

Linguistische Informatik

Linguistische Ebenen, Subsprachen

Gerhard Heyer
Universität Leipzig
heyer@informatik.uni-leipzig.de

UNIVERSITÄT LEIPZIG

Institut für Informatik



Automatische Sprachverarbeitung

Beispiel einer Einzelsprache

Um welche Sprache handelt es sich?

Was fällt Ihnen auf?

Keväällä 1992 Nokia vahvisti asemaansa kulutuselektronikka-markkinoilla hankkimalla omistukseensa televisioita valmistavan Finluxin. Nokian tietoliikennekaapeliteollisuus vahvistui, kun yhtymään kuuluva hollantilainen NKF osti saksalaisen Philips Kommunikations Industrien kaapeliliiketoiminnan elokuussa 1993.

... ..

Was sehen wir?

- **Buchstaben**
- **Buchstabenkombinationen**
- **Wortformen**
- **Sätze**
-

Keväällä 1992 Nokia vahvisti asemaansa kulutuselektronikka-markkinoilla hankkimalla omistukseensa televisioita valmistavan Finluxin. Nokian tietoliikennekaapeliteollisuus vahvistui, kun yhtymään kuuluva hollantilainen NKF osti saksalaisen Philips Kommunikations Industrien kaapeliliiketoiminnan elokuussa 1993.

... ..

Was sehen wir?

- **Buchstaben**
- **Buchstabenkombinationen**
- **Wortformen**
- **Sätze**
-

- **Einzelne Wortformen sind vermutlich Eigennamen**

Keväällä 1992 **Nokia** vahvisti asemaansa kulutuselektronikka-markkinoilla hankkimalla omistukseensa televisioita valmistavan **Finluxin**. **Nokian** tietoliikennekaapeliteollisuus vahvistui, kun yhtymään kuuluva hollantilainen **NKF** osti saksalaisen **Philips Kommunikations Industrien** kaapeliliiketoiminnan elokuussa 1993.

... ..

Was sehen wir?

- Buchstaben
- Buchstabenkombinationen
- Wortformen
- Sätze
-

- Einzelne Wortformen sind vermutlich Eigennamen
- Außerdem erkennen wir Zahlen, vermutlich Jahreszahlen

Keväällä **1992** **Nokia** vahvisti asemaansa kulutuselektronikka-markkinoilla hankkimalla omistukseensa televisioita valmistavan **Finluxin**. **Nokian** tietoliikennekaapeliteollisuus vahvistui, kun yhtymään kuuluva hollantilainen **NKF** osti saksalaisen **Philips Kommunikations Industrien** kaapeliliiketoiminnan elokuussa **1993**.

... ..

Linguistische Ebenen

Explanandum

Laute (tokens)

Lautgruppen

*Phonem: kleinste bedeutungs-
unterscheidende Einheit*

Gruppen von Phonemen

*Morphem: kleinste bedeutungstragende
Einheit*

*Allomorphe: bedeutungsäquivalente
Morpheme*

Bsp.: sprech={sprech, sprich, sprach, ...}

Explanans

Phonetik

Phonologie

Morphologie

Linguistische Ebenen

Gruppen von Morphemen

Wortformen: flektierte Formen eines Wortes

Wort: Äquivalenzklasse v. Wortformen

Gruppen von Wörtern

Phrasen: zulässige Kombination von Wortformen

Sätze: grammatisch vollständige

Sequenz von Phrasen

Aussagen: wahrheitsfähige Sätze

Sprechakte: zustandsverändernd

Lexikon

Syntax

Semantik /
Pragmatik

Beispiel für linguistische Ebenen

Buchstaben oder Phoneme

E i n k u r z e r B e i s p i e l s a t z w ü r d e h e l f e n

Morphe und Morpheme

Ein kurz er Bei spiel satz würd e helf en.

Komposition

i, e, u

Wortformen und Wörter

wird, werd, ...

Abstraktion

Ein kurzer Beispielsatz würde helfen.

Phrasen, Sätze

wurde, werde, wird, ...

Ein kurzer Beispielsatz würde helfen. Denn ohne Beispiele ist
alles ein wenig Abstrakt. ...

Ein Beispielsatz wird helfen.
Ein Beispielsatz wurde selten ignoriert.

Absätze

Texte

...

Merkmale für die ASV im allgemeinen

- **Features im Sinne des maschinellen Lernens finden sich auf allen linguistischen Ebenen**
- **Für Anwendungen der Sprachverarbeitung sind neben den rein „linguistischen“ features oftmals aber auch statistische oder Muster basierte features nützlich**

Beispiele

- **Uni-, Bi-, Tri-, n-Gramme**
- **(gewichtete) Wort- und/oder n-Gramm-Frequenzen**
- **Kookkurrenzen**
- **Zeitstempel und andere Metadaten**

Grundlegende Definitionen

1. Alphabet

Sei NL eine natürliche Sprache und sei A eine Menge von Zeichen, $A = \{l_1, l_2, \dots, l_k\}$. Wir nennen A ein Alphabet von NL der Größe k .

Beispiel: $A_E = \{a, b, \dots, z\}$ $k_E = 26$

2. Zeichenkette

Seien l_1, l_2, \dots, l_n Buchstaben aus A . Das Tupel t

mit $t = \langle l_1, l_2, \dots, l_n \rangle$ wird Zeichenkette genannt und n ist die Länge von t .

3. Menge von Zeichenketten

Sei A^n das kartesische Produkt des Alphabets A . A^n wird Menge von Zeichenketten der Länge n genannt.

Beispiel: $A^3 = \{$
 $(a,a,a), (a,a,b), \dots (a,a,z),$
 $(b,a,a), (b,a,b), \dots (b,a,z),$
 .
 .
 $(z,z,z) \}$

Definitionen 2

4. Lexikon einer Sprache

Sei NL eine natürliche Sprache und L eine Teilmenge von A^+ ($A^+ = \cup_{n>0} A^n$). Wir nennen $L \subseteq A^+$ ein **Lexikon** von NL.

5. Wortform, Menge von Wortformen der Länge n

Jedes Element W des Lexikons L wird eine (natürliche) Wortform genannt. W^n ist die Schnittmenge von A^n mit L und wird Menge von Wortformen der Länge n genannt.

6. token

Vorkommen einer Zeichenkette (Wortform) in einem Text.

(Gesamtzahl der tokens in einem Text = Textumfang)

7. type

Äquivalenzklasse gleicher Zeichenketten (Wortformen) in einem Text. (Gesamtzahl der types in einem Text = Vokabularumfang)

Wie viele Wörter?

- *Ich lese ab und zu gerne Zeitung, aber ich habe keine Zeit zu lesen.*
- **Wortform**
 - flektierte Form, wie sie im Text vorkommt
 - *lese* and *lesen* ... sind verschiedene Wortformen
- **Lemma**
 - Wortformen mit demselben Stamm, gleiche Wortkategorie und Bedeutung
 - *lese* und *lesen* ... haben dasselbe Lemma (“lesen”, Stamm “les-”)
- **token:**
 - Tatsächliches Vorkommen einer Wortform
 - *Ich lese ab und zu gerne Zeitung, aber ich habe keine Zeit zu lesen.*
 - *14 tokens* (ohne Punctuation)
- **type:**
 - Muster eines tokens
 - *Ich lese ab und zu gerne Zeitung, aber ich habe keine Zeit zu lesen.*
 - *13 types* (ohne Punctuation)

Tokenisierung

- **Zerlegung von Zeichenketten (einer Sprache!) in Wortformen**
- **Benötigt werden Festlegungen zu**
 - Sonderzeichen
 - Zeichensetzung
 - Worttrennung
 - Wortkombinationen
- **NLTK Tokenisierung**
 - <http://nltk.googlecode.com/svn/trunk/doc/api/nltk.tokenize-module.html>

Ein nützliches Maß: token - type Quotient

- **Mark Twain's *Tom Sawyer***
 - 71,370 tokens
 - 8,018 word types
 - tokens/type Verhältnis = 8.9
- **Alle Werke Shakespeares**
 - 884,647 word tokens
 - 29,066 word types
 - tokens/type Verhältnis = 30.4
- **Wie ist dieses Maß zu interpretieren?**

Definitionen 3

8. *Trigramme einer Wortform*

Sei $t \in A^+$ mit $t = (l_1, l_2 \dots l_n)$, $0 =$ leeres Element.

Die Menge T von Trigrammen aus t ist die Menge von 3-Tupeln, so dass

$$T = \{ \langle 0, 0, l_1 \rangle, \langle 0, l_1, l_2 \rangle, \langle l_1, l_2, l_3 \rangle, \langle l_2, l_3, l_4 \rangle, \dots \\ \dots \langle l_{n-2}, l_{n-1}, l_n \rangle, \langle l_{n-1}, l_n, 0 \rangle, \langle t_n, 0, 0 \rangle \}$$

Bigrammhäufigkeiten (Auszug)

<i>en</i>	4.47	<i>nw</i>	0.55	<i>ab</i>	0.25	<i>nm</i>	0.16
<i>er</i>	3.40	<i>us</i>	0.54	<i>il</i>	0.25	<i>pe</i>	0.16
<i>ch</i>	2.80	<i>nn</i>	0.53	<i>mm</i>	0.25	<i>rl</i>	0.16
<i>nd</i>	2.58	<i>nt</i>	0.52	<i>nz</i>	0.25	<i>sm</i>	0.16
<i>ei</i>	2.26	<i>ta</i>	0.51	<i>sg</i>	0.25	<i>sp</i>	0.16
<i>de</i>	2.14	<i>eg</i>	0.50	<i>sw</i>	0.25	<i>th</i>	0.16
<i>in</i>	2.04	<i>eh</i>	0.50	<i>rn</i>	0.24	<i>wo</i>	0.16
<i>es</i>	1.81	<i>zu</i>	0.50	<i>ro</i>	0.24	<i>af</i>	0.15
<i>te</i>	1.78	<i>al</i>	0.49	<i>ea</i>	0.23	<i>lu</i>	0.15
<i>ie</i>	1.76	<i>ed</i>	0.48	<i>fr</i>	0.23	<i>mu</i>	0.15
<i>un</i>	1.73	<i>ru</i>	0.48	<i>sd</i>	0.23	<i>no</i>	0.15
<i>ge</i>	1.68	<i>rs</i>	0.47	<i>tt</i>	0.23	<i>nv</i>	0.15
<i>st</i>	1.24	<i>ig</i>	0.45	<i>tw</i>	0.23	<i>rf</i>	0.15
<i>ic</i>	1.19	<i>ts</i>	0.45	<i>gr</i>	0.22	<i>ut</i>	0.15
<i>he</i>	1.17	<i>ma</i>	0.43	<i>tz</i>	0.22	<i>br</i>	0.14
<i>ne</i>	1.17	<i>sa</i>	0.43	<i>fe</i>	0.21	<i>ez</i>	0.14
<i>se</i>	1.17	<i>wa</i>	0.43	<i>gt</i>	0.21	<i>ho</i>	0.14
<i>ng</i>	1.07	<i>ac</i>	0.42	<i>rh</i>	0.21	<i>ka</i>	0.14
<i>re</i>	1.07	<i>eu</i>	0.42	<i>ds</i>	0.20	<i>os</i>	0.14
<i>au</i>	1.04	<i>so</i>	0.41	<i>du</i>	0.20	<i>bl</i>	0.13
<i>di</i>	1.02	<i>ar</i>	0.40	<i>mi</i>	0.20	<i>dw</i>	0.13
<i>be</i>	0.96	<i>tu</i>	0.40	<i>nb</i>	0.20	<i>ep</i>	0.13
<i>ss</i>	0.94	<i>ck</i>	0.37	<i>nk</i>	0.20	<i>hm</i>	0.13
<i>ns</i>	0.93	<i>or</i>	0.37	<i>rk</i>	0.20	<i>hw</i>	0.13
<i>an</i>	0.92	<i>rt</i>	0.36	<i>rz</i>	0.20	<i>pr</i>	0.13

Trigrammhäufigkeiten (Auszug)

<i>ein</i>	1.22	<i>ese</i>	0.27	<i>hre</i>	0.18	<i>nne</i>	0.14	<i>auc</i>	0.11
<i>ich</i>	1.11	<i>auf</i>	0.26	<i>hei</i>	0.18	<i>nes</i>	0.14	<i>als</i>	0.11
<i>nde</i>	0.89	<i>ben</i>	0.26	<i>lei</i>	0.18	<i>ond</i>	0.14	<i>alt</i>	0.11
<i>die</i>	0.87	<i>ber</i>	0.26	<i>nei</i>	0.18	<i>oen</i>	0.14	<i>eic</i>	0.11
<i>und</i>	0.87	<i>eit</i>	0.26	<i>nau</i>	0.18	<i>sdi</i>	0.14	<i>esc</i>	0.11
<i>der</i>	0.86	<i>ent</i>	0.26	<i>sge</i>	0.18	<i>sun</i>	0.14	<i>enh</i>	0.11
<i>che</i>	0.75	<i>est</i>	0.26	<i>tte</i>	0.18	<i>von</i>	0.14	<i>eil</i>	0.11
<i>end</i>	0.75	<i>sei</i>	0.26	<i>wei</i>	0.18	<i>bei</i>	0.13	<i>fen</i>	0.11
<i>gen</i>	0.71	<i>and</i>	0.25	<i>abe</i>	0.17	<i>chl</i>	0.13	<i>gan</i>	0.11
<i>sch</i>	0.66	<i>ess</i>	0.25	<i>chd</i>	0.17	<i>chn</i>	0.13	<i>hte</i>	0.11
<i>cht</i>	0.61	<i>ann</i>	0.24	<i>des</i>	0.17	<i>chw</i>	0.13	<i>iea</i>	0.11
<i>den</i>	0.57	<i>esi</i>	0.24	<i>nte</i>	0.17	<i>ech</i>	0.13	<i>ieb</i>	0.11
<i>ine</i>	0.53	<i>ges</i>	0.24	<i>rge</i>	0.17	<i>edi</i>	0.13	<i>nli</i>	0.11
<i>nge</i>	0.52	<i>nsc</i>	0.24	<i>tes</i>	0.17	<i>enk</i>	0.13	<i>rda</i>	0.11
<i>nun</i>	0.48	<i>nwi</i>	0.24	<i>uns</i>	0.17	<i>eun</i>	0.13	<i>rsc</i>	0.11
<i>ung</i>	0.48	<i>tei</i>	0.24	<i>vor</i>	0.17	<i>enz</i>	0.13	<i>std</i>	0.11
<i>das</i>	0.47	<i>eni</i>	0.23	<i>dem</i>	0.16	<i>hau</i>	0.13	<i>sst</i>	0.11
<i>hen</i>	0.47	<i>ige</i>	0.23	<i>hin</i>	0.16	<i>ite</i>	0.13	<i>tre</i>	0.11
<i>ind</i>	0.46	<i>aen</i>	0.22	<i>her</i>	0.16	<i>ief</i>	0.13	<i>uss</i>	0.11
<i>enw</i>	0.45	<i>era</i>	0.22	<i>lle</i>	0.16	<i>imm</i>	0.13	<i>all</i>	0.10
<i>ens</i>	0.44	<i>ern</i>	0.22	<i>nan</i>	0.16	<i>ihr</i>	0.13	<i>aft</i>	0.10
<i>ies</i>	0.44	<i>rde</i>	0.22	<i>tda</i>	0.16	<i>iss</i>	0.13	<i>bes</i>	0.10
<i>ste</i>	0.44	<i>ren</i>	0.22	<i>tel</i>	0.16	<i>kei</i>	0.13	<i>dei</i>	0.10
<i>ten</i>	0.44	<i>tun</i>	0.22	<i>ueb</i>	0.16	<i>mei</i>	0.13	<i>erf</i>	0.10
<i>ere</i>	0.43	<i>ing</i>	0.21	<i>ang</i>	0.15	<i>nsi</i>	0.13	<i>ess</i>	0.10
<i>lic</i>	0.42	<i>sta</i>	0.21	<i>cha</i>	0.15	<i>nem</i>	0.13	<i>esw</i>	0.10

„künstliche Sprache“ (nach Kupfmüller)

Einergruppen (Buchstabenhäufigkeit)

EME GKNEET ERS TITBL BTZENFNDGBGD EAI E LASZ
BETEATR IASMIRCH EGEOM

Zweiergruppen (Paarhäufigkeit)

AUSZ KEINU WONDINGLIN DUFNRN ISAR STEISBERER ITEHM
ANORER

Dreiergruppen

PLANZEUNDGES PHIN INE UNDEN ÜBBEICHT GES AUF ES SO
UNG GAN DICH WANDERSO

Vierergruppen

ICH FOLGEMÄSZIG BIS STEHEN DISPONIN SEELE NAMEN

Stringähnlichkeiten von Wortformen

- **N-Gramm-Ähnlichkeit oftmals interessantes *feature***
- **Nützlich für u.a.**
 - **Rechtschreibprüfung**
 - **Zitations- und Plagiatserkennung**
 - **(semi-)automatisches Lernen von sprachlichen Strukturen**
- **Verfahren u.a.**
 - **Anzahl gleicher Trigramme, z. B. Dice-Koeffizient**

$$d_w(a, b) = \frac{2 |T(a) \cap T(b)|}{|T(a)| + |T(b)|}$$

- **Editierdistanz, Kosten für die Transformation einer Zeichenkette in eine andere, z. B. Levenshtein Matrix**

Levenshtein Matrix

$$d(\alpha, \beta) = \phi + \min_{i, j, k, l, a, b} \left\{ \begin{array}{l} \alpha = \beta \\ d(\text{sub}(\alpha, i, a), \beta) \quad 1 \leq i \leq |\alpha| \\ d(\text{del}(\alpha, j), \beta) \quad 1 \leq j \leq |\alpha| \\ d(\text{ins}(\alpha, k, b), \beta) \quad 0 \leq k \leq |\alpha| + 1 \\ d(\text{tra}(\alpha, l), \beta) \quad 1 \leq l \leq |\alpha| - 1 \end{array} \right.$$

1.) $d(\alpha, \beta) = \phi$ dann und nur dann wenn $\alpha = \beta$, sonst $d(\alpha, \beta) > \phi$

2.) $d(\alpha, \beta) = d(\beta, \alpha)$

3.) $d(\alpha, \beta) + d(\beta, \gamma) \geq d(\alpha, \gamma)$

4.) Maß ist Länge

Beispiel

$$d(\text{abc}, \text{axc}) = 1 + (\text{sub}(\text{abc}, 2, \text{x}), \text{axc}) = 1 + d(\text{axc}, \text{axc}) = 1$$

Definitionen 4

9. Teilzeichenkette

Seien $t, u \in A^+$ Zeichenketten mit

$t = (l_1, l_2, \dots, l_n)$ und

$u = (u_1, u_2, \dots, u_m)$.

Wir nennen u eine Teilzeichenkette von t ,

falls $1 \leq m \leq n$ und es gibt ein i und ein j so daß $u = l_i$ bis
für alle $1 \leq j \leq m$,

Beispiel: „satis“ ist Teilzeichenkette von „Mensatisch“.

Me n **s a t i s** c h **s a t i s**

i= 1 2 3 **4 5 6 7 8** 9 10 j= 1 2 3 4 5

Definitionen 5

8. Wortformkombinationen der Länge r

Sei L ein Tupel von Wortformen, $L=(W_1, W_2, \dots, W_r)$ mit $W_i \in L$. Wir nennen L eine Wortformkombination der Länge r .

9. Menge von Wortformkombinationen

Sei L^r das kartesische Produkt von L . L^+ wird Menge von Wortformkombinationen der Länge r genannt. ($L^+ = \cup_{r>0} L^r$)

10. Menge von Sätzen

SYN sei eine Menge von syntaktischen Restriktionen. Die Menge S , mit $S \leq L^+$, die SYN folgen, wird Menge von Sätzen genannt. (Ableitung !)

Definitionen 6

11. Wort

Äquivalenzklasse von morphologisch zusammengehörigen Wortformen.

12. Konzept

**Äquivalenzklasse von semantisch zusammengehörigen Wörtern
(z. B. globaler Kontext)**

Akquisition Linguistischen Wissens

Gegeben ein Möglichkeitsraum von N-Grammen (Buchstaben und Wortformen) einer Sprache:

- **Bestimme das Lexikon der Wortformen**
- **Bestimme das Lexikon der Wörter**
- **Bestimme die Menge der syntaktischen Restriktionen SYN**

Was heisst hier *Sprache* ?

Einzel sprache vs. Sprachfähigkeit

Einzel sprache: Deutsch, Englisch, Finnisch, ...

Empirisch finden wir *Realisierungen* (Äußerungen, Texte) von Einzel sprachen. Menschen besitzen *Kenntnisse* einer (oder mehrerer) Einzel sprachen.

Sprachfähigkeit: Für den Menschen spezifische Möglichkeit, eine Einzel sprache zu lernen (verstehen, anwenden).

Sprache i.S. von Sprachfähigkeit: abstraktes *System*

Einzel Sprachen und ihre Realisierungen

Einzel Sprachen begegnen uns in einer Vielzahl von Realisierungen:

- **Alltagssprache**
- **Zeitungen**
- **wissenschaftliche Aufsätze, Fachbücher, Lexika**
- **Internet**
- **Korrespondenzen und email**
- **technische Dokumentationen, Benutzerhandbücher und Produktbeschreibungen**
- **Normen, Gesetze, Kommentare und Verträge**
- **Organisationsanweisungen**
- **usw.**

Subsprachen

Z. Harris (1968): Untermenge der vom Sprachsystem erzeugbaren Strukturen

1. *syntaktische und semantische Beschränkungen*

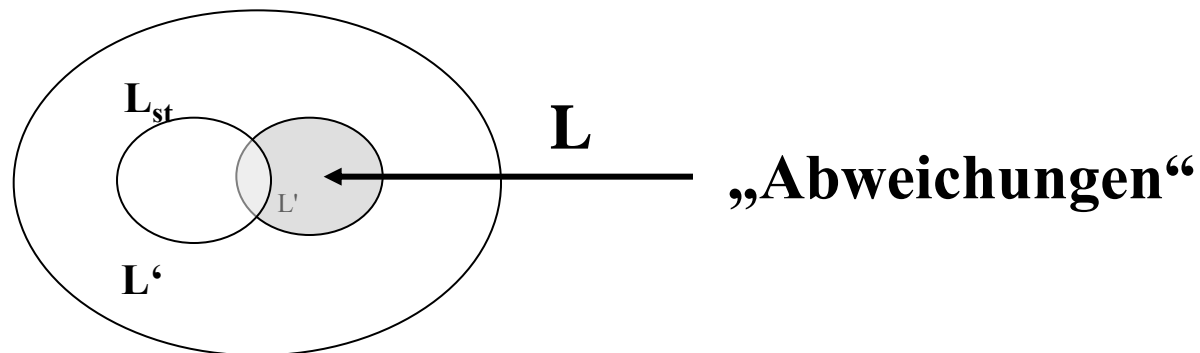
- *abweichende Grammatik*
- *hohe Wahrscheinlichkeit bestimmter Konstruktionen*

2. *lexikalische Beschränkungen*

- *Medizin, Wetterberichte, juristische Texte, technische Anleitungen ...*

3. *charakteristische Morpheme*

- *Medizin, Chemie, technische Anleitungen*



Beispiel Wetterbericht

Mit Tief „Fritz“ über der Biskaya gelangt feuchte, aber allmählich mildere Luft nach Mitteleuropa. In Spanien und Frankreich bringt das Tief kräftige Regengüsse. Über dem Baltikum und in weiten Teilen von Skandinavien sorgt Hoch „Birgit“ für trockenes, aber kaltes Winterwetter. Im östlichen Mittelmeerraum liegt ein weiteres Tiefdruckgebiet mit Regen.

Beispiel Wetterbericht

- **Syntaktische und semantische Beschränkungen**
- **Lexikalische Beschränkungen**

Mit Tief „Fritz“ über der Biskaya gelangt **feuchte, aber** allmählich mildere Luft nach Mitteleuropa. In Spanien und Frankreich bringt das Tief kräftige Regengüsse. Über dem Baltikum und in weiten Teilen von Skandinavien sorgt Hoch „Birgit“ für **trockenes, aber** kaltes Winterwetter. Im östlichen Mittelmeerraum liegt ein weiteres Tiefdruckgebiet mit Regen.

Beispiel Wetterbericht

- **Syntaktische und semantische Beschränkungen**
- **Lexikalische Beschränkungen**

Mit **Tief „Fritz“** über der Biskaya gelangt *feuchte, aber allmählich mildere* Luft nach Mitteleuropa. In Spanien und Frankreich bringt das Tief kräftige Regengüsse. Über dem Baltikum und in weiten Teilen von Skandinavien sorgt **Hoch „Birgit“** für *trockenes, aber kaltes* Winterwetter. Im östlichen Mittelmeerraum liegt ein weiteres Tiefdruckgebiet mit Regen.

Beispiel Wetterbericht

- **Syntaktische und semantische Beschränkungen**
- **Lexikalische Beschränkungen**

Mit **Tief** „**Fritz**“ über der **Biskaya** gelangt *feuchte, aber allmählich mildere* Luft nach **Mitteleuropa**. In **Spanien und Frankreich** bringt das **Tief** kräftige Regengüsse. Über dem **Baltikum** und in weiten Teilen von **Skandinavien** sorgt **Hoch** „**Birgit**“ für *trockenes, aber kaltes* Winterwetter. Im **östlichen Mittelmeerraum** liegt ein weiteres Tiefdruckgebiet mit **Regen**.

Eigenschaften von Subsprachen

1. SL bilden thematische Gruppen

konstante Lexika

konstante Syntax

2. SL (= Sublanguage)-Merkmale können über verschiedene Sprachen hinweg identisch bleiben.

Bsp.: a) Passiv in techn. Anleitung

b) Frequenzen von Satzbauplänen und Begriffen

c) Auslassen des Artikels

3. SL Merkmale sind gradierbar und verändern sich

N. Sager (NY), "LinguisticString Project", 1967

Language registers

***Language registers* refer to linguistic variations that correlate with different *situations* of language use**

Different situations

- 1. channel factors (e.g. mobile phone, Internet, telephone,...)
restricting the mode and quantity of linguistic action**
- 2. different social roles requiring specific modes of linguistic action (e.g. politeness)**

Linguistic variations

- 1. constraints**
- 2. amendmends**
- 3. deviant frequencies**
 - *syntax*
 - *lexicon*

Example – Language register *Computer Talk*

- **When people communicate with a computer system they adapt their language register assuming that the system will require „simpler“ input**
- **Computer talk comparable to *baby talk* or *foreigner talk***
- **Empirical validation (Krause et. al.) using simulations of 4 systems:**

System 1: optimal **human-human** information system

System 2: optimal **human-computer** IS

System 3: *restricted* **human-computer** IS

System 4: *strongly restricted* **human-computer** IS

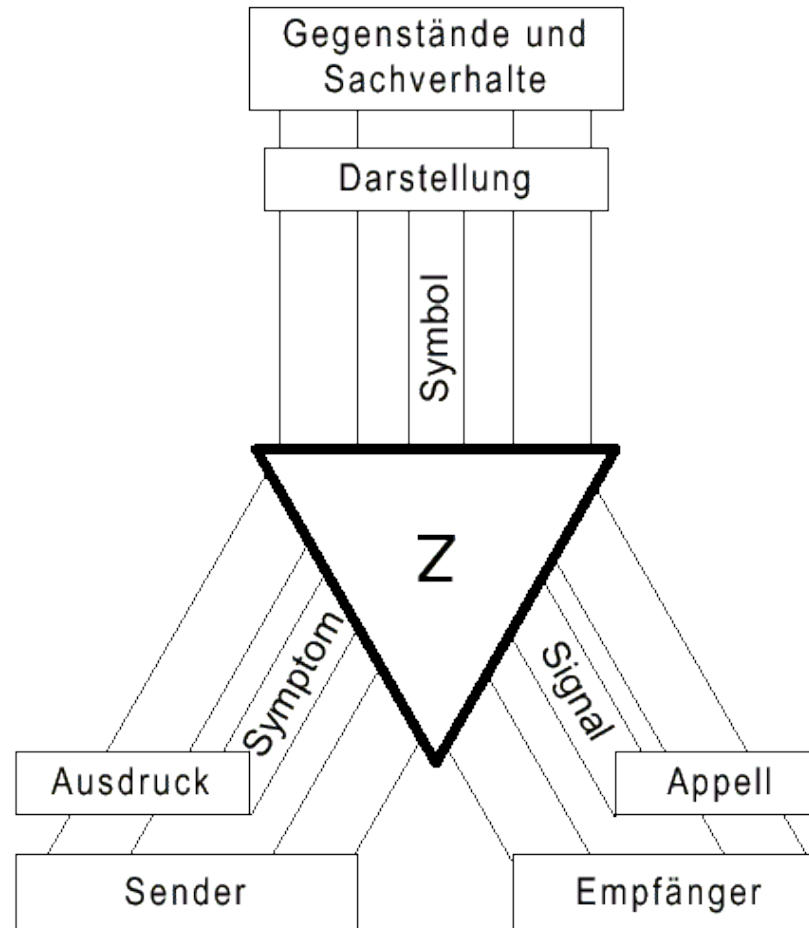
Computer Talk – Share of 10 most frequent sentence patterns at total input

	Typed input	Spoken input
S1	63,04 %	35,75 %
S2	88,57 %	77,12 %
S3	97,41 %	82,14 %
S4	99,26 %	99,24 %

Das Organonmodell von Bühler (Sprachtheorie)

- **Grundlage sprachlicher Kommunikation sind sprachliche Ausdrücke**
- **Sprachliche Ausdrücke haben drei Dimensionen:**
 - den Sender (Sprecher, Schreiber)
 - den Empfänger (Hörer, Leser)
 - die referenzierte Sache (Objekte und Ereignisse, Eigenschaften, Tatsachen, ...)
- **In Bezug auf einen Sender, (intendierten) Empfänger und die referenzierte Sache haben sprachliche Ausdrücke daher eine dreifache Funktion:**
 - Symptom
 - Apell
 - Symbol

Organonmodell – *Schema*



Linguistische Aspekte

- **Wir finden daher sprachliche Ausdrücke in Form von**

- Ausrufen (Symptom)
- Bewertungen (Apell)
- Aussagen (Symbol)

- **Beispiel**

ECKHARD BERGER (stockend)

Ich weiß nicht, was passieren wird... aber ich habe Angst... Angst vor meinen Kollegen: Jürgen Wiesehöfer... Michael Nauen... und Sven Lienecke. Wenn mir etwas zustößt, dann... (eine quälende Pause, dann) diese drei Männer sind gefährlich... (leise) möglicherweise Mörder.

**Drehbuch SoKo
Leipzig Folge 6,
2004**

Literaturempfehlung

Heyer, Quasthoff, Wittig, Text Mining - Wissensrohstoff
Text, W3L-Verlag, Bochum 2006

Krause, Hitzenberger, Computer Talk, Olms 1992