

CSE401: Lexical Analysis

Larry Snyder
Spring 2003

Slides by Chambers, Eggers, Notkin, Ruzzo, Snyder and others

Objectives (today and tomorrow)

- Define overall theory and practical structure of lexical analysis
- Briefly recap regular languages, expressions, finite state machines, and their relationships
- How to define tokens with regular expressions
- How to leverage this to implement a lexer

2

Lexical analysis (scanning)

```

graph LR
    A[Source Program] --> B[stream of characters]
    B -- Lexical analysis --> C[sequence of tokens]
  
```

- The scanner/lexer groups characters into tokens
- A *token* is a basic, atomic chunk of syntax, e.g.
 - Literals: 17, 42, 3.1415, "Hello.", ...
 - Punctuation & operators: },),], :=, <, <=, ...
 - Reserved words: if, then, else, for, while, int, char, ...
 - Identifiers: snork, x, dogbert, sqrt, printf, ...
- The lexer also removes whitespace
 - Whitespace: characters that are ignored between tokens
 - Ex: spaces, tabs, newlines, comments
- Definitions of tokens and whitespace vary among languages

3

Separation of lexing & parsing

- A universal separation:
 - Lexer: character stream to token stream
 - Parser: token stream to syntax tree
- Advantages:
 - Simpler design
 - Based on related but distinct theoretical underpinnings
 - Compartmentalizes some low-level issues, e.g., I/O, internationalization, ...
 - Faster
 - Lexing is time-consuming in many compilers (40-60% ?)
 - By restricting the job of the lexer, a faster implementation is usually feasible

4

Overall approach to scanning

- Define language tokens using regular expressions
 - Natural representation for tokens
 - But difficult to produce a scanner from REs
- Convert the regular expressions into a non-deterministic finite state automaton (NFA)
 - Straightforward conversion
 - Can produce a scanner from NFA, but an inefficient one
- Convert the NFA into a deterministic finite state automaton (DFA)
 - Straightforward conversion
- Convert the DFA into an efficient scanner implementation

5

Language & automata theory: a speedy reminder

- Alphabet: a finite set of symbols
- String: a finite, possibly empty, sequence of symbols from an alphabet
- Language: a set, often infinite, of strings
- Finite specifications of (possibly infinite) languages:
 - Automaton – a recognizer; a machine that accepts all strings in the language (and rejects all other strings)
 - Grammar – a generator; a system for producing all strings in the language (and no other strings)
- A language may be specified by many different grammars and automata
- A grammar or automaton specifies only one language

6

Definitions: token vs lexeme

- Token: an "atom of syntax"; set of lexemes
 - Ex: int literal, string literal, identifier, keyword-if
- Lexeme: the character string forming a token
 - Ex: 17, 42, "Hello", "Goodbye", x, dogbert, if
- A token may have attributes, if the set has more than a single lexeme
 - "int literal" token might have attribute "17" or "42"
 - "keyword-if" token probably needs no attributes

7

Regular expressions: a notation for defining tokens

- Regular expressions (REs) are defined inductively:
 - Base cases
 - The empty string (ϵ)
 - A symbol from the alphabet
 - Inductive cases
 - Choice of two REs: $E_1 | E_2$
 - Sequence of two REs: $E_1 E_2$
 - Kleene closure (zero or more occurrences) of an RE: E^*
- Use parentheses for grouping
- Whitespace is not significant

Increasing precedence

8

Examples

- a
- a b
- (a | b)
- (a | b) c
- a | b c
- a b*
- (a | b)(0 | 1)*

9

Notational conveniences: no additional expressive power

- E^+ means one or more occurrences of E
- E^k means k occurrences of E (k a literal constant)
- $[E]$ means 0 or 1 occurrences of E (it's optional)
- $\{E\}$ means E^+
- $\text{not}(x)$ means any character in the alphabet but x

rarely implemented (potentially expensive)

- $\text{not}(E)$ means any strings in the alphabet but those matching E
- $E_1 - E_2$ means any strings matching E_1 except those matching E_2

10

Naming regular expressions: simplify RE definitions

- Can assign names to regular expressions
- Can use these names in the definition of another regular expression
- Examples
 - letter ::= a | b | ... | z
 - digit ::= 0 | 1 | ... | 9
 - alphanumeric ::= letter | digit
- Can eliminate names by macro expansion
- No recursive definitions are allowed! Why?

11

Regular expressions for PL/0

```

Digit ::= 0 | ... | 9
Letter ::= a | ... | z | A | ... | Z
Integer ::= Digit*
AlphaNum ::= Letter | Digit
Id ::= Letter AlphaNum*
Keyword ::= module | procedure | begin | end | const
           | var | if | then | while | do | input
           | output | odd | int
Punct ::= ; | : | . | , | ( | )
Operator ::= = | * | / | + | - | = | < > | <= | < | >= | >
Token ::= Id | Integer | Keyword | Operator | Punct
White ::= <space> | <tab> | <newline>
Program ::= (Token | White)*
  
```

12

Generate scanner from regular expressions?

- n This would be ideal: REs as input to a scanner generator, and a scanner as output
 - n Indeed, some tools can mostly do this
- n But it's not straightforward to do this
 - n One reason: there is a lot of non-determinism — choice — inherent in most regular expressions
 - n Choice can be implemented using backtracking, but it's generally very inefficient
- n In any case, these tools go through a process like the one we'll look at

13

Next steps

- n Convert regular expressions to non-deterministic finite state automata (NFA)
- n Then convert the NFA to deterministic finite state automata (DFA)
- n Then convert DFA into code

14

Finite state automaton

- n A finite set of states
 - n One marked as the initial state
 - n One or more marked as final states
- n A set of transitions from state to state
 - n Each transition is marked with a symbol from the alphabet or with ϵ
- n Operate by reading symbols in sequence
 - n A transition can be taken if it labeled with the current symbol
 - n An ϵ -transition can be taken at any point, without consuming a symbol
- n Accept if no more input and in a final state
- n Reject if no transition can be taken or if no more input and not in a final state (DFA case)

15

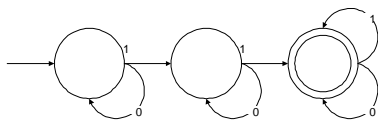
DFA vs. NFA

- n A deterministic finite state automaton (DFA) is one in which there is no choice of which transition to take under any condition
- n A non-deterministic finite state automaton (NFA) is one in which there is a choice of which transition to take in at least one situation
 - n "Accept" == some way } to reach final state
 - n "Reject" == all ways fail } at end of input

16

Example

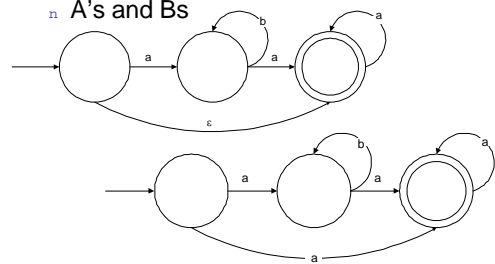
- n Describe the language



17

Examples

- n A's and Bs



18

Plan of attack

- Convert from regular expressions to NFAs because there is an easy construction
 - However, NFAs encode choice, and choice implies backtracking, which is slow
- Convert from NFAs to DFAs, because there is a well-defined procedure
 - And DFAs lay the foundation for an efficient scanner implementation

19

Exercise

- Consider the language that includes only those binary strings that have odd parity
- For this language, define
 - the alphabet
 - a grammar
 - an automaton

20

Converting REs to NFAs: base cases

21

$E_1 E_2$

22

$E_1 \mid E_2$

23

E^*

?

24

RE to NFA

- Those rules are sufficient for constructing an equivalent NFA from a regular expression

25

Exercise

- Define a regular expression that recognizes comments of the form
 - `/* ... */`
 - Be careful in defining "..."
- Then convert that regular expression to an NFA

26

Building lexers from regular expressions

- Convert the regular expressions into deterministic finite state automata (DFA)
 - Manually
 - Mechanically by converting first to non-deterministic finite state automata (NFA) and then into DFA
- Convert DFA into scanner implementation
 - By hand into a collection of procedures
 - Mechanically into a table-driven parser

27

Why convert to DFAs?

- Because
 - they are equivalent in power to NFAs
 - they are deterministic, which makes them a terrific basis for an efficient implementation of a scanner

28

NFA => DFA

- Basic problem
 - NFA can choose among alternative paths
 - either ϵ transitions or
 - multiple transitions from a state with the same label
 - But a DFA cannot have this kind of choice
- Solution: subset construction
 - In the newly constructed DFA, each state represents a *set* of states in the NFA,
- Key Idea:
 - the* state of the DFA after reading $x_1 x_2 \dots x_k$ is the set of *all* states that the NFA might reach after reading the same input

29

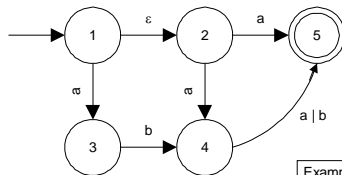
Subset construction algorithm

initial step

- Create start state of new DFA
 - Label it with the set of NFA states that can be reached without consuming any input
 - I.e., NFA's start state, or reachable by ϵ transitions
 - Think of it as all possible start states in the NFA, since there could be more than one, given the ϵ transitions
- Then "process" this new start state
 - Details below

30

Example



Example from
Crafting a Compiler,
Fischer & LeBlanc

31

Example (cont.)

32

Subset construction algorithm *processing a state*

- To process a state S in the new DFA with label $\{s_1, \dots, s_n\}$
- For each symbol x in the alphabet
 - Compute the set T of NFA states reached from any of the NFA states s_1, \dots, s_n by *one* x transition followed by *any number* of ϵ transitions
 - If T is not empty
 - If there is not already a DFA state with T as a label, create one, and add T to the list of states to be processed
 - Add a transition labeled x from S to T
- Repeat until no unprocessed states

33

Subset construction algorithm *defining final states*

- After the algorithm terminates
- Mark every DFA state as final if *any* of the NFA states in its label is final

34

Subset construction: notes

- It is provable that this works and produces an equivalent DFA (c.f. CSE 322)
- This activity can be automated
- Question: What can be said about the number of states in the DFA relative to the NFA?
 - In theory? In practice?

35

Minimizing DFAs

- There is also an algorithm for minimizing the number of states in a DFA
- Given an arbitrary DFA, one can find a unique DFA with a minimum number of states that is equivalent to the original DFA
 - Except for a renaming of the states
 - Essentially, try to merge states

36

Constructing scanners from DFAs

- n Use a table-driven scanner
- n Write disciplined procedures that encode the DFA
- n We'll talk about both (the first briefly)
- n The second approach is used in the PL/0 compiler
 - n Because it's generally easier to handle a few practical issues (but may be slower?)

37

Approach 1: Table-driven

- n Represent the DFA as an adjacency matrix
 - n One row per state
 - n One column per character in the alphabet
 - n Entry is state to transition to
- n Mechanically walk the input, taking appropriate transitions
 - n Rules for termination remain unchanged

	a	b
{1,2}	{3,4,5}	
{3,4,5}	{5}	{4,5}
{4,5}	{5}	{5}
{5}		

38

Approach 2: Procedural

- n Define a procedure for each state in the DFA
- n Use conditionals to check the input character and then make the appropriate transition
- n A transition is a call to the procedure for the next state
- n (Call overhead optimizable)

```

procedure {3,4,5} begin
  if nextChar() == 'a'
    call {5}
  elsif nextChar() == 'b'
    call {4,5}
  else
    reject("no transition
           out of this
           state")
end
    
```

39

The heart of the PL/0 scanner *it's not quite as clean (but it's not bad!)*

```

Token ::= Id |
         Integer |
         Keyword |
         Operator |
         Punct

if (isalpha(CurrentCh)) {
  T = GetIdent()
} else if (isdigit(CurrentCh)) {
  T = GetInt()
} else {
  T = GetPunct();
}
    
```

- n Where's the DFA?
- n How come five kinds of tokens and only three branches?

40

PL/0's GetIdent method

- n Is PL/0 case-sensitive?
- n What does SearchReserved return?

```

Token* Scanner::GetIdent() {
  char ident[MaxIdLength+1];
  int LengthOfId = 0;
  while (isalnum(CurrentCh)) {
    ident[LengthOfId] =
      tolower(CurrentCh);
    LengthOfId++;
    GetCh();
  }
  ident[LengthOfId] = '\0';
  return SearchReserved(ident);
}
    
```

41

PL/0's GetInt method

```

Token* Scanner::GetInt() {
  int integer = 0;
  while (isdigit(CurrentCh)) {
    integer = 10 * integer + (CurrentCh - '0');
    GetCh();
  }
  return new IntegerToken(integer);
}
    
```

42

PL/0's GetPunct method

```

Token* Scanner::GetPunct() {
    Token* T;
    switch (CurrentCh) {
        case ' ':
            GetCh();
            if (CondReadCh('=')) {
                T = new Token(GETS);
            } else {
                T = new Token(COLON);
            }
            break;
        case '<':
            GetCh();
            if (CondReadCh('=')) {
                T = new Token(LEQ);
            } else if (CondReadCh('>')) {
                T = new Token(NEQ);
            } else {
                T = new Token(LSS);
            }
            break;
        ...
    }
}

```

43

A few PL/0 scanner notes

- n There is a Scanner class
 - n There is only one instance of this class
 - n This is an example of the *Singleton* design pattern
- n The high-level structure we showed has the scanner scan before the parser parses
 - n Study the compiler to figure out what really happens
- n Make sure (for this and all other phases) to read the interface (the .h file) very, very carefully

44

Language design issues (lexical)

- n Most languages are now free-form
 - n Layout doesn't matter
 - n Use whitespace to separate tokens, if needed
 - n Alternatives include
 - n Fortran, Algol68: whitespace is ignored
 - n Haskell: use layout to imply grouping
- n Most languages now have reserved words
 - n Cannot be used as identifiers
 - n Alternative: PL/1 has keywords that are treated specially only in certain contexts, but may be used as identifiers, too
- n Most languages separate scanning & parsing
 - n Alternative: C/C++ *type* vs *ident*

```

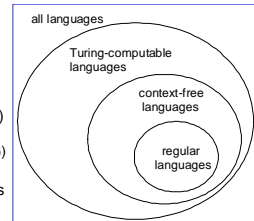
typedef int mytype;
int myvar;
mytype i,j,k;

```

47

Classes of languages

- n Regular languages can be specified by
 - n regular expressions
 - n regular grammars
 - n finite-state automata (FSA)
- n Context-free languages (CFL) can be specified by
 - n context-free grammars (CFG)
 - n push-down automata (PDA)
- n Turing-computable languages can be specified by
 - n arbitrary grammars
 - n Turing machines



Strict inclusion of these classes of languages

46

Objectives: next lectures

- n Understand the theory and practice of parsing
- n Describe the underlying language theory of parsing (CFGs, etc.)
- n Understand and be able to perform top-down parsing
- n Understand bottom-up parsing

47