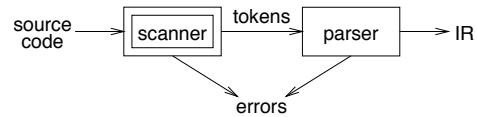


Scanner



- maps characters into *tokens* – the basic unit of syntax
`x = x + y;`
 becomes
`<id, x> = <id, x> + <id, y> ;`
- character string value for a *token* is a *lexeme*
- typical tokens: *number, id, +, -, *, /, do, end*
- eliminates white space (*tabs, blanks, comments*)
- a key issue is speed
 ⇒ use specialized recognizer (as opposed to `lex`)

Copyright ©2007 by Antony L. Hosking. Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and full citation on the first page. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or fee. Request permission to publish from hosking@cs.purdue.edu.

CS502

Scanning

1

Specifying patterns

A scanner must recognize the units of syntax

Some parts are easy:

white space

```

<ws> ::= <ws> ' '
        | <ws> '\t'
        | ' '
        | '\t'
  
```

keywords and operators

specified as literal patterns: `do, end`

comments

opening and closing delimiters: `/* ... */`

CS502

Scanning

2

Specifying patterns

A scanner must recognize the units of syntax

Other parts are much harder:

identifiers

alphanumeric followed by *k* alphanumerics (`_, $, &, ...`)

numbers

integers: 0 or digit from 1-9 followed by digits from 0-9

decimals: integer `'.'` digits from 0-9

reals: (integer or decimal) `'E'` (+ or -) digits from 0-9

complex: `'('` real `'.'` real `)'`

We need a powerful notation to specify these patterns

CS502

Scanning

3

Operations on languages

Operation	Definition
union of <i>L</i> and <i>M</i> written <i>L</i> ∪ <i>M</i>	$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$
concatenation of <i>L</i> and <i>M</i> written <i>LM</i>	$LM = \{st \mid s \in L \text{ and } t \in M\}$
Kleene closure of <i>L</i> written <i>L</i> [*]	$L^* = \bigcup_{i=0}^{\infty} L^i$
positive closure of <i>L</i> written <i>L</i> ⁺	$L^+ = \bigcup_{i=1}^{\infty} L^i$

CS502

Scanning

4

Regular expressions

Patterns are often specified as *regular languages*

Notations used to describe a regular language (or a regular set) include both *regular expressions* and *regular grammars*

Regular expressions (over an alphabet Σ):

- ϵ is a RE denoting the set $\{\epsilon\}$
- if $a \in \Sigma$, then a is a RE denoting $\{a\}$
- if r and s are REs, denoting $L(r)$ and $L(s)$, then:

(r) is a RE denoting $L(r)$

$(r) | (s)$ is a RE denoting $L(r) \cup L(s)$

$(r)(s)$ is a RE denoting $L(r)L(s)$

$(r)^*$ is a RE denoting $L(r)^*$

If we adopt a *precedence* for operators, the extra parentheses can go away. We assume *closure*, then *concatenation*, then *alternation* as the order of precedence.

CS502

Scanning

5

Algebraic properties of REs

Axiom	Description
$r s = s r$	$ $ is commutative
$r (s t) = (r s) t$	$ $ is associative
$(rs)t = r(st)$	concatenation is associative
$r(s t) = rs rt$ $(s t)r = sr tr$	concatenation distributes over $ $
$\epsilon r = r$ $r \epsilon = r$	ϵ is the identity for concatenation
$r^* = (r \epsilon)^*$ $r^{**} = r^*$	relation between $*$ and ϵ $*$ is idempotent

CS502

Scanning

7

Examples

identifier

$letter \rightarrow (a | b | c | \dots | z | A | B | C | \dots | Z)$

$digit \rightarrow (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)$

$id \rightarrow letter (letter | digit)^*$

numbers

$integer \rightarrow (+ | - | \epsilon) (0 | (1 | 2 | 3 | \dots | 9) digit^*)$

$decimal \rightarrow integer . (digit)^*$

$real \rightarrow (integer | decimal) \mathbb{E} (+ | -) digit^*$

$complex \rightarrow ' (' real , real ') '$

Numbers can get much more complicated

Most programming language tokens can be described with REs

We can use REs to build scanners automatically

CS502

Scanning

6

Examples

Let $\Sigma = \{a, b\}$

- $a|b$ denotes $\{a, b\}$
- $(a|b)(a|b)$ denotes $\{aa, ab, ba, bb\}$
i.e., $(a|b)(a|b) = aa|ab|ba|bb$
- a^* denotes $\{\epsilon, a, aa, aaa, \dots\}$
- $(a|b)^*$ denotes the set of all strings of a 's and b 's (including ϵ)
i.e., $(a|b)^* = (a^*b^*)^*$
- $a|a^*b$ denotes $\{a, b, ab, aab, aaab, aaaab, \dots\}$

CS502

Scanning

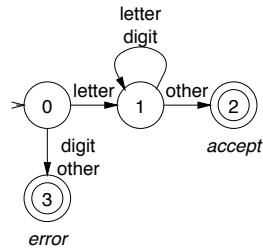
8

Recognizers

From a regular expression we can construct a

deterministic finite automaton (DFA)

Recognizer for *identifier*:



identifier

letter → (a | b | c | ... | z | A | B | C | ... | Z)

digit → (0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9)

id → *letter* (*letter* | *digit*)*

CS502

Scanning

9

Code for the recognizer

```

char ← next_char();
state ← 0; /* code for state 0 */
done ← false;
token_value ← "" /* empty string */
while( not done ) {
  class ← char_class[char];
  state ← next_state[class,state];
  switch(state) {
    case 1: /* building an id */
      token_value ← token_value + char;
      char ← next_char();
      break;
    case 2: /* accept state */
      token_type = identifier;
      done = true;
      break;
    case 3: /* error */
      token_type = error;
      done = true;
      break;
  }
}
return token_type;

```

CS502

Scanning

10

Tables for the recognizer

Two tables control the recognizer

char_class:		a-z	A-Z	0-9	other
	value	letter	letter	digit	other

	class	0	1	2	3
next_state:	letter	1	1	-	-
	digit	3	1	-	-
	other	3	2	-	-

To change languages, we can just change tables

CS502

Scanning

11

Automatic construction

Scanner generators automatically construct code from RE-like descriptions

- construct a DFA
- use state minimization techniques
- emit code for the scanner (table driven or direct code)

A key issue in automation is an interface to the parser

lex is a scanner generator supplied with UNIX

- emits C code for scanner
- provides macro definitions for each token (used in the parser)

CS502

Scanning

12

Grammars for regular languages

Can we place a restriction on the *form* of a grammar to ensure that it describes a regular language?

Provable fact:

For any RE r , \exists a grammar g such that $L(r) = L(g)$

Grammars that generate regular sets are called *regular grammars*:

They have productions in one of 2 forms:

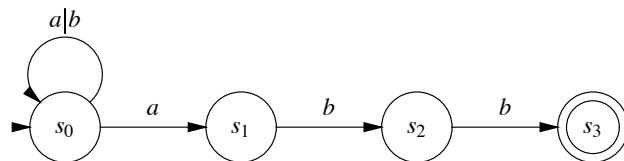
1. $A \rightarrow aA$
2. $A \rightarrow a$

where A is any non-terminal and a is any terminal symbol

These are also called *type 3* grammars (Chomsky)

More regular expressions

What about the RE $(a|b)^*abb$?



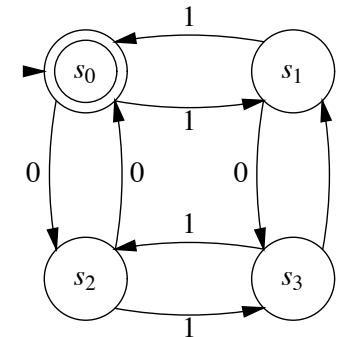
State s_0 has multiple transitions on a !

\Rightarrow *nondeterministic finite automaton*

	a	b
s_0	$\{s_0, s_1\}$	$\{s_0\}$
s_1	-	$\{s_2\}$
s_2	-	$\{s_3\}$

More regular languages

Example: the set of strings containing an even number of zeros and an even number of ones



The RE is $(00 | 11)^*((01 | 10)(00 | 11)^*(01 | 10)(00 | 11)^*)^*$

Finite automata

A *non-deterministic finite automaton* (NFA) consists of:

1. a set of *states* $S = \{s_0, \dots, s_n\}$
2. a set of input symbols Σ (the alphabet)
3. a transition function *move* mapping state-symbol pairs to sets of states
4. a distinguished *start state* s_0
5. a set of distinguished *accepting or final states* F

A *Deterministic Finite Automaton* (DFA) is a special case of an NFA:

1. no state has a ϵ -transition, and
2. for each state s and input symbol a , there is at most one edge labelled a leaving s

A DFA *accepts* x iff. \exists a *unique* path through the transition graph from s_0 to a final state such that the edges spell x .

DFAs and NFAs are equivalent

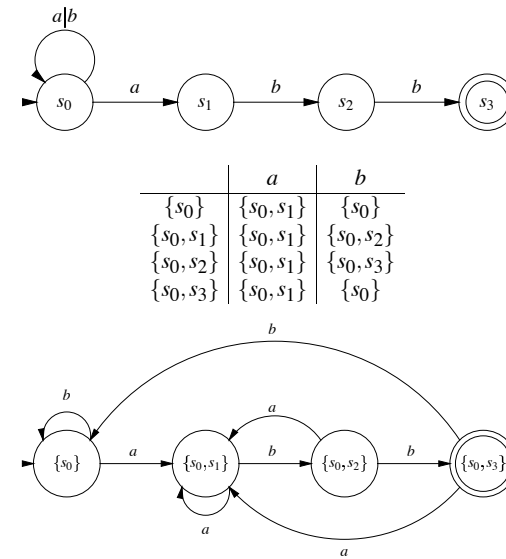
1. DFAs are clearly a subset of NFAs
2. Any NFA can be converted into a DFA, by simulating sets of simultaneous states:
 - each DFA state corresponds to a set of NFA states
 - possible exponential blowup

CS502

Scanning

17

NFA to DFA using the subset construction: example 1

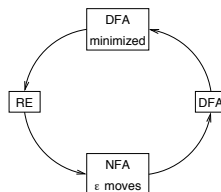


CS502

Scanning

18

Constructing a DFA from a regular expression



RE \rightarrow NFA w/ ϵ moves
 build NFA for each term
 connect them with ϵ moves

NFA w/ ϵ moves to DFA
 construct the simulation
 the "subset" construction

DFA \rightarrow minimized DFA
 merge compatible states

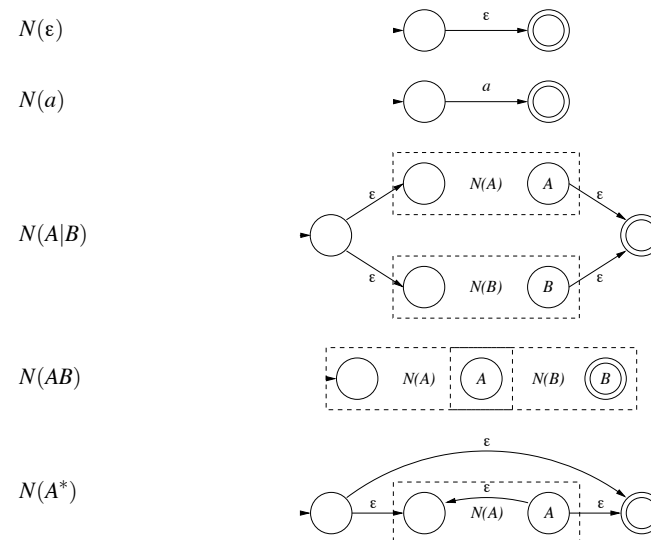
DFA \rightarrow RE
 construct $R_{ij}^k = R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1} \cup R_{ij}^{k-1}$

CS502

Scanning

19

RE to NFA

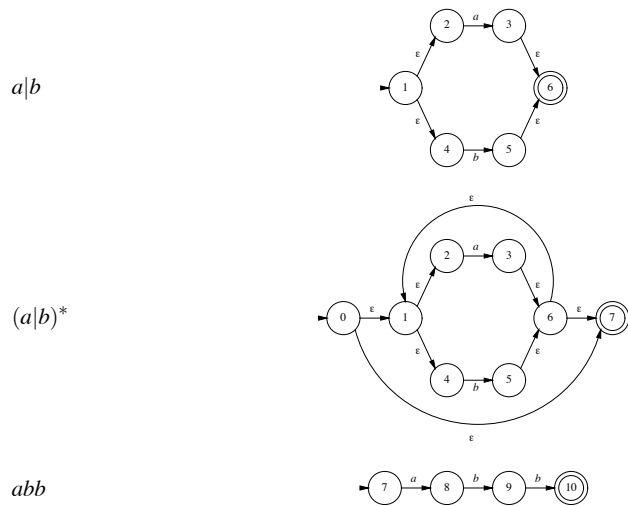


CS502

Scanning

20

RE to NFA: example



CS502

Scanning

21

NFA to DFA: the subset construction

Input: NFA N
 Output: DFA D with states $Dstates$ and transitions $Dtrans$ such that $L(D) = L(N)$
 Method: Let s be a state in N and T be a set of states, define:

Operation	Definition
$\epsilon\text{-closure}(s)$	set of NFA states reachable from NFA state s on ϵ -transitions alone
$\epsilon\text{-closure}(T)$	set of NFA states reachable from some NFA state s in T on ϵ -transitions alone
$move(T, a)$	set of NFA states to which there is a transition on input symbol a from some NFA state s in T

```

add state  $T = \epsilon\text{-closure}(s_0)$  unmarked to  $Dstates$ 
while  $\exists$  unmarked state  $T$  in  $Dstates$ 
  mark  $T$ 
  for each input symbol  $a$ 
     $U = \epsilon\text{-closure}(move(T, a))$ 
    if  $U \notin Dstates$  then add  $U$  to  $Dstates$  unmarked
     $Dtrans[T, a] = U$ 
  endfor
endwhile
    
```

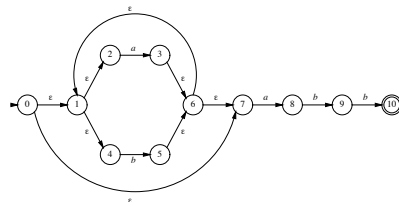
$\epsilon\text{-closure}(s_0)$ is the start state of D
 A state of D is final if it contains at least one final state in N

CS502

Scanning

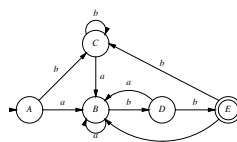
22

NFA to DFA using subset construction: example 2



$A = \{0, 1, 2, 4, 7\}$ $D = \{1, 2, 4, 5, 6, 7, 9\}$
 $B = \{1, 2, 3, 4, 6, 7, 8\}$ $E = \{1, 2, 4, 5, 6, 7, 10\}$
 $C = \{1, 2, 4, 5, 6, 7\}$

	a	b
A	B	C
B	B	D
C	B	C
D	B	E
E	B	C



CS502

Scanning

23

Limits of regular languages

Not all languages are regular

One cannot construct DFAs to recognize these languages:

- $L = \{p^k q^k\}$
- $L = \{wcw^r \mid w \in \Sigma^*\}$

Note: neither of these is a regular expression!
 (DFAs cannot count!)

But, this is a little subtle. One can construct DFAs for:

- alternating 0's and 1's
 $(\epsilon \mid 1)(01)^*(\epsilon \mid 0)$
- sets of pairs of 0's and 1's
 $(01 \mid 10)^+$

CS502

Scanning

24

So what is hard?

Language features that can cause problems:

reserved words

PL/I had no reserved words

```
if then then then = else; else else = then;
```

significant blanks

FORTRAN and Algol68 ignore blanks

```
do 10 i = 1,25
```

```
do 10 i = 1.25
```

string constants

special characters in strings

newline, tab, quote, comment delimiter

finite closures

some languages limit identifier lengths

adds states to count length

FORTRAN 66 → 6 characters

These can be swept under the rug in the language design

CS502

Scanning

25

How bad can it get?

```
1      INTEGERFUNCTIONA
2      PARAMETER (A=6,B=2)
3      IMPLICIT CHARACTER*(A-B) (A-B)
4      INTEGER FORMAT(10),IF(10),D09E1
5      100  FORMAT(4H)=(3)
6      200  FORMAT(4 )=(3)
7      D09E1=1
8      D09E1=1,2
9          IF(X)=1
10         IF(X)H=1
11         IF(X)300,200
12      300  CONTINUE
13      END
        C    this is a comment
        $    FILE(1)
14      END
```

Example due to Dr. F.K. Zadeck of IBM Corporation

CS502

Scanning

26

Scanning MiniJava

White space:

- ' ', '\t', '\n', '\r', '\f'

Tokens:

- Operators, keywords (straightforward; I've done them for you)
- Identifiers (straightforward)
- Integers (straightforward)
- Strings (tricky for escapes)

CS502

Scanning

27