

Languages, Automata, Regular Expressions & Scanners

CSE 413
Autumn 2007

11/14/2007

1

Agenda

- Basic concepts of formal grammars
- Scanner Theory
 - Regular expressions
 - Finite automata (to recognize regular expressions)
- Scanner Implementation

11/14/2007

2

Programming Language Specs

- Since the 1960s, the syntax of every significant programming language has been specified by a formal grammar
 - First done in 1959 with BNF (Backus-Naur Form or Backus-Normal Form) used to specify the syntax of ALGOL 60
 - Borrowed from the linguistics community (Chomsky)

11/14/2007

3

Grammar for a Tiny Language

- $program ::= statement \mid program\ statement$
- $statement ::= assignStmt \mid ifStmt$
- $assignStmt ::= id = expr ;$
- $ifStmt ::= if (expr) stmt$
- $expr ::= id \mid int \mid expr + expr$
- $id ::= a \mid b \mid c \mid i \mid j \mid k \mid n \mid x \mid y \mid z$
- $int ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

11/14/2007

4

Productions

- The rules of a grammar are called *productions*
- Rules contain
 - Nonterminal symbols: grammar variables (*program*, *statement*, *id*, etc.)
 - Terminal symbols: concrete syntax that appears in programs (a, b, c, 0, 1, if, (, ...)
- Meaning of
 $nonterminal ::= \langle \text{sequence of terminals and nonterminals} \rangle$
In a derivation, an instance of *nonterminal* can be replaced by the sequence of terminals and nonterminals on the right of the production
- Often, there are two or more productions for a single nonterminal – can use either at different times

11/14/2007

5

Alternative Notations

- There are several syntax notations for productions in common use; all mean the same thing

$ifStmt ::= if (expr) stmt$

$ifStmt \rightarrow if (expr) stmt$

$\langle ifStmt \rangle ::= if (\langle expr \rangle) \langle stmt \rangle$

11/14/2007

6

Example Derivation

```

program ::= statement | program statement
statement ::= assignStmt | ifStmt
assignStmt ::= id = expr;
ifStmt ::= if ( expr ) stmt;
expr ::= id | int | expr + expr
Id ::= a | b | c | i | j | k | n | x | y | z
int ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

```

a = 1 ; if (a + 1) b = 2 ;

11/14/2007

Parsing

- Parsing: reconstruct the derivation (syntactic structure) of a program
- In principle, a single recognizer could work directly from the concrete, character-by-character grammar
- In practice this is never done

11/14/2007

Parsing & Scanning

- In real compilers the recognizer is split into two phases
 - **Scanner**: translate input characters to tokens
 - Also, report lexical errors like illegal characters and illegal symbols
 - **Parser**: read token stream and reconstruct the derivation

```

source → [Scanner] → tokens → [Parser]

```

11/14/2007

Characters vs Tokens (review)

- Input text


```
// this statement does very little
if ( x >= y ) y = 42;
```
- Token Stream

IF	LPAREN	ID(x)	GEQ	ID(y)
RPAREN	ID(y)	BECOMES	INT(42)	SCOLON

11/14/2007

Why Separate the Scanner and Parser?

- Simplicity & Separation of Concerns
 - Scanner hides details from parser (comments, whitespace, etc.)
 - Parser is easier to build; has simpler input stream (tokens)
- Efficiency
 - Scanner can use simpler, faster design
 - (But still often consumes a surprising amount of the compiler's total execution time)

11/14/2007

Tokens

- Idea: we want a distinct token kind (lexical class) for each distinct terminal symbol in the programming language
 - Examine the grammar to find these
- Some tokens may have attributes
 - Examples: integer constant token will have the actual integer (17, 42, ...) as an attribute; identifiers will have a string with the actual id

11/14/2007

Typical Tokens in Programming Languages

- Operators & Punctuation
 - + - * / () { } [] ; : :: < <= == != ! ...
 - Each of these is a distinct lexical class
- Keywords
 - if while for goto return switch void ...
 - Each of these is also a distinct lexical class (*not* a string)
- Identifiers
 - A single ID lexical class, but parameterized by actual id
- Integer constants
 - A single INT lexical class, but parameterized by int value
- Other constants, etc.

11/14/2007

13

Principle of Longest Match

- In most languages, the scanner should pick the longest possible string to make up the next token if there is a choice
- Example

```
return foobar != hohum;
```

should be recognized as 5 tokens

RETURN ID(foobar) NEQ ID(hohum) SCOLON

not more (i.e., not parts of words or identifiers, or ! and = as separate tokens)

11/14/2007

14

Formal Languages & Automata Theory (in one slide)

- 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 a language (and rejects all other strings)
 - Grammar** – a *generator*, a system for producing all strings in the language (and no other strings)
- A particular language may be specified by many different grammars and automata
- A grammar or automaton specifies only one language

11/14/2007

15

Regular Expressions and FAs

- The lexical grammar (structure) of most programming languages can be specified with regular expressions
 - Aside: Difficulties with Fortran
- Tokens can be recognized by a deterministic finite automaton
 - Can be either table-driven or built by hand based on lexical grammar

11/14/2007

16

Regular Expressions

- Defined over some alphabet Σ
 - For programming languages, commonly ASCII or Unicode
- If re is a regular expression, $L(re)$ is the language (set of strings) generated by re

11/14/2007

17

Fundamental REs

re	$L(re)$	Notes
a	$\{ a \}$	Singleton set, for each a in Σ
ϵ	$\{ \epsilon \}$	Empty string
\emptyset	$\{ \}$	Empty language

11/14/2007

18

Operations on REs

<i>re</i>	$L(re)$	Notes
rs	$L(r)L(s)$	Concatenation
$r s$	$L(r) \cup L(s)$	Combination (union)
r^*	$L(r)^*$	0 or more occurrences (Kleene closure)

- Precedence: * (highest), concatenation, | (lowest)
- Parentheses can be used to group REs as needed

11/14/2007

19

Abbreviations

- The basic operations generate all possible regular expressions, but there are common abbreviations used for convenience. Typical examples:

Abbr.	Meaning	Notes
$r+$	(r^*)	1 or more occurrences
$r?$	$(r \epsilon)$	0 or 1 occurrence
$[a-z]$	$(a b \dots z)$	1 character in given range
$[abxyz]$	$(a b x y z)$	1 of the given characters

11/14/2007

20

Examples

<i>re</i>	Meaning
$+$	single + character
$!$	single ! character
$=$	single = character
$!=$	2 character sequence
$<=$	2 character sequence
hogwash	7 character sequence

11/14/2007

21

More Examples

<i>re</i>	Meaning
$[abc]^+$	
$[abc]^*$	
$[0-9]^+$	
$[1-9][0-9]^*$	
$[a-zA-Z][a-zA-Z0-9_]^*$	

11/14/2007

22

Abbreviations

- Many systems allow abbreviations to make writing and reading definitions easier

name ::= *re*

- Restriction: abbreviations may not be circular (recursive) either directly or indirectly

11/14/2007

23

Example

- Possible syntax for numeric constants

```
digit ::= [0-9]
digits ::= digit+
number ::= digits ( . digits ) ?
           ( [eE] ( + | - ) ? digits ) ?
```

11/14/2007

24

Recognizing REs

- Finite automata can be used to recognize strings generated by regular expressions
- Can build by hand or automatically
 - Not totally straightforward, but can be done systematically
 - Tools like Lex, Flex, and JLex do this automatically, given a set of REs

11/14/2007

25

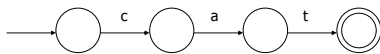
Finite State Automaton (FSA)

- A finite set of states
 - One marked as initial state
 - One or more marked as final states
 - States sometimes labeled or numbered
- A set of transitions from state to state
 - Each labeled with symbol from Σ , or ϵ
- Operate by reading input symbols (usually characters)
 - Transition can be taken if labeled with current symbol
 - ϵ -transition can be taken at any time
- Accept when final state reached & no more input
 - Scanner slightly different – accept longest match each time called, even if more input; i.e., run the FSA each time the scanner is called
- Reject if no transition possible or no more input and not in final state (DFA)

11/14/2007

26

Example: FSA for "cat"



11/14/2007

27

DFA vs NFA

- **Deterministic Finite Automata (DFA)**
 - No choice of which transition to take under any condition
- **Non-deterministic Finite Automata (NFA)**
 - Choice of transition in at least one case
 - Accept - if *some way* to reach final state on given input
 - Reject - if *no possible way* to final state

11/14/2007

28

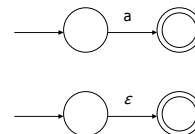
FAs in Scanners

- Want DFA for speed (no backtracking)
- Conversion from regular expressions to NFA is easy
- There is a well-defined procedure for converting a NFA to an equivalent DFA

11/14/2007

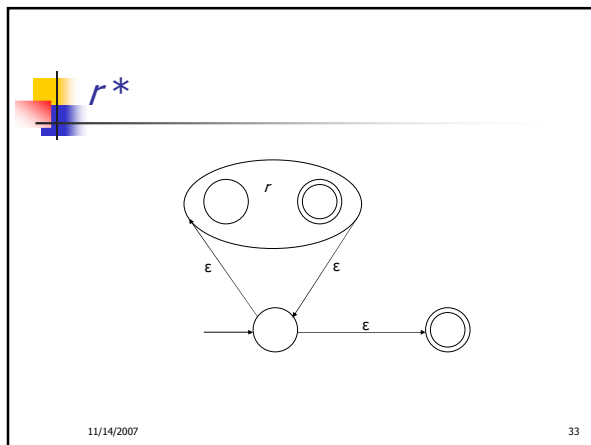
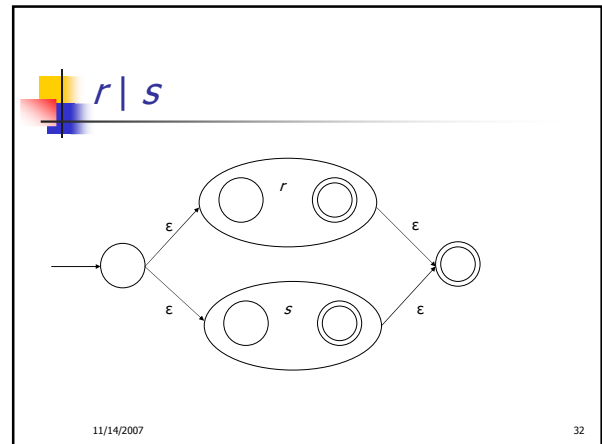
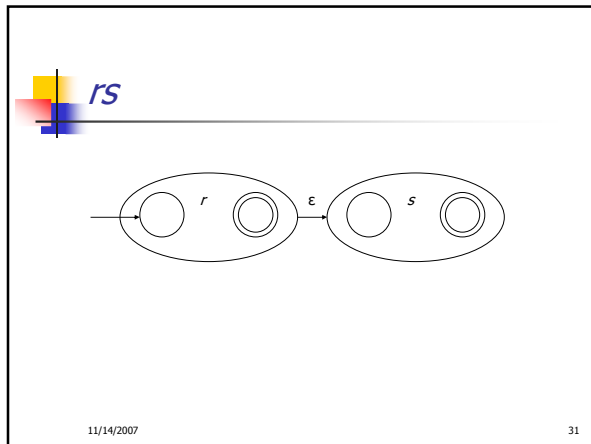
29

From RE to NFA: base cases



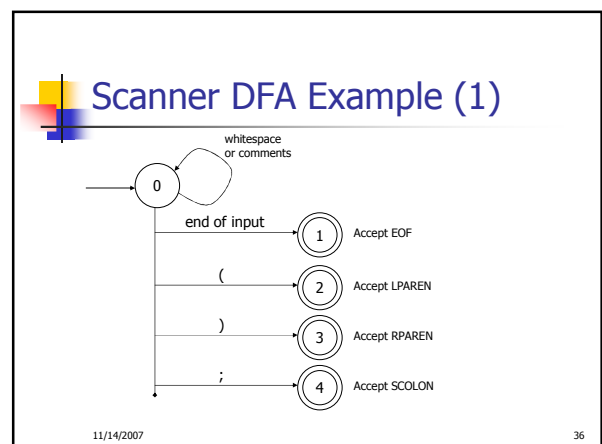
11/14/2007

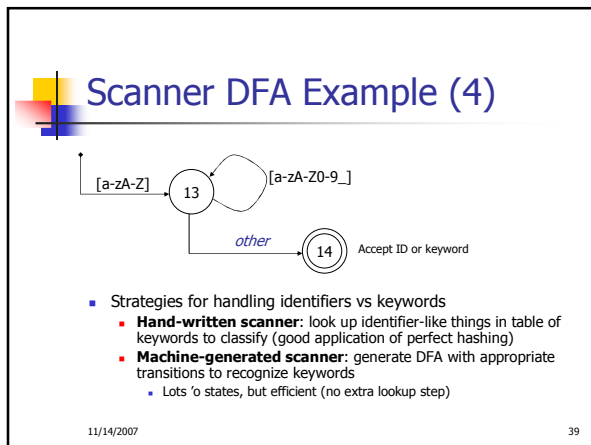
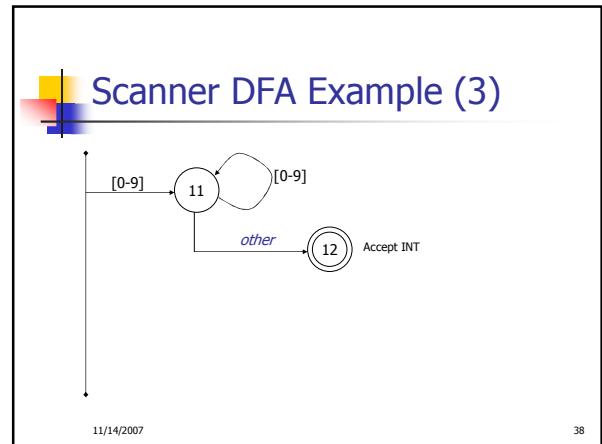
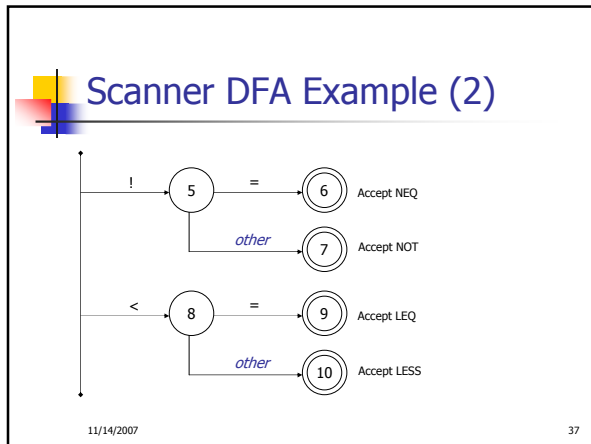
30



- ### From NFA to DFA
- Subset construction
 - Construct a DFA from the NFA, where each DFA state represents a *set of NFA states*
 - Key idea
 - The state of the DFA after reading some input is the set of *all* states the NFA could have reached after reading the same input
 - Algorithm: example of a fixed-point computation
 - If NFA has n states, DFA has at most 2^n states
 - => DFA is finite, can construct in finite # steps
 - Resulting DFA may have more states than needed
 - See books for construction and minimization details
- 11/14/2007 34

- ### Example: DFA for hand-written scanner
- **Idea:** show a hand-written DFA for some typical programming language constructs
 - Then use to construct hand-written scanner
 - **Setting:** Scanner is called whenever the parser needs a new token
 - Scanner stores current position in input file
 - Starting there, use a DFA to recognize the longest possible input sequence that makes up a token and return that token
- 11/14/2007 35





Implementing a Scanner by Hand – Token Representation

- A token is a simple, tagged structure

```
public class Token {
    public int kind;           // token's lexical class
    public int intVal;        // integer value if class = INT
    public String id;         // actual identifier if class = ID
    // lexical classes
    public static final int EOF = 0; // "end of file" token
    public static final int ID = 1;  // identifier, not keyword
    public static final int INT = 2; // integer
    public static final int LPAREN = 4;
    public static final int SCOLN = 5;
    public static final int WHILE = 6;
    // etc. etc. etc. ...
}
```

11/14/2007 40

Simple Scanner Example

```
// global state and methods
static char nextch; // next unprocessed input character

// advance to next input char
void getch() { ... }

// skip whitespace and comments
void skipWhitespace() { ... }
```

11/14/2007 41

Scanner getToken() method

```
// return next input token
public Token getToken() {
    Token result;

    skipWhiteSpace();

    if (no more input) {
        result = new Token(Token.EOF); return result;
    }

    switch(nextch) {
        case '(': result = new Token(Token.LPAREN); getch(); return result;
        case ')': result = new Token(Token.RPAREN); getch(); return result;
        case ';': result = new Token(Token.SCOLON); getch(); return result;
    }

    // etc. ...
}
```

11/14/2007 42

getToken() (2)

```
case '!': // ! or !=
    getch();
    if (nextch == '=') {
        result = new Token(Token.NEQ); getch(); return result;
    } else {
        result = new Token(Token.NOT); return result;
    }

case '<': // < or <=
    getch();
    if (nextch == '=') {
        result = new Token(Token.LEQ); getch(); return result;
    } else {
        result = new Token(Token.LESS); return result;
    }

// etc. ...
```

11/14/2007

43

getToken() (3)

```
case '0': case '1': case '2': case '3': case '4':
case '5': case '6': case '7': case '8': case '9':
    // integer constant
    String num = nextch;
    getch();
    while (nextch is a digit) {
        num = num + nextch; getch();
    }
    result = new Token(Token.INT, Integer(num).intValue());
    return result;

...
```

11/14/2007

44

getToken (4)

```
case 'a': ... case 'z':
case 'A': ... case 'Z': // id or keyword
    string s = nextch; getch();
    while (nextch is a letter, digit, or underscore) {
        s = s + nextch; getch();
    }
    if (s is a keyword) {
        result = new Token(keywordTable.getKind(s));
    } else {
        result = new Token(Token.ID, s);
    }
    return result;
```

11/14/2007

45