



# Outil semi-automatique d'aide au balisage de texte pour le projet TMG

Florence Le Priol, Pôle des Systèmes d'Information,  
Université Paris-Sorbonne

# Plan

- Objectifs
- Choix techniques
- Fonctionnalités
- Démo
- Améliorations

# Objectifs

- Construire un outil pour l'aide au balisage des textes
  - Automatiser ce qui peut l'être (nettoyage, balisage des termes connus)
  - Outil d'aide pour le reste (nouveaux termes)
- Répondre aux normes de la TEI.
- multiplateforme (windows, mac os)
- pas nécessairement disponible en ligne
- Installation et utilisation simples y compris pour des non-informaticiens.

# Choix techniques

- Pourquoi choisir PYTHON
  - Python est bon dans tous les domaines
  - Il est excellent dans l'analyse des textes et les interfaces utilisateurs
  - Langage jamais bloquant quelque soit les besoins : nombreuses bibliothèques disponibles
  - Interfaçage natif avec de nombreux autres langages
  - Langage bien classé et en progression
  - Utilisé en production dans des développements commerciaux importants

# Choix techniques

- Et pourquoi pas JAVA ?
  - Excellent langage
  - Le plus utilisé
- Néanmoins
  - Python est passé en tête des langages d'apprentissage en Europe et aux Etats-Unis
  - La syntaxe de Python est plus simple et bien plus lisible que celle de Java
  - Python est plus simple à apprendre et à comprendre que Java

# Choix techniques

- Distribution
  - Stand-alone
  - Version Mac os
  - Version Windows
  - Utilisation de cx\_freeze

# Fonctionnalités

- Nettoyage automatique :
  - Suppression des balises issues de la numérisation
  - Utilisation de listes de balises ou d'expressions régulières
- Entrée :
  - Texte tel que récupéré de la numérisation, au format txt (dossier TMG\_sources)
  - Liste des balises (TMG\_data/nettoyage.txt)
  - Liste des expressions régulières (TMG\_data/nettoyageregex.txt)
- Sortie : texte nettoyé, au format txt (dossier TMG\_textes)

# Fonctionnalités

- Liste de balises
- Forme initiale fixe, Forme souhaitée

- Par exemple

#.s#.s est

remplacé

par ff

```
%--B, B
%--A, A
%--a, a
#.sz, f
#.si, fi
#.s#.s, ff
#.s, f
#.r, r
#. , , "
#;ou, ô
#;eæ, æ
{v_er_pu}, ð
```

- Liste d'expressions régulières
- Forme initiale variable, Forme souhaitée

```
<SE[0-9][0-9]?[0-9]?>,
<filename:[a-zA-Z0-9_.*]*/> ,
```

- Par exemple

<SE suivi de 1 à 3 chiffres  
> est supprimé (remplacé  
par rien)



# Fonctionnalités

Extrait de texte  
issu de la numérisation

```
<br type="page"/>
<p>@@@<filename:TMG_Anonymus_1496_0005.tif/></p>
<p>@@1@<SE5></p>
<p><P>E{v_ee}erpta mu#.sice omnis</P></p>
<p><P>cantus G#.rego#.riani Pita/</P></p>
<p><P>go#.rici {et1$} contrap%-ucti #.simplicis</P></p>
<p/>
<br type="page"/>
<p>@@@<filename:TMG_Anonymus_1496_0006.tif/></p>
<p>@@1@<SE6></p>
<p><tspb>ela</p>
<p>b</p>
<p>dd dd dd ddla#.sol</p>
<p>cc#.solfa</p>
<p>b bbfahmi</p>
<p>b aalamire</p>
<p>g g g g#.solreut</p>
<p>ffaut</p>
<p>b elami</p>
<p>dla#.solre</p>
<p>c c c c#.solfaut</p>
<p>b bfahmi</p>
<p>b alamire</p>
<p>G#.solreut</p>
<p>f f f Effaut</p>
<p>b Elami</p>
<p>D#.solre</p>
<p>Cfaut</p>
<p>b Bmi</p>
```

nettoyé

```
@@@
@@1@
Eyerpta mufice omnis
cantus Gregoriani Pita/
gorici 7 contrapūcti fimplicis

@@@
@@1@
<tspb>ela
b
dd dd dd ddlafol
ccfolfa
b bbfahmi
b aalamire
g g g gfolreut
ffaut
b elami
dlafolre
c c c cfolfaut
b bfahmi
b alamire
Gfolreut
f f f Effaut
b Elami
Dfolre
Cfaut
b Bmi
```

# Fonctionnalités

- Balisage automatique :
  - Balisage des termes dont la correspondance est déjà connue
  - Utilisation de la liste des ponctuations
  - Utilisation du dictionnaire
  
- Entrée :
  - texte nettoyé, au format txt (dossier TMG\_textes)
  - Liste des ponctuations (TMG\_data/balises\_ponctuations.txt)
  - Dictionnaire (TMG\_dictionnaires/dico.dico → URL serveur)
- Sortie : texte balisé, au format xml (dossier TMG\_textes)

# Fonctionnalités

- Liste des ponctuations
- Ponctuation, Balise

```
/, <pc type="S2"> / </pc>  
., <pc type="S1">. </pc>  
:, <pc type="S3">: </pc>  
,, <pc type="S4">, </pc>  
?, <pc type="S5">? </pc>
```

- Par exemple

/

est balisé

```
<pc type="S2"> / </pc>
```

- Dictionnaire
- Expression, Balise
- Par exemple

Componift

est balisé

```
<choice n =  
"gr"><orig>Componift</orig><reg>Componi  
st</reg></choice>
```

```
Capellmeister, <choice n = "gr"><orig>Capellmeister</orig><reg>Capellmeister</reg></choice>  
Componift, <choice n = "gr"><orig>Componift</orig><reg>Componist</reg></choice>  
Componiften, <choice n = "gr"><orig>Componiften</orig><reg>Componisten</reg></choice>  
Compositiōn, <choice n = "gr"><orig>Compositiōn</orig><reg>Composition</reg></choice>  
ConcertGefänge, <choice n = "gr"><orig>ConcertGefänge</orig><reg>ConcertGesänge</reg></choice>  
Consequentz, <choice n = "gr"><orig>Consequentz</orig><reg>Consequentz</reg></choice>  
Confonantia, <choice n = "gr"><orig>Confonantia</orig><reg>Consonantia</reg></choice>  
Confonantiam, <choice n = "gr"><orig>Confonantiam</orig><reg>Consonantiam</reg></choice>  
Confonantias, <choice n = "gr"><orig>Confonantias</orig><reg>Consonantias</reg></choice>  
Confonantien, <choice n = "gr"><orig>Confonantien</orig><reg>Consonantien</reg></choice>  
Confonatias, <choice n = "gr"><orig>Confonatias</orig><reg>Consonatias</reg></choice>
```

# Fonctionnalités

## Extrait de texte nettoyé

## balisé

Qm̄ magnos studiosis afferat fructus antiquissima illa nobilissimaq;  
musicalis scientia. quantaue sit autoritate probata. quātum deniq; lau=  
dis honorifue attigerit. haud facile dictu existimem. Naz oēs fere sancti  
patres: prophete: philosophi: ac romani pōtīfices. aut musici erant: aut  
musicis delectauere. Jd nimirū pfecto cum nec creditum esset quēpiam  
fatis eruditum fuisse: ni adipe perdulcis cātilene foret enutritus. Quod  
et ipi lacedemonij maxima ope seruauere. dū apud eos Thaletas Cretē=  
fis gortinus magno precio accitus: pueros disciplina musice artis im=  
bueret. Quoniā tanta fuit apud eos musice diligentia: vt eam animos  
q̄qz obtinere arbitrarentur. ac fane quidez rati. vulgatū em̄ est: q̄sepe ira  
cundias cantilena reprefferit: q̄multa vel in corporum vel in animor um  
affectioni[us] miranda p̄fecerit. Legit̄ em̄ eam nonnūq; sua dulcedine de  
monia fugasse: Infanos quoq; ad rationis v̄sum: incontinētes ad casti=  
tatem: tristes ad leticiā: egros ad fanitatem. inordinatas imaginationes  
ad constantiam deliberationēq; pigros ḡ inutiles ad agibilitatē reduxi=  
fe: nec nō etiam melanc oliam tyrannidemq; a magnati[us] depulisse: Jdq;  
non mō sese in singulis vel studijs vel etati[us] tenet: verū p̄ cuncta diffun  
ditur studia. vt infantes. inuenes. fenes. ḡ viri ḡ mulieres. ita naturaliter  
affectu quodā spontaneo modis musicis adiunguntur: vt nulla omnino  
sit etas que a cantilene dulcis delectatione seiuncta sit. Nonne etiā illud  
manifestum est in bellum pugnantium animos tubarum carmine incēdi ḡ  
fopitos excitari: vigilantū quoq; animis cum dulcior esset melodia fopo  
rem peruaderi. Fuit em̄ pitagoricis id in morem diuq; permanfit. vt cū  
diuturnas in sōno resoluere curas quibusdam cantilenis vtrentur. vt  
eis lenis ḡ quietus esset sopor. itaq; expereti alijs quibusdā modis fu=  
porem somni confusionemq; purgabant Jd nimirū scientes siquidez eā

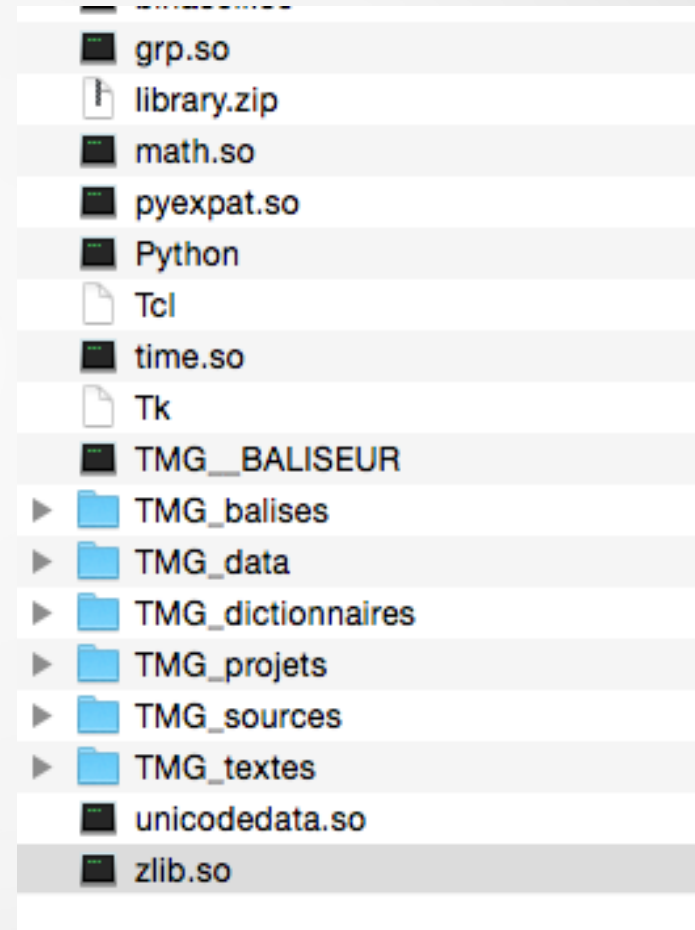
105 Qm̄ magnos studiosis afferat fructus antiquissima illa nobilissimaq;  
106 musicalis scientia<pc type="S1">.</pc> quantaue sit autoritate probata<pc type="S1">.</pc> quātum deniq; lau=  
107 dis honorifue attigerit<pc type="S1">.</pc> haud facile dictu existimem<pc type="S1">.</pc> Naz oēs fere sancti  
108 patres<pc type="S3">.</pc> prophete<pc type="S3">.</pc> philosophi<pc type="S3">.</pc> ac romani pōtīfices<pc type="S1">.</pc> aut  
musicis erant<pc type="S3">.</pc> aut  
109 musicis delectauere<pc type="S1">.</pc> Jd nimirū pfecto cum nec creditum esset quēpiam  
110 fati eruditum fuisse<pc type="S3">.</pc> ni adipe perdulcis cātilene foret enutritus<pc type="S1">.</pc> Quod  
111 et ipi lacedemonij maxima ope seruauere<pc type="S1">.</pc> dū apud eos Thaletas Cretē=  
112 fis gortinus magno precio accitus<pc type="S3">.</pc> pueros disciplina musice artis im=  
113 bueret<pc type="S1">.</pc> Quoniā tanta fuit apud eos musice diligentia<pc type="S3">.</pc> vt eam animos  
114 q̄qz obtinere arbitrarentur<pc type="S1">.</pc> ac fane quidez rati<pc type="S1">.</pc> vulgatū em̄ <choice n = "gr"><orig>eft</  
orig><reg>est</reg></choice><pc type="S3">.</pc> q̄sepe ira  
115 cundias cantilena reprefferit<pc type="S3">.</pc> q̄multa vel in corporum vel in animor um  
116 affectioni[us] miranda p̄fecerit<pc type="S1">.</pc> Legit̄ em̄ eam nonnūq; sua dulcedine de  
117 monia fugasse<pc type="S3">.</pc> Infanos quoq; ad rationis v̄sum<pc type="S3">.</pc> incontinētes ad casti=  
118 tatem<pc type="S3">.</pc> tristes ad leticiā<pc type="S3">.</pc> egros ad fanitatem<pc type="S1">.</pc> inordinatas imaginationes  
119 ad constantiam deliberationēq; pigros ḡ inutiles ad agibilitatē reduxi=  
120 fe<pc type="S3">.</pc> nec nō etiam melanc oliam tyrannidemq; a magnati[us] depulisse<pc type="S3">.</pc> Jdq;  
121 non mō sese in singulis vel studijs vel etati[us] tenet<pc type="S3">.</pc> verū p̄ cuncta diffun  
122 ditur studia<pc type="S1">.</pc> vt infantes<pc type="S1">.</pc> inuenes<pc type="S1">.</pc> fenes<pc type="S1">.</pc> ḡ viri ḡ mulieres<  
pc type="S1">.</pc> ita naturaliter  
123 affectu quodā spontaneo modis musicis adiunguntur<pc type="S3">.</pc> vt nulla omnino  
124 sit etas que a cantilene dulcis delectatione seiuncta sit<pc type="S1">.</pc> Nonne etiā illud  
125 manifestum <choice n = "gr"><orig>eft</orig><reg>est</reg></choice> in bellum pugnantium animos tubarum carmine incēdi ḡ  
126 fopitos excitari<pc type="S3">.</pc> vigilantū quoq; animis cum dulcior esset melodia fopo  
127 rem peruaderi<pc type="S1">.</pc> Fuit em̄ pitagoricis id in morem diuq; permanfit<pc type="S1">.</pc> vt cū

# Fonctionnalités

- Balisage semi-automatique
  - Balisage des termes nouveaux (non présent dans le dictionnaire)
  - A l'aide de l'interface de balisage
- Entrée :
  - Texte balisé, au format xml (dossier TMG\_textes)
  - Dictionnaire (TMG\_dictionnaires/dico.dico → URL serveur)
  - Dictionnaire temporaire (TMG\_dictionnaires/dico.dico.temp)
  - Schéma de balises (TMG\_balises/balises.bal)
  - Liste des termes comportant des caractères spéciaux (à partir de TMG\_data/caracteresspeciaux.txt)
- Sortie : Texte balisé, au format xml (dossier TMG\_textes)

# Démo

- Lancement : `TMG__BALISEUR`



# Améliorations

- Traitement de tous les schémas de balise
- Possibilité d'individualiser le traitement d'un terme
- Navigation à travers les termes balisés / à baliser
- Rédiger les menus d'aide
- Visuel / ergonomie
- Version windows
  
- ... et en fonction des retours des utilisateurs...