974).

Il s'agit de définir un modèle adapté

- Aux documents que l'on édite
- À nos questions de recherche
- Aux moyens (techniques, financiers. . .) dont on dispose

Attention ! Il est souvent coûteux et compliqué de revenir sur certains choix. Il s'agit donc de bien réfléchir !

Concrètement, pour un philologue, les premières questions sont les suivantes ?

- Quels passages du textes doivent être balisés ? Les noms ? les locutions étrangères tous les mots ? Doit-on mettre la catégorie morpho-syntaxique et le lemme ?
- Doit-on représenter la structure physique du document (folios, pages. . .) ou la structure logique (chapitres, parties. . .)
- **Attention, il est (presque) impossible de tout faire : il faut choisir !**

Modéliser

Deux encodages différents d'un même texte Structure logique

```
1 <document>
2 <paragraphe>
3 <phrase>
4   On emploie <locutionÉtrangère>a priori</locutionÉtrangère> les
5   italiques pour les termes empruntés à d'autres langues.
6 </phrase>
7 <phrase>
8   On emploie les petites capitales pour les noms propres,
9   comme <nom>Léopold Delisle</nom> ou <nom>Jules Quicherat</nom>.
10 </phrase>
11 </paragraphe>
12 </document>
```

Structure physique

```
1 <document>
2 <pb n="1"/>
3   On emploie <locutionÉtrangère>a priori</locutionÉtrangère> les
4   italiques pour les termes empruntés à d'autres langues. On
5 <pb n="2"/>
6   emploie les petites capitales pour les noms propres, comme
7   <nom>Léopold Delisle</nom> ou <nom>Jules Quicherat</nom>.
8 </document>
```

La granularité

Définition

Degré de finesse ou précision d'un modèle, conçu comme le niveau de son plus petit composant. Plus la granularité est grande, plus on descend dans la modélisation (niveau phrase, mot, graphème, etc.) – et plus on ajoute de balises.

Modéliser

Deux encodages différents d'un même texte

Faible granularité

```
1 <document>
2 <paragraphe>
3 <phrase>
4   On emploie a priori les
5   italiques pour les termes empruntés à d'autres langues.
6 </phrase>
7 <phrase>
8   On emploie les petites capitales pour les noms propres,
9   comme <nom>Léopold Delisle</nom> ou <nom>Jules Quicherat</nom>.
10 </phrase>
11 </paragraphe>
12 </document>
```

Moyenne granularité

```
1 <document>
2 <paragraphe>
3 <phrase>
4   On emploie <locutionÉtrangère>a priori</locutionÉtrangère> les
5   italiques pour les termes empruntés à d'autres langues.
6 </phrase>
7 <phrase>
8   On emploie les petites capitales pour les noms propres,
9   comme <nom>Léopold Delisle</nom> ou <nom>Jules Quicherat</nom>.
10 </phrase>
11 </paragraphe>
```

Modéliser

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

On peut aller très loin

```
1 <document>
2 <paragraphe>
3 <phrase>
4 <w lemme="on" POS="PROper">On</w>
5 <w lemme="employer" POS="VERcjk">emploie</w>
6 ...
```

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

**Pourquoi la
TEI ?**

Conclusion

Pourquoi la TEI ?

Défauts de la TEI

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

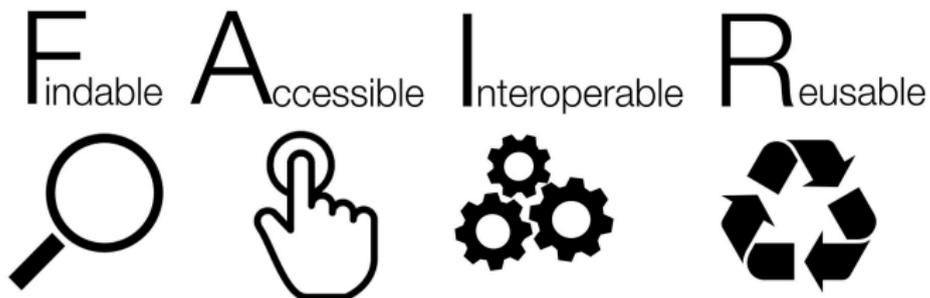
La TEI pose des problèmes

- ▶ Elle force à utiliser un standard, par définition générique, et qui ne convient pas forcément exactement à nos données
- ▶ Elle nécessite un apprentissage, notamment pour respecter le sémantisme du vocabulaire

Alors pourquoi la TEI ?

Fair Data

La TEI permet de répondre (en partie) aux principes du Fair Data



(Img. SangyaPundir CC BY-SA 4.0)

Les principes soutenus par les États et agences de recherche (ANR, ERC, FNS, DFG...)

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Les quatre grands principes du Fair Data :

- ▶ *Findable* : Les données sont faciles à trouver en ligne, y compris sur le long terme (archivage pérenne) ;
- ▶ *Accessible* : Les métadonnées accompagnent les données, y compris quand ces dernières ne peuvent être diffusées (*open science*)
- ▶ *Interoperable* : Le format doit être ouvert, libre, documenté et compréhensible ;
- ▶ *Reusable* : La modalités techniques et légales de réutilisation sont claires.

Bonnes pratiques

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

La TEI encourage (force ?) à avoir des bonnes pratiques

- ▶ Utilisation de métadonnées (des données sur les données)
- ▶ Identification claire des auteurs et de leurs responsabilités
- ▶ Propose une documentation rigoureuse

Format pivot

Encoder en
TEI

Simon Gabay

Edition
numérique

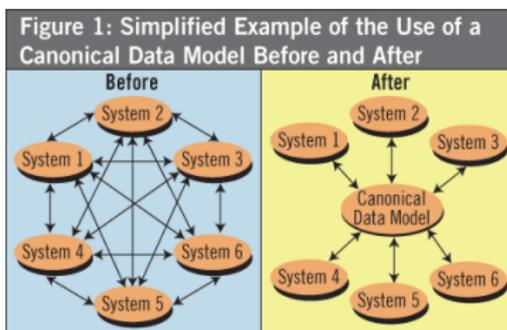
XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

La TEI est un format pivot



extr. de Steve Hoberman, « Canonical Data Model », *Information Management Magazine*, en ligne : http://www.information-management.com/issues/2007_50/10001733-1.html.

Chaîne de traitement

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

La TEI est une étape de la chaîne de traitement

- ▶ La chaîne de traitement numérique des documents (*Workflow*) est une suite de tâches ou opérations
- ▶ Elle implique l'utilisation de différentes solutions informatiques, par exemple :
 - XML-TXM pour l'analyse linguistique
 - HTML pour la publication web
 - LaTeX pour l'édition papier ou pdf
 - ...
- ▶ LA TEI n'est pas la réponse à tout
- ▶ Mais une attention particulière est portée à son intégration dans la chaîne de traitement

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Conclusion

Conclusion

Encoder en
TEI

Simon Gabay

Edition
numérique

XML-TEI

Avant toute
chose

Pourquoi la
TEI ?

Conclusion

Conclusion

- ▶ La TEI n'est qu'un (petit) morceau de la philologie numérique
- ▶ La philologie numérique est avant tout de la philologie
- ▶ La philologie numérique nécessite des compétences en humanités numériques
- ▶ D'autres formations existent