



Lexical Analysis - An Introduction

Lecture 4

Spring 2005

Department of Computer Science

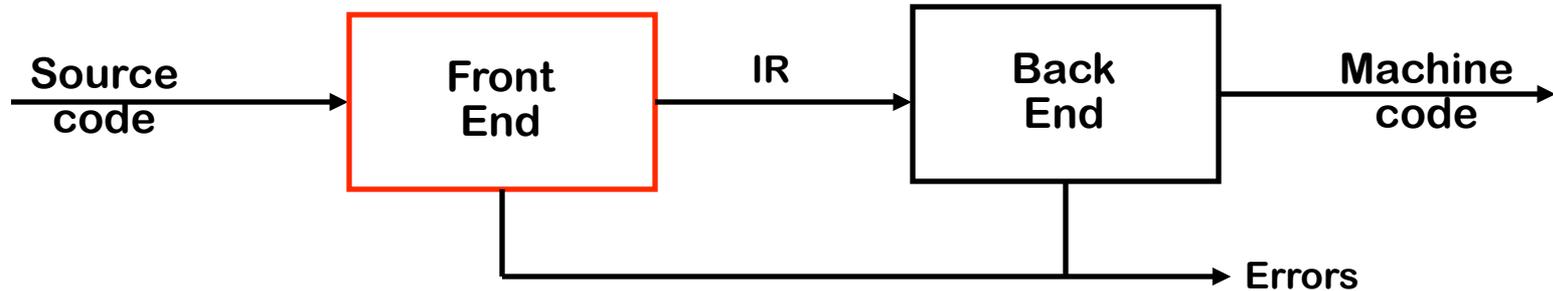
University of Alabama

Joel Jones

Copyright 2003, Keith D. Cooper, Ken Kennedy & Linda Torczon, all rights reserved.
Students enrolled in Comp 412 at Rice University have explicit permission to make copies of these materials for their personal use.



The Front End

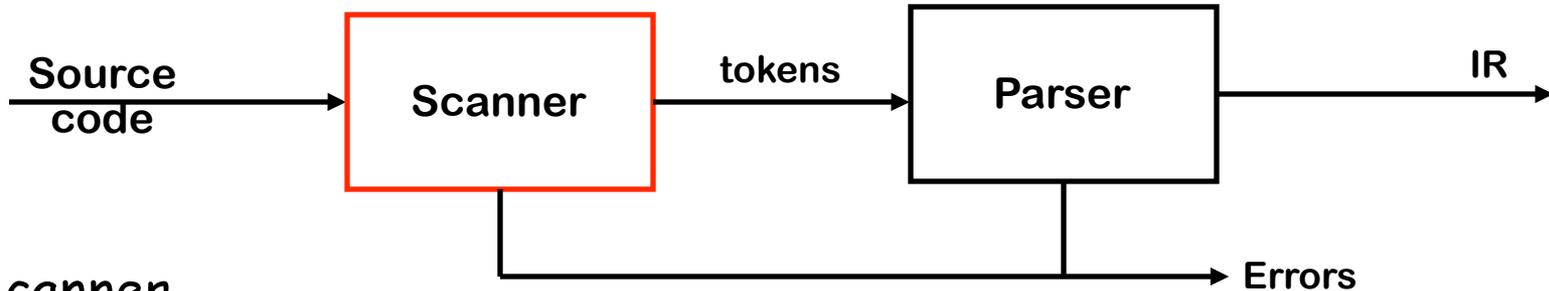


The purpose of the front end is to deal with the input language

- Perform a membership test: $\text{code} \in \text{source language?}$
- Is the program well-formed (semantically)?
- Build an IR version of the code for the rest of the compiler

The front end is not monolithic

The Front End



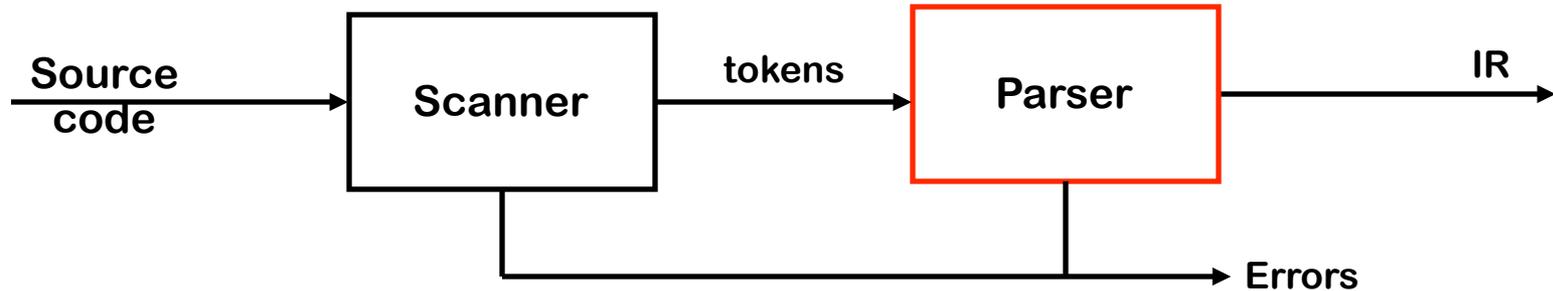
Scanner

- Maps stream of characters into words
 - Basic unit of syntax
 - $x = x + y ;$ becomes
 $\langle id, x \rangle \langle eq, = \rangle \langle id, x \rangle \langle pl, + \rangle \langle id, y \rangle \langle sc, ; \rangle$
- Characters that form a word are its **lexeme**
- Its part of speech (or syntactic category) is called its **token type**
- Scanner discards white space & (often) comments

Speed is an issue in scanning
⇒ use a specialized recognizer



The Front End



Parser

- Checks stream of classified words (parts of speech) for grammatical correctness
- Determines if code is syntactically well-formed
- Guides checking at deeper levels than syntax
- Builds an IR representation of the code

We'll come back to parsing in a couple of lectures

The Big Picture



- Language syntax is specified with parts of speech, not words
- Syntax checking matches parts of speech against a grammar

1. goal \rightarrow expr
2. expr \rightarrow expr op term
3. | term
4. term \rightarrow number
5. | id
6. op \rightarrow +
7. | -

S = goal
T = { number, id, +, - }
N = { goal, expr, term, op }
P = { 1, 2, 3, 4, 5, 6, 7 }



The Big Picture

- Language syntax is specified with parts of speech, not words
- Syntax checking matches parts of speech against a grammar

1. goal \rightarrow expr
2. expr \rightarrow expr op term
3. | term
4. term \rightarrow number
5. | id
6. op \rightarrow +
7. | -

S = goal
T = { number, id, +, - }
N = { goal, expr, term, op }
P = { 1, 2, 3, 4, 5, 6, 7 }

No words here!

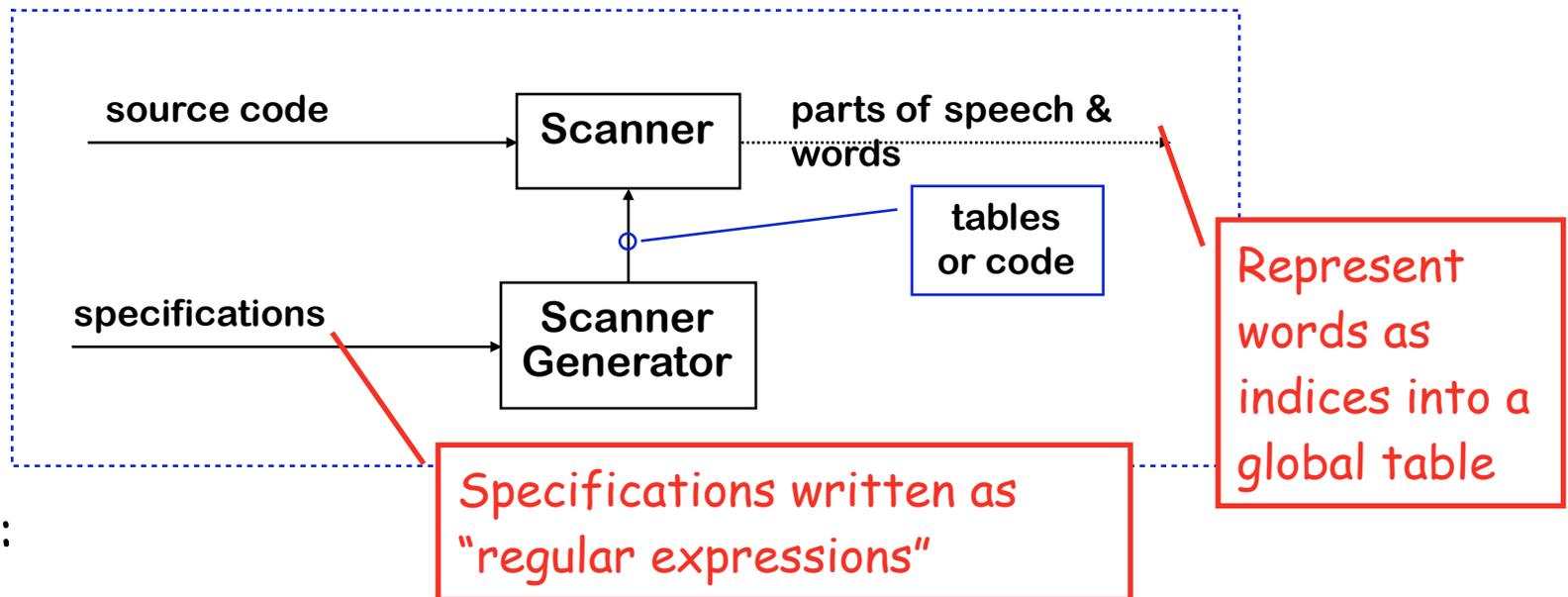
Parts of speech,
not words!



The Big Picture

Why study lexical analysis?

- We want to avoid writing scanners by hand



Goals:

- To simplify specification & implementation of scanners
- To understand the underlying techniques and technologies



Regular Expressions

Lexical patterns form a **regular language**

*** any finite language is regular ***

Ever type
"rm *.o a.out" ?

Regular expressions (REs) describe regular languages

Regular Expression (over alphabet Σ)

- ϵ is a RE denoting the set $\{\epsilon\}$
- If a is in Σ , then a is a RE denoting $\{a\}$
- If x and y are REs denoting $L(x)$ and $L(y)$ then
 - $x \mid y$ is an RE denoting $L(x) \cup L(y)$
 - xy is an RE denoting $L(x)L(y)$
 - x^* is an RE denoting $L(x)^*$

Precedence is closure, then concatenation, then alternation

Set Operations

(review)



Operation	Definition
<i>Union of L and M</i> <i>Written $L \cup M$</i>	$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$
<i>Concatenation of L and M</i> <i>Written LM</i>	$LM = \{st \mid s \in L \text{ and } t \in M\}$
<i>Kleene closure of L</i> <i>Written L^*</i>	$L^* = \bigcup_{0 \leq i \leq \infty} L^i$
<i>Positive Closure of L</i> <i>Written L^+</i>	$L^+ = \bigcup_{1 \leq i \leq \infty} L^i$

These definitions should be well known



Examples of Regular Expressions

Identifiers:

Letter \rightarrow (a|b|c ... |z|A|B|C ... |Z)

Digit \rightarrow (0|1|2 ... |9)

Identifier \rightarrow Letter (Letter | Digit)^{*}

0 or number
with no leading
zeros

Numbers:

Integer \rightarrow (+|-|ε) (0 | (1|2|3 ... |9)(Digit^{*}))

Decimal \rightarrow Integer . Digit^{*}

Real \rightarrow (Integer | Decimal) E (+|-|ε) Digit^{*}

Complex \rightarrow (Real , Real)

Numbers can get much more complicated!

Regular Expressions

(the point)



Regular expressions can be used to specify the words to be translated to parts of speech by a lexical analyzer

Using results from automata theory and theory of algorithms, we can automatically build recognizers from regular expressions

⇒ We study REs and associated theory to automate scanner construction !



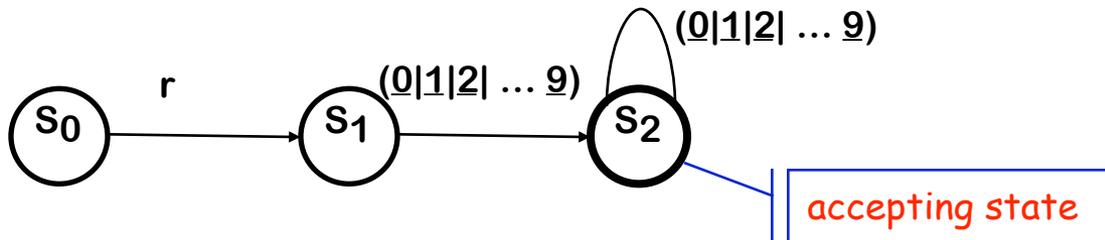
Example

Consider the problem of recognizing ILOC register names

Register $\rightarrow r (0|1|2| \dots | 9) (0|1|2| \dots | 9)^*$

- Allows registers of arbitrary number
- Requires at least one digit

RE corresponds to a recognizer (or DFA)



Recognizer for Register

Transitions on other inputs go to an error state, s_e

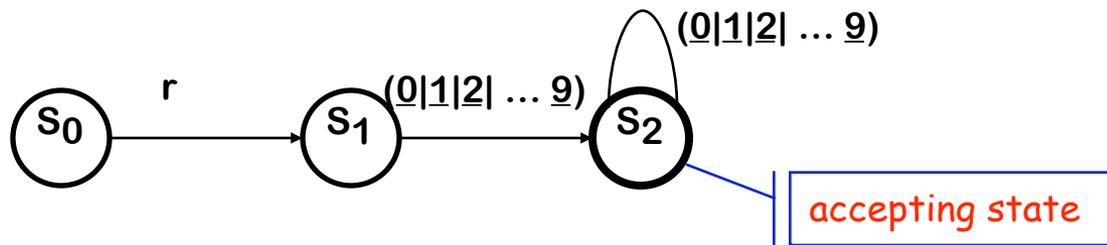
Example

(continued)



DFA operation

- Start in state S_0 & take transitions on each input character
- DFA accepts a word \underline{x} iff \underline{x} leaves it in a final state (S_2)



Recognizer for Register

So,

- r17 takes it through s_0, s_1, s_2 and accepts
- r takes it through s_0, s_1 and fails
- a takes it straight to s_e

Example

(continued)



To be useful, recognizer must turn into code

```
Char ← next character
State ← s0

while (Char ≠ EOF)
    State ← δ(State,Char)
    Char ← next character

if (State is a final state )
    then report success
    else report failure
```

Skeleton recognizer

δ	r	0,1,2,3,4, 5,6,7,8,9	All others
s0	s1	s _e	s _e
s1	s _e	s2	s _e
s2	s _e	s2	s _e
s _e	s _e	s _e	s _e

Table encoding RE

Example

(continued)



To be useful, recognizer must turn into code

```
Char ← next character
State ← s0

while (Char ≠ EOF)
    State ← δ(State,Char)
    perform specified action
    Char ← next character

if (State is a final state )
    then report success
    else report failure
```

Skeleton recognizer

δ	r	0,1,2,3,4, 5,6,7,8,9	All others
s0	s1 start	s _e error	s _e error
s1	s _e error	s2 add	s _e error
s2	s _e error	s2 add	s _e error
s _e	s _e error	s _e error	s _e error

Table encoding RE



What if we need a tighter specification?

\underline{r} Digit Digit* allows arbitrary numbers

- Accepts r00000
- Accepts r99999
- What if we want to limit it to r0 through r31 ?

Write a tighter regular expression

→ Register → $\underline{r} ((\underline{0}|\underline{1}|\underline{2}) (\text{Digit} | \epsilon) | (\underline{4}|\underline{5}|\underline{6}|\underline{7}|\underline{8}|\underline{9}) | (\underline{3}|\underline{30}|\underline{31}))$

→ Register → $\underline{r0}|\underline{r1}|\underline{r2} | \dots | \underline{r31}|\underline{r00}|\underline{r01}|\underline{r02} | \dots | \underline{r09}$

Produces a more complex DFA

- Has more states
- Same cost per transition
- Same basic implementation

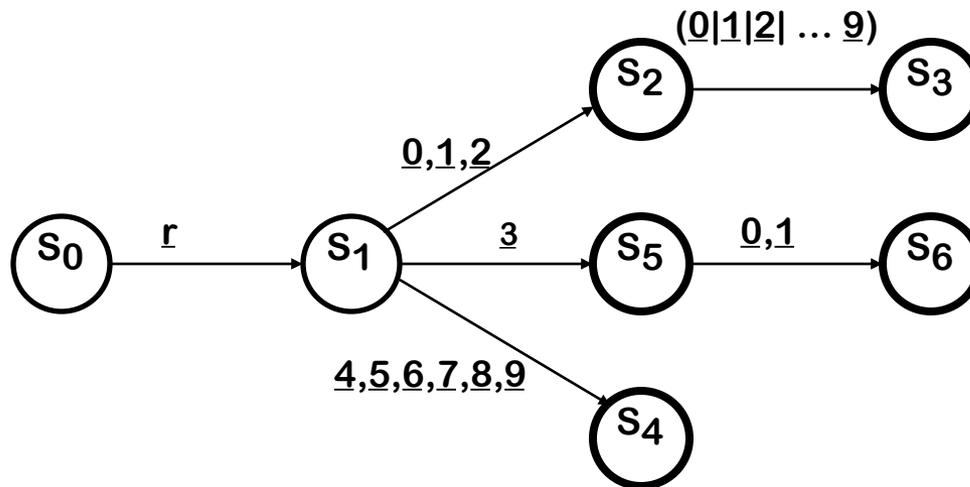
Tighter register specification

(continued)



The DFA for

Register $\rightarrow \underline{r} ((\underline{0}|\underline{1}|\underline{2}) (\text{Digit} | \varepsilon) | (\underline{4}|\underline{5}|\underline{6}|\underline{7}|\underline{8}|\underline{9}) | (\underline{3}|\underline{30}|\underline{31}))$



- Accepts a more constrained set of registers
- Same set of actions, more states

Tighter register specification

(continued)



δ	r	0,1	2	3	4-9	All others
s0	s1	s _e				
s1	s _e	s2	s2	s5	s4	s _e
s2	s _e	s3	s3	s3	s3	s _e
s3	s _e					
s4	s _e					
s5	s _e	s6	s _e	s _e	s _e	s _e
s6	s _e					
s _e						

Runs in the same skeleton recognizer

Table encoding RE for the tighter register specification



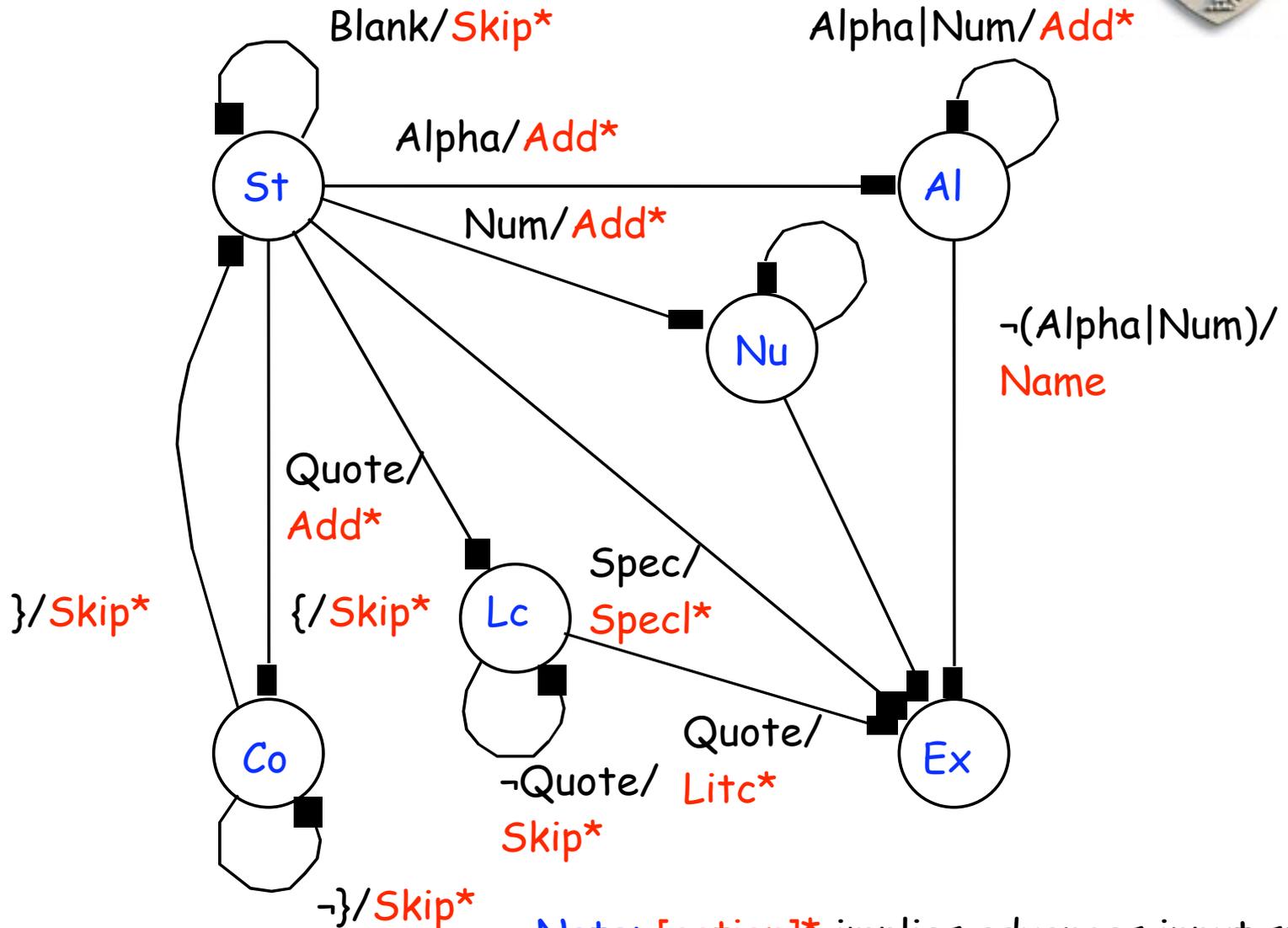
Extra Slides Start Here

Principles of Scanners



- Lexical Analysis Strategy: **Simulation of Finite Automaton**
 - States, characters, actions
 - State transition $\delta(\text{state}, \text{charclass})$ determines next state
- **Next character** function
 - Reads next character into buffer
 - Computes **character class** by fast table lookup
- Transitions from state to state
 - Current state and next character determine (via δ)
 - Next state and action to be performed
 - Some actions preload next character
- Identifiers distinguished from keywords by hashed lookup
 - This differs from EAC advice (discussion later)
 - Permits translation of identifiers into **<type, symbol_index>**
 - Keywords each get their own type

A Lexical Analysis Example



Note: [action]* implies advances input stream

Example Lexical Scan Code



```
current = START_STATE;
token = "";
// assume next character has been preloaded into a buffer
while (current != EX)
{
    int charClass = inputstream->thisClass();
    switch (current->action(charClass))
    {
        case SKIP:
            inputstream->advance();break;
        case ADD:
            char* t = token; int n = ::strlen(t);
            token = new char[n + 2]; ::strcpy(token, t);
            token[n] = inputstream->thisChar(); token[n+1] = 0;
            delete [] t; inputstream->advance(); break;
        case NAME:
            Entry * e = symTable->lookup(token);
            tokenType = (e->type==NULL_TYPE ? NAME_TYPE : e->type);
            break;
        ...
    }
    current = current->nextState(charClass);
}
```

Tighter register specification

(continued)



state action	r	0,1	2	3	4,5,6 7,8,9	other
0	1 <u>start</u>	e	e	e	e	e
1	e	2 <u>add</u>	2 <u>add</u>	5 <u>add</u>	4 <u>add</u>	e
2	e	3 <u>add</u>	3 <u>add</u>	3 <u>add</u>	3 <u>add</u>	x <u>exit</u>
3,4	e	e	e	e	e	x <u>exit</u>
5	e	6 <u>add</u>	e	e	e	x <u>exit</u>
6	e	e	e	e	e	x <u>exit</u>
e	e	e	e	e	e	e