

# Compiler Construction

## Lecture 5: Lexical Analysis IV (Practical Aspects)

Thomas Noll

Lehrstuhl für Informatik 2  
(Software Modeling and Verification)

RWTH Aachen University

`noll@cs.rwth-aachen.de`

`http://www-i2.informatik.rwth-aachen.de/i2/cc10/`

Winter semester 2010/11

# Lehrpreis der Fachgruppe Informatik

**Wir suchen:**

## Beste selbstständige Lehre

Professor/-innen, Juniorprofessor/-innen und (habilitierte) Mitarbeiter/-innen, die selbstständig Vorlesungen, Seminare, Praktika oder ähnliches anbieten.

## Beste Unterstützung in der Lehre

Wissenschaftliche Mitarbeiter/-innen und studentische Hilfskräfte, die maßgeblich an der Betreuung von Vorlesungen, Übungen, Seminaren und Praktika oder ähnlichem mitwirken.

Vorschläge (mit Begründung) bitte bis zum 10. November an:

**[lehrpreis@informatik.rwth-aachen.de](mailto:lehrpreis@informatik.rwth-aachen.de)**

Die Preisverleihung erfolgt im Rahmen des diesjährigen Tages der Informatik.

Weitere Informationen:

**[www.informatik.rwth-aachen.de](http://www.informatik.rwth-aachen.de)**

**RWTH**AACHEN  
UNIVERSITY

**FACH**INFORMATIK  
GRUPPE

Due to the **Fachschaftsvollversammlung** (whatever this means in English...), the Tuesday lecture on November 2 will start at **14:15**.

- 1 Repetition: First-Longest-Match Analysis
- 2 First-Longest-Match Analysis with NFA
- 3 Longest Match in Practice
- 4 Regular Definitions
- 5 Generating Scanners Using `[f]lex`
- 6 Further Problems in Lexical Analysis

# The Extended Matching Problem

## Problem (Extended matching problem)

*Given  $\alpha_1, \dots, \alpha_n \in RE_\Omega$  and  $w \in \Omega^*$ , decide whether there exists a decomposition of  $w$  w.r.t.  $\alpha_1, \dots, \alpha_n$  and determine a corresponding analysis.*

To ensure **uniqueness**:

- ❶ **Principle of the longest match** (“maximal munch tokenization”)
  - for uniqueness of decomposition
  - make lexemes as long as possible
  - motivated by applications: usually every (nonempty) prefix of an identifier is also an identifier
- ❷ **Principle of the first match**
  - for uniqueness of analysis
  - choose first matching regular expression (in the order given)

# Implementation of FLM Analysis

## Algorithm (FLM analysis)

**Input:** expressions  $\alpha_1, \dots, \alpha_n \in RE_\Omega$ , tokens  $\{T_1, \dots, T_n\}$ ,  
input word  $w \in \Omega^*$

**Procedure:**

- ➊ for every  $i \in [n]$ , construct  $\mathfrak{A}_i \in DFA_\Omega$  such that  $L(\mathfrak{A}_i) = \llbracket \alpha_i \rrbracket$  (see *DFA method*)
- ➋ construct the *product automaton*  $\mathfrak{A} \in DFA_\Omega$  such that  $L(\mathfrak{A}) = \bigcup_{i=1}^n \llbracket \alpha_i \rrbracket$
- ➌ *partition the set of final states* of  $\mathfrak{A}$  to follow the first-match principle
- ➍ extend the resulting DFA to a *backtracking DFA* which implements the longest-match principle, and let it run on  $w$

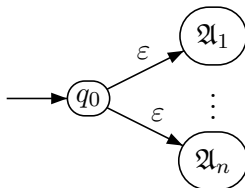
**Output:** FLM analysis of  $w$  (if it exists)

- 1 Repetition: First-Longest-Match Analysis
- 2 First-Longest-Match Analysis with NFA
- 3 Longest Match in Practice
- 4 Regular Definitions
- 5 Generating Scanners Using `[f]lex`
- 6 Further Problems in Lexical Analysis

# A Backtracking NFA

A similar construction is possible for the **NFA method**:

- 1  $\mathfrak{A}_i = \langle Q_i, \Omega, \delta_i, q_0^{(i)}, F_i \rangle \in NFA_\Omega$  ( $i \in [n]$ ) by NFA method
- 2 “Product” automaton:  $Q := \{q_0\} \uplus \biguplus_{i=1}^n Q_i$



- 3 Partitioning of final states:
  - $M \subseteq Q$  is called a  **$T_i$ -matching** if
$$M \cap F_i \neq \emptyset \text{ and for all } j \in [i-1] : M \cap F_j = \emptyset$$
  - yields set of  $T_i$ -matchings  $F^{(i)} \subseteq 2^Q$
  - $M \subseteq Q$  is called **productive** if there exists a productive  $q \in M$
  - yields productive state sets  $P \subseteq 2^Q$
- 4 Backtracking automaton: similar to DFA case



- 1 Repetition: First-Longest-Match Analysis
- 2 First-Longest-Match Analysis with NFA
- 3 Longest Match in Practice**
- 4 Regular Definitions
- 5 Generating Scanners Using `[f]lex`
- 6 Further Problems in Lexical Analysis

# Longest Match in Practice

- In general: **lookahead of arbitrary length** required
  - that is,  $|v|$  unbounded in configurations  $(T, vqw, W)$
  - see Example 4.15:  $\alpha_1 = a$ ,  $\alpha_2 = a^*b$ ,  $w = a \dots a$
- “Modern” programming languages (Pascal, ...):  
**lookahead of one or two characters** sufficient
  - separation of keywords, identifiers, etc. by spaces
  - Pascal: two-character lookahead required to distinguish **1.5** (real number) from **1..5** (integer range)

**However:** principle of longest match not always applicable!

## Example 5.1 (Longest Match in FORTRAN)

### ① Relational expressions

- valid lexemes: `.EQ.` (relational operator), `EQ` (identifier), `12` (integer), `12.`, `.12` (reals)
  - input string: `12.EQ.12`  $\rightsquigarrow$  `12.EQ.12` (ignoring blanks!)
  - intended analysis: `(int, 12)(rel, eq)(int, 12)`
  - LM yields: `(real, 12.0)(id, EQ)(real, 0.12)`
- $\Rightarrow$  wrong interpretation

### ② DO loops

- (correct) input string: `DO 5 I = 1, 20`  $\rightsquigarrow$  `D05I=1, 20`
  - intended analysis:  
`(do, )(label, 5)(id, I)(gets, )(int, 1)(comma, )(int, 20)`
  - LM analysis (wrong): `(id, )(gets, )(int, 1)(comma, )(int, 20)`
- (erroneous) input string: `DO 5 I = 1. 20`  $\rightsquigarrow$  `D05I=1. 20`
  - LM analysis (correct): `(id, D05I)(gets, )(real, 1.2)`

## Example 5.2 (Longest Match in C)

- valid lexemes:
    - `x` (identifier)
    - `--` (decrement operator; ANSI-C: `--`)
    - `1`, `-1` (integers)
  - input string: `x=-1`
  - intended analysis: `(id, x)(gets, )(int, -1)`
  - LM yields: `(id, x)(dec, )(int, 1)`
- ⇒ wrong interpretation

Possible solutions:

- Hand-written (non-FLM) scanners
- Lookahead (later)

- 1 Repetition: First-Longest-Match Analysis
- 2 First-Longest-Match Analysis with NFA
- 3 Longest Match in Practice
- 4 Regular Definitions
- 5 Generating Scanners Using `[f]lex`
- 6 Further Problems in Lexical Analysis

# Regular Definitions I

**Goal:** modularizing the representation of regular sets by introducing additional identifiers

## Definition 5.3 (Regular definition)

Let  $\{R_1, \dots, R_n\}$  be a set of symbols disjoint from  $\Omega$ . A **regular definition** (over  $\Omega$ ) is a sequence of equations

$$\begin{array}{c} R_1 = \alpha_1 \\ \vdots \\ R_n = \alpha_n \end{array}$$

such that, for every  $i \in [n]$ ,  $\alpha_i \in RE_{\Omega \uplus \{R_1, \dots, R_{i-1}\}}$ .

**Remark:** since no recursion is involved, every  $R_i$  can (iteratively) be substituted by a regular expression  $\alpha \in RE_{\Omega}$   
(otherwise  $\implies$  context-free languages)

## Example 5.4 (Symbol classes in Pascal)

Identifiers:  $Letter = A \mid \dots \mid Z \mid a \mid \dots \mid z$   
 $Digit = 0 \mid \dots \mid 9$   
 $Id = Letter (Letter \mid Digit)^*$

Numerals:  
(unsigned)  $Digits = Digit^+$   
 $Empty = \Lambda^*$   
 $OptFrac = . Digits \mid Empty$   
 $OptExp = E (+ \mid - \mid Empty) Digits \mid Empty$   
 $Num = Digits OptFrac OptExp$

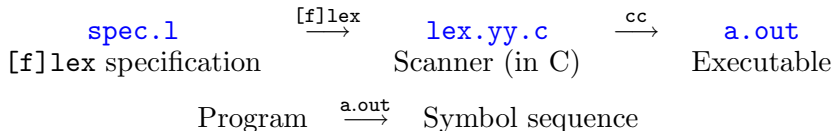
Rel. operators:  $RelOp = < \mid <= \mid = \mid > \mid >=$

Keywords:  $If = if$   
 $Then = then$   
 $Else = else$

- 1 Repetition: First-Longest-Match Analysis
- 2 First-Longest-Match Analysis with NFA
- 3 Longest Match in Practice
- 4 Regular Definitions
- 5 Generating Scanners Using `[f]lex`
- 6 Further Problems in Lexical Analysis



Usage of [f]lex (“[fast] lexical analyzer generator”):



A [f]lex **specification** is of the form

*Definitions (optional)*

%%

*Rules*

%%

*Auxiliary procedures (optional)*

- Definitions:
- C code for declarations etc.: `%{ Code %}`
  - **Regular definitions:** *Name RegExp ...*  
(non-recursive!)

Rules: of the form *Pattern { Action }*

- *Pattern*: regular expression, possibly using *Names*
- *Action*: C code for computing  
**symbol** = (token, attribute)
  - **token**: integer **return** value, 0 = EOF
  - **attribute**: passed in global variable `yylval`
  - **lexeme**: accessible by `yytext`
- matching rule found by **FLM strategy**
- **lexical errors** caught by `.` (any character)

# Example [f]lex Specification

```
%{
#include <stdio.h>
typedef enum {EOF, IF, ID, RELOP, LT, ...} token_t;
unsigned int yylval;    /* attribute values */
}%

LETTER      [A-Za-z]
DIGIT       [0-9]
ALPHANUM    {LETTER}|{DIGIT}
SPACE       [ \t\n]

%%

"if"        { return IF; }
"<"         { yylval = LT; return RELOP; }
{LETTER}{ALPHANUM}*  { yylval = install_id(); return ID; }
{SPACE}+     /* eat up whitespace */
.            { fprintf (stderr, "Invalid character '%c'\n", yytext[0]); }

%%

int main(void) {
    token_t token;
    while ((token = yylex()) != EOF)
        printf ("(Token %d, Attribute %d)\n", token, yylval);
    exit (0);
}

unsigned int install_id () {...}    /* identifier name in yytext */
```

# Regular Expressions in [f]lex

Syntax	Meaning
printable character	this character
<code>\n</code> , <code>\t</code> , <code>\123</code> , etc.	newline, tab, octal representation, etc.
<code>.</code>	any character except <code>\n</code>
<code>[Chars]</code>	one of <i>Chars</i> ; ranges possible (" <code>0-9</code> ")
<code>[^Chars]</code>	none of <i>Chars</i>
<code>\\</code> , <code>\.</code> , <code>\[</code> , etc.	<code>\</code> , <code>.</code> , <code>[</code> , etc.
<code>"Text"</code>	<i>Text</i> without interpretation of <code>.</code> , <code>[</code> , <code>\</code> , etc.
<code>^α</code>	$\alpha$ at beginning of line
<code>α\$</code>	$\alpha$ at end of line
<code>{Name}</code>	<i>RegExp</i> for <i>Name</i>
<code>α?</code>	zero or one $\alpha$
<code>α*</code>	zero or more $\alpha$
<code>α+</code>	one or more $\alpha$
<code>α{n,m}</code>	between <i>n</i> and <i>m</i> times $\alpha$ (" <i>m</i> " optional)
<code>(α)</code>	$\alpha$
<code>α<sub>1</sub>α<sub>2</sub></code>	concatenation
<code>α<sub>1</sub> α<sub>2</sub></code>	alternative
<code>α<sub>1</sub>/α<sub>2</sub></code>	$\alpha_1$ but only if followed by $\alpha_2$ (lookahead)

## Example 5.5 (Lookahead in FORTRAN)

### ❶ DO loops (cf. Example 5.1)

- input string: `DO 5 I = 1, 20`
- LM yields: `(id, )(gets, )(int, 1)(comma, )(int, 20)`
- observation: decision for `do` only possible after reading “,”
- specification of `DO` keyword in `[f]lex`, using lookahead:  
`DO / ({LETTER}|{DIGIT})* = ({LETTER}|{DIGIT})* ,`

### ❷ IF statement

- problem: FORTRAN keywords not reserved
- example: `IF(I,J) = 3` (assignment to an element of matrix `IF`)
- conditional: `IF (condition) THEN ...` (`IF` keyword)
- specification of `IF` keyword in `[f]lex`, using lookahead:  
`IF / \( .* \) THEN`

# Longest Match and Lookahead in [f]lex

```
%{
#include <stdio.h>
typedef enum {EoF, AB, A} token_t;
}%
%%
ab      { return AB; }
a/bc    { return A; }
.       { fprintf (stderr, "Invalid character '%c'\n", yytext[0]); }
%%
int main(void) {
    token_t token;
    while ((token = yylex()) != EoF) printf ("Token %d\n", token);
    exit (0);
}
```

returns on input

- a: Invalid character 'a'
- ab: Token 1
- abc: Token 2   Invalid character 'b'   Invalid character 'c'

⇒ lookahead counts for length of match

- 1 Repetition: First-Longest-Match Analysis
- 2 First-Longest-Match Analysis with NFA
- 3 Longest Match in Practice
- 4 Regular Definitions
- 5 Generating Scanners Using `[f]lex`
- 6 Further Problems in Lexical Analysis

## Identifiers that influence subsequent parsing

- **Example:** (type definitions in C)

The program fragment `(T *)` is

- valid in the scope of declaration `typedef int T;`  
(casting to type “pointer to `T`”)
- invalid in the scope of declaration `int T;`  
(incorrect expression with missing second multiplication factor)

- **Solution:** exploit symbol table information

## Macro processing

- macro substitution
- parameter substitution
- file inclusion
- conditional compilation
- **Solution:** separate preprocessing phase between reading and lexical analysis