

Teil VII: Modelle, Algorithmen, Eigenschaften

1. Motivation
2. Grammatiken und Formale Sprachen
3. Formale Sprachen und Endliche Automaten
4. Suchen in Zeichenketten

© K. Murmann, H. Neumann, Fakultät für Informatik, Universität Ulm, 2000

Motivation

Suche nach Zeichenkombinationen

Beispiel

„Richard Feynman war sehr amerikanisch. Wie das große Land selbst, an dessen Ostküste er 1918 geboren [] und an dessen Ostküste er 1988 starb, so steckt [] ch voller Gegensätze. Während sich die USA zum Beispiel als Land von Mickey Mouse und Mondfahrt charakterisieren lassen, kann man Feynman als genialen Physiker und großen Kindsopf zugleich beschreiben. So wie sich in Amerika sowohl die größte Zahl der Analphabeten in der westlichen Welt als auch die meisten Nobelpreisträger in den Naturwissenschaften finden lassen, so stellt man in Feynman die höchste Originalität in der Physik und die platteste Banalität in Kunst und Philosophie fest. ...“

(aus E.P. Fischer, Einstein & Co. Piper, München, 1995)

Aufgaben und Fragen:

- Wir haben uns vertippt und **Feymann** statt **Feynman** geschrieben. Wir wollen die Fehler korrigieren !
Lösungsidee : Überprüfung des Textes von Anfang bis Ende auf Zeichenfolgen, die mit **Feymann** übereinstimmen
- Die Jahreszahlen **19xx** müssen auf Richtigkeit überprüft werden !
Lösungsidee : Überprüfung des Textes auf 4-ziffrige Zeichenfolgen mit **19** am Anfang sowie **2** weiteren beliebigen Ziffern

Arithmetische Ausdrücke

Aufgabe: Für geklammerte arithmetische Ausdrücke muss entschieden werden, ob sie korrekt gebildet sind !

- Es sei:
- S ist eine **Zeichenmenge (Alphabet)**, $S = \{ (,), +, -, *, /, a \}$
(a ist Platzhalter für beliebige Konstanten oder Variablen)
 - eine **(formale) Sprache L** (über S) ist eine Teilmenge aus der Menge der **Konkationen** der Zeichen aus S (bezeichnet als S^*)
- hier: $EXPR \subseteq S^*$

Einige arithmetische Ausdrücke:

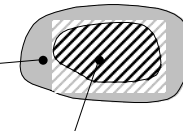
$$(a - a) * a + a / (a + a) - a \in EXPR$$

$$(((a))) \in EXPR$$

$$((a+) - a) \notin EXPR$$

$$a * / a \notin EXPR$$

Menge aller Sprachen
(~ freie Sprache)

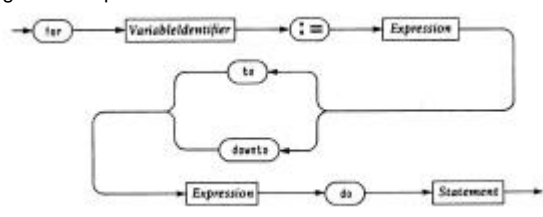


Zusammenfassung aller
korrekten Konstruktionen

Problem: Sprachen sind allgemein unendliche Objekte (i.S. beliebig langer und strukturierter Zeichenfolgen), die selbst endlich beschrieben werden müssen !
→ **Grammatiken** und **Automaten**

Interpretation / Entwurf syntaktischer Beschreibungen

geg.: 1. **Syntaxdiagramm** für Wiederholungsschleifen in der Programmiersprache **Pascal**



2. **EBNF-Syntax** für die Spezifikation von Zahlen in **Pascal**

Definition:

```

UnsignedNumber = UnsignedInteger | UnsignedReal .
UnsignedInteger = DigitSequence .
UnsignedReal = UnsignedInteger "." DigitSequence
               [ "*" ScaleFactor ] |
               UnsignedInteger "e" ScaleFactor .
ScaleFactor = [ Sign ] UnsignedInteger .
Sign = "+" | "-" .
DigitSequence = Digit | Digit | .
    
```

Beispiele:

Examples of unsigned integers:
1 100 00100

Examples of unsigned reals:
0.1 0.1e0 87.35e+8 1E2

Einige Fragen und Anwendungen:

1. Dem Buch über die (Programmier-) Sprache „Y“ liegt eine **Syntaxbeschreibung** bei. **Wie können wir diese lesen und verstehen ?**
2. Sie haben für eine **interaktive Software eine Benutzerschnittstelle** entwickelt und müssen nun für die Benutzer die Formate der erlaubten Kommandos mit ihren Parametern spezifizieren. **Wie kann man dies sauber (am besten) beschreiben ?**
3. Sie wollen eine **Editierhilfe für ein Softwarewerkzeug** entwickeln und nehmen dabei Anleihen aus dem bekannten vi-Editor (so soll z.B. in einer Eingabe ein Textstück gesucht werden können mit `:/<string>/`). **Wie lassen sich derartige Kommandos beschreiben und anschliessend realisieren ?**
4. Wir haben eine **eigene (Programmier-) Sprache** entwickelt. **Wie kann man einen Compiler für diese Sprache definieren, mit dem Programme in meiner neuen Sprache automatisch in ausführbaren Code übersetzt werden ?**

Grammatiken und Formale Sprachen

Grammatik und Sprache

Def.:

Grammatik

$$G = (S_N, S_T, P, w_s)$$

S_N : Menge der **Nichtterminal-Zeichen (Variablen)**

S_T : Menge der **Terminal-Zeichen**

$P = \{p_i\}$, $|P| = n < \infty$: **Produktions- / Ableitungs- / Ersetzungsregeln**

$$P \subseteq (S_N \cup S_T)^{\pm} \times (S_N \cup S_T)^*$$

w_s : **Startvariable**, $w_s \in S_N$ ($|w_s| = 1$)

es können nur Zeichen echt ersetzt werden !



Def.:

Sprache

$$L(G) = \{w \in S_T^* \mid S \xrightarrow{*}_G w\}$$

mit $\xrightarrow{*}_G$: endliche Folgen von Ableitungen (= (Einzel-) Ersetzungen) \rightarrow_G

Erläuterungen:

Achtung: Die Bezeichnung hat nichts mit den **Variablen** (im Gegensatz zu **Konstanten**) eines **Programms** zu tun !

1. Die **Variablen** legen die Menge der **nicht-terminalen Symbole** fest
2. Die **Terminalsymbole** bilden die Menge der „eigentlichen“ Symbole; jede Zeichenfolge einer mittels der Produktionen abgeleiteten Konstruktionen enthält nur Terminalsymbole !
3. Die **Regeln** der schrittweisen **Zeichenersetzungen** werden durch die **Produktionen** (Ableitungen) festgelegt; sie haben allgemein die Form:

linke Seite \longrightarrow rechte Seite

die erste Ersetzung beginnt stets bei der **Startvariablen** !

Ableitung einer Zeichenfolge (Wortproblem)

Problem :Es liegt eine Zeichenfolge w vor !

Frage : **Gehört w zur Sprache, d.h. ist** $w \in R(S)$

\longrightarrow Dieses Entscheidungsproblem ist **nicht** für beliebige Grammatiken lösbar, jedoch für solche mit bestimmten Ersetzungsregeln !

Strategien :

1. **Top-Down**
 \rightarrow ausgehend von der Startvariable wird solange abgeleitet, bis die gegebene Folge produziert wird !
2. **Bottom-Up**
 \rightarrow ausgehend von der gegebenen Folge werden die Regeln „rückwärts“ angewandt, bis die Startvariable erreicht ist !

Bsp. :

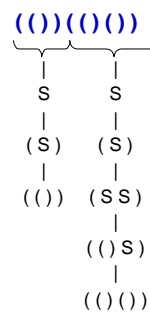
$G = (\{ S \},$
 $\{ (,) \},$
 $\{ S \xrightarrow{1} (), S \xrightarrow{2} (S), S \xrightarrow{3} SS \},$
 $S)$

„Ausgeglichene Klammern“

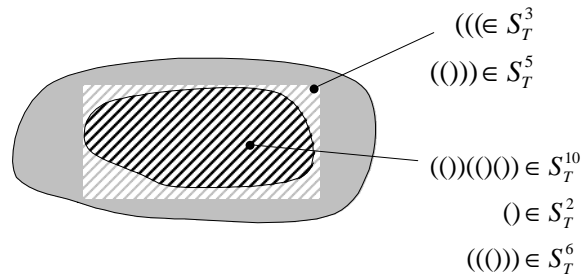
Frage : Gehört $((()))((()())$ zur Sprache ?

Top-Down Ableitung :

$S \xrightarrow{3} SS$
 $\rightarrow \xrightarrow{2} (S) S$
 $\rightarrow \xrightarrow{1} (() S$
 $\rightarrow \xrightarrow{2} (() (S)$
 $\rightarrow \xrightarrow{3} (() (SS)$
 $\rightarrow \xrightarrow{1} (() (() S)$
 $\rightarrow \xrightarrow{1} (() (() ())$



Aus der Menge aller durch Konkatination aus Klammern über $S_T = \{ (,) \}$ gebildeter Zeichenketten ($S_T^* = \{ ()^* \}$) werden **nur die ausgeglichenen (balancierten) Konstruktionen** generiert !



Konkrete Anwendungen :

- Prüfung der Gültigkeit geschachtelter Ausdrücke (**BEGIN** \equiv (, **END** \equiv))
- Editor mit automatischer Prüfung ausgewogener Ausdrücke während der Eingabe (... ((...)

Chomsky-Hierarchie

Def.:

Typ 0 – Grammatik

- jede Grammatik ist zunächst automatisch Typ 0 !
- Produktionen $(w_1 \rightarrow w_2 \in P)$ **keine Einschränkungen**
(→ allgemeine Phrasenstruktur-Grammatiken)



Typ 1 – Grammatik

- **kontextsensitiv**
- Produktionen $(w_1 \rightarrow w_2 \in R)$ $|w_1| \leq |w_2|$

Kontextsensitivität (Bsp.) :

Produktionsregeln der Art

$$\begin{array}{l} aAb \rightarrow agb \\ dAh \rightarrow dfh \end{array} \quad \text{mit} \quad |w_1| = |w_2|$$

Def.:

Typ 2 – Grammatik

- **kontextfrei**
- Produktionen $(w_1 \rightarrow w_2 \in P)$ $|w_1| = 1, w_1 \in S_N$
d.h. Produktionen sind Einzelzeichenersetzungen

Typ 3 – Grammatik

- **regulär**
- Produktionen $(w_1 \rightarrow w_2 \in R)$ $|w_1| = 1, w_1 \in S_N$

und zusätzlich: $w_2 \in S_T \cup S_T S_N$

d.h. die Ableitungen sind alle von der Form

- $r : B \rightarrow g$ einzelne Terminalzeichen, oder
- $r : B \rightarrow gA$ ein Terminalzeichen, gefolgt von einer Variablen
(Linksableitung)

(→ bei Rechtsableitungen ist die 2. Ersetzungsregel stets $r : B \rightarrow Ag$)

Eine **Sprache** $L \subseteq S_T^*$ ist vom Typ 0, 1, 2, 3, falls es eine Typ 0, 1, 2, 3 – Grammatik G mit $L(G) = L$ gibt !

Syntaxbäume (Ableitungsbäume)

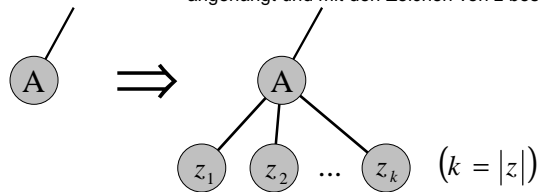
- Ableitungen** eines Wortes w aus einer **Typ 2** (oder **Typ 3**) Grammatik G kann man einen **Ableitungsbaum** zuordnen !

geg.: $w \in L(G)$,

$$S = w_0 \Rightarrow w_1 \Rightarrow w_2 \Rightarrow \dots \Rightarrow w_n = w$$

ist eine Ableitung des Wortes w

$i = 0$ Startvariable S ist die Wurzel des Ableitungsbaums
 $i = 1, 2, \dots, n$ für j -ten Ableitungsschritt (Übergang $w_{j-1} \Rightarrow w_j$)
 mit Ersetzung $(A \rightarrow z) \in P$
 werden im Baum am Knoten A $|z|$ Söhne
 angehängt und mit den Zeichen von z beschriftet

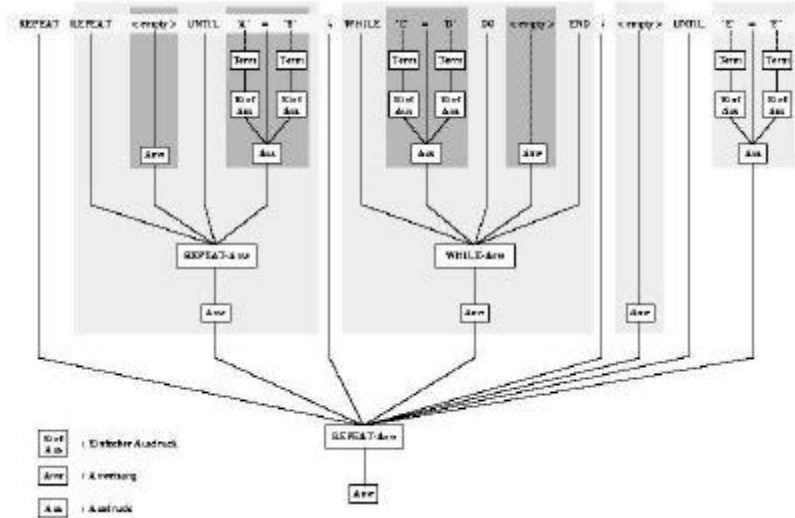


Bsp.: Ableitung einer MODULA-2 Programmsequenz
(zur Anschauung ...)

```

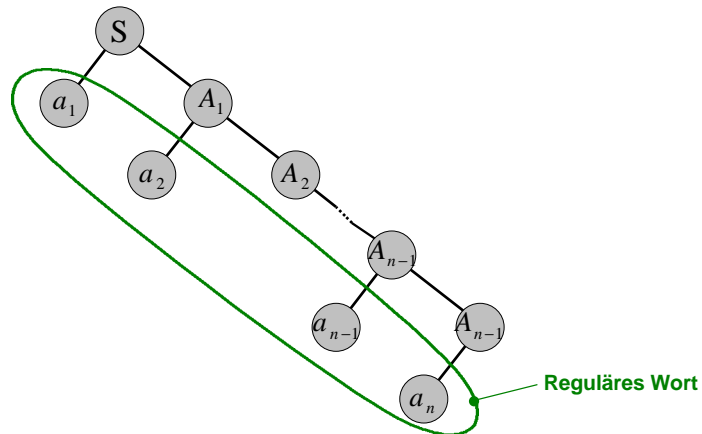
REPEAT
  REPEAT
  UNTIL A = B;
  WHILE C = D DO
  END;
UNTIL E = F
  
```

Achtung
Endlosschleife



- Ableitungsbäume bei **regulären Grammatiken** sind stets „entartet“, d.h. es werden lineare Ketten mit geeigneten Bäumen erzeugt !

hier: Linksableitung



Backus-Naur-Form

Backus-Naur-Form (BNF)

→ Formalismus zur Darstellung einer **kontextfreien Grammatik** (Typ 2)

Die Darstellung von Regeln mit derselben linken Seite wird zusammengefasst

$$\left. \begin{array}{l} A \rightarrow \mathbf{b}_1 \\ A \rightarrow \mathbf{b}_2 \\ \vdots \\ A \rightarrow \mathbf{b}_n \end{array} \right\} \text{wird zu einer Meta-Regel} \quad A \rightarrow \mathbf{b}_1 \mid \mathbf{b}_2 \mid \dots \mid \mathbf{b}_n$$

mithilfe des Meta-Symbols | (Alternative, Auswahl)

(Hinweis: Backus / Naur verwenden statt \rightarrow das Zeichen $::=$)

Erweiterte BNF (EBNF)

Einführung meta-sprachlicher Nicht-Terminalsymbole, insbesondere **Klammerkonstrukte** :

→	Ableitungssymbol (Original ::=)
.	Markierung des Endes einer Ableitungsregel
	Alternative (Auswahl)
[...]	Optionale Ausdrücke : 0- oder 1-mal auftretend
{...}	Optionale Ausdrücke : 0-, 1- oder mehrmals auftretend
(...)	Präzedenzregelnde (Strukturierungs-) Klammern

Erweiterungen durch die EBNF !

Beispiel: Syntax für ganze Dezimalzahlen (evtl. mit Vorzeichen)

Herkömmliche Notation :

$\langle \text{GanzeZahl} \rangle \rightarrow \langle \text{Zahl} \rangle .$
 $\langle \text{GanzeZahl} \rangle \rightarrow \langle \text{Vorzeichen} \rangle \langle \text{Zahl} \rangle .$
 $\langle \text{Zahl} \rangle \rightarrow \langle \text{Ziffer} \rangle .$
 $\langle \text{Zahl} \rangle \rightarrow \langle \text{Ziffer} \rangle \langle \text{Zahl} \rangle .$
 $\langle \text{Vorzeichen} \rangle \rightarrow + .$
 $\langle \text{Vorzeichen} \rangle \rightarrow - .$
 $\langle \text{Ziffer} \rangle \rightarrow 0 .$
 $\langle \text{Ziffer} \rangle \rightarrow 1 .$
 $\langle \text{Ziffer} \rangle \rightarrow 2 .$
 $\langle \text{Ziffer} \rangle \rightarrow 3 .$
 $\langle \text{Ziffer} \rangle \rightarrow 4 .$
 $\langle \text{Ziffer} \rangle \rightarrow 5 .$
 $\langle \text{Ziffer} \rangle \rightarrow 6 .$
 $\langle \text{Ziffer} \rangle \rightarrow 7 .$
 $\langle \text{Ziffer} \rangle \rightarrow 8 .$
 $\langle \text{Ziffer} \rangle \rightarrow 9 .$

BNF :

$\langle \text{GanzeZahl} \rangle \rightarrow \langle \text{Zahl} \rangle | \langle \text{Vorzeichen} \rangle \langle \text{Zahl} \rangle .$
 $\langle \text{Zahl} \rangle \rightarrow \langle \text{Ziffer} \rangle | \langle \text{Ziffer} \rangle \langle \text{Zahl} \rangle .$
 $\langle \text{Vorzeichen} \rangle \rightarrow + | - .$
 $\langle \text{Ziffer} \rangle \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 .$

EBNF :

$\langle \text{GanzeZahl} \rangle \rightarrow [+ | -] \langle \text{Ziffer} \rangle \{ \langle \text{Ziffer} \rangle \}$
 $\langle \text{Ziffer} \rangle \rightarrow 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9$

Syntaxdiagramme

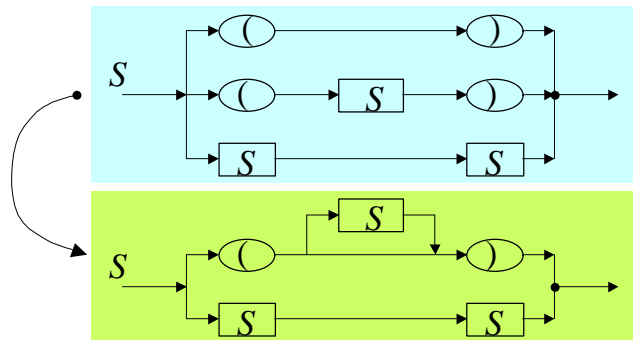
Umsetzung : Produktionen \rightarrow Syntaxdiagramm

Bsp.: „Ausgeglichene Klammern“

1. Produktionen und EBNF

$$P \rightarrow \{S \rightarrow (), S \rightarrow (S), S \rightarrow SS\} \iff P \rightarrow (\langle S \rangle) | \langle S \rangle \langle S \rangle$$

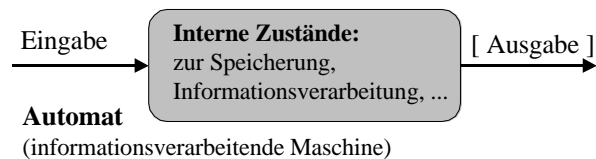
2. Diagramme



Formale Sprachen und Endliche Automaten

Endliche Automaten

Allgemeines Modell eines Automaten



Def.:

Deterministischer endlicher Automat M (DEA)
 („discrete finite automaton“, DFA) **ohne Ausgabe**

$$M = (Z, S, d, z_0, T)$$

Z : Menge der Zustände , $|Z| = n + 1 < \infty$

S : Eingabe-Alphabet , $|S| = m < \infty$

$d : Z \times S \rightarrow Z$: Zustandsübergangs-Funktion

$z_0 \in Z$: Anfangs- / Startzustand

$T \subseteq Z$: Menge der End- / Terminalzustände

Zustandstafel und Zustandsübergangsfunktion (Tabellenform) :

Eingabe \ Zustand	s_1	...	s_k	...	s_m
z_0					
\vdots					
z_i			z_j		
\vdots					
z_n					

→ Automat sei im Zustand z_i und erhält die Eingabe s_k , dann geht der Automat in den Zustand z_j über!

Graphische Darstellung : Zustandsdiagramm / -graph

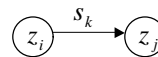
▪ **Knoten (= Zustände)**

- Anfangs- / Startzustand (offener Eingang)
- End- / Terminalzustand (Doppelkreis)

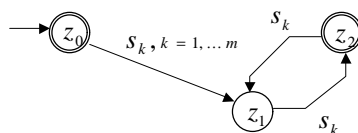


▪ **Kanten**

- Zustandsübergang $d(z_i, s_k) = z_j$



Einfaches Beispiel:



Für welche Eingaben gelangt der Automat in einen Endzustand ?

Akzeptor einer Sprache

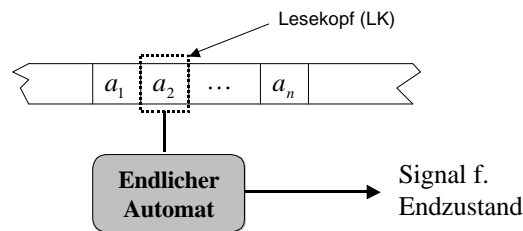
Ein endlicher Automat M beschreibt (erkennt, **akzeptiert**) eine Sprache $L \subseteq S^*$

- ein Wort $a_1 a_2 \dots a_n$ wird elementweise gelesen
- der Automat durchläuft schrittweise eine Folge von Zuständen z_1, z_2, \dots, z_n

z_0 : Startzustand
 $d(z_{i-1}, a_i) = z_i$, für $i = 1, \dots, n$

- es muss abschliessend gelten

$z_n \in T$



Def.: Erweiterung der Zustandsübergangsfunktion von Einzelzeichen zu Wörtern

geg.: $M = (Z, S, d, z_0, T)$

Induktive Festlegung einer neuen Funktion

$$\hat{d} : Z \times S^* \rightarrow Z$$

mit $\hat{d}(z, \epsilon) = z$

$$\hat{d}(z, ax) = \hat{d}(d(z, a), x) \text{ wobei } z \in Z, x \in S^*, a \in S$$

Es gilt : $\hat{d}(z, a) = d(z, a), a \in S$

$$\hat{d}(z, a_1 a_2 \dots a_n) = d(\dots d(d(z, a_1), a_2) \dots, a_n)$$

→ vgl. Ableitungsregeln von **regulären Sprachen** !

Die von M akzeptierte Sprache ist

$$T(M) = \{x \in S^* \mid \hat{d}(z_0, x) \in T\}$$

Satz: Jede durch endliche Automaten erkennbare Sprache ist regulär (also Typ 3)



(ohne Beweis !)

Bsp.: geg.: $M = (S, Z, d, z_0, T)$

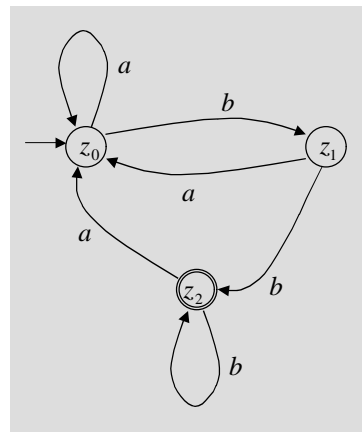
$$S = \{a, b\}$$

$$Z = \{z_0, z_1, z_2\}$$

$$T = \{z_2\}$$

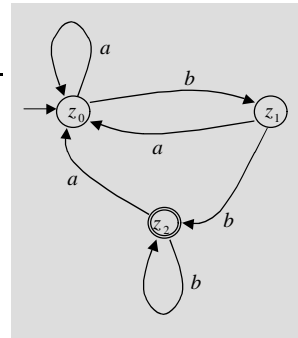
d :

	a	b
z_0	z_0	z_1
z_1	z_0	z_2
z_2	z_0	z_2



Bsp. (cont'd)

Eingaben	End-Zustand
<i>aab</i>	$\hat{d}(z_0, aab) = z_1$
<i>baba</i>	$\hat{d}(z_0, baba) = z_0$
<i>abb</i>	$\hat{d}(z_0, abb) = z_2$
<i>abbabb</i>	$\hat{d}(z_0, abbabb) = z_2$
<i>bba</i>	$\hat{d}(z_0, bba) = z_0$



Frage : Was leistet der Automat ?

⇒ Akzeptor für ... bb !

Simulation

Algorithmus eines deterministischen endlichen Automaten (DEA):

BEGIN

Bringe EA in Zustand z_0 ;

Setze LK über linkes Zeichen des Eingabewortes (falls vorhanden);

WHILE Zeichen unter LK vorhanden **DO**

Gehe in Folgezustand gemäss $d : Z \times S \rightarrow Z$;

Bewege LK um ein Feld nach rechts

END (* WHILE*)

END.

Reguläre Ausdrücke

Def.:

Reguläre Ausdrücke (induktiv)

- die leere Menge $\{\}$ ist ein regulärer Ausdruck
- das leere Wort ϵ ist ein regulärer Ausdruck
- jedes Terminalsymbol $a \in \Sigma$ ist ein regulärer Ausdruck
- geg. : reguläre Ausdrücke a, b
dann sind ab (Konkatenation), $(a|b)$ (Auswahl), $(a)^*$ (Folge) reguläre Ausdrücke



→ **Meta-Symbole:** (), * , |

Zuordnung einer **regulären Sprache $L(\gamma)$** (γ ist ein regulärer Ausdruck)

- $g = \{\}$ $\Rightarrow L(g) = \{\}$
- $g = \epsilon$ $\Rightarrow L(g) = \{\epsilon\}$
- $g = a$ $\Rightarrow L(g) = \{a\}$
- $g = ab$ $\Rightarrow L(g) = L(a)L(b)$
- $g = (a|b)$ $\Rightarrow L(g) = L(a) \cup L(b)$
- $g = (a)^*$ $\Rightarrow L(g) = L(a)^*$

Bsp.: geg. : regulärer Ausdruck $(0|(0|1)^*00)$

In der Sprache enthalten sind alle Formeln mit

- $x = 0$ oder
- x endet mit $\dots 00$

Grammatik :

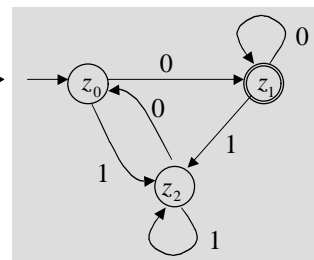
$G = (\{A, B, C, D\}, \{0, 1\},$
 $\{A \rightarrow 0 | B, B \rightarrow C0, C \rightarrow D0, D \rightarrow \lambda, D \rightarrow D1, D \rightarrow D0\},$
 $A)$

Automat als Akzeptor :

$M = (\{z_0, z_1, z_2\}, \{0,1\}, d, z_0, \{z_1\})$

$d :$

	0	1
z_0	z_1	z_2
z_1	z_1	z_2
z_2	z_0	z_2



Konkrete Anwendungen :

- Zur Erinnerung (Teil IV / 3 – Modularisierung und Programmentwicklung) : Telegramme

Eingabe-Format als regulärer Ausdruck :

```
( zzzzzz (b)+ ((c)*(b))* <STOP> (b)+ <STOP> (b)+ )+
```

Ende-Marke

mit $z \in \{0, \dots, 9\}$, b : Leerzeichen, c : Druckzeichen,
 (.)+ : Metazeichen zur 1- oder mehrmaligen Wiederholung,
 (.)* : Metazeichen zur 0- oder mehrmaligen Wiederholung

Beispiel :

```
635444 WEIHNACHTEN IST NAHE STOP STOP 634560 IM
WINTERSEMESTER GIBT ES EINE ERHOLUNGSPAUSE
IN DER MITTE DER VORLESUNGSZEIT STOP UND BALD
IST ES SOWEIT STOP STOP
```

- Reguläre Ausdrücke und UNIX-Kommandos

c	Terminalzeichen c
\c	Escapesequenz, oder falls c Meta-Zeichen, c als Terminalsymbol
^	Anfang (Zeile, Wort) – Meta-Zeichen
\$	Ende (Zeile, Wort) – Meta-Zeichen
.	beliebiges Einzel-Zeichen – Meta-Zeichen
[abc...]	ein Terminalzeichen aus der Aufzählung zwischen []
[^abc...]	alle Zeichen als Terminalzeichen außer den Aufgezählten in []
[a-f]	alle Zeichen zwischen a und f nach ASCII-Ordnung als Terminalzeichen
[^a-f]	alle Zeichen außer denen zwischen a und f als Terminalzeichen
r1 r2	r1 oder r2
(r1)(r2)	r1 gefolgt von r2
(r)*	beliebige viele r's hintereinander (auch keins)
(r)+	beliebig viele, aber mindestens ein r
(r)?	ein oder kein r
(r)	Klammern dienen der Gruppierung – Meta-Zeichen

Metazeichen sind also: () [] \ . ^ + * ? \$

- UNIX-Kommando **egrep** : Sucht in einer Datenmenge (Standardeingabe, Datei)

Aufruf: **egrep options pattern {file}**

pattern regulärer Ausdruck (mit Meta-Symbolen gebildet)

Ausgabe der Zeilennummer Umlenken der Standardausgabe In Datei

Bsp.1: Der Beispieltext vom Anfang (plus einiger Zeilen ...) – Datei „SuchText.txt“

```
Richard Feynman war sehr amerikanisch. Wie das grosse Land selbst, an dessen Ostkueste er 18 geboren wurde und an dessen Westkueste er 88 starb, so steckte auch er voller Gegensaeetze. Waehrend sich die USA zum Beispiel als Land von Mickey Mouse und Mondfahrt charakterisieren lassen, kann man Feynman als genialen Physiker und grossen Kindskopf zugleich beschreiben. So wie sich in Amerika sowohl die groesste Anzahl der Analphabeten in der westlichen Welt als auch die meisten Nobelpreistraeger in den Naturwissenschaften finden lassen, so stellt man in Feynman hoechste Originalitaet in der Physik und die platteste Banalitaet in Kunst und Philosophie fest.
-----
Noch'n paar Zeilen ...
Ab 15 Uhr gibt's Bier bis 23 Uhr -- danach Kaffee und Saft.
```

UNIX-Kommando: `egrep -n SuchText.txt > 19.doc`

```
2:Land selbst, an dessen Ostkueste er 18 geboren wurde
3:und an dessen Westkueste er 88 starb, so steckte auch
15:Ab 15 Uhr gibt's Bier bis 23 Uhr -- danach Kaffee und Saft.
```

UNIX-Kommando: `egrep -n SuchText.txt > 1900.doc`

```
2:Land selbst, an dessen Ostkueste er 18 geboren wurde
3:und an dessen Westkueste er 88 starb, so steckte auch
```

Bsp.2: Ein Modula-2 Programm – Datei „power.m2“

```
MODULE power;
FROM InOut IMPORT WriteString, WriteLn;
FROM RealInOut IMPORT WriteFloat;

VAR
  n          : INTEGER;
  power, base : REAL;
BEGIN
  [redacted]
  WHILE n >= 1 DO
    IF ODD(n) THEN
      [* -- n ungerade -- *]
      [redacted]
    ELSE
      [* == n gerade == *]
      [redacted]
    END
  END
END;

WriteString("Ergebnis : "); WriteFloat(power, 5, 5); WriteLn
END power.
```

Programm berechnet den Wert $1.1^{15} = 4.17725$

Aufgabe: **Alle Zuweisungen anzeigen !**

UNIX-Kommando: `egrep -n power.m2 > powerAssign.doc`

```
10: base := 1.1;
11: n := 15;
12: power := 1.0;
18: power := power * base;
19: n := n - 1;
22: n := n DIV 2;
23: base := base * base
```

```

17: (* -- n ungerade -- *)
21: (* -- n gerade -- *)

```

Aufgabe: **Alle (nicht-leeren) String-Konstanten anzeigen !**

UNIX-Kommando: `egrep -n [] power.m2 > powerStringConst.doc`

```

28: WriteString("Ergebnis : "); WriteFloat(power, 5, 5); WriteLn

```

Aufgabe: **Wieviele Leerzeilen ?**

UNIX-Kommando: `egrep -n [] power.m2 | wc -l > powerLeer.doc`

```

6

```

wc options (file)
options: clw - Zeichen, Zeilen, Worte

Aufgabe: **(Reine) Kommentarzeilen ?**

UNIX-Kommando: `egrep -n [] power.m2 > powerComment.doc`

```

17: (* -- n ungerade -- *)
21: (* -- n gerade -- *)

```

Suchen in Zeichenketten

Fragestellung

Beispiel

geg.: String (Text)

„Als Strukturwissenschaften wird man nicht nur die reine und angewandte Mathematik bezeichnen, sondern auch das in seiner Gliederung noch nicht voll durchschaute Gebiet der Wissenschaften, die man mit Namen wie Systemanalyse, Informationstheorie, Kybernetik, Spieltheorie bezeichnet. Sie sind geometrische Mathematik zeitlicher Vorgänge, die durch menschliche Entscheidung, durch Planung, durch Strukturen, die sich darstellen lassen, als seien sie geplant, oder schließlich durch Zufall gesteuert werden. Sie sind also Strukturtheorien zeitlicher Veränderung. Ihr wichtigstes praktisches Hilfsmittel ist der Computer, dessen Theorie selbst eine der Strukturwissenschaften ist. ...“

(aus C.F. v. Weizsäcker, Die Einheit der Natur, dtv, München, 1974)

Pattern (Muster) : Kybernetik

Aufgabe :

Wir wollen ein Programm entwickeln, das in dem angegebenen Text ein bestimmtes (Vergleichs-) Muster sucht und im Falle des Vorhandenseins die Position (Index) des ersten Zeichens des Musters ausgibt !

Randbedingungen für die Realisierung – Spezifikation und Entwurf

a) Wie wird das Ende einer Zeichenkette („string“) gekennzeichnet ?

Realisierung : 0-Character (0C, CHR(0)) (Modula-2)

b) Wie wird die Zeichenkette repräsentiert ?

Vorschlag : Feld (ARRAY) fester maximaler Länge, A[1..maxLENGTH]

c) Was liefert die Vergleichs-Prozedur als Ergebnis ?

Vorschlag : wenn gefunden : pos ∈ {1, 2, ..., maxLENGTH}
wenn nicht gefunden : pos := UNDEFINED (≡ -1)

d) Wie soll das Ergebnis lauten, wenn der gesuchte Text mehrmals vorkommt ?

Vorschlag : Position des 1. Auftretens wird geliefert
(wenn $pos < maxLENGTH - LENGTH(pattern)$, dann kann für den restlichen Text nochmals das Muster gesucht werden !)

Einfacher Lösungsansatz

Struktur

geg.: String (Text)

PETER BAUER HOERT BACH UND SINGT (0C)

Muster

BACH

▪ **Regulärer Ausdruck**

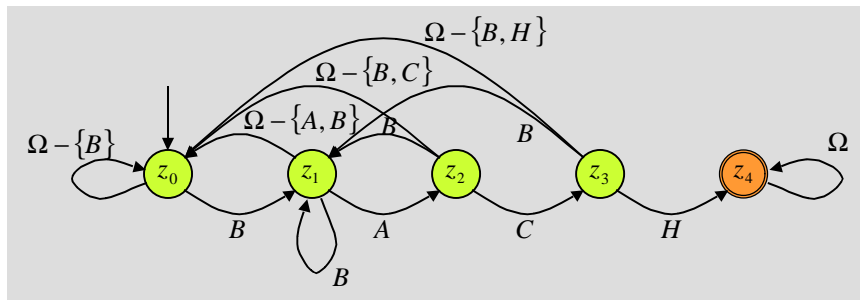
→ vgl. Struktur der „String-Anfragen“ in UNIX

UNIX-Kommando: `egrep 'BACH' string.txt`

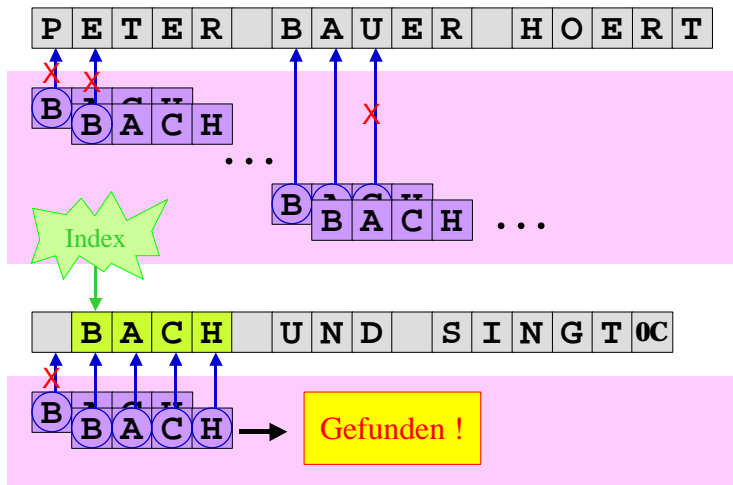
▪ **Automat** $M = (S, Z, d, z_0, T)$

Ziel : Definition eines Akzeptors für eine Zeichenkette, die als Teilkette mindestens 1-mal die Folge „BACH“ enthält !

Eingabe : $S = \{A, B, \dots, Z\} = \Omega$
 Zustände : $Z = \{z_0, z_1, z_2, z_3, z_4\}$
 Terminalzustand : $T = \{z_4\}$



▪ Ablauf



Programm – Realisierung des Automaten

Sequentielle Lösung der „Pattern matching“ Aufgabe

Besonderheiten :

- File I/O (Dateien)
 - neuer Datentyp FILE
 - Dateioperationen :
 - Fopen(...)
 - Fclose(...)
 - Eingabe :
 - FreadChar(...)
 - Done (~ Flag zur Rückmeldung der erfolgreichen Leseoperation)
- Ein-/Ausgabe-Ströme als Dateien
 - stdin

```
MODULE Suchmuster;
FROM SIO2 IMPORT s2;
FROM InOut IMPORT ReadString, WriteInt, WriteInt, WriteString, Write, Read;
FROM FILE2 IMPORT Done, FreadChar;
FROM SIO1 IMPORT s1In, read, FILE, Fopen, Fclose;
FROM Strings IMPORT StrLen;

CONST
  MAXLENGTH = 10000;
  notDone = -1;

TYPE
  string = ARRAY [1..MAXLENGTH] OF CHAR;

VAR
  filename, text, pattern : STRING;
  index : INTEGER;
  file : FILE;
  result : BOOLEAN;
  ch : CHAR;

PROCEDURE ReadLine(
  file : FILE;
  VAR s : ARRAY OF CHAR; : BOOLEAN);
VAR
  i : INTEGER;
  ch : CHAR;
BEGIN
  i := 0;
  FreadChar(file, ch);
  WHILE Done & i <= MAXLENGTH DO
    s[i] := ch;
    INC(i);
    FreadChar(file, ch);
  END;
  s[i] := 0;
  IF Done & i <= MAXLENGTH THEN
    RETURN TRUE;
  ELSE
    RETURN FALSE;
  END;
END ReadLine;

PROCEDURE ReadText(
  file : FILE;
  VAR s : ARRAY OF CHAR; : BOOLEAN);
VAR
  i : INTEGER;
  ch : CHAR;
BEGIN
  i := 0;
  FreadChar(file, ch);
  WHILE Done & i <= MAXLENGTH DO
    s[i] := ch;
    INC(i);
    FreadChar(file, ch);
  END;
  s[i] := 0;
  IF Done THEN
    RETURN TRUE;
  ELSE
    RETURN FALSE;
  END;
END ReadText;

PROCEDURE MatchSimple;
VAR
  s : string;
  pattern : ARRAY OF CHAR;
  VAR index : INTEGER;

VAR
  gpos, spos, spos : INTEGER;
  index := notDone;
  gpos := 0;
  spos := 0;
  spos := 0;
  WHILE | pattern[gpos] <> 0C | & ( s[pos] <= 0C | 0C ) DO
    IF | s[pos] = pattern[gpos] | THEN
      gpos := gpos + 1;
      spos := spos + 1;
    ELSE
      gpos := gpos - spos + 1;
      spos := 0;
    END;
  END;
  IF | pattern[gpos] = 0C | THEN
    index := spos - gpos;
  ELSE
    index := notDone;
  END;
END MatchSimple;

BEGIN (* -- Hauptprogramm "Suchmuster" -- *)
  WriteString("Bitte den Namen der zu durchsuchenden Datei eingeben : ");
  ReadLine(filename);
  IF Fopen(file, filename, read, TRUE) THEN
    WriteString("Bitte Suchmuster eingeben : ");
    IF ReadLine(text, pattern) THEN
      MatchSimple(text, pattern, index);
      IF index = notDone THEN
        WriteString("**** Muster nicht gefunden : ****");
        WriteLn;
      ELSE
        WriteString("Muster an Position ");
        WriteInt(index, 5);
        WriteString(" in Text gefunden !");
        WriteLn;
      END;
    ELSE
      WriteString("**** Fehler: konnte Text nicht lesen : ****");
      WriteLn;
    ELSE
      WriteString("**** Fehler: konnte Datei nicht öffnen : ****");
      WriteLn;
    END;
  ELSE
    WriteString("**** Fehler: konnte Datei nicht öffnen : ****");
    WriteLn;
  END;
END Suchmuster.
```

Aufwand („Laufzeit“)

Länge des **Strings** N

Länge des **Musters** M

a) Im schlechtesten Fall („worst case“):

P	E	T	E	R		B	A	U	E	R		H	O	E	R	T
B	A	C	H													

Für jede Position des **Strings** muss jedes Zeichen des **Musters** getestet werden !

⇒ $M * N$ Operationen

b) Normal-Fall („average case“ ~ mittlerer Aufwand):

In den meisten Fällen wird der Vergleich des **Strings** mit dem **Muster** (~ innere Schleife) frühzeitig abgebrochen !

⇒ $M + N$ Operationen

Boyer-Moore Algorithmus

Quelle : R.S. Boyer, J.S. Moore (1977). A fast string searching algorithm. Communications of the ACM, **20**, 762-772.

Grundidee

geg.: String (Text)

PETER BAUER HOERT BACH UND SINGT (0C)

Muster

BACH

Einfacher Ansatz : String wird elementweise mit **Muster** verglichen, beginnend mit dem ersten Element

Nachteil : viele Fehlvergleiche wg. fehlender Nutzung der Struktur-Information im Such-**Muster**

Neuer Ansatz

- Anwendung des **Musters** im Vergleich mit **String** jeweils **von rechts nach links** !
- Struktur-Information im **Muster** wird in Form einer **Distanz-Tabelle** gespeichert



- Für den Entwurf von Algorithmen ist es sinnvoll, **Modelle** anzuwenden, deren Funktionalität formal beschrieben werden kann!
Automaten (und ihr Bezug zu den **formalen Sprachen**) dienen als Formalismus zur Beschreibung einfacher Mechanismen als Folge von Zustandsübergängen
- **Grammatiken** bilden Systeme, mit denen Folgen von Zeichenketten **regelmäßig** erzeugt werden können; diese Zeichenketten bilden gültige Worte einer durch die Grammatik definierten **Sprache**;
Grammatiken (und die hieraus definierten Sprachen) bilden je nach Format der **Ersetzungs-** (**Produktions-**) **Regeln** eine Hierarchie (**Chomsky-Hierarchie**)
- Ableitungen **kontextfreier Sprachen** (Typ 2 & 3) lassen sich in Form von **Syntax-** oder **Ableitungsbäumen** darstellen; kompakte Darstellungen sind die **Backus-Naur-Form (BNF)** und **EBNF**
- **Endliche Automaten** (deterministisch und ohne Ausgabe) definieren einfach(st)e Maschinen mit Zuständen und Übergängen zwischen Zuständen; diese können **graphisch** mittels **Zustandsdiagrammen (-graphen)** dargestellt werden
- Endliche Automaten sind Akzeptoren **regulärer** (Typ 3) **Sprachen**
- Mit **regulären Ausdrücken** lassen sich komplexe Formate beschreiben, die als Muster (z.B. bei Such- und Vergleichsaufgaben, Textersetzungen in Editoren, etc.) dienen