



HAL
open science

TreeCloud & Unitex: an increased synergy

Claude Martineau

► **To cite this version:**

Claude Martineau. TreeCloud & Unitex: an increased synergy. ECLAVIT Workshop, Nov 2017, Champs-sur-Marne, France. 2017. hal-01702091

HAL Id: hal-01702091

<https://hal.science/hal-01702091v1>

Submitted on 19 Apr 2019

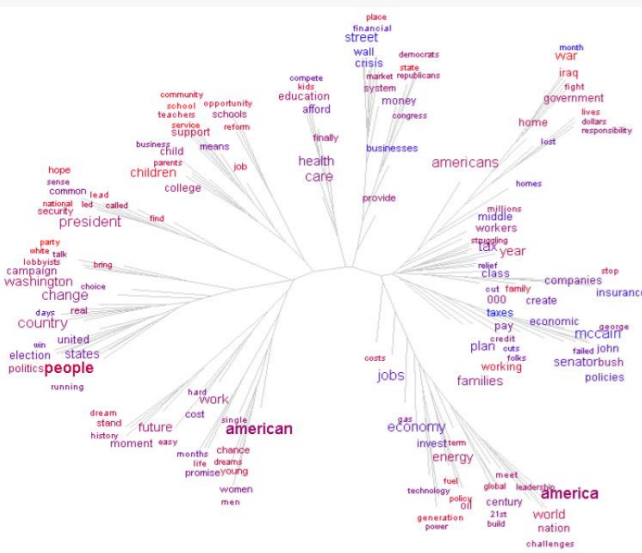
HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

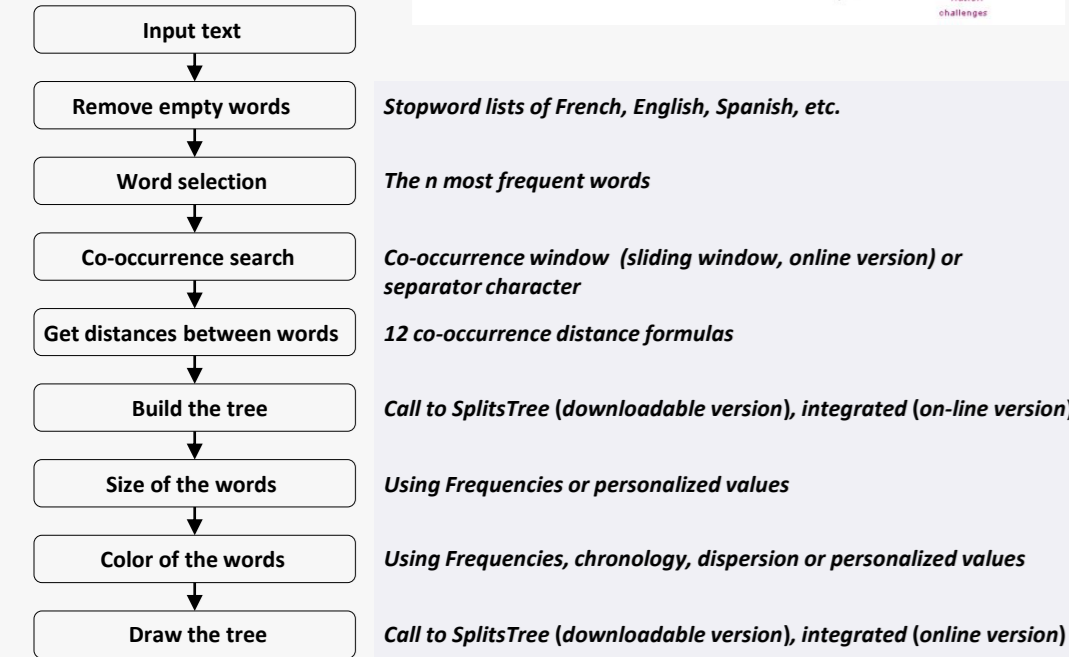
TreeCloud is a tree cloud visualization of a text

TreeCloud builds a tree cloud visualization of a text, which looks like a tag cloud where the tags are displayed around a tree to reflect the co-occurrence distance between the words in the text.

The adjacent treecloud gives an overview of Barack Obama's 2008 presidential campaign speeches.



Tree construction process



The stopword list contains grammatical words and auxiliary or modal verbs.

The stopwords are usually the most common words in a language. To build a meaningful tree, all these words must be firstly removed from the input text.

Whenever a sliding window of analysis is used, its size (i.e. a number of words) must be given as parameter.

Co-occurrence and distance calculation

- Given two words A and B and:
• O11, observed number of sliding windows containing both A and B
• O12, observed number of sliding windows containing A but not B
• O21, observed number of sliding windows not containing A but B
• O22, observed number of sliding windows containing neither A nor B

- the following variables are defined:
• R1 = O11 + O12, number of sliding windows containing A
• R2 = O21 + O22, number of sliding windows not containing A
• C1 = O11 + O21, number of sliding windows containing B
• C2 = O12 + O22, number of sliding windows not containing B
• N = R1 + R2 = C1 + C2, number of sliding windows
• E11 = (R1C1/N), expected number of sliding windows containing both A and B
• E12 = (R1C2/N), expected number of sliding windows containing A but not B
• E21 = (R2C1/N), expected number of sliding windows not containing A but B
• E22 = (R2C2/N), expected number of sliding windows containing neither A nor B

- The definitions of co-occurrence formulas are the following:
• jaccard: 1 - (O11 / (O11 + O12 + O21))
• liddell: 1 - (O11 * O22 - O12 * O21) / (C1 * C2)
• dice: 1 - 2 * O11 / (R1 + C1)
• hyperlex: 1 - max(O11 / R1, O11 / C1)
• poissonstirling: O11 (log O11 - log E11 - 1)
• chisquared: 1000 - N(O11 - E11)^2 / (E11 * E22)
• zscore: 1 - (O11 - E11) / sqrt(E11)
• ms: 1 - min(O11 / R1, O11 / C1)
• oddsratio: 1 - log((O11 * O22) / (O12 * O21))
• loglikelihood: 1 - 2(O11 * log(O11 / E11) + O12 * log(O12 / E12) + O21 * log(O21 / E21) + O22 * log(O22 / E22))
• gmean: 1 - O11 * sqrt(R1 * C1) = 1 - O11 * sqrt(N * E11)
• mi (mutual information): 1 - log(O11 / E11)
• ngd (normalized Google distance): (max(log R1, log C1) - log O11) / (N - min(log R1, log C1))

Several versions of TreeCloud

Downloadable version in Python
• 2009: TreeCloud 1.3 for Windows, Linux, Mac developed by Philippe Gambette
Use of SplitsTree 4.10 to draw the tree

On-line version in C

- 2009: 1st C version developed by Jean-Charles Bontemps
• 2012: Transition to Unicode developed by Claude Martineau
• 2014: 1st implementation of Unitex developed by Claude Martineau

Unitex/GramLab is a corpus analyser and annotation tool

- Based on Automata and RTNs with outputs
• Multilingual: Up to 22 languages (French, English,..., Greek, ... , Korean, Thai)
• Unicode 3.0 (UTF8, UTF16LE, UTF16BE)
• Cross-platform: Linux, macOS, Windows
• Open source: https://github.com/UnitexGramLab
• Website and binary installers: http://unitexgramlab.org
• Under development since 2001 by a group of passionate volunteers

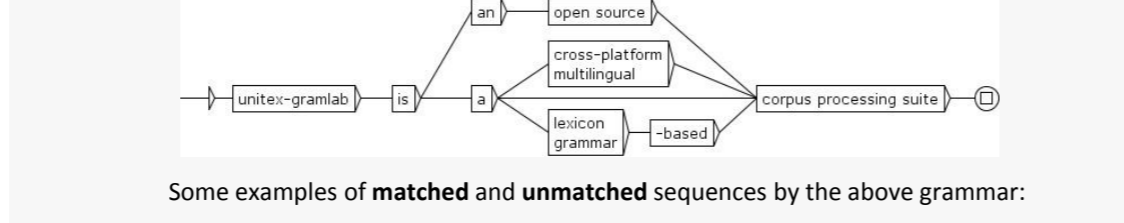
Unitex/GramLab uses linguistic resources:

• DELA (LADL electronic dictionaries)
A typical DELA entry is composed by a simple or compound inflected form, followed by a lemma and grammatical information. Each entry can be associated with syntactic and semantic attributes and inflection rules:

inflected_form, lemma.grammatical_information+attributes:inflection_rule
Example: Given the French simple word "avocat" (lawyer) and the compound word "avocat d'affaires" (business lawyer), a DELA representation would be:

avocat,avocat.N+Hum+Prof:ms avocat d'affaires,avocat d'affaires.N+Hum+Prof:ms
avocate,avocat.N+Hum+Prof:fs avocate d'affaires,avocate d'affaires.N+Hum+Prof:fs
avocats,avocat.N+Hum+Prof:mp avocats d'affaires,avocats d'affaires.N+Hum+Prof:mp
avocates,avocat.N+Hum+Prof:fp avocates d'affaires,avocate d'affaires.N+Hum+Prof:fp

- Syntactic or semantic rules called «local grammars» represented by graphs
• Graphical representations of local grammars are composed by a set of linked boxes.
• A successful path is a path between initial and final states.



Some examples of matched and unmatched sequences by the above grammar:
Unitex-GramLab is a corpus processing suite [MATCH]
Unitex-GramLab is an open source corpus processing suite [MATCH]
Unitex-GramLab is a hard to learn corpus processing suite [FAIL]
Unitex-GramLab is [FAIL]

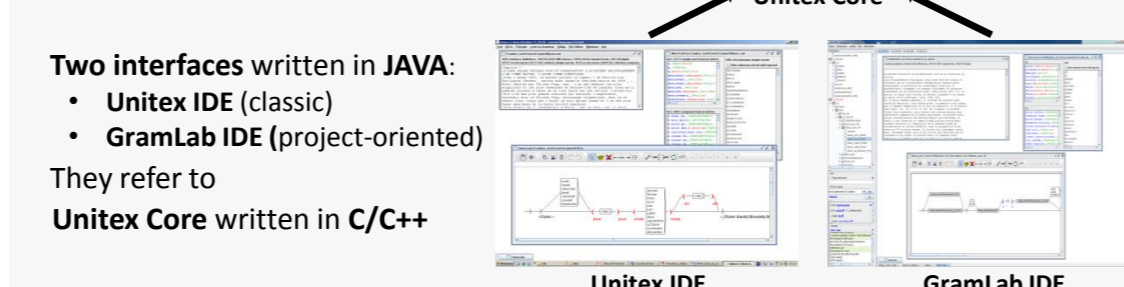
An example of analysis

Application of a dictionary; the result is the text dictionary, then application of a local grammar
The grammar below contains two search paths:
• an adverb (<ADV>) ending in -ly followed by a past participle (<V:K>)
• a noun (<N>) followed by a verb in progressive form (<be.V> <V:G>)
A lexical mask like <V:K> refers to the text dictionary.
The recognized sequences are surrounded by the tag <pattern>. The results are represented in the form of concordances.

Word Lists in ChildrenPublic/Documents/...
Unitex-GramLab is a corpus processing suite [MATCH]
Unitex-GramLab is an open source corpus processing suite [MATCH]
Unitex-GramLab is a hard to learn corpus processing suite [FAIL]
Unitex-GramLab is [FAIL]

Notice that Unitex/GramLab can also produce an annotated text that includes all the tagged patterns.

Several ways to use Unitex/GramLab



Use the API C and JAVA (JNI) that provides access to
• a virtual file system
• a persistence layer for resources (alphabets, dictionaries and corpora)

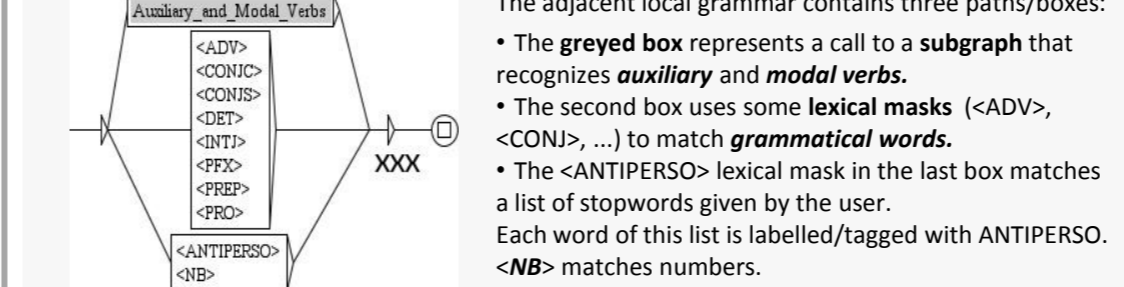
Command lines or system calls with Perl, Python, etc.

How and Why to plug Unitex into TreeCloud?

Construction of the tree with Unitex

- 1) Unitex transforms the input text into a new text with all the forbidden/stopwords replaced by the XXX « word »
2) The new text is sent to TreeCloud with XXX as the unique forbidden word (the unique word in the stop list)

Get a larger and more accurate coverage of forbidden words



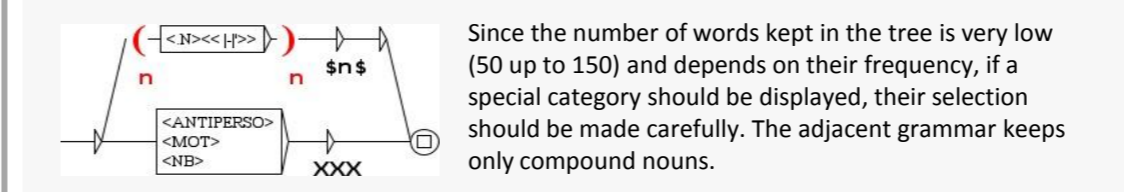
Insert multiwords into the tree

If a dictionary contains compound words, these words can be kept in the tree but all multiwords cannot be listed in a dictionary.

The adjacent grammar contains a path that recognizes person names. The variable \$P\$ contains the name of a person captured by the subgraph Name_of_Person.

The .{Pers} label is added to the output (2017 online version).

Make a strong selection of the words kept in the tree



In order to display verbs into the tree, it could be useful to get their LEMMA (since they have many inflected forms). The adjacent grammar is designed for this goal. A similar process is used for languages with nouns or adjectives with cases (e.g. Serbian or Greek).

2017 on-line version of TreeCloud: an improved implementation

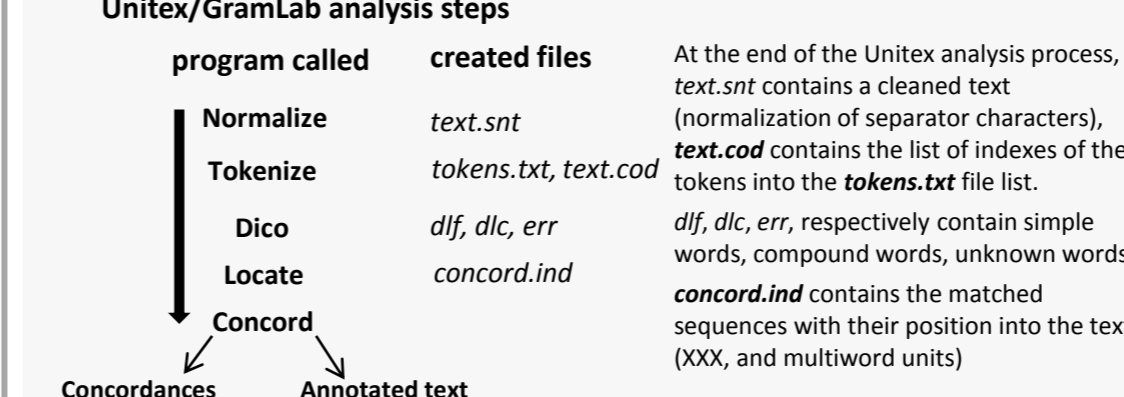
Introduction of the concept of file processing
In the online version of 2014, there is only a single Unitex processing for each language. The necessary resources were hard-coded into the program.

In the new 2017 version, there can be several processing procedures for each language. Furthermore, the Serbian language (Latin and Cyrillic) has been added. In order to manage these pairs (language, processing) a concept of processing file has been set up.

For example, in the processing file below, the first line indicates the path for French resources, then three dictionaries (French general dictionary, first name dictionary, toponym dictionary) have to be applied to the text. In the end, the local grammar is applied in the replace mode.

```
REP:TreeCloud_WS/French/src
NB_DICOS=3
FICHER=Dela/dela-fr-public.bin
FICHER=Dela/prenom-s.bin
FICHER=Dela/Prolex-Unitex_1_2_TOPONYMES.bin
GRAMMAIRE:
FICHER=Graphs/Treecloud_N_Pers_Top_v1_FR.fst2
MODE=REPLACE
```

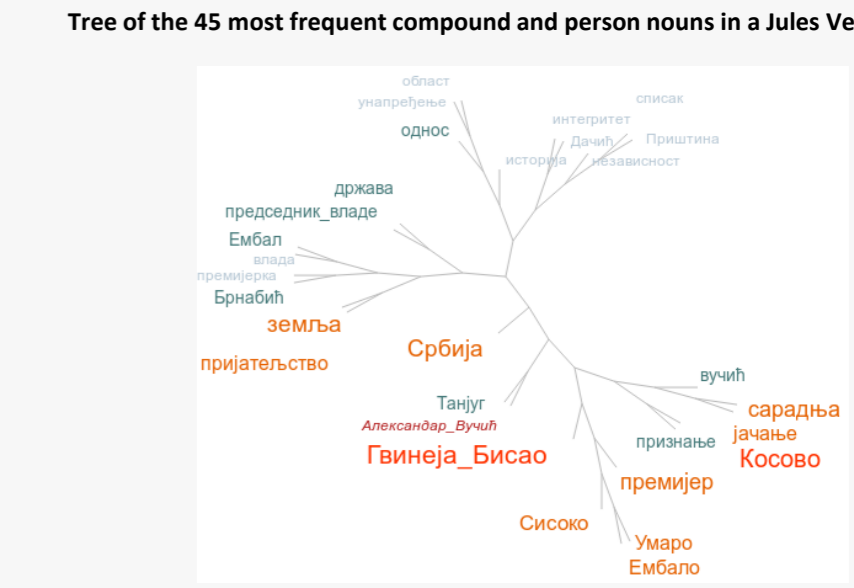
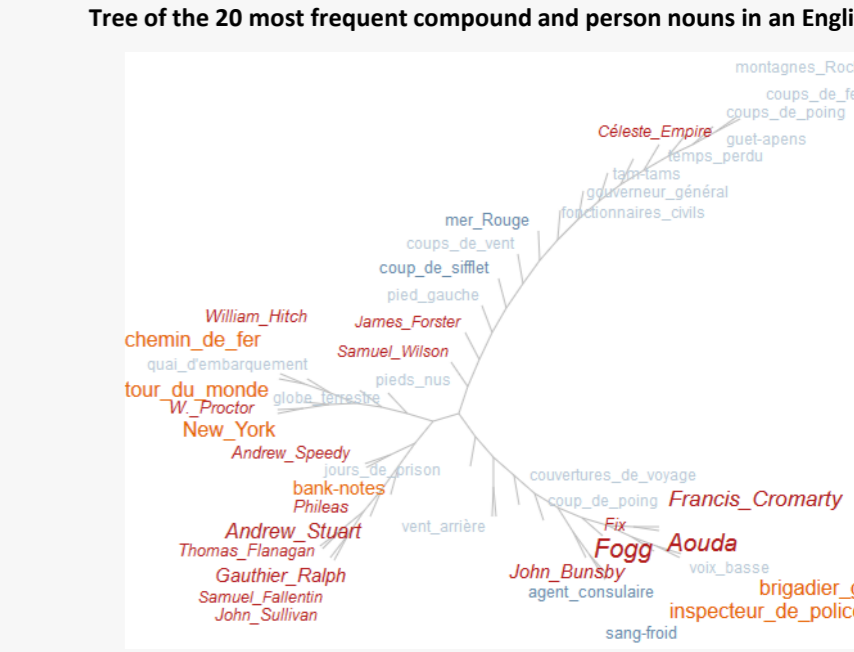
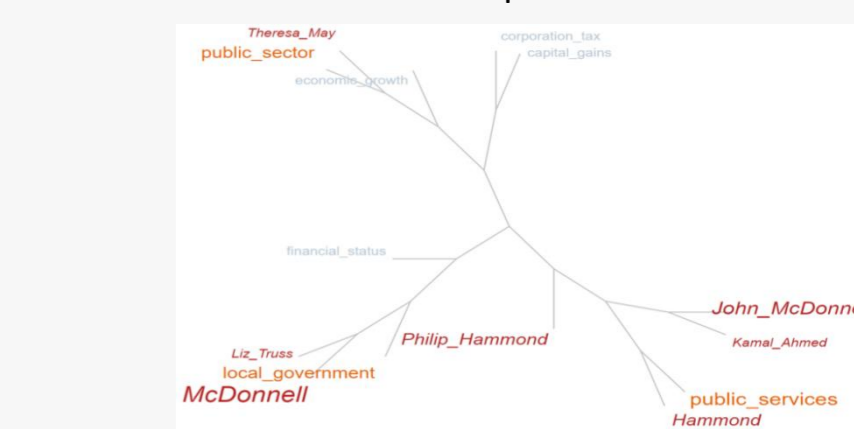
Take advantage of the work already done by Unitex



To get the «new text», we retokenize the text with matched sequences of the concord.ind file as the new tokens of the text. New token.txt and text.cod files are created. This process prevents double reading of the text and double division into words.

Thanks to the Unitex API and virtual file system, all this work is done in memory.

Some examples of trees in different languages



Conclusion

Plug-in of Unitex into TreeCloud provides:

- A more accurate representation of forbidden words
• All kinds of multiwords to be recognized in the text and presented in the tree
• A visual representation of some grammatical or semantic categories of the words.

• A faster construction of the tree (via a careful use of the Unitex API)

http://treecloud.univ-mlv.fr