# Semantic Word Sketches

Diana McCarthy,[†] Adam Kilgarriff,[◇]
Miloš Jakubíček,[◇‡] Siva Reddy[⋆]

DTAL University of Cambridge[†], Lexical Computing[◇],
University of Edinburgh[⋆], Masaryk University[‡]

July 2015

# Outline

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Concordances
Word Sketches

# The Sketch Engine

- concordances, word lists, collocations
- word sketches
  - create and examine syntactic profiles and collocations of words
  - input automatic part-of-speech tags and a bespoke 'sketch grammar'
- automatic thesauruses: which other words have similar profiles?
- sketch differences between words

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Concordances
Word Sketches

# The Sketch Engine

## for viewing corpora



Semantic Word Sketches

# The Sketch Engine

Word Sketches: syntactic profiles

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Concordances
Word Sketches

# Sketch Grammars

Under the hood

- Definitions: define('any_noun','"N.."')

  ...

- Relations

  =subject/subject_of
  2:any_noun rel_start? adv_aux_string_incl_be 1:verb_not_pp
  2:any_noun rel_start? adv_aux_string_incl_be aux_have adv_string 1:past_part

  1:past_part adv_string [word="by"] long_np

The Sketch Engine
**Semantic Tagging**
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Super Sense Tagger (SST)
SST Supersenses

# Semantic Class Tagging

- aim to build word sketches on syntactic and semantic information
- automatic 'superclass' tagging technology
- superclass: a coarse grained semantic class that is applicable to multiple words (e.g. **animal** for *cat*, *fly*, *hare*, *pig* etc...
- allow search and analysis with these classes and
- semantic word sketches: basic semantic frame with semantic preferences for arguments

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Super Sense Tagger (SST)
SST Supersenses

# Semantic Class Tagging

Super Sense Tagger (SST) Ciaramita and Altun (2006)
(http://sourceforge.net/projects/supersensetag/)

- semantic tags are WordNet Fellbaum (1998) lexicographer classes
- supervised word sense disambiguation (i.e. it requires hand labelled data for training) using a Hidden Markov Model e.g. labels *mouse* as **animal**, **artifact**)
- SemCor (Landes et al., 1998) used as training data
- Named Entity Recognition e.g. < *RHM Technology Ltd.*> **organization**
- Multiword tagging using multiwords from WordNet e.g. *couch potato*

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Super Sense Tagger (SST)
SST Supersenses

## SST WordNet Noun Classes (25)

| | |
|---|---|
| act | acts or actions |
| object | natural objects (not man-made) |
| animal | animals |
| quantity | quantities and units of measure |
| artifact | man-made objects |
| phenomenon | natural phenomena |
| attribute | attributes of people and objects plant plants |
| food | food and drinks |
| . . . | . . . |

# SST WordNet Verb Classes (15)

| | |
|---|---|
| body | grooming, dressing and bodily care |
| emotion | feeling |
| change | size, temperature change, intensifying |
| motion | walking, flying, swimming |
| cognition | thinking, judging, analyzing, doubting |
| perception | seeing, hearing, feeling |
| communication | telling, asking, ordering, singing |
| possession | buying, selling, owning |
| creation | sewing, baking, painting, performing |
| . . . | . . . |

# Experiments

- just over 25% of the UKWaC Ferraresi et al. (2008)
- SST tagged with
    - part-of-speech tags (Penn TreeBank)
    - supersenses (WordNet labels)
    - Named Entity Labels
    - WordNet multiwords

# Semantic Tags in the Concordance

First | Previous    Page [3] of 789    Go    Next | Last

| #1493049 | credibility that surrounds the <mwe><ne> Mickey | Mouse | /NNP/other.n | </ne></mwe><n |
| #1496976 | hovering , <mwe> seeking out </mwe> a vole or | mouse | /NN/animal.n | . INFORMATION F |
| #1545720 | audio visuals and the <mwe><ne> Soviet Spy | Mouse | /NNP/other.n | Trail </ne></mw |
| #1561637 | even after disruption . Control through the | mouse | /NN/animal.n | Interactive stori |
| #1561653 | keyboard . They are controlled entirely by the | mouse | /NN/animal.n | which moves the |
| #1561673 | of effects . In a very real sense , the | mouse | /NN/animal.n | represents cont |
| #1561693 | importance of access to the controlling device ( | mouse | /NN/animal.n | or keyboard , <n |
| #1561728 | </mwe> , 1993 ) . Seen in this light , the | mouse | /NN/animal.n | might be conside |
| #1561813 | different methods were used when passing the | mouse | /NN/animal.n | to the <ne> next |
| #1561826 | Sometimes , <mwe> for example </mwe> , when the | mouse | /NN/animal.n | had been left in |
| #1561862 | remaining member of the group moved the | mouse | /NN/animal.n | towards the nex |
| #1561928 | within the group . If simply leaving the | mouse | /NN/animal.n | is seen as unhel |
| #1562003 | members who were not in possession of the | mouse | /NN/animal.n | issued a significa |
| #1562048 | in <mwe> actual possession </mwe> . While the | mouse | /NN/animal.n | gave the undisp |
| #1562081 | directed the same commands to the holder of the | mouse | /NN/animal.n | , adding extra ps |
| #1608486 | and dragging the control points using the | mouse | /NN/animal.n | ( left button ) o |
| #1608563 | dialog or dragging the points using left <mwe> mouse | /NN/artifact.n | button </mwe> . |
| #1664217 | with a cenral finger-ball ( e.g. marble | mouse | /NN/animal.n | ) . Just more ph |
| #1768071 | students attending workshops ( eg. <mwe> mouse | /NN/artifact.n | mats </mwe> , |
| #1831217 | </ne></mwe> include the <mwe><ne> Yellow Necked | Mouse | /NNP/other.n | </ne></mwe> , |

First | Previous    Page [3] of 789    Go    Next | Last

# Semantic Tags in the Word Sketch (selected)

**eat** *(verb)*   **UKWaC super sensed freq = 26329** (71.2 per million)

| transframe | 1241 | 6.6 | intransframe | 1021 | 2.7 |
|---|---|---|---|---|---|
| person.n_*consumption.v_food.n | 178 | 11.42 | animal.n_*consumption.v | 85 | 10.84 |
| group.n_*consumption.v_food.n | 57 | 10.25 | person.n_*consumption.v | 382 | 10.25 |
| person.n_*consumption.v_plant.n | 37 | 9.77 | group.n_*consumption.v | 145 | 9.82 |
| person.n_*consumption.v_animal.n | 35 | 9.62 | 0_*consumption.v | 61 | 9.73 |
| animal.n_*consumption.v_animal.n | 30 | 9.56 | state.n_*consumption.v | 20 | 8.89 |
| animal.n_*consumption.v_plant.n | 25 | 9.32 | time.n_*consumption.v | 19 | 8.74 |
| animal.n_*consumption.v_food.n | 24 | 9.21 | communication.n_*consumption.v | 30 | 8.57 |
| person.n_*consumption.v_person.n | 52 | 8.97 | artifact.n_*consumption.v | 30 | 8.51 |
| 0_*consumption.v_food.n | 20 | 8.92 | food.n_*consumption.v | 19 | 8.4 |
| animal.n_*consumption.v_artifact.n | 19 | 8.79 | other.n_*consumption.v | 15 | 8.1 |

# Semantic Tags in the Word Sketch (selected)



laugh (verb)    UKWaC super sensed freq = 6489 (17.5 per million)

| V_PP | 148 | 9.1 |
|---|---|---|
| *body.v_at_cognition.n | 14 | 11.02 |
| *body.v_at_communication.n | 12 | 10.94 |
| *body.v_at_person.n | 7 | 9.77 |
| *body.v_as_person.n | 6 | 9.78 |
| *communication.v_at_location.n | 6 | 7.8 |
| *body.v_in_cognition.n | 5 | 9.58 |
| *communication.v_at_communication.n | 5 | 8.39 |
| *communication.v_at_person.n | 5 | 8.14 |

| intransframe | 1101 | 8.9 |
|---|---|---|
| person.n_*body.v | 556 | 10.49 |
| group.n_*body.v | 143 | 10.02 |
| 0_*body.v | 102 | 10.33 |
| artifact.n_*body.v | 49 | 9.19 |
| time.n_*body.v | 23 | 8.49 |
| location.n_*body.v | 16 | 7.88 |
| event.n_*body.v | 9 | 7.56 |
| other.n_*body.v | 8 | 7.38 |
| cognition.n_*body.v | 7 | 7.2 |
| communication.n_*body.v | 7 | 6.89 |

The Sketch Engine
Semantic Tagging
**Semantic Tags in Sketch Engine**
Comparison to FrameNet
Conclusions
References

In the Concordance
**Semantic Word Sketches**
Other Possibilities from SST Output

# Semantic Word Sketch Grammar

An example for the intransitive frame

=intransframe
*COLLOC "%(2.sense)_*%(1.sense)-x"
        2:any_noun rel_start? adv_aux_string_incl_be 1:verb_not_pp
not_np_start
        2:any_noun rel_start? adv_aux_string_incl_be aux_have
adv_string 1:past_part not_np_start

# MWEs: Sketch Diff chip (green) vs chips (red)

| mwe | 325 | 314 | -2.0 | -2.0 |
|---|---|---|---|---|
| fish_and_chip_food.n | 0 | 70 | 0.0 | 10.4 |
| tortilla_chip_food.n | 0 | 14 | 0.0 | 8.3 |
| potato_chip_food.n | 10 | 40 | 7.7 | 9.7 |
| memory_chip_artifact.n | 25 | 32 | 9.0 | 9.4 |
| gene_chip_artifact.n | 11 | 10 | 7.9 | 7.8 |
| poker_chip_artifact.n | 14 | 12 | 8.3 | 8.0 |
| silicon_chip_artifact.n | 40 | 27 | 9.7 | 9.1 |
| bargaining_chip_attribute.n | 29 | 7 | 9.3 | 7.2 |
| chocolate_chip_food.n | 14 | 0 | 8.3 | 0.0 |
| chip_shot_act.n | 15 | 0 | 8.4 | 0.0 |

# Portion of Sketch Diff laugh (green) vs cry (red)

| V_PP | 148 | 117 | 9.1 | 7.0 |
|---|---|---|---|---|
| *communication.v_in_location.n | 0 | 6 | 0.0 | 7.3 |
| *stative.v_for_act.n | 0 | 11 | 0.0 | 6.6 |
| *communication.v_at_location.n | 6 | 0 | 7.8 | 0.0 |
| *communication.v_at_person.n | 5 | 0 | 8.1 | 0.0 |
| *communication.v_at_communication.n | 5 | 0 | 8.4 | 0.0 |
| *body.v_in_cognition.n | 5 | 0 | 9.6 | 0.0 |
| *body.v_at_person.n | 7 | 0 | 9.8 | 0.0 |
| *body.v_as_person.n | 6 | 0 | 9.8 | 0.0 |
| *body.v_at_communication.n | 12 | 0 | 10.9 | 0.0 |
| *body.v_at_cognition.n | 14 | 0 | 11.0 | 0.0 |

# Semantic Word Lists: CQL + Word Frequency (Communication Verbs)

# Semantic Word Lists: FindX (communication verbs)

| | | | |
|---|---|---|---|
| 0.1 | 122.9 | say.-v | 271 |
| 0.1 | 119.4 | tell.-v | 62 |
| 0.1 | 112.0 | ask.-v | 75 |
| 0.2 | 101.9 | out.-v | 53 |
| 0.2 | 100.0 | humour-v | 53 |
| 0.2 | 100.0 | critique-v | 142 |
| 0.2 | 100.0 | underline-v | 2166 |
| 0.2 | 100.0 | stammer-v | 50 |
| 0.2 | 100.0 | reintroduce-v | 501 |
| 0.2 | 100.0 | re-introduce-v | 109 |
| 0.2 | 100.0 | shriek-v | 116 |
| 0.2 | 100.0 | exhort-v | 88 |
| 0.2 | 100.0 | publicise-v | 1244 |
| 0.2 | 100.0 | chide-v | 116 |
| 0.3 | 100.0 | interrogate-v | 730 |
| 0.3 | 100.0 | fate-v | 67 |
| 0.3 | 100.0 | bemoan-v | 277 |
| 0.3 | 100.0 | absolve-v | 136 |
| 0.3 | 100.0 | signpost-v | 160 |
| 0.3 | 100.0 | unrated-v | 321 |
| 0.3 | 100.0 | chronicle-v | 487 |
| 0.3 | 100.0 | telegraph-v | 75 |
| 0.3 | 100.0 | spam-v | 218 |
| 0.3 | 100.0 | misquote-v | 119 |
| 0.4 | 100.0 | extol-v | 87 |
| 0.4 | 100.0 | eschew-v | 322 |
| 0.4 | 100.0 | nominates-v | 71 |
| 0.4 | 100.0 | evince-v | 156 |
| 0.4 | 100.0 | spoof-v | 66 |
| 0.4 | 100.0 | rejuvenate-v | 205 |
| 0.4 | 100.0 | symbolise-v | 292 |
| 0.4 | 100.0 | pardon-v | 445 |

# Comparing to FrameNet (Ruppenhofer et al., 2010)

- FrameNet contains lots of useful information e.g. [FRAME **employing**:
  Frame Elements: Employer Employee Position Tasks Compensation ...
  Definition: An Employer *employs* an Employee whose Position entails that the Employee perform certain Tasks in exchange for Compensation
- lots of other information
  - lexical units *employ.v commision.v staff.n employment.n*
  - precedes frame **firing**
  - with corpus examples, *I employed him as Chief Gardener for ten years*
- but manually produced so low coverage
- Semantic word sketches can provide additional information and high coverage

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
**Conclusions**
References

# Summary

- semantic tagging alongside part-of-speech for semantic word sketches
- provide syntactic and semantic profiling for
  - semantic queries and word lists
  - semantic and syntactic profiling in the word sketch
  - comparing words by the profiles

# Future Possibilities

- try other semantic tagsets, taggers and tools
- sketch grammar could be developed further
- no identification of semantic roles as yet in contrast to FrameNet (Ruppenhofer et al., 2010), Propbank (Palmer et al., 2005) and VerbNet (Kipper-Schuler, 2005)
- Semantic word sketches could be used to provide selectional preferences and corpus information to such resources

# Thank You

| | | |
|---|---|---|
| http://www... | see , a little more really means a lot - | **Thank** /communication. |
| http://www... | Jamie and Lynne Reilly 29 June , 2002 `` | **Thank** /communication. |
| http://www... | Team , ( Jack , Billy and young Andy , ) | **THANK** /communication. |
| http://www... | important member of the family , and we | **thank** /communication. |
| http://www... | B. I would have liked the opportunity to | **thank** /communication. |
| http://www... | obtained some very interesting information . I | **thank** /communication. |
| http://bee... | again for eternity in heaven . I want to | **thank** /communication. |
| http://www... | to cover , but thank you for coming and | **thank** /communication. |
| http://www... | ZENAB ( ) posted : 08.03.2006 message : | **Thank** /communication. |
| http://www... | Joanne.gowing@pizzaexpress.com Happy donating and | **thank** /communication. |
| http://www... | Platform : GameCube Sent by : Andrew Bernish ( | **Thank** /communication. |
| http://www... | to keep Nicholas in all of your prays and | **thank** /communication. |
| http://www... | top Thank you June 2006 I would like to | **thank** /communication. |
| http://spo... | Peter MacLeod . Hello to you too , sir , and | **thank** /communication. |
| http://www... | , I do . You seem so damn sure . HOMER : | **Thank** /communication. |

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
References

Ciaramita, M. and Altun, Y. (2006). Broad-coverage sense disambiguation and information extraction with a supersense sequence tagger. In *Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing*, pages 594–602, Sydney, Australia. Association for Computational Linguistics.

Fellbaum, C., editor (1998). *WordNet, An Electronic Lexical Database*. The MIT Press, Cambridge, MA.

Ferraresi, A., Zanchetta, E., Baroni, M., and Bernardini, S. (2008). Introducing and evaluating ukwac, a very large web-derived corpus of english. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC 2008)*, Marrakech, Morocco.

Kipper-Schuler, K. (2005). *VerbNet: A broad-coverage, comprehensive verb lexicon*. PhD thesis, Computer and

The Sketch Engine
Semantic Tagging
Semantic Tags in Sketch Engine
Comparison to FrameNet
Conclusions
**References**

Information Science Dept., University of Pennsylvania.
Philadelphia, PA.

Landes, S., Leacock, C., and Randee, I. T. (1998). Building
  semantic concordances. In Fellbaum, C., editor, *WordNet: an
  Electronic Lexical Database*, pages 199–237. MIT Press.

Palmer, M., Gildea, D., and Kingsbury, P. (2005). The proposition
  bank: A corpus annotated with semantic roles. *Computational
  Linguistics*, 31(1):71–106.

Ruppenhofer, J., Ellsworth, M., Petruck, M. R. L., Johnson, C. R.,
  and Scheffczyk, J. (2010). FrameNet II: Extended theory and
  practice. Technical report, International Computer Science
  Institute, Berkeley. http://framenet.icsi.berkeley.edu/.