



# **Lexikalische Semantik: WordNet**

Skript, Version 2019-11-10

**Timm Lichte**



- 1 Intro
- 2 Fakten und Zahlen
- 3 Struktur von WordNet
- 4 Erstellung von WordNet
- 5 Unzulänglichkeiten von WordNet
- 6 Relationspfade und Ähnlichkeit
- 7 Erweiterungen
- 8 Begriffsdiagramm



# **Lexikalische Semantik: WordNet**

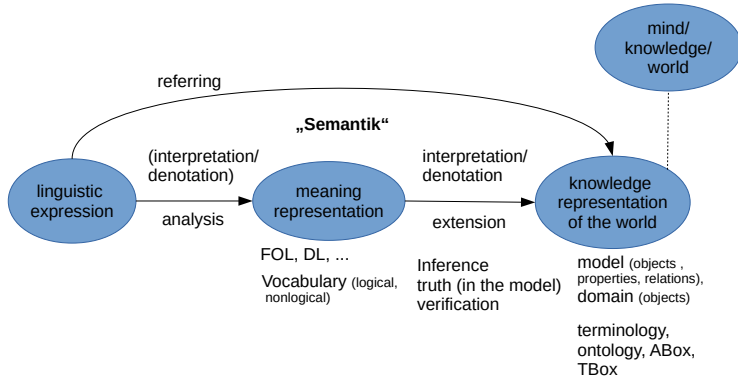
Skript, Version 2019-11-10

**Timm Lichte**



Die **SEMANTIK** beschäftigt sich mit der Bedeutung von Symbolen bzw. Zeichen (Symbole, die eine Bedeutung haben). Sind diese Symbole linguistische Wortformen wie **Bank**, dann haben wir es genauer mit der **LEXIKALISCHEN SEMANTIK** zu tun.

Vielleicht die wichtigste Ressource im Bereich der lexikalischen Semantik ist **WORDNET**. Dank der relativ einfachen, theorieneutralen Grundprinzipien genießt WordNet auch nach 30 Jahren immer noch einen Sonderstatus in der NLP.



Semantik/Bedeutung stellt eine **Verbindung** her zwischen einem Zeichen und der “Welt”.



- Miller, George A. 1995. WordNet: A lexical database for English. Communications of the ACM 38(11). 39–41.
- Jurafsky, Daniel & James H. Martin. 2018. Speech and Language Processing. An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition. Third edition draft of September 23, 2018. <https://web.stanford.edu/~jurafsky/slp3/ed3book.pdf>.

Wird jeweils angezeigt!



## WordNet:

- entwickelt seit 1986 unter der Leitung von George A. Miller (†)
- entstanden im Bereich der Künstlichen Intelligenz und Psychologie
- Umfang Version 3.1 von 2012 (Quelle Wikipedia):
  - 155 327 Lemmata
  - 175 979 Synsets (~Bedeutungen)
  - 207 016 Lemma-Synset-Paare
  - durchschn. Lemma-Ambiguität: 1,34
  - durchschn. Synset-Umfang: 1,176





## Synset

Each member of a given synset expresses the same concept, though not all synset members are interchangeable in all contexts.<sup>[2]</sup>

Tests für (Fast-)Synonymie (“near synonymy”):

- “interchangeability in some contexts”
- *Mein Wagen / Auto ist kaputt.*
- *Ich brauche eine Münze für einen Wagen / #ein Auto .*





Jedes Synset hat die folgenden Attribute:

- Synset-Label hello.n.01
- Lemmata *hello, hullo, hi, howdy*
- Definition/Glosse 'an expression of greeting'
- Examples *every morning they exchanged polite hellos*
- Frequenz (in SemCor)



Semantische Relationen bestehen zwischen **nominalen** Synsets:

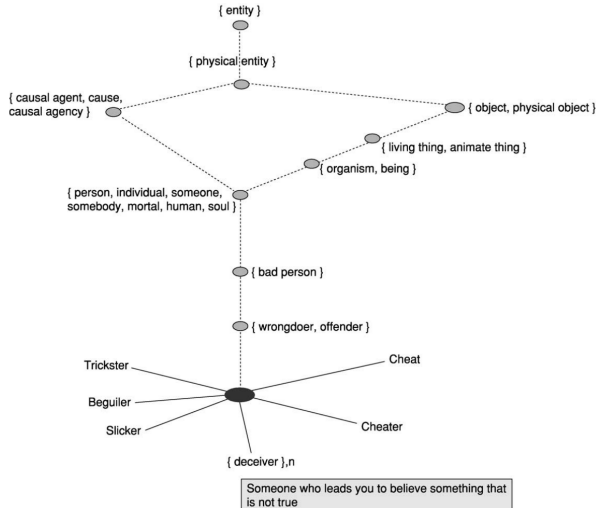
Relation	Also Called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> <sup>1</sup> → <i>meal</i> <sup>1</sup>
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> <sup>1</sup> → <i>lunch</i> <sup>1</sup>
Instance Hypernym	Instance	From instances to their concepts	<i>Austen</i> <sup>1</sup> → <i>author</i> <sup>1</sup>
Instance Hyponym	Has-Instance	From concepts to concept instances	<i>composer</i> <sup>1</sup> → <i>Bach</i> <sup>1</sup>
Member Meronym	Has-Member	From groups to their members	<i>faculty</i> <sup>2</sup> → <i>professor</i> <sup>1</sup>
Member Holonym	Member-Of	From members to their groups	<i>copilot</i> <sup>1</sup> → <i>crew</i> <sup>1</sup>
Part Meronym	Has-Part	From wholes to parts	<i>table</i> <sup>2</sup> → <i>leg</i> <sup>3</sup>
Part Holonym	Part-Of	From parts to wholes	<i>course</i> <sup>7</sup> → <i>meal</i> <sup>1</sup>
Substance Meronym		From substances to their subparts	<i>water</i> <sup>1</sup> → <i>oxygen</i> <sup>1</sup>
Substance Holonym		From parts of substances to wholes	<i>gin</i> <sup>1</sup> → <i>martini</i> <sup>1</sup>
Antonym		Semantic opposition between lemmas	<i>leader</i> <sup>1</sup> ↔ <i>follower</i> <sup>1</sup>
Derivationally Related Form		Lemmas w/same morphological root	<i>destruction</i> <sup>1</sup> ↔ <i>destroy</i>

**Figure C.2** Noun relations in WordNet.

(aus Jurafsky & Martin [3] )



(aus Fellbaum [2])



**Figure 1** A WordNet noun tree.



Semantische Relationen bestehen zwischen **verbalen** Synsets:

Relation	Definition	Example
Hypernym	From events to superordinate events	<i>fly</i> <sup>9</sup> → <i>travel</i> <sup>5</sup>
Troponym	From events to subordinate event (often via specific manner)	<i>walk</i> <sup>1</sup> → <i>stroll</i> <sup>1</sup>
Entails	From verbs (events) to the verbs (events) they entail	<i>snore</i> <sup>1</sup> → <i>sleep</i> <sup>1</sup>
Antonym	Semantic opposition between lemmas	<i>increase</i> <sup>1</sup> ⇔ <i>decrease</i> <sup>1</sup>
Derivationally Related Form	Lemmas with same morphological root	<i>destroy</i> <sup>1</sup> ⇔ <i>destruction</i> <sup>1</sup>

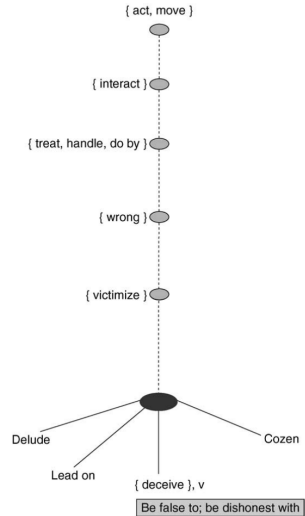
**Figure C.3** Verb relations in WordNet.

(aus Jurafsky & Martin [3] )

# Beispiel: *deceive*



(aus Fellbaum [2])



**Figure 2** A WordNet verb tree.

# Beispiel: *dry* versus *wet*



(aus Fellbaum [2])

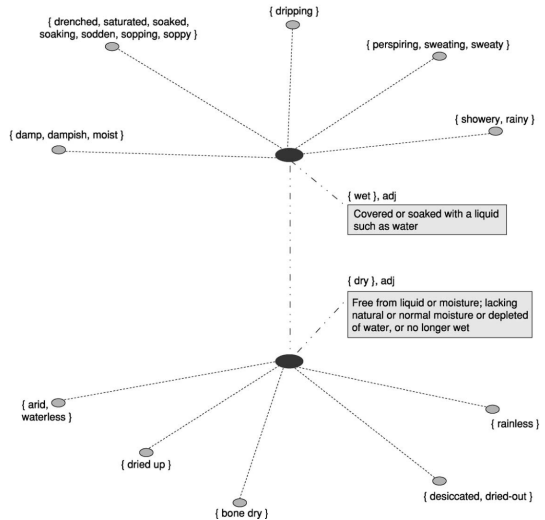


Figure 3 An adjective cluster in WordNet.



Unterschiedliche Dateien je nach syntaktischer Kategorie (Zahlen aus Jurafsky & Martin [3]):

- Nomen: 117 798 (1,23 Bedeutungen/Lemma)
- Verben: 11 529 (2,16 Bedeutungen/Lemma)
- Adjektive: 22 479
- Adverbien: 4 481



# Unique Beginners

Die WordNet-Hierarchie für Nomen kann in mehrere Hierarchien partitioniert werden, die “relatively distinct semantic fields representieren”

Die Wurzel (“top”) einer solchen Teilhierarchy wird auch “unique beginner” genannt (aus Fellbaum [1: §1.2]):

**Table 1.1**

List of 25 unique beginners for noun source files

{act, activity}	{food}	{possession}
{animal, fauna}	{group, grouping}	{process}
{artifact}	{location}	{quantity, amount}
{attribute}	{motivation, motive}	{relation}
{body}	{natural object}	{shape}
{cognition, knowledge}	{natural phenomenon}	{state}
{communication}	{person, human being}	{substance}
{event, happening}	{plant, flora}	{time}
{feeling, emotion}		





# Unique Beginners – Taxonomic Sisters – Semantic Atoms

---

Das entspricht in gewisser Weise den “**taxonomic sisters**” bei Kroeger (ohne “mutually exclusive” zu sein), oder den **semantischen Atomen** in der Komponentenanalyse.

Die Herausforderung ist offensichtlich:

*The problem, of course, is to **decide what the primitive semantic components should be**. Different workers make different choices; one important criterion is that, collectively, they should provide a place for every English noun. (Miller in Fellbaum [1])*



(aus Fellbaum [1])



**Figure 1.1**

Diagrammatic representation of relations that reduce the 25 noun source files to 11 unique beginners. The unique beginners are italicized.



Fellbaum [1: §1.2]:

- seldom more than **10 or 12 levels deep**
- The deepest examples usually contain technical distinctions that are not part of the everyday vocabulary.



<http://wordventure.eti.pg.gda.pl/wne/wne.html>





ganz viel Handarbeit . . . .

⇒ teuer und zeitaufwendig, aber hohe Qualität

Bemühungen, den Erstellungsprozess so weit wie möglich zu automatisieren.



Miller in Fellbaum [1: §1.4]

- is-not-a-Relationen / Ausnahmen:
  - *A whale is not a fish.*
- Unterschiedliche Arten von Hyponymie:
  - **taxonomic** (is-a-kind-of)
  - **functional** (is-used-as-a-kind-of)
  - {*chicken*} @→ {*bird*} versus {*chicken*} @→ {*food*}
- Unterscheidung zwischen Eigennamen und Gattungsname, Maßnamen und zählbaren Nomen
- andere semantische Relationen
  - implizit verfügbar über Glossen



Fellbaum [2]:

- syntagmatische Relationen:
  - thematische Rollen (AGENS, PATIENS, ...)
  - grammatische Funktionen (Subjekt, Objekt, ...)

## The Tennis Problem<sup>[1]</sup>

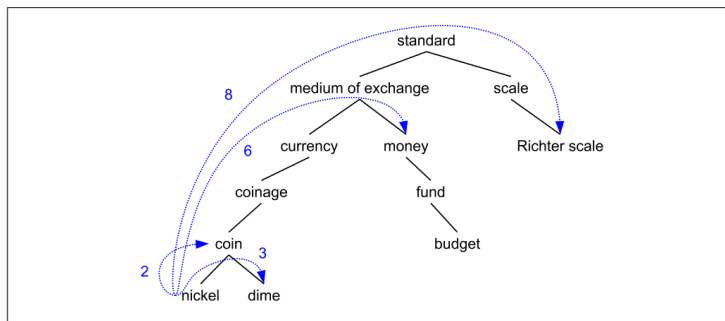
- Semantische Felder wie *racquet*, *ball*, and *net* werden nicht direkt repräsentiert.
- ad-hoc/diskursgetriggerte Relationen: *racquet*, *ball*, *net*, *court game*, *physician*, *hospital* ...

Aber vielleicht indirekt? → Relationspfade



## Intuition

Synsets sind konzeptuell um so ähnlicher, je näher sie in WordNet beieinander liegen, d.h. je **kürzer** der Relationspfad zwischen ihnen ist



**Figure C.5** A fragment of the WordNet hypernym hierarchy, showing path lengths (number of edges plus 1) from *nickel* to *coin* (2), *dime* (3), *money* (6), and *Richter scale* (8).

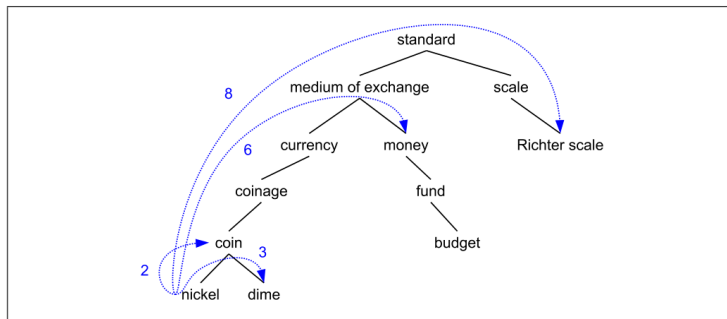




# Pfadlänge und Wortähnlichkeit

Das lässt sich sehr einfach auf Worte/Lemmata übertragen:

- $pathlen(c_1, c_2) = 1 + \text{the number of edges in the shortest path between the sense nodes } c_1 \text{ and } c_2$
- $sim_{path}(c_1, c_2) = \frac{1}{pathlen(c_1, c_2)}$
- $wordsim(w_1, w_2) = \max_{c_1 \in senses(w_1), c_2 \in senses(w_2)} sim(c_1, c_2)$



**Figure C.5** A fragment of the WordNet hypernym hierarchy, showing path lengths (number of edges plus 1) from *nickel* to *coin* (2), *dime* (3), *money* (6), and *Richter scale* (8).



# Gewichtung von Relationspfaden

## Problem

Die Ähnlichkeit zwischen *nickel* und *money* ist intuitiv viel größer als die zwischen *nickel* und *standard*.

**Lösung:** unterschiedliche Gewichtung von Relationspfaden

- abhängig von der Tiefe der Einbettung (Wu & Palmer)
- abhängig von der Auftretenshäufigkeit in einem Korpus

**Resnik/Information-based Word Similarity:**

- $P(c) = \frac{\sum_{w \in \text{words}(c)} \text{count}(w)}{N}$
- $IC(c) = -\log P(c)$
- $LCS(c_1, c_2)$  = the lowest common subsumer
- $sim_{resnik}(c_1, c_2) = -\log P(LCS(c_1, c_2))$



Mit solchen Ähnlichkeitsmaßen für Bedeutungen in WordNet lassen WSD-Verfahren umsetzen.

## Idee:

- Gegeben ein Target-Lemma in einem Satz mit möglichen Bedeutungen  $\sigma_1, \dots, \sigma_n$ .
- Man wähle die Bedeutung  $\sigma_i$  so, dass die Ähnlichkeit zu den Bedeutungen der anderen Lemmata im Satz maximiert wird.
- *Ich ging zur Bank, um Geld abzuheben.*



- Deutsch: GermaNet (in Tübingen)
- Polnisch
- Japanisch
- Hindu
- ...

<http://globalwordnet.org/wordnets-in-the-world/>



## IMAGENET: Photodatenbank mit Verknüpfungen zu WordNet-Bedeutungen

### Moneybag

A drawstring bag for holding money

389  
pictures

32.16%  
Popularity  
Percentile

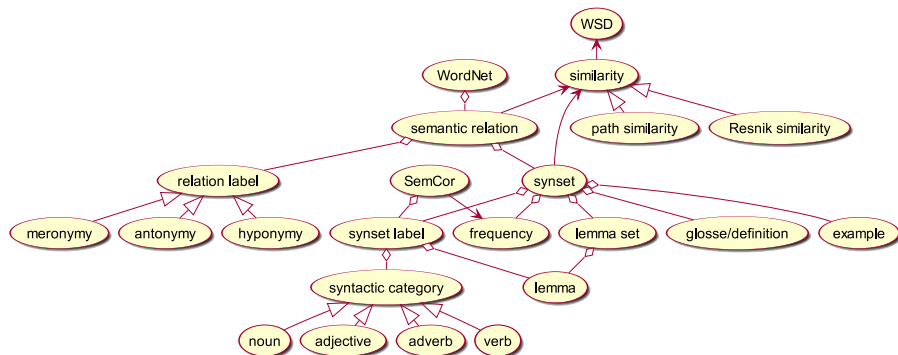
Wordnet  
IDs

- tucker-bag (0)
- backpack, back pack
- sweat bag (0)
- shopping bag (0)
- burn bag (0)
- bladder (0)
- sandbag (0)
- skin (2)
- ragbag (0)
- rosin bag (0)
- sick bag, sickbag (0)
- ice pack, ice bag (0)
- purse (0)
- sleeping bag (0)
- saddlebag (0)
- drawstring bag (3)
- moneybag (0)
- duffel bag, duffle
- seabag (0)
- dust bag, vacuum ba
- gunnysack, gunny sa
- plastic bag (1)
- schoolbag (0)
- sack, poke, paper ba
- nosebag, feedbag (0)
- body bag, personnel
- pouch (13)
- manger, trough (1)
- case, display case, sho
- spoon (9)
- powder horn, powder fla
- wastepaper basket, was
- sack (0)

[Treemap Visualization](#)
[Images of the Synset](#)
[Downloads](#)

\*Images of children synsets are not included. All images shown are thumbnails. Images may be subject to copyright.

[Prev](#)
[1](#)
[2](#)
[3](#)
[4](#)
[5](#)
[6](#)
[7](#)
[8](#)
[9](#)
[10](#)
[11](#)
[12](#)
[Next](#)





- [1] Fellbaum, Christiane (ed.). 1998. *WordNet: An electronic lexical database*. (Language, Speech, and Communication). Cambridge, MA: MIT Press.
- [2] Fellbaum, Christiane. 2006. WordNet(s). *Encyclopedia of Language & Linguistics*. 665–670.
- [3] Jurafsky, Daniel & James H. Martin. 2018. *Speech and Language Processing. An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Third edition draft of September 23, 2018.  
<https://web.stanford.edu/~jurafsky/slp3/ed3book.pdf>.
- [4] Miller, George A. 1995. WordNet: A lexical database for english. *Communications of the ACM* 38(11). 39–41.