

# Finding commonest match for word sketch collocations

**Vít Baisa**

Lexical Computing

SKEW3, 2012, Brno

- intended as helpful feature for WS evaluators:  
**see** has\_obj **final**

- intended as helpful feature for WS evaluators:  
**see** `has_obj` **final** → **saw** *world cup* **finals**

- intended as helpful feature for WS evaluators:  
**see** `has_obj` **final** → **saw** *world cup* **finals**
- discover a common string for collocates
- using left and right context
- include intervening words

# Algorithm

**Input:** a pair of lemmas from a word sketch

**Init:** from corresponding concordance lines, use strings from the left lemma to the right lemma (*match*) plus 3 left and 3 right words

```
if the most freq. match occurs > N/4 times:
    commonest_match = the most frequent match
    n = freq(commonest_match)
    if the most freq. L/R extension occurs > n/4 times:
        extended_commonest_match = the most freq. L/R extension
        n = freq(extended_commonest_match)
        if the most freq. L/R extension occurs < n/4:
            return extended_commonest_match
        else:
            extend recursively L/R, up to 3 words
    else:
        return commonest_match
else:
    return ""
```

# Example I

## Word sketch for lemma **morning**

<i>modifier</i>	tomorrow	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">1422</a>	tomorrow morning [93.0 %]
<i>modifier</i>	next	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">8722</a>	next morning [87.5 %]
<i>modifier</i>	yesterday	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">1267</a>	yesterday morning [84.1 %]
<i>modifier</i>	coffee	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">1254</a>	coffee morning [49.2 %]
<i>modifier</i>	early	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">4387</a>	early morning [77.3 %]
<i>modifier</i>	sunny	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">490</a>	sunny morning [85.1 %]
<i>modifier</i>	weekday	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">289</a>	weekday morning [62.3 %]
<i>modifier</i>	mid	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<a href="#">297</a>	mid morning [86.9 %]

## Example II

Extended commonest matches for collocates from various gramrels.

lemma	collocate	commonest	%
evaluation	histological	eosin for histological evaluation	37.5
evaluation	determinant	evaluation of determinant	40.0
key	backspace	TAB and BACKSPACE keys can be used	30.0
key	keyboard	key on the keyboard	27.3
squirrel	golden-mantled	golden-mantled ground squirrels	33.3
squirrel	popper	squirrels in a popcorn popper	27.3
republic	USSR	republics of the USSR	43.4
bride	toast	toast to the bride and groom	55.6
marry	subsequently	whom he subsequently married	40.0

# What it is good for?

- provide evaluators with explanatory examples
- provide the most common utterance of the collocates
- Good Dictionary EXample for the collocates
- discover obvious multiword expressions



# Conclusion, future plans

- accessible on `beta.sketchengine.co.uk` (**Sketch-Eval**)
- computed for top N collocates, on the fly
- make it faster
- pre-compute commonest matches
- show within WS by default
- get feedback from users
- tune parameters, or provide a way to change them