

# Automatic Induction of Semantic Classes for German Verbs

Sabine Schulte im Walde  
Institut für Maschinelle Sprachverarbeitung (IMS)  
Universität Stuttgart

DGfS Jahrestagung 2003  
AG 6 : *Semantisches Wissen im Lexikon*  
Universität München

## Goals

- Automatic acquisition of high-quality and large-scale lexical resource for NLP applications: German semantic verb classes
- Theoretical investigation of relationship between verb behaviour and meaning components
- Development of clustering methodology suitable for the demands of natural language

## Overview

1. Theoretical aspects of verb classes
2. German verbs and verb classes
3. Clustering methodology
4. Corpus-based statistical lexical acquisition → verb description
5. Clustering experiments and results
6. NLP applications of lexical verb class information

## Verb Classes

*verb meaning components* ↔ *verb behaviour*

To a certain extent, the lexical meaning components of a verb determine its behaviour, particularly with respect to the choice of its arguments.

## Existing Verb Classes

- German (Schumacher 1986; Ballmer and Brennenstuhl 1986)
- English (Levin 1993)
- French (Saint-Dizier 1996)
- Spanish (Vázquez et al. 2000)

## German Verb Classes: Constitution

- 168 German verbs → 43 semantic classes
- Class size: 2-7 verbs (∅ 3.9 verbs per class)
- Ambiguity: 8 verbs with 2 senses
- High and low frequency verbs:  $8 \leq \text{freq} \leq 71,604$
- Basis: semantic intuition
- Relation to Levin (1993), consistency with Schumacher (1986)

## German Semantic Verb Classes (1)

1. *Aspect* : anfangen, aufhören, beenden, beginnen, enden
2. *Propositional Attitude* : ahnen, denken, glauben, vermuten, wissen
3. *Desire*
  - (a) *Wish* : erhoffen, wollen, wünschen
  - (b) *Need* : bedürfen, benötigen, brauchen
4. *Transfer of Possession (Obtaining)* : bekommen, erhalten, erlangen, kriegen
5. *Transfer of Possession (Giving)*
  - (a) *Gift* : geben, leihen, schenken, spenden, stiften, vermachen, überschreiben
  - (b) *Supply* : bringen, liefern, schicken, vermitteln<sub>1</sub>, zustellen
6. *Manner of Motion*
  - (a) *Manner of Locomotion* : gehen, klettern, kriechen, laufen, rennen, schleichen, wandern
  - (b) *Rotation* : drehen, rotieren
  - (c) *Rush* : eilen, hasten
  - (d) *Means* : fahren, fliegen, rudern, segeln
  - (e) *Flotation* : fließen, gleiten, treiben

## German Semantic Verb Classes (2)

7. *Emotion*
  - (a) *Origin* : ärgern, freuen
  - (b) *Expression* : heulen<sub>1</sub>, lachen<sub>1</sub>, weinen
  - (c) *Objection* : ängstigen, ekeln, fürchten, scheuen
8. *Manner of Look on the Face* : gähnen, grinsen, lachen<sub>2</sub>, lächeln, starren
9. *Perception* : empfinden, erfahren<sub>1</sub>, fühlen, hören, riechen, sehen, wahrnehmen
10. *Manner of Articulation* : flüstern, rufen, schreien
11. *Moaning* : heulen<sub>2</sub>, jammern, klagen, lamentieren
12. *Communication* : kommunizieren, korrespondieren, reden, sprechen, verhandeln
13. *Statement*
  - (a) *Announcement* : ankündigen, bekanntgeben, eröffnen, verkünden
  - (b) *Constitution* : anordnen, bestimmen, festlegen
  - (c) *Promise* : versichern, versprechen, zusagen
14. *Observation* : bemerken, erkennen, erfahren<sub>2</sub>, feststellen, realisieren, registrieren
15. *Description* : beschreiben, charakterisieren, darstellen<sub>1</sub>, interpretieren
16. *Presentation* : darstellen<sub>2</sub>, demonstrieren, präsentieren, veranschaulichen, vorführen
17. *Speculation* : grübeln, nachdenken, phantasieren, spekulieren



## German Semantic Verb Classes (3)

18. *Insistence* : beharren, bestehen<sub>1</sub>, insistieren, pochen
19. *Teaching* : beibringen, lehren, unterrichten, vermitteln<sub>2</sub>
20. *Position*
  - (a) *Bring into Position* : legen, setzen, stellen
  - (b) *Be in Position* : liegen, sitzen, stehen
21. *Production* : bilden, erzeugen, herstellen, hervorbringen, produzieren
22. *Renovation* : dekorieren, erneuern, renovieren, reparieren
23. *Support* : dienen, folgen<sub>1</sub>, helfen, unterstützen
24. *Quantum Change* : erhöhen, erniedrigen, senken, steigern, vergrößern, verkleinern
25. *Opening* : öffnen, schließen<sub>1</sub>
26. *Existence* : bestehen<sub>2</sub>, existieren, leben
27. *Consumption* : essen, konsumieren, lesen, saufen, trinken
28. *Elimination* : eliminieren, entfernen, exekutieren, töten, vernichten
29. *Basis* : basieren, beruhen, gründen, stützen
30. *Inference* : folgern, schließen<sub>2</sub>
31. *Result* : ergeben, erwachsen, folgen<sub>2</sub>, resultieren
32. *Weather* : blitzen, donnern, dämmern, nieseln, regnen, schneien

## Verb Class Scene Description: *Aspect Verbs*

Verbs: *anfangen, aufhören, beenden, beginnen, enden*

**Scene** : [*E* An event] begins or ends, either internally caused or externally caused by [*I* an initiator]. The event may be specified with respect to [*T* tempus], [*L* locus], [*X* an experiencer], or [*R* a result].

Frame Roles: **I**nitiator, **E**vent

Modification Roles: **T**emporal, **L**ocal, e**X**periencer, **R**esult

Levin class: 55.1 (*Aspectual Verbs* → *Begin Verbs*)

Schumacher class: 2.1 (*Verben der speziellen Existenz*  
→ *Verben der Existenzsituierung*)  
*beenden* is not classified.

## Verb Class Frame Variants: Aspect Verbs (1)

Frame	Participating Verbs & Corpus Examples
<p><b>n<sub>E</sub></b></p>	<p>+ anfangen, aufhören, beginnen / +<sub>adv</sub> enden / ¬ beenden</p> <p>Nun aber muß [<sub>E</sub> der Dialog] <b>anfangen</b>.</p> <p>... bevor [<sub>E</sub> der Golfkrieg] <b>angefangen</b> hatte ...</p> <p>... damit [<sub>E</sub> die Kämpfe] <b>aufhören</b>.</p> <p>Erst muß [<sub>E</sub> das Morden] <b>aufhören</b>.</p> <p>[<sub>E</sub> Der Gottesdienst] <b>beginnt</b>.</p> <p>[<sub>E</sub> Das Schuljahr] <b>beginnt</b> [<sub>T</sub> im Februar].</p> <p>[<sub>X</sub> Für die Flüchtlinge] <b>beginnt</b> nun [<sub>E</sub> ein Wettlauf gegen die Zeit].</p> <p>[<sub>E</sub> Sein kurzes Zwischenspiel] bei der Wehrmacht <b>endete</b> ... [<sub>R</sub> glimpflich].</p> <p>[<sub>E</sub> Die Ferien] <b>enden</b> [<sub>R</sub> mit einem großen Fest].</p> <p>[<sub>E</sub> Druckkunst] ... <b>endet</b> [<sub>R</sub> beim guten Buch].</p> <p>[<sub>E</sub> Die Partie] <b>endete</b> [<sub>R</sub> 0:1].</p> <p>[<sub>L</sub> An einem Baum] <b>endete</b> in Höchst [<sub>E</sub> die Flucht] ...</p> <p>[<sub>E</sub> Der Informationstag] ... <b>endet</b> [<sub>T</sub> um 14 Uhr].</p>

## Verb Class Frame Variants: Aspect Verbs (2)

Frame	Participating Verbs & Corpus Examples
<b>n<sub>I</sub></b>	<p>+ anfangen, aufhören / ¬ beenden, beginnen, enden</p> <p>[<sub>I</sub> Die Hauptstadt] muß <b>anfangen</b>.</p> <p>... daß [<sub>I</sub> er] [<sub>T</sub> pünktlich] <b>anfing</b>.</p> <p>Jetzt können [<sub>I</sub> wir] nicht einfach <b>aufhören</b>.</p> <p>Vielleicht sollte [<sub>I</sub> ich] <b>aufhören</b> und noch studieren.</p>
<b>n<sub>I</sub> a<sub>E</sub></b>	<p>+ anfangen, beenden, beginnen / ¬ aufhören, enden</p> <p>Nachdem [<sub>I</sub> wir] [<sub>E</sub> die Sache] <b>angefangen</b> haben,</p> <p>... [<sub>I</sub> er] versucht, [<sub>E</sub> ein neues Leben] <b>anzufangen</b>.</p> <p>[<sub>I</sub> Die Polizei] <b>beendete</b> [<sub>E</sub> die Gewalttätigkeiten].</p> <p>[<sub>T</sub> Nach dem Abi] <b>beginnt</b> [<sub>I</sub> Jens] [<sub>L</sub> in Frankfurt] [<sub>E</sub> seine Lehre] ...</p>
<b>n<sub>I</sub> a<sub>E</sub> [P]</b>	<p>+ anfangen, beenden, beginnen / ¬ aufhören, enden</p> <p>Wenn [<sub>E</sub> die Arbeiten] [<sub>T</sub> vor dem Bescheid] <b>angefangen</b> werden ...</p> <p>Während [<sub>X</sub> für Senna] [<sub>E</sub> das Rennen] <b>beendet</b> war ...</p> <p>[<sub>E</sub> Das Durcheinander, das es zu CDU-Zeiten gegeben habe,] sei <b>beendet</b> worden.</p> <p>... ehe [<sub>E</sub> eine militärische Aktion] <b>begonnen</b> wird ...</p>

## Verb Class Frame Variants: Aspect Verbs (3)

Frame	Participating Verbs & Corpus Examples
<b>n<sub>I</sub> i<sub>E</sub></b>	<p>+ anfangen, aufhören, beginnen / ¬ beenden, enden</p> <p>[<sub>I</sub> Ich] habe <b>angefangen</b>, [<sub>E</sub> Hemden zu schneiden].</p> <p>[<sub>I</sub> Die Bahn] will [<sub>T</sub> 1994] <b>anfangen</b> [<sub>E</sub> zu bauen].</p> <p>... daß [<sub>I</sub> der Alkoholiker] <b>aufhört</b> [<sub>E</sub> zu trinken].</p> <p>... daß [<sub>I</sub> die Säuglinge] einfach <b>aufhören</b> [<sub>E</sub> zu atmen].</p> <p>In dieser Stimmung <b>begannen</b> [<sub>I</sub> Männer] [<sub>E</sub> auf den Straßen den Tango zu tanzen] ...</p> <p>[<sub>I</sub> Tausende von Pinguinen] <b>beginnen</b> [<sub>E</sub> dort zu brüten].</p>
<b>n<sub>I</sub> p<sub>E</sub>:mit<sub>Dat</sub></b>	<p>+ anfangen, aufhören, beginnen / ¬ beenden, enden</p> <p>Erst als [<sub>I</sub> der versammelte Hofstaat] [<sub>E</sub> mit Klatschen] <b>anfang</b>,</p> <p>Aber [<sub>I</sub> wir] müssen endlich [<sub>E</sub> damit] <b>anfangen</b>.</p> <p>[<sub>I</sub> Der Athlet] ... kann ... [<sub>E</sub> mit seinem Sport] <b>aufhören</b>.</p> <p>... müßten noch [<sub>I</sub> viel mehr Frauen] [<sub>E</sub> mit ihrer Arbeit] <b>aufhören</b> ...</p> <p>Schließlich zog [<sub>I</sub> er] einen Trennstrich, <b>begann</b> [<sub>E</sub> mit dem Selbstentzug] ...</p> <p>[<sub>I</sub> Man] <b>beginne</b> [<sub>E</sub> mit eher katharsischen Werken].</p>
<b>n<sub>I</sub> p<sub>E</sub>:mit<sub>Dat</sub> [P]</b>	<p>+anfangen, aufhören, beginnen / ¬ beenden, enden</p> <p>Und [<sub>E</sub> mit den Umbauarbeiten] könnte <b>angefangen</b> werden.</p> <p>[<sub>E</sub> Mit diesem ungerechten Krieg] muß sofort <b>aufgehört</b> werden.</p> <p>[<sub>T</sub> Vorher] dürfe [<sub>E</sub> mit der Auflösung der Szene] nicht <b>begonnen</b> werden.</p>

## Clustering Methodology

1. Statistical acquisition of lexical verb information
2. Automatic verb clustering by standard technique k-Means
3. Clustering evaluation against manual verb classification

## Lexical Acquisition Framework

- Lexicalised probabilistic context-free grammar (Carroll/Rooth 1998)
- Unsupervised training by the *EM-Algorithm* (Baum 1972)
- Robust statistical parser **LoPar** (Schmid 2000)
- 35 million words of German newspaper corpora
- Lexicalised grammar rules and lexical choice parameters
- Corpus-based statistical lexical acquisition

## Corpus-Based Statistical Lexical Acquisition

- D1** purely syntactic definition of subcategorisation
- D2** syntactico-semantic definition of subcategorisation with prepositional preferences
- D3** syntactico-semantic definition of subcategorisation with prepositional and selectional preferences



## Subcategorisation Frame Elements

<b>n</b>	noun phrase: NP <sub>Nom</sub>
<b>a</b>	noun phrase: NP <sub>Akk</sub>
<b>d</b>	noun phrase: NP <sub>Dat</sub>
<b>r</b>	reflexive pronoun
<b>p</b>	prepositional phrase
<b>x</b>	expletive <i>es</i>
<b>i</b>	subordinated non-finite clause
<b>s-2</b>	subordinated finite verb second clause
<b>s-dass</b>	subordinated finite <i>dass</i> -clause
<b>s-ob</b>	subordinated finite <i>ob</i> -clause
<b>s-w</b>	indirect <i>wh</i> -questions
<b>k</b>	copula constructions

Examples: **nad**, **np:auf<sub>Akk</sub>**, **nai**

## Lexical Verb Information (1): Subcategorisation Frame Definition

Frame	Freq	Prob
ns-dass	1,929	0.27945
ns-2	1,888	0.27358
np	687	0.09951
n	608	0.08811
na	555	0.08046
ni	346	0.05015
nd	234	0.03392
nad	160	0.02325
nds-2	70	0.01011

probability distribution for *glauben* 'to think/believe'  
(probability values >1%)

## Lexical Verb Information (2): Subcategorisation Frame Definition + PPs

Refined Frame		Freq	Prob
<b>np:über</b> <sub>Akk</sub>	acc / 'about'	480	0.11981
<b>np:von</b> <sub>Dat</sub>	dat / 'about'	463	0.11568
<b>np:mit</b> <sub>Dat</sub>	dat / 'with'	280	0.06983
<b>np:in</b> <sub>Dat</sub>	dat / 'in'	81	0.02031

refined **np** probability distribution for *reden* 'to talk'  
with total joint probability  $p(\textit{reden}, \mathbf{np}) = 0.35820$   
(probability values >1%)

## Selectional Preferences: Nominal Level

Noun		Freq
Ziel	'goal'	86.30
Strategie	'strategy'	27.27
Politik	'policy'	25.30
Interesse	'interest'	21.50
Konzept	'concept'	16.84
Entwicklung	'development'	15.70
Kurs	'direction'	13.96
Spiel	'game'	12.26
Plan	'plan'	10.99
Spur	'trace'	10.91
Programm	'program'	8.96
Weg	'way'	8.70

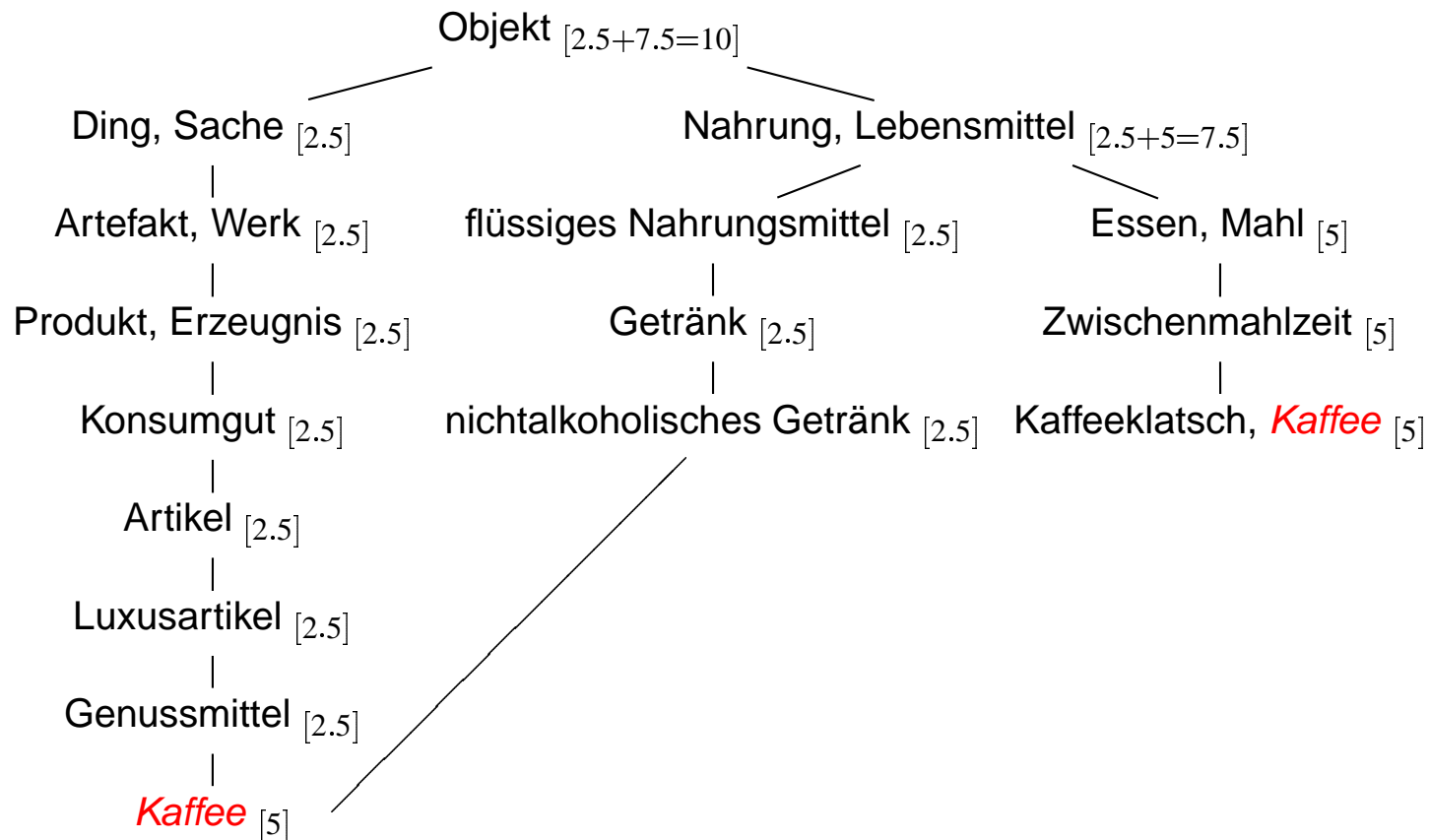
Nominal arguments for *verfolgen* 'to follow' in **na**

## Selectional Preferences: Nominal Level

Noun		Freq
Uhr	'o'clock'	85.38
Prozeß	'process'	77.14
Kampf	'fight'	69.77
Verhandlung	'negotiation'	65.51
Krieg	'war'	64.39
Tag	'day'	52.12
Zeit	'time'	51.94
Arbeit	'work'	46.86
Geschichte	'history'	46.10
Karriere	'career'	42.49
Spiel	'game'	37.79
Diskussion	'discussion'	31.25

Nominal arguments for *beginnen* 'to begin' in n

## Selectional Preferences: GermaNet Nodes (1)



GermaNet hierarchy for noun *Kaffee* 'coffee'

## Selectional Preferences: GermaNet Nodes (2)

Synset		Freq
Objekt	'object'	261.25
Nahrung, Lebensmittel, Esswaren, Speisen	'food'	127.98
festes Nahrungsmittel	'solid food'	100.28
Ding, Sache, Gegenstand, Gebilde	'thing'	66.24
Lebewesen, Kreatur, Wesen	'creature'	50.06
natürliches Lebewesen, Organismus	'organism'	49.14
Fleischware, Fleisch	'meat'	37.52
höheres Lebewesen	'higher creature'	34.51
Tier	'animal'	26.18
Backware	'pastry'	25.96
Gericht, Speise, Essen	'food'	22.36
Grünzeug	'vegetables' (coll.)	20.78

Selectional preferences for *essen* 'to eat' in na

## GermaNet Top Level Nodes

Lebewesen	'creature'
Sache	'thing'
Besitz	'property'
Substanz	'substance'
Nahrung	'food'
Mittel	'means'
Situation	'situation'
Zustand	'state'
Struktur	'structure'
Physis	'body'
Zeit	'time'
Ort	'space'
Attribut	'attribute'
Kognitives Objekt	'cognitive object'
Kognitiver Prozess	'cognitive process'



## Selectional Preferences: GermaNet Top Level Nodes

Synset	Freq	Prob
Situation	140.99	0.244
Kognitives Objekt	109.89	0.191
Zustand	81.35	0.141
Sache	61.30	0.106
Attribut	52.69	0.091
Lebewesen	46.56	0.081
Ort	45.96	0.080
Struktur	14.25	0.025
Kognitiver Prozess	11.77	0.020
Zeit	4.58	0.008
Besitz	2.86	0.005
Substanz	2.08	0.004
Nahrung	2.00	0.003
Physis	0.50	0.001

GermaNet arguments for *verfolgen* 'to follow' in na

## Selectional Preferences: GermaNet Top Level Nodes

Synset	Freq	Prob
Situation	1,102.26	0.425
Zustand	301.82	0.116
Zeit	256.64	0.099
Sache	222.13	0.086
Kognitives Objekt	148.12	0.057
Kognitiver Prozess	139.55	0.054
Ort	107.68	0.041
Attribut	101.47	0.039
Struktur	87.08	0.034
Lebewesen	81.34	0.031
Besitz	36.77	0.014
Physis	4.18	0.002
Substanz	3.70	0.001
Nahrung	3.29	0.001

GermaNet arguments for *beginnen* 'to begin' in n

## Verb Description

- D1** purely syntactic definition of subcategorisation
- D2** syntactico-semantic definition of subcategorisation with prepositional preferences
- D3** syntactico-semantic definition of subcategorisation with prepositional and selectional preferences

**$D1 \rightarrow D2 \rightarrow D3$** 

<i>D1</i>		<i>D2</i>		<i>D3</i>	
np	0.43	n	0.28	<u>n</u> (Situation)	0.12
n	0.28	np:um <sub>Akk</sub>	0.16	<u>np:um</u> <sub>Akk</sub> (Situation)	0.09
ni	0.09	ni	0.09	<u>np:mit</u> <sub>Dat</sub> (Situation)	0.04
na	0.07	np:mit <sub>Dat</sub>	0.08	<u>ni</u> (Lebewesen)	0.03
nd	0.04	na	0.07	<u>n</u> (Zustand)	0.03
nap	0.03	np:an <sub>Dat</sub>	0.06	<u>np:an</u> <sub>Dat</sub> (Situation)	0.03
nad	0.03	np:in <sub>Dat</sub>	0.06	<u>np:in</u> <sub>Dat</sub> (Situation)	0.03
nir	0.01	nd	0.04	<u>n</u> (Zeit)	0.03
ns-2	0.01	nad	0.02	<u>n</u> (Sache)	0.02
xp	0.01	np:nach <sub>Dat</sub>	0.01	<u>na</u> (Situation)	0.02

Example: *beginnen* 'to begin'

***D1* → *D2* → *D3***

<i>D1</i>		<i>D2</i>		<i>D3</i>	
na	0.42	na	0.42	<u>na</u> (Lebewesen)	0.33
n	0.26	n	0.26	<u>na</u> (Nahrung)	0.17
nad	0.10	nad	0.10	<u>na</u> (Sache)	0.09
np	0.06	nd	0.05	<u>n</u> (Lebewesen)	0.08
nd	0.05	ns-2	0.02	<u>na</u> (Lebewesen)	0.07
nap	0.04	np:auf <sub>Dat</sub>	0.02	<u>n</u> (Nahrung)	0.06
ns-2	0.02	ns-w	0.01	<u>n</u> (Sache)	0.04
ns-w	0.01	ni	0.01	<u>nd</u> (Lebewesen)	0.04
ni	0.01	np:mit <sub>Dat</sub>	0.01	<u>nd</u> (Nahrung)	0.02
nas-2	0.01	np:in <sub>Dat</sub>	0.01	<u>na</u> (Attribut)	0.02

Example: *essen* 'to eat'

## *D1* → *D2* → *D3*

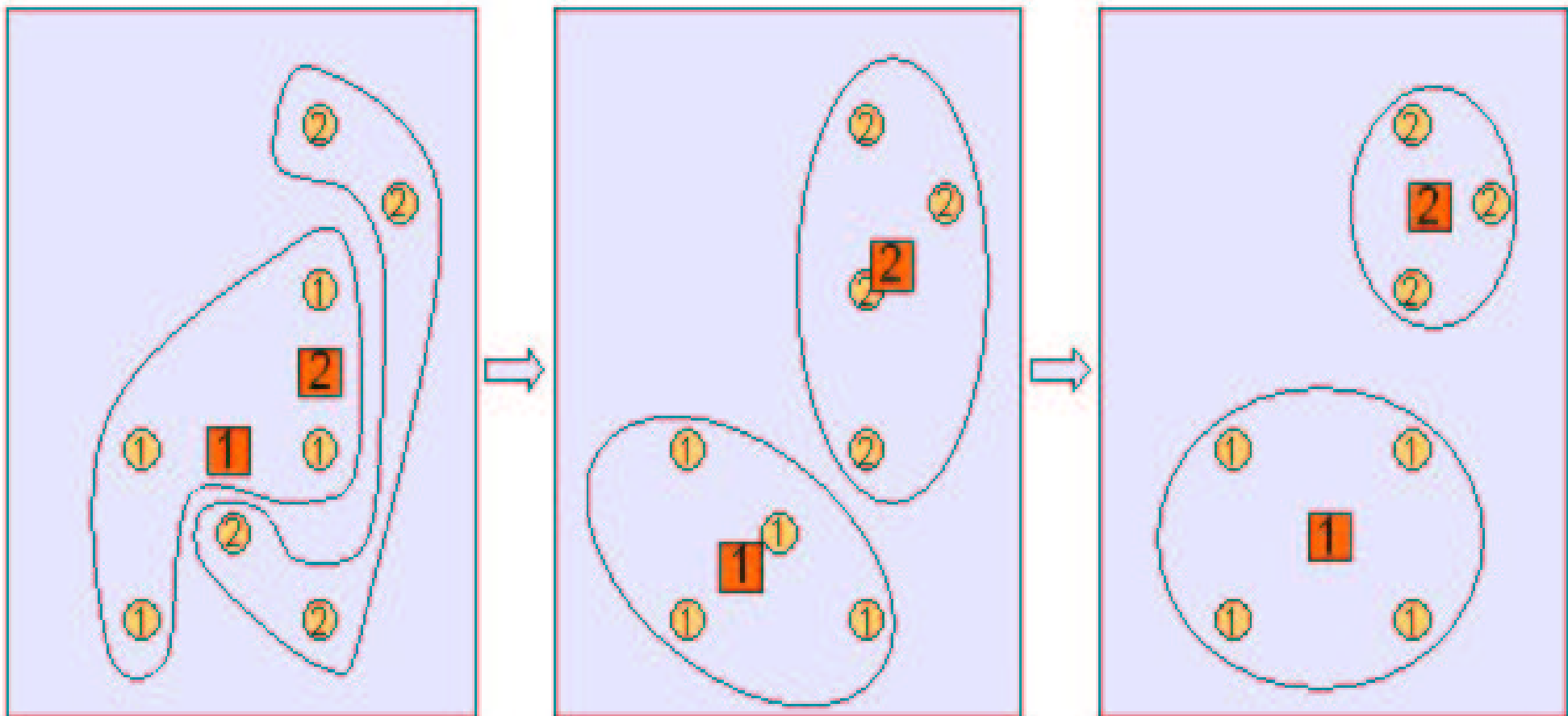
<i>D1</i>		<i>D2</i>		<i>D3</i>	
n	0.34	n	0.34	<u>n</u> (Sache)	0.12
np	0.29	na	0.19	<u>n</u> (Lebewesen)	0.10
na	0.19	np:in <sub>Akk</sub>	0.05	<u>na</u> (Lebewesen)	0.08
nap	0.06	nad	0.04	<u>na</u> (Sache)	0.06
nad	0.04	np:zu <sub>Dat</sub>	0.04	<u>n</u> (Ort)	0.06
nd	0.04	nd	0.04	<u>na</u> (Sache)	0.05
ni	0.01	np:nach <sub>Dat</sub>	0.04	<u>np:in<sub>Akk</sub></u> (Sache)	0.02
ns-2	0.01	np:mit <sub>Dat</sub>	0.03	<u>np:zu<sub>Dat</sub></u> (Sache)	0.02
ndp	0.01	np:in <sub>Dat</sub>	0.03	<u>np:in<sub>Akk</sub></u> (Lebewesen)	0.02
ns-w	0.01	np:auf <sub>Dat</sub>	0.02	<u>np:nach<sub>Dat</sub></u> (Sache)	0.02

Example: *fahren* 'to drive'

## k-Means Clustering

- k-Means algorithm (Forgy 1965)
- Unsupervised hard clustering:  $n$  objects  $\rightarrow k$  clusters
- Iterative re-organisation of cluster membership:
  1. Initial cluster assignment: agglomerative hierarchical clusters
  2. Calculation of cluster centroids
  3. Determining closest cluster (centroid)
  4. Re-arrangement of cluster membership
- Similarity measure: skew divergence
$$d(v_1, v_2) = D(p \parallel w * q + (1 - w) * p) \quad \text{with weight } w \text{ set to } 0.9$$
- Number of clusters: 43 (= manual)

## k-Means Clustering: Illustration





## Clustering Evaluation

**Gold standard** : hand-constructed verb classes

- Evaluation = Similarity between two partitions on verb set
- Evaluation measure: adjusted Rand index (Hubert/Arabie 1985):
  - Agreement between verb pairs in the partitions
  - Correction for chance in comparison to random partition
  - Typical range:  $0 \leq R_{adj} \leq 1$

$$R_{adj} = \frac{\sum_{i,j} \binom{t_{ij}}{2} - \frac{\sum_i \binom{t_{i\cdot}}{2} \sum_j \binom{t_{\cdot j}}{2}}{\binom{n}{2}}}{\frac{1}{2} (\sum_i \binom{t_{i\cdot}}{2} + \sum_j \binom{t_{\cdot j}}{2}) - \frac{\sum_i \binom{t_{i\cdot}}{2} \sum_j \binom{t_{\cdot j}}{2}}{\binom{n}{2}}} \quad (1)$$

## Clustering Results on *D1* and *D2*

Distribution		$R_{adj}$
<i>D1</i>		0.094
<i>D2</i>	$pp_{arg}$	0.151
	$pp_{chosen}$	0.151
	$pp_{all}$	0.160

## Clustering Results on Varying *D3*

Single Slots		Slot Combinations	
<u>n</u>	0.125	<u>na</u>	0.137
<u>na</u>	0.176	<u>n/na</u>	0.128
<u>na</u>	0.164	<u>nad</u>	0.088
<u>nad</u>	0.144	<u>n/na/nad</u>	0.118
<u>nad</u>	0.115	<u>nd</u>	0.150
<u>nad</u>	0.161	<u>n/na/nd</u>	0.124
<u>nd</u>	0.152	<u>n/na/nad/nd</u>	0.161
<u>nd</u>	0.143	<u>n/na/nd/nad/ns-dass</u>	0.182
<u>np</u>	0.133	<u>np/ni/nr/ns-2/ns-dass</u>	0.131
<u>ni</u>	0.148	all NP	0.158
<u>nr</u>	0.136	all NPs+PPs	0.176
<u>ns-2</u>	0.121		
<u>ns-dass</u>	0.156		

## Cluster Analysis (D3) – Examples

- (a) beginnen<sub>1</sub> bestehen<sub>37</sub> enden<sub>1</sub> existieren<sub>37</sub> laufen<sub>8</sub> liegen<sub>31</sub>  
sitzen<sub>31</sub> stehen<sub>31</sub>
- (b) eilen<sub>10</sub> gleiten<sub>12</sub> kriechen<sub>8</sub> rennen<sub>8</sub> starren<sub>16</sub>
- (c) fahren<sub>11</sub> fliegen<sub>11</sub> fließen<sub>12</sub> klettern<sub>8</sub> segeln<sub>11</sub> wandern<sub>8</sub>
- (d) bilden<sub>32</sub> erhöhen<sub>35</sub> festlegen<sub>22</sub> senken<sub>35</sub> steigern<sub>35</sub> vergrößern<sub>35</sub>  
verkleinern<sub>35</sub>
- (e) töten<sub>39</sub> unterrichten<sub>29</sub>
- (f) nieseln<sub>43</sub> regnen<sub>43</sub> schneien<sub>43</sub>
- (g) dämmern<sub>43</sub>

## Large-Scale Cluster Analysis (D3) – Examples ↓ →

- (a) *anhören* ‘to listen’, *auswirken* ‘to affect’, *einigen* ‘to agree’, *lohnen* ‘to be worth’, *verhalten* ‘to behave’, *wandeln* ‘to promenade’
- (b) *beschleunigen* ‘to speed up’, ***bilden***, *darstellen* ‘to illustrate’, *decken* ‘to cover’, *erfüllen* ‘to fulfil’, ***erhöhen*** ‘to raise’, *erledigen* ‘to fulfil’, *finanzieren* ‘to finance’, *füllen* ‘to fill’, *lösen* ‘to solve’, *rechtfertigen* ‘to justify’, ***reduzieren*** ‘to reduce’, ***senken*** ‘to lower’, ***steigern*** ‘to increase’, ***verbessern*** ‘to improve’, ***vergrößern*** ‘to enlarge’, ***verkleinern*** ‘to make smaller’, ***verringern*** ‘to decrease’, ***verschieben*** ‘to shift’, ***verschärfen*** ‘to intensify’, ***verstärken*** ‘to intensify’, ***verändern*** ‘to change’

## Large-Scale Cluster Analysis (D3) – Examples ↑

- (a) **ahnen** ‘to guess’, **bedauern** ‘to regret’, **befürchten** ‘to fear’, **bezweifeln** ‘to doubt’, **merken** ‘to notice’, **vermuten** ‘to assume’, **weißen** ‘to whiten’, **wissen** ‘to know’
- (b) **basieren** ‘to be based on’, **beruhen** ‘to be based on’, **resultieren** ‘to result from’, **stammen** ‘to stem from’
- (c) **befragen** ‘to interrogate’, **entlassen** ‘to release’, **ermorden** ‘to assassinate’, **erschießen** ‘to shoot’, **festnehmen** ‘to arrest’, **töten** ‘to kill’, **verhaften** ‘to arrest’
- (d) **beziffern** ‘to amount to’, **schätzen** ‘to estimate’, **veranschlagen** ‘to estimate’

## Discussion

- Clustering methodology for automatic basis of large-scale verb classification
- Step-wise refinement of features improves clustering
- Linguistic intuition and clustering results do not necessarily align
- Difficulty of selecting and encoding verb features  
(→ selectional preferences)
- Idiosyncratic properties of verbs
- Specific properties of desired verb classification

## Verb Class Usage

- Redundancy reduction in verb descriptions by generalising over the common properties of verbs
  - Prediction of properties and refinement of vague properties
- Possible NLP applications: parsing, language modelling, machine translation, information extraction



## Verb Classes for NLP Applications

- Machine translation (Dorr 1997)
- Document classification (Klavans and Kan 1998)
- Word sense disambiguation  
(Dorr and Jones 1996, Prescher et al. 2000)
- Subcategorisation acquisition and filtering (Korhonen 2002)