

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 1 of 40

CS421 COMPILERS AND INTERPRETERS

## **Example: Source Code**

#### A Sample Toy Program:

CS421 COMPILERS AND INTERPRETERS

## **The Lexical Structure**

Output after the Lexical Analysis ----- token + associated value

<b>LET</b> 51	FUNCTION 56	<pre>ID(do_nothing1) 65</pre>
LPAREN 76	<b>ID</b> (a) 77	COLON 78
<b>ID</b> (int) 80	COMMA 83	<b>ID</b> (b) 85
COLON 86	ID(string) 88	RPAREN 94
<b>EQ</b> 95	<pre>ID(do_nothing2)</pre>	99
LPAREN 110	<b>ID</b> (a) 111	<b>PLUS</b> 112
<b>INT</b> (1) 113	RPAREN 114	FUNCTION 117
<pre>ID(do_nothing2)</pre>	126	LPAREN 137
<b>ID</b> (d) 138	<b>COLON</b> 139	<b>ID</b> (int) 141
RPAREN 144	<b>EQ</b> 146	
<pre>ID(do_nothing1)</pre>	150	LPAREN 161
<b>ID</b> (d) 162	<b>COMMA</b> 163	STRING(str) 165
RPAREN 170	<b>IN</b> 173	
<pre>ID(do_nothing1)</pre>	177	LPAREN 188
<b>INT</b> (0) 189	<b>COMMA</b> 190	STRING(str2) 192
RPAREN 198	<b>END</b> 200	EOF 203

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 3 of 40

CS421 COMPILERS AND INTERPRETERS

### Tokens

• **Tokens** are the <u>atomic unit</u> of a language, and are usually <u>specific</u> <u>strings</u> or <u>instances</u> of <u>classes</u> of strings.

Tokens	Sample Values	Informal Description
LET	let	keyword LET
END	end	keyword END
PLUS	+	
LPAREN	(	
COLON	:	
STRING	"str"	
RPAREN	)	
INT	49, 48	integer constants
ID	do_nothingl, a, int, string	letter followed by letters, digits, and under-scores
EQ	=	
EOF		end of file

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 4 of 40



Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 5 of 40

CS421 COMPILERS AND INTERPRETERS

# **Regular Expressions**

 regular expressions are concise, linguistic characterization of regular languages (regular sets)

identifier = letter (letter | digit | underscore)\*

- each regular expression define a regular language ---- a <u>set of strings</u> over some alphabet, such as ASCII characters; each member of this set is called a **sentence**, or a **word**
- we use regular expressions to define each category of tokens

For example, the above identifier specifies a set of strings that are a sequence of letters, digits, and underscores, starting with a letter.

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 6 of 40

" 0 or more"

CS421 COMPILERS AND INTERPRETERS

## Regular Expressions and Regular Languages

- Given an alphabet ? ??the regular expressions over ? and their corresponding regular languages are
  - a) ? 2denotes??? ?? , the empty string, denotes the language { ? }.
  - b) for each a in ?, a denotes { a } --- a language with one string.
  - c) if *R* denotes  $L_R$  and *S* denotes  $L_S$  then *R* / *S* denotes the language  $L_R$ ?  $\mathcal{L}_S$ , i.e, { x | x ?  $\mathcal{L}_R$  or x ?  $\mathcal{L}_S$  }.
  - d) if *R* denotes  $L_R$  and *S* denotes  $L_S$  then *RS* denotes the language  $L_RL_S$ , that is, { xy | x ?  $L_R$  and y ?  $L_S$  }.
  - e) if R denotes L<sub>R</sub> then R<sup>\*</sup> denotes the language L<sub>R</sub><sup>\*</sup> where L<sup>\*</sup> is the union of all L<sup>i</sup> (i=0,...,? ?and L<sup>i</sup> is just {x<sub>1</sub>x<sub>2</sub>...x<sub>i</sub> | x<sub>1</sub>? 1, ..., x<sub>i</sub>? 1}.
  - f) if R denotes  $L_R$  then (R) denotes the same language  $L_R$ .

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 7 of 40

#### CS421 COMPILERS AND INTERPRETERS

### Example

Regular Expression	Explanation
a*	0 or more a's
a <sup>+</sup>	1 or more a's
(a b)*	all strings of a's and b's (including ?)
(aa ab ba bb) <sup>*</sup>	all strings of a's and b's of even length
[a-zA-Z]	shorthand for "a b  z A  Z"
[0-9]	<i>shorthand for</i> "0 1 2  9"
0([0-9])*0	numbers that start and end with 0
(ab aab b) <sup>*</sup> (a aa ?)	?
?	all strings that contain foo as substring

the following is **not** a regular expression:

 $a^{n}b^{n}$  (n > 0)

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

#### CS421 COMPILERS AND INTERPRETERS

## **Lexical Specification**

• Using regular expressions to specify tokens

```
keyword = begin | end | if | then | else
identifier = letter (letter | digit | underscore)*
integer = digit*
relop = < | <= | = | <> | > | >=
letter = a | b | ... | z | A | B | ... | Z
digit = 0 | 1 | 2 | ... | 9
```

- Ambiguity : is "begin" a keyword or an identifier ?
- Next step: to construct a token recognizer for languages given by regular expressions --- by using finite automata !

given a string *x*, the token recognizer says "yes" if *x* is a sentence of the specified language and says "no" otherwise

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 9 of 40

# **Transition Diagrams**

- Flowchart with states and edges; each edge is labelled with characters; certain subset of states are marked as "final states"
- Transition from state to state proceeds along edges according to the next input character



- Every string that ends up at a **final state** is accepted
- If get "stuck", there is no transition for a given character, it is an error
- Transition diagrams can be easily translated to programs using **case** statements (in C).

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 10 of 40

CS421 COMPILERS AND INTERPRETERS

# Transition Diagrams (cont'd)

The token recognizer (for identifiers) based on transition diagrams:

#### CS421 COMPILERS AND INTERPRETERS

# Finite Automata

- Finite Automata are similar to transition diagrams; they have states and labelled edges; there are one unique start state and one or more than one final states
- Nondeterministic Finite Automata (NFA) :
  - a) ? can label edges (these edges are called ?-transitions)
  - b) some character can label 2 or more edges out of the same state

#### • Deterministic Finite Automata (DFA) :

- a) no edges are labelled with ?
- b) each charcter can label at most one edge out of the same state
- NFA and DFA accepts string x if there exists a path from the start state to a final state labeled with characters in x

NFA: multiple paths

DFA: one unique path

```
Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University
```

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University









CS421 COMPILERS AND INTERPRETERS **Transition Table** • Finite Automata can also be represented using transition tables For NFA, each entry is a set of For DFA, each entry is a unique states: state: STATE STATE а b а b 0 {0,1} {0} 0 1 0 1 -{2} 1 1 2 2 {3} 2 1 3 -3 3 1 0 --

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 15 of 40



CS421 COMPILERS AND INTERPRETERS



<section-header>

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 19 of 40





Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 21 of 40

CS421 COMPILERS AND INTERPRETERS



Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 22 of 40

CS421 COMPILERS AND INTERPRETERS

# NFA -> DFA

- NFA are non-deterministic; need DFA in order to write a deterministic prorgam !
- There exists an algorithm ("subset construction") to convert any NFA to a DFA that accepts the same language
- States in DFA are **sets of states** from NFA; DFA simulates "in parallel" all possible moves of NFA on given input.
- Definition: for each state s in NFA,

?-CLOSURE(s) = { s } ? { t | s can reach t via ?-transitions }

• Definition: for each set of states S in NFA,

?-CLOSURE(**S**) = ? ; ?-CLOSURE(**s**) for all  $\mathbf{s}_i$  in **S** 

CS421 COMPILERS AND INTERPRETERS

# NFA -> DFA (cont'd)

- each DFA-state is a set of NFA-states
- suppose the start state of the NFA is s, then the start state for its DFA is ?-CLOSURE(s); the final states of the DFA are those that include a NFA-final-state
- Algorithm : converting an NFA N into a DFA D ----

```
Dstates = {?-CLOSURE(s<sub>0</sub>),s<sub>0</sub> is N's start state}
Dstates are initially "unmarked"
while there is an unmarked D-state X do {
    mark X
    for each a???? do {
        T = {states reached from any s<sub>i</sub> in X via a}
        Y = ?-CLOSURE(T)
        if Y? Dstates then add Y to Dstates "unmarked"
        add transition from X to Y, labelled with a
    }
}
```

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 24 of 40

# Example : NFA -> DFA

• converting NFA for (a|b)\*abb to a DFA -----

The start state  $A = ?-CLOSURE(0) = \{0, 1, 2, 4, 7\};$  Dstates={A}

1st iteration: A is unmarked; mark A now; a-transitions: T = {3, 8} a new state B= ?-CLOSURE(3) ? ??-CLOSURE(8) = {3, 6, 1, 2, 4, 7} ? ??8) = {1, 2, 3, 4, 6, 7, 8} add a transition from A to B labelled with a

> b-transitions: T =  $\{5\}$ a new state C = ?-CLOSURE(5) =  $\{1, 2, 4, 5, 6, 7\}$ add a transition from A to C labelled with b **Dstates** =  $\{A, B, C\}$

2nd iteration: B, C are unmarked; we pick B and mark B first; B =  $\{1, 2, 3, 4, 6, 7, 8\}$ B's a-transitions: T =  $\{3, 8\}$ ; T's ?-CLOSURE is B itself. add a transition from B to B labelled with a

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 25 of 40

CS421 COMPILERS AND INTERPRETERS

### Example : NFA -> DFA (cont'd)

B's b-transitions: T =  $\{5, 9\}$ ; a new state D = ?-CLOSURE( $\{5, 9\}$ ) =  $\{1, 2, 4, 5, 6, 7, 9\}$ add a transition from B to D labelled with b **Dstates** =  $\{A, B, C, D\}$ 

then we pick C, and mark C C's a-transitions: T = {3, 8}; its ?-CLOSURE is B. add a transition from C to B labelled with a C's b-transitions: T = {5}; its ?-CLOSURE is C itself. add a transition from C to C labelled with b

next we pick D, and mark D D's a-transitions: T = {3, 8}; its ?-CLOSURE is B. add a transition from D to B labelled with a D's b-transitions: T = {5, 10}; a new state E = ?-CLOSURE({5, 10}) = {1, 2, 4, 5, 6, 7, 10} **Dstates** = {A, B, C, D, E}; E is a **final state** since it has 10;

next we pick E, and mark E

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 26 of 40

CS421 COMPILERS AND INTERPRETERS

# Example : NFA -> DFA (cont'd)

E's a-transitions:  $T = \{3, 8\}$ ; its ?-CLOSURE is B. add a transition from E to B labelled with a E's b-transitions:  $T = \{5\}$ ; its ?-CLOSURE is C itself. add a transition from E to C labelled with b

all states in Dstates are marked, the DFA is constructed !



Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 27 of 40

CS421 COMPILERS AND INTERPRETERS

# **Other Algorithms**

- How to minimize a DFA ? (see Dragon Book 3.9, pp141)
- How to convert RE to DFA directly ? (see Dragon Book 3.9, pp135)
- How to prove two Regular Expressions are equivalent ? (see Dragon Book pp150, Exercise 3.22)

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 28 of 40









CS421 COMPILERS AND INTERPRETERS **ML-Lex** • ML-Lex is like Lex ------ it takes lexical specification as input, and produces a lexical processor written in Standard ML. Lex Specification ML-Lex foo.lex.sml foo.lex ML Compiler foo.lex.sml module Mlex M lex input text sequence of tokens • Implementation of ML-Lex is similar to implementation of Lex Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University Lexical Analysis : Page 31 of 40



Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

C S 4 2 1	COMPILERS	AND	INTERPRETERS



<b>ML-Lex Definitions</b>				
• Things you can write inside the "ml-lex definitions" section (2nd part):				
%s COMMENT STRING	define new start states			
%reject %count %structure {identifier}	<b>REJECT()</b> to reject a match count the line number the resulting structure name (the default is Mlex)			
<pre>(hint: you probably don't need use %reject, %count,or %structure for assignment 2.)</pre>				
Definition of named regular expressions :				
identifier = regular	expression			
SPACE=[ \t\n\012] IDCHAR=[_a-zA-Z0-9]				
Convright 1984 - 2001 Cursten Schürmann Zhong Shao Yale University	Lexical Analysis - Proc 34 of 40			

CS421 COMPILERS AND INTERPRETERS

# **ML-Lex Translation Rules**

• Each translation rule (3rd part) are in the form

<start-state-list> regular expression => (action);

• Valid ML-Lex regular expressions: (see ML-Lex-manual pp 4-6)

a character stands for itself except for the reserved chars: ? \* + | ( ) ^ \$ / ; . = < > [ { " to use these chars, use backslash! for example, \\\" represents the string \"

using square brackets to enclose a set of characters  $( \ - \ ^{\circ}$  are reserved)

[abc]	<b>char</b> a, <b>or</b> b, <b>or</b> c
[^abc]	<b>all chars except</b> a, b, c
[a-z]	all chars from a to z
[\n\t\b]	new line, tab, or backspace
[-abc]	char - or a or b or c

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

Lexical Analysis : Page 35 of 40

CS421 COMPILERS AND INTERPRETERS

# ML-Lex Translation Rules (cont'd)

• Valid ML-Lex regular expressions: (cont'd)

escape sequences: (can be used inside or outside square brackets) backspace ∖b \n newline \t tab \ddd any ascii char (ddd is 3 digit decimal) any char except newline (equivalent to  $[^{n}]$ ) match string x exactly even if it contains reserved chars "x" an optional x x? 0 or more x's x\* 1 or more x's x+ xy x or y if at the beginning, match at the beginning of a line only ^x  $\{x\}$ substitute definition x (defined in the lex definition section) same as regular expression x (x) repeating x for n times x{n}  $x\{m-n\}$  repeating x from m to n times Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University Lexical Analysis : Page 36 of 40



Lexical Analysis : Page 39 of 40

Lexical Analysis : Page 38 of 40

CS421 COMPILERS AND INTERPRETERS

## Start States (or Start Conditions)

- start states permit multiple lexical analyzers to run together.
- each <u>translation rule</u> can be prefixed with <start-state>
- the lexer is initially in a predefined start stae called INITIAL
- define new start states (in ml-lex-definitions): %s COMMENT STRING
- to switch to another start states (in *action*): YYBEGIN COMMENT
- example: multi-line comments in C

Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University

```
ୢୡୢୡ
%s COMMENT
ୢଌୢଌ
<INITIAL>"/*"
                 => (YYBEGIN COMMENT; continue());
<COMMENT>"*/"
                 => (YYBEGIN INITIAL; continue());
<COMMENT>. | "\n" => (continue());
<INITIAL> .....
```

CS421 COMPILERS AND INTERPRETERS

## Implementation of Lex

- construct NFA for <u>sum</u> of Lex translation rules (regexp/action);
- convert NFA to DFA. then minimize the DFA
- to recognize the input, simulate DFA to termination; find the last DFA state that includes NFA final state, execute associated action (this pickes *longest* match). If the last DFA state has >1 NFA final states, pick one for rule that appears first
- how to represent DFA, the transition table:



```
Copyright 1994 - 2001, Carsten Schürmann, Zhong Shao, Yale University
```

Lexical Analysis : Page 40 of 40