

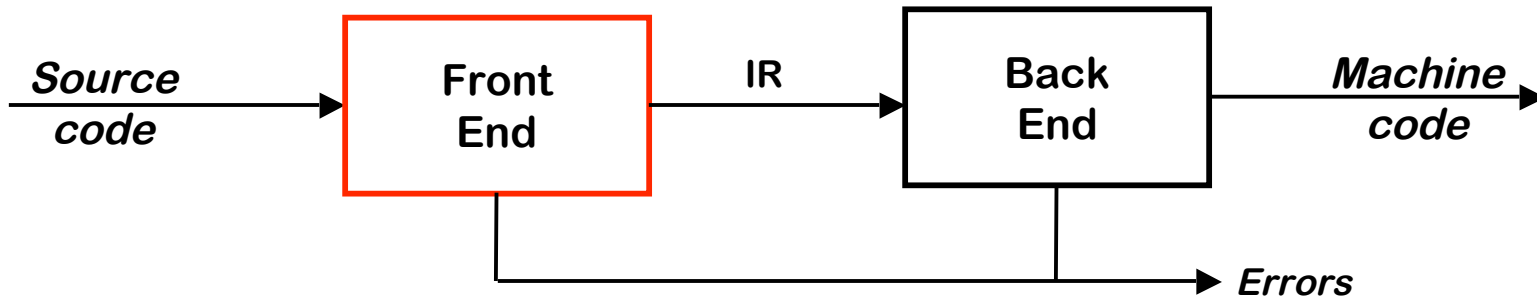


Lexical Analysis - An Introduction

COMP 412
Fall 2005

Copyright 2005, 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

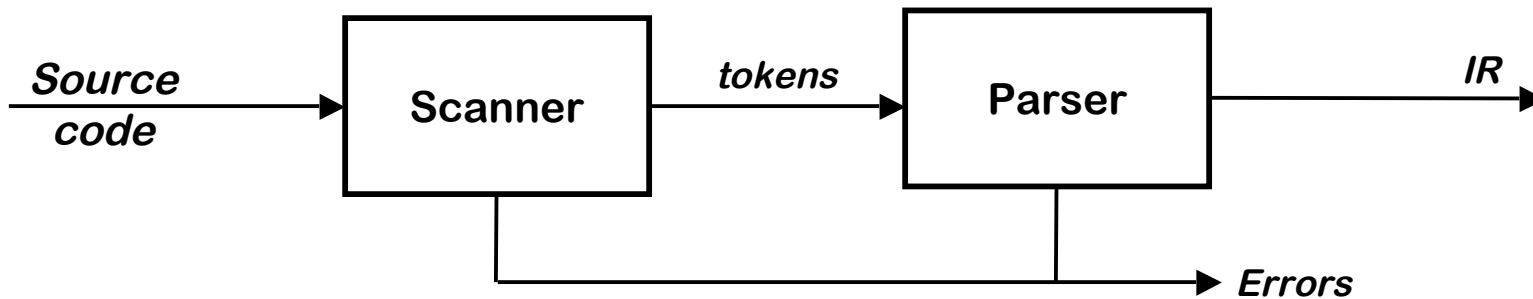


The purpose of the front end is to deal with the input language

- 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

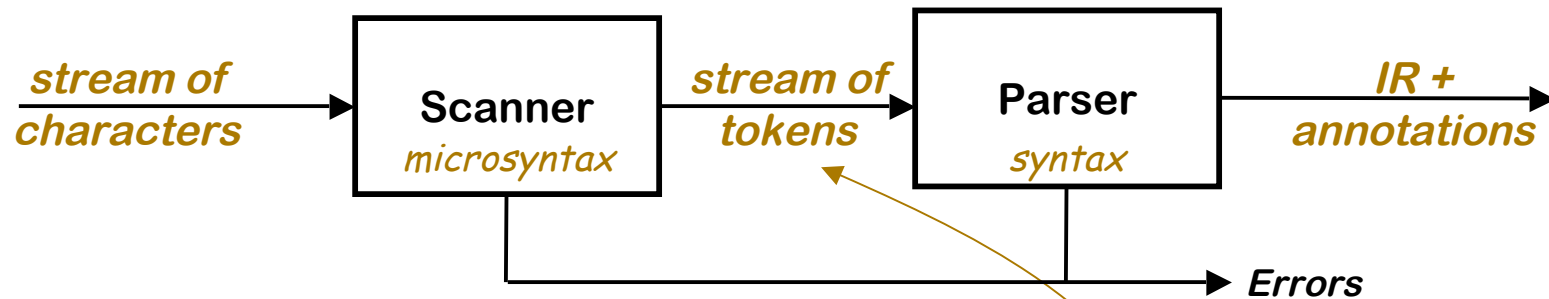


Implementation Strategy

- Specify syntax in a formal notation
 - regular expressions in scanning, context-free grammars in parsing
- Simulate an automaton to recognize valid strings
 - finite automata, push-down automata
- Automate construction of the simulations
 - table-driven simulations or direct-coded simulations
- Add "actions" to automaton to create representations

(COMP 481)

The Front End



Why separate the scanner and the parser?

- Scanner classifies words
- Parser constructs grammatical derivations
- Parsing is harder and slower
- Separation simplifies implementation
 - smaller grammar for parser
 - faster front end

token is a pair
<part of speech, lexeme>

also called *syntactic categories* or *tokentypes*



The Big Picture

The front end deals with syntax

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

Simple expression grammar from lecture 2

1. *goal* → *expr*

2. *expr* → *exp*

3. | *ter*

4. *term* → *num*

5. | *id*

6. *op* → *+*

7. | *-*

The scanner turns a stream of characters into a stream of words, classified with their part of speech.

$N = \{ \textit{goal}, \textit{expr}, \textit{term}, \textit{op} \}$

$P = \{ 1, 2, 3, 4, 5, 6, 7 \}$

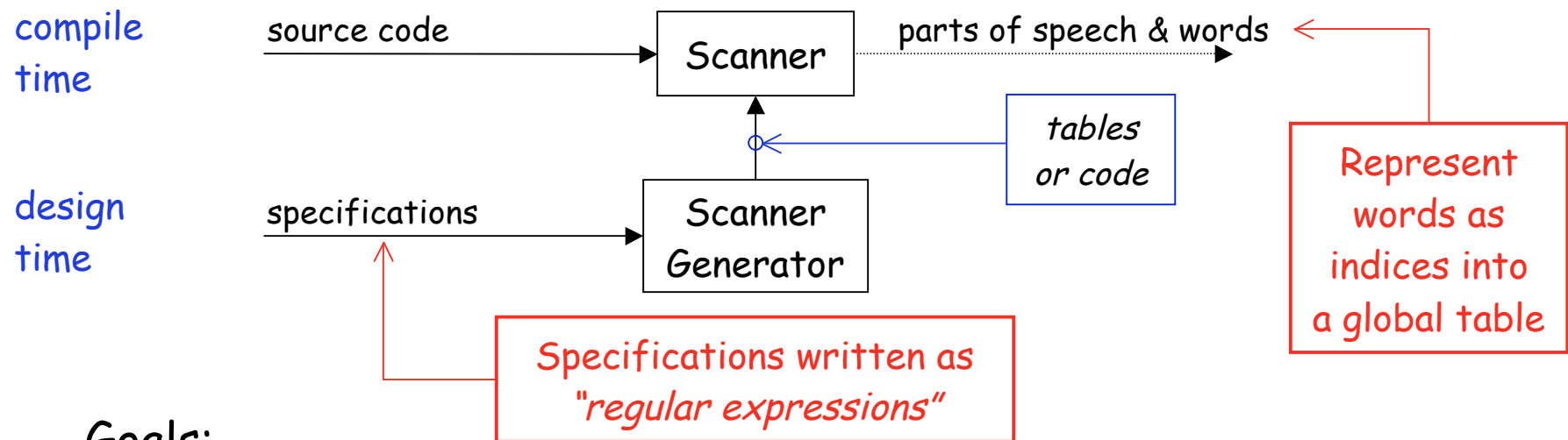
parts of speech
syntactic variables

The Big Picture



Why study scanner construction?

- We want to avoid writing scanners by hand
- We want to harness the theory from classes like COMP 481



Goals:

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

Set Operations

(review)



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

These definitions should be well known

Regular Expressions



We constrain programming languages so that the spelling of a word always implies its part of speech *(few exceptions)*

The rules or patterns that impose this mapping form a *regular language*

Regular expressions (REs) describe regular languages

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 $\{\underline{a}\}$
- 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*

Regular Expressions



How do these operators help?

Regular Expression (over alphabet Σ)

- ε is a RE denoting the set $\{\varepsilon\}$
- If \underline{a} is in Σ , then \underline{a} is a RE denoting $\{a\}$
 - the spelling of a word is an RE
- If x and y are REs denoting $L(x)$ and $L(y)$ then
 - $x | y$ is an RE denoting $L(x) \cup L(y)$
 - any finite list of words can be written as an RE $(w_0 / w_1 / \dots / w_n)$
 - xy is an RE denoting $L(x)L(y)$
 - x^* is an RE denoting $L(x)^*$
 - we can use concatenation & closure to write more concise patterns and to specify infinite sets that have finite descriptions

Examples of Regular Expressions



Identifiers:

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

Digit → (0|1|2 | ... |9)

Identifier → *Letter* (*Letter* | *Digit*)*

shorthand
for

$(\underline{a}|\underline{b}|\underline{c} | \dots | \underline{z}|\underline{A}|\underline{B}|\underline{C} | \dots | \underline{Z}) ((\underline{a}|\underline{b}|\underline{c} | \dots | \underline{z}|\underline{A}|\underline{B}|\underline{C} | \dots | \underline{Z}) | (\underline{0}|\underline{1}|\underline{2} | \dots | \underline{9}))^*$

Numbers:

Integer → (+|-|ε) (0 | (1|2|3 | ... |9)(*Digit**))

Decimal → *Integer* . *Digit**

Real → (*Integer* | *Decimal*) E (+|-|ε) *Digit**

Complex → (*Real* , *Real*)

underlining indicates
a letter in the input
stream

Numbers can get much more complicated!

Regular Expressions

(the point)



We use regular expressions to specify the mapping of words to parts of speech for the lexical analyzer

Using results from automata theory and theory of algorithms, we can automate construction of recognizers from REs

- ⇒ We study REs and associated theory to automate scanner construction !
- ⇒ Fortunately, the automatic techniques lead to fast scanners
 - used in text editors, URL filtering software, ...

Example

(from Lab 1)

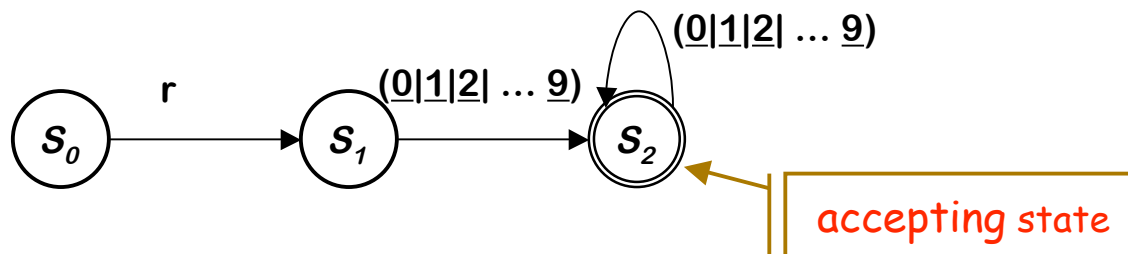


Consider the problem of recognizing ILOC register names

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

- 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

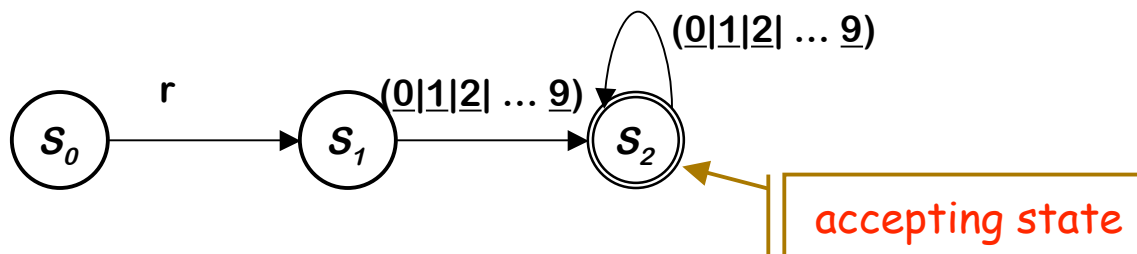
Example

(continued)



DFA operation

- Start in state S_0 & make 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, the recognizer must be converted 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 the RE

Example

(continued)



We can add "actions" to each transition

```
Char ← next character
State ← s0

while (Char ≠ EOF)
  Next ← δ(State,Char)
  Act ← α(State,Char)
  perform action Act
  State ← Next
  Char ← next character

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

Skeleton recognizer

δ	α	0,1,2,3,4, 5,6,7,8,9	All others
s_0	s_1 <i>start</i>	s_e <i>error</i>	s_e <i>error</i>
s_1	s_e <i>error</i>	s_2 <i>add</i>	s_e <i>error</i>
s_2	s_e <i>error</i>	s_2 <i>add</i>	s_e <i>error</i>
s_e	s_e <i>error</i>	s_e <i>error</i>	s_e <i>error</i>

Table encoding RE



What if we need a tighter specification?

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* \rightarrow r ((0|1|2) (*Digit* | ϵ) | (4|5|6|7|8|9) | (3|30|31))
- *Register* \rightarrow r0|r1|r2| ... |r31|r00|r01|r02| ... |r09

Produces a more complex DFA

- DFA has more states
- DFA has **same cost** per transition *(or per character)*
- DFA has same basic implementation

