

# New features in Corpus Architect (corpus management)

Vít Suchomel



`vit.suchomel@sketchengine.co.uk`

4<sup>th</sup> Sketch Engine Workshop  
Tallinn, October 16, 2013

# Nicer user interface

- Creation of user corpora simplified – automated steps:
  - The most complete processing chain available
  - The most up-to-date sketch grammar

## Create new corpus

The screenshot shows a web form titled "Create new corpus". It contains the following elements:

- Corpus name:** A text input field with a cursor at the start.
- Language:** A dropdown menu showing "-----" with a downward arrow.
- Configuration template:** A section with two sub-sections:
  - Preloaded configuration templates:** A box containing a radio button labeled "Default for selected language".
  - My configuration templates:** A box containing two radio buttons:
    - The first is labeled "Farsi word-POS-lemma (Persian)" and is accompanied by a visual representation of a sketch grammar: **V** (red), **T** (red), **L** (yellow), **S** (red).
    - The second is labeled "my Basque test template (Basque)" and is accompanied by the same visual representation: **V** (red), **T** (red), **L** (yellow), **S** (red).

# Nicer user interface

- Better display of corpora and file names
  - The original file and corpus names displayed
  - No more underscores to replace non latin characters

**ნარმატებული 14 წელი**  
carmatebuli\_14\_celi

 [Add new file](#) /  [Add data from web using WebBootCaT](#) /  [Compile corpus](#) / 

#	Original file	Plain text	Vertical	Tokens 
1	<a href="#">ჰელაინ ინტერნეიშნლ</a>			256

# Keywords, terms

- Term extraction for English, French, Chinese, Japanese, Korean, Spanish, Russian
  - Corpus processing chain reviewed and separated from CA
  - New/reviewed term grammars applied to reference corpora
- Keyword and term extraction results merged
  - One page, keywords + multi word terms
- API in development
  - input text, output keywords

# Keywords & terms extracted from a domain specific corpus

## Keywords

- |  |   |
|--|---|
| <input type="checkbox"/> dioxide (415.2, <a href="#">427</a> )       | <input type="checkbox"/> mutualism (75.6, <a href="#">8</a> )     |
| <input type="checkbox"/> trophic (264.9, <a href="#">33</a> )        | <input type="checkbox"/> radiative (75.0, <a href="#">12</a> )    |
| <input type="checkbox"/> greenhouse (238.4, <a href="#">282</a> )    | <input type="checkbox"/> gasses (75.0, <a href="#">12</a> )       |
| <input type="checkbox"/> ecology (237.7, <a href="#">196</a> )       | <input type="checkbox"/> lca (74.4, <a href="#">10</a> )          |
| <input type="checkbox"/> methane (233.5, <a href="#">108</a> )       | <input type="checkbox"/> biotic (74.2, <a href="#">10</a> )       |
| <input type="checkbox"/> arrhenius (232.2, <a href="#">25</a> )      | <input type="checkbox"/> acidification (74.1, <a href="#">9</a> ) |
| <input type="checkbox"/> photosynthesis (230.6, <a href="#">46</a> ) | <input type="checkbox"/> above-ground (73.6, <a href="#">9</a> )  |
| <input type="checkbox"/> callendar (215.4, <a href="#">22</a> )      | <input type="checkbox"/> holism (73.5, <a href="#">9</a> )        |
| <input type="checkbox"/> ecosystems (211.4, <a href="#">114</a> )    | <input type="checkbox"/> felzer (73.5, <a href="#">7</a> )        |
| <input type="checkbox"/> warming (193.8, <a href="#">504</a> )       | <input type="checkbox"/> carbonic (72.4, <a href="#">9</a> )      |
| <input type="checkbox"/> keeling (192.5, <a href="#">23</a> )        | <input type="checkbox"/> loa (71.5, <a href="#">10</a> )          |
| <input type="checkbox"/> carbon (186.8, <a href="#">558</a> )        | <input type="checkbox"/> biogeography (71.2, <a href="#">9</a> )  |
| <input type="checkbox"/> n't (177.1, <a href="#">17</a> )            | <input type="checkbox"/> organisms (70.4, <a href="#">86</a> )    |
| <input type="checkbox"/> gases (173.9, <a href="#">159</a> )         | <input type="checkbox"/> mauna (69.7, <a href="#">10</a> )        |
| <input type="checkbox"/> -oct- (169.3, <a href="#">28</a> )          | <input type="checkbox"/> flowering (68.4, <a href="#">23</a> )    |
| <input type="checkbox"/> vapor (151.3, <a href="#">72</a> )          | <input type="checkbox"/> emitted (68.2, <a href="#">27</a> )      |

## Terms

- ☐ carbon dioxide (567.1)
- ☐ greenhouse effect (515.0)
- ☐ water vapor (486.8)
- ☐ global warming (298.8)
- ☐ industrial ecology (261.6)
- ☐ infrared radiation (170.9)
- ☐ carbon cycle (169.0)
- ☐ surface temperature (161.0)
- ☐ elevated carbon (156.4)
- ☐ elevated carbon dioxide (156.4)
- ☐ greenhouse gas (135.8)
- ☐ climate system (134.1)
- ☐ food web (124.3)
- ☐ amount of carbon dioxide (116.8)
- ☐ other greenhouse (114.2)
- ☐ global temperature (109.1)

## Configuration templates in Corpus Architect administration

Name	Default grammar	Reference corpus	Terms grammar	Terms corpus
<b>TreeTagger for English</b>	english-penn_tt-2.5.wsdef.txt	enTenTen08 (trial)	english-penn_tt-terms-2.3.wsdef.m4	enTenTen12 [sample 40M] with term grammar (trial)
<b>TreeTagger for French</b>	french-tt-1.0.wsdef.txt	frTenTen12 (trial)	french-tt-terms-1.0.wsdef.m4	frTenTen12 [sample] with term grammar (trial)
<b>Stanford Chinese Segmenter and Tagger</b>	chinese-universal-with-tags.1.0.wsdef.m4	zhTenTen11 (trial)	chinese-stanford-terms-1.0.wsdef.m4	zhTenTen11 [sample 10M] with term grammar (trial)
<b>Mecab (+Unidic2 +LUW) for Japanese</b>	japanese-mecab-unidic2-1.7.wsdef.txt	jpTenTen11 [LUW, sample] (trial)	japanese-mecab_unidic2_luw-terms-1.0.wsdef.m4	jpTenTen11 [LUW, sample] with term grammar (trial)
<b>HanNanum for Korean</b>	korean-universal-with-tags.1.0.wsdef.m4	koTenTen12 (trial)	korean-hannanum_simplified-terms-1.0.wsdef.m4	koTenTen12 [sample] with term grammar (trial)
<b>TreeTagger for Russian</b>	russian-tt-1.0.wsdef.txt	ruTenTen11 [sample 50M] (main)	russian-tt-terms-1.0.wsdef.m4	ruTenTen11 [sample 50M] with term grammar (restricted)
<b>Freeling for Spanish</b>	spanish-freeling-1.0.wsdef.txt	esTenTen11 [Eu + Am, Freeling, sample 50M] (trial)	spanish-freeling-terms-1.0.wsdef.m4	esTenTen11 [Eu + Am, Freeling, sample 50M] with term grammar

# Single sign on

- Service provider software set up
  - Use case: University students are granted access based on authentication at their university, no additional authentication with Sketch Engine required
- Ready in UK academic identity federation

# WebBootCaT search engine

- New word search service in MS Azure Marketplace used
- Service availability improvements



# Future work

- improve usability
- finish term extraction
- finish corpus comparison
- implement corpus homogeneity

# Conclusion

Nicer interface, parallel corpora support, keyword & term extraction,...  $\Rightarrow$  much work done!