



RICE

COMP 412
FALL 2011

Lexical Analysis – Introduction

Comp 412

The slides assume some familiarity with finite automata.
For a different (& more intuitive?) introduction to finite automata & recognizers, see Section 2.2 of EaC

Copyright 2011, Keith D. Cooper & 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.
Faculty from other educational institutions may use these materials for nonprofit educational purposes, provided this copyright notice is preserved.

Lab 1 Discussion



1. Get started - two weeks left
2. The lab manipulates the syntax of the input program, not its meaning or value
 - Read the code into some data structure
 - Rename values so that each operation defines a unique name
 - *only valid in straightline code; makes names ≡ to live ranges*
 - Run the allocator
 - Write the results back out as valid ILOC
3. First step - write something to read and write the ILOC
 - Keep it simple
4. If you look at the example code in § 13.2 of EaC, keep in mind that you only have 1 register class in your machine
 - Substantial simplification to the code

↙ defines ≡ assigns

Lab 1 Discussion

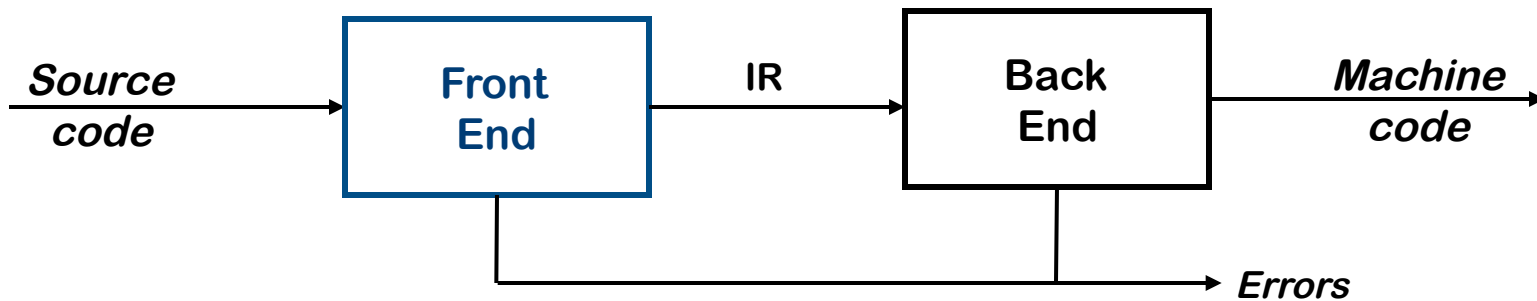


Reading and writing ILOC

- Use a simple data structure - say a $k \times 4$ array of integers
 - Turn opcode into an integer
 - *Simple function using string comparisons*
 - Represent registers and immediates as integers
 - Prettyprinter can interpret those ints contextually

- This part of the lab should take an evening or two (The allocators are the tricky part, not the I/O.)

The Front End

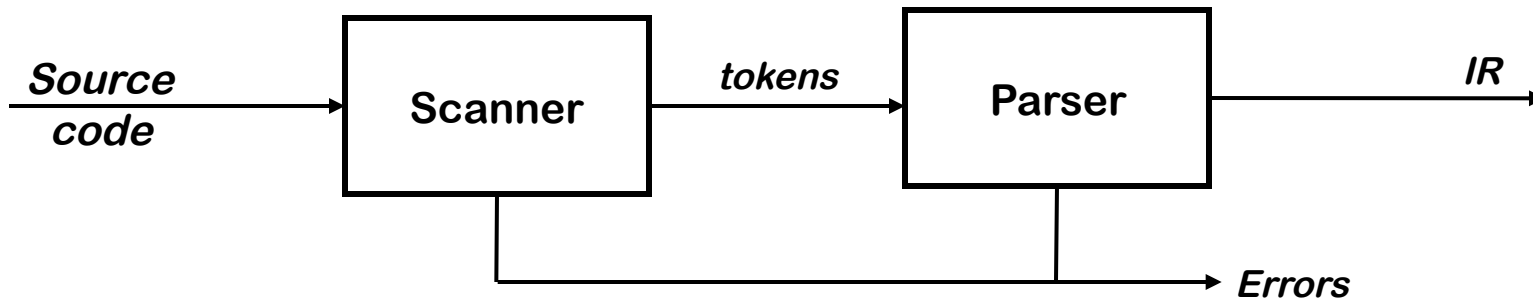


The purpose of the front end is to deal with the input language

- 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 deals with form (syntax) & meaning (semantics)

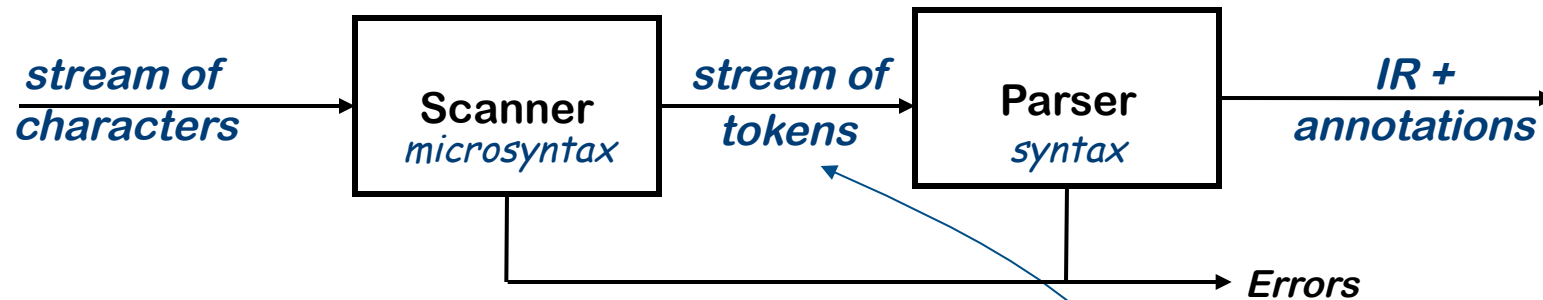
The Front End



Implementation Strategy

	Scanning	Parsing
Specify Syntax	regular expressions	context-free grammars
Implement Recognizer	deterministic finite automaton	push-down automaton
Perform Work	Actions on transitions in automaton	

The Front End



Why separate the scanner and the parser?

- Scanner classifies words
- Parser constructs grammatical derivations
- Parsing is harder and slower

Separation simplifies the implementation

- Scanners are simple
- Scanner leads to a faster, smaller parser

Scanner is only pass that touches every character of the input.

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* → *expr op term*
3. | *term*
4. *term* → number
5. | id
6. *op* → +
7. | -

$S = \text{goal}$

$T = \{ \text{number}, \text{id}, +, - \}$

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

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

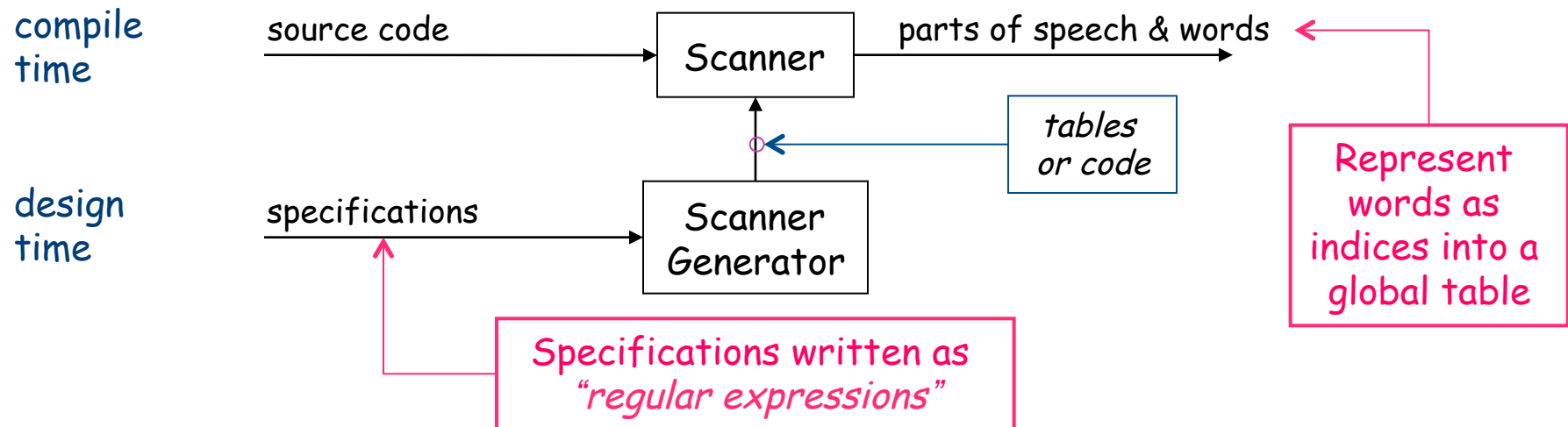
parts of speech
syntactic variables

The Big Picture

In practice, many scanners are hand coded. Even if you build a hand-coded scanner, it is useful to write down the specification and the automaton.

Why study automatic scanner construction?

- Avoid writing scanners by hand
- 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 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 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*



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 $\{\underline{a}\}$
 - the spelling of any specific 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

(a|b|c | ... |z|A|B|C | ... |Z) (a|b|c | ... |z|A|B|C | ... |Z) | (0|1|2 | ... |9)*

Numbers:

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

Decimal → *Integer* . *Digit**

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

Complex → (*Real* , *Real*)

Numbers can get much more complicated!

Using symbolic names
does not imply recursion

underlining indicates
a letter in the input
stream

Regular Expressions

So what's 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)

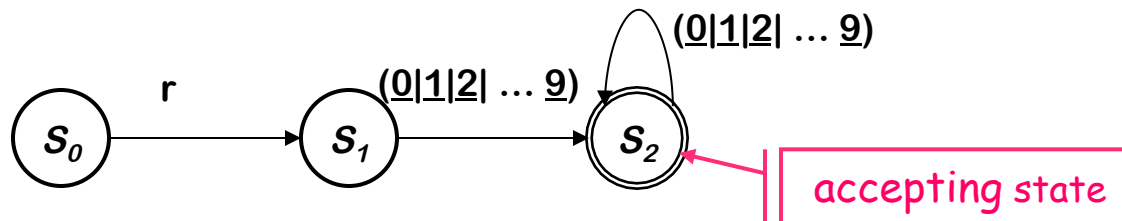


Consider the problem of recognizing ILOC register names

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

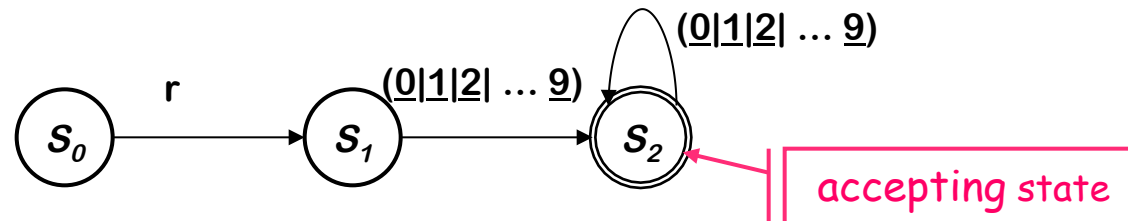
Example

(continued)



DFA operation

- Start in state S_0 & make transitions on each input character
- DFA accepts a word 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, 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	r	0,1,2,3,4, 5,6,7,8,9	All others
α			
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* → r (0|1|2) (*Digit* | ϵ) | (4|5|6|7|8|9) | (3|30|31))
- *Register* → 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

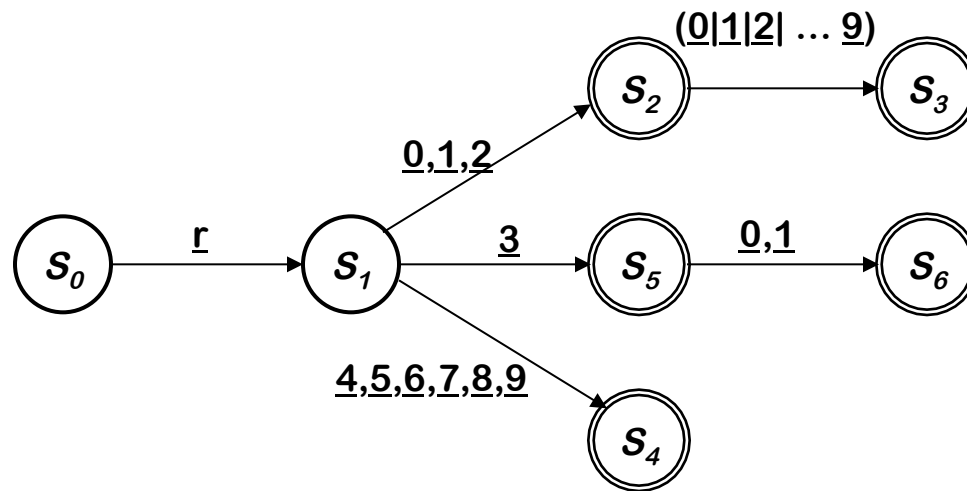
Tighter register specification

(continued)



The DFA for

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



- Accepts a more constrained set of register names
- Same set of actions, more states

Tighter register specification

(continued)



δ	r	0,1	2	3	4-9	All others
s_0	s_1	s_e	s_e	s_e	s_e	s_e
s_1	s_e	s_2	s_2	s_5	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_6	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

This table runs in the same skeleton recognizer

Table encoding RE for the tighter register specification

Tighter register specification

(continued)



State Action	r	0,1	2	3	4,5,6 7,8,9	other
0	1 <i>start</i>	e	e	e	e	e
1	e	2 <i>add</i>	2 <i>add</i>	5 <i>add</i>	4 <i>add</i>	e
2	e	3 <i>add</i>	3 <i>add</i>	3 <i>add</i>	3 <i>add</i>	e <i>exit</i>
3,4	e	e	e	e	e	e <i>exit</i>
5	e	6 <i>add</i>	e	e	e	e <i>exit</i>
6	e	e	e	e	e	x <i>exit</i>
e	e	e	e	e	e	e

