

# **Comparable Corpora BootCaT (CCBC)**

**Adam Kilgarriff, Avinesh PVS,  
Jan Pomikalek**

Lexical Computing Ltd.



# Just-in-time corpora

- Krista Varantola
- Translators, terminologists
- *In-domain terminology:*
  - Domain dictionaries
    - Don't exist
    - Out of date
    - Not accessible
- Collect in-domain web pages
- Instant corpus



# BootCaT (Bootstrapping Corpora and Terms)

- Baroni and Bernardini 2004
- User: input ‘seed terms’
- Send 3-at-a-time to a search engine
  - Returns search hits page
- Retrieve those pages
- A corpus!
  - Cleaning, deduplicating, linguistic processing
- Extract terms
  - Can use extracted terms as seeds, iterate



# Works well

- Widely used
- More implementations
  - SkE has WebBootCaT, web front end
- Secret:
  - *piggybacks on search engines*
  - They do the donkey-work
    - on-domain, text-rich pages, no spam, ...



## Also in use for

- General language corpus
  - Long list of general seed words
    - Pioneer: Serge Sharoff
    - LCL: Corpus Factory
- ‘Varieties of Learner English’
  - General English, same queries *except*
    - Region=UK, US, Canada, Aus, China, Japan, Korea
  - Validation under way



# The Sketch Engine

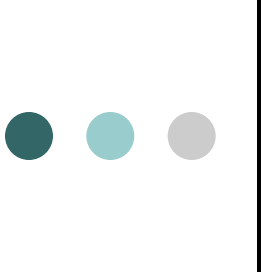
- Corpus query tool, since 2003
- Widely used by lexicographers
  - Commercial
    - OUP, CUP, Collins, Macmillan, Le Robert, Cornelsen, Shogakukan
  - National dictionary projects
    - Bulgaria, Czech Republic, Estonia, Netherlands, Slovakia, Slovenia
- Universities
  - Linguistics, language research, NLP, language teaching, teaching translation



# 55 languages and counting

Large corpora ready-to-use for

*Arabic Bengali Bulgarian Chinese Czech  
Croatian Danish Dutch English Estonian Finnish  
French German Greek Gujarati Hebrew Hindi  
Indonesian Irish Italian Japanese Korean Latin  
Malay Malayalam Norwegian Persian Polish  
Portuguese Romanian Russian Serbian  
Setswana Slovak Slovene Spanish Swahili  
Swedish Tamil Telugu Thai Turkish Urdu  
Vietnamese*

- 
- Handles large corpora
    - Largest to date: 8 billion words
  - Fast
  - Web-based: no software to install
  - Build ‘instant corpora’ from the web
  - Load your own corpus
    - Quota of space on SkE server
  - ***Word sketches***
    - One-page, automatic accounts of a word’s grammatical and collocational behaviour
  - Free 30-day trial: [sketchengine.co.uk](http://sketchengine.co.uk)



# goal *(noun)* ukWaC freq = 168345

<u>object of</u>	<u>58924</u>	<u>3.0</u>
score	<a href="#">8390</a>	11.42
achieve	<a href="#">9422</a>	10.05
concede	<a href="#">1421</a>	9.43
accomplish	<a href="#">585</a>	8.05
reach	<a href="#">1924</a>	7.84
pursue	<a href="#">648</a>	7.5
net	<a href="#">337</a>	7.43
set	<a href="#">2413</a>	7.42
attain	<a href="#">400</a>	7.42
grab	<a href="#">406</a>	7.39
pull	<a href="#">501</a>	7.11
disallow	<a href="#">190</a>	6.67
bag	<a href="#">186</a>	6.64
meet	<a href="#">1335</a>	6.62

<u>and/or</u>	<u>16213</u>	<u>0.9</u>
objective	<a href="#">858</a>	7.01
aspiration	<a href="#">159</a>	6.73
ambition	<a href="#">151</a>	6.42
appearance	<a href="#">216</a>	5.67
penalty	<a href="#">102</a>	5.32
target	<a href="#">320</a>	5.3
goal	<a href="#">315</a>	5.3
dream	<a href="#">129</a>	5.27
motivation	<a href="#">67</a>	5.22
aim	<a href="#">227</a>	5.15
try	<a href="#">34</a>	5.13
vision	<a href="#">154</a>	5.1
ideal	<a href="#">52</a>	5.1
expectation	<a href="#">98</a>	5.0

<u>pp after-i</u>	<u>336</u>	<u>3.6</u>
minute	<a href="#">150</a>	4.03
break	<a href="#">11</a>	1.55
goal	<a href="#">17</a>	1.18
ball	<a href="#">9</a>	0.85

<u>possessor</u>	<u>1934</u>
poacher	<a href="#">14</a>
opponent	<a href="#">39</a>
striker	<a href="#">15</a>
defender	<a href="#">11</a>
visitor	<a href="#">67</a>
opposition	<a href="#">16</a>
government	<a href="#">169</a>
charity	<a href="#">38</a>
organization	<a href="#">21</a>
client	<a href="#">61</a>
administration	<a href="#">18</a>
learner	<a href="#">10</a>
organisation	<a href="#">71</a>
researcher	<a href="#">10</a>



# WebBootCaT

- BootCaT integrated in SkE
- BootCaT a corpus
  - Clean, de-dupe, POS-tag, then
  - *Load into Sketch Engine*

# WebBootCaT: Create corpus

Corpus ID

Unique identifier of your corpus. May only contain letters, numbers, underscores and hyphens.

Language

Creating BootCaT corpora is available only for those language which we can at least tokenise. All such languages are listed here.

Build word sketches

This option only has effect if a pre-loaded sketch grammar is available for the selected language.

Input type  Seed words  
 URLs

Select "URLs" to download data from specified URLs rather than use seed words for finding the URLs.

Seed words

Use space as separator. Enclose multiword expressions into quotes (").

Please select URLs which you would like to process.

Cancel

< Back

OK

Query: "volcanic eruption" Eyjafjallajokull geodic

no results found

Query: stratigraphic tephra volcanology

- <http://en.wikipedia.org/wiki/Volcanology>
- <http://www.unige.ch/sciences/terre/mineral/volcano/Recherche/tephraChaiten.html>
- <http://www.earth-prints.org/bitstream/2122/2883/3/bozza%20IC-%20giaccio.pdf>
- <http://tripatlas.com/Volcanology>
- <http://publications.esc.cam.ac.uk:8080/509/>
- <http://www.answers.com/topic/volcanology>
- <http://www.therafoundation.org/articles/volcanology/>
- <http://www.earth-prints.org/handle/2122/4827>
- <http://www.geo.mtu.edu/~raman/papers2/CentralAmer/Alvarado/soto%20arenal.pdf>
- <http://maps.thefullwiki.org/Volcanology>

Invert selection

Query: "volcanic eruption" stratigraphic tephra

- <http://www.archaeologywordsmith.com/lookup.php?category=&where=headword&terms=tephra>
- [http://en.wikipedia.org/wiki/List\\_of\\_potentially\\_dangerous\\_volcanic\\_eruptions](http://en.wikipedia.org/wiki/List_of_potentially_dangerous_volcanic_eruptions)

## volcano\_en: WebBootCaT: Downloading data...

		31%
Successfully processed files	8	
Files remaining	46	
Data downloaded	669 kB	
Tokens retrieved	25,336	
Tokens per file (avg)	3,167	
Time elapsed	0:43	
Estimated time remaining	1:34	
Average file processing time	2.1 s	
Errors		13
- unable to retrieve		8
- invalid content-type		0
- file size out of range		2
- number of words out of range		3
- keywords filter applied		0
- unable to convert to text		0
- duplicate		0

 [Cancel processing](#)

```

Processing http://www.answers.com/topic/tephrochronology-1
- Content-type: text/html;charset=UTF-8
- File type: html
- Bytes read: 43468
- Number of words (approx.): 138
  - Too few words (min: 300)
Processing http://geology.geoscienceworld.org/cgi/content/abstract/15/9/809
- Content-type: text/html
- File type: html
- Bytes read: 29773
- Number of words (approx.): 172
  - Too few words (min: 300)
Processing http://tripatlas.com/Vesuvius

```

<a href="#">ukWaC</a>	English	1,565,
<a href="#">FinnishWaC</a>	Finnish	144,
<a href="#">frWaC</a>	French	1,628,

SH

## My corpora

Corpus ID	Corpus name	Language	Size	
<a href="#">brewing</a>	brewing	English	351,078	
<a href="#">camera</a>	camera	English	80,435	
<a href="#">ergative_en2</a>	ergative_en2	English	1,037,229	
<a href="#">tufs_anatomy</a>	tufs_anatomy	English	94,111	
<a href="#">ukwac_10m</a>	ukwac_10m	English	1,002,520	
<a href="#">volcano2_en</a>	volcano2_en	English	946,819	
<a href="#">volcano_en</a>	volcano_en	English	77,031	
<a href="#">volcanology</a>	volcanology	English	73,906	
<a href="#">ergative_fr3</a>	ergative_fr3	French	663,014	
<a href="#">volcan_fr</a>	volcan_fr	French	660,178	
<a href="#">heidelberg-anatomy</a>	heidelberg-anatomy	German	18,094	
<a href="#">sanskrit</a>	sanskrit	Hindi	5	
<a href="#">turkish1</a>	turkish1	Turkish	17,554	

[Create corp](#)

Sketch grammar development corpora

# How big a corpus do we get?

<b>Lg</b>		<b>Q's sent</b>	<b>Volcanoes</b>			<b>Stradivarius</b>		
			<b>Url</b>	<b>Doc</b>	<b>KW</b>	<b>Url</b>	<b>Doc</b>	<b>KW</b>
<b>En</b>	<b>B</b>	10	84	51	244	70	46	230
		50	318	180	679	230	150	1230
		250	941	515	1580	808	483	5326
	<b>Y</b>	10	67	39	152	60	47	148
		50	281	176	445	267	196	1071
		250	867	527	1232	937	649	3700
<b>Fr</b>	<b>B</b>	10	79	45	150	74	52	264
		50	246	152	461	225	145	1020
		250	755	506	1445	612	379	3815
	<b>Y</b>	10	79	36	118	82	60	720
		50	285	154	695	257	156	1155
		250	994	527	1737	843	510	2317



# Observation

- Specialist domain, L1
- Specialist domain, L2
- *Matching terminology*





# Going multilingual

- Translate seeds

- **English:** volcanology volcanologist "volcanic eruption" seismographs Eyjafjallajokull geodic "deformation monitoring" tephra magma stratigraphic tephrochronology geochronological "volcanic ash" ablation rhyolitic
- **French:** vulcanologue volcanologie "éruption volcanique" sismographes Eyjafjallajokull "surveillance de la déformation" géodiques tephra magma téphrochronologie stratigraphique géochronologiques "de cendres volcaniques" ablation rhyolitiques

- BootCaT for French

# volcanic *(adjective)*

volcano2\_en freq = 2207

and/or	531	2.0	modifies	2110	3.8
deadly	<a href="#">18</a>	10.0	eruption	<a href="#">414</a>	10.94
explosive	<a href="#">24</a>	9.82	ash	<a href="#">241</a>	10.9
active	<a href="#">23</a>	9.57	activity	<a href="#">160</a>	10.57
major	<a href="#">22</a>	9.44	rock	<a href="#">115</a>	9.6
recent	<a href="#">17</a>	9.41	hazard	<a href="#">45</a>	9.18
Icelandic	<a href="#">12</a>	9.07	field	<a href="#">47</a>	9.0
felsic	<a href="#">9</a>	8.99	cloud	<a href="#">42</a>	8.95
subglacial	<a href="#">7</a>	8.43	glass	<a href="#">36</a>	8.83
geothermal	<a href="#">7</a>	8.42	complex	<a href="#">29</a>	8.69
hot	<a href="#">7</a>	8.3	center	<a href="#">30</a>	8.68
Quaternary	<a href="#">7</a>	8.27	vent	<a href="#">29</a>	8.61
igneous	<a href="#">7</a>	8.18	gas	<a href="#">33</a>	8.49
porous	<a href="#">5</a>	8.17	event	<a href="#">33</a>	8.47
future	<a href="#">5</a>	8.0	plume	<a href="#">27</a>	8.47
tectonic	<a href="#">6</a>	7.98	system	<a href="#">27</a>	8.38

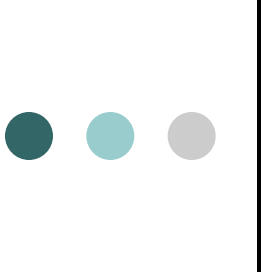
# volcanique *(adjective)* volcan\_fr

modifié	869	5.6	pp_sous-x	2	5.6
éruption	<a href="#">105</a>	10.87	pression	<a href="#">2</a>	8.8
roche	<a href="#">98</a>	10.72	pp_à-x	<a href="#">8</a>	1.9
activité	<a href="#">70</a>	10.7	travers	<a href="#">2</a>	10.1
édifice	<a href="#">44</a>	10.51	l	<a href="#">2</a>	7.4
cendre	<a href="#">59</a>	10.43	pp_au-x	<a href="#">10</a>	1.1
explosivité	<a href="#">15</a>	9.09	cour cours	<a href="#">2</a>	8.4
risque	<a href="#">16</a>	9.06			
île	<a href="#">17</a>	9.02			
verre	<a href="#">15</a>	9.0			
manifestation	<a href="#">12</a>	8.74			
bombe	<a href="#">12</a>	8.73			
arc	<a href="#">13</a>	8.71			
faciès	<a href="#">12</a>	8.66			
gaz	<a href="#">14</a>	8.64			
zone	<a href="#">15</a>	8.55			



# CCBC

- Input: L1, L1 seeds, L2
- Choose dictionary
  - Google as default
    - Google dictionary (25 lg pairs, limited API)
    - Google translate (1225 lg pairs, only 1 transl)
  - Option: edit translations
- Bootcat 2 corpora
- Bilingual word sketches



# Bilingual word sketches (very first pass)

- For L1 nodeword  $n$ 
  - For each of its translations  $n_1, n_2, \dots$ 
    - For each collocate  $c$  in word sketch
      - For each of its translations  $c_1, c_2, \dots$ 
        - Does  $c_i$  occur as collocate in word sketch for  $n_i$ ?
        - If yes: output  $\langle c; n_i, c_i \rangle$
        - Add L1 and L2 examples sentences

**volcano** (*noun*) volcano2\_en freq = 2490

volcan (*noun*) volcan\_fr freq = 1268

- be** [224](#) **Close view of the eruption column of Mount St. Helens on May 18 , 1980 ; the volcano photograph .**
- être [87](#) Le sommet du volcan est dégagé et l'on peut apercevoir une épaisse colonne de vapeur blanche monter droit dans le ciel.
- aller [6](#) Ce gonflement du volcan va générer des microséismes , une augmentation de l'inclinaison des pentes du volcan , une activité de la caldeira sommitale .
- produire [4](#) Le volcan a produit l'un des plus grands volumes de lave de partout dans le monde au cours du dernier millénaire , soit plus de 100 km<sup>3</sup> de lave.
- venir [3](#) Les dynamismes : Les principales différences entre volcans viennent de la puissance de leur éruption , ou de leur explosivité.
- active** [103](#) **Is an ongoing eruption the only way to tell if a volcano is active ?**
- actif [64](#) Les scientifiques considèrent habituellement un volcan être actif s'il est actuellement éclatant ou montrant des signes de tremblement de terre ou les nouvelles émissions de gaz significatives .
- eruption** [89](#) **Photo showing a cross section of the volcano before the 1980 eruption and the 330 Mt .**
- éruption [2](#) Lave La lave est une roche issue d'un magma qui est émise sous une forme plus ou moins fluide par les volcans en éruption.
- erupt** [85](#) **Volcanic ash : how far will it fall downwind from an erupting volcano ?**
- have** [61](#) **Pinatubo in June 1991 , this volcano has been an ideal location for remote sensing a volcano .**
- prendre [3](#) Puisque les coulées pyroclastiques successives se sont accumulées à la base de la montagne , le volcan a pris une forme conique.
- avoir [37](#) Des parties du premier cône du volcan ont été déplacées par des glaciers durant la Glaciation qui eut lieu il y a 14 000 ans.
- Icelandic** [37](#) **Well folks , it was only a matter of time .... " Ice cap thaw may awaken Icelandic volcano**
- islandais [2](#) Les effets géographiques des éruptions à l'échelle régionale Perturbations de la circulation aérienne 13 Pendant l'éruption du volcan Eyjafjallajökull en 2010 , les autorités de régulation aéronautique ont convenablement joué leur rôle d'informateur par le biais des METARs ( METARs )



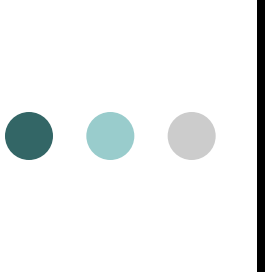
# Matching seeds – how?

- User translates
  - Yes but limited
- Bilingual dictionary
  - Yes but finding them??
  - Google dictionary
- Machine translations
- **Wikipedia**
  - Matching articles



# Evaluation

- Extract terms for L1, L2
- Ask expert
  1. Are they terms
  2. Do the L1, L2 lists contain translations of each other?

- 
- 3 Ig-pairs
    - En-Fr, En-De, En-Cz
    - One expert for each pair
  - 3 domains
    - Volcanoes
    - Stradivarius
    - Pancreatic cancer
      - Wikipedia: En and De only





# Results

Who	Wds	Trans	Mwds
<i>Volcanoes, En</i>			
E-Cz	29/30		10/85
E-De	29/30	10/29	16/85
E-Fr	29/30	19/29	24/85
<i>Stradivarius, En</i>			
E-Cz	19/29		13/85
E-De	26/30	3/26	9/85
<i>Stradivarius, De</i>			
E-De	16/30	2/16	6/84
<i>Cancer, En</i>			
E-De	27/30	9/27	
<i>Cancer, De</i>			
E-De	22/30	10/22	8/90
<i>Volcanoes, Fr</i>			
E-De <sup>25</sup>	27/30	19/27	5/83

In brief

- Words good
- Multiwords bad



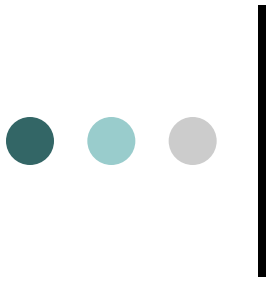
# Unithood and termhood

- To find terms
  - For multiwords only
    - Does it hang together?
    - *Unithood*
  - It it distinctive?
    - Keywords
    - *Termhood*
- We didn't use termhood for multiwords but need to



# Next steps

- Termhood for multiwords
- WebBootCaT from wikipedia
- From collocations to terms
  - More-than-2-word collocations
    - ... deadline next Tuesday



Thank you

<http://www.sketchengine.co.uk>