Analyse100010010100010011000110101010100011
010011et01010011000111001101010010101011
1Traitement0101000110001010110011001
0100Informatique0101001011001001001100100
de0101la0100011101010001
0101Langue010111001
Française01010011
0101010Analyse
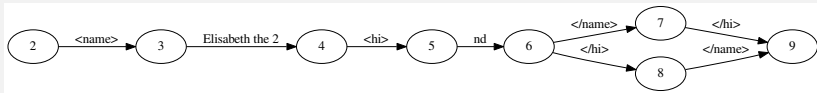100110001101010

# A Tool for adding word Annotations into TEI Files

## Bertrand Gaiffe

# A small example

\<**p**\>The queen Elisabeth the 2\<**hi** rend="sup"\>nd\</**hi**\>
celebrates her 93\<**hi** rend="sup"\>rd\</**hi**\> birthday.\</**p**\>

red :  \<name xmlns="http://www.tei-c.org/ns/1.0"\>

blue : \<num xmlns="http://www.tei-c.org/ns/1.0"\>

An extract of the graph:



Parsing with the grammar of the xml language, we get:

<**p**>The queen <name>Elisabeth the
2<**hi** rend="sup">nd</**hi**></name> celebrates her
<num>93<**hi** rend="sup">rd</**hi**></num> birthday.</**p**>

**Mixer**

The program (Mixer) takes as arguments:
- an XML file
- a "companion file" of the form:

| index1 | index2 | XML elements |
|--------|--------|--------------|
| 10 | 27 | <name xmlns="http://www.tei-c.org/ns/1.0"> |
| 43 | 47 | <num xmlns="http://www.tei-c.org/ns/1.0"> |

It proved reliable enough to POS-tags a corpus of about 5000 french novels.

**Annotations alignment**

The annotation data in practice:

```
The        O
queen      O
Elisabeth  B-<name xmlns="http://www.tei-c.org/ns/1.0">
the        I-<name>
2          I-<name>
nd         I-<name>
celebrates O
her        O
93         B-<num xmlns="http://www.tei-c.org/ns/1.0">
rd         I-<num>
birthday   O
```
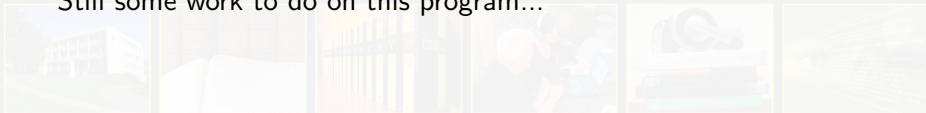
- Takes as input an XML file and a tabular file and adds two columns corresponding to the indexes of the companion file.
- Of course, the text contents of the XML file has to match the text contents of the tabular file...

Still some work to do on this program...

# The whole process

- ▶ fold the non annotated parts (ex <teiHeader xml:id="something"/>)
- ▶ extract the text and feed some annotation process (pos, ner...). If the result is not tabular, transform it into a tabular file.
- ▶ align the tabular result with the "folded xml"
- ▶ produce the "companion file"
- ▶ apply mixer
- ▶ unfold the xml result

# Conclusion and perspectives

- ▶ we do not modify the text contents of the XML file
- ▶ no schema so far (so the result is not necesserally TEI) ($\longrightarrow$ to do, but...
    - ▶ some more work on corrections...
    - ▶ it may happend that no valid result
  )
- ▶ so far, gives "a" result somewhat randomly (could output the intersection grammar...)
- ▶ the alignment process should be improved...

to test:

https://github.com/bgaiffe/Annotations