

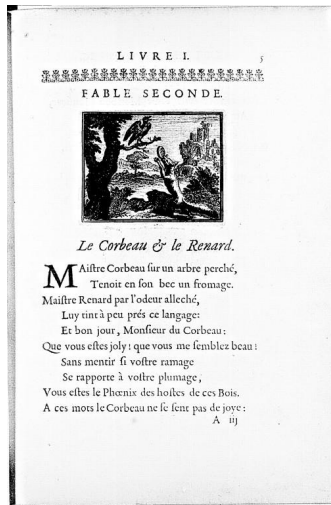
Construction des corpus pour l'analyse lexicale

Lou Burnard Consulting

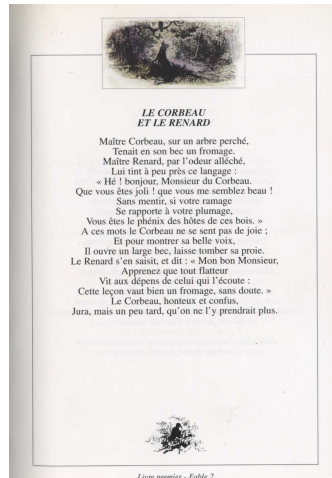
Objectifs de cette presentation

- ① Préciser ce que c'est que l'encodage textuel
- ② Un cas d'étude: le *Petit Comtois*
- ③ Présenter quelques concepts fondamentaux de TEI et d'XML

Est-ce que ceux-ci représente la meme chose ?



Source gallica.bnf.fr / Bibliothèque nationale de France



Un texte n'est pas un document...

En quoi consiste l'essentiel d'un texte ?

- en l'apparence des lettres et leur mise-en-page?
- en la version originelle (pretendue) de cette copie?
- en les interpretations/lectures apportées ou trouvées? en les intentions (supposées) de son auteur?

Un "texte" est quelque chose d'abstrait: la construction d'une communauté de lecteurs.

L'encodage explicite cette abstraction à fin de la mieux gérer

Un texte n'est pas un document...

En quoi consiste l'essentiel d'un texte ?

- en l'apparence des lettres et leur mise-en-page?
- en la version originelle (pretendue) de cette copie?
- en les interpretations/lectures apportées ou trouvées? en les intentions (supposées) de son auteur?

Un "texte" est quelque chose d'abstrait: la construction d'un communauté de lecteurs.

L'encodage explicite cette abstraction à fin de la mieux gérer

Un texte n'est pas un document...

En quoi consiste l'essentiel d'un texte ?

- en l'apparence des lettres et leur mise-en-page?
- en la version originelle (pretendue) de cette copie?
- en les interpretations/lectures apportées ou trouvées? en les intentions (supposées) de son auteur?

Un "texte" est quelque chose d'abstrait: la construction d'un communauté de lecteurs.

L'encodage explicite cette abstraction à fin de la mieux gérer

Un texte n'est pas un document...

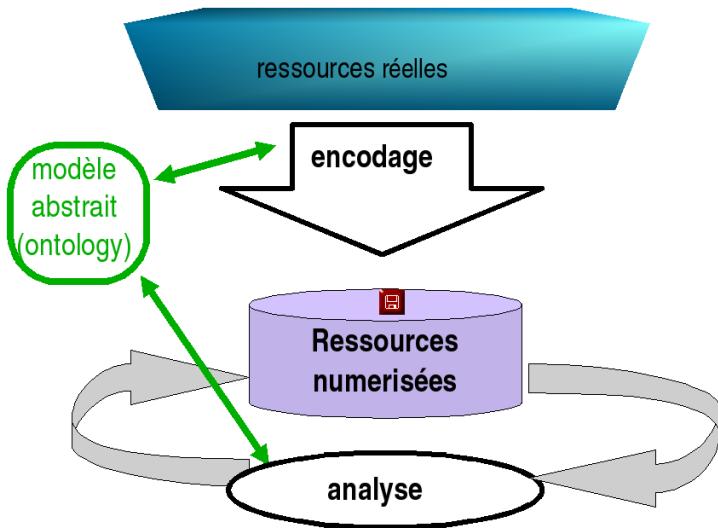
En quoi consiste l'essentiel d'un texte ?

- en l'apparence des lettres et leur mise-en-page?
- en la version originelle (pretendue) de cette copie?
- en les interpretations/lectures apportées ou trouvées? en les intentions (supposées) de son auteur?

Un "texte" est quelque chose d'abstrait: la construction d'une communauté de lecteurs.

L'encodage explicite cette abstraction à fin de la mieux gérer

Qu'est-ce qu'on fait en numérisant un texte?

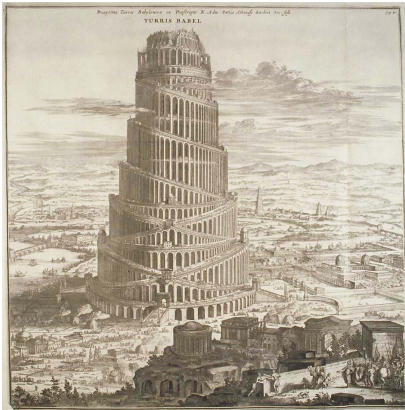


L'encodage

- Un texte est plus qu'une séquence de caractères encodés!
- Un text est plus qu'une séquence de formes lexicaux!
 - Il a une **structure** et une **signification**
 - Un texte peut avoir plusieurs **lectures** variantes
 - La portée d'un texte peut être **enrichie** par des annotations
- L'encodage explicite les lectures
- Sans explicitation, on ne peut rien traiter

L'effet Babel

Bien sûr il existe plusieurs lectures possibles pour la plupart des textes...



... et (malheureusement) plusieurs manières d'expression pour ces lectures!

Par exemple...

I

Loomings

Call me Ishmael. Some years ago – never mind how long precisely – having little or no money in my purse, and nothing particular to interest me on shore, I thought I would sail about a little and see the watery part of the world. It is a way I have of driving off the

Encodage ou babel?

```
|chap1
<C 1> Loomings
\chapter[1]{Loomings}
:h1.1. Loomings
MOBY001001LOOMINGS
|C1
.chapter Loomings
.cp;.sp 6 a;.ce .bd 1. Loomings
<h1>Loomings</h1>
<p class="h1">Loomings
~x
```

- Bonne nouvelle: il existe des logiciels capables de traduire entre 500 formats divers
- Mauvaise nouvelle: on en a besoin

Encodage ou babel?

```
|chap1
<C 1> Loomings
\chapter[1]{Loomings}
:h1.1. Loomings
MOBY001001LOOMINGS
|C1
.chapter Loomings
.cp;.sp 6 a;.ce .bd 1. Loomings
<h1>Loomings</h1>
<p class="h1">Loomings
~x
```

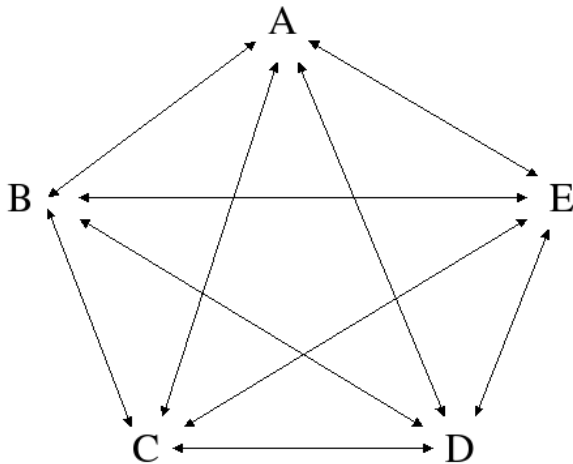
- Bonne nouvelle: il existe des logiciels capables de traduire entre 500 formats divers
- Mauvaise nouvelle: on en a besoin

Encodage ou babel?

```
|chap1
<C 1> Loomings
\chapter[1]{Loomings}
:h1.1. Loomings
MOBY001001LOOMINGS
|C1
.chapter Loomings
.cp;.sp 6 a;.ce .bd 1. Loomings
<h1>Loomings</h1>
<p class="h1">Loomings
~x
```

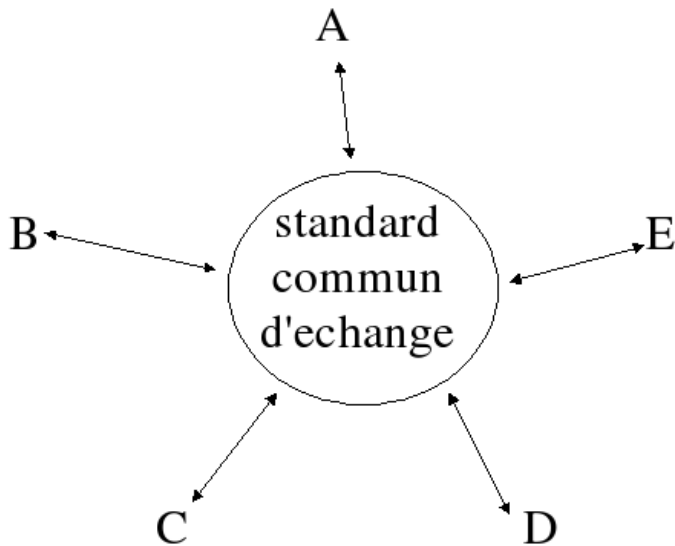
- Bonne nouvelle: il existe des logiciels capables de traduire entre 500 formats divers
- Mauvaise nouvelle: on en a besoin

Echange d'informations (1)



*sans format pivot: 20 passerelles requises ($n*n-n$)*

Echange d'informations (2)



Objectif d'un format pivot

- permettre la representation de toute information requise
- d'une maniere independente
 - de logiciel specifique
 - de plateforme specifique
 - d'usage specifique
- sans hausse improbable d'effort

Un objectif assez classique, qui s'exprime avec un **balisage** du texte

Objectif du balisage

Le balisage constitue un langage d'encodage, qui devrait...

- spécifier les caractères d'un texte
- expliciter la/les structures aperçue/s dans un texte
- linéariser le texte
- spécifier les méta-informations, renseignements contextuels etc.

Le balisage: un technique bien-compris

- Un balisage explicite les distinctions qu'on désire faire en traitant une chaîne de caractères
- Le balisage est une manière de nommer et de caractériser les composants d'une structure textuelle, d'une manière quasiment formelle
- Quel genre de composants va-t-on baliser?

Qu'est ce qu'on balisera?

Comparer:

```
<lb/>LIVRE I.  
<lb/>FABLE SECONDE.  
<lb/>Le Corbeau & le Renard  
<lb/>MAistre Corbeau sur un arbre perché,  
<lb/>Tenoit en son bec un fromage.  
<lb/>Maistre Renard par l'odeur alleché,  
<lb/>Luy tint à peu près ce langage: ...  
<lb/>A ces mots, le Corbeau ne se sent pas de joye:  
<lb/>A iii
```

avec

```
<div type="fable" n="I2">  
  <head>Le Corbeau & le Renard</head>  
  <l>  
    <hi rend="lettrine">M</hi>aistre Corbeau sur un arbre perché,</l>  
  <l>  
    <reg>Tenait</reg> en son bec un fromage.</l>  
  <l>Maistre Renard par l'odeur alleché,</l>  
  <l>  
    <reg>Lui</reg> tint à peu près ce langage:</l>  
<!-- ... -->  
  <l>A ces mots, le Corbeau ne se sent pas de joye:</l>  
  <fw type="sig">A iii</fw> ...  
</div>
```

Ou bien

```
<lg>
  <milestone unit="voix" n="narr"/>
  <l>Maître corbeau, sur un arbre perché</l>
  <l>Tenait en son bec un fromage.</l>
  <l>Maître renard par l'odeur alléché</l>
  <l>Lui tint à peu près ce langage:</l>
  <l>
    <milestone unit="voix" n="renard"/>Hé! bonjour Monsieur du Corbeau</l>
  <l>Que vous êtes joli! que vous me semblez beau!</l> ...
</lg>
```

ou bien

```
<lg>
  <sp who="#narr">
    <ab>
      <lb/>Maître corbeau, sur un arbre perché
      <lb/>Tenait en son bec un fromage.
      <lb/>Maître renard par l'odeur alléché
      <lb/>Lui tint à peu près ce langage:</ab>
    </sp>
  <sp who="#renard">
    <ab>
      <lb/>Hé! bonjour Monsieur du Corbeau
      <lb/>Que vous êtes joli! que vous me semblez beau! ... </ab>
    </sp>
  </lg>
```

N° 6 — Première année

CINQ CENTIMES

Lundi 6 août 1883

LE PETIT COMTOIS

JOURNAL RÉPUBLICAIN DÉMOCRATIQUE QUOTIDIEN

RÉDACTEUR EN CHEF, JULES GROS

ABONNEMENTS

	1989	1990	1991
France et colonies	12,1	12,1	12,1
Pays étrangers	15,1	15,1	15,1
Total hexagonal	27,2	27,2	27,2

RÉDACTION ET ADMINISTRATION

BESANCON. — 7: Semore Saint-Amour, 7 — BESANCON

Т.бул айырымдардын жарыгында 4-мартта 1941-жылдын 1-апрелине чейин

INSERTIONS

Ammoniac, 100 grs. 500
 Nivellures, 100 grs. 500
 Poudre d'Alun, 100 grs. 500

— **Service d'été.** — 1 heure de la soirée au début de 10 minutes.

Géant de Bassano: Point Dûn, Bassano, 3 h. 45, 8 h.
Bois, 2 h. 10, 5 h. 20, 8 h. 24.

Pour Boivin & Martin, 5 h. 22, 5 h. 22, 8 h. 00; Pour, 2 h. 35.
 Pour Jimmy & Martin, 2 h. 45; Boiv, 2 h. 44, 4 h. 40 pour Monchaud et son.

son: Maria, 1 p. 52; | Four Great War Mitzvah: Ma-
 9, 10, 40; Rev. 20, 24, | No. 11. 10; Son: 2 p. 40

Beobachtung : Mai, 3 h. 15 : h. 21, 90.

De Hout : 2 h. 25.	De Hout : 2 h. 25.
De Hout : 2 h. 25.	De Hout : 2 h. 25.

and r. Martin, 9 h, 24, +40 Grey; 30.5 h, 24, 30.5

LA POSITION COLONIALE

Lorsque suivait le duel, deux dignitaires, de son ministre (Gortchakoff), et *le Russe se soulevait* — après ses désastres de Crimée, et que, renfermé dans sa chambre chez lui dans l'attente de la mort, il se soulevait et se levait, quoique son allié respecté, monarque, au contraire, lui eût adressé ses vœux dans sa grande courtoisie aristocratique et s'abandonnait tout entier à sa magnifique mouvement d'expansion coloniale, et que, dans ce moment, le principe puissance européenne de la Russie, alors il ne manquait pas chez elle d'hommes qui — comme aujourd'hui chez nous — blâmaient cette politique, prêchaient l'absorption et dénonçaient comme une infamie tout empire

Valées terribles ! L'avenir, en nous montrant les précédents, nous avertit, s'il est permis, contre eux, et l'histoire officielle de la Russie qui, après s'être retirée dans l'ombre immense de la nuit, et du la colonisation de l'Asie, se réveille, s'est relevée depuis de ses désastres européens, nous démontre tout suffisamment que ces hommes livrés à qui dévouera à leur propos « le peuple d'aujourd'hui » sont eux-mêmes les

l'absence de motivations et des défaillances du Parlement, qui ont fait compromettre irrémédiablement nos deux expéditions de Tunisie et de l'Algerie, comme elles nous ont fait perdre l'Algérie. C'est pourquoi nous ne devons pas nous laisser aller à penser que dans une politique d'insécurité d'expressions politiques, et deux ans après cela, tout va bien le dire, « pas de doute ».

Enfin, adjoint aux dix maîtres à penser, et sans avoir de leur attitude, ne sommes-nous pas en train d'élaborer ce vaste empire colonial à partir d'un mot nous avons déjà entrepris les bases par le développement de la colonisation algérienne, par l'occupation de l'Algérie, par la Tunisie, par le Maroc, par les pays méditerranéens, les Indes, le Soudan, le Congo, et par la conquête de Madagascar, cette « France équinoxiale » dont nous sommes en train, nous le savons, de servir l'implémentation comme rose.

Nous ne passons pas un anneau escarpé des Andes, dont nous avons tout récemment cessé la poursuite et déjà presque réalisé la conquête sur le continent asiatique, et nous arrivons à la ville de décapiter notre influence et notre puissance coloniale par l'adjonction de deux vastes et riches provinces, l'Araucan et le Tenkin, à nos belles provinces de Cochinchine et du Cambodge !

vous avons suivie sur le terrain colonial.

Malheureusement guidée dans cette oeuvre de vulgarisation par la préoccupation de faire campagne, avec l'ingénierie dont notre indépendance est la meilleure garantie, ce que nous croyons être la vérité sur notre politique coloniale, nous estimons supérieur de répondre aux arguments sans portée forcée de conjectures et d'hypothèses et de chercher la solution de problèmes dont l'avenir seul dira le succès.

Qu'importe, en effet, de savoir si, pour servir cette politique, nous encourageons la jalousie de l'Angleterre et les tentatives de sa réprobation, qu'elle ne mène jamais à tout voisin qui s'enrichit.

Il s'agit de savoir si, prenant en compte, abandonnant, cédant quatre milliards de la conquête du globe à des associations anglo-saxones, nous rappellerons en France ces hommes de courage et du courage : Bruma, Lucio et Piero, Bergami-Darbois, tous ces luttards pieux qui, portant avec nous notre langue et notre civilisation, ont accordé hier sur tous les points du globe l'émancipation du pays.

Il s'agit de savoir si, arrivés au terme
de notre développement et, de notre ex-
pansion humaine nous ne souhaitons
explorer dès à présent toute possibilité

mande, et si, conscients de notre vitalité comme race, nous nous efforçons encore au cœur l'autrui, de longs espoirs et de vastes pensées.

ЖАҢА ДӘҢГЕЛ СӨЗІНДІГІ:

DÉPÊCHES DE NUIT

Service de votre entreprise est-il satisfaisant?

INFORMATIONS GÉNÉRALES

Los magistrados des commissions

La nouvelle loi sur la magistrature a pu pour elle d'émanciper les magistrats en leur faisant perdre les commissions mixtes et leur enlever ainsi leur activité de service.

M. Camarero, conseiller à la cour de cassation et père du préfet de police, est, comme par le fait sous l'application de la loi, M. Camarero a demandé sa mise en retraite et il l'a obtenue, mais la disposition qui lui était faite a sa demande ne lui confère aucune pension.

Accident & Philhellenism

Un «*star*» Hard'ien, professeur au lycée de Phéligueville, ayant son premier signe dans une feuille locale des articles toujours contre les colons de l'Algérie, a été la cause d'une manifestation tumultueuse.

Un petit exercice intellectuel...

Imaginez vous: on a des milliers des pages à baliser....

- Quels traits existent dans ces matériaux?
- Quels traits faut-il baliser? lesquels seraient (in)utiles?
- Comment justifier l'étiquetage choisi?
- Comment garantir la consistance et l'étendue du balisage?

Maintenant, on réduit le budget de 50%. Répétez l'exercice!

Par exemple...

On commence par noter la structuration:

- la page contient des entêtes, des colonnes, des lignes, etc.
- le journal contient des titres, des regroupements de notices, une episode de feuilleton, des notices de publicité, etc.
- dans le texte du journal, il y a partout des noms d'autre journaux, de personnages, des lieux, etc.
- le texte est (principalement) en français, d'un style d'intérêt historique
- on note également des références aux événements historiques
- et on a bien-sûr des infos supplémentaires regardant la production, la dissémination, la bibliographie etc. de cette source...

Focalisons...

taurien, recteur de l'Académie d'Alger, a demandé la révocation de ce professeur imprudent.

L'organisation de l'artillerie de forteresse.

Le général Tricoche, directeur de l'artillerie au ministère de la guerre, sera nommé général de division et chargé de l'organisation des batteries d'artillerie de forteresse, avec le titre d'inspecteur général.

Il sera remplacé dans ses fonctions au ministère par le général Lavocat.

Tremblement de terre en Grèce.

Aujourd'hui, à deux heures du matin, un fort tremblement de terre a été ressenti au Pirée.

On peut se passer des balises ?

Un journal algérien annonce que M. Constantin, recteur de l'académie d'Alger, a demandé la révocation de ce professeur imprudent. Le général Tricoche, directeur de l'artillerie au ministère de la guerre, sera nommé général de division et chargé des l'organisation des batteries d'artillerie de forteresse, avec le titre d'inspecteur général. Il sera remplacé dans ses fonctions au ministère par le général Lavocat. Aujourd'hui, à deux heures du matin, un fort tremblement de terre a été ressenti au Pirée. On télégraphie d'Athènes qu'aucun accident de personnes n'est à déplorer.

- on ne mets en evidence que les mots des notices
- il n'y a pas de frontiere entre les deux notices parce que les parties structurantes sont supprimées
- la ponctuation est retenue telle quelle, y compris les liaisons
- tous les mots semblent pareilles en terme de leur fonctionnement dans le texte

Balisage ASTARTEX

```
<ASTX_ATTR SOURCE=PETIT_COMTOIS DATE=1883-08-06 NUMERO=0006
PAGE=1 LARGEUR=1 EMBLACEMENT=4M
RUBRIQUE=DEPECHES_DE_NUIT_INFORMATIONS_GENERALES
RUB_GR=DEPECHES_DE_NUIT
TITRE=L'organisation de l'artillerie de forteresse
SIGNATURE=Service_de_notre_correspondant_spezial. > Le general
Tricoche, directeur de l'artillerie au ministere de la guerre, sera
nomme general de division et charge des l'organisation des batteries
d'artillerie de forteresse, avec le titre d'inspecteur general. Il
sera remplace dans ses fonctions au ministere par le general Lavocat.
<ASTX_ATTR SOURCE=PETIT_COMTOIS DATE=1883-08-06 NUMERO=0006 PAGE=1_2
LARGEUR=1 EMBLACEMENT=4M_1H
RUBRIQUE=DEPECHES_DE_NUIT_INFORMATIONS_GENERALES
RUB_GR=DEPECHES_DE_NUIT
TITRE=Tremblement_de_terre_en_Grece.
SIGNATURE=Service_de_notre_correspondant_spezial. >
Aujourd'hui, a deux heures du matin, un fort tremblement de terre a
ete ressenti au Pirée. On telegraphie d'Athenes qu'aucun accident de
personnes n'est a deplorer.
```

- on ne mets en evidence que les mots des notices
- les parties structurantes sont transformées en annotation
- les annotations s'expriment comme proprietés des mots

Un balisage TEI

```
<div type="groupe" decls="DDN" n="4">
  <head type="RUBRIQUE">DEPECHEs DE NUIT</head>
  <head type="SIGNATURE">Service de notre correspondant spécial.</head>
  <head type="RUBRIQUE">INFORMATIONS GENERALES</head>
  <div type="NOTICE" rend="EMP.4M LG.1"
    n="4">
    <head type="TITRE">L'organisation de l'artillerie de forteresse</head>
    <p>Le général Tricoche, directeur de l'artillerie au ministère de la guerre, sera
      nommé général de division et chargé des l'organisation des batteries d'artillerie de
      forteresse, avec le titre d'inspecteur général. </p>
    <p>Il sera remplacé dans ses fonctions au ministère par le général Lavocat.</p>
  </div>
  <div type="NOTICE" rend="EMP.4M LG.1"
    n="5">
    <head>Tremblement de terre en Grèce.</head>
    <p>Aujourd'hui, à deux heures du matin, un fort tremblement de terre a été ressenti au
      Pirée.</p>
    <p>On télégraphie d'Athènes qu'aucun accident de personnes n'est à déplorer.</p>
  </div>
</div>
```

A noter...

- on met en evidence les mots des notices mais également les objets qu'ils constituent (paragaphes, titres...)
- (une selection des) objets textuels sont balisés explicitement avec leur fonction
- les annotations s'expriment comme propriétés des objets auxquels elles sont attachées
- on se sert de XML, un syntaxe standardisé

XML: ce que c'est et pourquoi on devrait le connaître

- XML est une manière de représenter les **données structurées** en forme de chaîne de caractères
- XML fournit une syntaxe tout à fait générique, à appliquer dans toute domaine de ressources numérisées
- un document XML doit être **bien formé**
- un document XML peut être **valide**
- XML est indépendant de l'application, de la plateforme et du vendeur
- XML rend le pouvoir aux fournisseurs de données, et facilite l'intégration des ressources diverses et polyglottes

Syntaxe XML

Un document XML contient:-

- des *éléments*, qui portent (facultativement) des *attributs*, marqués par *balises*
- des *commentaires*
- des *instructions de traitement*
- des *references à entité* (interne ou externe)
- des **sections CDATA**
- ...et des caractères Unicode

C'est tout!

XML: règles du jeu

- Un document XML représente une arborescence composée de **noeuds**
- il y a un seul noeud racine qui contient tous les autres
- chaque noeud peut être
 - une arborescence
 - un **élément** (qui porte facultativement des **attributs**)
 - une chaîne de **caractères**
- Chaque élément porte un nom ou **identification générique**
- Chaque attribut porte un nom et une valeur
- les noms sont liés avec un **namespace** (espace de noms)

Donc... on a des contenants emboîtés, avec des propriétés.

Representation d'une arborescence XML

- Un document XML linéarisé commence par une instruction de traitement special
- Les occurrences d'élément sont marqués entre **balises ouvrantes** et **balises fermantes**
- Les caractères < et & sont Magiques et doivent être cachés au moyen de références entité (< et & respectivement)
- Les paires nom/valeurs qui constituent les attributs d'un élément peuvent apparaître sans ordre à l'intérieur d'une balise ouvrante
- L'espace de noms auquel appartient un élément peut être signalé par un **namespace-prefix** (p.e. xml:) prédéfini

Syntaxe XML: le "fine print"

Pour qu'un document soit *bien formé*, il faut que:

- ① une seule racine contienne le document entier
- ② chaque arborescence soit proprement imbriquée
- ③ tout les noms soient sensibles à la casse
- ④ chaque balise ouvrante ait sa balise fermante (sauf qu'on peut combiner les deux, le noeud étant vide)
- ⑤ les valeurs d'attribut soient présentées correctement entre guillemets

! pas de chevauchement

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Bien formé? Oui ou non?

- `<seg>some text</seg>`
- `<seg><foo>some</foo> <bar>text</bar></seg>`
- `<seg><foo>some <bar></foo> text</bar></seg>`
- `<seg type="text">some text</seg>`
- `<seg type='text'>some text</seg>`
- `<seg type=text>some text</seg>`
- `<seg type = "text">some text</seg>`
- `<seg type="text">some text<seg/>`
- `<seg type="text">some text<gap/></seg>`
- `<seg type="text">some text< /seg>`
- `<seg type="text">some text</Seg>`

Validation XML

Un document XML *valide* est (bien sûr) bien formé, et en plus conforme à des règles supplémentaires, qui constituent un *schéma*

Un schéma peut spécifier:

- le nom de l'élément racine
- les noms de tous les éléments légaux
- les noms et les types des attributs
- des règles concernant l'imbrication et le contenu des éléments
- et quelques autres menus propos...

n.b. Un schéma ne spécifie point la signification sémantique des éléments

Continuons notre balisage TEI ...

```
<p>
  <s n="42">
    <w ana="#NOM" lemma="XXXX">0n</w>
    <w ana="#VER" lemma="télégraphier">télégraphie</w>
    <w ana="#PRP" lemma="de">d'</w>
    <w ana="#NOM" lemma="XXXX">Athènes</w>
    <w ana="#KON" lemma="que">qu'</w>
    <w ana="#PRO" lemma="aucun">aucun</w>
    <w ana="#NOM" lemma="accident">accident</w>
    <w ana="#PRP" lemma="de">de</w>
    <w ana="#NOM" lemma="personne">personnes</w>
    <w ana="#ADV" lemma="ne">n'</w>
    <w ana="#VER" lemma="être">est</w>
    <w ana="#PRP" lemma="à">à</w>
    <w ana="#VER" lemma="déplorer">déplorer</w>
    <c ana="#SENT">.</c>
  </s>
</p>
```


Encore de balisage TEI

```
<p>  
  <persName key="#TRIC01">Le général Tricoche</persName>,  
  directeur de l'artillerie au  
  ministère de la guerre, sera nommé général de division et  
  chargé des l'organisation des  
  batteries d'artillerie de forteresse, avec le titre  
  d'inspecteur général.  
</p>  
<p>Il sera remplacé dans ses fonctions au ministère par  
<persName ref="#LAVOC32">le  
  général Lavocat</persName>.</p>  
<!-- ... -->  
<person xml:id="LAVOC32">  
  <persName>  
    <forename>Jean-Louis</forename>  
    <surname>Lavocat</surname>  
  </persName>  
  <birth when="1823-02-09">ne le 9 fevrier 1823 a  
  Besancon</birth>  
<!-- ... -->  
</person>  
<relation type="remplacement"  
  active="#LAVOC32" passive="#TRIC01"/>
```

Et encore de balisage TEI...

```
<teiHeader>
  <fileDesc>
    <titleStmt>
      <title>Le Petit Comtois: TEI-XML Edition</title>
    </titleStmt>
    <publicationStmt>
      <p>Prepared for MISAT</p>
    </publicationStmt>
    <sourceDesc>
      <bibl>
        <title>Le Petit Comtois: Journal Républicain démocratique quotidien
      </title>
      <ref>http://laseldi.univ-fcomte.fr/petit\_comtois/archives.php</ref>
      </bibl>
    </sourceDesc>
    <encodingDesc>
      <classDecl>
        <taxonomy xml:id="notice-class">
          <category xml:id="DDN">
            <catDesc>Dépêches de nuit</catDesc>
          </category>
          <category xml:id="CR">
            <catDesc>Chroniques régionales</catDesc>
          </category>
        </taxonomy>
      </classDecl>
    </encodingDesc>
    <revisionDesc>
      <change when="2009-06-15">Taxinomies révisées</change>
    </revisionDesc>
  </fileDesc>
</teiHeader>
```

Le paysage actuel de la TEI

- Structuration basique des textes continus
- Transcription diplomatique, images, multimédia, annotations...
- Données formelles : dates, noms de lieux ou de personnes...
- Données paratextuelles et "meta"
- Analyses linguistiques à tout niveau (y compris l'oral)
- Documentation de balisage
- Et cetera : voir

<http://www.tei-c.org/P5/Guidelines/>

... Bref : une sorte d'encyclopédie du balisage !

Composants de TEI ALL: les modules TEI P5

nom	chapitre P5
analysis	Simple Analytic Mechanisms
certainity	Certainty and Responsibility
core	Elements Available in All TEI Documents
corpus	Language Corpora
dictionaries	Dictionaries
drama	Performance Texts
figures	Tables, Formulae, and Graphics
gaiji	Representation of Non-standard Characters and Glyphs
header	The TEI Header
iso-fs	Feature Structures
linking	Linking, Segmentation, and Alignment
msdescription	Manuscript Description
namesdates	Names, Dates, People, and Places
nets	Graphs, Networks, and Trees
spoken	Transcriptions of Speech
tagdocs	Documentation Éléments
tei	The TEI Infrastructure
textcrit	Critical Apparatus
textstructure	Default Text Structure
transcr	Representation of Primary Sources
verse	Verse

Les TEI Guidelines

Dans ses 1 400 pp imprimées, vous trouverez :

un lexique et une grammaire 22 'modules' regroupant en total 564 éléments qui sont d'ailleurs classifiés en 189 classes

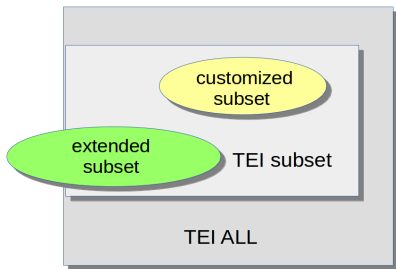
des contraintes additionnelles 31 types de données, 63 contrainte/regles Schematron

des conventions d'utilisation *beaucoup* de prose

plusieurs exemples d'utilisation dont au moins un par élément

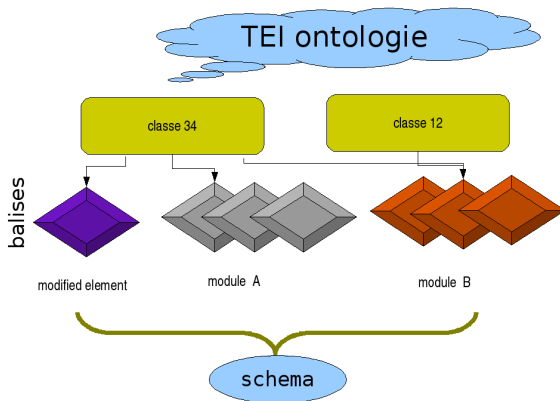
Comment régler ces richesses?

La TEI est conçue pour soutenir une variété d'approches



- on peut simplement utiliser un sous-ensemble de ses propositions (TEI subset)
- on peut y ajouter des contraintes supplémentaires (customized subset)
- on peut y ajouter de nouveaux composants (extended subset)

Architecture de la TEI



1

L'esprit TEI

Qu'est-ce que cela veut dire : utiliser la TEI ?

- une pratique de balisage consensuelle
- un lexique commun
- un respect de l'autonomie

La standardisation ne devrait pas signifier « fais comme moi » mais plutôt « explique-moi ce que tu fais. »

Etre conforme à la TEI veut dire quoi?

- **être honnet** : Les éléments XML qui se déclarent comme appartenant au namespace TEI doivent respecter les définitions TEI de ces éléments
- **être explicite** : Pour valider un document TEI, un ODD est fortement conseillé, parce que cela mettra en évidence toutes les modifications effectuées
- Tout document valide TEI est valide par rapport au schéma TEI-ALL

L'objet de ces règles est de faciliter le "blind interchange" des documents.

"One Document Does it all (ODD)"

- Les Guidelines TEI sont produits à partir d'une même ressource XML qui contient:
 - de la prose descriptive (une grande quantité)
 - des exemples d'utilisations (plusieurs)
 - des déclarations formelles pour les constituants du modèle abstrait de TEI
 - les éléments et leurs attributs
 - les modules
 - les classes, et les macros
- On appelle cette ressource un ODD (bien qu'elle consiste en des centaines de petits fichiers)

Et alors?

- Vu ses visées ambitieuses, on ne peut se servir du système TEI qu'en le personnalisant
- Les personnalisations s'expriment également en langue ODD
- Cela permet de communiquer la manière où on a appliqué la TEI, sans perturber les normes

Exemple...

Un fichier TEI ODD est un document TEI comme les autres, avec autant de prose discursif qu'il faut, mais en plus un élément spécialisé: un `<schemaSpec>` qui définit le schéma documenté

```
<text>
  <body>
    <div>
      <head>Notre manuel d'encodage</head>
      <p>Dans ce projet on ....</p>
      <schemaSpec ident="myTei"
        start="TEI">
        <moduleRef key="tei"/>
        <moduleRef key="header"/>
        <moduleRef key="core"
          except="sp stage speaker"/>
        <moduleRef key="textstructure"/>
        <elementSpec ident="head"
          mode="change" module="core">
          <content>
            <textNode/>
          </content>
        </elementSpec>
      </schemaSpec>
    </div>
  </body>
</text>
```

L'outillage TEI

Une des raisons fortes pour lesquelles se servir de la TEI est l'existence des outils TEI p.e.

roma <http://www.tei-c.org/Roma> interface web permettant de construire des schemas TEI

oxGarage services web pour transformation entre TEI/docx/html/epub etc.

TEI publisher systeme d'édition numérique

Versioning Machine systeme de gestion des variance

txm outil textometrique capable de gerer directement TEI

Une autre est la possibilite de se servir de n'importe quelle outillage XML – parce que la TEI est 100% standard XML!

oxygen editeur (etc.) XML générique

saxon moteur de transformation XML générique

basex; exist systemes de base de données XML

Des références...

Site web du Consortium TEI <http://www.tei-c.org>

Depot github <https://github.com/TEIC>

Listes de discussion tei-l@listserv.brown.edu;
tei-fr@groupe.renater.fr (francophone)

Une excellente introduction francophone
<http://books.openedition.org/oep/1237>