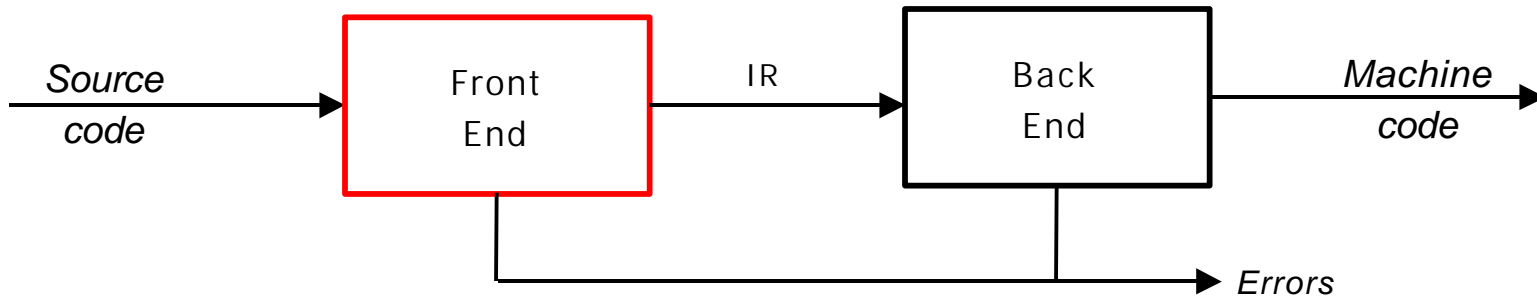


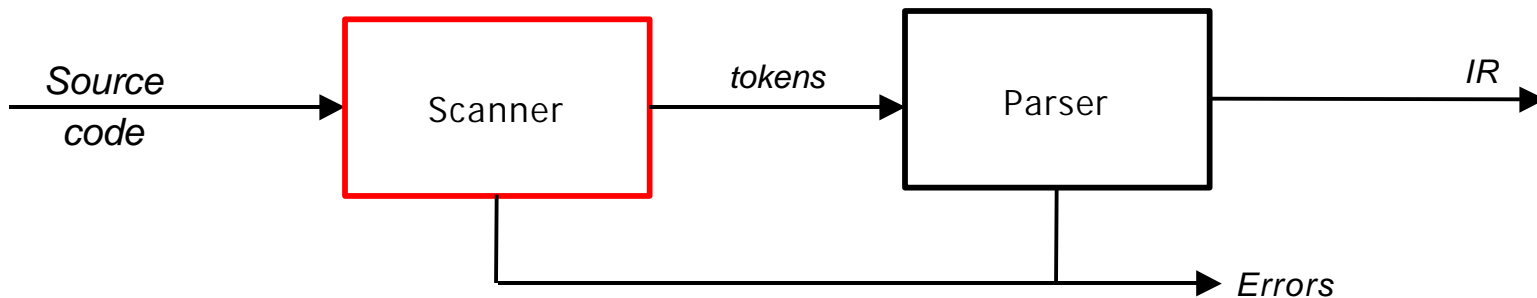
The Front End



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 (syntactically) ?
- Build an IR version of the code for the rest of the compiler

The Front End

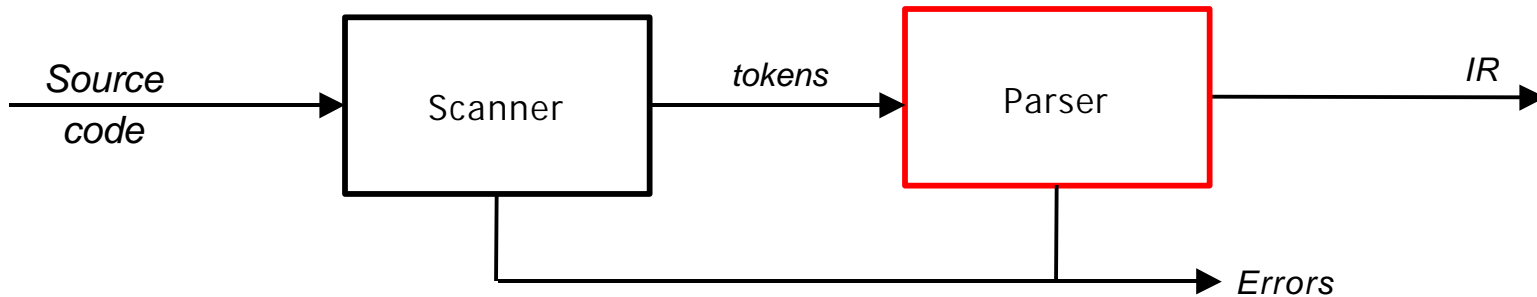


Scanner

- Maps stream of characters into words
 - > Basic unit of syntax
 - > `x = x + y ;` becomes `<id,x> <assignop,=> <id,x> <arithop,+> <id,y> ;`
- Characters that form a word are its *lexeme*
- Its *part of speech* (or *syntactic category*) is called its *token*
- 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

In natural languages, *word* @ *part of speech* is idiosyncratic

- > Based on connotation & context
- > Typically done with a table lookup

In formal languages, *word* @ *part of speech* is syntactic

- > Based on denotation
- > Makes this a matter of syntax, or *micro-syntax*
- > We can recognize this micro-syntax efficiently
- > Reserved keywords are critical (no context!)

⇒ Fast recognizers can map *words* into their *parts of speech*

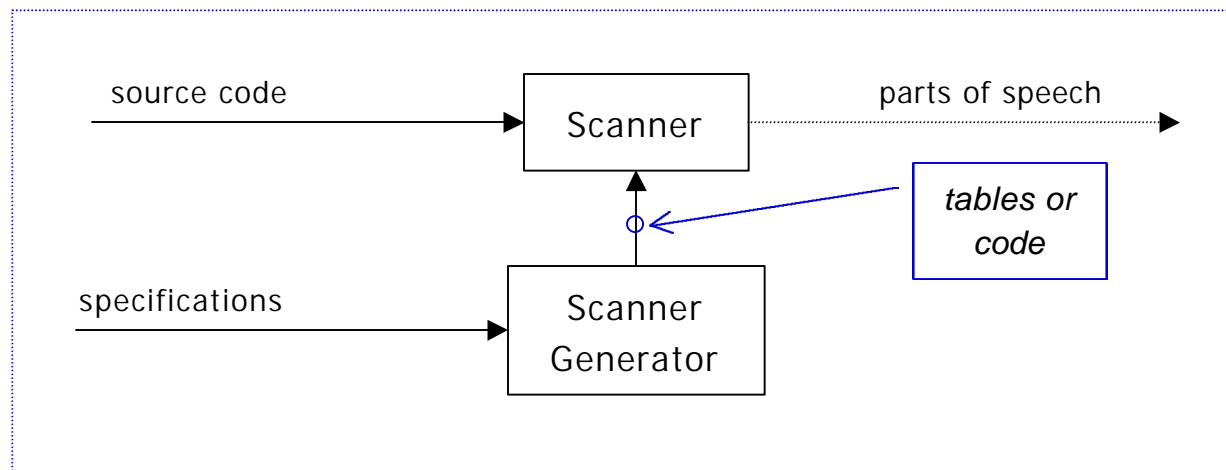
⇒ Study formalisms to automate construction of recognizers



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

Specifying Lexical Patterns

(micro-syntax)



A scanner recognizes the language's parts of speech

Some parts are easy

- White space
 - > *WhiteSpace* → blank | tab | *WhiteSpace* blank | *WhiteSpace* tab
- Keywords and operators
 - > Specified as literal patterns: if, then, else, while, =, +, ...
- Comments
 - > Opening and (*perhaps*) closing delimiters
 - > /* followed by */ in C
 - > // in C++
 - > % in LaTeX

Specifying Lexical Patterns

(micro-syntax)



A scanner recognizes the language's parts of speech

Some parts are more complex

- Identifiers
 - > Alphabetic followed by alphanumerics + `_`, `&`, `$`, ...
 - > May have limited length
- Numbers
 - > Integers: 0 *or* a digit from 1-9 followed by digits from 0-9
 - > Decimals: integer `.` digits from 0-9, *or* `.` digits from 0-9
 - > Reals: (integer or decimal) `E` (`+` *or* `-`) digits from 0-9
 - > Complex: (`(` real `.` real `)`

We need a notation for specifying these patterns

We would like the notation to lead to an implementation

Regular Expressions



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 is a RE denoting $L(x)$
 - > $x|y$ is a RE denoting $L(x) \dot{\cup} L(y)$
 - > xy is a RE denoting $L(x)L(y)$
 - > x^* is a RE denoting $L(x)^*$

Precedence is
closure, then
concatenation, then
alternation

Set Operations

(refresher)



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

You need to know these definitions



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*)^{*}

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)



To make scanning tractable, programming languages differentiate between parts of speech by controlling their spelling (as opposed to dictionary lookup)

Difference between *Identifier* and *Keyword* is entirely lexical

- > While is a *Keyword*
- > Whilst is an *Identifier*

The lexical patterns used in programming languages are regular

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

⇒ We study REs to automate scanner construction !



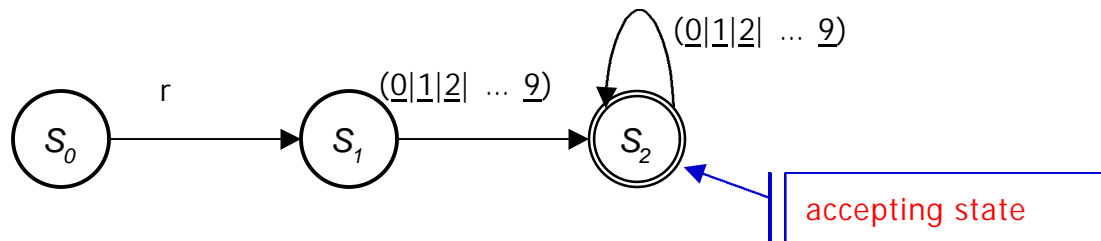
Example

Consider the problem of recognizing register names

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

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

RE corresponds to a recognizer (or DFA)



Recognizer for *Register*

With implicit transitions on other inputs to an error state, s_e

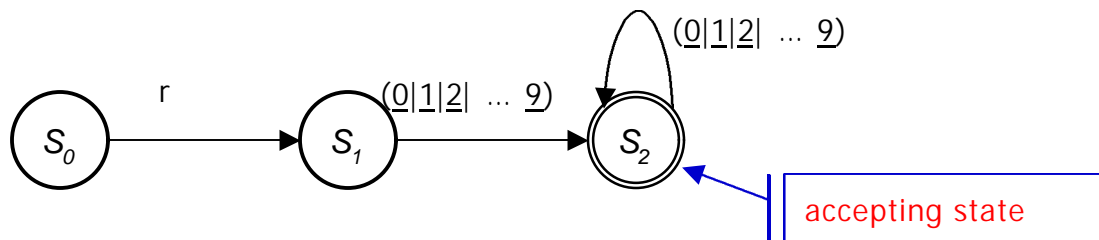


Example

(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)

```

char  $\rightarrow$  next character;
state  $\rightarrow$   $s_0$ ;
call action(state,char);
while (char  $\neq$  eof)
    state  $\rightarrow$   $\delta$ (state,char);
    call action(state,char);
    char  $\rightarrow$  next character;
  
```

```

if  $T$ (state) = final then
    report acceptance;
else
    report failure;
  
```

```

action(state,char)
switch( $T$ (state) )
  case start:
    word  $\rightarrow$  char;
    break;
  case normal:
    word  $\rightarrow$  word + char;
    break;
  case final:
    word  $\rightarrow$  char;
    break;
  case error:
    report error;
    break;
end;
  
```

T	<i>a c t i o n</i>
S_0	<u>s t a r t</u>
S_1	<u>n o r m a l</u>
S_2	<u>f i n a l</u>
S_e	<u>e r r o r</u>

δ	r	0,1, 2,3, 4,5, 6, 7,8, 9	<i>o t h e r</i>
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

- The recognizer translates directly into code
- To change DFAs, just change the tables



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* \textcircled{R} $\underline{r} ((0|1|2) (Digit | \epsilon) | (4|5|6|7|8|9) | (3|30|31)$
- > *Register* \textcircled{R} r0|r1|r2 | ... |r31|r00|r01|r02 | ... |r09

Produces a more complex DFA

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

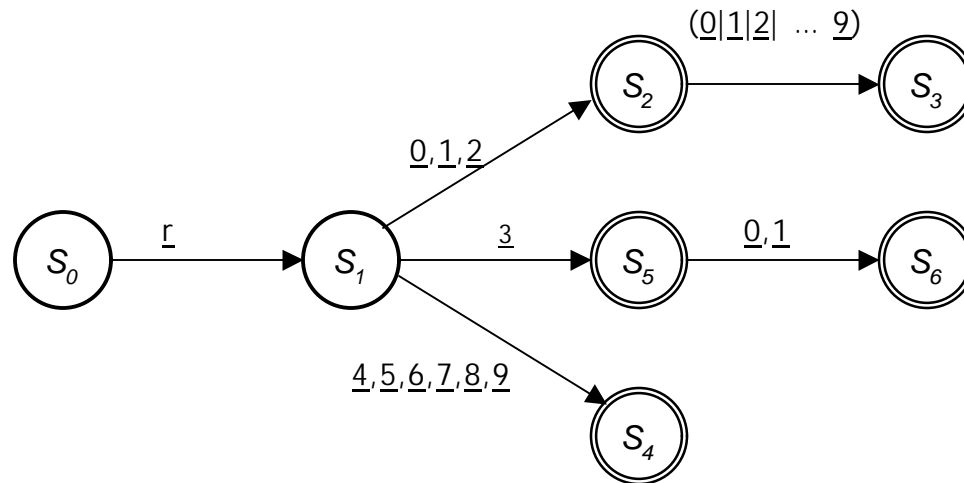
Tighter register specification

(continued)



The DFA for

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



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



Tighter register specification

(continued)

To implement the recognizer

- Use the same code skeleton
- Use transition and action tables for the new RE

δ	r	0,1	2	3	4,5,6 7,8,9	other
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

T	action
S_0	<u>start</u>
S_1	<u>normal</u>
$S_{2,3,4,5,6}$	<u>final</u>
S_e	<u>error</u>

- Bigger tables, more space, same asymptotic costs
- Better (micro-)syntax checking at the same cost