

Komplexität und natürliche Sprache

Einordnung natürlicher Sprache in die Chomsky-Hierarchie

Timm Lichte & Christian Wurm

HHU Düsseldorf

SS 2017, 03.05.2015



Vier Formen der Komplexität:

- **Extensionale Komplexität**
- Algorithmische Komplexität / Verarbeitungskomplexität
- Beschreibungskomplexität
- Informatische Komplexität

Wir betrachten Stringsprachen:

$$a \triangleleft \bar{w} \in L_G \in \{L_G \mid G \in \hat{G}\}$$

Wir betrachten Stringsprachen:

$$a \prec \bar{w} \in L_G \in \{L_G | G \in \hat{G}\}$$

a ist ein **Buchstabe** und Element eines Alphabets Σ .

\bar{w} ist ein **Wort** und Element einer Sprache L .

L_G ist die **Sprache** der Grammatik G .

\hat{G} ist eine Menge/Klasse von Grammatiken.

Wir betrachten Stringsprachen:

$$a \prec \bar{w} \in L_G \in \{L_G | G \in \hat{G}\}$$

a ist ein **Buchstabe** und Element eines Alphabets Σ .

\bar{w} ist ein **Wort** und Element einer Sprache L .

L_G ist die **Sprache** der Grammatik G .

\hat{G} ist eine Menge/Klasse von Grammatiken.

Beispiele:

$$\Sigma = \{a, b\}$$

Wir betrachten Stringsprachen:

$$a \triangleleft \bar{w} \in L_G \in \{L_G | G \in \hat{G}\}$$

a ist ein **Buchstabe** und Element eines Alphabets Σ .

\bar{w} ist ein **Wort** und Element einer Sprache L .

L_G ist die **Sprache** der Grammatik G .

\hat{G} ist eine Menge/Klasse von Grammatiken.

Beispiele:

$$\Sigma = \{a, b\}$$

$$\bar{w} = abba$$

Wir betrachten Stringsprachen:

$$a \triangleleft \bar{w} \in L_G \in \{L_G | G \in \hat{G}\}$$

a ist ein **Buchstabe** und Element eines Alphabets Σ .

\bar{w} ist ein **Wort** und Element einer Sprache L .

L_G ist die **Sprache** der Grammatik G .

\hat{G} ist eine Menge/Klasse von Grammatiken.

Beispiele:

$$\Sigma = \{a, b\}$$

$$\bar{w} = abba$$

$$L = \{abba, abab, bbb, \dots\}$$

Wir betrachten Stringsprachen:

$$a \prec \bar{w} \in L_G \in \{L_G | G \in \hat{G}\}$$

a ist ein **Buchstabe** und Element eines Alphabets Σ .

\bar{w} ist ein **Wort** und Element einer Sprache L .

L_G ist die **Sprache** der Grammatik G .

\hat{G} ist eine Menge/Klasse von Grammatiken.

Beispiele:

$$\Sigma = \{a, b\}$$

$$\bar{w} = abba$$

$$L = \{abba, abab, bbb, \dots\}$$

$$G = \langle \mathcal{N}, \Sigma, S, \mathcal{P} \rangle \quad (\text{z.B. eine CFG})$$

Wir betrachten Stringsprachen:

$$a \triangleleft \bar{w} \in L_G \in \{L_G | G \in \hat{G}\}$$

a ist ein **Buchstabe** und Element eines Alphabets Σ .

\bar{w} ist ein **Wort** und Element einer Sprache L .

L_G ist die **Sprache** der Grammatik G .

\hat{G} ist eine Menge/Klasse von Grammatiken.

Beispiele:

$$\Sigma = \{a, b\}$$

$$\bar{w} = abba$$

$$L = \{abba, abab, bbb, \dots\}$$

$$G = \langle \mathcal{N}, \Sigma, S, \mathcal{P} \rangle \quad (\text{z.B. eine CFG})$$

$$L_G = \{\bar{w}' | S \vdash_{\mathcal{P}_G}^* \bar{w}'\}$$

Intuition

Je komplexer desto umfangreicher.

Intuition

Je komplexer desto umfangreicher.

Gegeben zwei Grammatikklassen \hat{G} und \hat{G}' :

$K(\hat{G}) > K(\hat{G}') \iff$ Für jede $G' \in \hat{G}'$ gibt es eine $G \in \hat{G}$ und $\mathcal{P}_{G'} = \mathcal{P}_G$

Intuition

Je komplexer desto umfangreicher.

Gegeben zwei Grammatikklassen \hat{G} und \hat{G}' :

$K(\hat{G}) > K(\hat{G}') \leftrightarrow$ Für jede $G' \in \hat{G}'$ gibt es eine $G \in \hat{G}$ und $\mathcal{P}_{G'} = \mathcal{P}_G$

$K(\hat{G}) > K(\hat{G}') \leftrightarrow$ Für jede $G' \in \hat{G}'$ gibt es eine $G \in \hat{G}$ und $L_{G'} = L_G$

Intuition

Je komplexer desto umfangreicher.

Gegeben zwei Grammatikklassen \hat{G} und \hat{G}' :

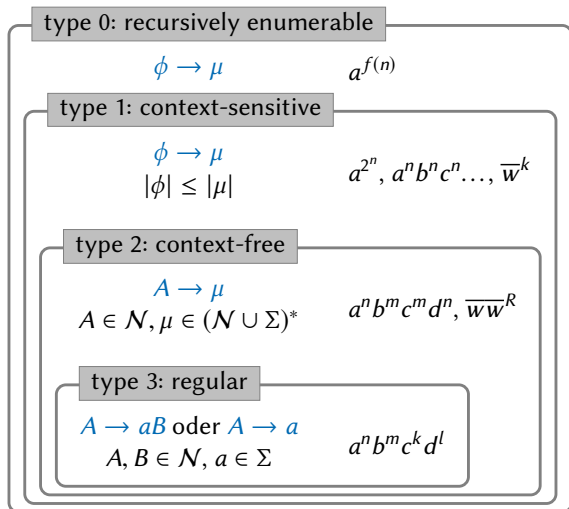
$K(\hat{G}) > K(\hat{G}') \leftrightarrow$ Für jede $G' \in \hat{G}'$ gibt es eine $G \in \hat{G}$ und $\mathcal{P}_{G'} = \mathcal{P}_G$

$K(\hat{G}) > K(\hat{G}') \leftrightarrow$ Für jede $G' \in \hat{G}'$ gibt es eine $G \in \hat{G}$ und $L_{G'} = L_G$

\Rightarrow berühmtestes Beispiel: Chomsky(-Schützenberger)-Hierarchie^[6]

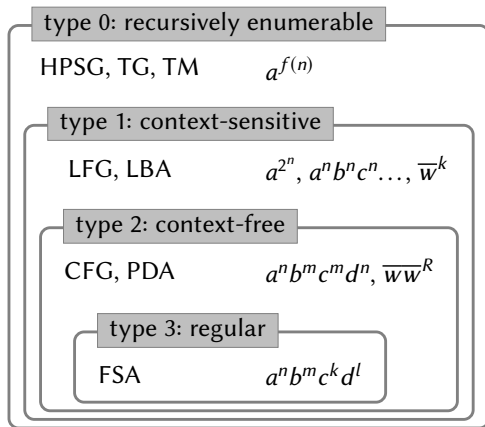
Relative extensionale Komplexität von Grammatikklassen

Chomsky(-Schützenberger)-Hierarchie^[6]



Relative extensionale Komplexität von Grammatikklassen

Chomsky(-Schützenberger)-Hierarchie^[6]



Wo liegt die Klasse
der natürlichen
Sprachen?

(Einzelne natürliche
Sprachen können in
unterschiedlichen
Klassen sein.)

NL \subset Type 3 / regular?

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

- Q: endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

- Q : endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)
- δ : endliche Menge von Übergängen zwischen Zuständen

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

- Q: endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)
- δ : endliche Menge von Übergängen zwischen Zuständen
 - Bei jedem Übergang wird ein Symbol von der Eingabe gelesen.

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

- Q: endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)
- δ : endliche Menge von Übergängen zwischen Zuständen
 - Bei jedem Übergang wird ein Symbol von der Eingabe gelesen.
 - kein Speicher (außer den Zuständen)!

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

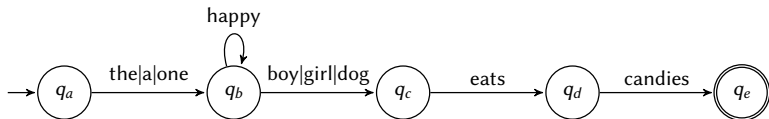
- Q: endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)
- δ : endliche Menge von Übergängen zwischen Zuständen
 - Bei jedem Übergang wird ein Symbol von der Eingabe gelesen.
 - kein Speicher (außer den Zuständen)!

NL \subset Type 3 / regular?

Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

- Q: endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)
- δ : endliche Menge von Übergängen zwischen Zuständen
 - Bei jedem Übergang wird ein Symbol von der Eingabe gelesen.
 - kein Speicher (außer den Zuständen)!



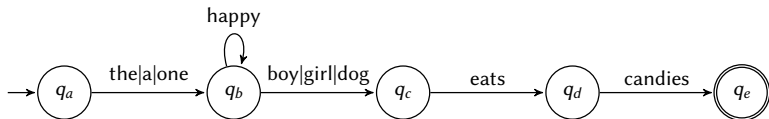
Hypothese

Für jede natürliche Sprache L gibt es einen **endlichen Automaten**, der genau die Wörter aus L akzeptiert.

Q: endliche Menge von Zuständen (inklusive einem Startzustand und mindestens einem Endzustand)

δ : endliche Menge von Übergängen zwischen Zuständen

- Bei jedem Übergang wird ein Symbol von der Eingabe gelesen.
- kein Speicher (außer den Zuständen)!



Wie **beweist** man die Inadäquatheit von endlichen Automaten auf der Ebene der Stringsprache?

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^n y} \in L$.

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny}$ $\in L$.

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^n y} \in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny}$ $\in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman** **whom another woman hired**
hired hired another woman

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny}$ $\in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman** **whom another woman hired**
hired hired another woman
d. ...

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny}$ $\in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman whom another woman hired**
hired hired another woman
d. ...

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny}$ $\in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman** **whom another woman hired**
hired hired another woman
d. ...

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^n y} \in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman** **whom another woman hired**
hired hired another woman
d. ...

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny}$ $\in L$.

Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman** **whom another woman hired**
hired hired another woman
d. ...

Formaler Beweis durch Widerspruch (mittels **Abschlusseigenschaften**):

- Jede reguläre Sprache L erfüllt das **Pump-Lemma** für reguläre Sprachen, d.h. für jedes Wort ab einer bestimmten Wortlänge gilt: Es gibt eine Dekomposition \overline{xy} , so dass $\overline{xa^ny} \in L$.

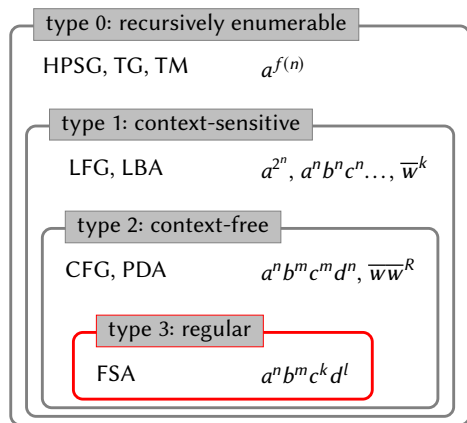
Einbettende Abhängigkeiten (Chomsky 1957)

- (1) a. a woman hired another woman
b. a woman **whom another woman hired** hired another woman
c. a woman **whom another woman whom another woman hired**
hired hired another woman
d. ...

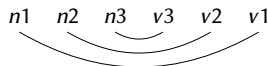
- Homomorphismus f : $f(\text{a woman}) = x$, $f(\text{whom another woman}) = a$,
 $f(\text{hired}) = b$, $f(\text{hired another woman}) = y$
 - xa^*b^*y ist eine reguläre Sprache; und
 - $f(\text{Englisch}) \cap xa^*b^*y = xa^nb^ny$ sollte auch regulär sein.
- xa^nb^ny widerspricht dem Pump-Lemma für reguläre Sprachen.

NL \notin Type 3 / regular!

Chomsky(-Schützenberger)-Hierarchie^[6]



NL ist nicht regulär!^[3,4]
einbettende Abhängigkeiten
bei Relativsätzen



NL \subset Type 2 / context-free?

- eine langwierige Diskussion

NL \subset Type 2 / context-free?

- eine langwierige Diskussion

Chomsky (1957):34

“Of course there are languages (in our general sense) that cannot be described in terms of phrase structure, but I do not know whether or not English is itself literally outside the range of such analysis.”

NL \subset Type 2 / context-free?

- eine langwierige Diskussion

Chomsky (1957):34

“Of course there are languages (in our general sense) that cannot be described in terms of phrase structure, but I do not know whether or not English is itself literally outside the range of such analysis.”

- mehrere fehlerhafte Versuche (siehe Pullum & Gazdar [17]), e.g.:

- eine langwierige Diskussion

Chomsky (1957):34

“Of course there are languages (in our general sense) that cannot be described in terms of phrase structure, but I do not know whether or not English is itself literally outside the range of such analysis.”

- mehrere fehlerhafte Versuche (siehe Pullum & Gazdar [17]), e.g.:

Bresnan (1978):37–38

“in many cases the number of a verb agrees with that of a noun phrase at some distance from it ... this type of syntactic dependency can extend as memory or patience permits ... the distant type of agreement ... cannot be adequately described even by context-sensitive phrase-structure rules, for the possible context is not correctly describable as a finite string of phrases.”

NL \subset Type 2 / context-free?

- eine langwierige Diskussion

Chomsky (1957):34

“Of course there are languages (in our general sense) that cannot be described in terms of phrase structure, but I do not know whether or not English is itself literally outside the range of such analysis.”

- mehrere fehlerhafte Versuche (siehe Pullum & Gazdar [17]), e.g.:

Bresnan (1978):37–38

“in many cases the number of a verb agrees with that of a noun phrase at some distance from it ... this type of syntactic dependency can extend as memory or patience permits ... the distant type of agreement ... cannot be adequately described even by context-sensitive phrase-structure rules, for the possible context is not correctly describable as a finite string of phrases.”

- richtige Beweistechniken: Pump-Lemma und Abschlusseigenschaften
- Was ist ein nicht-kontextfreies Phänomen in natürlicher Sprache?

NL \subset Type 2 / context-free?

Nahe dran: Niederländisch (Bresnan et al. 1982)

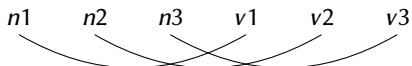
- (2) dat Jan Piet de kinderen zag helpen zwemmen.
that Jan Piet the children saw help swim
'that Jan saw Piet helping the children to swim.'

NL \subset Type 2 / context-free?

Nahe dran: Niederländisch (Bresnan et al. 1982)

- (2) dat Jan Piet de kinderen zag helpen zwemmen.
that Jan Piet the children saw help swim
'that Jan saw Piet helping the children to swim.'

Linguistische Abhängigkeiten sind kreuzend:

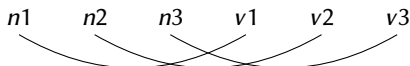


NL \subset Type 2 / context-free?

Nahe dran: Niederländisch (Bresnan et al. 1982)

- (2) dat Jan Piet de kinderen zag helpen zwemmen.
that Jan Piet the children saw help swim
'that Jan saw Piet helping the children to swim.'

Linguistische Abhängigkeiten sind kreuzend:



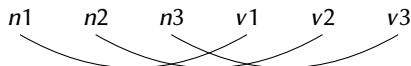
Aber: Kein Reflex an der Oberfläche, d.h. in der Stringsprache!

NL \subset Type 2 / context-free?

Nahe dran: Niederländisch (Bresnan et al. 1982)

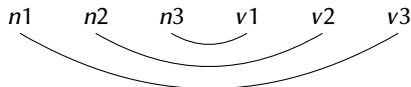
- (2) dat Jan Piet de kinderen zag helpen zwemmen.
that Jan Piet the children saw help swim
'that Jan saw Piet helping the children to swim.'

Linguistische Abhängigkeiten sind kreuzend:



Aber: Kein Reflex an der Oberfläche, d.h. in der Stringsprache!

- ⇒ Der String kann durch eine CFG erzeugt werden, auch wenn die Abhängigkeiten verloren gehen.



NL \subset Type 2 / context-free?

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

(3) a. wulu

'dog'

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- (3) a. wulu
b. wulu-lela

'dog'

'dog watcher'

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- (3) a. wulu *'dog'*
b. wulu-lela *'dog watcher'*
c. wulu-lela-nyinila *'dog watcher hunter'*

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- | | | |
|--------|-------------------|-----------------------------|
| (3) a. | wulu | <i>'dog'</i> |
| b. | wulu-lela | <i>'dog watcher'</i> |
| c. | wulu-lela-nyinila | <i>'dog watcher hunter'</i> |
| d. | wulu-o-wulu | <i>'whatever dog'</i> |

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- | | | |
|--------|-----------------------|-----------------------------|
| (3) a. | wulu | <i>'dog'</i> |
| b. | wulu-lela | <i>'dog watcher'</i> |
| c. | wulu-lela-nyinila | <i>'dog watcher hunter'</i> |
| d. | wulu-o-wulu | <i>'whatever dog'</i> |
| e. | wulu-lela-o-wulu-lela | |

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- | | | |
|--------|-----------------------|-------------------------------|
| (3) a. | wulu | <i>'dog'</i> |
| b. | wulu-lela | <i>'dog watcher'</i> |
| c. | wulu-lela-nyinila | <i>'dog watcher hunter'</i> |
| d. | wulu-o-wulu | <i>'whatever dog'</i> |
| e. | wulu-lela-o-wulu-lela | <i>'whatever dog watcher'</i> |

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- | | | |
|--------|---------------------------------------|--------------------------------------|
| (3) a. | wulu | <i>'dog'</i> |
| b. | wulu-lela | <i>'dog watcher'</i> |
| c. | wulu-lela-nyinila | <i>'dog watcher hunter'</i> |
| d. | wulu-o-wulu | <i>'whatever dog'</i> |
| e. | wulu-lela-o-wulu-lela | <i>'whatever dog watcher'</i> |
| f. | wulu-lela-nyinila-o-wulu-lela-nyinila | <i>'whatever dog watcher hunter'</i> |

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- | | | |
|--------|---------------------------------------|-------------------------------|
| (3) a. | wulu | 'dog' |
| b. | wulu-lela | 'dog watcher' |
| c. | wulu-lela-nyinila | 'dog watcher hunter' |
| d. | wulu-o-wulu | 'whatever dog' |
| e. | wulu-lela-o-wulu-lela | 'whatever dog watcher' |
| f. | wulu-lela-nyinila-o-wulu-lela-nyinila | 'whatever dog watcher hunter' |

Muster: $a^n b^m a^n b^m$ oder \overline{ww} (copy language)

⇒ nicht kontextfrei!

NL \subset Type 2 / context-free?

Ein weiterer Versuch von Culy (1985): Duplikation in der Morphologie des Bambara

- | | | |
|--------|---------------------------------------|-------------------------------|
| (3) a. | wulu | 'dog' |
| b. | wulu-lela | 'dog watcher' |
| c. | wulu-lela-nyinila | 'dog watcher hunter' |
| d. | wulu-o-wulu | 'whatever dog' |
| e. | wulu-lela-o-wulu-lela | 'whatever dog watcher' |
| f. | wulu-lela-nyinila-o-wulu-lela-nyinila | 'whatever dog watcher hunter' |

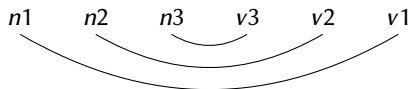
Muster: $a^n b^m a^n b^m$ oder \overline{ww} (copy language)

⇒ nicht kontextfrei!

Aber: Beweis für die Morphologie, nicht für die Syntax (Stringsprache)!

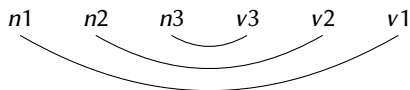
■ Deutsch: **einbettende Abhängigkeiten** (VE-Sätze)

- (4) (dass) er die Kinder dem Hans das Haus streichen helfen ließ
(that) he the children the Hans the house paint help let
'(that) he let the children help Hans paint the house'



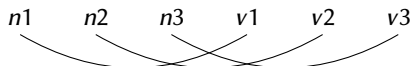
■ Deutsch: **einbettende Abhängigkeiten** (VE-Sätze)

- (4) (dass) er die Kinder dem Hans das Haus streichen helfen ließ
(that) he the children the Hans the house paint help let
'(that) he let the children help Hans paint the house'



■ Schwyzerdütsch: **kreuzende Abhängigkeiten**

- (5) ...mer d'chind em Hans es huus lönd hälfe aastriiche
...we children.ACC the Hans.DAT the house.ACC let help paint
'...we let the children help Hans paint the house'



NL \subset Type 2 / context-free?

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastriiche.

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastrichte.

■ Homomorphismus f :

$$\begin{array}{llll} f(\text{Jan säit das mer}) = x & f(\text{d'chind}) = a & f(\text{em Hans}) = b & f(\text{es huus}) = y \\ f(\text{lönd}) = c & f(\text{hälfe}) = d & & f(\text{aastriche}) = z \end{array}$$

NL \subset Type 2 / context-free?

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastrichte.

■ Homomorphismus f :

$$\begin{array}{llll} f(\text{Jan säit das mer}) = x & f(\text{d'chind}) = a & f(\text{em Hans}) = b & f(\text{es huus}) = y \\ & f(\text{lönd}) = c & f(\text{hälfe}) = d & f(\text{aastriche}) = z \end{array}$$

■ $f(\text{Schwyzerdütsch}) \cap xa^*b^*yc^*d^*z = xa^m b^n yc^m d^n z$

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastriiche.

- Homomorphismus f :

$$\begin{array}{llll} f(\text{Jan säit das mer}) = x & f(\text{d'chind}) = a & f(\text{em Hans}) = b & f(\text{es huus}) = y \\ f(\text{lönd}) = c & f(\text{hälfe}) = d & & f(\text{aastriiche}) = z \end{array}$$

- $f(\text{Schwyzerdütsch}) \cap xa^*b^*yc^*d^*z = xa^m b^n yc^m d^n z$
 - Kontextfreie Sprachen sind **abgeschlossen** unter Schnitt mit regulären Sprachen.

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastriiche.

- Homomorphismus f :

$$\begin{array}{llll} f(\text{Jan säit das mer}) = x & f(\text{d'chind}) = a & f(\text{em Hans}) = b & f(\text{es huus}) = y \\ f(\text{lönd}) = c & f(\text{hälfe}) = d & & f(\text{aastriiche}) = z \end{array}$$

- $f(\text{Schwyzerdütsch}) \cap xa^*b^*yc^*d^*z = xa^mb^nc^md^nz$
 - Kontextfreie Sprachen sind **abgeschlossen** unter Schnitt mit regulären Sprachen.
 - $xa^*b^*yc^*d^*z$ ist regulär, also sollte $xa^mb^nc^md^nz$ kontextfrei sein.

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastriiche.

- Homomorphismus f :

$$\begin{array}{llll} f(\text{Jan säit das mer}) = x & f(\text{d'chind}) = a & f(\text{em Hans}) = b & f(\text{es huus}) = y \\ f(\text{lönd}) = c & f(\text{hälfe}) = d & & f(\text{aastriiche}) = z \end{array}$$

- $f(\text{Schwyzerdütsch}) \cap xa^*b^*yc^*d^*z = xa^mb^ncy^md^nz$
 - Kontextfreie Sprachen sind **abgeschlossen** unter Schnitt mit regulären Sprachen.
 - $xa^*b^*yc^*d^*z$ ist regulär, also sollte $xa^mb^ncy^md^nz$ kontextfrei sein.
 - Aber mittels **Pump-Lemma für kontextfreie Sprachen**:
 $xa^mb^ncy^md^nz$ ist nicht kontextfrei.

Beweis von Shieber (1985):

(6) Jan säit das mer d'chind em Hans es huus lönd hälfe aastriiche.

- Homomorphismus f :

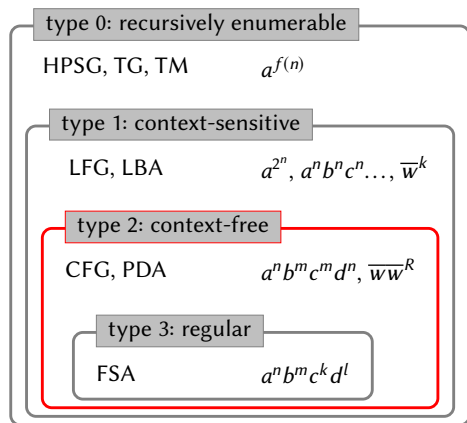
$$\begin{array}{llll} f(\text{Jan säit das mer}) = x & f(\text{d'chind}) = a & f(\text{em Hans}) = b & f(\text{es huus}) = y \\ f(\text{lönd}) = c & f(\text{hälfe}) = d & & f(\text{aastriiche}) = z \end{array}$$

- $f(\text{Schwyzerdütsch}) \cap xa^*b^*yc^*d^*z = xa^mb^ncy^md^nz$
 - Kontextfreie Sprachen sind **abgeschlossen** unter Schnitt mit regulären Sprachen.
 - $xa^*b^*yc^*d^*z$ ist regulär, also sollte $xa^mb^ncy^md^nz$ kontextfrei sein.
 - Aber mittels **Pump-Lemma für kontextfreie Sprachen**:
 $xa^mb^ncy^md^nz$ ist nicht kontextfrei.

⇒ Schwyzerdütsch ist nicht kontextfrei!

NL $\not\subset$ Type 2 / context-free!

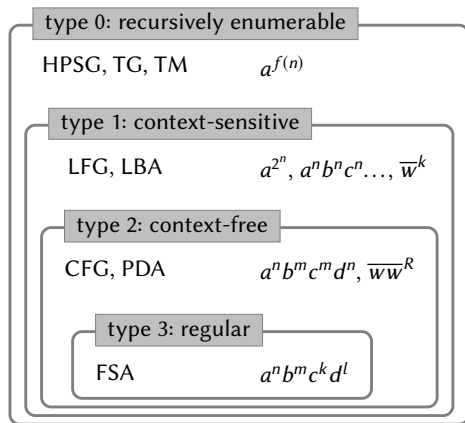
Chomsky(-Schützenberger)-Hierarchie^[6]



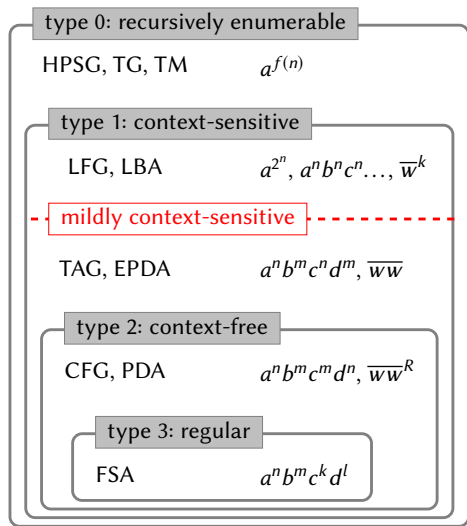
NL ist nicht kontextfrei!^[21]
kreuzende Abhängigkeiten
im Niederländischen und
Schwyzerdütsch

$n1 \quad n2 \quad n3 \quad v1 \quad v2 \quad v3$

NL \subset Type 1–2 / mildly context-sensitive?



NL \subset Type 1-2 / mildly context-sensitive?



NL ist **schwach kontextsensitiv**? (Joshi [13])

- \supset CFL
- kreuzende Abh.
- semi-linear
- in PTIME

Argumente gegen Semilinearität:

- Kasusstapel (“Suffixaufnahme”) im Alt-Georgischen (Michaelis & Kracht 1997)

N_1 -NOM N_2 -GEN-NOM N_3 -GEN²-NOM ... N_k -GEN ^{$k-1$} -NOM

“The number of all genitive suffixes of all nouns within a complex NP consisting of k stacked NPs, where $k \in \mathbb{N}^+$, is bounded by $k^2/2 - k/2$.”

Argumente gegen Semilinearität:

- Kasusstapel (“Suffixaufnahme”) im Alt-Georgischen (Michaelis & Kracht 1997)

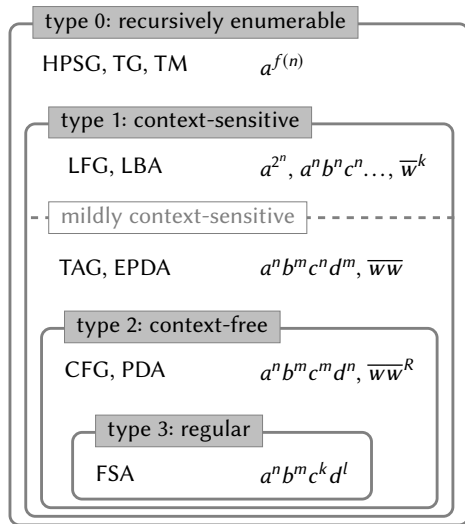
N_1 -NOM N_2 -GEN-NOM N_3 -GEN²-NOM ... N_k -GEN ^{$k-1$} -NOM

“The number of all genitive suffixes of all nouns within a complex NP consisting of k stacked NPs, where $k \in \mathbb{N}^+$, is bounded by $k^2/2 - k/2$.”

- Chinesische Zahlterme (Radzinski 1991)
- Koordination im Niederländischen (Groenink 1997)
- Relativsätze im Yoruba (Kobele 2006)

⇒ Im Allgemeinen in der Literatur nicht berücksichtigt.

NL \subset Type 1 / context-sensitive?



NL unentscheidbar?

- Hintikka (1974)
- Chomsky (1980)
- Pollard (1996)

Subreguläre Typen: Type 4, 5, ...

Wie Reguläre Sprachen weiter beschränken?

- Regelschema: Anzahl der Nichtterminale
- Automaten: Anzahl der Zustände, Struktur der Übergänge

Subreguläre Typen: Type 4, 5, ...

Wie Reguläre Sprachen weiter beschränken?

- Regelschema: Anzahl der Nichtterminale
- Automaten: Anzahl der Zustände, Struktur der Übergänge

Endliche Sprachen: Naheliegend aber formal eher uninteressant

- $S \rightarrow \Sigma^+$
- $S \rightarrow ab$
- $S \rightarrow a$

Subreguläre Typen: Type 4, 5, ...

Wie Reguläre Sprachen weiter beschränken?

- Regelschema: Anzahl der Nichtterminale
- Automaten: Anzahl der Zustände, Struktur der Übergänge

Endliche Sprachen: Naheliegend aber formal eher uninteressant

- $S \rightarrow \Sigma^+$
- $S \rightarrow ab$
- $S \rightarrow a$

Interessantere subreguläre Sprachen:

- Strictly Local Languages (SL_k)^[12,19]

Subreguläre Typen: Type 4, 5, ...

Wie Reguläre Sprachen weiter beschränken?

- Regelschema: Anzahl der Nichtterminale
- Automaten: Anzahl der Zustände, Struktur der Übergänge

Endliche Sprachen: Naheliegend aber formal eher uninteressant

- $S \rightarrow \Sigma^+$
- $S \rightarrow ab$
- $S \rightarrow a$

Interessantere subreguläre Sprachen:

- Strictly Local Languages (SL_k)^[12,19]

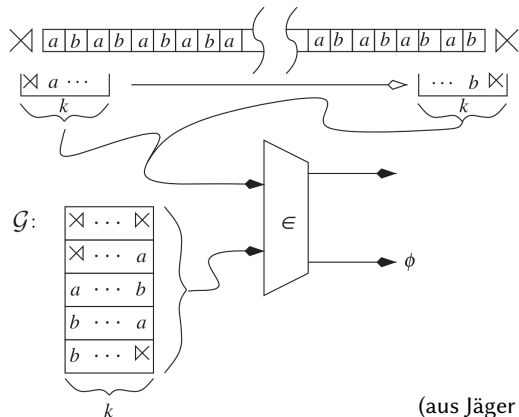
Aber was hat das mit natürlicher Sprache zu tun?

⇒ Syntax/Typologie^[11], Phonologie

Strictly Local Languages (SL_k)

Die Sprache wird nur anhand zusammenhängender “Blöcke” der Größe k beschränkt.

Automatenmodell: “scanner” mit “look-up list”



(aus Jäger & Rogers [12])

(siehe z.B. Rogers & Pullum [19])

Strictly Local Languages (SL_k)

Die Sprache wird nur anhand zusammenhängender “Blöcke” der Größe k beschränkt.

Automatenmodell: “scanner” mit “look-up list”

- Die Grammatik G besteht aus einer Menge zulässiger Blöcke, z.B. $G = \{\times A, AB, BA, B \times\}$ mit $k = 2$.
- $F_k(\bar{w})$ ist die Menge der k -Faktoren in \bar{w} .
$$F_k(\bar{w}) = \begin{cases} \{y \mid \bar{w} = x \cdot y \cdot z, \ x, y, z \in \Sigma^*, |y| = k\} & \text{if } |\bar{w}| > k, \\ \{\bar{w}\} & \text{otherwise.} \end{cases}$$
- $L(G) = \{\bar{w} \mid F_k(\times \cdot \bar{w} \cdot \times) \subseteq G, \bar{w} \text{ finite}\}$

(siehe z.B. Rogers & Pullum [19])

Strictly Local Languages (SL_k)

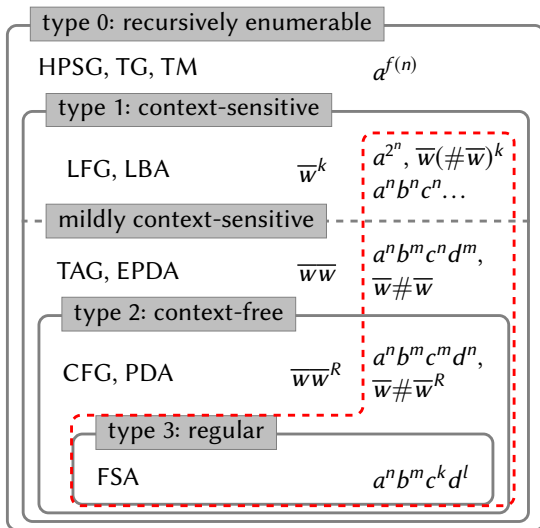
Die Sprache wird nur anhand zusammenhängender “Blöcke” der Größe k beschränkt.

Automatenmodell: “scanner” mit “look-up list”

- Die Grammatik G besteht aus einer Menge zulässiger Blöcke, z.B. $G = \{\times A, AB, BA, B\times\}$ mit $k = 2$.
- $F_k(\bar{w})$ ist die Menge der k -Faktoren in \bar{w} .
$$F_k(\bar{w}) = \begin{cases} \{y \mid \bar{w} = x \cdot y \cdot z, \ x, y, z \in \Sigma^*, |y| = k\} & \text{if } |\bar{w}| > k, \\ \{\bar{w}\} & \text{otherwise.} \end{cases}$$
- $L(G) = \{\bar{w} \mid F_k(\times \cdot \bar{w} \cdot \times) \subseteq G, \bar{w} \text{ finite}\}$
- $(ab)^n \in SL_2$
- aber: $a^*ba^* \notin SL_k$

(siehe z.B. Rogers & Pullum [19])

Orthogonale Hierarchien



orthogonal classes:^[9,22]
deterministic real-time
definable (Rosenberg 1967)

(siehe auch die Sitzung
über Lernbarkeit)

- extensionale Komplexität von (formalen) Sprachklassen
- Wo liegt natürliche Sprache in der Chomsky-Hierarchie?
- Die **Klasse** der natürlichen Sprachen ist zumindest schwach kontextsensitiv.
- formale Beweise mittels Abschlusseigenschaften und Pump-Lemma

Appendix: Abschlusseigenschaften

	$L_1 \cap L_2$	$L_1 \cup L_2$	L^c	$L_1 \bullet L_2$	L^*
Type 0	+	+	-	+	+
Type 1	+	+	+	+	+
Type 2	-	+	-	+	+
Type 3	+	+	+	+	+
Type 2/3	+	+	+	+	+

Ist eine Sprache L kontextfrei, dann existiert eine Mindestwortlänge k für L , ab der ein Wort $\bar{w} \in L$ in Teilworte $\bar{x}\bar{a}\bar{y}\bar{b}\bar{z}$ dekomponiert werden kann, so dass auch $\bar{x}\bar{a}^n\bar{y}\bar{b}^n\bar{z} \in L$ für $n \geq 0$, wobei gelten muss:

1 $|\bar{xyz}| \leq k,$

2 $|\bar{ab}| \geq 1.$

- [1] Bresnan, Joan. 1978. A realistic transformational grammar. In Morris Halle, Joan Bresnan & George A. Miller (eds.), *Linguistic theory and psychological reality*, 1–59. Cambridge, MA: The MIT Press.
- [2] Bresnan, Joan, Ronald M. Kaplan, Stanley Peters & Annie Zaenen. 1982. Cross-serial dependencies in dutch. *Linguistic Inquiry* 13(4). 613–634.
- [3] Chomsky, Noam. 1956. Three models for the description of language. *IRE Transactions on Information Theory* 2. 113–124.
- [4] Chomsky, Noam. 1957. *Syntactic structures*. Den Haag: Mouton.
- [5] Chomsky, Noam. 1980. *Rules and representations*. Oxford, UK: Basil Blackwell.
- [6] Chomsky, Noam & Marcel-Paul Schützenberger. 1963. The algebraic theory of context-free languages. In P. Braffort & D. Hirschberg (eds.), *Computer programming and formal systems* (Studies in Logic and the Foundations of Mathematics 35), 118–161. Elsevier.
- [7] Culy, Christopher. 1985. The complexity of the vocabulary of Bambara. *Linguistics and Philosophy* 8(3). 345–351. <http://www.jstor.org/stable/25001211>.
- [8] Groenink, Annius V. 1997. Mild context-sensitivity and tuple-based generalizations of context-grammar. *Linguistics and Philosophy* 20(6). 607–636. <http://dx.doi.org/10.1023/A%3A1005376413354>.
- [9] Hausser, Roland. 1992. Complexity in Left-Associative Grammar. *Theoretical Computer Science* 106(2). 283–308. <http://www.sciencedirect.com/science/article/pii/030439759290253C>.
- [10] Hintikka, Jaakko. 1974. Quantifiers vs. quantification theory. *Linguistic Inquiry* 5(2). 153–177. <http://www.jstor.org/stable/4177815>.

- [11] Jackendoff, Ray & Eva Wittenberg. 2014. What you can say without syntax: A hierarchy of grammatical complexity. In Frederick J. Newmeyer & Laurel B. Preston (eds.), *Measuring grammatical complexity*, 65–82. Oxford: Oxford University Press.
- [12] Jäger, Gerhard & James Rogers. 2012. Formal language theory: Refining the Chomsky hierarchy. *Philosophical Transactions of the Royal Society B: Biological Sciences* 367(1598). 1956–1970.
- [13] Joshi, Aravind K. 1985. Tree adjoining grammars: how much context-sensitivity is required to provide reasonable structural descriptions. In David Dowty, Lauri Karttunen & Arnold Zwicky (eds.), *Natural language parsing*, 206–250. Cambridge University Press.
- [14] Kobele, Gregory M. 2006. *Generating copies: an investigation into structural identity in language and grammar*. Los Angeles: University of California dissertation.
<http://home.uchicago.edu/~gkobe/~/files/Kobe06GeneratingCopies.pdf>.
- [15] Michaelis, Jens & Marcus Kracht. 1997. Semilinearity as a syntactic invariant. In Christian Retoré (ed.), *Logical aspects of computational linguistics* (Lecture Notes in Computer Science 1328), 329–345. Berlin: Springer.
<http://dx.doi.org/10.1007/BFb0052165>.
- [16] Pollard, Carl. 1996. The nature of constraint grammar. Paper presented at the 11th Pacific Conference of Language, Information and Computation (PACLIC).
- [17] Pullum, Geoffrey K. & Gerald Gazdar. 1982. Natural languages and context-free languages. *Linguistics and Philosophy* 4(4). 471–504.
<http://www.jstor.org/stable/25001071>.
- [18] Radzinski, Daniel. 1991. Chinese number-names, tree adjoining languages, and mild context-sensitivity. *Computational Linguistics* 17. 277–299.

- [19] Rogers, James & Geoffrey K. Pullum. 2007. Aural pattern recognition experiments and the subregular hierarchy. *Journal of Language, Logic and Information* 20. 329–342.
- [20] Rosenberg, Arnold L. 1967. Real-time definable languages. *Journal of the Association for Computing Machinery* 14(4). 645–662.
- [21] Shieber, Stuart. 1985. Evidence against the context-freeness of natural language. *Linguistics and Philosophy* 8. 333–343.
- [22] Wurm, Christian. 2012. Regular Growth Automata: properties of a class of finitely induced infinite machines. In Makoto Kanazawa, Markus Kracht, Hiroyuki Seki & Andras Kornai (eds.), *Proceedings of the 12th conference on the mathematics of language (MOL 2012)* (LNCS 6878), 192–208. Nara, Japan: Springer.