

Lexical tokens

Type	Examples
ID	foo n14 last
NUM	73 0 00 515 082
REAL	66.1 .5 10. 1e67 5.5e-10
IF	if
COMMA	,
LPAREN	(
ASSIGN	:=

- Some tokens eg ID NUM REAL have *semantic values* attached to them, eg ID(n14), NUM(515)
- Reserved words: Punctuation tokens e.g. IF, VOID, RETURN constructed from alphanumeric characters, cannot be used as identifiers

1 February, 2010

IN2009 Language Processors - Session 2

7

Example informal specification

Identifiers in C or Java:

- “An identifier is a sequence of letters and digits; the first character must be a letter. The underscore (`_`) counts as a letter. Upper- and lowercase letters are different.”
- “If the input stream has been divided into tokens up to a given character, the next token is taken to include the longest string of characters that could possibly constitute a token.”

1 February, 2010

IN2009 Language Processors - Session 2

8

Identifiers in C or Java

- *Blanks, tabs, newlines and comments are ignored except as they serve to separate tokens. Some white space is required to separate otherwise adjacent identifiers, keywords and symbols.*

How do we formally specify the description above?

1 February, 2010

IN2009 Language Processors - Session 2

9

Formal specifications of tokens

- Approach:
 - Specify lexical tokens using the formal language of regular expressions.
 - Transform regular expressions into deterministic finite automata (DFA)
 - Implement lexical analysers (lexers) using DFA.
 - Fortunately for our sanity, there are automatic conversion tools...

1 February, 2010

IN2009 Language Processors - Session 2

10

Formal specifications of tokens

“a formal language is a language that is defined by precise mathematical or machine processable formulas” (Wikipedia)

- Languages
 - A language is a set of strings
 - A string is a finite sequence of symbols
 - Symbols are taken from a finite alphabet
- Many types of formal languages, at this stage we will consider **regular languages**.

1 February, 2010

IN2009 Language Processors - Session 2

11

Regular languages

- Regular languages can be:
 - Specified by regular expressions
 - Accepted by deterministic/non-deterministic finite state automata (DFA,NFA)
 - A word is accepted by the state machine if and only if the word belongs to the language i.e. it satisfies the regular expression specification
- A regular expression can be translated into an NFA.
- DFAs can be efficiently implemented.
- An NFA can be converted into a DFA.

1 February, 2010

IN2009 Language Processors - Session 2

12

Regular expressions

- Regular expressions (regex, regexp) - strings that describe the syntax of other strings.
- These are simply *patterns*, constructed by syntactical rules and defining a set of strings.
- Components
 - Alphabet: Valid symbols of the language
 - Operators: \cdot $|$ $*$ $+$

1 February, 2010

IN2009 Language Processors - Session 2

13

Regular expressions

- There is a useful connection to real life software development here, too.
 - Regexp can be used in a number of programming languages and tools
 - TextPad, Emacs, etc. allow for searching for text given a regular expression - much more powerful than plain text searching.
 - Java contains a class RegExp, too.
- On to the rules...

1 February, 2010

IN2009 Language Processors - Session 2

14

Regular expressions

Symbol

- For each symbol a in the alphabet of the language, the regexp a denotes the language containing just the string a
- $\text{Lang}(a) = \{ "a" \}$

1 February, 2010

IN2009 Language Processors - Session 2

15

Regular expressions

Alternation

- Given 2 regular expressions M and N then $M | N$ is a new regexp.
- A string is in $\text{lang}(M|N)$ if it is in $\text{lang}(M)$ or $\text{lang}(N)$.
- $\text{lang}(M|N) = \text{lang}(M) \cup \text{lang}(N)$
- E.g. $\text{lang}(a | b) = \{ "a", "b" \}$ contains the 2 strings a and b .

1 February, 2010

IN2009 Language Processors - Session 2

16

Regular expressions

Concatenation

- Given 2 regexes M and N then $M \cdot N$ is a new regexp.
- A string is in $\text{lang}(M \cdot N)$ if it is the concatenation of 2 strings α and β s.t. α in $\text{lang}(M)$ and β in $\text{lang}(N)$.
- $\text{lang}(M \cdot N) = \{ \alpha\beta \text{ s.t. } \alpha \text{ in } \text{lang}(M) \text{ and } \beta \text{ in } \text{lang}(N) \}$
- Thus regexp $(a|b) \cdot a = \{ "aa", "ba" \}$ defines the language containing the 2 strings aa and ba

1 February, 2010

IN2009 Language Processors - Session 2

17

Regular expressions

Epsilon

- The regexp ϵ represents the language whose only string is the empty string.
- $\text{Lang}(\epsilon) = \{ "" \}$
- Thus $(a \cdot b) | \epsilon$ represents the language $\{ "", "ab" \}$

1 February, 2010

IN2009 Language Processors - Session 2

18

Regular expressions

Repetition

- M^* - Kleene closure (or Kleene Star) of M
- A string is in M^* if it is the concatenation of ≥ 0 strings, all in M .
- Thus $((a|b) \cdot a)^*$ represents the infinite set
 $\{ "", "aa", "ba", "aaaa", "baaa", "aaba", "baba", "aaaaaa", \dots \}$

1 February, 2010

IN2009 Language Processors - Session 2

19

Examples

- $(0|1)^* \cdot 0$
- $b^*(abb^*)^*(a|\epsilon)$
- $(a|b)^*aa(a|b)^*$
- Conventions:
 - omit \cdot and ϵ , assume Kleene closure binds tighter than \cdot and concatenation binds tighter than $|$
 - $a \cdot b^*$ means $a(b^*)$
 - $ab|c$ means $(a \cdot b)|c$
 - $(a|)$ means $(a|\epsilon)$

1 February, 2010

IN2009 Language Processors - Session 2

20

Abbreviations (extensions)

$[abcd]$	means $(a b c d)$
$[b-g]$	means $[bcdefg]$
$\sim[b-g]$	means everything but $[bcdefg]$
$[b-gM-Qkr]$	means $[bcdefgMNO PQkr]$
$M?$	means $(M \epsilon)$
M^+	means $M(M)^*$

1 February, 2010

IN2009 Language Processors - Session 2

21

Regular expression summary

a or $"a"$	ordinary character, stands for itself
ϵ	the empty string
	another way to write the empty string
$M N$	alternation
$M \cdot N$	concatenation (often written simply as MN)
M^*	repetition (zero or more times)
M^+	repetition (one or more times)
$M?$	Optional, zero or one occurrence of M

1 February, 2010

IN2009 Language Processors - Session 2

22

Regular expression summary

$["a"- "z" "A" - "Z"]$	Character set alternation (JavaCC)
$\sim[]$	Any single character ($\sim[]$ is JavaCC form)
$"\n" "\t" "\""$	newline, tab, double quote (quoted special characters)
$"a. +"$	quotation, string stands for itself (in this case $a. +$)
$[a-zA-Z]$	Character set alternation
$.$	Any single character except newline

1 February, 2010

IN2009 Language Processors - Session 2

23

Regular expressions for some tokens

$(" " "\n" "\t")$	<i>whitespace; ignore</i>
if	IF
$[a-z][a-z0-9]^*$	ID
$[0-9]^+$	NUM
$([0-9]^+ \cdot "." [0-9]^*) ([0-9]^* "." [0-9]^+)$	REAL
$("--" [a-z]^* "\n")$	<i>comment starting --; ignore</i>
$.$	Error

Completeness: there must always be a match for some initial substring of the input.
 We always want the **longest match**.
Rule Priority: If more than one regexp matches the string, pick the first one from the lexical spec.

1 February, 2010

IN2009 Language Processors - Session 2

24

Example: Lexical Specification

We'll assume "whitespaces" are skipped.

```
if                IF
[a-z][a-z0-9]*    ID
.                 Error
```

Sample input: "a2", "if", "if2", "02a12".

What happens if we put **ID** first?

What happens if we put **.** first?

Can we use **.***?

1 February, 2010

IN2009 Language Processors - Session 2

25

Finite Automata

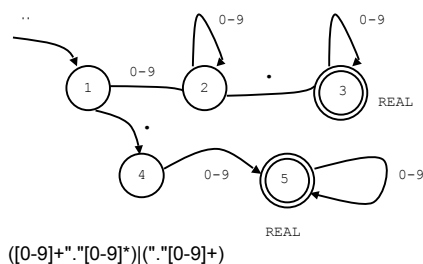
- A finite set of states and edges.
- Edges lead from one state to another. Each edge is labeled by a symbol.
- A single start state and one or more final states.
- Each final state identifies a TOKEN from the language
- A string is accepted if, after all characters have been matched by a transition, the last state is a final state.
- Deterministic or Non-deterministic (NFA) .

1 February, 2010

IN2009 Language Processors - Session 2

26

Finite automaton



1 February, 2010

IN2009 Language Processors - Session 2

27

DFA Implementation

```
method Token lex() { //example automaton for REAL only
    int state = 1; String text = ""; char ch;
    while (true) {
        ch = nextchar(); // Get the next input character
        text = text + ch; // Collect text of the token
        if (state == 1) {
            if (ch >= '0' && ch <= '9') {
                state = 2;
            } else if (ch == '.') {
                state = 4;
            } else {
                lexerror(ch);
            }
        }
        ...
    }
}
```

1 February, 2010

IN2009 Language Processors - Session 2

28

DFA Implementation

```
...
else if (state == 5) {
    if ch >= '0' && ch <= '9' {
        state = 5;
    } else {
        return new Token(REAL, new Double(text));
    }
} else {
    error ("Illegal state: shouldn't happen");
}
}
```

- Inefficient and verbose..use transition matrix instead

1 February, 2010

IN2009 Language Processors - Session 2

29

Efficient DFA implementation

- Create a table that represents the transition matrix
 - `int edges[NumStates][NumCharacters]`
 - `Edges[s][c] = sn`
If there exists a transition labeled with character `c` that joins states `s` and `sn`
 - Index `c` can be extracted from the ASCII code
 - Create a "dead" state that loops to itself on all characters (to denote no match)
- Create an array that marks states as final and maps them to tokens.
- This table can be generated automatically

1 February, 2010

IN2009 Language Processors - Session 2

30

From RegExp to Implementation

- A RegExp can be converted into an NFA.
- However, the implementation of NFAs is inefficient.
- But NFAs can be translated into DFAs!

All these steps can be automated
There are tools that generate lexical analysers from lexical specifications
E.g. JavaCC

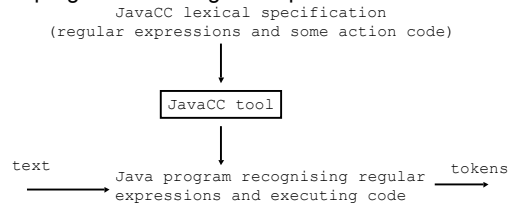
1 February, 2010

IN2009 Language Processors - Session 2

31

JavaCC compiler-compiler

- Fortunately, tools can produce finite automata programs from regular expressions...



1 February, 2010

IN2009 Language Processors - Session 2

32

JavaCC token regexps

- Characters and strings must be quoted, eg:
 - “,” “int” “while” “\n” “\hello\”
- Character lists [...] provide a shorthand for |, eg:
 - “[a”-”z”] matches “a” through “z”, “[a”, “e”, “i”, “o”, “u”] matches any single vowel, ~ “[a”, “e”, “i”, “o”, “u”] any non-vowel, ~[] any character
- Repetition with + and *, eg:
 - “[a”-”z”, “A”-”Z”]+ matches one or more letters
 - “[a”-”z”][“0”-”9”]* matches a letter followed by zero or more digits

1 February, 2010

IN2009 Language Processors - Session 2

33

JavaCC token regexps

- Shorthand with ? provides for optional expr, eg:
 - “(“+”|“-”)?(“0”-”9”)+ matches signed and unsigned integers
- Tokens can be named
 - TOKEN : {
 < IDENTIFIER: < LETTER> (< LETTER> | < DIGIT>)* >
 }
– TOKEN : {
 < LETTER: [“a”-”z”, “A”-”Z”] >
 | < DIGIT: [“0”-”9”] >
 }
– now < IDENTIFIER> can be used later in defining syntax

1 February, 2010

IN2009 Language Processors - Session 2

34

JavaCC lexical analysis example

```
/* file MyParser.jj */
PARSER_BEGIN(MyParser)
class MyParser {
  PARSER_END(MyParser)

  TOKEN : {
    < IF: "if" >
    | < #DIGIT: [“0”-”9”] >
    | < ID: [“a”-”z”] ([“a”-”z”] | < DIGIT>)* >
    | < NUM: (< DIGIT>)+ >
    | < REAL: ( (< DIGIT>)+ “.” (< DIGIT>)* ) |
      ( (< DIGIT>)* “.” (< DIGIT>)+ ) >
  }

  SKIP : { // skipped during lexical analysis
    < “ ” > ([“a”-”z”])* (“\n” | “\r” | “\r\n”) >
    | “ ” | “\t” | “\n” | “\r”
  }
}
```

means can use **only** in
TOKEN definitions (local)

1 February, 2010

IN2009 Language Processors - Session 2

35

JavaCC lexical analysis example

```
void Start() : // the grammar
{Token t;}
{ ( ( t=<IF> | t=<ID> | t=<NUM> | t=<REAL> )
  { System.out.println("token found: "
    + tokenImage[t.kind]
    + " (" + t.image + ")"); }
  )* <EOF>
} /* end of file MyParser.jj */

/* file MyParser.java */
class MyParser {
  public static void main(String args[]) throws
    ParseException {
    ...
    MyParser parser = new MyParser(System.in);
    parser.Start();
  }
}
```

1 February, 2010

IN2009 Language Processors - Session 2

36

JavaCC introduction

- Generates a combined lexical analyser and parser (a Java class). Here the class is called `MyParser`
- This session we're learning about lexical analysis
 - JavaCC does this and it is the focus of our first example.
 - JavaCC-generated parsers (later).
- You can see another example in today's lab.

1 February, 2010

IN2009 Language Processors - Session 2

37

JavaCC introduction

- all you need to know now is that it also generates a Java method to recognize the things we've labelled under 'Start()' – this is the default.
- to make it work,
 - create an object
`MyParser parser = new MyParser(inputstream)`
 - and then call the method...`parser.Start()`
- a class `Token` is also generated
 - field **`image`** is the string matched for the token
 - field **`kind`** can be used to index array **`tokenImage`** to see the token type

1 February, 2010

IN2009 Language Processors - Session 2

38