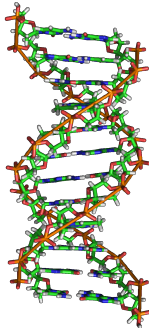


Grammatik der Genome

Andreas de Vries

Fachhochschule Südwestfalen, Campus Hagen

SS 2011



1 Formale Sprachen

2 Genomgrammatik

Übersicht

1 Formale Sprachen

2 Genomgrammatik

Formale Sprachen

- Die Theorie der formalen Sprachen behandelt die systematische Analyse, Klassifizierung und Konstruktion von Wortmengen über endlichen Alphabeten.
- Entstanden in den 1950er Jahren in der Linguistik.
- Formale Sprachen (Prädikatenkalkül, Java, HTML, ...) haben wohldefinierte und eindeutige **Grammatiken**.
- Natürliche Sprachen haben dagegen **Mehrdeutigkeiten** auf jeder Beschreibungsebene, von der phonetischen bis zur soziologischen.
- Sprachliche **Witze** verwenden diese Mehrdeutigkeiten, insbesondere Wortspiele: „*Bis dann, wir sehen uns! – Wir sind ja nicht blind.*“

Formale Sprachen

- Die Theorie der formalen Sprachen behandelt die systematische Analyse, Klassifizierung und Konstruktion von Wortmengen über endlichen Alphabeten.
- Entstanden in den 1950er Jahren in der Linguistik.
- Formale Sprachen (Prädikatenkalkül, Java, HTML, ...) haben wohldefinierte und eindeutige **Grammatiken**.
- Natürliche Sprachen haben dagegen **Mehrdeutigkeiten** auf jeder Beschreibungsebene, von der phonetischen bis zur soziologischen.
- Sprachliche **Witze** verwenden diese Mehrdeutigkeiten, insbesondere Wortspiele: „*Bis dann, wir sehen uns! – Wir sind ja nicht blind.*“
- Merke: **Formale Sprachen eignen sich nicht für Späße!**

Formale Grammatiken

Example (*Simple English*)

Mit dem Lateinischen Alphabet mit den 26 Kleinbuchstaben, dem Leerzeichen und dem Punkt ist durch die \langle Variablen \rangle und Regeln für die Wortbildung („Produktionen“) eine **Grammatik** gegeben:

\langle sentence $\rangle \rightarrow \langle$ subject $\rangle \langle$ predicate $\rangle \langle$ object \rangle .
 \langle subject $\rangle \rightarrow \langle$ article $\rangle \langle$ adjective $\rangle \langle$ noun \rangle
 \langle article $\rangle \rightarrow a \mid the$
 \langle adjective $\rangle \rightarrow sweet \mid quick \mid small$
 \langle noun $\rangle \rightarrow duck \mid frog \mid mouse \mid hippopotamus$
 \langle predicate $\rangle \rightarrow likes \mid catches \mid eats$
 \langle object $\rangle \rightarrow cookies \mid chocolate \mid pizza$

Die durch diese Grammatik erzeugte Sprache beinhaltet Sätze wie “the small duck eats pizza.” und “a quick mouse catches cookies.”

Chomsky-Hierarchie

General languages

\mathcal{L}_0 : Phrase-structured languages
(Type-0 languages)

\mathcal{L}_1 : Context-sensitive languages
(Type-1 languages)

\mathcal{L}_2 : Context-free languages
(Type-2 languages)

\mathcal{L}_3 : Regular languages
(Type-3 languages)

Chomsky-Hierarchie

General languages

\mathcal{L}_0 : Phrase-structured languages
(Type-0 languages)

\mathcal{L}_1 : Context-sensitive languages
(Type-1 languages)

\mathcal{L}_2 : Context-free languages
(Type-2 languages)

\mathcal{L}_3 : Regular languages
(Type-3 languages)

General languages

\mathcal{L}_0 : Phrase-structured languages

\mathcal{L}_1 : Context-sensitive languages

\mathcal{L}_{ind} : Indexed languages

\mathcal{L}_2 : Context-free languages

\mathcal{L}_{lin} : Linear languages

\mathcal{L}_3 : Regular languages

Übersicht

1 Formale Sprachen

2 Genomgrammatik

Genomsequenzen

Genomsequenzen sind auffassbar als Strings über dem DNA Alphabet

$$\Sigma_{\text{DNA}} = \{A, C, G, T\}. \quad (1)$$



■ Phosphatrückgrat (-) ■ Adenin (A) ■ Thymin (T) ■ Guanin (G) ■ Cytosin (C)

Angelehnt an: http://upload.wikimedia.org/wikipedia/commons/e/e7/DNA_simple.svg

DNA-Komplement

Definition

Sei die Abbildung $\bar{\cdot} : \Sigma_{\text{DNA}} \cup \{\varepsilon\} \rightarrow \Sigma_{\text{DNA}} \cup \{\varepsilon\}$ mit

$$\bar{\varepsilon} = \varepsilon, \quad \bar{A} = T, \quad \bar{C} = G, \quad \bar{G} = C, \quad \bar{T} = A$$

gegeben. Für ein Wort $w \in \Sigma_{\text{DNA}}^*$ ist dann das **Komplement** \bar{w} gegeben durch

$$\overline{w_1 w_2 \cdots w_n} = \bar{w}_n \cdots \bar{w}_2 \bar{w}_1, \quad (2)$$

mit $w = w_1 w_2 \cdots w_n$, wobei $w_i \in \Sigma$. Die so definierte Abbildung $\bar{\cdot} : \Sigma_{\text{DNA}}^* \rightarrow \Sigma_{\text{DNA}}^*$ heißt **DNA-Komplementierung**, und \bar{w} das **DNA-Komplement** des Wortes w .

- Für ein DNA-Wort $w \in \Sigma_{\text{DNA}}^*$ modelliert das Komplement \bar{w} den gegenüberliegenden Helixstrang.
- Komplementäre Wörter können auch als Resultate der DNA-Replikation angesehen werden.

Sätze über komplementäre Genomsequenzen

Theorem (Watson und Crick 1953)

Die Replikation des Wortes des gegenüberliegenden DNA-Strangs ergibt wieder das originale Wort.

Lemma (Dyadische Symmetrie)

Für ein Wort $w \in \Sigma_{\text{DNA}}^$ gilt $w = \bar{w}$ genau dann, wenn $w = u\bar{u}$ für ein $u \in \Sigma_{\text{DNA}}^*$. Insbesondere hat w eine gerade Länge.*

Lemma (Kopiereigenschaft)

Für ein Wort $w \in \Sigma_{\text{DNA}}^$ mit $w = \bar{w}$ folgt $w^n = \overline{w^n}$ für alle $n \in \mathbb{N}$.*

DNA-Grammatik: Haarnadel-Moleküle

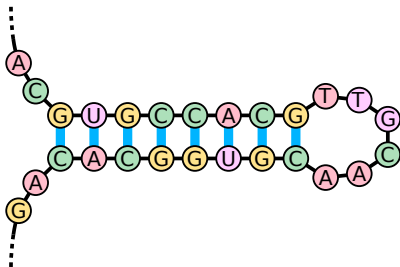
Haarnadel-Moleküle sind Einzelstränge, die sich auf sich selbst falten und sich wie getrennte Stränge verhalten. Sie bilden die formale Sprache

$$L_h = \{w \in \Sigma_{\text{DNA}}^* \mid w = \bar{w}\}. \quad (3)$$

L_h ähnelt einer Palindrom-Sprache. Sie wird erzeugt von der kontextfreien Grammatik

$$S \rightarrow x S \bar{x} \mid \varepsilon \quad (4)$$

mit $x \in \Sigma_{\text{DNA}}$.



Angelehnt an: <http://upload.wikimedia.org/wikipedia/commons/3/3f/Stem-loop.svg>

DNA-Grammatik: Haarnadel-Moleküle

Example

Die Wörter $v = \text{GATC}$ und $w = \text{TGGCCA}$,

$$v = \begin{array}{c} \boxed{\begin{array}{c} 5' \rightarrow \text{GA} \rightarrow 3' \\ 3' \leftarrow \text{CT} \leftarrow 5' \end{array}} \quad w = \begin{array}{c} \boxed{\begin{array}{c} 5' \rightarrow \text{TGG} \rightarrow 3' \\ 3' \leftarrow \text{GGT} \leftarrow 5' \end{array}} \end{array} \quad (5)$$

sind in L_h , da die Produktionen (4) ergeben:

$$S \Rightarrow \text{GSC} \Rightarrow \text{GASTC} \Rightarrow \text{GATC}$$

und

$$S \Rightarrow \text{GSC} \Rightarrow \text{GGSCC} \Rightarrow \text{TGGSCCA} \Rightarrow \text{TGGCCA}.$$

Zu welcher Sprachfamilie gehört die DNA-Grammatik?

Theorem (Searl 1999)

L_h ist kontextfrei, aber nicht regulär.

Zu welcher Sprachfamilie gehört die DNA-Grammatik?

Theorem (Searl 1999)

L_h ist kontextfrei, aber nicht regulär.

Theorem (Head, Păun und Pixton 2010)

Falls die Sprache der DNA nicht kontextfrei ist, so ist sie auch nicht kontextsensitiv.

Zu welcher Sprachfamilie gehört die DNA-Grammatik?

Theorem (Searl 1999)

L_h ist kontextfrei, aber nicht regulär.

Theorem (Head, Păun und Pixton 2010)

Falls die Sprache der DNA nicht kontextfrei ist, so ist sie auch nicht kontextsensitiv.

Vermutung

Die Sprache der DNA ist nicht kontextsensitiv.

Diskussion

Noch Fragen 