

## TP2 : oXygen

TEI@Oxford

Sept 2012

### 1 À la découverte d'Oxygen - 2

Dans ce deuxième exercice, nous allons approfondir votre expérience avec Oxygen pour :

- l'encodage XML d'un document existant
- la représentation des structures aperçues dans un texte littéraire
- l'affichage en mode « auteur »

### 2 Sample text

Dans cet exercice, nous allons baliser ce sonnet de Du Bellay, un écrivain du XVI<sup>e</sup> siècle :

**Al'Ambicieux,  
ET AVARE ENNEMY  
DES BONNES LETTRES.**

Sonnet.

*Serf de Faueur, Esclave d'Avarice,  
Tu n'as jamais sur toy mesmes pouuoir,  
Et ie me veux d'un tel Maître pouruoir,  
Que l'Esprit libre en plaisir se nourrisse.  
L'Air, la Fortune, & l'humaine Folice  
Ont en leurs Mains ton malheureux Auoir.  
Le Iuge auare icy n'a rien à voir.  
Ny les troys Seurs, ny du Tens la malice.  
Regarde donc qui est plus souhaitable  
L'ayse, ou l'ennuy, le certain, ou l'instable.  
Quand à l'honneur, s'effere estre immortel:  
Car en cler Nom soubz Mort jamais ne tombe.  
Le tien obscur ne te promet rien tel.  
Ainsi, tous deux ferez soubz mesme Tumbes.*

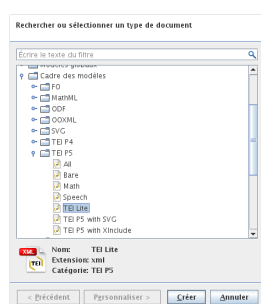
JO CARLO MVSA BEAT.

Pour vous simplifier la vie, on a déjà préparé une transcription prête à l'emploi... ne la regardez pas tout de suite !

### 3 Création du document

Ouvrez Oxygen. Cliquez sur l'icône Nouveau en haut à gauche, (ou sélectionnez Nouveau du menu Fichier, ou tapez CTRL-N.) pour ouvrir la boîte de dialogue Nouveau

- Cette fois, choisissez Cadre des modèles, ensuite TEI P5, et ensuite TEI-Lite, pour sélectionner un schéma TEI plus complet.



- Comme auparavant, Oxygen vous propose un squelette du document ; vous observerez que le schéma TEI lite décrit un document TEI plus élaboré avec un choix de balises plus important.
- Complétez l'entête TEI avec des indications de titre ('Un sonnet du Du Bellay numérisé'), distribution, et source ('encodé à partir d'un facsimilé de la première édition produit au CESR, Tours')
- Enlevez toute la partie proposée par défaut pour le `<body>` (mais retenez les balises `<body>` et `</body>`).

### 4 Ajout de texte au document

Rien ne vous empêche de taper à la main tout le poème directement. Mais pour gagner du temps on vous propose la démarche suivante :

- Assurez-vous que le curseur soit toujours entre les deux balises `<body>` et `</body>` de votre document vide
- Dans le menu Document, sélectionnez Fichier, et ensuite Insérer un fichier
- Naviguez jusqu'au fichier duBellay.txt dans le dossier Travaux qui devrait se trouver sur votre Bureau, et insérez-le. Ce fichier est aussi disponible sur le web, à cet URL )

### 5 Attention, Will Robinson !

Votre document est rempli de taches rouges!

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <div>
3   A l'antebellum,
4   ET AVANT D'AVANT
5   DES BONNES LETTRES
6   Sonnet.
7
8
9   Serf de Faveur, Esclave d'Avance,
10  Tu n'as jamais sur toi-même pouvoir,
11  Si je me veux d'un tel Maître pourvoir,
12  Que l'Esprit libre en plaisir se nourisse.
13
14
15  L'Air, la Fortune, & l'humaine Police
16  Ont en leurs flancs son malheureux Avoir
17  Le Jugement n'y n'a rien à voir,
18  N'y les trois Seurs, n'y du Tens la malice.
19
20
21  Regarde donc, qui est plus souhaitable
22  L'aveir, ou l'enfer, le certain, ou l'instable,
23  Quand à l'homme, j'espère estre immortel.
24
25
26  Car un cher Nom souz Mort jamais ne tombe,
27  Le tien obscur ne te promet rien tel,
28  Dans tous deux sans souz même Tumba,
29
30
31  C'ESTO MYSA BEAT.
32
33 </div>
34
```

[X] (Error) The entity name must immediately follow the '@' in the entity reference.

Texte Grille Auteur

On va régler cela petit à petit. Commençons avec la partie la plus sérieuse : où le texte lui-même est affiché en rouge.

- Oxygen vous indique que la partie en rouge n'est pas bien formée
- Il y a un message en bas de l'écran qui vous aidera (un peu) à identifier le problème. Sauriez vous comment le résoudre ?
- Conseil : souvenez-vous du fait que l'esperluette est un caractère magique en XML, qui doit donc être représentée de manière indirecte.

### 6 Structuration du document

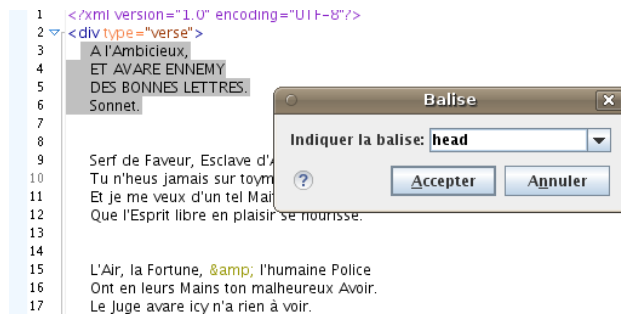
Il est possible d'identifier dans ce document plusieurs sous-parties. En particulier, il contient :

- un titre ... à baliser avec la balise `<head>`
- quatre strophes ... à baliser avec la balise `<lg>`
- des vers ... à baliser avec la balise `<l>`

- une citation d'Horace ... à baliser avec la balise `<quote>`

Allons-y !

- Avec la souris, sélectionnez les quatre premières lignes du document 'À l'ambitieux et avare ... sonnet'; ceux-ci constituent son titre.
- Tapez CTRL-E (ou sélectionnez XML Refactoring et ensuite Entourer des balises dans le menu Document
- Oxygen vous propose toutes les balises disponibles à cet emplacement.
- Sélectionnez `<head>` et cliquez Accepter



- Ce sonnet est bicéphale : coupez en deux parties le `<head>` existant.. La première partie est le titre. La seconde partie est le mot sonnet. Chacune de ces deux parties devrait- être balisée avec un `<head>` au début et à la fin, en vous plaçant avant ou après le mot 'Sonnet' et en vous servant de CTRL-ALT-D ou d'une autre manière comme lors du premier exercice.

## 7 Balisage du poème

Enfin, il faut travailler un peu sur l'étiquetage du poème lui-même. Notre but serait de baliser chaque vers, se servant de la balise TEI `<l>` ('line') et en amont de baliser chaque ensemble de vers constituant une strophe avec la balise TEI `<lrg>`.

Bien sur, il y a plusieurs manières de faire. On va vous décrire une méthode assez rapide, mais si vous préférez une autre démarche, on n'insiste pas !

- Sélectionnez avec la souris tout le texte du poème, jusqu'à la citation d'Horace 'caelo musa beat'.
- Entourez toute le poème d'une balise `<l>`.
- Entourez la citation d'Horace d'une balise `<quote>`
- Le petit carré vert apparaît ! Avons-nous terminé ? Hélas non : un document peut être valide, tout en contenant des mensonges ! Voyez-vous des mensonges ? Cliquez sur le bouton Indentation (ou tapez CTRL-SHIFT-P).

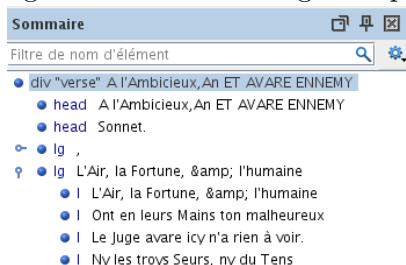
```
<body>
<head>A l'Ambitieux, ET AVARE ENNEMY DES BONNES LETTRES. Sonnet.</head>
<l>Serf de Faveur, Esclave d'Avarice, Tu n'heus jamais sur toymesmes pouvoir, Et je me veux
d'un tel Maître pourvoir, Que l'Esprit libre en plaisir se nourrisse. L'Air, la Fortune,
& l'humaine Police Ont en leurs Mains ton malheureux Avoir. Le Juge avare icy n'a rien à
voir. Ny les troys Seurs, ny du Tens la malice. Regarde donc qui est plus souhaitable
L'ayse, ou l'ennuy, le certain, ou l'instable. Quand à l'honneur, j'espere estre immortel:
Car un cler Nom souzb Mort jamais ne tombe. Le tien obscur ne te promet rien tel. Ainsi,
tous deux serez souzb mesme Tumbé. </l>
<quote> CAELO MUSA BEAT.</quote>
</body>
```

Parce qu'elles ne sont pas explicités par le balisage, la distinction entre chaque vers et entre chaque strophe ne sont plus affichées. Il faut donc les baliser. (Si vous ne voyez pas bien où indiquer ces séparations, vous pouvez ré-établir l'affichage originel en tapant CTRL-Z)

- Mettez le curseur au début de chaque vers, et distinguez-le du précédent, en vous servant du bouton **division** ou en tapant ALT-SHIFT-D
- Si cela vous semble plus rapide, vous pouvez aussi ‘tricher’ : copiez la séquence `</l><l>`, et collez-la au commencement de chaque vers.
- Quand vous aurez fini, cliquez encore sur le bouton **Indentation**. Ou sont les strophes ?
- Entourez tous les vers (tous les `<l>`) avec une balise `<lg>`.
- Divisez cette strophe unique en quatre strophes distinctes. Attention : il faut effectuer cette distinction au bon endroit — entre la fin d’un vers et le début du suivant, (i.e. après une balise fermante et avant une balise ouvrante).

## 8 Affichage du texte balisé

Est-ce que la disparition graduelle de notre beau sonnet au dessus d’un nuage de balises commence à vous inquiéter un peu ? Cela nous permet au moins de visualiser la structure émergente du sonnet – regardez par exemple dans la fenêtre ‘Sommaire’ à gauche :



Pour contrôler l’affichage d’un text balisé, il faut préciser comment nous désirons visualiser chaque balise. Ces précisions se font avec ce qu’on appelle une *feuille de style*.

- Sélectionnez **XML Document** et ensuite **Associer une feuille de style** à partir du menu **Document**.
- Naviguez jusqu’au fichier **duBellay.css** dans le dossier **Travaux** qui devrait se trouver sur votre Bureau, et associez-le (ce fichier est aussi disponible à l’URL sur le web ).
- Notez qu’une ligne contenant un ”processing instruction” est ajouté au début de votre document.
- Contrôlez l’effet sur l’affichage de votre document au mode **Auteur**.

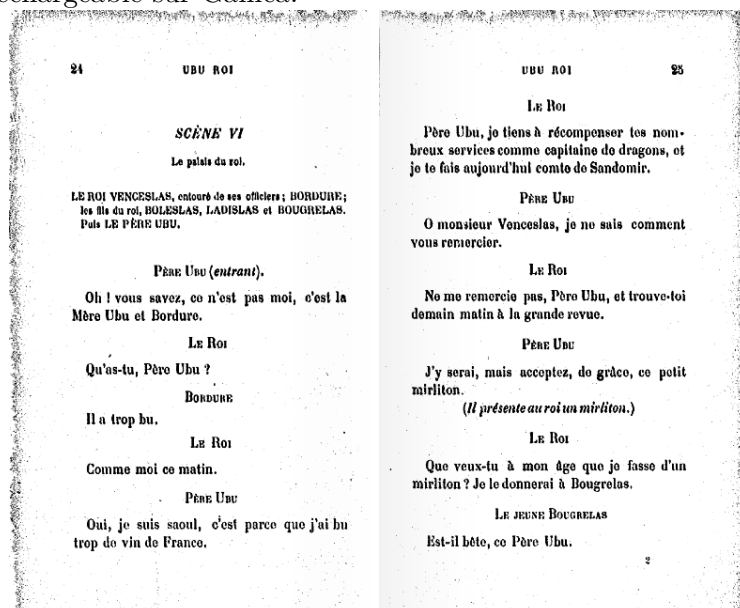
On peut changer cet affichage. Ouvrez le fichier **duBellay.css** avec Oxygen et modifiez-le, si cela vous intéresse. Mais il y a des limites :

- CSS n’est pas un langage XML, et ne permet que de modifier l’affichage du contenu d’un document XML.
- On peut changer la taille de police, l’étendu des marges, les couleurs etc.
- Changez la couleur du titre de **red** en **purple**.

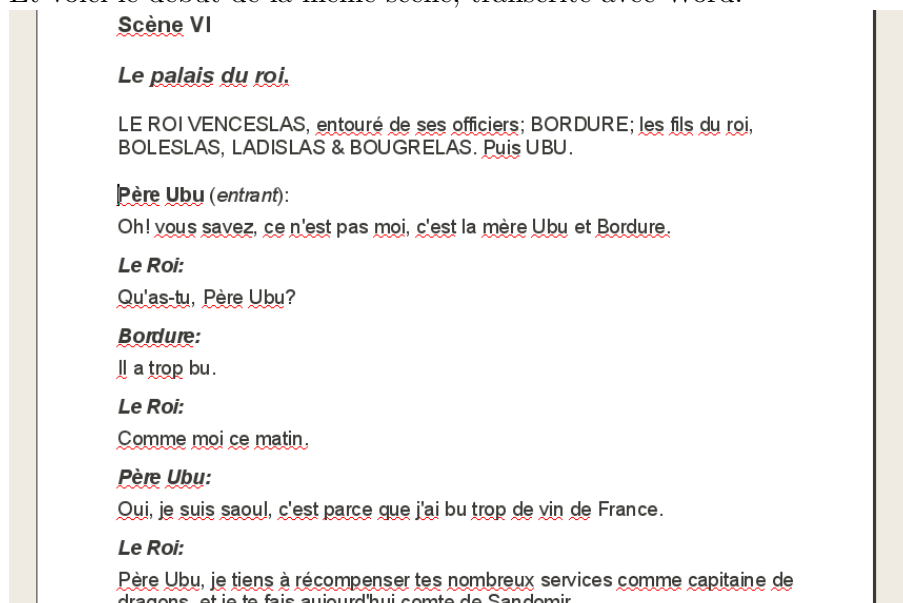
Des transformations beaucoup plus intéressantes sont possibles avec un autre langage de feuille de style : le XSLT.

## 9 Par ma chandelle verte...

Voici le commencement d'une scène extraite de *Ubu Roi* de Alfred Jarry (1896), dans la version téléchargeable sur Gallica.



Et voici le début de la même scène, transcrite avec Word:



## 10 Transformation d'un fichier Word

Vous avez vu comment baliser un fichier « text ». Comment pouvons-nous faire de même avec notre fichier Word ?

Nous pourrions exporter le fichier Word au format « plein text ». Nous pourrions aussi faire un copier-coller. Mais dans chaque cas, nous perdriions le formatage qui distingue (par exemple) le nom de chaque locuteur de ces énoncés.

Si le fichier Word est bien stylé, on peut le transformer en XML, sans perte d'information. L'outil OxGarage nous aidera.

- Ouvrez votre Internet browser, et allez sur le site <http://www.tei-c.org/ege-webclient/>
- Cliquez sur Text Documents. Dans la liste des formats de saisi proposée, cochez la case pour Microsoft Word Document.

- Une liste des formats cibles proposés apparaît. Cochez la case pour TEI P5 XML.
- En haut de l'écran un bouton **Browse** apparaît. Cliquez sur ceci, et naviguez jusqu'au fichier `ubu.doc` dans votre dossier **Travaux**.
- Cliquez sur le grand bouton **Convert** et patientez.
- Le site vous renvoie (après un bref délai) un fichier `ubu.xml`. Enregistrez-le dans votre dossier **Travaux**.

## 11 Structuration des pièces de théâtre

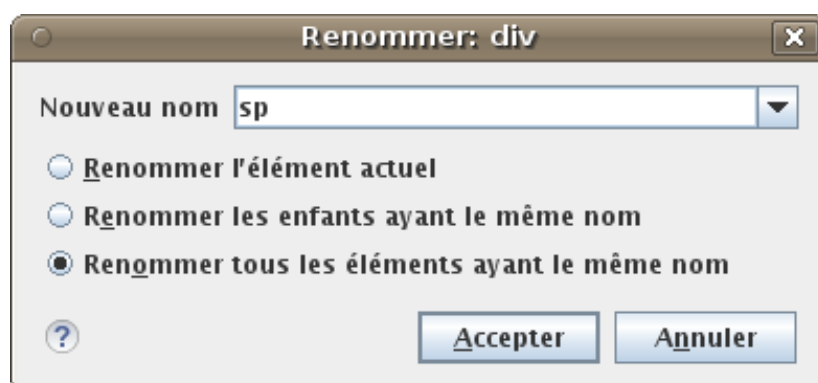
Une pièce de théâtre contient :

- des *didascalies* de plusieurs types... à baliser avec `<stage>` (en se servant de l'attribut `type` pour les distinguer au cas ou)
- des *énoncés* ou discours ... à baliser avec `<sp>` ('speech')
- des titres ou marques de *locuteur* ... à baliser avec `<speaker>`
- des paragraphes, ou des vers, balisés comme d'habitude.

## 12 Transformation des balises

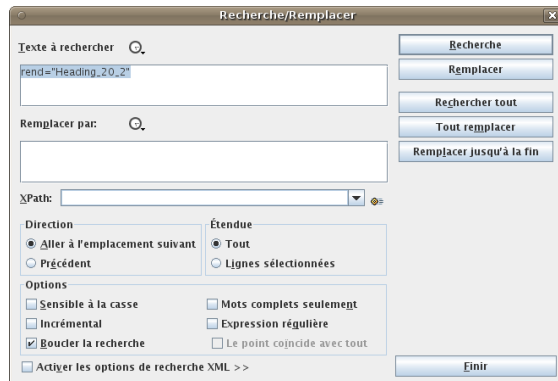
Voyons ce que le Garage aura fait de notre fichier Word. En tout cas, il est devenu un document TEI valide, même s'il est plein de mensonges...

- Lancez Oxygen, et ouvrez le fichier `ubu.xml` que vous venez de créer.
- On laisse passer pour le moment le TEI Header. Vous pouvez le compléter plus tard.
- D'abord notez que chaque énoncé est devenu un `<div>`, contenant un `<head>` et un `<p>`. Cela va beaucoup nous aider, et c'est une conséquence du fait que la version Word s'est servi des styles (chaque énoncé était précédé d'un objet style 'heading 2')
- Mettez le curseur d'abord sur une des balises `<div>`
- Dans le menu **Document** sélectionnez la commande **XML-Refactoring**, puis **Renommer l'élément** (ou tapez ALT-MAJ-R).
- Oxygen vous propose ce dialogue



- Dans le champs **Nouveau nom** entrez `sp`, le nom de la balise souhaitée, et cocher la case **Renommer tous les éléments ayant le même nom**.
- Cliquez sur **Accepter** et tous les `<div>` deviennent `<sp>`. (Le document n'est plus valide, mais on va rectifier cela tout à l'heure.)

- Faites de même pour les balises **<head>** : transformez les tous en **<speaker>**.
- Avant de rectifier la structuration du document, on va supprimer les **rend="Heading\_20\_2"** qui n'ont plus de sens. Sélectionnez cette chaîne de caractères n'importe où sur l'écran.
- Tapez CRL-F, ou sélectionnez Rechercher/Remplacer sur le menu Recherche.



- Tapez **Tout remplacer**.

Enfin, regardez les trois morceaux de texte mal-balisés au début de la scène : le plus simple serait d'enlever leur balisage actuel et ensuite d'ajouter les balises qu'il nous faut.

- Mettez le curseur sur le texte 'Scène Vi'. Dans le menu Document, sélectionnez XML Refactoring, et ensuite Effacer les balises (ALT-MAJ-X). Répétez jusqu'à ce que toutes les balises qui l'entourent soient effacées.
- Même jeu pour les deux lignes suivantes : 'Le palais...' et 'LE ROI VENCESLAUS...'
- Maintenant, entourez tout le texte souligné en rouge, et balisez-le avec un **<head>**.
- Votre document redevient valide ! à vous maintenant de le rendre correcte :
  - diviser les trois parties du **<head>** au début de la scène.
  - les phrases balisées **<emph>** et l'un des **<p>** sont des didascalies : changez-les donc en **<stage>**.
  - il reste des valeurs de l'attribut **rend** à enlever
  - Dans ce **<body>** nous n'avons qu'une des scènes du texte d'Ubu. Il serait préférable (plus honnête !) d'indiquer cela en entourant le tout avec un **<div type="scene">**

S'il vous reste du temps, pensez à compléter le TEI Header. Quand vous aurez fini, n'oubliez pas d'enregistrer votre belle version TEI du fichier!