

Scanners and Parsers

Summary

- What is a programming language
- What makes a language successful
- Why study programming languages
- Programming language views
- Types of languages
- Phases of compilation

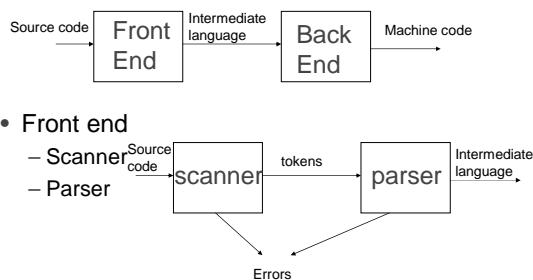
Types of languages

- Imperative
 - Von Neumann (Fortran, Pascal, BASIC, C, ...)
 - Object-oriented (C++/Java, Smalltalk, ...)
- Declarative
 - Functional (Scheme, Lisp, ML, ...)
 - Logic, constraint-based (Prolog, VisiCalc, ...)

Phases of Compilation

- Scanning (lexical analysis)
- Parsing (syntax analysis)
- Semantic analysis
- Intermediate code generation
- Machine-independent code improvement
- Machine-specific code improvement
- Target code generation

Overview



Specifying Patterns

A scanner must recognize various parts of the language's syntax

Some parts are easy –

e.g., white space, keywords, operators, comments

Other parts require more notation –

e.g., identifiers, numbers

→ Need a powerful notation to specify these patterns

Regular Expressions

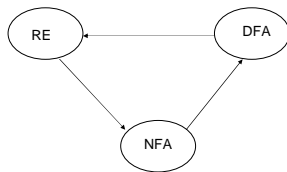
- Defines a set of strings over the characters contained in some alphabet, defining a language
- Three basic operations
 - Union – e.g., $a|b$
 - Concatenation – e.g., ab
 - Closure – (Kleene closure) – e.g., a^* - where a can be a set concatenated with itself zero or more times

Finite Automata

- A finite collection of states
- A set of transitions between those states
- An alphabet
- A start state, S_0
- One or more final states

Deterministic vs. Non-deterministic – single-valued vs. multi-valued transitions, no epsilon transitions

RE \leftrightarrow DFA



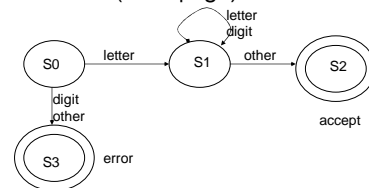
Recognizers

Identifier

letter $\rightarrow (a|b|c| \dots |z|A|B|C| \dots |Z)$

digit $\rightarrow (0|1|2|3|4|5|6|7|8|9)$

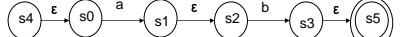
id $\rightarrow \text{letter} (\text{letter}|\text{digit})^*$



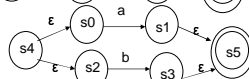
Automated RE \rightarrow NFA

- Build NFA for each term, connect them with ϵ moves

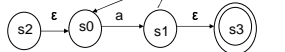
– Concatenate - ab



– Union - $a|b$



– Kleene Closure - a^*



Thompson's Construction

- Each NFA has a single start state and a single final state
- An ϵ -move always connects two states that were start or final states
- A state has at most 2 entering and 2 exiting ϵ -moves, and at most 1 entering and 1 exiting move on a symbol in the alphabet

Automated DFA->RE

```

for i = 1 to N
  for j = 1 to N
     $R_{ij}^0 = \{a | \delta(s_i, a) = s_j\}$ 
    if (i == j)
       $R_{ij}^0 = R_{ij}^0 \cup \epsilon$ 
  for k = 1 to N
    for i = 1 to N
      for j = 1 to N
         $R_{ij}^k = R_{ik}^{k-1} R_{kk}^{k-1} R_{kj}^{k-1} \cup R_{ij}^{k-1}$ 
L =  $\bigcup_{i,j \in S_F} R_{ij}^N$ 

```

NFA->DFA

- Subset construction algorithm
 - Each state in DFA corresponds to a set of states in NFA
- ```

q0 ← ε-closure(n0)
initialize Q with {q0}
while (Q is still changing)
 for each qi ∈ Q
 for each character α ∈ Σ
 t ← ε-closure(move(qi, α))
 T[qi, α] ← t
 if t ∉ Q then
 add t to Q

```

## DFA Minimization

```

P ← P {S_F, {S - S_F}}
while (P is still changing)
 T ← ∅
 for each set p ∈ P
 for each α ∈ Σ
 partition p by α
 into p1, p2, p3, ... pk
 T ← T ∪ p1, p2, p3, ... pk
 if T ≠ P then
 P ← T

```

## Implementing Scanners

- Ad-hoc
- Semi-mechanical pure DFA
- Table-driven DFA

## Code for Recognizer/Scanner

```

char = next_char();
state = S0; /* code for S0 */
token_value = "" /* empty string */
while (not done) {
 class = char_class[char];
 state = next_state[class, state];
 switch(state) {
 case S1: /* building an id */
 token_value = token_value + char;
 char = next_char; break;
 case S2: /* accept state */
 token_type = identifier;
 done = true; break;
 case S3: /* error */
 token_type = error;
 done = true; break;
 }
}
return(token_type);

```

## Tables Driving the Recognizer

| char_class | a-z    | A-Z    | 0-9   | other |
|------------|--------|--------|-------|-------|
| value      | letter | letter | digit | other |

|            |        |    |    |    |    |
|------------|--------|----|----|----|----|
|            | class  | S0 | S1 | S2 | S3 |
| next_state | letter | S1 | S1 | -- | -- |
|            | digit  | S3 | S1 | -- | -- |
|            | other  | S3 | S2 | -- | -- |

To change language, we can just change tables

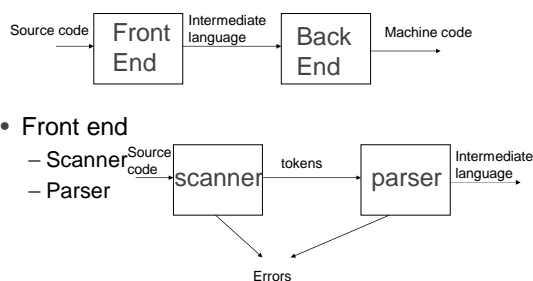
## Error Recovery

- E.g., illegal character
- What should the scanner do?
  - Report the error
  - Try to correct it?
- Error correction techniques
  - Minimum distance corrections
  - Hard token recovery
  - Skip until match

## Scanner Summary

- Break up input into tokens
- Catch lexical errors
- Difficulty affected by language design
- Issues
  - Input buffering
  - Lookahead
  - Error recovery
- Scanner generators
  - Tokens specified by regular expressions
  - Regular expression  $\rightarrow$  DFA
  - Highly efficient in practice

## Overview



## Parsers

- Check input for grammatical correctness
- Determine syntax of token stream
- Constructs intermediate representation
- Produces meaningful error messages

Needs mathematical model of syntax  $\rightarrow$  grammar  
 Needs algorithm for testing syntax  $\rightarrow$  parsing

## Context Free Grammars

- Adds recursion/allows non-terminals to be expressed in terms of themselves
- Can be used to count/impart structure – e.g., nested parentheses
- Notation – grammar  $G(S, N, T, P)$ 
  - $S$  is the start symbol
  - $N$  is a set of non-terminal symbols (LHS)
  - $T$  is a set of terminal symbols (tokens)
  - $P$  is a set of productions or rewrite rules ( $P : N \rightarrow N \cup T$ )

## Advantages of CFGs

- Precise syntactic specification of a programming language
- Easy to understand, avoids ad hoc definition
- Easier to maintain, add new language features
- Can automatically construct efficient parser
- Parser construction reveals ambiguity, other difficulties
- Imparts structure to language
- Supports syntax-directed translation

## Derivations

- A sequence of application of the rewrite rules is a derivation or a parse (e.g., deriving the string  $x + 2 - y$ )
- The process of discovering a derivation is called parsing

## Types of derivations

- Leftmost derivation
  - The leftmost non-terminal is replaced at each step
- Rightmost derivation
  - The rightmost non-terminal is replaced at each step
- Ambiguous grammar – one with multiple leftmost (or multiple rightmost) derivations for a single sentential form

## Types of parsers

- Top-down (LL) parsers
  - Left to right, leftmost derivation
  - Starts at the root of the derivation tree and fills in
  - Predicts next state with  $n$  lookahead
- Bottom-up (LR) parsers
  - Left to right, rightmost derivation
  - Starts at the leaves and fills in
  - Start with input string, end with start symbol
  - Starts in a state valid for legal first tokens
  - Change state to encode possibilities as input is consumed
  - Use a stack to store both state and sentential form

## Top-Down Parsing

- At a node labeled  $A$ , select a production with  $A$  on its LHS and for each symbol on its RHS, construct the appropriate child
- When a terminal is added that does not match the input, backtrack
- Find the next node to be expanded (must have a label in NT)

## Calculator Example

- All variables are integers
- There are no declarations
- The only statements are assignments, input, and output
- Expressions use one of four arithmetic operators and parentheses
- Operators are left associative, with the usual precedence
- There are no unary operators

## Regular Expressions

$\text{id} \rightarrow \text{letter} ( \text{letter} | \text{digit} )^*$   
 $\text{literal} \rightarrow \text{digit digit}^*$   
read, write, ":", "+", "-", "\*", "/", "(", ")"  
\$\$ /\* end of input \*/

## Grammar for Calculator

```

<pgm> → <stmtlist> $$
<stmtlist> → <stmtlist> <stmt> | ε
<stmt> → id := <expr> | read <id> | write <expr>
<expr> → <term> | <expr> <add op> <term>
<term> → <factor> | <term> <multop> <factor>
<factor> → (<expr>) | id | literal
<addop> → + | -
<multop> → * | /

```

## Eliminating Left Recursion

- A grammar is left recursive if there exists A in NT such that  $A \rightarrow A\delta$  for some string  $\delta$
- Transform the grammar to remove left recursion

```

<foo> → <foo> δ
 | μ

```

→

```

<foo> → μ <bar>
<bar> → δ <bar> | ε

```

where <bar> is a new non-terminal

## Eliminating Common Prefixes

```

foo → bar δ
 → bar (μ)

```

→

```

foo → bar footail
footail → δ | (μ)

```

## LL Grammar for Calculator

```

<pgm> → <stmtlist> $$
<stmtlist> → <stmt> <stmtlist> | ε
<stmt> → id := <expr> | read <id> | write <expr>
<expr> → <term> <termtail>
<termtail> → <addop> <term> <termtail> | ε
<term> → <factor> <factortail>
<factortail> → <multop> <factor> <factortail> | ε
<factor> → (<expr>) | id | literal
<addop> → + | -
<multop> → * | /

```

## Example non-LL Grammar Construct

```

stmt → if condition then_clause else_clause
 | other_stmt
then_clause → then stmt
else_clause → else stmt | ε
→ Ambiguous – allows dangling else to be
 paired with either then in
 if A then if B then C else D

```

## Fix – Bottom-up parsing (LR)

```

stmt → balanced_stmt | unbalanced_stmt
balanced_stmt → if condition then balanced_stmt
 | else balanced_stmt
 | other_stmt
unbalanced_stmt → if condition then stmt
 | if condition then balanced_stmt
 | else unbalanced_stmt

```

OR

Use special disambiguating rules → use production that occurs first in case of conflict

## Parser Construction

- Recursive descent parsing
  - Top-down parsing algorithm
  - Built on procedure calls (may be recursive)
  - Write procedure for each non-terminal, turning each production into clause
  - Insert call to procedure  $A()$  for non-terminal  $A$  and to  $\text{match}(x)$  for terminal  $x$
  - Start by invoking procedure for start symbol  $S$

## Predictive (Table-Driven) Parsing

- Actions –
  - Match a terminal
  - Predict a production
  - Announce a syntax error
- Push as yet unseen portions of productions onto a stack
- Use –
  - $\text{FIRST}(A)$
  - $\text{FOLLOW}(A)$

## The FIRST Set

- $\text{FIRST}(\alpha)$  is the set of terminal symbols that begin strings derived from  $\alpha$
- To build  $\text{FIRST}(X)$ :
  - If  $X$  is a terminal,  $\text{FIRST}(X)$  is  $\{X\}$
  - If  $X \leftarrow \epsilon$ , then  $\epsilon \in \text{FIRST}(X)$
  - If  $X \leftarrow Y_1 Y_2 \dots Y_k$  then put  $\text{FIRST}(Y_1)$  in  $\text{FIRST}(X)$
  - If  $X$  is a non-terminal and  $X \leftarrow Y_1 Y_2 \dots Y_k$ , then a  $a \in \text{FIRST}(X)$  if  $a \in \text{FIRST}(Y_i)$  and  $\epsilon \in \text{FIRST}(Y_j)$ , for all  $1 \leq j < i$

## The Follow Set

- For a non-terminal  $A$ ,  $\text{FOLLOW}(A)$  is the set of terminals that can appear immediately to the right of  $A$  in some sentential form
- To build  $\text{FOLLOW}(B)$  for all  $B$  -
  - Starting with goal, place eof in  $\text{FOLLOW}(<\text{goal}>)$
  - If  $A \leftarrow \alpha B \beta$ , then put  $\{\text{FIRST}(\beta) - \epsilon\}$  in  $\text{FOLLOW}(B)$
  - If  $A \leftarrow \alpha B$ , then put  $\text{FOLLOW}(A)$  in  $\text{FOLLOW}(B)$
  - If  $A \leftarrow \alpha B \beta$  and  $\epsilon \in \text{FIRST}(\beta)$ , then put  $\text{FOLLOW}(A)$  in  $\text{FOLLOW}(B)$

## Using FIRST and FOLLOW

- For each production  $A \leftarrow \alpha$  and lookahead token
  - Expand  $A$  using the production if token  $\epsilon \in \text{FIRST}(\alpha)$
  - If  $\epsilon \in \text{FIRST}(\alpha)$ , expand  $A$  using the production if token  $\epsilon \in \text{FOLLOW}(A)$
  - All other tokens return error
- If there are multiple choices, the grammar is not LL(1) (predictive)

## LL(1) Grammars

A Grammar  $G$  is LL(1) if and only if, for all non-terminals  $A$ , each distinct pair of productions  $A \leftarrow \alpha$  and  $A \leftarrow \beta$  satisfy the condition  $\text{FIRST}(\alpha) \cap \text{FIRST}(\beta) = \Phi$ , i.e.,

For each set of productions  $A \leftarrow \alpha_1 | \alpha_2 | \dots | \alpha_n$

- $\text{FIRST}(\alpha_1), \text{FIRST}(\alpha_2), \dots, \text{FIRST}(\alpha_n)$  are all pairwise disjoint
- If  $\alpha_i \leftarrow \epsilon$  for any  $i$ , then  $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \Phi$ , for all  $j \neq i$

## Bottom-Up Parsers

- End with start symbol
- Apply productions in reverse to input, replacing RHS of production with LHS non-terminal
- Final result is a rightmost derivation, in reverse
- Bottom-up parsers can use stack and lookahead to choose production
- LR(k) parsers see the entire RHS before choosing a production → more powerful than LL(k)

### Definitions –

The *handle* is defined as the combination of

- 1) RHS to be replaced
- 2) Its position

Replacement step is called a *reduction*

## Shift-Reduce Parsers

### • Actions

- Shift – next input symbol shifted onto stack
- Reduce – right end of handle is on top of stack; locate left end of handle within stack; pop handle off stack and push appropriate non-terminal LHS
- Accept – terminate parsing and signal success
- Error – initiate/call an error recovery routine

## Skeleton Parser

```
push s0
token = next_token()
while (1)
 s = top of stack
 if action[s,token] = "shift si" then
 push token
 push si
 token = next_token()
 else if action[s,token] = "reduce A ←" then
 pop 2 * |β| symbols
 s = top of stack
 push A
 push goto[s,A]
 else if action[s,token] = "accept" then
 return
 else error()
```

Results in  $k$  shifts,  $r$  reduces, and 1 accept for a valid string, where  $k$  is the length of the input string and  $r$  is the length of the reverse rightmost derivation

## LR(1) Items

- An LR(k) item is a pair  $[\alpha, \beta]$ , where  $\alpha$  is a production from  $G$  with a  $\bullet$  at some position in the rhs and  $\beta$  is a lookahead string containing  $k$  symbols (terminals or eof)
- The  $\bullet$  indicates how much of an item we have seen at a given state in the parse
- $[A \leftarrow XY \bullet Z, a]$  indicates that the parser has seen a string derived from  $XY$  and is looking for one derivable from  $Z$ , and that  $a$  is in  $\text{FOLLOW}(A)$
- $A \leftarrow XYZ$  generates 4 LR(0) items

## LR(1) Machine

- Definitions
  - Closure of  $[A \leftarrow \alpha \bullet B \beta, a]$  contains itself and any items of form  $[B \leftarrow \bullet \gamma, \text{FIRST}(\beta a)]$ , repeating for new items
  - $\text{goto}(X)$  of  $[A \leftarrow \alpha \bullet X \beta, a]$  contains the closure of  $[A \leftarrow \alpha X \bullet \beta, a]$
- LR(1) DFA construction
  - Begin with closure of start symbol  $[S \leftarrow \bullet \alpha, \text{eof}]$
  - For each state, calculate  $\text{goto}(X)$  for all grammar symbols  $X$ , generating states
    - Each state is a set of LR(1) items
  - Repeat above step for newly generated states

## LR(1) Table Construction

- If  $S$  appears on RHS, create augmented grammar  $G'$  by adding  $S' \leftarrow S$
- Construct the collection of sets of LR(1) items for  $G'$
- For each state  $S_i$  and each item in state  $S_i$ 
  - If  $[A \leftarrow \alpha \bullet a \beta, b]$  in  $S_i$  and  $\text{goto}[S_i] = S_j$ , then set  $\text{action}[i, a]$  to "shift  $j$ " ( $a$  must be a terminal)
  - If  $[A \leftarrow \alpha \bullet a]$  in  $S_i$ , then set  $\text{action}[i, a]$  to "reduce  $A \leftarrow \alpha$ "
  - If  $[S' \leftarrow S \bullet, \text{eof}]$  in  $S_i$ , then set  $\text{action}[i, \text{eof}]$  to "accept"
- For all non-terminals
  - if  $\text{goto}(S_i, A) = S_j$ , then set  $\text{goto}[i, A]$  to  $j$
- All other entries in action and goto are set to "error"
- The initial state of the parser is the state constructed from the set containing the item  $[S' \leftarrow \bullet S, \text{eof}]$



## Error Recovery

- Phrase-level
  - Delete tokens until something likely to follow expression (in FIRST or FOLLOW for expression or eof)
  - Context-specific FOLLOW sets
  - Avoid deleting statically defined “starter” symbols
- Exception-based
  - Small set of contexts to which to back out in case of error
  - Re-raise exception to pop to next handler on seeing starter symbol that it shouldn't delete
  - Delete input tokens until parsing can recommence

## Error Recovery Example

```
stmt_list ← stmt
 | stmt_list; stmt
```

*can be augmented with error*

```
stmt_list ← stmt
 | error
 | stmt_list; stmt
```

- Throw out the erroneous statement, synchronize at “;” or eof, invoke `yyerror(“syntax error”)`
- Other natural places for errors – lists, missing parentheses or brackets, extra operator or missing operator