# Programming Languages and Compilers Symbol Specification and Lexical Analysis Dr. Peter Pfahler Based on the lecture by Prof. Dr. Uwe Kastens Universität Paderborn Fakultät EIM Institut für Informatik Winter 2016/2017 PLaC Symbol Specification and Lexical Analysis

**Input:** Program represented by a sequence of characters

Tasks	Compiler Module
Read character stream	Input reader
Recognize and classify tokens;	Scanner (central phase, finite state
skip irrelevant characters	machine)
Encode tokens; Conversion;	Identifier Module, Literal Module
Store token information	

**Output:** Program represented by a sequence of *tokens* encoded by triples:



Lexical Analysis

## Token Representation

```
Symbol Specification and Lexical Analysis
```

```
Example SetLan Program
```

```
{
  set one;
  set two;
  one = [100];
  two = one + [77];
}
```

Typical token classes are

- Identifiers
- Keywords
- Literals
- Special Symbols

P. Pfahler (upb)

Ident :

Attribute values are relevant for identifiers and literals.

le SetLan Program	F	0	(1 1)	r	
one; two; = [100]; = one + [77];	5 35 36 14 35 36 14 36 32 21	0 2 0 0 3 0 2 0 0	(1, 1) (2, 3) (2, 7) (2, 10) (3, 3) (3, 7) (3, 10) (4, 3) (4, 7) (4, 9)	t set ; set two ; one = Γ	
token classes are ntifiers ywords erals ecial Symbols re values are relevant for	26 22 14 36 32 36 4 21 26 22 14 6 1	100 0 3 0 2 0 0 77 0 0 0 0 0	(4, 10) (4, 13) (4, 14) (5, 3) (5, 7) (5, 13) (5, 15) (5, 16) (5, 18) (5, 19) (6, 1) (7, 1)	100 ] ; two = one + [ 77 ] ; } EOF	
e values are relevant for	I	0	(7, 1)	EUF	
rs and literals.					
Pfahler (uph) P	1 2 (			Winter 2016/	017 3 / 1
Symbol Specification and Lexical Analysis	Luc				
	Tok	ken Sp	ecification a	and Recog	gnition
	Ident	t:	Let	ter	
ent : Letter (Letter   Digit)*				git	
		Letter	~		

## **Typical Notation: Regular Expression**

Ident	:	Letter X
Х	:	Letter X
Х	:	Digit X
Х	:	

# **Formal Notation: Regular Grammar**





### **Implementation: Finite Automaton**



# Transform the NFA into a DFA (1)

The "subset construction" uses two functions:

- The ε-closure function takes a state and returns the set of states reachable from it based on (one or more) ε-transitions.
- The function *move* takes a state and a character, and returns the set of states reachable by one transition on this character.

These functions are generalized to apply to sets of states by taking the union of the application to individual states.

#### Example:

```
If A, B and C are states,
then move(\{A, B, C\}, `a') = move(A, `a') \cup move(B, `a') \cup move(C, `a')
```



# Recognize the Longest Match

An automaton may contain transitions from final states: When does the automaton stop?

#### Rule of the longest match

- The automaton continues as long as there is a transition with the next character.
- After having stopped it sets back to the most recently passed final state. The input position is adjusted.
- If no final state has been passed an error message is issued.

Consequence: Some kinds of tokens have to be separated explicitly.

P. Pfahler (upb)	PLaC			Winter 2016/2017	9 / 1
	Symbol Specification and Lexical Analysis				
		Aspects of	Scanner	Implementat	ion

Scanner runtime is proportional to the number of characters in the program  $\Rightarrow$  Operations per character must be fast, otherwise the scanner dominates compilation time.

Two implementation techniques:

• Table driven scanning:

Scanner control loop interprets state transition table.

 Directly programmed scanning: State transitions are coded by control flow (switches, branches, loops).

Directly programmed scanner implementations are usually faster than table-driven scanners.

# Aspects of Lexical Level Language Design



Scanner generators generate the central function of lexical analysis from token specifications given as regular expressions.

- Lex / Flex
  standard Unix tool and its successor
  scanner implemented in C/C++
  table driven
  - GLA (Generator for Lexical Analysis)
    - University of Boulder, Colorado
    - Part of the Eli system, interfaces with other components
    - library of RE for typical tokens ("canned descriptions")
    - scanner implemented in C/C++
    - directly programmed implementation
  - JLex / JFlex
    - Java Scanner Generator and its successor, GNU license
    - scanner implemented in Java
    - table driven

P. Pfahler (upb)	PLaC	Winter 2016/2017	7 13 / 1
	Symbol Specification and Lexical Analysis		

# Example: SetLan Token Specification for Eli

Written in GLA notation using canned descriptions.

Literal tokens (keywords and special symbols) are automatically taken from concrete syntax and do not have to be specified in Eli.

SetLan Toker	n Specification
Identifier:	C_IDENTIFIER
Number:	C_INTEGER
	C_COMMENT

Equivalent specification after expanding canned descriptions:

SetLan Token Specification without Canned Descriptions

The specification uses token processors like mkidn to access identifier and literal modules and so-called *auxilliary scanners* like auxCComment for special scanning purposes, like e.g. skipping (possibly nested) comments.