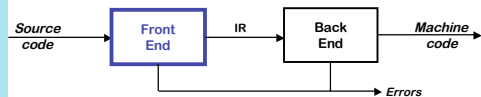


Lexical Analysis

Introduction

Copyright 2009, Pedro C. Diniz, all rights reserved.
Students enrolled in the Compilers class at Instituto Superior Técnico (IST/UTL) have explicit permission to make copies of these materials for their personal use.

The Front End

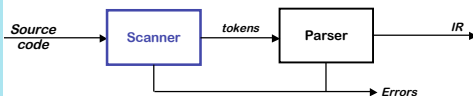


The purpose of the front end is to deal with the input language

- 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

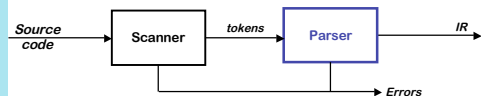


Scanner

- Maps stream of characters into words
 - Basic unit of syntax
 - $x = x + y$; becomes $\langle \text{id}, x \rangle \langle \text{eq}, = \rangle \langle \text{id}, x \rangle \langle \text{pl}, + \rangle \langle \text{id}, y \rangle \langle \text{sc}, ; \rangle$
- 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
 \Rightarrow 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 → expr
2. expr → expr op term
3.   | term
4. term → number
5.   | id
6. op  → +
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 → expr
2. expr → expr op term
3.   | term
4. term → number
5.   | id
6. op  → +
7.   | -
    
```

```

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

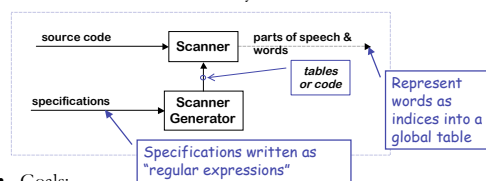
No words here!

Parts of speech,
not words!

The Big Picture

Why study Lexical Analysis?

- We want to avoid writing scanners by hand
- We want to harness the theory classes



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

What is a Lexical Analyzer?

Source program text → Tokens

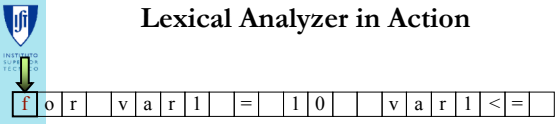
Example of Tokens

- Operators: `= + - > ({ := == <>`
- Keywords: `if while for int double`
- Numeric literals: `43 6.035 -3.6e10 0x13F3A`
- Character literals: `'a' '~' '\'`
- String literals: `"6.891" "Fall 06" "\\\" = empty"`

Example of non-tokens

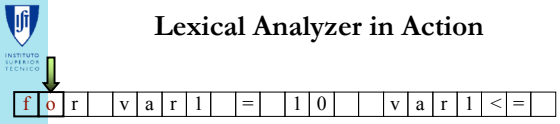
- White space: `space(' ') tab('\t') end-of-line('\n')`
- Comments: `/*this is not a token*/`

Lexical Analyzer in Action



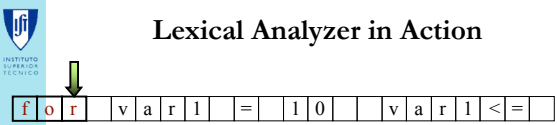
The diagram shows a lexical analyzer's initial state. A green arrow points to the first character 'f' in the input string 'for var l = l 0 var l < ='. The input is displayed in a sequence of boxes: f, o, r, space, v, a, r, space, l, space, =, space, l, space, 0, space, v, a, r, space, l, space, <, space, =.

Lexical Analyzer in Action



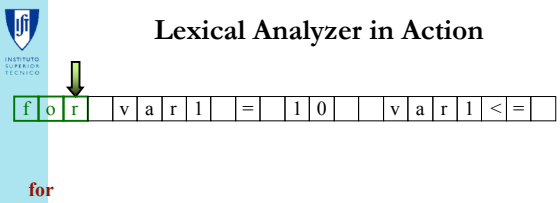
The diagram shows the lexical analyzer moving to the second character 'o'. The green arrow now points to 'o', and the 'f' is highlighted in red.

Lexical Analyzer in Action




The diagram shows the lexical analyzer moving to the third character 'r'. The green arrow now points to 'r', and both 'f' and 'o' are highlighted in red.

Lexical Analyzer in Action




The diagram shows the lexical analyzer identifying the keyword 'for'. The green arrow points to the space after 'r', and the entire sequence 'for' is highlighted in green. Below the input sequence, the word 'for' is printed in red, indicating it is a recognized keyword.



Lexical Analyzer in Action

f	o	r		v	a	r	l	=		l	0		v	a	r	l	<	=	
---	---	---	--	---	---	---	---	---	--	---	---	--	---	---	---	---	---	---	--


for



Lexical Analyzer in Action

f	o	r		v	a	r	l	=		l	0		v	a	r	l	<	=	
---	---	---	--	---	---	---	---	---	--	---	---	--	---	---	---	---	---	---	--


for



Lexical Analyzer in Action

f	o	r		v	a	r	l	=		l	0		v	a	r	l	<	=	
---	---	---	--	---	---	---	---	---	--	---	---	--	---	---	---	---	---	---	--


for



Lexical Analyzer in Action

f	o	r		v	a	r	l	=		l	0		v	a	r	l	<	=	
---	---	---	--	---	---	---	---	---	--	---	---	--	---	---	---	---	---	---	--


for



Lexical Analyzer in Action

f	o	r	v	a	r	l	=	l	0		v	a	r	l	<=	
---	---	---	---	---	---	---	---	---	---	--	---	---	---	---	----	--


for



Lexical Analyzer in Action

f	o	r	v	a	r	l	=	l	0		v	a	r	l	<=	
---	---	---	---	---	---	---	---	---	---	--	---	---	---	---	----	--


for



Lexical Analyzer in Action

f	o	r	v	a	r	l	=	l	0		v	a	r	l	<=	
---	---	---	---	---	---	---	---	---	---	--	---	---	---	---	----	--


for ID("varl")



Lexical Analyzer in Action

f	o	r	v	a	r	l	=	l	0		v	a	r	l	<=	
---	---	---	---	---	---	---	---	---	---	--	---	---	---	---	----	--


for ID("varl")



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		l	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1")



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		l	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1")



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		l	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1") eq_op



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		l	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1") eq_op



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		1	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1") eq_op



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		1	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1") eq_op



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		1	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--


for ID("var1") eq_op



Lexical Analyzer in Action

f	o	r		v	a	r		l		=		1	0			v	a	r		l		<	=	
---	---	---	--	---	---	---	--	---	--	---	--	---	---	--	--	---	---	---	--	---	--	---	---	--

for ID("var1") eq_op Num(10)




INSTITUTO
SUPERIOR
TÉCNICO

Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l	<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	---	---	--

for ID("var1") eq_op Num(10)




INSTITUTO
SUPERIOR
TÉCNICO

Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l	<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	---	---	--

for ID("var1") eq_op Num(10)




INSTITUTO
SUPERIOR
TÉCNICO

Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l	<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	---	---	--

for ID("var1") eq_op Num(10)



INSTITUTO
SUPERIOR
TÉCNICO

Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l	<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	---	---	--

for ID("var1") eq_op Num(10)

Lexical Analyzer in Action

f	o	r		v	a	r		=		1	0		v	a	r		<	=	
---	---	---	--	---	---	---	--	---	--	---	---	--	---	---	---	--	---	---	--

for ID("var1") eq_op Num(10)

Lexical Analyzer in Action

f	o	r		v	a	r		=		1	0		v	a	r		<	=	
---	---	---	--	---	---	---	--	---	--	---	---	--	---	---	---	--	---	---	--

for ID("var1") eq_op Num(10)

Lexical Analyzer in Action


f	o	r		v	a	r		=		1	0		v	a	r		<	=	
---	---	---	--	---	---	---	--	---	--	---	---	--	---	---	---	--	---	---	--

for ID("var1") eq_op Num(10)

Lexical Analyzer in Action

f	o	r		v	a	r		=		1	0		v	a	r		<	=	
---	---	---	--	---	---	---	--	---	--	---	---	--	---	---	---	--	---	---	--

for ID("var1") eq_op Num(10)




Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l		<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	--	---	---	--

↓

for ID("var1") eq_op Num(10) ID("var1")




Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l		<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	--	---	---	--

↓

for ID("var1") eq_op Num(10) ID("var1")




Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l		<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	--	---	---	--

↓

for ID("var1") eq_op Num(10) ID("var1")



Lexical Analyzer in Action

f	o	r		v	a	r	l		=		1	0			v	a	r	l		<	=	
---	---	---	--	---	---	---	---	--	---	--	---	---	--	--	---	---	---	---	--	---	---	--

↓

for ID("var1") eq_op Num(10) ID("var1") leq_op

Lexical Analyzer in Action

The diagram shows a sequence of characters in a grid: f, o, r, space, v, a, r, 1, space, =, space, 1, 0, space, v, a, r, 1, space, <, =, space. A green arrow points to the space after the second '='. Below the grid, the tokens are listed: **for**, ID("var1"), eq_op, Num(10), ID("var1"), leq_op.

Lexical Analyzer in Action

The diagram shows a sequence of characters in a grid: f, o, r, space, v, a, r, 1, space, =, space, 1, 0, space, v, a, r, 1, space, <, =, space. A green arrow points to the space after the second '<='. Below the grid, the tokens are listed: **for**, ID("var1"), eq_op, Num(10), ID("var1"), leq_op.

Lexical Analyzer needs to...

- Partition Input Program Text into Subsequence of Characters Corresponding to Tokens
- Attach the Corresponding Attributes to the Tokens
- Eliminate White Space and Comments

Lexical Analysis: Basic Issues

- How to Precisely Match Strings to Tokens
- How to Implement a Lexical Analyzer

Regular Expressions

Lexical patterns form a *regular language*

*** any finite language is regular ***

Regular expressions (REs) describe regular languages

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

Regular Expression (over alphabet Σ)

- ϵ is a RE denoting the set $\{\epsilon\}$
- If \underline{a} is in Σ , then \underline{a} is a RE denoting $\{a\}$
- If x and y are REs denoting $L(x)$ and $L(y)$ then
 - $x | y$ 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
Union of L and M Written $L \cup M$	$L \cup M = \{s \mid s \in L \text{ or } s \in M\}$
Concatenation of L and M Written LM	$LM = \{st \mid s \in L \text{ and } t \in M\}$
Kleene closure of L Written L^*	$L^* = \bigcup_{0 \leq i < \infty} L^i$
Positive Closure of L Written L^+	$L^+ = \bigcup_{1 \leq i < \infty} L^i$

These definitions should be well known

Examples of Regular Expressions

Identifiers:

Letter $\rightarrow (a|b|c| \dots |z|A|B|C| \dots |Z)$
 Digit $\rightarrow (0|1|2| \dots |9)$
 Identifier $\rightarrow \text{Letter} (\text{Letter} | \text{Digit})^*$

Numbers:

Integer $\rightarrow (\pm | \epsilon) (0|1|2| \dots |9)^+$
 Decimal $\rightarrow \text{Integer} . \text{Digit}^*$
 Real $\rightarrow (\text{Integer} | \text{Decimal}) E (\pm | \epsilon) \text{Digit}^*$
 Complex $\rightarrow \text{Real} . \text{Real}$

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

Some of you may have seen this construction for string pattern matching

\Rightarrow We study REs and associated theory to automate scanner construction !

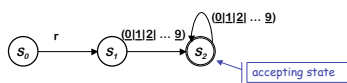
Example

Consider the problem of recognizing 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)



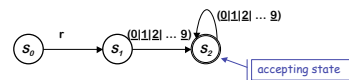
Recognizer for Register

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 x iff 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
s_0	s_1	s_e	s_e
s_1	s_e	s_2	s_e
s_2	s_e	s_2	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
s_0	s_1	s_e	s_e
s_1	start	error	error
s_2	error	add	error
s_e	error	error	error

Table encoding RE

What about a Tighter Specification?

\mathbb{R} *Digit Digit*^{*} allows arbitrary numbers

- Accepts $\mathbb{R}00000$
- Accepts $\mathbb{R}99999$
- What if we want to limit it to $\mathbb{R}0$ through $\mathbb{R}31$?

Write a tighter regular expression

- $\text{Register} \rightarrow \mathbb{R} (0|1|2) (\text{Digit} | \epsilon) | (4|5|6|7|8|9) | (3|30|31)$
- $\text{Register} \rightarrow \mathbb{R} | \mathbb{R}1 | \mathbb{R}2 | \dots | \mathbb{R}31 | \mathbb{R}00 | \mathbb{R}01 | \mathbb{R}02 | \dots | \mathbb{R}09$

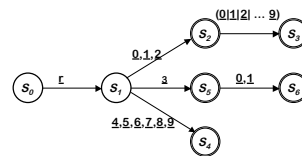
Produces a more complex DFA

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

Tighter Register Specification (cont'd)

The DFA for

$\text{Register} \rightarrow \mathbb{R} (0|1|2) (\text{Digit} | \epsilon) | (4|5|6|7|8|9) | (3|30|31)$



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

Tighter Register Specification (cont'd)

δ	\mathbb{R}	0,1	2	3	4-9	All other s
s_0	s_1	s_e	s_e	s_e	s_e	s_e
s_1	s_e	s_2	s_2	s_3	s_4	s_e
s_2	s_e	s_3	s_3	s_3	s_3	s_e
s_3	s_e	s_e	s_e	s_e	s_e	s_e
s_4	s_e	s_e	s_e	s_e	s_e	s_e
s_5	s_e	s_e	s_e	s_e	s_e	s_e
s_6	s_e	s_e	s_e	s_e	s_e	s_e
s_e	s_e	s_e	s_e	s_e	s_e	s_e

Runs in the same skeleton recognizer

Table encoding RE for the tighter register specification

Summary

- The role of the lexical Analyzer
 - Partition input stream into tokens
- Regular Expressions
 - Used to describe structure of tokens
- DFA: Deterministic Finite Automata
 - Machinery to recognize Regular Languages