

# D4.1: Methodology of Sketch Grammar evaluation

The evaluation of word sketch grammars can be performed in a number of ways depending on the particular purpose word sketches are used for. In the following we however focus on a general-purpose approach to the word sketch evaluation and related issues.

## Quantitative evaluation

Quantitative evaluation of word sketches can be useful especially when developing a new sketch grammar or trying to improve the coverage of an existing one. The quantitative evaluation in this case consists merely of estimating what percentage of corpus tokens is covered as headwords with meaningful grammatical relations.

This information is useful for potential improvements of individual grammatical relations as defined using the CQL queries (Jakubicek, 2010). The Sketch Engine (Kilgarriff, 2004, 2014) has been recently enhanced by feature that directly shows the proportion of corpus positions covered by the displayed word sketch table for a given headword.

**test** (*noun*) Alternative PoS: *verb* (6,896)

British National Corpus freq = **15,789** (140.74 per million) Coverage: **92.85%**

modifier		modifies			object_of		and/or		subject_of					
	8,554	1.80		3,517	0.90		3,461	2.10		1,476	0.70		1,123	1.20
driving	<u>127</u>	8.64	cricket	<u>105</u>	8.65	pass	<u>294</u>	8.29	international	<u>13</u>	7.58	reveal	<u>51</u>	6.48
blood	<u>204</u>	8.19	match	<u>170</u>	8.16	undergo	<u>53</u>	7.78	histology	<u>7</u>	7.09	match	<u>19</u>	6.14
fitness	<u>82</u>	7.93	statistic	<u>32</u>	8.07	satisfy	<u>49</u>	7.75	x ray	<u>11</u>	7.09	prove	<u>44</u>	6.00
beta	<u>69</u>	7.91	kit	<u>60</u>	7.95	fail	<u>134</u>	7.55	examination	<u>28</u>	6.23	consist	<u>13</u>	5.42
u	<u>49</u>	7.52	ban	<u>58</u>	7.83	devise	<u>39</u>	7.39	exam	<u>8</u>	6.17	detect	<u>9</u>	5.41
diagnostic	<u>50</u>	7.44	score	<u>76</u>	7.81	conduct	<u>69</u>	7.39	Southampton	<u>6</u>	6.12	<b>show</b>	<u>124</u>	5.37
screening	<u>46</u>	7.22	tube	<u>60</u>	7.79	administer	<u>31</u>	7.24	board	<u>42</u>	5.34	confirm	<u>17</u>	5.23
means	<u>97</u>	7.15	wicket	<u>35</u>	7.71	perform	<u>77</u>	7.23	interview	<u>12</u>	5.11	measure	<u>13</u>	5.17
breath	<u>75</u>	7.14	pilot	<u>56</u>	7.50	apply	<u>101</u>	6.82	check	<u>6</u>	5.04	check	<u>13</u>	5.09
nuclear	<u>75</u>	7.10	debut	<u>33</u>	7.18	boycott	<u>13</u>	6.68	measurement	<u>6</u>	4.86	assess	<u>7</u>	4.34
exact	<u>43</u>	7.04	cricketer	<u>20</u>	7.17	carry	<u>120</u>	6.32	trial	<u>13</u>	4.81	determine	<u>14</u>	4.28
laboratory	<u>62</u>	7.01	meal	<u>60</u>	7.13	stand	<u>80</u>	6.24	blood	<u>12</u>	4.64	result	<u>8</u>	4.28
statistical	<u>42</u>	6.96	certificate	<u>40</u>	6.92	standardise	<u>8</u>	6.11	assessment	<u>11</u>	4.57	apply	<u>15</u>	4.19
liver	<u>42</u>	6.96	suite	<u>22</u>	6.76	withstand	<u>9</u>	6.10	laboratory	<u>6</u>	4.50	demonstrate	<u>6</u>	4.05
acid	<u>70</u>	6.94	series	<u>85</u>	6.59	face	<u>57</u>	6.06	curriculum	<u>6</u>	4.25	run	<u>25</u>	4.00
Hiv	<u>39</u>	6.92	item	<u>61</u>	6.52	screen	<u>10</u>	5.99	procedure	<u>12</u>	4.16	involve	<u>20</u>	3.98
written	<u>45</u>	6.91	innings	<u>12</u>	6.44	use	<u>227</u>	5.44	exercise	<u>7</u>	4.00	fail	<u>10</u>	3.96
intelligence	<u>52</u>	6.91	result	<u>115</u>	6.25	repeat	<u>16</u>	5.41	analysis	<u>13</u>	3.92	drive	<u>6</u>	3.31
rank	<u>51</u>	6.86	rig	<u>10</u>	6.16	undertake	<u>18</u>	5.37	result	<u>21</u>	3.85	indicate	<u>6</u>	3.25
antibody	<u>39</u>	6.81	treaty	<u>30</u>	6.16	survive	<u>13</u>	5.36	treatment	<u>10</u>	3.70	identify	<u>6</u>	3.13
reasonableness	<u>31</u>	6.80	strip	<u>17</u>	6.12	introduce	<u>36</u>	5.35	culture	<u>8</u>	3.69	provide	<u>24</u>	3.06
forensic	<u>31</u>	6.80	igg	<u>8</u>	6.10	play	<u>66</u>	5.33	datum	<u>6</u>	2.76	lay	<u>6</u>	2.96
smear	<u>31</u>	6.80	batsman	<u>10</u>	6.10	score	<u>12</u>	5.32	patient	<u>7</u>	2.59	follow	<u>16</u>	2.96
positive	<u>58</u>	6.78	flight	<u>29</u>	6.07	develop	<u>48</u>	5.27	case	<u>10</u>	2.10	remain	<u>8</u>	2.89
iq	<u>29</u>	6.75	drive	<u>27</u>	6.02	win	<u>37</u>	5.12	study	<u>6</u>	1.73	suggest	<u>9</u>	2.87

Figure 1: Sample word sketch table for the English noun *test* with its coverage displayed.

In Figure 1 the screenshot from Sketch Engine for the word sketch of the English noun *test* shows a coverage of 92.85 %, meaning that this percentage of all occurrences of *test* as noun in the corpus is covered by some grammatical relation defined in the word sketch grammar and hence the collocations are captured by the word sketch table as presented.

## Qualitative evaluation

### Extrinsic evaluation

The qualitative evaluation of word sketches is of course a much more complex task -- as it is usually the case. Given the practical nature of word sketches, ideally they should be evaluated extrinsically by measuring their contribution to improvements given a particular task.

A candidate task would be a real lexicographic setting where e.g. one group of lexicographers would be working on dictionary entries while using Sketch Engine with sketch grammar A and another group would be using the same setup with sketch grammar B. One could evaluate the group performance in terms of speed but also in terms of quality of the entries, number of examples etc.

Such evaluations have been performed as in-house trials by major UK publishing houses (e.g. Oxford University Press, Cambridge University Press, Macmillan etc.) who are using Sketch Engine but only in very small scale and only to compare with a null variant -- i.e. not using Sketch Engine.

### Intrinsic evaluation

First experiments with both quantitative and qualitative evaluation of word sketches have been performed in 2010 for Dutch, English, Japanese and Slovene, Japanese (see Kilgarriff et al. 2010). In these experiments only precision of the collocations automatically found by word sketches was evaluated, not the recall. Therefore, one was able to answer the question: *how many of the automatically found collocations are good?* but not *which good collocates are missing?* For this some form of a gold standard of collocations for a set of headwords would be necessary so as to have a reference set to be compared with.

In 2011 we therefore setup (and published in Kilgarriff, 2012) an intrinsic evaluation task where we have collected a gold standard of collocates (disregarding which grammatical relation they would fit in) by hiring a group of professional lexicographers (for English) and linguist students (for Czech) as annotators, providing them with collocation candidates and asking the question *Would you put this word into a collocation dictionary?* The answers were binary yes or no.

Of course this very much depends on the actual definition of a collocation (which later turned out to be the biggest methodological problem). The definition that the annotators were given was the one given in the Oxford Collocations Dictionary:

*Collocation is the way words combine in a language to produce natural-sounding speech and writing. ... Combinations of words in a language can be ranged on a cline from the totally free — see a man/car/book — to the totally fixed and idiomatic — not see the wood for the trees. ... All these combinations, apart from those at the very extremes of the cline, can be called collocation. And it is combinations such as these — particularly in the ‘medium-strength’ area — that are vital to communicative competence in English. (Crowther et al., 2002, vii)*

When preparing the collocation candidates we collected a number of text corpora (both web-based and edited-content-based) and existing collocations dictionaries. The headwords were only nouns, adjectives and verbs chosen randomly from three frequency bands:

- high: top 100–2999 words by frequency
- mid: top 3000–9999 words by frequency
- low top 10,000–30,000 words by frequency

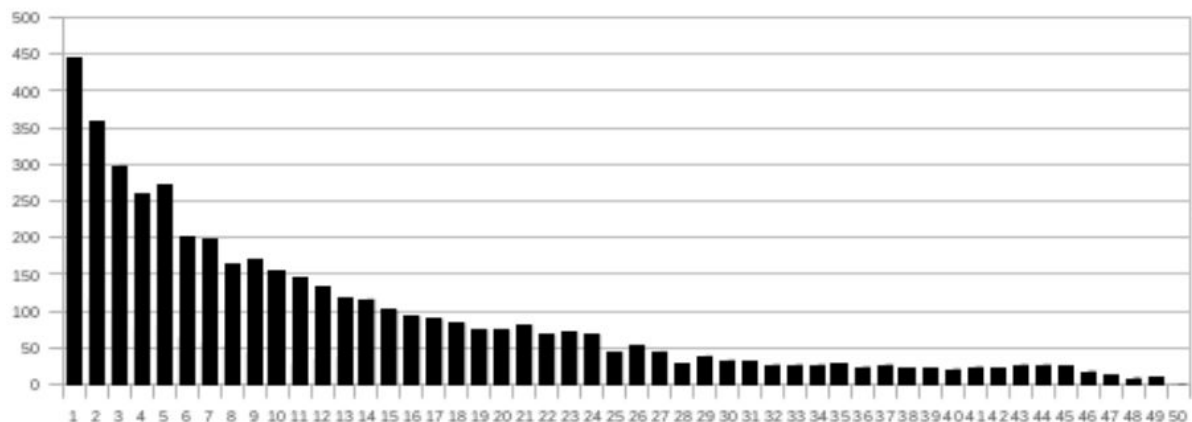


Figure 2: Distribution of good collocations in fiftieths, ordered by score.

We used various sources to achieve close to 100 % coverage for the collocation candidates and we clearly showed that there is a strong correlation between the word sketch score and the goodness of a collocation (as provided in Figure 2) nevertheless our later experiments have shown that a full coverage was not achieved, which of course represents a methodological issue from the point of view of using such a dataset as a gold standard. With respect to this, the findings were blatant in saying:

- *if we showed judges more candidates from the same corpora, they found more collocations (though with diminishing returns)*
- *if we showed judges more candidates from new corpora, they found more collocations.*

It is worth noting that the task as presented in (Kilgarriff, 2014) was setup an extrinsic corpus evaluation using word sketches -- of course, this becomes then intrinsic evaluation from the perspective of a word sketch grammar.

The underlying assumption was that, having a processing pipeline for a corpus as shown in Figure 3, one can evaluate any of the components by interchanging it but keeping the rest the same.



Figure 3: sample processing pipeline from corpus source to collocations

We have evaluated a large number of corpora with different processing tools and various parameters -- including different taggers and word sketch grammar exploiting their annotation.

The most cumbersome issue remained the definition of collocation which materialized in rather low inter-annotator agreement: for Czech, the pairwise agreement varied between 73.6% and 82.3%, for English, between 81.1% and 85.8% (with three annotator for both languages). Later investigations revealed that for some of the headwords (such as *plutonium*), no pair of annotators agreed on a good collocation. We have found out that the guidance to the annotators was clearly adopted with a printed dictionary on mind, where the set of collocations is rather small and focuses only on very strong collocations.

However, what we were actually most interested in, was an evaluation which would clearly distinguish between obvious errors (i.e. non-collocations) and remaining mild, medium or strong collocations. Such type of evaluation could be used to improvements alongside the whole processing pipeline, including word sketch grammar, in a rather straightforward way. In some cases, we have also noticed that, since the words were chosen randomly, there were terms from domains which the lexicographers were not familiar with, and hence were unable to identify good and strong collocations.

Recently, we have therefore reviewed the methodology and we are now in preparation of a new gold standard collocation set following a revised methodology for annotation that aims to be more inclusive. The annotators are now classifying into five categories:

- strong collocation
- weak collocation
- correct word combination but not a significant collocation
- error
- I don't understand

To help reducing the number of unknown collocations, the word sketches have been enhanced by the so called longest-commonest match (LCM) string -- the most common headword-collocation combination (Kilgarriff et al. 2015).

We have also devised a detailed annotation manual that would help the annotators to decide among these categories which is provided in Figure 4:

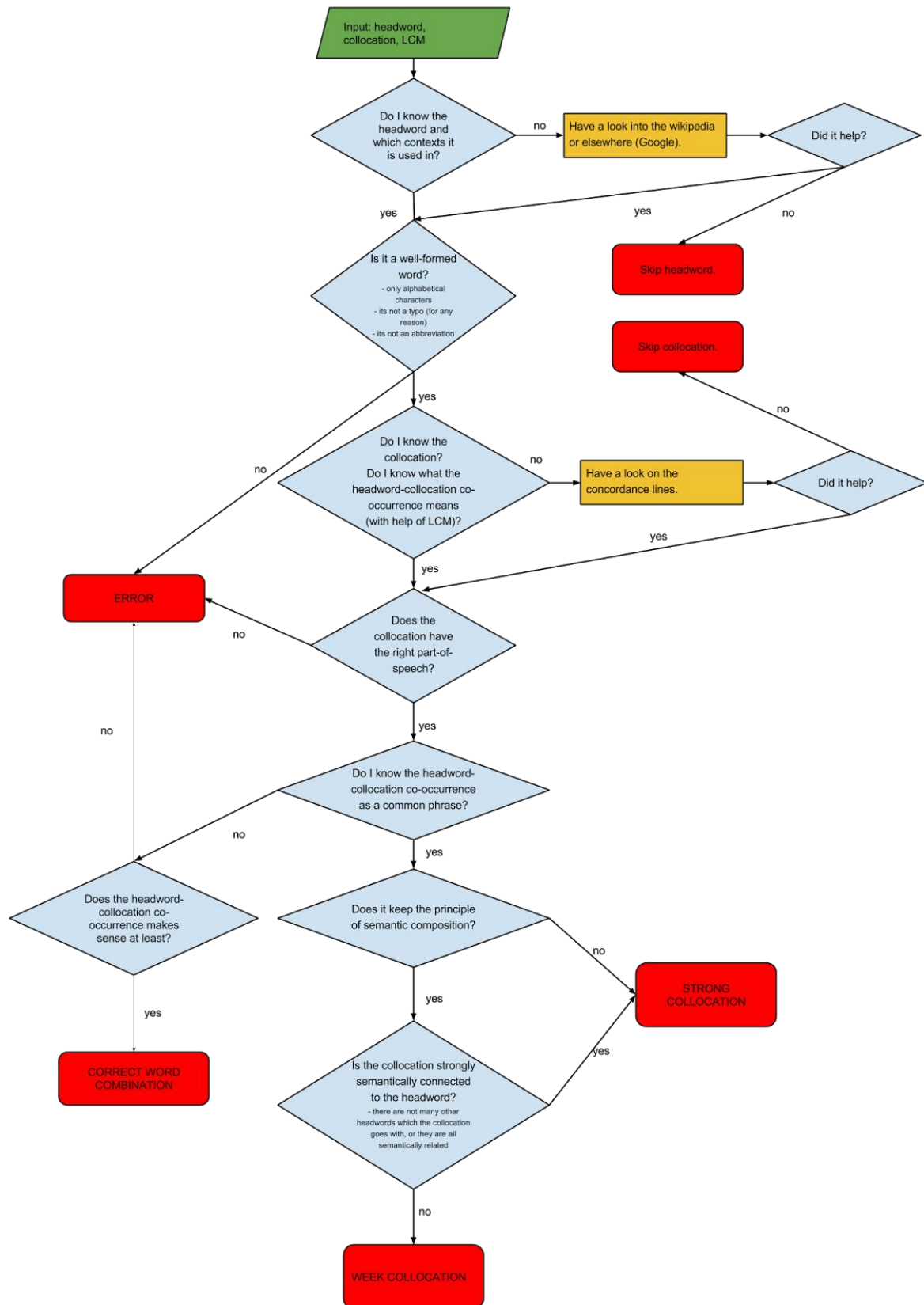


Figure 4: New annotation manual for collocations classification

## References

JAKUBÍČEK, Miloš, et al. Fast Syntactic Searching in Very Large Corpora for Many Languages. In: PACLIC. 2010. p. 741-747.

KILGARRIFF, Adam, et al. A quantitative evaluation of word sketches. In: *Proceedings of the 14th EURALEX International Congress, Leeuwarden, The Netherlands*. 2010.

KILGARRIFF, Adam, et al. Extrinsic corpus evaluation with a collocation dictionary task. In: *Proceedings of the 9th International Conference on Language Resources and Evaluation (LREC 2014)*. 2014.

KILGARRIFF, Adam, et al. Longest–commonest Match. 2015.