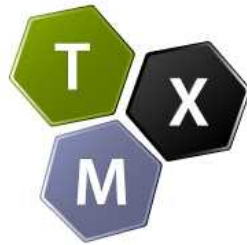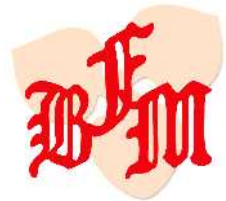# Problems in Linguistic Annotation and Analysis of XML TEI Editions of Textual Heritage Works

ICAR Research laboratory
Université de Lyon, ENS de Lyon, CNRS

El'Manuscript 2012
Petrozavodsk, 3 – 9 September 2012

# Outline

- Reasons for using TEI in digital editions and in corpus building

- Possible conflicts between philological markup and linguistic annotation

- Analyzing and enriching critical editions with TXM platform
  - TEI-BFM customization
  - tokenization and applying NLP tools
  - TEI-TXM extension
  - corpus query and analysis

# Reasons for using TEI XML

- Detailed documentation
  - elaborate text structure and markup theory
- Wide community (TEI < XML)
- Constant evolution
  - relatively long backward compatibility
  - tools for markup upgrade
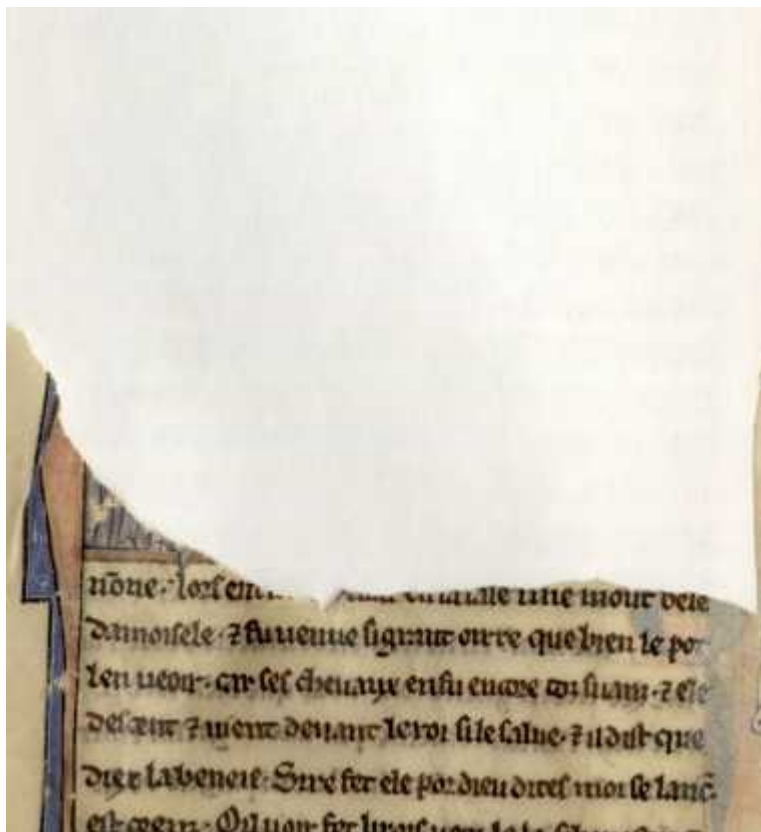- Opportunities for interchange and using external tools

# Possible conflicts

- XML does not allow overlapping structures:
  - elements must nest
- Only one hierarchy can be marked up directly
- "Competing" hierarchies can be marked up using
  - joining mechanisms
  - "milestones" (empty tags)
  - stand-off mechanisms
- Sophisticated tools may be necessary to re-build elements of alternative hierarchies
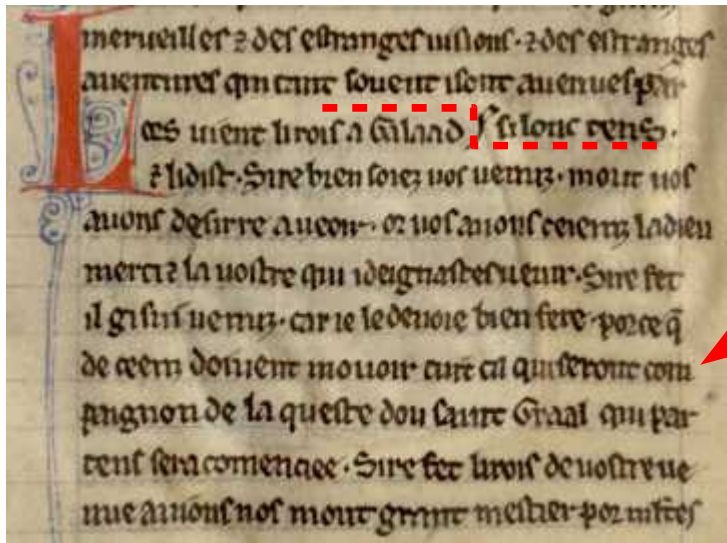
# Multiple hierarchies

- **Physical (geometrical)**
  - ☐ book > folio > page > column>  line...
- **Logical (semantic)**
  - ☐ work > part > paragraph / verse group...
- **Rythmic**
  - ☐ verse...
- **Linguistic**
  - ☐ sentences > phrases > words > morphemes...
- **Genetic**
  - ☐ additions, deletions...

# Examples of competing hierarchies
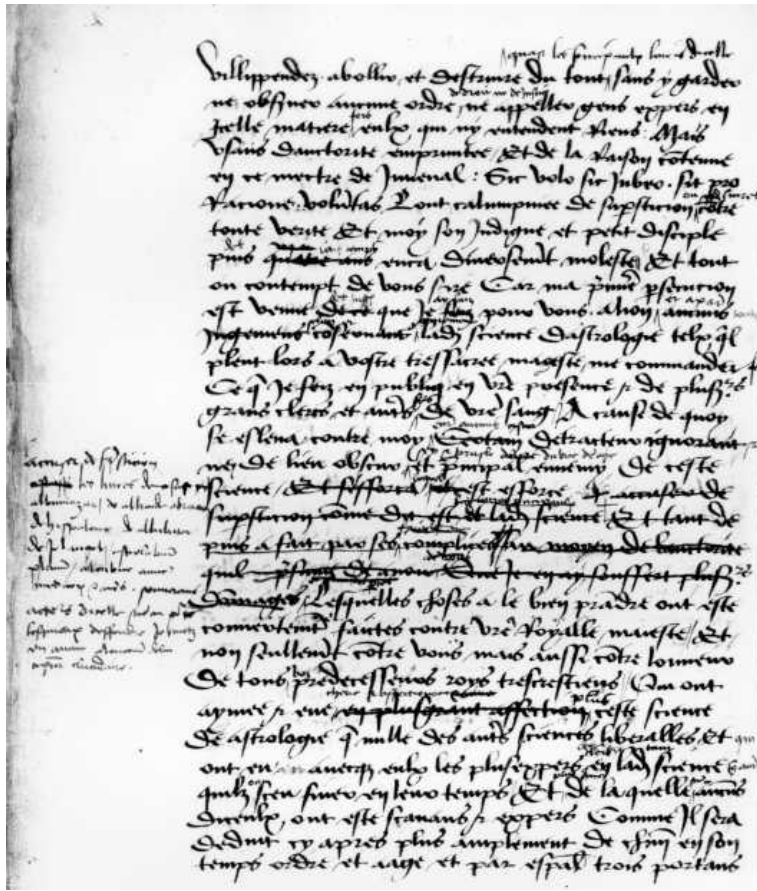


- Linguistic vs. editorial structures
  - `<supplied>`A la veille de la Pentecoste […] tables a heure de`</supplied>` nonne lors en`<supplied>`tra a cheval`</supplied>` `<unclear>`en la s`</unclear>`ale

# Compatibility problems



- Linguistic vs. physical structures
  - "overlapping" lines in the end of divisions
  - line/pagebreaks "inside" words
  - "irregular" word segmentation
    - lirois, idegnastesuenir
    - es paule, en mena

# Manuscript transcriptions for linguistic research



- If a manuscript contains additions, deletions, corrections, etc. it may be interesting to analyse different states of the text (genetic edition)

Ms. BnF fr. 1357, f° 2v, Simon de Phares, *Recueil des astrologues*

# Processing texts with NLP tools

- Lexical items recognition (tokenisation)
  - w
- Sentence splitting
  - s
- POS annotation and lemmatization
  - e.g. TreeTagger
  - export / import
- Inline / standoff annotation

# Compatibility with NLP tools

- Tokenizers usually rely on "separator" characters
  - xml editing tools can easily reformat text (insert line breaks and tabs) to make the xml structure more visible ==> spacing characters around xml tags are not reliable
- Sentence parsers rely on "strong" punctuation
  - they can also use structural tags to force sentence breaks (e.g. end of heading)
  - it may be tricky to place s tags correctly in a text with a rich editorial markup

# Compatibility with NLP tools

- Rich manuscript transcriptions are likely to be regularly (constantly?) updated
- Standoff markup solutions require stable sources
  - character position for simple text alignment
  - identifier for xml alignment

# Corpus query tools

- **IMS Open corpus workbench (CWB)**
  - ☐ CQL query language ← index tables
- **TXM**
  - ☐ TEI (extended)
  - ☐ CQP concordances and indexes
  - ☐ statistical analysis tools
- **TigerSearch**
  - ☐ querying syntactic structures ← index tables

# BFM TEI customization

- Define tag classes with respect to their position in linguistic hierarchy (project-specific)
  - word-internal (c, g, am, ex...)
  - = word (num, abbr)
  - one or more words within a sentence
  - at least one sentence (div, p, head...)
  - out-of-text for NLP tools (note, foreign...)
  - floating (most editorial tags)
    - but may have a "privileged" position
    - additional markup for "unusual" positions (see next slide)

# BFM TEI customization

- ■ Deal with the floating tags
    - ■ split, if a tag starts in the middle of a word and ends after a few other words
        - □ en\<supplied\>tra\</supplied\> \<supplied\>a cheval\</supplied\>
    - ■ use special attribute values if a tag is not in its default position (e.g. "word_part" and "multi_s")
        - □ en\<supplied txm:range="word_part"\>tra\</supplied\>
        - □ cf. lb/@type
            - ▪ The type attribute may be used to characterize the line break in any respect, but its most common use is to specify that the presence of the line break does not imply the end of the word in which it is embedded. A value such as inWord or nobreak is recommended for this purpose, but encoders are free to choose whichever values are appropriate. (TEI Guidelines)
    - ■ or pre-tag words containing floating tags
        - □ \<w\>en\<supplied\>tra\</supplied\>\</w\> \<supplied\>a cheval\</supplied\>

# TEI TXM Extension

- **Systematic tagging of words and sentences**
  - `<s>` / `<w>` with `@xml:id`
- **Inline markup**
  - `<txm:form>` +
  - `<txm:ana>` *
- **Standoff**
  - referring to w/`@xml:id`
  - updating
    - "inject" annotation in-line before updates (if possible)
    - use diff algorithms to re-align text and markup

# Generalization

- Check the tag usage in the document
- Use an XSLT filter
  - □ before the tokenization
  - □ maybe also after the tokenization
- Remove xml elements unlikely to used in corpus query
- Deal with typical problems
  - □ e.g. full stop in the end of a sentence-internal element

# Generalization

- Tested on
  - Bouvard et Pécuchet digital edition
    - http://dossiers-flaubert.ish-lyon.cnrs.fr/
  - Bibliothèques virtuelles des humanistes
    - http://www.bvh.univ-tours.fr/
  - Frantext
    - http://www.frantext.fr

# *Queste del saint Graal* digital edition

- XML markup (tag usage)

| Tag | Occurs |
|---|---:|
| w | 118 885 |
| lb | 10 751 |
| me:facs, me:dipl; me:norm, choice, s | 2000 - 3500 |
| q | 622 |
| milestone, orig, pb, seg, p | 200 - 350 |
| corr, sic, ex, bfm:mdvAbbr | 150 - 199 |
| note, bfm:hyphen, supplied | 10 - 50 |
| del, hi, bfm:sb, bfm:lettrine | 5 - 8 |
| text, subst, damage | 1 |

# *Queste del saint Graal* digital edition

■ Interface http://txm.bfm-corpus.org/txm/

# *Queste del saint Graal* digital edition

■ Interface http://txm.bfm-corpus.org/txm/

# *Queste del saint Graal* digital edition

- Interface http://txm.bfm-corpus.org/txm/

# Thank you!