

# Manuel d'encodage XML-TEI étendu des transcriptions de manuscrits dans le projet BFM-Manuscrits

Dernier enregistrement le 26 juin 2008

Alexei Lavrentiev (Alexei.Lavrentev@ens-lsh.fr)

CNRS / ENS-LSH, UMR 5191 ICAR



(VERSION 2.1- JUIN 2008)  
(VERSION 2.0- MARS 2008)  
(VERSION 1.0- DECEMBRE 2005)

Ce document de travail est élaboré dans le cadre des opérations de saisie et de traitement automatique des transcriptions de manuscrits de la Base de Français Médiéval (<http://bfm.ens-lsh.fr>). Le protocole de balisage est une extension du schéma de la TEI (P5) (<http://www.tei-c.org>) et se base sur les principes exposés dans le *Manuel d'encodage XML-TEI des textes de la Base de Français Médiéval* ([http://bfm.ens-lsh.fr/IMG/pdf/Manuel\\_Encodage\\_TEI.pdf](http://bfm.ens-lsh.fr/IMG/pdf/Manuel_Encodage_TEI.pdf)) et dans le *Manuel de description de textes pour la Base de Français Médiéval*, sur l'expérience du projet *Charrette* (<http://www.princeton.edu/~lancelot/>), sur les recommandations présentées dans le *Handbook* (v. 2) du projet Menota (<http://www.menota.org/>) et sur les propositions de l'initiative MIFI (<http://gandalf.aksis.uib.no/mifi/>).



## Sommaire

Introduction : Principes du codage XML, application et extension de la TEI.....	3
1. Descripteur (entête TEI) .....	4
<teiHeader> .....	7
<fileDesc>.....	7
<titleStmt> .....	8
<title>.....	8
<author>.....	8
<editor>.....	8
<respStmt>.....	8
<resp> .....	8
<name> .....	8
<extent> .....	8
<publicationStmt> .....	9
<sourceDesc>.....	9

<msDesc> .....	9
<msIdentifier> .....	9
<msPart> .....	9
<physDesc> .....	9
<handDesc> .....	9
<handNote> .....	9
<history> .....	9
<origin> .....	10
<profileDesc> .....	10
<creation> .....	10
<date> .....	10
<langUsage> .....	10
<textDesc> .....	10
<textClass> .....	11
<catRef> .....	11
<encodingDesc> .....	11
<classDecl> .....	11
<taxonomy> .....	11
<category> .....	11
<catDesc> .....	11
<revisionDesc> .....	11
<change> .....	11
2. Niveaux de représentation des données (diffraction) .....	12
Représentation normalisée : <me:norm> .....	13
Représentation diplomatique : <me:dipl> .....	13
Représentation imitative : <me:fac> .....	14
Représentation paléographique (facultative) : <me:pal> .....	14
3. Corrections sribales et éditoriales ; passages difficilement lisibles ou illisibles .....	15
Partie d'un mot .....	15
Un mot entier ou plusieurs mots .....	15
Segment commençant ou finissant au milieu d'un mot et s'étendant sur plusieurs mots .....	16
Notes et commentaires du relecteur/encodeur .....	16
4. Définition et description des éléments .....	16
<ab> .....	16
<add> .....	17
<am> .....	17
<bfm:headlb> .....	17
<bfm:hyphen> .....	19
<bfm:lettrine> .....	19
<bfm:mdvAbbr> .....	20
<bfm:punct> .....	21
<bfm:sb> .....	21
<cb> .....	21
<choice> .....	22
<corr> .....	22
<damage> .....	22
<del> .....	22
<div> .....	23
<ex> .....	23
<gap> .....	24

<head> .....	24
<hi> .....	24
<lb> .....	24
<me:dipl> .....	25
<me:fac> .....	25
<me:norm> .....	25
<me:pal> .....	25
<note> .....	25
<p> .....	26
<pb> .....	26
<q> .....	26
<seg> .....	27
<space> .....	27
<subst> .....	27
<supplied> .....	28
<unclear> .....	28
<w> .....	28
5. Saisie et correction des transcriptions (syntaxe compacte).....	29
Caractères de raccourci, balises simplifiées.....	29
6. Relations dans la base des descripteurs BFM .....	34
7. Tableaux d'encodage des caractères « spéciaux » .....	35
Lettres initiales (letrines) .....	35
Variantes de caractères .....	35
Ligatures .....	35
Caractères exponctués.....	35
Abréviations .....	36
Marques de ponctuation et de mise en page .....	37
8. Projets cités : .....	37
9. Index des éléments (entête).....	37
10. Index des éléments (hors entête).....	38

## Introduction :

# Principes du codage XML, application et extension de la TEI

Comme c'est le cas de l'ensemble des textes constituant la Base de Français Médiéval (BFM), le corpus de transcriptions de manuscrits, d'incunables et de livres imprimés du XVI<sup>e</sup> siècle (BFM-MSS) est encodé selon la norme XML (<http://www.w3c.org/xml>), en conformité avec les recommandations de la Text Encoding Initiative (TEI, version P5, <http://www.tei-c.org/Guidelines/P5>).

Les principes de base du langage XML et de la TEI sont exposés dans la documentation de référence publiée sur les sites indiqués ci-dessus.

Les recommandations de la TEI étant très générales, une adaptation et une sous-spécification sont nécessaires dans le cadre d'un projet particulier. Le protocole d'encodage du corps des textes d'éditions modernes numérisées pour la BFM est exposé dans le *Manuel d'encodage XML-TEI des textes de la Base de Français Médiéval* ([http://bfm.ens-lsh.fr/IMG/pdf/Manuel\\_Encodage\\_TEI.pdf](http://bfm.ens-lsh.fr/IMG/pdf/Manuel_Encodage_TEI.pdf)). Le protocole de description et de caractérisation

des textes de la BFM est exposé dans le *Manuel de description de textes pour la Base de Français Médiéval*. Ces deux documents servent de base pour des extensions éventuelles dans le cadre d'un projet particulier.

La transcription fine de manuscrits et d'incunables a nécessité l'usage de balises spéciales, définies pour la plupart dans les modules **transcr** et **msdescription** de la TEI. Dans quelques cas précis, l'introduction de balises supplémentaires à celles proposées par la TEI a été jugée utile. Quelques unes de ces balises sont empruntées au schéma du projet Menota (Medieval Nordic Text Archive, <http://www.menota.org>), d'autres sont définies par le projet BFM-MSS. L'introduction de ces balises étrangères à la TEI s'est effectuée en conformité avec les mécanismes d'extension prévue par la TEI même.

L'espace de nommage des balises XML par défaut est celui de la TEI (<http://www.tei-c.org/ns/1.0>). Toutes les balises non-TEI appartiennent soit à l'espace Menota (<http://www.menota.org/ns/1.0>, préfixe **me:**), soit BFM (<http://bfm.ens-lsh.fr/ns/1.0>, préfixe **bfm:**).

## 1. Descripteur (entête TEI)

Le corpus est composé de transcriptions de fragments de manuscrits français médiévaux, d'incunables et de livres du XVI<sup>e</sup> siècle. Chaque transcription forme une unité du corpus. Chaque unité du corpus porte un entête TEI.

Comme c'est le cas de l'ensemble des textes de la BFM, les descripteurs des unités textuelles sont stockés dans une base de données relationnelle, les entêtes TEI sont générés à partir des données de cette base au moment de la publication du corpus ou d'un document particulier. Par la publication nous entendons l'intégration à l'outil d'exploitation, toute mise en ligne, cession à un partenaire, etc.

Le descripteur d'une unité du corpus BFM-MSS inclut tous les champs caractérisant une œuvre dans la Base de Français Médiéval (titre, auteur, date de composition, forme, domaine, genre, etc.). La liste complète de ces champs est donnée dans le *Manuel de description de textes pour la BFM*.

La description du manuscrit et de sa transcription sont basées sur les dernières recommandations de la TEI, qui sont à leur tour, le fruit d'un long travail d'un groupe d'experts.

Un manuscrit médiéval (ou codex) peut en effet être composé de plusieurs parties hétérogènes réunies sous une reliure longtemps après l'écriture du texte. Chacune de ses parties peut à son tour inclure plusieurs unités textuelles (« œuvres » dans la terminologie de la BFM). Une transcription représente enfin une ou plusieurs parties d'une unité textuelle dans un manuscrit. Nous avons donc affaire à une structure hiérarchique, qui, dans une représentation XML prend forme d'arborescence d'éléments : **<msDesc>** (description du volume physique) peut contenir plusieurs **<msPart>** (parties physiquement homogènes) qui, à leur tour, contiennent un certain nombre de **<msItem>** (unités textuelles dans un manuscrit). Bien entendu, l'entête TEI d'une transcription donnée ne contiendra qu'un seul **<msItem>**.

Dans notre base de données relationnelle, cette hiérarchie est traduite par un système de tables entretenant une relation « un à plusieurs » et portant, pour la commodité, les mêmes noms que les balises TEI correspondantes (cf. le schéma des relations ci-joint). Enfin, une table spéciale ('MsTranscription') est dédiée aux données sur la transcription.

Les incunables et les livres du XVI<sup>e</sup> siècle occupent une place intermédiaire entre les manuscrits et les éditions modernes : la description de l'exemplaire (objet physique localisé dans une bibliothèque) est pour eux aussi importante que les références de générale de

l'édition. Dans notre base de données ces imprimés anciens sont décrits dans des tables spécialisées ('Incunable', 'IncItem' et 'IncTranscription', cf. le schéma des relations). Dans l'entête TEI, les incunables sont décrits de la même façon que les manuscrits, avec quelques différences légères (par exemple, il n'y a pas de description des mains).

Un exemple d'un entête complet d'une transcription d'un manuscrit est présenté ci-dessous. Si la structure générale de l'entête est imposée par la TEI, les projets individuels disposent d'une certaine liberté dans le choix d'informations renseignées et dans leur disposition. Nous commenterons ensuite les éléments de l'entête dans l'ordre de leur apparition en indiquant nos choix aux cas où le schéma de la TEI est flexible.

```

<teiHeader type="text">
  <fileDesc>
    <titleStmt>
      <title type="normal">Chronique</title>
      <title type="source">Manuscrit Paris, BnF fr. 2682</title>
      <title type="reference" n="1012">monstre2682</title>
      <title type="medium">transcription électronique</title>
      <author>Enguerrand de Monstrelet</author>
      <editor role="editor">Équipe diachronie et bases textuelles d'ancien
et moyen français - UMR5191 ICAR, CNRS/ENS-LSH</editor>
      <respStmt>
        <resp>Transcription du manuscrit, encodage XML</resp>
        <name type="person" xml:id="AL">Alexei Lavrentiev</name>
      </respStmt>
      <respStmt>
        <resp>Datation du manuscrit</resp>
        <name type="document" xml:id="catBNF">Catalogue de la BnF</name>
      </respStmt>
    </titleStmt>
    <extent>- Taille approximative du fichier encodé en TEI non compressé :
104827 octets
- Nombre d'occurrences-mots : 847
- Numéros de pages saisies : de 233r à 233v
</extent>
    <publicationStmt>
      <distributor><name>Projet BFM, UMR 5191 ICAR, CNRS/ENS LSH</name>
      <address>
        <addrLine>15 parvis René-Descartes</addrLine>
        <addrLine>B.P. 7000 </addrLine>
        <addrLine>69342 Lyon Cedex 07</addrLine>
        <addrLine>France</addrLine>
        <addrLine>Tél : 04 37 37 63 10</addrLine>
        <addrLine>Fax : 04 37 37 62 65</addrLine>
        <addrLine>http://bfm.ens-lsh.fr/</addrLine>
      </address>
      <email>bfm@ens-lsh.fr</email>
    </distributor>
    <availability status="restricted">
      <p>(c) 2007, Equipe BFM, CNRS/ENS-LSH.
      <hi>Conditions d'utilisation</hi> : usage interne équipe BFM
Pour d'autres conditions d'usage, contacter :
      <address>
        <addrLine><name type="person">Céline Guillot</name></addrLine>
        <addrLine>Mèl : Celine.Guillot@ens-lsh.fr</addrLine>
      </address>
    </p>
    </availability>
  </publicationStmt>

```

```

<sourceDesc>
  <msDesc>
    <msIdentifier>
      <country>France</country>
      <settlement>Paris</settlement>
      <repository>BnF</repository>
      <idno>fr. 2682</idno>
    </msIdentifier>
    <msPart>
      <altIdentifier>
        <idno>bnf_fr2682</idno>
      </altIdentifier>
      <physDesc>
        <handDesc hands="2">
          <handNote xml:id="scribe1">Premier scribe :
            <name type="person">inconnu</name>
          </handNote>
          <handNote xml:id="rubricateur1">Premier rubricateur :
            <name type="person">inconnu</name>
          </handNote>
        </handDesc>
      </physDesc>
      <history>
        <origin>
          Date du manuscrit : <origDate value="1450-01-01"
                                notBefore="1400-01-01"
                                notAfter="1499-01-01"
                                resp="#catBNF"
                                evidence="undefined"
                                cert="low">15e s.
                                </origDate>
          Région : <region>inconnue</region>
          Ville : <settlement>inconnue</settlement>
          Atelier : <origPlace>inconnu</origPlace>
        </origin>
      </history>
    </msPart>
  </msDesc>
</sourceDesc>
</fileDesc>
<profileDesc>
  <creation>
    <date when="1441-01-01">1441</date>
  </creation>
  <langUsage>
    <language ident="fr" usage="100">Le texte est entièrement écrit en
vers en français médiéval</language>
  </langUsage>
  <textDesc n="chronique">
    <channel mode="w"></channel>
    <constitution type="single"></constitution>
    <derivation type="original"></derivation>
    <domain type="historique"></domain>
    <factuality type="inapplicable"></factuality>
    <interaction type="none"></interaction>
    <preparedness type="prepared"></preparedness>
    <purpose type="inform"></purpose>
  </textDesc>
  <textClass>
    <catRef target="#forme_prose"/>
  </textClass>

```

```

</profileDesc>
<encodingDesc>
  <classDecl>
    <taxonomy xml:id="forme">
      <category xml:id="forme_vers">
        <catDesc>vers</catDesc>
      </category>
      <category xml:id="forme_prose">
        <catDesc>prose</catDesc>
      </category>
      <category xml:id="forme_mixte">
        <catDesc>mixte</catDesc>
      </category>
    </taxonomy>
  </classDecl>
  <projectDesc>
    <p>Projet : BFM-MSS - Base de Français Médiéval - Manuscrits
    Resp. : <name type="person">Alexei Lavrentiev</name>
    Équipe : Diachronie et bases textuelles d'ancien et de moyen
français
    Laboratoire : UMR5191 ICAR
    Institution : CNRS / ENS LSH, Lyon
    <address>
      <addrLine>15 parvis René-Descartes</addrLine>
      <addrLine>B.P. 7000 </addrLine>
      <addrLine>69342 Lyon Cedex 07</addrLine>
      <addrLine>France</addrLine>
      <addrLine>Tél : 04 37 37 63 10</addrLine>
      <addrLine>Fax : 04 37 37 62 65</addrLine>
      <addrLine>http://bfm.ens-lsh.fr/</addrLine>
    </address>
    <email>bfm@ens-lsh.fr</email>
  </p>
  </projectDesc>
  <editorialDecl>
    <p>See BFM-MSS Encoding Guidelines at http://bfm.ens-lsh.fr/IMG/pdf/BFM-Mss\_Encodage-XML.pdf.
    See also BFM XML-TEI Encoding Guidelines at http://bfm.ens-lsh.fr/IMG/pdf/Manuel\_Encodage\_TEI.pdf
    and BFM XML-TEI Encoding Instructions at http://bfm.ens-lsh.fr/IMG/pdf/Consignes\_BFM.pdf</p>
  </editorialDecl>
</encodingDesc>
<revisionDesc>
  <change who="#AL" when="2007-07-22">Création de l'entête TEI et de la
  représentation diffractée</change>
</revisionDesc>
</teiHeader>

```

## <teiHeader>

L'entête TEI peut contenir 4 éléments : <fileDesc>, <encodingDesc>, <profileDesc> et <revisionDesc>. Tous ces éléments sont présents dans l'entête BFM-MSS.

## <fileDesc>

Description du document électronique. Dans le projet BFM-MSS cette description est composée de <titleStmt>, <extent>, <publicationStmt> et <sourceDesc>.

### **<titleStmt>**

Déclaration du titre. Contient une ou plusieurs occurrences des éléments suivants : **<title>**, **<author>**, **<editor>**.

### **<title>**

Balise générique pour un titre. Son attribut **type** est utilisé pour les catégoriser.

- **type="normal"** – titre conventionnel de l'œuvre, tel que donné dans les ouvrages de référence ;
- **type="source"** – brève indication de la source ;
- **type="reference"** – identifiant BFM (l'attribut **n** contient l'identifiant numérique dans le corpus de transcriptions) ;
- **type="medium"** – sert à préciser qu'il s'agit d'une transcription électronique.

### **<author>**

Le nom de l'auteur.

### **<editor>**

Balise générique pour un auteur « secondaire ». Son attribut **role** permet de préciser le rôle.

- **role="editor"** – institution responsable de l'édition du document numérique ;
- **role="ed\_sci"** – éditeur scientifique (responsable des corrections dans la représentation normalisée) ;
- **role="imprimeur"** – scribe (utilisé dans **<msDesc>** pour les incunables).

Dans les deux derniers cas, cet élément porte également un attribut **xml:id** permettant de se référer à la personne correspondante.

### **<respStmt>**

Déclaration des responsabilités, contient un ou plusieurs éléments **<resp>** et **<name>**.

### **<resp>**

Description d'un rôle d'un responsable de l'édition.

### **<name>**

Balise générique pour un nom propre. Son attribut **type** permet de les catégoriser. Son attribut **xml:id** permet de créer un identifiant auquel d'autres éléments pourront se référer (par exemple, d'indiquer le responsable d'une correction dans le texte).

### **<extent>**

Informations sur les dimensions du document numérique.

### **<publicationStmt>**

Déclaration de la publication : informations sur la distribution et l'accessibilité du document numérique. Contient les éléments **<distributor>** et **<availability>** que nous ne détaillerons pas davantage.

### **<sourceDesc>**

Conteneur de la description de la source du document numérique. Dans l'entête d'un texte du corpus BFM-MSS, il contient toujours un élément **<msDesc>**.

### **<msDesc>**

Conteneur de la description du codex ou de l'incunable.

Il contient un **<msIdentifier>** et

- **<msPart>**, si le codex est hétérogène ;
- **<physDesc>** et **<history>**, si le manuscrit est homogène.

### **<msIdentifier>**

Identifiant du manuscrit. Contient les éléments permettant de localiser le manuscrit ou l'exemplaire de l'incunable transcrit : **<country>** (pays), **<settlement>** (ville), **<repository>** (bibliothèque ou autre lieu de dépôt), **<idno>** (cote).

### **<msPart>**

Contient un identifiant **<altIdentifier>** et les mêmes éléments que **<msDesc>** pour les manuscrits homogènes.

### **<physDesc>**

Description physique du manuscrit ou de l'incunable. Contient l'élément **<objectDesc>** (description du volume en « prose libre ») et, pour les manuscrits, **<handDesc>**.

### **<handDesc>**

Description des mains du manuscrit. Contient plusieurs éléments **<handNote>**. Son attribut **hands** sert à indiquer le nombre de mains. On distingue habituellement au moins deux mains : celles du scribe et du rubricateur (qui réalise les lettrines et les titres en couleur).

### **<handNote>**

Une note caractérisant une main particulière du manuscrit. Son attribut **xml:id** permet de créer l'identifiant auquel on peut se référer pour indiquer la responsabilité de corrections sribales. La valeur de cet attribut indique le rôle et un indice numérique. Par exemple, **xml:id="scribe1"** pour le premier (ou unique) scribe identifié. Cet élément peut contenir de la « prose libre » avec un élément **<name>** qui indique éventuellement le nom de la personne.

### **<history>**

Dans le corpus BFM-MSS, il contient un seul élément **<origin>**.

## **<origin>**

Contient les informations sur l'origine du manuscrit ou de l'incunable :

- date du manuscrit (élément **<origDate>**, cf. élément **<date>** ci-dessous) ;
- région (élément **<region>**) ;
- ville (élément **<settlement>**) ;
- atelier pour les manuscrits (élément **<origPlace>**) ;
- nom de l'imprimeur pour les incunables (élément **<editor role="imprimeur">**).

## **<profileDesc>**

Description du texte représenté par la transcription, indépendante de sa forme matérielle. Le même schéma est utilisé pour l'ensemble des textes de la BFM. Contient les éléments **<creation>**, **<langUsage>**, **<textDesc>** et **<textClass>**. Le modèle de cette description est commun pour l'ensemble des textes de la BFM.

## **<creation>**

Informations sur la composition du texte, contient un élément **<date>** indiquant la date de composition.

## **<date>**

Élément générique permettant d'encoder une date. D'autres éléments, plus spécialisés peuvent être utilisés selon les cas (par exemple, **<origDate>**), mais le mécanisme est toujours le même.

Le contenu de cet élément est l'indication de la date sous forme libre. Ses attributs permettent de formaliser la date pour la rendre exploitable par des outils informatiques :

- **when** : date au format ISO (aaaa-mm-jj) ;
- **notBefore** : la date la plus ancienne possible au format ISO ;
- **notAfter** : la date la plus tardive possible au format ISO ;
- **cert** : certitude de la datation ;
- **resp** : responsable de la datation.

L'algorithme du calcul de la date formalisée du projet BFM est décrit dans le *Manuel de description de textes pour la BFM*.

## **<langUsage>**

Cet élément défini dans le module indique l'usage des différentes langues dans le texte. Contient un ou plusieurs éléments **<language>**.

## **<textDesc>**

Caractérisation typologique du texte, contient les éléments **<channel>**, **<constitution>**, **<derivation>**, **<domain>**, **<factuality>**, **<interaction>**, **<preparedness>** et **<purpose>** dont l'usage est précédemment défini par la TEI. Dans la BFM, le domaine du texte occupe une place centrale parmi les variables de la typologie textuelle. Une liste d'autorité des domaines a été établie et intégrée à la personnalisation du schéma de la TEI pour la BFM. Cette liste d'autorité figure dans le *Manuel de description de textes pour la BFM*.

D'autres variables importantes pour la BFM (telles que la forme et le genre du texte), elles ont donc été introduites à l'aide d'un mécanisme plus générique utilisant les éléments **<catRef>** et **<catDesc>** présentés ci-dessous.

### **<textClass>**

Conteneur de la référence à la caractérisation du texte selon une taxonomie quelconque. Contient un ou plusieurs éléments **<catRef>**. Dans le projet BFM ce mécanisme est utilisé pour caractériser la forme (vers, prose ou mixte) et le genre du texte (roman, chronique, etc.). La liste d'autorité des genres et des formes figure dans le *Manuel de description de textes pour la BFM*.

### **<catRef>**

Référence à une catégorie. Son attribut **target** assure la connexion avec la description de la catégorie **<catDesc>** dans le cadre d'une taxonomie.

### **<encodingDesc>**

Description des normes d'encodage. Contient les déclarations des classes **<classDecl>**, ainsi que les éléments **<projectDesc>** et **<editorialDecl>** que nous ne détaillerons pas davantage.

### **<classDecl>**

Déclaration d'une classe. Utilisée dans la BFM pour présenter les taxonomies de catégories de textes. Contient plusieurs éléments **<taxonomy>**.

### **<taxonomy>**

Contient plusieurs éléments **<category>**.

### **<category>**

Contient la description d'une catégorie **<catDesc>**. Son attribut **xml:id** permet de se référer à cette description à partir de **<catRef>**.

### **<catDesc>**

Contient une description de la catégorie en « prose libre ».

### **<revisionDesc>**

Historique des modifications du document numérique. Contient plusieurs éléments **<change>**.

### **<change>**

Contient la description d'une opération de modification. Ses attributs **who** et **when** permettent d'indiquer la personne intervenante (qui doit figurer dans la déclaration des responsabilités) et la date de la modification.

## 2. Niveaux de représentation des données (diffraction)

Un document XML est organisé comme une arborescence d'éléments imbriqués. Ainsi, un roman en prose peut être représenté comme un élément `<text>` composé d'une partie liminaire (titre, épigraphe, prologue, etc.) et du « corps » du texte, composé à son tour de plusieurs chapitres. Les transcriptions des sources du projet BFM-Manuscrits visent à représenter de façon fine et fidèle les données des originaux, en y ajoutant des couches distinctes de normalisation et de « filtrage » permettant d'en faciliter la lecture et l'analyse. Pour atteindre cet objectif, il a été nécessaire de créer des représentations multiples des données et de choisir le niveau de la structure textuelle auquel ces représentations divergent.

Puisque la BFM est destinée avant tout à servir à des recherches linguistiques, l'unité principale de représentation qui s'impose est le mot. Le mot est en effet « porteur » des propriétés morpho-syntaxiques et l'unité d'indexation pour une grande partie des outils de l'informatique linguistique. Au même niveau que les mots se situent les marques de ponctuation et certaines balises de mise en page et de référence (sauts de ligne et de page, certains commentaires).

La structure de la transcription se complique dans les cas où le manuscrit présente des corrections sribales (suppressions ou ajouts) ou des passages illisibles (détérioration matérielle ou écriture ambiguë).

Nous utilisons le terme de *représentation diffractée*<sup>1</sup> pour désigner la séparation systématique des niveaux de transcription. C'est à partir de cette représentation que sont générées les différentes formes d'éditions (visualisations) et c'est à cette représentation que sont rattachées les annotations linguistiques. En revanche, une *syntaxe compacte* est utilisée au stade de saisie et de correction des transcriptions (cf. section 5 ci-dessous).

### Illustration de la structure arborescente

```
<text>
  <front>
  </front>
  <body>
    <div>
      <p>
        <w><choice>
          <me:norm>chevaliers</me:norm>
          <me:dipl>ch<ex>eualie</ex>rs</me:dipl>
          <me:facs>
            <bfm:mdvAbbr>ch&apos;&rrot;</bfm:mdvAbbr>s
          </me:facs>
          <me:pal> ... </me:pal>*2
        </choice></w>
        <bfm:punct><choice>
          <me:norm></me:norm>
          <me:dipl>,</me:dipl>
          <me:facs>&punctelev;</me:facs>
          <me:pal> ... </me:pal>*
        </choice></bfm:punct>
      </p>
    </div>
  </body>
```

<sup>1</sup> Terme proposé par Nicolas Mazziotta.

<sup>2</sup> Dans les exemples de balisage, un astérisque après une balise signifie que l'usage de cette balise est optionnel.

### **Représentation normalisée : <me:norm>**

Ce niveau de représentation est conforme aux normes appliquées dans les éditions critiques modernes (cf. Viellard, Guyotjeannin 2001). La seule différence est qu'il s'agit toujours de la transcription d'un manuscrit, par conséquent toutes les corrections basées sur des leçons des autres manuscrits ou sur des données extratextuelles (corrections de noms géographiques, personnages historiques, etc.) ne sont pas prises en compte. Celles-ci peuvent éventuellement figurer dans des commentaires (élément <note>).

#### **Traits caractéristiques**

- Usage de caractères modernes (suppression des variantes calligraphiques et des diacritiques médiévaux) ;
- Résolution tacite des abréviations ;
- Introduction de la distinction *u/v, i/j* ;
- Usage de diacritiques modernes désambiguïsants (*é, ï, ü*, etc.) ;
- Séparation de ligatures ;
- Notation des noms propres avec les majuscules ;
- Segmentation des unités-mots normalisée ;
- Ponctuation modernisée.

#### **Contenu :**

```
text (PCDATA)
<corr>
<supplied>
```

### **Représentation diplomatique : <me:dipl>**

Cette représentation se situe dans la lignée des éditions dites « diplomatiques », même si les pratiques varient beaucoup selon les éditeurs. Les normes appliquées dans le projet BFM-MSS se basent sur le principe de transcription graphématique : les variantes de graphèmes sont neutralisées, mais on n'introduit pas de marques désambiguïsantes.

#### **Traits caractéristiques**

- Usage de caractères modernes (suppression des variantes calligraphiques et des diacritiques médiévaux) ;
- Résolution des abréviations avec une mise en évidence typographique des caractères restitués ;
- Les caractères *i, j, u* et *v* du manuscrit sont respectés ;
- Pas d'usage de diacritiques modernes désambiguïsants (*é, ï, ü*, etc.) ;
- Séparation de ligatures ;
- Respect des majuscules et des minuscules du manuscrit ;
- Segmentation des unités-mots normalisée ;
- Représentation de la ponctuation de la source : une virgule pour une ponctuation faible et un point pour une ponctuation forte.

## Contenu :

```
texte (PCDATA)
<bfm:expan>
<ex>
<sic>
```

## **Représentation imitative : <me:facs>**

Cette représentation s'inspire de certaines éditions dites « hyper-diplomatiques » ou « imitatives ». Son principe est de représenter fidèlement toutes les données potentiellement pertinentes pour une analyse linguistique du système graphique. Cette représentation ne peut remplacer ni la représentation normalisée, ni la représentation photographique de la source, mais elle les complète en répondant à des besoins de recherches spécifiques.

## Traits caractéristiques

- Distinction de variantes calligraphiques clairement identifiables (allographes), maintien de diacritiques médiévaux ;
- Représentation des marques d'abréviation médiévales ;
- Les caractères *i*, *j*, *u* et *v* du manuscrit sont respectés ;
- Pas d'usage de diacritiques modernes désambiguïsants (*é*, *ï*, *ü*, etc.) ;
- Maintien de certaines ligatures ;
- Représentation des majuscules et des minuscules du manuscrit (éventuellement, des « grandes minuscules » et des « petites majuscules ») ;
- Segmentation des unités-mots originale (maintien des agglutinations et des déglutinations) ;
- Marques de ponctuation médiévales ;
- Représentation des marques de correction sribales.

## Contenu (mots) :

```
texte (PCDATA)
<add>
<am>
<bfm:mdvAbbr>
<bfm:lettrine>
<del>
<unclear>
```

## Contenu (ponctuation)

```
texte (PCDATA)
<bfm:eol>
```

## **Représentation paléographique (facultative) : <me:pal>**

Niveau de représentation prévu pour les annotations paléographiques (distinction fine des tracés de caractères, des traits d'écriture individuels, des espaces de différentes tailles, etc.). Il n'est pas implémenté à présent dans les transcriptions BFM-MSS.

### 3. Corrections sribales et éditoriales ; passages difficilement lisibles ou illisibles

Les transcriptions du projet BFM-Manuscrits visent à représenter les manuscrits sans corriger les erreurs éventuelles des scribes ou copistes. Par conséquent, les seules corrections y intégrées de façon systématique sont celles effectuées sur le manuscrit par le scribe lui-même ou par un (re)lecteur médiéval. Des corrections éditoriales de coquilles évidentes peuvent toutefois être ajoutées au niveau de représentation normalisée. Les mêmes principes de balisage s'appliquent aux passages du manuscrit difficilement lisible ou illisible à cause de l'état matériel du manuscrit. Les balises XML recommandées par la TEI sont utilisées pour toutes les corrections et passages illisibles.

Le niveau hiérarchique sur lequel apparaissent les balises de corrections et assimilées dépend de l'étendu de celles-ci. Trois cas de figure se présentent :

#### *Partie d'un mot*

Les éléments correspondants (<corr>, <add>, <del>, <unclear>, etc.) sont placés à l'intérieur de la représentation facsimilaire, diplomatique ou normalisée selon les cas.

#### Exemple

```
<w><choice>
  <me:norm>Platon</me:norm>
  <me:dipl>platon</me:norm>
  <me:facs>p<del rend="dots_below"
hand="#scribe1">i</del>laton</me:facs>
</choice></w>
```

#### *Un mot entier ou plusieurs mots*

Les mots (<w>) concernés sont placés dans les éléments correspondants (<corr>, <add>, <del>, <unclear>, etc.). Dans les passages supprimés, seul le niveau de représentation facsimilé sera présent.

#### Exemple

```
<del rend="dots_below" hand="#scribe1">
  <w type="num">
    <me:facs>.ii.j.</me:facs>
  </w>
  <w>
    <bfm:mdvAbbr><me:facs>phi<am>&apos;. </am><bfm:mdvAbbr></me:facs>
  </w>
  <space dim="horizontal" unit="chars" quantity="9"/>
  <w>
    <me:facs>in</me:facs>
  </w>
  <w>
    <me:facs>p&rrot;ologo</me:facs>
  </w>
</del>
```

## Segment commençant ou finissant au milieu d'un mot et s'étendant sur plusieurs mots

Ce cas, assez rare, présente un problème de structure hiérarchique XML. La technique de « fragmentation et reconstitution d'éléments virtuels » est utilisée pour le résoudre (cf. [TEI Guidelines, chapitre 20, section 3](#)), qui semble satisfaisante pour le traitement de cas rares.

### Exemple

```
<w>
  <norm>lairons</me:norm>
  <me:dipl>l<unclear xml:id="ucl123"
next="#ucl124">airons</unclear></me:dipl>
  <me:fac>l<unclear corresp="#ucl123" next="#ucl124">airons</unclear>
  </me:fac>
</w>
<unclear xml:id="ucl124" prev="#ucl123">
  <w>
    <norm>nous</me:norm>
    <me:dipl>nous</me:dipl>
    <me:fac>nous</me:fac>
  </w>
</unclear>
```

### Notes et commentaires du relecteur/encodeur

Une `<note>` peut être utilisée pour ajouter un commentaire à l'endroit du texte où l'encodage pose un problème ou si le choix de balisage nécessite une explication. L'élément `<note>` est normalement placé au même niveau d'arborescence XML que les mots (`<w>`), à la fin du passage commenté.

## 4. Définition et description des éléments

Par défaut, les éléments ci-dessous sont définis par la TEI. Les éléments définis par le projet Menota portent un préfixe `me:` (`xmlns:me="http://www.menota.org/ns/1.0"`). Les éléments définis par le projet BFM-MSS portent le préfixe `bfm:` (`xmlns:bfm="http://bfm.ens-lsh.fr/ns/1.0"`).

### **<ab>**

Balise générique pour une unité textuelle de niveau équivalent au paragraphe. Dans la BFM, cette balise est utilisée dans les textes en vers.

La structure linguistique (en particulier, la division en phrases ou « unités ponctuables ») est prioritaire dans la BFM. Pour éviter le conflit de « hiérarchies concurrentes », les éléments `<l>` (vers) et `<lg>` (groupe de vers) ne sont donc pas utilisés.

### Position (« parents » ou « frères et sœurs »)

A l'intérieur d'une division (`<div type="chapitre">` ou autre).

## Contenu (« enfants »)

<w>  
<bfm:punct>  
<lb>, etc.

### Attributs

- **type** : type d'unité ("**gv**" pour groupe de vers) ;
- **n** : indique éventuellement le numéro du paragraphe.

## <add>

Lettres, mots ou groupes de mots ajoutés au texte initial dans le manuscrit par le scribe ou autre personne à l'époque proche de la création du manuscrit.

### Attributs

- **place** : indique l'endroit où le texte ajouté est placé. La TEI propose une liste de valeurs non restrictive ("**infralinear**", "**supralinear**", "**margin-right**", "**inline**", etc.). Une valeur "**overwritten**" est ajoutée par le projet BFM-MSS. Attribut obligatoire ;
- **hand** : indique la main ayant effectué l'ajout. Si absent, on suppose que c'est le premier scribe ;
- **resp** : indique le responsable du balisage et d'attribution de l'ajout
- **seq** : indique éventuellement l'ordre séquentiel en cas de plusieurs « couches » de corrections dans un manuscrit.

## Position et contenu

Cf. section 3 ci-dessus.

## <am>

Marque d'abréviation

### Position

Dans l'élément <bfm:mdvAbbr> (représentation facsimilaire).

### Contenu

PCDATA (un caractère MUFI représenté par une entité, cf. section 6).

### Attributs

- **resp** – responsable du balisage (identification de la marque).

### Exemples

Cf. <bfm:mdvAbbr> plus bas.

## <bfm:headlb>

Saut de ligne « supplémentaire » dans un titre.

Cet élément est utilisé dans un cas particulier de la mise en page médiévale. Dans certains manuscrits les titres ou résumés de chapitres sont écrits en encre de couleur différente et placés immédiatement après la fin du chapitre précédent, souvent sans retour à la ligne. Parfois ces titres continuent sur les fins des premières lignes du chapitre suivant, ce qui pose un problème de continuité linéaire du texte (cf. illustration ci-dessous). La solution proposée est de placer la totalité du titre avant le début du chapitre annoncé et d'utiliser une balise spéciale là où une partie du titre se trouve sur la même ligne physique que le contenu d'une autre unité sémantique du texte (fin de la division précédente ou début de la division suivante). En revanche une simple balise `<lb>` est utilisée lorsqu'une ligne du titre occupe la totalité d'une ligne physique.



Fragment de *l'Image du monde*, ms. Paris, BNF, fr. 574, f° 83r (cf. exemple de code ci-dessous)

## Position

Dans un titre (ou résumé) d'un chapitre ou autre unité textuelle (`<head>`).

Comme `<lb>`, cet élément marque en effet le début (et non la fin d'une ligne). La première balise se trouve donc au début du titre.

## Contenu

Élément vide.

## Attributs

- **lines** – nombre de lignes de décalage par rapport à la position normale (fin de la dernière ligne de la division précédente. Sa valeur est "0" si le titre commence à la même ligne, elle est négative si le titre commence avant la dernière ligne de la division précédente et positive si le titre continue sur des fins de ligne du début de la division suivante.

## Exemple

(cf. illustration ci-dessus)

```
<lb/>est desus ia ne sera si fo<bfm:lb/>rt
<lb/>tour . </p></div>
<div><head><bfm:headlb lines="0"/>comment la terre
<bfm:headlb lines="1"/>crolle et fent . </head>
<p>
<lb/>ORe
<lb/>entendez donques
<lb/>du mouuement que
```

## <bfm:hyphen>

Marque de coupure de mot en fin de ligne

### Position dans la structure

Représentation facsimilaire : <me:facs>.

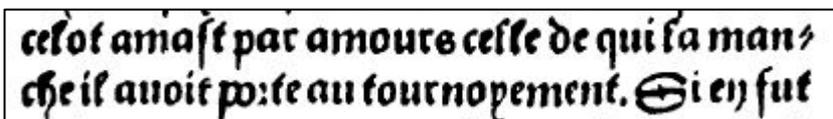
### Contenu

Un caractère (représenté souvent par une entité MUFI `&sol;`, `&dbloblhyphen;`)

### Attributs

- `place` – indique éventuellement la position de la marque de coupure ("`margin`" ou, par défaut, "`inline`").

### Exemples



```
la man<bfm:hyphen>&dbloblhyphen;</bfm:hyphen><lb/>che
```

## <bfm:lettrine>

Lettre initiale (lettrine).

### Position dans la structure

Représentation imitative : <me:facs>.

### Contenu

Un caractère (PCDATA)

### Attributs

- `color` – couleur (principale) de l'initiale ;
- `colorSuppl` – couleur secondaire ;
- `size` – taille prévue : nombre de lignes d'alinéa ;
- `sizeAct` – taille réelle de l'initiale en nombre de lignes ;
- `decoration` – détails de la décoration ;
- `ref` – référence de la description complète de l'initiale ;
- `resp` – responsable du balisage.

### Exemples



```
<bfm:lettrine color="blue" colorSuppl="red" size="6"
sizeAct="6">O</bfm:lettrine>R
```

### <bfm:mdvAbbr>

Abréviation médiévale. Cette balise marque le segment graphique « affecté » par l'abréviation. Selon le type d'abréviation, il peut s'étendre de la marque d'abréviation proprement dite à une ou plusieurs lettres associées à la marque, voire à un mot entier. Le fait que les limites de l'abréviation ne coïncident pas forcément avec celles d'un mot nous empêche d'utiliser l'élément <abbr> de la TEI.

### Position dans la structure

Dans la représentation facsimilaire : <me:fac>

### Contenu

caractères (PCDATA)

<am>

### Attributs

- **range** – précision de l'étendu de l'abréviation ("chars", "word" or "phrase") ;
- **type** – type d'abréviation ("regular", "frequent\_word", etc.) ;
- **resp** – responsable du balisage (en particulier, de la délimitation de l'abréviation) ;
- **cert** – certitude de l'identification de l'abréviation.

### Exemples<sup>3</sup>

```
<bfm:mdvAbbr><am>&et ; </am></bfm:mdvAbbr>
h<bfm:mdvAabbr>o<am>&bar ; </am></bfm:mdvAbbr>me
<bfm:mdvAbbr>ml<am>&apos ; </am>t</bfm:mdvAbbr>
```

<sup>3</sup> Tous les exemples cités sont tirés du corpus BFM-MSS. Pour améliorer la lisibilité, seules les balises pertinentes pour l'exemple en question sont conservées. Sauf indication contraire, les graphies sont représentées au niveau diplomatique.

## **<bfm:punct>**

Marque de ponctuation ou toute « frontière ponctuable ».

Cette balise sert à la fois à représenter les marques de ponctuation utilisées dans le manuscrit, les marques ajoutées au niveau « normalisé » pour faciliter la lecture du texte et à repérer les « frontières ponctuables ».

Une balise **<me:punct>** est proposée par le projet Menota, mais elle n'a pas exactement les mêmes propriétés que la balise **<bfm:punct>**.

## **Position**

Au même niveau que les mots (**<w>**).

## **Contenu**

Un élément **<choice>** ou l'une des représentations multiples : **<me:norm>**, **<me:dipl>**, **<me:fac>**.

## **Attributs**

- **force** : indique la force de la ponctuation ("strong" ou "weak") ;
- **mark** : donne le nom conventionnel de la marque de ponctuation médiévale, sa valeur est "none" si aucune marque de ponctuation médiévale n'est présente dans la source (une marque peut apparaître en revanche au niveau normalisé ou diplomatique) ;
- **ana** : le code du type de frontière ponctuable.

## **<bfm:sb>**

Espace blanc à l'intérieur d'un mot (« déglutination »).

## **Position dans la structure**

Dans la représentation facsimilaire : **<me:fac>** d'un mot.

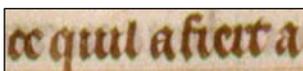
## **Contenu**

Élément vide

## **Attributs**

- **cert** : certitude de la présence d'un espace ;  
NB : L'absence de certitude de déglutination signifie normalement qu'il y a un « petit blanc ».
- **resp** : responsable du balisage.

## **Exemples**



```
ce quil a<mdv_sb cert="no"/>fiert a
```

## **<cb>**

Saut de colonne.

## Position et contenu

Élément vide, situé entre des mots ou des divisions de niveaux supérieurs.

## Attributs

- **n** : indique le numéro de la colonne (habituellement "a", "b", "c" ou "d").

## <choice>

Balise TEI servant à délimiter les représentations alternatives d'un même segment du texte. Dans le projet BFM-MSS, à l'exemple du projet Menota, cette balise entoure les représentations alternatives d'un mot ou d'une marque de ponctuation.

## Position

A l'intérieur d'un <w> ou d'un <bfm:punct>.

## Contenu

<me:norm>, <me:dipl> et <me:fac>.

## Attributs

Pas d'attributs.

## <corr>

Lettres, mots ou groupes de mots corrigés par l'éditeur ou par un relecteur/encodeur dans la représentation normalisée.

## Position et contenu

Cf. section 3 ci-dessus.

## Attributs

- **cert** : indique la certitude de la correction ;
- **resp** : indique le responsable du balisage et d'attribution de la correction.

## <damage>

Fragment du manuscrit endommagé.

## Position et contenu

Cf. section 3 ci-dessus.

## Attributs

- **extent** : indication de la taille du fragment endommagé ;
- **resp** : indique le responsable du balisage et d'attribution de la correction.

## <del>

Lettres, mots ou groupes de mots supprimés dans le manuscrit.

## Position et contenu

Cf. section 3 ci-dessus.

### Attributs

- **rend** : indique la méthode de suppression utilisée. Les valeurs utilisées sont les suivantes : "dots\_below", "horbar", "diagbar", etc. ;
- **hand** : indique la main ayant effectué l'ajout, si identifiable ;
- **resp** : indique le responsable du balisage et d'attribution de l'ajout.

### <div>

Division structurelle du texte.

### Position

A l'intérieur du corps du texte <body> ou d'une autre <div>.

### Contenu

Eventuellement <head> au début ;  
<p> ou <ab>.

### Attributs

- **rend** : indique la méthode de suppression utilisée. Les valeurs utilisées sont les suivantes : "dots\_below", "horbar", "diagbar", etc. ;
- **hand** : indique la main ayant effectué l'ajout, si identifiable ;
- **resp** : indique le responsable du balisage et d'attribution de l'ajout.

### <ex>

Lettres restituées dans la résolution d'une abréviation.

### Position dans la structure

Dans la représentation diplomatique.

### Contenu

Une ou plusieurs lettres (PCDATA)

### Attributs

- **cert** : certitude de la résolution ;
- **resp** : responsable de la restitution des lettres.

### Exemples

```
ch<ex>eualie</ex>r  
ih<ex>e</ex>r<ex>usa</ex>l<ex>e</ex>m
```

## **<gap>**

Fragment de texte omis : soit parce qu'il est illisible dans le manuscrit, soit pour des raisons d'échantillonnage.

### **Position et contenu**

Cf. section 3 ci-dessus.

### **Attributs**

- **unit** : indique l'unité de mesure : "**chars**", "**words**", "**pages**", etc. ;
- **extent** : indication numérique de l'étendu (en unités choisies) ;
- **reason** : indique la raison de l'omission : "**illegible**", "**deleted**" ou "**sampling**" ;
- **hand** : en cas d'effacement volontaire, indique le responsable, si identifiable ;
- **resp** : indique le responsable du balisage et d'attribution de l'ajout.

## **<head>**

Titre ou autre élément liminaire d'une division textuelle (chapitre, partie, etc.).

### **Position et contenu**

Début d'une division **<div>**.

### **Contenu**

Mots **<w>** et autres unités du même niveau ;  
**<bfm:head1b>** dans un cas particulier (cf. la définition de cet élément).

### **Attributs**

Attribut globaux seulement.

## **<hi>**

Segment de texte mise en relief.

### **Position et contenu**

Cf. section 3 ci-dessus.

### **Attributs**

- **rend** : indique la forme de mise en relief : "**exp**", "**color:red**", etc.

## **<lb>**

Saut de ligne.

### **Position et contenu**

Élément vide, situé entre des mots (sauf car particulier).

Contrairement à l'apparence, cette balise marque un début et non la fin d'une ligne. Elle est par conséquent placée au début de la première ligne et est absente à la fin de la dernière.

## Attributs

- **n** : indique le numéro de la ligne dans la page ou dans la colonne.

### Cas particulier : saut de ligne coupant un mot.

L'attribut **ed="facs"** est utilisé à l'intérieur de la représentation imitative à l'endroit précis du saut de ligne dans le manuscrit. Une deuxième balise **<lb>** est placée après le mot coupé, portant un attribut **ed="norm"**.

## Exemples

```
<lb n="4"/> que parfaite proëce estoit entee et en<lb ed="facs"/>rachinee  
<lb ed="norm" n="5"/> au plus fort ez cuers des nobles
```

### **<me:dipl>**

Représentation diplomatique d'un mot ou d'une marque de ponctuation.

### Position et contenu

Cf. section 0 ci-dessus.

### **<me:facs>**

Représentation « facsimilé » d'un mot ou d'une marque de ponctuation. Élément défini par le projet Menota.

### Position et contenu

Cf. section 0 ci-dessus.

### **<me:norm>**

Représentation normalisée d'un mot ou d'une marque de ponctuation.

### Position et contenu

Cf. section 0 ci-dessus.

### **<me:pal>**

Représentation « paléographique » d'un mot ou d'une marque de ponctuation. Élément défini par le projet Menota.

Le projet BFM-MSS ne prend pas en charge ce niveau de représentation au stade actuel.

### **<note>**

Note ou commentaire du relecteur/encodeur. Élément défini par la TEI.

## Position dans la structure

A l'intérieur de paragraphes `<p>`.

## Contenu

Texte brut (métatexte par rapport à la transcription).

### `<p>`

Paragraphe (dans les textes en prose).

## Position

A l'intérieur d'une division (`<div type="chapitre">` ou autre).

## Contenu

`<w>`

`<bfm:punct>`

`<lb>`, etc.

## Attributs

- `n` : indique éventuellement le numéro du paragraphe.

### `<pb>`

Saut de page (ou de folio).

## Position

Entre des mots ou entre des divisions de niveau supérieur.

Contrairement à l'apparence, cette balise marque un début et non la fin d'une page. Elle est par conséquent placée au début de la première page et est absente à la fin de la dernière.

## Contenu

Élément vide.

## Attributs

- `n` : indique le numéro de la page ou du folio.

### `<q>`

Balise du discours direct.

## Position

A l'intérieur des paragraphes. Si le discours direct dépasse un paragraphe, le mécanisme de « fragmentation et reconstitution d'éléments virtuels » est utilisé (cf. section 3).

## Contenu

Éventuellement, des segments ponctuables `<seg type="sp">` ;

Mots `<w>` et autres éléments de ce niveau.

## Attributs

- **type** : type de discours direct ("**spoken**" par défaut, "**written**" ou autre éventuellement).

## Exemple

A quoy il respondy , sans dire wy , ou nanil , `<q>` " par foy &slong;ire il souffrira assez se je voy ce qui se fera , " `</q>` et ce disoit le conte...

## `<seg>`

Élément générique pour la segmentation linguistique du texte.

## Position

A l'intérieur des paragraphes.

## Contenu

Mots, marques de ponctuation, etc.

## Attributs

- **type** : type de segment ("**sp**" pour « segment ponctuable »).  
La limite de `<seg type="sp">` est placée avant la marque de ponctuation. La marque de ponctuation est ainsi intégrée à l'unité textuelle qui la suit.

## Exemple

```
<seg type="sp" ana="#d9 #z4"> ... voudroy je parler </seg> <seg type="sp" ana="#c4 #c3">, C'est de chace ... </seg>
```

## `<space>`

Espace blanc de taille significative.

## Position dans la structure

Entre des mots.

## Contenu

Élément vide.

## Attributs

- **dim** : dimension ("**vertical**" ou "**horizontal**");
- **unit** : unité de mesure ("**chars**", "**lines**", etc.);
- **quantity** : valeur numérique ;
- **resp** : responsable du balisage.

## `<subst>`

Regroupe les éléments `<del>` et `<add>` en cas de correction sribale.

## Position

Représentation imitative.

## Contenu

<del> et <add>.

## <supplied>

Lettres, mots ou groupes de mots ajoutés par l'éditeur ou par un relecteur/encodeur pour corriger les coquilles du manuscrit.

## Position et contenu

Représentation normalisée, cf. section 3 ci-dessus.

## <unclear>

Lettres, mots ou groupes de mots difficilement lisibles dans le manuscrit. Élément défini par la TEI.

## Position et contenu

Cf. section 3 ci-dessus.

## <w>

Mot lexical ou grammatical. Élément défini par la TEI ; modèle du contenu modifié par le projet Menota ; attributs ajoutés par le projet BFM-MSS.

## Position dans la structure

Unité principale de structuration linguistique. Utilisé à l'intérieur de paragraphes <p>.

## Contenu

Un élément <choice> ou l'une des représentations multiples : <me:norm>, <me:dipl>, <me:fac>.

## Attributs

- **bfm:aggl** – indique l'absence de blanc avant du mot suivant (valeurs : *elision*, *simple*) ;
- **bfm:agglCert** – certitude de l'agglutination ;
- **bfm:agglResp** – responsable de l'identification de l'agglutination ;
- **lemma** – référence à un lemme ;
- **ana** – étiquette morpho-syntaxique ;
- **type**.

## 5. Saisie et correction des transcriptions (syntaxe compacte)

Les transcriptions sont saisies au niveau imitatif, avec des méta-caractères non-XML permettant de générer les niveaux diplomatique et normalisé.

Si la transcription s'appuie sur le texte d'une édition de référence, celui-ci est traité par un script permettant de le « diplomatiser » tout en gardant les informations nécessaires pour le niveau de représentation normalisée.

Même si elles contiennent des annotations non-XML, les transcriptions compactes sont des documents XML validables contre un schéma conforme à la TEI.

### Caractères de raccourci, balises simplifiées

#### Segmentation

- + – agglutination simple ;
- +? – agglutination simple « pas sure » ou « petit blanc » ;
- ´ – élision ;
- \_ – déglutination ;
- \_? – déglutination « pas sure » ou « petit blanc ».

#### Abréviations

La portée de l'abréviation (voir la définition plus haut) est marquée par des parenthèses doubles.

Si la résolution de l'abréviation est régulière (automatique à partir de l'association de la marque et des lettres dans la portée de l'abréviation, seule la marque de l'abréviation est donnée (sous la forme d'une entité MUFI).

Syntaxe compacte	Représentation diffractée
((&et;))	<pre>&lt;norm&gt;&lt;/me:norm&gt; &lt;me:dipl&gt;&lt;ex&gt;et&lt;/ex&gt;&lt;/me:dipl&gt; &lt;me:facs&gt;   &lt;bfm:mdvAbbr&gt;&lt;am&gt;&amp;et;&lt;/am&gt;&lt;/bfm:mdvAbbr&gt; &lt;/me:facs&gt;</pre>
((o&bar;))	<pre>&lt;norm&gt;on&lt;/me:norm&gt; &lt;me:dipl&gt;o&lt;ex&gt;n&lt;/ex&gt;&lt;/me:dipl&gt; &lt;me:facs&gt;   &lt;bfm:mdvAbbr&gt;o&lt;am&gt;&amp;bar;&lt;/am&gt;&lt;/bfm:mdvAbbr&gt; &lt;/me:facs&gt;</pre>
...((&m&dblbar;t))	<pre>&lt;norm&gt;...ment&lt;/me:norm&gt; &lt;me:dipl&gt;...m&lt;ex&gt;en&lt;/ex&gt;t&lt;/me:dipl&gt; &lt;me:facs&gt;   ...&lt;bfm:mdvAbbr&gt;m&amp;dblbar;t&lt;/bfm:mdvAbbr&gt; &lt;/me:facs&gt;</pre>

Si la résolution n'est pas « automatique », celle-ci est ajoutée, précédée d'une marque de soulignement. Les lettres « restituées » sont placées entre crochets.

Syntaxe compacte	Représentation diffractée
((e&bar;_e[st]))	<pre>&lt;me: norm&gt;est&lt;/me: norm&gt; &lt;me: dipl&gt;e&lt;ex&gt;st&lt;/ex&gt;&lt;/me: dipl&gt; &lt;me: facs&gt;   &lt;bfm: mdvAbbr&gt;     e&lt;am&gt;&amp;bar;&lt;/am&gt;   &lt;/bfm: mdvAbbr&gt; &lt;/me: facs&gt;</pre>
((ch&apomod;&errot;_ch[evalie]r))	<pre>&lt;me: norm&gt;chevalier&lt;/me: norm&gt; &lt;me: dipl&gt;ch&lt;ex&gt;evalie&lt;/ex&gt;r&lt;/me: dipl&gt; &lt;me: facs&gt;   &lt;bfm: mdvAbbr&gt;     ch&lt;am&gt;&amp;apomod;&lt;/am&gt;&amp;errot;   &lt;/bfm: mdvAbbr&gt; &lt;/me: facs&gt;</pre>
((nr&apomod;e_n[ost]re))	<pre>&lt;me: norm&gt;nostre&lt;/me: norm&gt; &lt;me: dipl&gt;n&lt;ex&gt;ost&lt;/ex&gt;re&lt;/me: dipl&gt; &lt;me: facs&gt;   &lt;bfm: mdvAbbr&gt;     nr&lt;am&gt;&amp;apomod;&lt;/am&gt;e   &lt;/bfm: mdvAbbr&gt; &lt;/me: facs&gt;</pre>

## Grandes initiales (letrines)

Les grandes initiales sont entourées de doubles accolades. La lettre est séparée de ses attributs (nombre de lignes en retrait, taille réelle en nombre de lignes, couleur principale, éventuellement la décoration) par des deux-points :

Syntaxe compacte	Représentation diffractée
{{Q:2:7:blue}}uant	<pre>&lt;me: norm&gt;Quant&lt;/me: norm&gt; &lt;me: dipl&gt;&lt;hi rend="initiale" &gt;Q&lt;/hi&gt;uant&lt;/me: dipl&gt; &lt;me: facs&gt;   &lt;bfm: letrine size="2" sizeAct="7" color="blue"&gt;     Q&lt;/bfm: letrine&gt;uant &lt;/me: facs&gt;</pre>

## Variantes de lettres

Entités MUFI

(voir la liste complète dans le *Tableau d'encodage des caractères*, section 6)

## Lettres à normaliser

Les lettres qui doivent être modifiées dans la représentation normalisée sont précédées de symboles \* et # :

- \*u – u diplomatique, à transformer en v
- \*v – v diplomatique, à transformer en u
- \*i – i diplomatique, à transformer en j
- \*j – j diplomatique, à transformer en i

Le dièse (#) est utilisé devant les minuscules à transformer en majuscules (noms propres, débuts de phrases dans la représentation normalisée).

Si les deux marques se combinent, le dièse est placé avant l'astérisque.

Les diacritiques de la représentation normalisée (é, ï, ü, etc.) sont utilisés dans la saisie initiale.

## Ponctuation

Les marques de ponctuation sont balisées `<c type="punct">` (le préfixe `bfm:` est ajouté ultérieurement). Les trois représentations sont fournies, séparées par % :

```
<punct>.%,%&punctelev; </punct>
```

Selon les cas, une ou deux représentations peuvent être vides :

```
<punct>,%% </punct>
<punct>%,%.</punct>
```

Le premier exemple représente une de ponctuation ajoutée par l'éditeur, qui ne correspond à aucune marque de la source. Le deuxième exemple représente une marque présente dans la source que l'éditeur a choisi d'éliminer dans la représentation normalisée.

## Corrections sribales

La syntaxe compacte est applicable pour le balisage des corrections sribales (suppressions et ajouts de mots ou de caractères) qui concernent un mot ou une partie d'un mot. Les balises XML TEI correspondantes sont doivent être utilisées pour des corrections de taille supérieure. Les marques utilisées pour noter les corrections sribales en syntaxe compacte sont basées sur les recommandations de F. Masai (1950) avec quelques modifications visant à les rendre plus systématiques.

De façon générale, les doubles crochets servent à délimiter le segment concerné par la correction. Une barre oblique indique l'ajout sur la ligne (à la place du texte supprimé ou d'un espace blanc). Une barre oblique inversée `<\ >` indique l'ajout au-dessus de la ligne. Un tiret signifie une suppression par rature ; un point souscrit (entité de type `&xdotbl;`) signifie une suppression par exponctuation.

Syntaxe compacte	Représentation diffractée
<ul style="list-style-type: none"> <li>addition interlinéaire</li> </ul> <pre>[[\après/]]</pre>	<pre>&lt;me: norm&gt;après&lt;/me: norm&gt; &lt;me: dipl&gt;apres&lt;/me: dipl&gt; &lt;me: facts&gt;   &lt;add place="interlinear"&gt;     apres   &lt;/add&gt; &lt;/me: facts&gt;</pre>
<ul style="list-style-type: none"> <li>addition sur la ligne</li> </ul> <pre>[[/et\]]</pre>	<pre>&lt;me: norm&gt;et&lt;/me: norm&gt; &lt;me: dipl&gt;et&lt;/me: dipl&gt; &lt;me: facts&gt;   &lt;add place="inline"&gt;et&lt;/add&gt; &lt;/me: facts&gt;</pre>
<ul style="list-style-type: none"> <li>addition en marge</li> </ul>	<pre>&lt;me: norm&gt;qui&lt;/me: norm&gt;</pre>

[[\qui//]]	<pre>&lt;me:dipl&gt;qui&lt;/me:dipl&gt; &lt;me:facs&gt;   &lt;add place="margin"&gt;qui&lt;/add&gt; &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• suppression par rature (lettre encore lisible)</li> </ul> vi[[ - r ]]ent	<pre>&lt;me:norm&gt;vient&lt;/me:norm&gt; &lt;me:dipl&gt;vient&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;del rend="line-through"&gt;r&lt;/del&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• suppression par rature (lettre illisible)</li> </ul> vi[[ - ]]ent	<pre>&lt;me:norm&gt;vient&lt;/me:norm&gt; &lt;me:dipl&gt;vient&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;del rend="line-through"&gt;     &lt;gap/&gt;   &lt;/del&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• suppression par exponctuation</li> </ul> vi[[ &rdotbl; ]]ent	<pre>&lt;me:norm&gt;vient&lt;/me:norm&gt; &lt;me:dipl&gt;vient&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;del rend="dotbl"&gt;&amp;rdotbl;&lt;/del&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• suppression par tout autre procédé (lettre encore lisible)</li> </ul> vi[[ r ]]ent	<pre>&lt;me:norm&gt;vient&lt;/me:norm&gt; &lt;me:dipl&gt;vient&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;del&gt;r&lt;/del&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• suppression par tout autre procédé (lettre illisible)</li> </ul> vi[[ ]]ent	<pre>&lt;me:norm&gt;vient&lt;/me:norm&gt; &lt;me:dipl&gt;vient&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;del&gt;&lt;gap/&gt;&lt;/del&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• substitution au dessus d'une rature (lettre encore lisible)</li> </ul> vi[[ - r \ &slong; ]]ent	<pre>&lt;me:norm&gt;visent&lt;/me:norm&gt; &lt;me:dipl&gt;visent&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;subst&gt;     &lt;del rend="line-through"&gt;       r     &lt;/del&gt;     &lt;add place="interlinear"&gt;       &amp;slong;     &lt;/add&gt;   &lt;/subst&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• substitution au dessus d'une rature (lettre supprimée illisible)</li> </ul> vi[[ - \ &slong; ]]ent	<pre>&lt;me:norm&gt;visent&lt;/me:norm&gt; &lt;me:dipl&gt;visent&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;subst&gt;     &lt;del rend="line-through"&gt;       &lt;gap/&gt;     &lt;/del&gt;     &lt;add place="interlinear"&gt;       &amp;slong;     &lt;/add&gt;   &lt;/subst&gt;ent &lt;/me:facs&gt;</pre>
<ul style="list-style-type: none"> <li>• substitution au dessus d'une exponctuation</li> </ul> vi[[ &rdotbl; \ &slong; ]]ent	<pre>&lt;me:norm&gt;visent&lt;/me:norm&gt; &lt;me:dipl&gt;visent&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;subst&gt;</pre>

	<pre> &lt;del rend="dotbl"&gt;&amp;rdotbl;&lt;/del&gt; &lt;add place="interlinear"&gt;   &amp;slong; &lt;/add&gt; &lt;/subst&gt;ent &lt;/me:facs&gt; </pre>
<ul style="list-style-type: none"> <li>substitution par transformation</li> </ul> <pre>[[e &gt; i]]n</pre>	<pre> &lt;me:norm&gt;in&lt;/me:norm&gt; &lt;me:dipl&gt;in&lt;/me:dipl&gt; &lt;me:facs&gt;   &lt;subst&gt;     &lt;del rend="transform"&gt;e&lt;/del&gt;     &lt;add place="overwrite"&gt;i&lt;/add&gt;   &lt;/subst&gt;n &lt;/me:facs&gt; </pre>
<ul style="list-style-type: none"> <li>substitution par superposition d'une lettre écrite sur une autre</li> </ul> <pre>[[e + i]]n</pre>	<pre> &lt;me:norm&gt;in&lt;/me:norm&gt; &lt;me:dipl&gt;in&lt;/me:dipl&gt; &lt;me:facs&gt;   &lt;subst&gt;     &lt;del rend="unmarked"&gt;e&lt;/del&gt;     &lt;add place="overwrite"&gt;i&lt;/add&gt;   &lt;/subst&gt;n &lt;/me:facs&gt; </pre>
<ul style="list-style-type: none"> <li>substitution par superposition d'une lettre écrite au dessus d'une autre (non expressément supprimée)</li> </ul> <pre>[[e + \ i]]n</pre>	<pre> &lt;me:norm&gt;in&lt;/me:norm&gt; &lt;me:dipl&gt;in&lt;/me:dipl&gt; &lt;me:facs&gt;   &lt;subst&gt;     &lt;del rend="unmarked"&gt;e&lt;/del&gt;     &lt;add place="interlinear"&gt;i&lt;/add&gt;   &lt;/subst&gt;n &lt;/me:facs&gt; </pre>
<ul style="list-style-type: none"> <li>substitution sur grattage ou autre type de suppression (lettre supprimée encore lisible)</li> </ul> <pre>vi[[r / &amp;slong;]]ent</pre>	<pre> &lt;me:norm&gt;visent&lt;/me:norm&gt; &lt;me:dipl&gt;visent&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;subst&gt;     &lt;del&gt;r&lt;/del&gt;     &lt;add place="overwrite"&gt;&amp;slong;&lt;/add&gt;   &lt;/subst&gt;ent &lt;/me:facs&gt; </pre>
<ul style="list-style-type: none"> <li>substitution sur grattage ou autre type de suppression (lettre supprimée illisible)</li> </ul> <pre>vi[[ / &amp;slong;]]ent</pre>	<pre> &lt;me:norm&gt;visent&lt;/me:norm&gt; &lt;me:dipl&gt;visent&lt;/me:dipl&gt; &lt;me:facs&gt;   vi&lt;subst&gt;     &lt;del&gt;&lt;gap/&gt;&lt;/del&gt;     &lt;add place="overwrite"&gt;&amp;slong;&lt;/add&gt;   &lt;/subst&gt;ent &lt;/me:facs&gt; </pre>



## 7. Tableaux d'encodage des caractères « spéciaux »

### *Lettres initiales (letrines)*

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
				Encodé au moyen d'éléments

### *Variantes de caractères*

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
ı	= i	&inodot;	0131	<i>i</i> simple utilisé dans la transcription, le caractère sans point apparaît seulement lors de la visualisation facsimilaire
í	= i	&iacute;	00ED	
ð	= d	&drot;	F109	
ɱ	[final m]	&mrdes;	F223 (v2)	
ŋ	&n;	&nrdes;	F228 (v2)	
ʒ	&r;	&rrot;	F20E	
f	&s;	&slong;	017F	
ς	&s2;	--		~ &stigma; 03C2
	&s3;	--		→ &slong;
	&s4;	--		→ &slong;
ε	s (final)	&sclose;	F128	non utilisé
ÿ	&y-dot;	&ydot;	1E8F	

### *Ligatures*

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
pp	--	&pplig;	EED6	
ct		&ctlig;	EEC5	
ft		&ftlig;	EECB	
ft		&slongtlig;	FB05	

### *Caractères exponctués*

#### point souscrit

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
ḋ		&ddotbl;	1E0D	

## Abréviations

### caractères spéciaux

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
7	&et; &et1;	&et;	204A	
ε		&etslash;	F158	
7	&et2;	&ET ;	F142	
	&et3;			
3	--	&etfin;	F155	
÷	&est;	&est;	223B	
°	&nine;	&usmod;	F151	cf. &us; si sur une lettre
9	&com;	&condes;	F156	
f	--	&is;	F15A	
℥	--	&rum;	F154	

### caractères suscrits (combinés)

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
$\overset{a}{x}$	&x <sup>4</sup> -alpha;	x&asup;	0363	
$\overset{e}{x}$	e	x&esup;	0364	
$\overset{o}{x}$	o	x&osup;	0366	
$\overset{u}{x}$	u	x&usup;	0367	
$\overset{s}{x}$	["]	x&ssup;	F027	??? « stigma » suscrit

### caractères diacritiques

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
x <sup>l</sup>	&x-vbar;	x&verbarup;	02C8	
$\overline{x}$	&x-hbar;	x&bar;	0305	
$\overline{\overline{xx}}$		*&dblbar;	035E	Absent MUFI, police Cardo
$\tilde{x}$	&x-tilde;	&combtilde;	0303	
$\hat{x}$		*&combinvbreve;	0311	Absent MUFI, police Cardo
$\widehat{xx}$		*&combdblinvbreve;	0361	Absent MUFI, police Cardo
$\ddot{x}$	&x-omeg;	&ra;	F157	
x <sup>'</sup> x	&apost;	&apomod;	02BC	
x <sup>˘</sup> x		&combcomma;	0315	code plus correct pour l'apostrophe abréviatif
$\acute{x}$	&ier;	&combtildevert;	033E	
$\tilde{x}$	--	&er;	035B	à distinguer de tildevert ?
$\tilde{x}$	--	&ercurl;	F1C8	pb espace après Andron
$\hat{x}$	--	&ur;	F1C3	^ dans trscr 2001
$\overset{2}{x}$	--	&urrot;	F153	
$\overset{2}{x}$	--/ &nine;	&us;	F15B	cf. &usmod;

<sup>4</sup> « x » signifie ici n'importe quel caractère alphabétique.

## lettres modifiées

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
ø	d°	&de;	F159	
þ	&par;	&pbardes;	E670	
Ɔ	&pro; &pro1;	&pflour;	E67D	
	&pro2;			???
ø		&qslstrok;	E8B1	
ƒ	--	&slongslstrok;	E8B8	
ß	--	&szlig;	00DF	
Ÿ	&Ver;	&Vslstrok;	2123	
ÿ	&ver;	&vdiagstrok;	E8BC (v2)	
	&Vos;			
	&vos;			

## Marques de ponctuation et de mise en page

Glyphe	CHR	MUFI	CODE	BFM + Commentaire
·	.	&middot;	00B7	ASCII 183
¿	&punct;	&punctelev;	F161	
/		&punctelevdiag;	F1F0	Est-ce le même que ! ?
¿		&punctinter;	F160	
/	,	&sol;	002F	
/		&virgmin;	F1F7 (v2)	Cf. 05F3
∴		&tridotsdownw;	F1EE	
⋄		&lozengedot;	2058	
⚡		*&dbloblhyphen;	2E17	Absent MUFI Code 03DF et police Cardo utilisés temporairement
¶		&para;	00B6	pied de mouche... faute de mieux
ƒ		&parag;	F1E1	crochet adlinéaire
^		&logand;	2227	marque d'insertion

## 8. Projets cités :

Charrette : <http://www.princeton.edu/~lancelot>

Christine de Pizan : <http://www.pizan.lib.ed.ac.uk/>

Khartês (Université de Liège) : non accessible en ligne

Menota : <http://www.menota.org>

MUFI : <http://gandalf.aksis.uib.no/mufi/>

## 9. Index des éléments (entête)

author, 8  
 catDesc, 11  
 category, 11  
 catRef, 11  
 change, 11  
 classDecl, 11  
 creation, 10  
 date, 10  
 editor, 8  
 encodingDesc, 11  
 extent, 8  
 fileDesc, 7  
 handDesc, 9  
 handNote, 9  
 history, 9  
 langUsage, 10  
 msDesc, 9  
 msIdentifier, 9  
 msPart, 9  
 name, 8  
 origin, 10  
 physDesc, 9  
 profileDesc, 10  
 publicationStmt, 9  
 resp, 8  
 respStmt, 8  
 revisionDesc, 11  
 sourceDesc, 9  
 taxonomy, 11  
 teiHeader, 7  
 textClass, 11  
 textDesc, 10  
 title, 8  
 titleStmt, 8

## 10. Index des éléments (hors entête)

ab, 16  
 abbr, 20  
 add, 15, 17  
 am, 17  
 bfm:headlb, 18  
 bfm:hyphen, 19  
 bfm:lettrine, 17  
 bfm:mdvAbbr, 17, 20  
 bfm:punct, 21  
 bfm:sb, 21  
 cb, 21  
 choice, 22  
 corr, 15, 22  
 damage, 22  
 del, 15, 22  
 div, 23  
 ex, 23  
 gap, 24  
 head, 24  
 hi, 24  
 l, 16  
 lb, 24  
 lg, 16  
 me:dipl, 13, 21, 25, 28  
 me:facs, 14, 18, 19, 20, 21, 25, 28  
 me:norm, 13, 21, 25, 28  
 me:pal, 14, 25  
 note, 16, 25  
 p, 26, 28  
 pb, 26  
 q, 26  
 seg, 27  
 space, 27  
 subst, 27  
 supplied, 28  
 unclear, 15, 28  
 w, 15, 28