# New word sketch data format

Miloš Jakubíček

Lexical 文 Computing

`milos.jakubicek@sketchengine.co.uk`

6th Sketch Engine Workshop
Herstmonceux, August 10, 2015

# Word sketches

- you know them everybody

# Word sketch

## resource *(noun)* **British National Corpus freq = 12658** (112.8 per million)

| modifier | 6477 | 1.5 | object_of | 3285 | 2.2 | modifies | 1906 | 0.5 | subject_of | 512 | 0.6 |
|----------|------|-----|-----------|------|-----|----------|------|-----|------------|-----|-----|
| scarce | 163 | 9.53 | allocate | 194 | 9.58 | allocation | 135 | 9.42 | devote | 28 | 7.69 |
| natural | 321 | 8.94 | pool | 39 | 8.43 | implication | 46 | 7.09 | consume | 4 | 5.36 |
| limited | 187 | 8.86 | exploit | 64 | 8.23 | management | 153 | 6.98 | tie | 6 | 4.87 |
| financial | 249 | 8.3 | divert | 38 | 7.86 | defense | 7 | 6.68 | last | 4 | 4.6 |
| mineral | 89 | 8.19 | deploy | 31 | 7.67 | Stonier | 6 | 6.65 | back | 5 | 4.5 |
| additional | 107 | 7.92 | devote | 44 | 7.64 | utilisation | 7 | 6.63 | stretch | 4 | 4.29 |
| valuable | 74 | 7.86 | concentrate | 62 | 7.35 | committee | 132 | 6.49 | result | 6 | 3.93 |
| extra | 88 | 7.53 | utilise | 22 | 7.28 | centre | 158 | 6.4 | depend | 6 | 3.84 |
| human | 134 | 7.38 | conserve | 17 | 7.09 | allocator | 5 | 6.4 | limit | 5 | 3.59 |
| renewable | 33 | 7.31 | lack | 37 | 7.0 | depletion | 6 | 6.21 | match | 3 | 3.58 |
| adequate | 49 | 7.28 | reallocate | 13 | 6.98 | pack | 17 | 6.2 | share | 6 | 3.55 |
| non-renewable | 25 | 6.97 | mobilise | 13 | 6.83 | investigator | 8 | 6.17 | earn | 3 | 3.55 |
| existing | 53 | 6.68 | mobilize | 13 | 6.79 | column | 20 | 6.16 | enable | 7 | 3.54 |
| finite | 22 | 6.66 | distribute | 29 | 6.73 | constraint | 14 | 6.14 | remain | 12 | 3.5 |

# Word sketch format 1 (2003)

- triples: headword – relation – collocation
- in a corpus: include headword position
- $\Rightarrow$ 4-tuple $(h, r, c, h_{pos})$

# Word sketch format 2 (2010)

- include collocation position
- $\Rightarrow$ 5-tuple $(h, r, c, h_{pos}, c_{pos})$

# Word sketch format 3 (2014)

- include longest-commonest match (**Vít Baisa, Wed 10.20**)
- more scalable – for corpora over 100 billion words
- $\Rightarrow$ n-tuple $(h, r, c, h_{pos}, c_{pos}, LCM_1, \ldots, LCM_n)$

# Word sketch format 4 (2015)

- score computation changes
- see `https://www.sketchengine.co.uk/statistics-used-in-sketch-engine/`

# Conclusions

New word sketch format being introduced:

- biggest changes since 2003
- affects local installations
- changes score computations
- is backward compatible, part of Manatee 2.125
- smaller indices, faster access
- corpora in SkE now being rebuilt