

## Syntax and Parsing

COMS W4115

Prof. Stephen A. Edwards  
Spring 2003

Columbia University  
Department of Computer Science

# Lexical Analysis (Scanning)

## Lexical Analysis (Scanning)

Goal is to translate a stream of characters

i n t s p g c d ( i n t s p  
a , s p i n t s p b

into a stream of tokens

ID	ID	LPAREN	ID	ID	COMMA	ID	ID
int	gcd	(	int	a	,	int	b

Each token consists of a token type and its text.

Whitespace and comments are discarded.

## Lexical Analysis

Goal: simplify the job of the parser.

Scanners are usually much faster than parsers.

Discard as many irrelevant details as possible (e.g.,  
whitespace, comments).

Parser does not care that the the identifier is  
"supercalifragilisticexpialidocious."

Parser rules are only concerned with token types.

## The ANTLR Compiler Generator

Language and compiler for writing compilers

Running ANTLR on an ANTLR file produces Java source  
files that can be compiled and run.

ANTLR can generate

- Scanners (lexical analyzers)
- Parsers
- Tree walkers

## An ANTLR File for a Simple Scanner

```
class CalcLexer extends Lexer;
```

```
LPAREN : '(' ; // Rules for punctuation  
RPAREN : ')' ;  
STAR : '*';  
PLUS : '+' ;  
SEMI : ';' ;  
protected // Can only be used as a sub-rule  
DIGIT : '0'..'9' ; // Any character between 0 and 9  
INT : (DIGIT)+ ; // One or more digits  
  
WS : (' ' | '\t' | '\n' | '\r') // Whitespace  
    { setType(Token.SKIP); } ; // Action: ignore
```

## ANTLR Specifications for Scanners

Rules are names starting with a capital letter.

A character in single quotes matches that character.

```
LPAREN : '(' ;
```

A string in double quotes matches the string

```
IF : "if" ;
```

A vertical bar indicates a choice:

```
OP : '+' | '-' | '*' | '/' ;
```

## ANTLR Specifications

Question mark makes a clause optional.

```
PERSON : ("wo")? 'm' ('a' | 'e') 'n' ;
```

(Matches man, men, woman, and women.)

Double dots indicate a range of characters:

```
DIGIT : '0'..'9' ;
```

Asterisk and plus match "zero or more," "one or more."

```
ID : LETTER (LETTER | DIGIT)* ;
```

```
NUMBER : (DIGIT)+ ;
```

## Kleene Closure

The asterisk operator (\*) is called the Kleene Closure  
operator after the inventor of regular expressions, Stephen  
Cole Kleene, who pronounced his last name "CLAY-nee."

His son Ken writes "As far as I am aware this  
pronunciation is incorrect in all known languages. I believe  
that this novel pronunciation was invented by my father."

## Scanner Behavior

All rules (tokens) are considered simultaneously. The longest one that matches wins:

1. Look at the next character in the file.
2. Can the next character be added to any of the tokens under construction?
3. If so, add the character to the token being constructed and go to step 1.
4. Otherwise, return the token.

How to keep track of multiple rules matching simultaneously? Build an automata.

## Deterministic Finite Automata

A state machine with an initial state

Arcs indicate "consumed" input symbols.

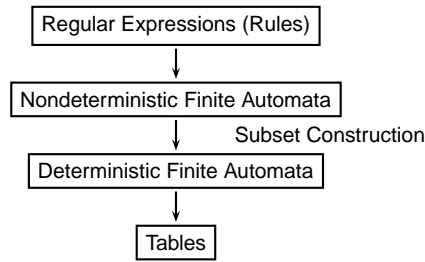
States with double lines are accepting.

If the next token has an arc, follow the arc.

If the next token has no arc and the state is accepting, return the token.

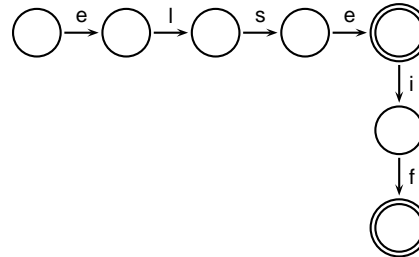
If the next token has no arc and the state is not accepting, syntax error.

## Implementing Scanners Automatically



## Deterministic Finite Automata

ELSE: "else" ;  
ELSEIF: "elseif" ;



## Nondeterministic Finite Automata

DFAs with  $\epsilon$  arcs.

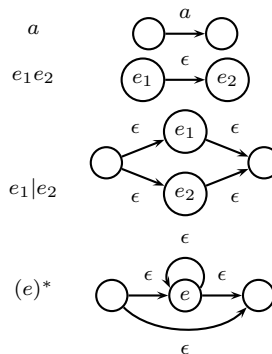
Conceptually,  $\epsilon$  arcs denote state equivalence.

$\epsilon$  arcs add the ability to make nondeterministic (schizophrenic) choices.

When an NFA reaches a state with an  $\epsilon$  arc, it moves into every destination.

NFAs can be in multiple states at once.

## Translating REs into NFAs



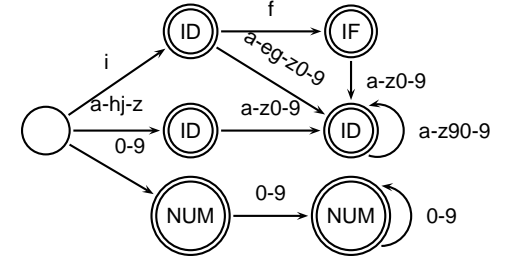
## Regular Expressions and NFAs

We are describing tokens with *regular expressions*:

- The symbol  $\epsilon$  always matches
- A symbol from an alphabet, e.g.,  $a$ , matches itself
- A sequence of two regular expressions e.g.,  $e_1e_2$  Matches  $e_1$  followed by  $e_2$
- An "OR" of two regular expressions e.g.,  $e_1|e_2$  Matches  $e_1$  or  $e_2$
- The Kleene closure of a regular expression, e.g.,  $(e)^*$  Matches zero or more instances of  $e_1$  in sequence.

## Deterministic Finite Automata

IF: "if" ;  
ID: 'a'..'z' ('a'..'z' | '0'..'9')\* ;  
NUM: ('0'..'9')+ ;

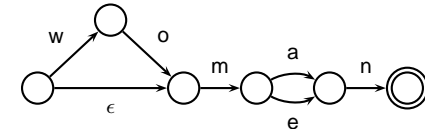


## RE to NFAs

Building an NFA for the regular expression

$(wo|\epsilon)m(a|e)n$

produces

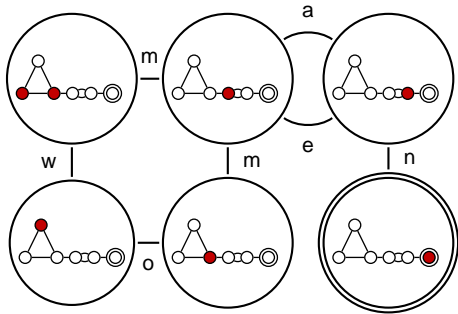


after simplification. Most  $\epsilon$  arcs disappear.

## Subset Construction

How to compute a DFA from an NFA.

Basic idea: each state of the DFA is a *marking* of the NFA



## Free-Format Languages

Java C C++ Algol Pascal

Some deviate a little (e.g., C and C++ have a separate preprocessor)

But not all languages are free-format.

## Syntax and Language Design

Does syntax matter? Yes and no

More important is a language's *semantics*—its meaning.

The syntax is aesthetic, but can be a religious issue.

But aesthetics matter to people, and can be critical.

Verbosity does matter: smaller is usually better.

Too small can be a problem: APL is a compact, cryptic language with its own character set (!)

```
E ← A TEST B;L
```

```
L ← 0.5
```

```
E ← ((A × A) + B × B) * L
```

## Subset Construction

An DFA can be exponentially larger than the corresponding NFA.

$n$  states versus  $2^n$

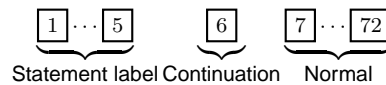
Tools often try to strike a balance between the two representations.

ANTLR uses a different technique.

## FORTRAN 77

FORTRAN 77 is not free-format. 72-character lines:

```
100  IF(IN .EQ. 'Y' .OR. IN .EQ. 'y' .OR.
      $  IN .EQ. 'T' .OR. IN .EQ. 't') THEN
```



When column 6 is not a space, line is considered part of the previous.

Fixed-length line works well with a one-line buffer.

Makes sense on punch cards.

## Syntax and Language Design

Some syntax is error-prone. Classic FORTRAN example:

```
DO 5 I = 1,25 ! Loop header (for i = 1 to 25)
```

```
DO 5 I = 1.25 ! Assignment to variable D05I
```

Trying too hard to reuse existing syntax in C++:

```
vector< vector<int> > foo;
vector<vector<int>> foo; // Syntax error
```

C distinguishes > and >> as different operators.

## Free-Format Languages

Typical style arising from scanner/parser division

Program text is a series of tokens possibly separated by whitespace and comments, which are both ignored.

- keywords (`if while`)
- punctuation (`, ( +`)
- identifiers (`foo bar`)
- numbers (`10 -3.14159e+32`)
- strings (`"A String"`)

## Python

The Python scripting language groups with indentation

```
i = 0
while i < 10:
    i = i + 1
    print i      # Prints 1, 2, ..., 10
```

```
i = 0
while i < 10:
    i = i + 1
print i          # Just prints 10
```

This is succinct, but can be error-prone.

How do you wrap a conditional around instructions?

## Keywords

Keywords look like identifiers in most languages.

Scanners do not know context, so keywords must take precedence over identifiers.

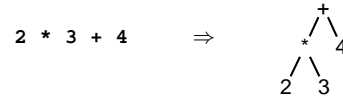
Too many keywords leaves fewer options for identifiers.

Languages such as C++ or Java strive for fewer keywords to avoid "polluting" available identifiers.

# Parsing

## Parsing

Objective: build an abstract syntax tree (AST) for the token sequence from the scanner.



Goal: discard irrelevant information to make it easier for the next stage.

Parentheses and most other forms of punctuation removed.

## Issues

Ambiguous grammars

Precedence of operators

Left- versus right-recursive

Top-down vs. bottom-up parsers

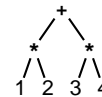
Parse Tree vs. Abstract Syntax Tree

## Operator Precedence

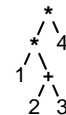
Defines how “sticky” an operator is.

$1 * 2 + 3 * 4$

\* at higher precedence than +:  
 $(1 * 2) + (3 * 4)$



+ at higher precedence than \*:  
 $1 * (2 + 3) * 4$



## Grammars

Most programming languages described using a *context-free grammar*.

Compared to regular languages, context-free languages add one important thing: recursion.

Recursion allows you to count, e.g., to match pairs of nested parentheses.

Which languages do humans speak? I'd say it's regular: I do not not not not not not not not not not understand this sentence.

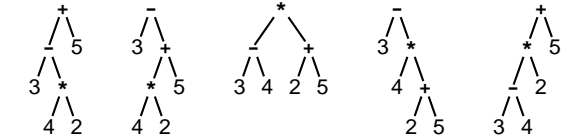
## Ambiguous Grammars

A grammar can easily be ambiguous. Consider parsing

$3 - 4 * 2 + 5$

with the grammar

$e \rightarrow e + e \mid e - e \mid e * e \mid e / e$



## Languages

Regular languages ( $t$  is a terminal):

$A \rightarrow t_1 \dots t_n B$

$A \rightarrow t_1 \dots t_n$

Context-free languages ( $P$  is terminal or a variable):

$A \rightarrow P_1 \dots P_n$

Context-sensitive languages:

$\alpha_1 A \alpha_2 \rightarrow \alpha_1 B \alpha_2$

“ $B \rightarrow A$  only in the ‘context’ of  $\alpha_1 \dots \alpha_2$ ”

## Operator Precedence and Associativity

Usually resolve ambiguity in arithmetic expressions

Like you were taught in elementary school:

“My Dear Aunt Sally”

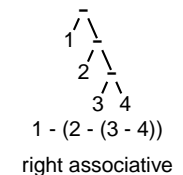
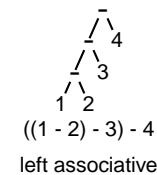
Mnemonic for multiplication and division before addition and subtraction.

## Associativity

Whether to evaluate left-to-right or right-to-left

Most operators are left-associative

$1 - 2 - 3 - 4$



## Fixing Ambiguous Grammars

Original ANTLR grammar specification

```
expr
: expr '+' expr
| expr '-' expr
| expr '*' expr
| expr '/' expr
| NUMBER
;
```

Ambiguous: no precedence or associativity.

## Parsing Context-Free Grammars

There are  $O(n^3)$  algorithms for parsing arbitrary CFGs, but most compilers demand  $O(n)$  algorithms.

Fortunately, the LL and LR subclasses of CFGs have  $O(n)$  parsing algorithms. People use these in practice.

## Writing LL(k) Grammars

Cannot have left-recursion

```
expr : expr '+' term | term ;
```

becomes

AST expr() –

```
switch (next-token) –
case NUMBER : expr(); /* Infinite Recursion */
```

## Assigning Precedence Levels

Split into multiple rules, one per level

```
expr : expr '+' expr
      | expr '-' expr
      | term ;

term : term '*' term
      | term '/' term
      | atom ;

atom : NUMBER ;
```

Still ambiguous: associativity not defined

## Parsing LL(k) Grammars

LL: Left-to-right, Left-most derivation

k: number of tokens to look ahead

Parsed by top-down, predictive, recursive parsers

Basic idea: look at the next token to predict which production to use

ANTLR builds recursive LL(k) parsers

Almost a direct translation from the grammar.

## Writing LL(1) Grammars

Cannot have common prefixes

```
expr : ID '(' expr ')'
      | ID '=' expr
```

becomes

AST expr() –

```
switch (next-token) –
case ID : match(ID); match('('); expr(); match(')');
case ID : match(ID); match('='); expr();
```

## Assigning Associativity

Make one side or the other the next level of precedence

```
expr : expr '+' term
      | expr '-' term
      | term ;

term : term '*' atom
      | term '/' atom
      | atom ;

atom : NUMBER ;
```

## A Top-Down Parser

```
stmt : 'if' expr 'then' expr
      | 'while' expr 'do' expr
      | expr ':'=' expr ;
```

```
expr : NUMBER | '(' expr ')';
```

AST stmt() {

```
switch (next-token) {
case "if" : match("if"); expr(); match("then"); expr();
case "while" : match("while"); expr(); match("do"); expr();
case NUMBER or "(" : expr(); match(":="); expr();
}
}
```

}

## Eliminating Common Prefixes

Consolidate common prefixes:

```
expr
: expr '+' term
| expr '-' term
| term
;
```

becomes

```
expr
: expr ('+' term | '-' term )
| term
;
```

## Eliminating Left Recursion

Understand the recursion and add tail rules

```
expr
  : expr ('+' term | '-' term )
  | term
  ;
```

becomes

```
expr : term exprt ;
exprt : '+' term exprt
      | '-' term exprt
      | /* nothing */
      ;
```

## The Dangling Else Problem

```
stmt : "if" expr "then" stmt iftail
      | other-statements ;
```

```
iftail
  : "else" stmt
  | /* nothing */
  ;
```

Problem comes when matching "iftail."

Normally, an empty choice is taken if the next token is in the "follow set" of the rule. But since "else" can follow an iftail, the decision is ambiguous.

## Statement separators/terminators

C uses ; as a statement terminator.

```
if (a<b) printf("a less");
else {
  printf("b"); printf(" less");
}
```

Pascal uses ; as a statement separator.

```
if a < b then writeln('a less')
else begin
  write('a'); writeln(' less')
end
```

Pascal later made a final ; optional.

## Using ANTLR's EBNF

ANTLR makes this easier since it supports \* and -:

```
expr : expr '+' term
      | expr '-' term
      | term ;
```

becomes

```
expr : term ('+' term | '-' term)* ;
```

## The Dangling Else Problem

ANTLR can resolve this problem by making certain rules "greedy." If a conditional is marked as greedy, it will take that option even if the "nothing" option would also match:

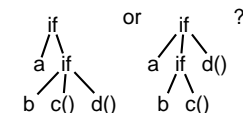
```
stmt
  : "if" expr "then" stmt
    ( options {greedy = true;}
    : "else" stmt
    )?
  | other-statements
  ;
```

## Bottom-up Parsing

## The Dangling Else Problem

Who owns the *else*?

```
if (a) if (b) c(); else d();
```



Grammars are usually ambiguous; manuals give disambiguating rules such as C's:

As usual the "else" is resolved by connecting an else with the last encountered elseless if.

## The Dangling Else Problem

Some languages resolve this problem by insisting on nesting everything.

E.g., Algol 68:

```
if a < b then a else b fi;
```

"fi" is "if" spelled backwards. The language also uses do-od and case-esac.

## Rightmost Derivation

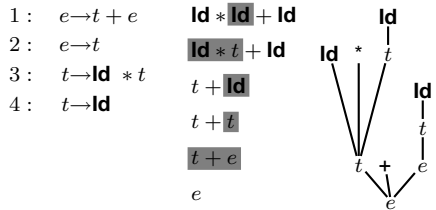
```
1: e → t + e
2: e → t
3: t → Id * t
4: t → Id
```

A rightmost derivation for **Id \* Id + Id**:

```
  e
t + e
t + Id
Id + Id
Id * Id + Id
Id * Id + Id
```

Basic idea of bottom-up parsing:  
construct this rightmost derivation  
**backward.**

## Handles

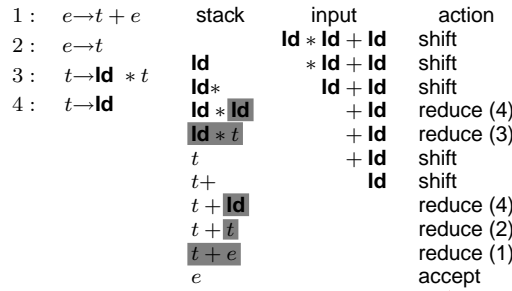


This is a reverse rightmost derivation for **ld \* ld + ld**.

Each highlighted section is a **handle**.

Taken in order, the handles build the tree from the leaves to the root.

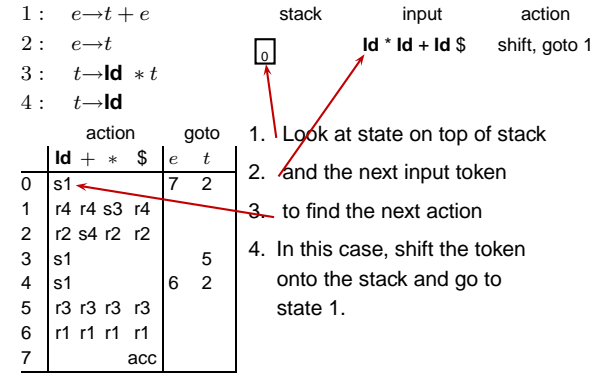
## Shift-reduce Parsing



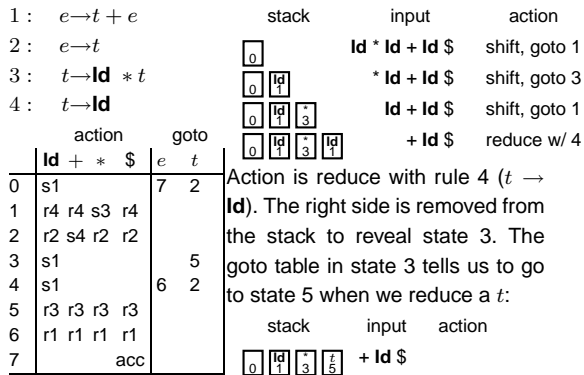
Scan input left-to-right, looking for handles.

An oracle tells what to do

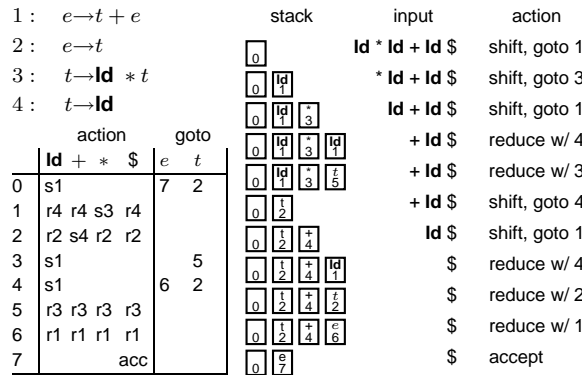
## LR Parsing



## LR Parsing



## LR Parsing



## Constructing the SLR Parse Table

The states are places we could be in a reverse-rightmost derivation. Let's represent such a place with a dot.

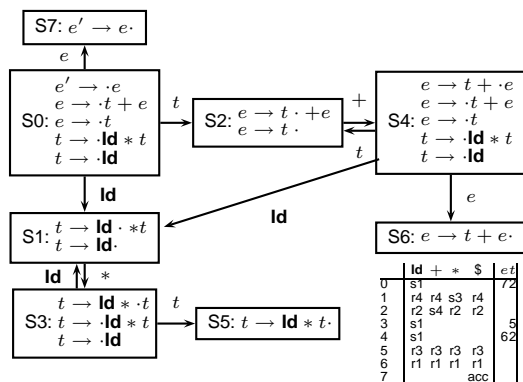
- $e \rightarrow t + e$
- $e \rightarrow t$
- $t \rightarrow ld * t$
- $t \rightarrow ld$

Say we were at the beginning ( $\cdot e$ ). This corresponds to

- $e' \rightarrow \cdot e$
- $e \rightarrow \cdot t + e$
- $e \rightarrow t \cdot$
- $t \rightarrow \cdot ld * t$
- $t \rightarrow \cdot ld$

The first is a placeholder. The second are the two possibilities when we're just before  $e$ . The last two are the two possibilities when we're just before  $t$ .

## Constructing the SLR Parsing Table



## The Punchline

This is a tricky, but mechanical procedure. The parser generators YACC, Bison, Cup, and others (but not ANTLR) use a modified version of this technique to generate fast bottom-up parsers.

You need to understand it to comprehend error messages:

Shift/reduce conflicts are caused by a state like

$t \rightarrow ld \cdot * t$

$t \rightarrow ld * \cdot t$

Reduce/reduce conflicts are caused by a state like

$t \rightarrow ld * t \cdot$

$e \rightarrow t + e \cdot$