

Script generated by TTT

Title: Simon: Compilerbau (28.04.2014)

Date: Mon Apr 28 14:22:08 CEST 2014

Duration: 45:18 min

Pages: 38

Chapter 5: Scanner design

51/61

Scanner design

Input (simplified): a set of rules:
$$\begin{array}{ll}
 e_1 & \{ \text{action}_1 \} \\
 e_2 & \{ \text{action}_2 \} \\
 \dots & \\
 e_k & \{ \text{action}_k \}
 \end{array}$$

Scanner design

Input (simplified): a set of rules:
$$\begin{array}{ll}
 e_1 & \{ \text{action}_1 \} \\
 e_2 & \{ \text{action}_2 \} \\
 \dots & \\
 e_k & \{ \text{action}_k \}
 \end{array}$$
Output: a program,

- ... reading a **maximal prefix** w from the input, that satisfies $e_1 \mid \dots \mid e_k$;
- ... determining the **minimal** i , such that $w \in \llbracket e_i \rrbracket$;
- ... executing **action_i** for w .

Implementation:

Idea:

- Create the DFA $\mathcal{P}(A_e) = (Q, \Sigma, \delta, q_0, F)$ for the expression $e = (e_1 | \dots | e_k)$;

- Define the sets:

$$F_1 = \{q \in F \mid q \cap \text{last}[e_1] \neq \emptyset\}$$

$$F_2 = \{q \in (F \setminus F_1) \mid q \cap \text{last}[e_2] \neq \emptyset\}$$

$$\dots$$

$$F_k = \{q \in (F \setminus (F_1 \cup \dots \cup F_{k-1})) \mid q \cap \text{last}[e_k] \neq \emptyset\}$$

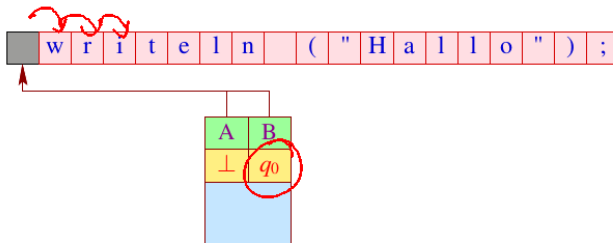
- For input w we find: $\delta^*(q_0, w) \in F_i$ iff the scanner must execute action_i for w

53 / 61

Implementation:

Idea (cont'd):

- The scanner manages two pointers $\langle A, B \rangle$ and the related states $\langle q_A, q_B \rangle \dots$
- Pointer A points to the last position in the input, after which a state $q_A \in F$ was reached;
- Pointer B tracks the current position.



54 / 61

Implementation:

Idea (cont'd):

- The scanner manages two pointers $\langle A, B \rangle$ and the related states $\langle q_A, q_B \rangle \dots$
- Pointer A points to the last position in the input, after which a state $q_A \in F$ was reached;
- Pointer B tracks the current position.

s t d o u t . w r i t e l n (" H a l l o ") ;



54 / 61

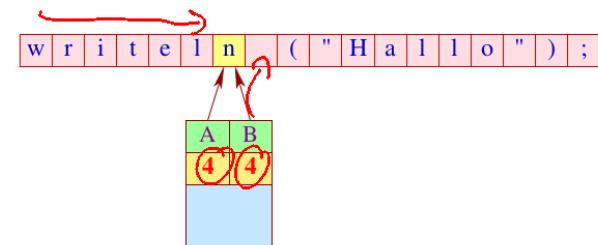
Implementation:

Idea (cont'd):

- The current state being $q_B = \emptyset$, we consume input up to position A and reset:

$$B := A; \quad A := \perp;$$

$$q_B := q_0; \quad q_A := \perp$$

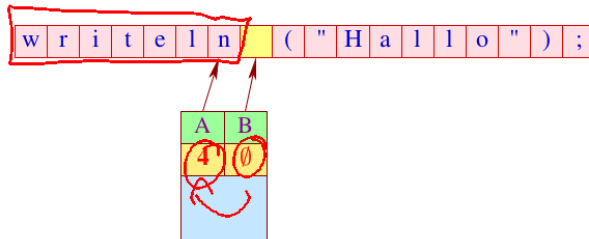


55 / 61

Implementation:

Idea (cont'd):

- The current state being $q_B = \emptyset$, we consume input up to position A and reset:

$$\begin{aligned} B &:= A; & A &:= \perp; \\ q_B &:= q_0; & q_A &:= \perp \end{aligned}$$


55 / 61

Input (generalized): a set of rules:

```

<state> { e1 { action1 yybegin(state1); }
          e2 { action2 yybegin(state2); }
          ...
          ek { actionk yybegin(statek); }
        }
    
```

- The statement `yybegin (statei);` resets the current state to `statei`.
- The start state is called (e.g. flex JFlex) `YYINITIAL`.

... for example:

```

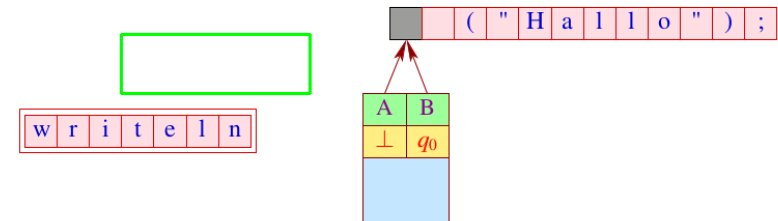
<YYINITIAL> { "/*" { yybegin(COMMENT); }
<COMMENT>  { "*" /" { yybegin(YYINITIAL); }
             . | \n { }
            }
    
```

57 / 61

Implementation:

Idea (cont'd):

- The current state being $q_B = \emptyset$, we consume input up to position A and reset:

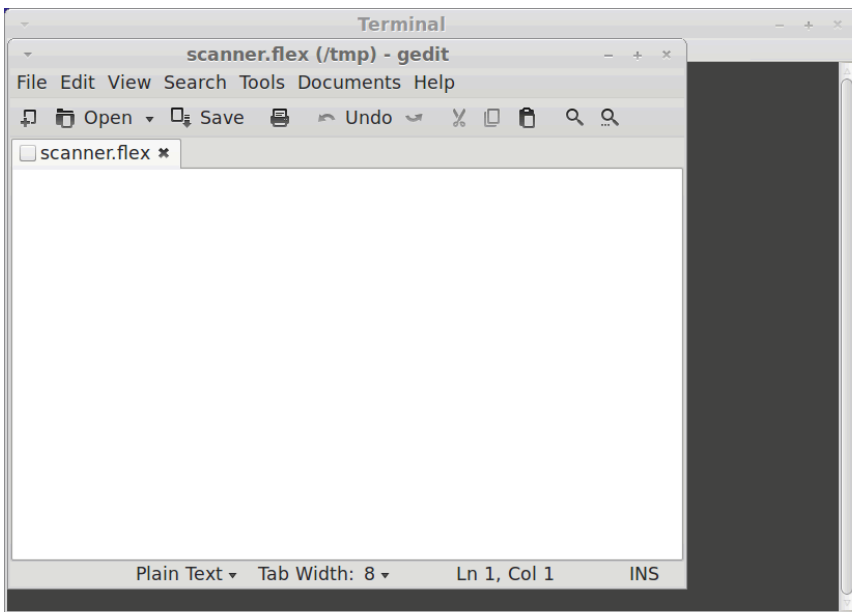
$$\begin{aligned} B &:= A; & A &:= \perp; \\ q_B &:= q_0; & q_A &:= \perp \end{aligned}$$



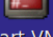
55 / 61

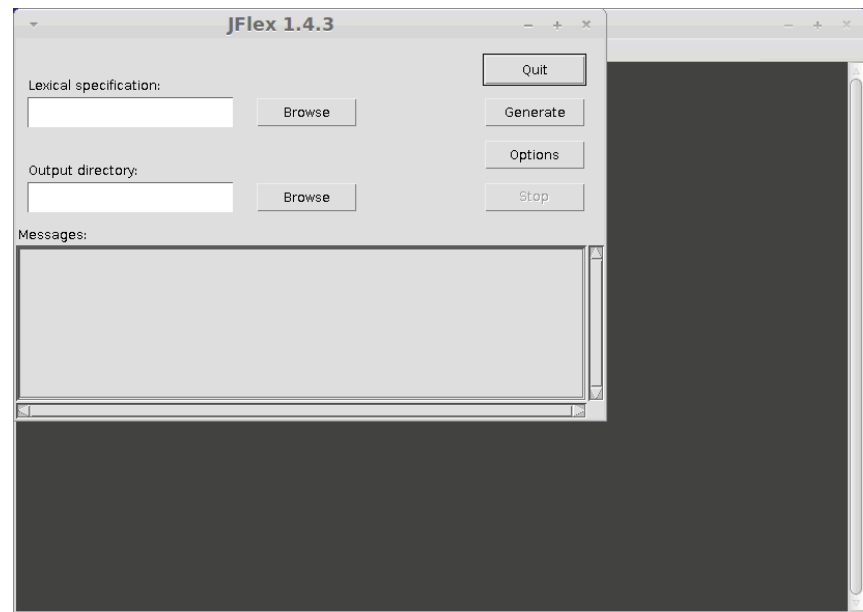
Remarks:



- “.” matches all characters different from “\n”.
- For every state we generate the scanner respectively.
- Method `yybegin (STATE);` switches between different scanners.
- Comments might be directly implemented as (admittedly overly complex) token-class.
- Scanner-states are especially handy for implementing **preprocessors**, expanding special fragments in regular programs.

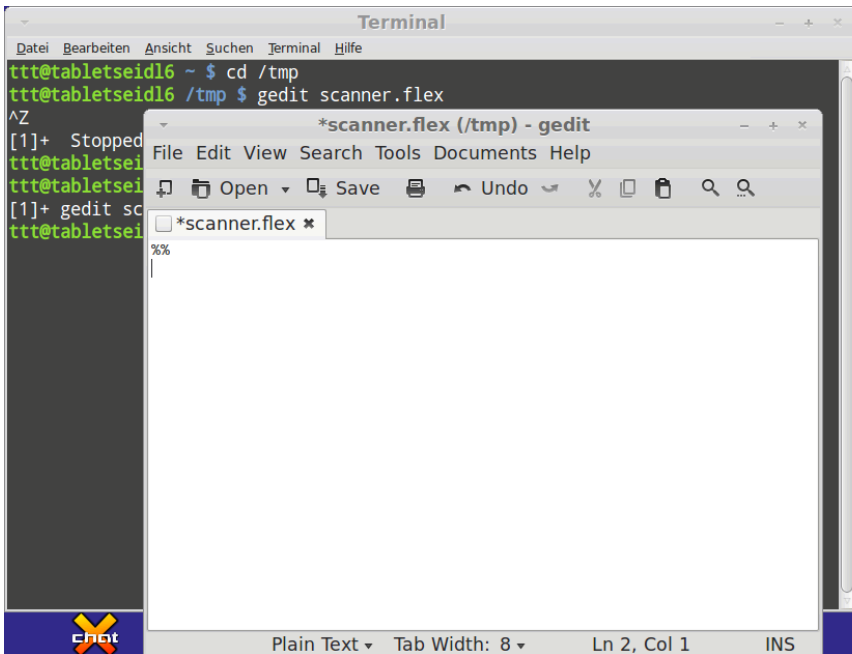
58 / 61



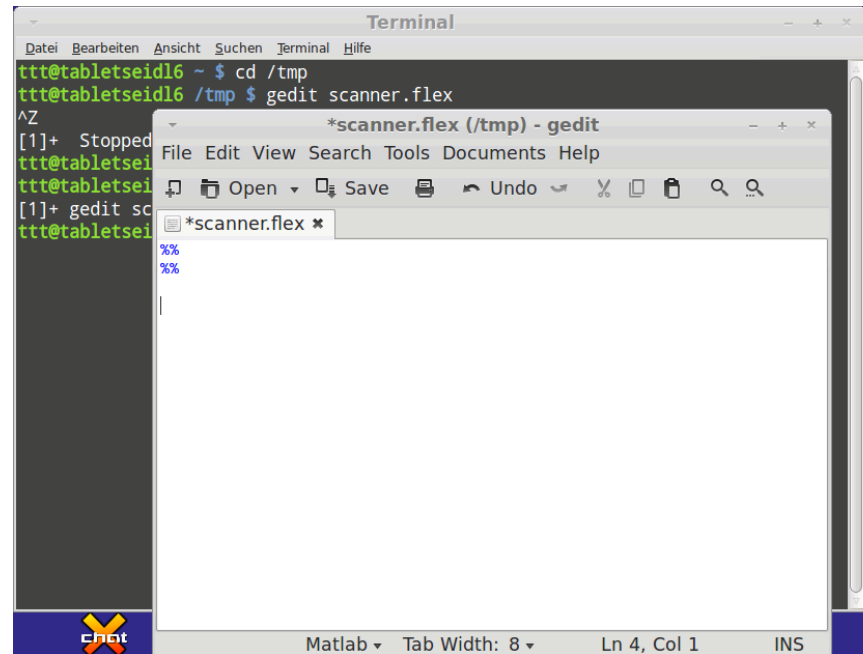
 Kill VNC-Server
 Start VNC-Server





 Kill VNC-Server
 Server



 Kill VNC-Server
 Server



 Kill VNC-Server
 Server


```

Terminal
Datei Bearbeiten Ansicht Suchen Terminal Hilfe
ttt@tabletseidl6 ~ $ cd /tmp
ttt@tabletseidl6 /tmp $ gedit scanner.flex
^Z
[1]+  Stopped                  gedit scanner.flex
ttt@tabletseidl6 ~ $ gedit scanner.flex
[1]+  gedit scanner.flex
ttt@tabletseidl6 ~ $ gedit scanner.flex
Reading "scanner.flex"
Constructing NFA
Converting NFA to LR(0) items
...
6 states before
Writing code
ttt@tabletseidl6 ~ $ gedit scanner.flex

scanner.flex (/tmp) - gedit
File Edit View Search Tools Documents Help
scanner.flex *
%%
Number = [0-9]+
%%
{Number} { System.out.println("Number detected"); }
"+ " { }
" _ " { }
"* " { }
Matlab Tab Width: 8 Ln 8, Col 14 INS
Kill VNC-Server Server

```

```

Yylex.java (/tmp) - gedit
File Edit View Search Tools Documents Help
scanner.flex * Yylex.java *
throw new Error(message);
}

/**
 * Pushes the specified amount of characters back into the input stream.
 * They will be read again by then next call of the scanning method
 * @param number the number of characters to be read again.
 * This number must not be greater than yylength()!
 */
public void yypushback(int number) {
    if ( number > yylength() )
        zzScanError(ZZ_PUSHBACK_2BIG);

    zzMarkedPos -= number;
}

/**
 * Resumes scanning until the next regular expression is matched,
 * the end of input is encountered or an I/O-Error occurs.
 * @return the next token
 * @exception java.io.IOException if any I/O-Error occurs
 */
public Yytoken yylex() throws java.io.IOException {
    int zzInput;
    int zzAction;
}
Java Tab Width: 8 Ln 440, Col 21 INS

```

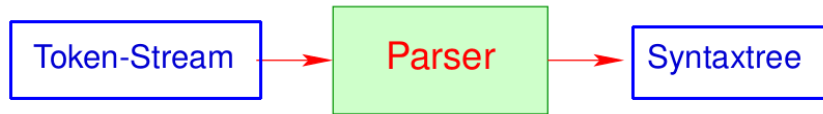
```

*scanner.flex (/tmp) - gedit
File Edit View Search Tools Documents Help
*scanner.flex * Yylex.java *
%%
Number = [0-9]+
%%
{Number} { System.out.println("Number detected"); }
"+ " { return | }
" _ " { }
"* " { }
Matlab Tab Width: 8 Ln 6, Col 19 INS

```

Topic: Syntactic Analysis

Syntactic Analysis



- Syntactic analysis tries to integrate Tokens into larger program units.

60 / 61

Syntactic Analysis



- Syntactic analysis tries to integrate Tokens into larger program units.
- Such units may possibly be:
 - Expressions;
 - Statements;
 - Conditional branches;
 - loops; ...

60 / 61

Discussion:

In general, parsers are not developed by hand, but **generated** from a specification:



61 / 61

Discussion:

In general, parsers are not developed by hand, but **generated** from a specification:



Specification of the hierarchical structure: contextfree grammars
Generated implementation: Pushdown automata + X

61 / 61

Chapter 1: Basics of Contextfree Grammars

62 / 61

Basics: Context-free Grammars

- Programs of programming languages can have arbitrary numbers of tokens, but only finitely many Token-classes.
- This is why we choose the set of Token-classes to be the finite alphabet of terminals T .
- The nested structure of program components can be described elegantly via context-free grammars...

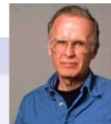
Definition:

A context-free grammar (CFG) is a 4-tuple $G = (N, T, P, S)$ with:

- N the set of nonterminals,
- T the set of terminals,
- P the set of productions or rules, and
- $S \in N$ the start symbol



Noam Chomsky



John Backus

63 / 61

Basics: Context-free Grammars

- Programs of programming languages can have arbitrary numbers of tokens, but only finitely many Token-classes.
- This is why we choose the set of Token-classes to be the finite alphabet of terminals T .
- The nested structure of program components can be described elegantly via context-free grammars...

63 / 61

Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \text{ with } A \in N, \alpha \in (N \cup T)^*$$

64 / 61

Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \text{ with } A \in N, \alpha \in (N \cup T)^*$$

... for example:

$$\begin{array}{l} S \rightarrow aSb \\ S \rightarrow \epsilon \end{array}$$

Specified language: $\{a^n b^n \mid n \geq 0\}$

... further examples:

$$\begin{array}{l} S \rightarrow \langle \text{stmt} \rangle \\ \langle \text{stmt} \rangle \rightarrow \langle \text{if} \rangle \mid \langle \text{while} \rangle \mid \langle \text{rexp} \rangle ; \\ \langle \text{if} \rangle \rightarrow \text{if} (\langle \text{rexp} \rangle) \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ \langle \text{while} \rangle \rightarrow \text{while} (\langle \text{rexp} \rangle) \langle \text{stmt} \rangle \\ \langle \text{rexp} \rangle \rightarrow \text{int} \mid \langle \text{lexp} \rangle \mid \langle \text{lexp} \rangle = \langle \text{rexp} \rangle \mid \dots \\ \langle \text{lexp} \rangle \rightarrow \text{name} \mid \dots \end{array}$$

Conventions

The rules of context-free grammars take the following form:

$$A \rightarrow \alpha \text{ with } A \in N, \alpha \in (N \cup T)^*$$

... for example:

$$\begin{array}{l} S \rightarrow aSb \\ S \rightarrow \epsilon \end{array}$$

Specified language: $\{a^n b^n \mid n \geq 0\}$

Conventions:

In examples, we specify nonterminals and terminals in general implicitly:

- nonterminals are: $A, B, C, \dots, \langle \text{exp} \rangle, \langle \text{stmt} \rangle, \dots;$
- terminals are: $a, b, c, \dots, \text{int, name, } \dots;$