

4. Die lexikalische Analyse

- zerlegt Folge von Zeichen in Eingabedatei in Folge von Symbolen (Token)
- Scanner-Sieber-Modul
- Token: Typ und Inhalt
- übliche Token-Typen:
 - reservierte Wörter (if, while, for,...)
 - Bezeichner (x, dauer, name, ..)
 - Literale für ganze Zahlen, Gleitkommazahlen, Strings, Zeichen
 - Trennzeichen (Komma, Semikolon)
 - Klammern
 - Operatoren (=, +, -, &&, ...)

4. Die lexikalische Analyse

var					a,	b:	real	;	NL	a:=	2.	1	;	NL	b:=	a*	a+	7	;													NL
-----	--	--	--	--	----	----	------	---	----	-----	----	---	---	----	-----	----	----	---	---	--	--	--	--	--	--	--	--	--	--	--	--	----

Scanner

Id(var) sep Id(a) com Id(b) col Id(real) sem sep Id(a) bec real(2.1) sem sep Id(b) bec Id(a) mul Id(a) add int(7) sem sep

Sieber

varSY id(1) com id(2) col realSY sem id(1) bec real(2.1) sem id(2) bec id(1) mul id(1) add int(7) sem

4. Die lexikalische Analyse

Tokenklassen:

- Menge aller möglichen Werte : reguläre Sprache
- beschrieben durch regulären Ausdruck

Wiederholung/Definition:

reguläre Sprachen über Alphabet Σ :

- leere Menge \emptyset , leere Sprache $\{\epsilon\}$
- für alle $a \in \Sigma$ ist $\{a\}$ reguläre Sprache
- R, S reguläre Sprachen über Σ , dann auch $R^*, RS, R \cup S$

4. Die lexikalische Analyse

Wiederholung/Definition:

reguläre Ausdrücke (RA) über Alphabet Σ :

- \emptyset ist RA, beschreibt Sprache \emptyset
- ε ist RA, beschreibt Sprache $\{\varepsilon\}$
- $a \in \Sigma$ ist RA, beschreibt $\{a\}$
- r, s RA, die R bzw. S beschreiben, dann
 - $(r | s)$ RA, der $R \cup S$ beschreibt
 - (rs) RA, der RS beschreibt
 - $(r)^*$ RA, der R^* beschreibt

Beispiel:

ab^*a	beschreibt $\{a\}\{b\}^*\{a\}$	$aa, aba, abba, \dots$
$(ab)^*$	beschreibt $\{ab\}^*$	$\varepsilon, ab, abab, \dots$

4. Die lexikalische Analyse

Aufgaben:

Erstellen eines regulären Ausdrucks für folgende Sprachen

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

1. Sprache aller Wörter, die mit a anfangen und mit b enden?
2. Sprache aller Wörter, die wenigstens 3 a enthalten?
3. Sprache aller Wörter, die mind. ein aa oder ein bb enthalten?

Welche Sprache wird durch folgenden RA beschrieben?

1. $a|b$
2. $a(aa)^*b^*$
3. $(ab)^*|(ba)^*$

Gegeben sei das Alphabet $\Sigma = \{0, \dots, 9, +, -, .\}$

Welcher RA beschreibt die Sprache aller Dezimalpunktzahlen?

4. Die lexikalische Analyse

Endliche Automaten:

- math. Maschinen zum Erkennen reg. Sprachen
- bestehen aus (Σ, Q, S, F, T)
 - Eingabealphabet Σ
 - endliche Zustandmenge Q
 - Startzustand $S \in Q$
 - Menge an Endzuständen F aus Q
 - Übergangsrelation T aus $(Q \times \Sigma \times Q)$
- Beispiel
 - $\Sigma = \{a, b\}$, $Q = \{z_0, z_1, z_2, z_3\}$,
 - $S = \{z_0\}$, $F = \{z_3\}$, $T =$

Q-alt	Zeichen	Q-neu
z0	a	z0
z0	b	z1
z1	a	z0
z1	b	z2
z2	a	z0
z2	b	z3
z3	a	z3
z3	b	z3

Welche Sprache erkennt der Automat?

Aufgabe: Implementierung in Java, C

4. Die lexikalische Analyse

Endliche Automaten als Scanner (1):

- reguläre Ausdrücke für alle Token(klassen) erstellen
- zugehörigen NEA konstruieren
- minimalen DEA ableiten (Verfahren von Thompson)
- DEA in Programmiersprache implementieren (s. letzte Übungsaufgabe)

4. Die lexikalische Analyse

Endliche Automaten als Scanner (2):

- startet auf dem ersten noch nicht gelesenen Zeichen
- Folge von Schritten abhängig vom Zustand und vom gelesenen Zeichen
- meldet Finden eines Symbols, wenn in Endzustand und zu nächstem Zeichen kein Zustandsübergang existiert
- existiert kein Übergang ohne Endzustand: letzte Übergänge rückgängig machen bis zuletzt durchlaufenen Endzustand

4. Die lexikalische Analyse

Endliche Automaten als Scanner (3):

Tokenklassen für „richtige“ Programmiersprache z.B.

- Bezeichner: $(a|\dots|z|A|\dots|Z)(a|\dots|z|A|\dots|Z|0|\dots|9)^*$
- Integer: $(1|2|\dots|9)(0|1|2|\dots|9)^*$
- Float: $(1|\dots|9)(0|\dots|9)^*(0|\dots|9)^*$
- Operatoren: $(+|-|^*|/|=|\dots)$
- Anweisungen: $\text{if} | \text{then} | \text{else} | \text{for} | \text{while} | \dots$
- Kommentar: $\#(a|\dots|z|A|\dots|Z|_|-|:|\dots)\#$
- Separator: $;$

Aufgabe: Schreiben Sie ein Java- o. C-Programm, welches eine Sprache mit diesen Token scannt und den Namen der jeweiligen Tokenklasse ausgibt.

4. Die lexikalische Analyse

lex, flex:

- Scanner-Generierungstool
- Input: Datei folgenden Aufbaus (bsp.l)

```
Definitionen
%%
lex-Regeln
%%
benutzerdefinierte
Routinen
```

- Generierung mittels `flex bsp.l`
- Output: Datei namens `lex.yy.c` mit Funktion `yylex()`
- `yylex()` in C-Dateien nutzen;
`gcc prog.c lex.yy.c -lfl -o prog`

4. Die lexikalische Analyse

Definitionen-Teil:

1. Festlegung der Wirtsprache (C, Ratfor): %C bzw. %R
2. C- oder Ratfor-Programmteile
3. C- oder Ratfor-Programmteile in %{%}
4. Reguläre Definitionen
5. Startbedingungen
6. Zeichensatz-Tabellen
7. Verstellen von lex-internen Tabellengrößen

4. Die lexikalische Analyse

C- bzw. Ratfor Programmteile:

entweder

- erstes Zeichen pro Zeile Leerstelle oder Tab
- dahinter C-/Ratfor-Deklarationen/Definitionen
- werden unverändert an globalster Stelle in generiertes Programm übernommen
- Beispiel:

```
int operand[10];  
char op;  
%%  
....
```