

COMPILER DESIGN

LECTURE NOTES ON COMPILER DESIGN

Prepared by
Dr. Subasish Mohapatra



Welcome

**Department of Computer Science and Application
College of Engineering and Technology, Bhubaneswar
Biju Patnaik University of Technology, Odisha**

SYLLABUS

Compiler Design (3-0-0)

MODULE – 1 (Lecture hours: 13)

Introduction: Overview and phases of compilation. (2-hours)

Lexical Analysis: Non-deterministic and deterministic finite automata (NFA & DFA), regular grammar, regular expressions and regular languages, design of a lexical analyzer as a DFA, lexical analyser generator. (3-hours)

Syntax Analysis: Role of a parser, context free grammars and context free languages, parse trees and derivations, ambiguous grammar.

Top Down Parsing: Recursive descent parsing, LL(1) grammars, non-recursive predictive parsing, error reporting and recovery.

Bottom Up Parsing: Handle pruning and shift reduces parsing, SLR parsers and construction of SLR parsing tables, LR(1) parsers and construction of LR(1) parsing tables, LALR parsers and construction of efficient LALR parsing tables, parsing using ambiguous grammars, error reporting and recovery, parser generator. (8-hours)

MODULE – 2 (Lecture hours: 14)

Syntax Directed Translation: Syntax directed definitions (SDD), inherited and synthesized attributes, dependency graphs, evaluation orders for SDD, semantic rules, application of syntax directed translation. (5-hours)

Symbol Table: Structure and features of symbol tables, symbol attributes and scopes. (2-hours)

Intermediate Code Generation: DAG for expressions, three address codes - quadruples and triples, types and declarations, translation of expressions, array references, type checking and conversions, translation of Boolean expressions and control flow statements, back patching, intermediate code generation for procedures. (7-hours)

MODULE – 3 (Lecture hours: 8)

Run Time Environment: storage organizations, static and dynamic storage allocations, stack allocation, handlings of activation records for calling sequences. (3-hours)

Code Generations: Factors involved, registers allocation, simple code generation using stack allocation, basic blocks and flow graphs, simple code generation using flow graphs. (3-hours)

Elements of Code Optimization: Objective, peephole optimization, concepts of elimination of local common sub-expressions, redundant and un-reachable codes, basics of flow of control optimization. (2-hours)

CONTENTS

| | |
|------------|---|
| Lecture-1 | Introduction to compiler & its phases |
| Lecture-2 | Overview of language processing system |
| Lecture-3 | Phases of a Compiler |
| Lecture-4 | Languages |
| Lecture-5 | Converting RE to NFA (Thomson Construction) |
| Lecture-6 | Lexical Analysis |
| Lecture-7 | Lexical Analyzer Generator |
| Lecture-8 | Basics of Syntax Analysis |
| Lecture-9 | Context-Free Grammar |
| Lecture-10 | Left Recursion |
| Lecture-11 | YACC |
| Lecture-12 | Top-down Parsing |
| Lecture-13 | Recursive Predictive Parsing |
| Lecture-14 | Non-recursive Predictive Parsing-LL(1) |
| Lecture-15 | LL(1) Grammar |
| Lecture-16 | Basics of Bottom-up parsing |
| Lecture-17 | Conflicts during shift-reduce parsing |
| Lecture-18 | Operator precedence parsing |
| Lecture-19 | LR Parsing |
| Lecture-20 | Construction of SLR parsing table |
| Lecture-21 | Construction of canonical LR(0) collection |
| Lecture-22 | Shift-Reduce & Reduce-Reduce conflicts |
| Lecture-23 | Construction of canonical LR(1) collection |
| Lecture-24 | Construction of LALR parsing table |
| Lecture-25 | Using ambiguous grammars |
| Lecture-26 | SYNTAX-DIRECTED TRANSLATION |
| Lecture-27 | Translation of Assignment Statements |
| Lecture-28 | Generating 3-address code for Numerical Representation of Boolean expressions |
| Lecture-29 | Statements that Alter Flow of Control |
| Lecture-30 | Postfix Translations |
| Lecture-31 | Array references in arithmetic expressions |
| Lecture-32 | SYMBOL TABLES |
| Lecture-33 | Intermediate Code Generation |
| Lecture-34 | Directed Acyclic Graph |
| Lecture-35 | Flow of control statements with Jump method |
| Lecture-36 | Backpatching |
| Lecture-37 | RUN TIME ADMINISTRATION |
| Lecture-38 | Storage Organization |
| Lecture-39 | ERROR DETECTION AND RECOVERY |
| Lecture-40 | Error Recovery in Predictive Parsing |
| Lecture-41 | CODE OPTIMIZATION |
| Lecture-42 | Local Optimizations |

Module-1

Lecture #1

INTRODUCTION TO COMPILERS AND ITS PHASES

A compiler is a program that takes a program written in a source language and translates it into an equivalent program in a target language. The source language is a high-level language and the target language is machine language.

Source program -> COMPILER -> Target program

Necessity of compiler

- Techniques used in a lexical analyzer can be used in text editors, information retrieval systems, and pattern recognition programs.
- Techniques used in a parser can be used in a query processing system such as SQL.
- Many software packages having a complex front-end may need techniques used in compiler design.
- A symbolic equation solver which takes an equation as input. That program should parse the given input equation.
- Most of the techniques used in compiler design can be used in Natural Language Processing (NLP) systems.

Properties of Compiler

- a) Correctness
 - i) Correct output in execution.
 - ii) It should report errors
 - iii) Correctly report if the programmer is not following language syntax.
- b) Efficiency
- c) Compile time and execution.
- d) Debugging / Usability.

| Compiler | Interpreter |
|---|--|
| <ol style="list-style-type: none">1. It translates the whole program at a time.2. Compiler is faster.3. Debugging is not easy.4. Compilers are not portable. | <ol style="list-style-type: none">1. It translates statement by statement.2. Interpreter is slower.3. Debugging is easy.4. Interpreters are portable. |

Types of compiler

1) Native code compiler

A compiler may produce binary output to run / execute on the same computer and operating system. This type of compiler is called as native code compiler.

2) Cross Compiler

A cross compiler is a compiler that runs on one machine and produces object code for another machine.

3) Bootstrap compiler

If a compiler has been implemented in its own language, it is a self-hosting compiler.

4) **One pass compiler**

The compilation is done in one pass over the source program, hence the compilation is completed very quickly. This is used for the programming language PASCAL, COBOL, FORTRAN.

5) **Multi-pass compiler(2 or 3 pass compiler)**

In this compiler, the compilation is done step by step. Each step uses the result of the previous step and it creates another intermediate result.

Example:- gcc, Turbo C++

6) **JIT Compiler**

This compiler is used for JAVA programming language and Microsoft .NET

7) **Source to source compiler**

It is a type of compiler that takes a high level language as an input and its output as high level language. Example Open MP

List of compiler

1. Ada compiler
2. ALGOL compiler
3. BASIC compiler
4. C# compiler
5. C compiler
6. C++ compiler
7. COBOL compiler
8. Smalltalk compiler
9. Java compiler

Lecture #2

OVERVIEW OF LANGUAGE PROCESSING SYSTEM

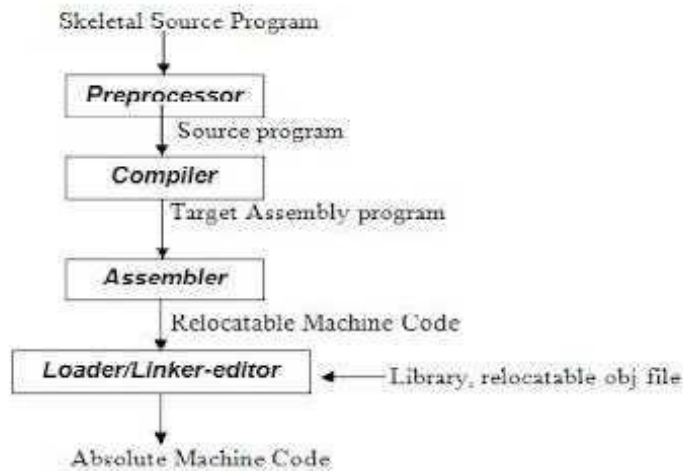


Fig 1.1 Language processing System

A source program may be divided into modules stored in separate files.

Preprocessor —collects all the separate files to the source program.

A preprocessor produce input to compilers. They may perform the following functions.

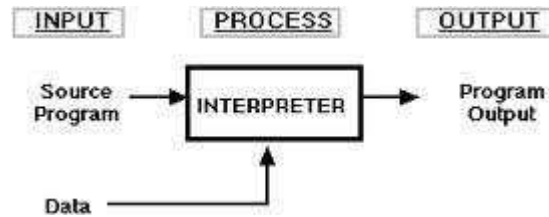
1. *Macro processing*: A preprocessor may allow a user to define macros that are short hands for longer constructs.
2. *File inclusion*: A preprocessor may include header files into the program text.
3. *Rational preprocessor*: these preprocessors augment older languages with more modern flow-of-control and data structuring facilities.
3. *Language Extensions*: These preprocessor attempts to add capabilities to the language by certain amounts to build-in macro

ASSEMBLER

Programmers found it difficult to write or read programs in machine language. They begin to use a mnemonic (symbols) for each machine instruction, which they would subsequently translate into machine language. Such a mnemonic machine language is now called an assembly language. Programs known as assembler were written to automate the translation of assembly language in to machine language. The input to an assembler program is called source program, the output is a machine language translation (object program).

INTERPRETER

An interpreter is a program that appears to execute a source program as if it were machine language



Languages such as BASIC, SNOBOL, LISP can be translated using interpreters. JAVA also uses interpreter. The process of interpretation can be carried out in following phases.

1. Lexical analysis
2. Syntax analysis
3. Semantic analysis
4. Direct Execution

Advantages

- Modification of user program can be easily made and implemented as execution proceeds.
- Type of object that denotes a various may change dynamically.
- Debugging a program and finding errors is simplified task for a program used for interpretation.
- The interpreter for the language makes it machine independent.

Disadvantages

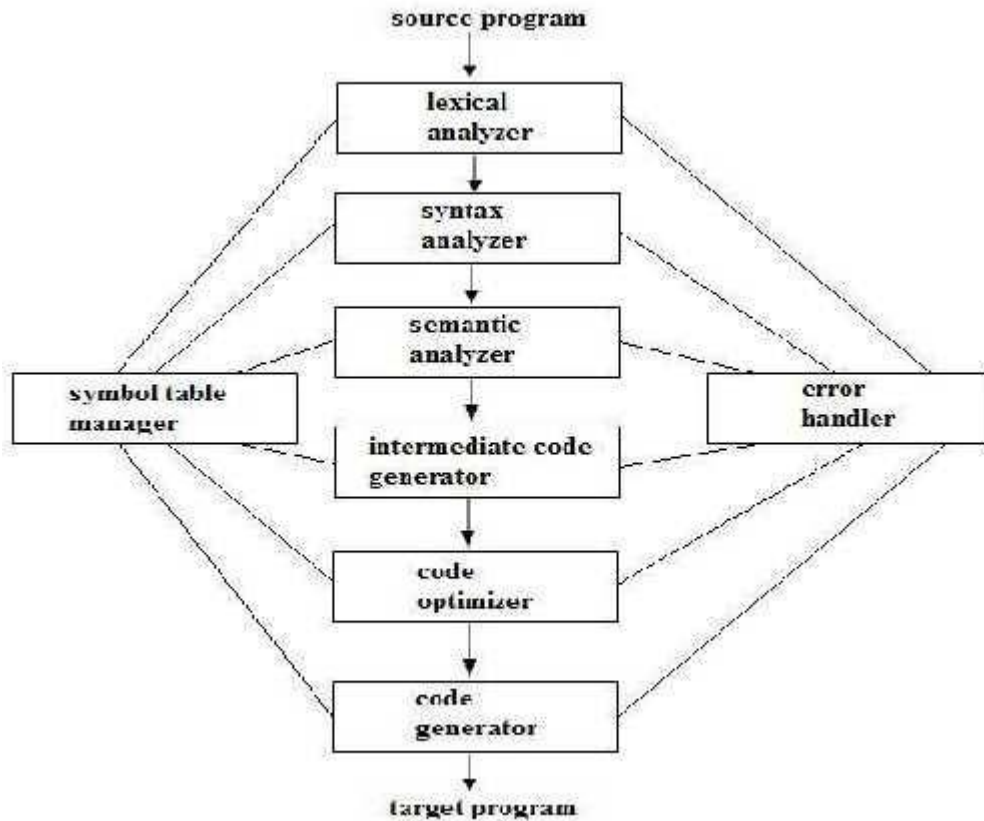
- The execution of the program is *slower*.
- Memory consumption is more.

Loader and Linker

Once the assembler procedures an object program, that program must be placed into memory and executed. The assembler could place the object program directly in memory and transfer control to it, thereby causing the machine language program to be execute. This would waste core by leaving the assembler in memory while the user's program was being executed. Also the programmer would have to retranslate his program with each execution, thus wasting translation time. To over come this problems of wasted translation time and memory. System programmers developed another component called Loader

“A loader is a program that places programs into memory and prepares them for execution.” It would be more efficient if subroutines could be translated into object form the loader could” relocate” directly behind the user's program. The task of adjusting programs o they may be placed in arbitrary core locations is called relocation. Relocation loaders perform four functions.

STRUCTURE OF THE COMPILER DESIGN



Major Parts of a Compiler

There are two major parts of a compiler: Analysis and Synthesis

- In analysis phase, an intermediate representation is created from the given source program. Lexical Analyzer, Syntax Analyzer and Semantic Analyzer are the phases in this part.
- In synthesis phase, the equivalent target program is created from this intermediate representation. Intermediate Code Generator, Code Generator, and Code Optimizer are the phases in this part.



Lecture #3

Phases of a Compiler

Each phase transforms the source program from one representation into another representation. They communicate with error handlers and the symbol table.

Lexical Analyzer

- Lexical Analyzer reads the source program character by character and returns the *tokens* of the source program.
- A *token* describes a pattern of characters having same meaning in the source program. (such as identifiers, operators, keywords, numbers, delimiters and so on)

Example:

In the line of code `newval := oldval + 12`, tokens are:

| | |
|---------------------|-----------------------|
| <code>newval</code> | (identifier) |
| <code>:=</code> | (assignment operator) |
| <code>oldval</code> | (identifier) |
| <code>+</code> | (add operator) |
| <code>12</code> | (a number) |

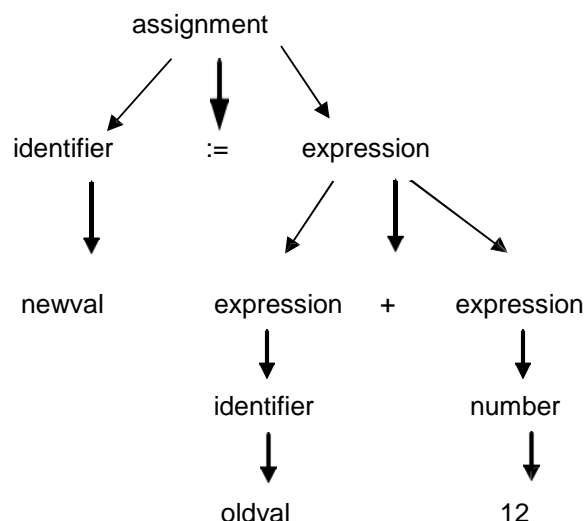
- Puts information about identifiers into the symbol table.
- Regular expressions are used to describe tokens (lexical constructs).
- A (Deterministic) Finite State Automaton can be used in the implementation of a lexical analyzer.

Syntax Analyzer

- A Syntax Analyzer creates the syntactic structure (generally a parse tree) of the given program.
- A syntax analyzer is also called a parser.
- A parse tree describes a syntactic structure.

Example:

For the line of code `newval := oldval + 12`, parse tree will be:



- The syntax of a language is specified by a context free grammar (CFG).
- The rules in a CFG are mostly recursive.
- A syntax analyzer checks whether a given program satisfies the rules implied by a CFG or not.

- If it satisfies, the syntax analyzer creates a parse tree for the given program.

Example:

CFG used for the above parse tree is:

assignment \rightarrow identifier := expression
expression \rightarrow identifier
expression \rightarrow number
expression \rightarrow expression + expression

- Depending on how the parse tree is created, there are different parsing techniques.
- These parsing techniques are categorized into two groups:
 - *Top-Down Parsing*,
 - *Bottom-Up Parsing*
- Top-Down Parsing:
 - Construction of the parse tree starts at the root, and proceeds towards the leaves.
 - Efficient top-down parsers can be easily constructed by hand.
 - Recursive Predictive Parsing, Non-Recursive Predictive Parsing (LL Parsing).
- Bottom-Up Parsing:
 - Construction of the parse tree starts at the leaves, and proceeds towards the root.
 - Normally efficient bottom-up parsers are created with the help of some software tools.
 - Bottom-up parsing is also known as shift-reduce parsing.
 - Operator-Precedence Parsing – simple, restrictive, easy to implement
 - LR Parsing – much general form of shift-reduce parsing, LR, SLR, LALR

Semantic Analyzer

- A semantic analyzer checks the source program for semantic errors and collects the type information for the code generation.
- Type-checking is an important part of semantic analyzer.
- Normally semantic information cannot be represented by a context-free language used in syntax analyzers.
- Context-free grammars used in the syntax analysis are integrated with attributes (semantic rules). The result is a syntax-directed translation and Attribute grammars

Example:

In the line of code `newval := oldval + 12`, the type of the identifier `newval` must match with type of the expression `(oldval+12)`.

Intermediate Code Generation

- A compiler may produce an explicit intermediate codes representing the source program.
- These intermediate codes are generally machine architecture independent. But the level of intermediate codes is close to the level of machine codes.

Example:

$newval := oldval * fact + 1$



$id1 := id2 * id3 + 1$



MULT id2, id3, temp1
ADD temp1, #1, temp2
MOV temp2, id1

The last form is the Intermediates Code (Quadruples)

Code Optimizer

The code optimizer optimizes the code produced by the intermediate code generator in the terms of time and space.

Example:

The above piece of intermediate code can be reduced as follows:

MULT id2, id3, temp1
ADD temp1, #1, id1

Code Generator

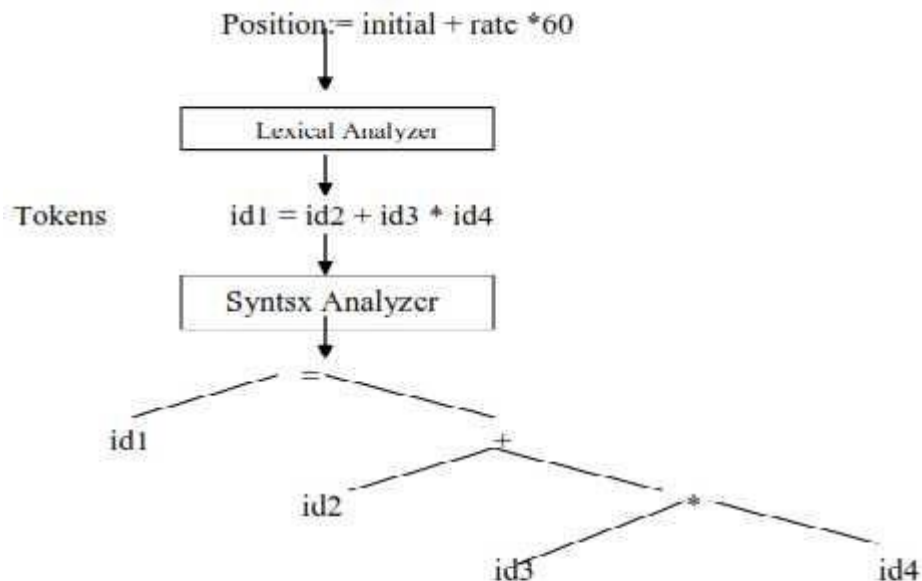
- Produces the target language in a specific architecture.
- The target program is normally is a relocatable object file containing the machine codes.

Example:

Assuming that we have architecture with instructions that have at least one operand as a machine register, the Final Code our line of code will be:

MOVE id2, R1
MULT id3, R1
ADD #1, R1
MOVE R1, id1

Ex:-



Phases of a compiler are the sub-tasks that must be performed to complete the compilation process. Passes refer to the number of times the compiler has to traverse through the entire program.

Symbol Table Management:

A symbol table is a data structure that contains a record for each identifier with field for attributes of the identifier.

The type information about the identifier is detected during the lexical analysis phases and is entered into the symbol table.

$$\text{Position} = \text{initial} + \text{rate} * 60;$$

| Address | Symbol | Location | attributes |
|---------|----------|----------|---------------|
| 1 | Position | 1000 | id, float |
| 2 | Initial | 2000 | id, float |
| 3 | Rate | 3000 | id, float |
| 4 | 60 | 4000 | constant, int |

Error Detection and Reporting:

Each phase detects/encounters errors after detecting errors.

This phase must deal with errors to continue with the process of compilation.

The following are some errors encountered in each phase:

- i) Lexical Analyzer- Miss spell token.
- ii) Semantic Analyzer- Type Mismatch.
- iii) Syntax Analyzer-Missing parenthesis , less no. of operands.
- iv) Intermediate code generation – In compatible operands for an operand.
- v) Code optimization- Unreachable statement.
- vi) Code Generation- Memory restriction to store a variable.

Lecture #4

Languages

Terminology

- Alphabet : a finite set of symbols (ASCII characters)
- String : finite sequence of symbols on an alphabet
- Sentence and word are also used in terms of string
- ϵ is the empty string
- $|s|$ is the length of string s .
- Language: sets of strings over some fixed alphabet
- \emptyset the empty set is a language.
- $\{\epsilon\}$ the set containing empty string is a language
- The set of all possible identifiers is a language.
- Operators on Strings:
- *Concatenation*: xy represents the concatenation of strings x and y . $s\epsilon = s$ $\epsilon s = s$
- $s^n = s s s \dots s$ (n times) $s^0 = \epsilon$

Operations on Languages

- Concatenation: $L_1 L_2 = \{ s_1 s_2 \mid s_1 \in L_1 \text{ and } s_2 \in L_2 \}$
- Union: $L_1 \cup L_2 = \{ s \mid s \in L_1 \text{ or } s \in L_2 \}$
- Exponentiation: $L^0 = \{\epsilon\}$ $L^1 = L$ $L^2 = LL$
- Kleene Closure: $L^* =$
- Positive Closure: $L^+ =$

Examples:

- $L_1 = \{a,b,c,d\}$ $L_2 = \{1,2\}$
- $L_1 L_2 = \{a1,a2,b1,b2,c1,c2,d1,d2\}$
- $L_1 \cup L_2 = \{a,b,c,d,1,2\}$
- $L_1^3 =$ all strings with length three (using a,b,c,d)
- $L_1^* =$ all strings using letters a,b,c,d and empty string
- $L_1^+ =$ doesn't include the empty string

Regular Expressions and Finite Automata

Regular Expressions

- We use regular expressions to describe tokens of a programming language.
- A regular expression is built up of simpler regular expressions (using defining rules)
- Each regular expression denotes a language.
- A language denoted by a regular expression is called as a regular set.

For Regular Expressions over alphabet Σ

| <u>Regular Expression</u> | <u>Language it denotes</u> |
|---------------------------|----------------------------|
| ϵ | $\{\epsilon\}$ |
| $a \in \Sigma$ | $\{a\}$ |
| $(r_1) \mid (r_2)$ | $L(r_1) \cup L(r_2)$ |
| $(r_1)(r_2)$ | $L(r_1)L(r_2)$ |
| $(r)^*$ | $L(r)^*$ |

- $(r)^+ = (r)(r)^*$
- $(r)? = (r) \mid \epsilon$
- We may remove parentheses by using precedence rules.
 - * highest
 - concatenation next
 - | lowest
- ab^*c means $(a(b^*))c$

Examples:

- $\Sigma = \{0,1\}$
- $0|1 = \{0,1\}$
- $(0|1)(0|1) = \{00,01,10,11\}$
- $0^* = \{\epsilon, 0, 00, 000, 0000, \dots\}$
- $(0|1)^* =$ All strings with 0 and 1, including the empty string

Finite Automata

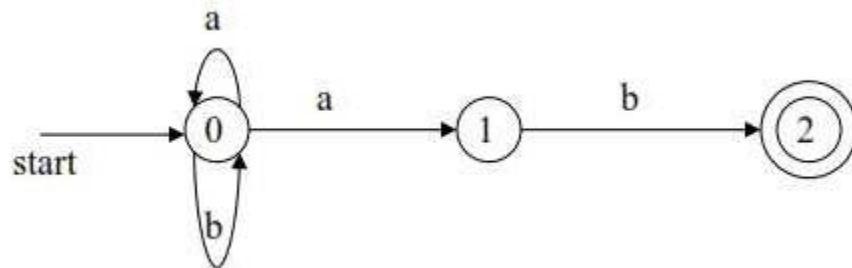
- A *recognizer* for a language is a program that takes a string x , and answers “yes” if x is a sentence of that language, and “no” otherwise.
- We call the recognizer of the tokens as a *finite automaton*.
- A finite automaton can be: *deterministic (DFA)* or *non-deterministic (NFA)*
- This means that we may use a deterministic or non-deterministic automaton as a lexical analyzer.

- Both deterministic and non-deterministic finite automaton recognize regular sets.
- Which one?
 - deterministic – faster recognizer, but it may take more space
 - non-deterministic – slower, but it may take less space
 - Deterministic automata are widely used lexical analyzers.
- First, we define regular expressions for tokens; Then we convert them into a DFA to get a lexical analyzer for our tokens.

Non-Deterministic Finite Automaton (NFA)

- A non-deterministic finite automaton (NFA) is a mathematical model that consists of:
 - S - a set of states
 - Σ - a set of input symbols (alphabet)
 - move - a transition function move to map state-symbol pairs to sets of states.
 - s_0 - a start (initial) state
 - F - a set of accepting states (final states)
- ϵ - transitions are allowed in NFAs. In other words, we can move from one state to another one
- without consuming any symbol.
- A NFA accepts a string x , if and only if there is a path from the starting state to one of accepting states such that edge labels along this path spell out x .

Example:



Transition Graph

0 is the start state s_0

{2} is the set of final states F

$\Sigma = \{a,b\}$

$S = \{0,1,2\}$

Transition Function:

| | a | b |
|---|-------|-----|
| 0 | {0,1} | {0} |
| 1 | {} | {2} |
| 2 | {} | {} |

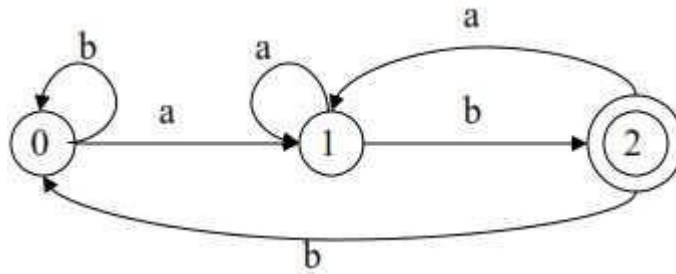
The language recognized by this NFA is $(a|b)^*ab$

Deterministic Finite Automaton (DFA)

- A Deterministic Finite Automaton (DFA) is a special form of a NFA.
- No state has ϵ - transition
- For each symbol a and state s , there is at most one labeled edge a leaving s . i.e. transition function is from pair of state-symbol to state (not set of states)

Example:

The DFA to recognize the language $(a|b)^* ab$ is as follows.



0 is the start state s_0

{2} is the set of final states F

$\Sigma = \{a,b\}$

$S = \{0,1,2\}$

Transition Function:

| | a | B |
|---|---|---|
| 0 | 1 | 0 |
| 1 | 1 | 2 |
| 2 | 1 | 0 |

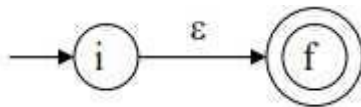
Note that the entries in this function are single value and not set of values (unlike NFA).

Lecture #5 \

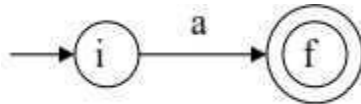
Converting RE to NFA (Thomson Construction)

- This is one way to convert a regular expression into a NFA.
- There can be other ways (much efficient) for the conversion.
- Thomson's Construction is simple and systematic method.
- It guarantees that the resulting NFA will have exactly one final state, and one start state.
- Construction starts from simplest parts (alphabet symbols).
- To create a NFA for a complex regular expression, NFAs of its sub-expressions are combined to create its NFA.

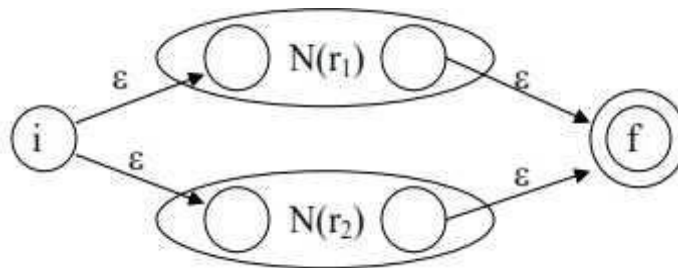
To recognize an empty string ϵ :



To recognize a symbol a in the alphabet Σ :



For regular expression $r_1 | r_2$:



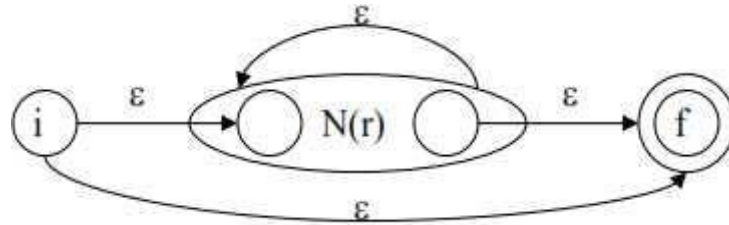
$N(r_1)$ and $N(r_2)$ are NFAs for regular expressions r_1 and r_2 .

For regular expression $r_1 \cdot r_2$



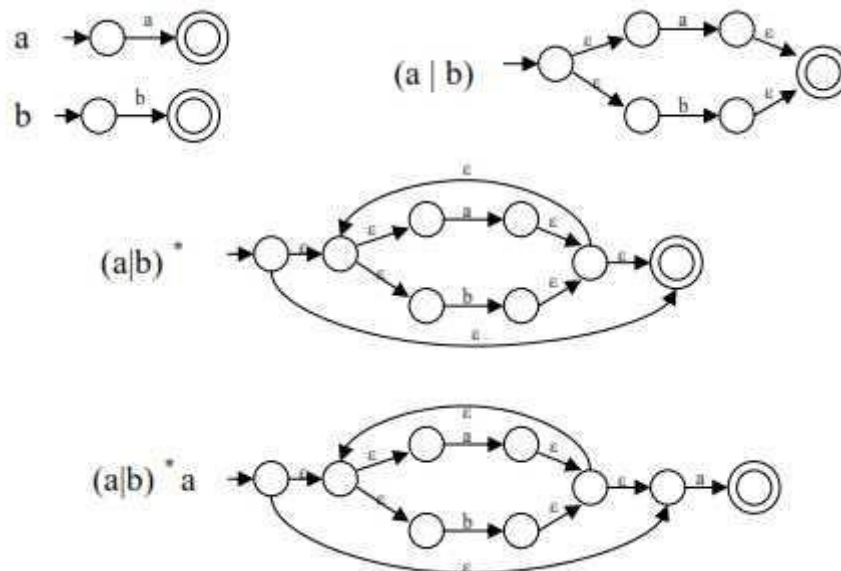
Here, final state of $N(r_1)$ becomes the final state of $N(r_1 r_2)$

For regular expression r^*



Example:

For a RE $(a|b)^* a$, the NFA construction is shown below.



Converting NFA to DFA (Subset Construction)

We merge together NFA states by looking at them from the point of view of the input characters:

From the point of view of the input, any two states that are connected by an ϵ -transition may as well be the same, since we can move from one to the other without consuming any character. Thus states which are connected by an ϵ -transition will be represented by the same states in the DFA.

If it is possible to have multiple transitions based on the same symbol, then we can regard a transition on a symbol as moving from a state to a set of states (ie. the union of all those states reachable by a transition on the current symbol). Thus these states will be combined into a single DFA state.

To perform this operation, let us define two functions:

- The ϵ -closure function takes a state and returns the set of states reachable from it based on (one or more) ϵ -transitions. Note that this will always include the state itself.

We should be able to get from a state to any state in its ϵ -closure without consuming any input.

- The function **move** takes a state and a character, and returns the set of states reachable by one transition on this character.

We can generalize both these functions to apply to sets of states by taking the union of the application to individual states.

For Example, if A, B and C are states, $\text{move}(\{A,B,C\}, 'a') = \text{move}(A, 'a') \cup \text{move}(B, 'a') \cup \text{move}(C, 'a')$.

The Subset Construction Algorithm is as follows:

put ϵ -closure($\{s_0\}$) as an unmarked state into the set of DFA (DS)

while (there is one unmarked S1 in DS) do

begin

mark S1

for each input symbol a do begin

S2 \leftarrow ϵ -closure(move(S1,a))

if (S2 is not in DS) then add S2 into DS as an unmarked state $\text{transfunc}[S1,a] \leftarrow S2$

end

End

ϵ -closure(move(S1,b)) = ϵ -closure({5}) = {1,2,4,5,6,7}

= S2 transfunc[S1,a]←S1 transfunc[S1,b]←S2

↓ mark S2

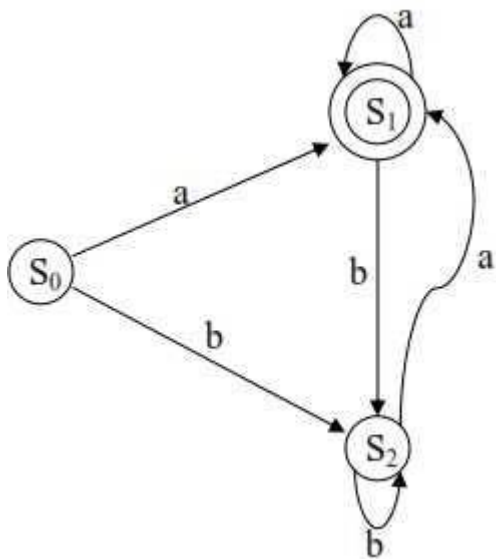
ϵ -closure(move(S2,a)) = ϵ -closure({3,8}) = {1,2,3,4,6,7,8} = S1

ϵ -closure(move(S2,b)) = ϵ -closure({5}) = {1,2,4,5,6,7}

= S2 transfunc[S2,a]←S1 transfunc[S2,b]←S2

S0 is the start state of DFA since 0 is a member of S0={0,1,2,4,7}

S1 is an accepting state of DFA since 8 is a member of S1 = {1,2,3,4,6,7,8}



Lecture #6

LEXICAL ANALYSIS

- Lexical Analyzer reads the source program character by character to produce tokens.
- Normally a lexical analyzer does not return a list of tokens at one shot; it returns a token when the parser asks a token from it.

Token, Pattern, Lexeme

- Token represents a set of strings described by a pattern. For example, an identifier represents a set of strings which start with a letter continues with letters and digits. The actual string is called as lexeme.
 - Since a token can represent more than one lexeme, additional information should be held for that specific lexeme. This additional information is called as the *attribute* of the token.
 - For simplicity, a token may have a single attribute which holds the required information for that token. For identifiers, this attribute is a pointer to the symbol table, and the symbol table holds the actual attributes for that token.
- Examples:
 - <identifier, attribute> where attribute is pointer to the symbol table
 - <assignment operator> no attribute is needed
 - <number, value> where value is the actual value of the number
 - Token type and its attribute uniquely identify a lexeme.
 - *Regular expressions* are widely used to specify patterns.

Pattern:

A pattern is a description of the form that the lexemes of a token may take. In the case of a keyword as a token, the pattern is just the sequence of characters that form the keyword. For identifiers and some other tokens, the pattern is a more complex structure that is matched by many strings.

Lexeme:

A lexeme is a sequence of characters in the source program that matches the pattern for a token and is identified by the lexical analyzer as an instance of that token.

| TOKEN | INFORMAL DESCRIPTION | SAMPLE LEXEMES |
|-------------------|---|----------------------------|
| if | characters i, f | if |
| else | characters e, l, s, e | else |
| comparison | < or > or <= or >= or == or != | <=, != |
| id | letter followed by letters and digits | pi, score, D2 |
| number | any numeric constant | 3.14159, 0, 6.02e23 |
| literal | anything but " , surrounded by " 's | "core dumped" |

Lexical Analysis versus parsing

There are a number of reasons the analysis portion of a compiler is separated into lexical analysis and parsing (syntax analysis) phases.

- Simplicity of design. The separation of lexical and syntactic analysis often allows us to simplify at least one of these tasks. For example, a parser that had to deal with comments and whitespace as syntactic units would be considerably more complex than one that can assume comments and whitespace have already been removed by the lexical analyzer.
- Compiler efficiency is improved. specialized buffering techniques for reading input characters can speed up the compiler significantly.
- Compiler portability is enhanced. Input-device-specific peculiarities can be restricted to the lexical analyzer.

Input Buffering(Lexical errors)

It is difficult to look one or more characters beyond the next lexeme before c o n f o r m the right lexeme.

There are many situations where we need to look at least one additional character ahead. For instance, we cannot be sure we've seen the end of an identifier until we see a character that is not a letter or digit, and therefore is not part of the lexeme for **id**.

In C, single-character operators like **-**, **=**, or **<** could also be the beginning of a two-character operator like **->**, **==**, or **<=**. Thus, we shall introduce a two-buffer scheme that handles large look aheads safely.

Buffer Pairs

Because of the amount of time taken to process characters and the large number of characters that must be processed during the compilation of a large source program, specialized buffering techniques have been developed to reduce the amount of overhead required to process a single input character. An important scheme involves two buffers that are alternately reloaded.

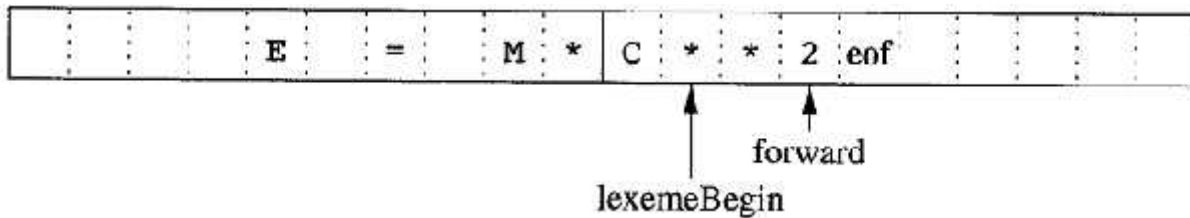


Fig: Using a Pair of Input Buffers

Two pointers to the input are maintained:

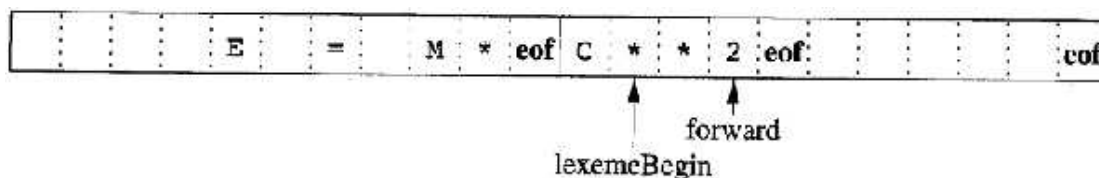
1. Pointer **lexemeBegin**, marks the beginning of the current lexeme
2. Pointer **forward** scans ahead until a pattern match is found.

Once the next lexeme is determined, forward is set to the character at its right end. Then, after the lexeme is recorded as an attribute value of a token returned to the parser, lexemeBegin is set to the character immediately after the lexeme just found. In Fig, forward has passed the end of the next lexeme, ** (the FORTRAN exponentiation operator), and must be retracted one position to its left.

Advancing forward requires that first test whether reached the end of one of the buffers, and if so, must reload the other buffer from the input, and move forward to the beginning of the newly loaded buffer.

Sentinels

for each character read, we make two tests: one for the end of the buffer, and one to determine what character is read. We can combine the buffer-end test with the test for the current character if we extend each buffer to hold a sentinel character at the end. The sentinel is a special character that cannot be part of the source program, and a natural choice is the character eof. Note that eof retains its use as a marker for the end of the entire input. Any eof that appears other than at the end of a buffer means that the input is at an end.



Panic Mode Error Recovery

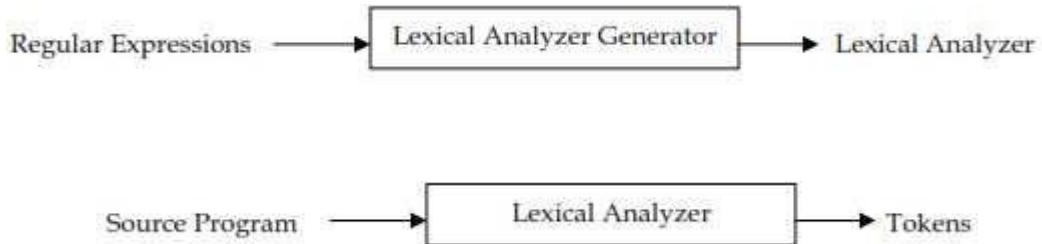
Suppose a situation arises in which a lexical analyzer is unable to proceed because none of the patterns for the token matches any prefix of the remaining input. It detects successive characters from the remaining input until the lexical analyzer can find a well

defined token at the beginning of the i/p.

The other possible error recovery actions are:-

- i) Delete 1 character from the remaining i/p.
- ii) Insert a missing character to the remaining i/p.
- iii) Replace a character by another character.
- iv) Transpose two adjacent characters.

Lecture #7
Lexical Analyzer Generator



Tokens

LEX is an example of Lexical Analyzer Generator.

Input to LEX

- The input to LEX consists primarily of *Auxiliary Definitions* and *Translation Rules*.
- To write regular expression for some languages can be difficult, because their regular expressions can be quite complex. In those cases, we may use *Auxiliary Definitions*.
- We can give names to regular expressions, and we can use these names as symbols to define other regular expressions.
- An *Auxiliary Definition* is a sequence of the definitions of the form:

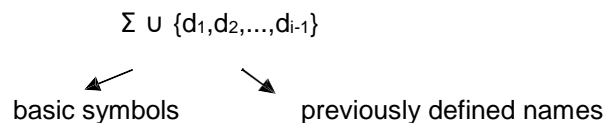
$d_1 \rightarrow r_1$

$d_2 \rightarrow r_2$

⋮

$d_n \rightarrow r_n$

where d_i is a distinct name and r_i is a regular expression over symbols in



Example:

For Identifiers in Pascal

letter $\rightarrow A | B | \dots | Z | a | b | \dots | z$ digit $\rightarrow 0 | 1 | \dots | 9$

id \rightarrow letter (letter | digit) *

If we try to write the regular expression representing identifiers without using regular definitions, that regular expression will be complex.

$(A|...|Z|a|...|z) ((A|...|Z|a|...|z) | (0|...|9)) ^ *$

Example:

For Unsigned numbers in Pascal digit $\rightarrow 0 | 1 | \dots | 9$ digits \rightarrow digit $^ +$

opt-fraction $\rightarrow (. \text{ digits }) ?$

opt-exponent $\rightarrow (E (+|-)? \text{ digits }) ?$

unsigned-num \rightarrow digits opt-fraction opt-exponent

- *Translation Rules* comprise of a ordered list Regular Expressions and the Program Code to be executed in case of that Regular Expression encountered.

| | |
|----------------|----------------|
| R ₁ | P ₁ |
| R ₂ | P ₂ |
| . | |
| R _n | P _n |

- The list is ordered i.e. the RE's should be checked in order. If a string matches more than one RE, the RE occurring higher in the list should be given preference and its Program Code is executed.

Implementation of LEX

- The Regular Expressions are converted into NFA's. The final states of each NFA correspond to some RE and its Program Code.
- Different NFA's are then converted to a single NFA with epsilon moves. Each final state of the NFA corresponds one-to-one to some final state of individual NFA's i.e. some RE and its Program Code. The final states have an order according to the corresponding RE's. If more than one final state is entered for some string, then the one that is higher in order is selected.
- This NFA is then converted to DFA. Each final state of DFA corresponds to a set of states (having at least one final state) of the NFA. The Program Code of each final state (of the DFA) is the program code corresponding to the final state that is highest in order out of all the final states in the set of states (of NFA) that make up this final state (of DFA).

Example:

AUXILIARY DEFINITIONS

(none)

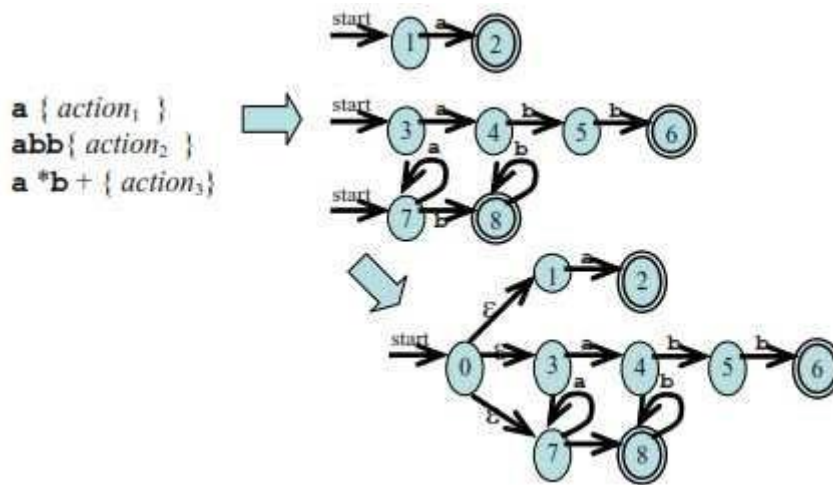
TRANSLATION RULES

a {Action₁}

abb{Action₂}

a*b⁺{Action₂}

First we construct an NFA for each RE and then convert this into a single NFA:



This NFA is now converted into a DFA. The transition table for the above DFA is as follows:

| State | A | b | Token found |
|-------|-----|----|-------------|
| 0137 | 247 | 8 | None |
| 247 | 7 | 58 | a |
| 8 | - | 8 | a^*b^+ |
| 7 | 7 | 8 | None |
| 58 | - | 68 | a^*b^+ |
| 68 | - | 8 | abb |

Lecture #8

BASICS OF SYNTAX ANALYSIS

- *Syntax Analyzer* creates the syntactic structure of the given source program.
- This syntactic structure is mostly a *parse tree*.
- Syntax Analyzer is also known as *parser*.
- The syntax of a programming is described by a *context-free grammar (CFG)*. We will use BNF (Backus-Naur Form) notation in the description of CFGs.
- The syntax analyzer (parser) checks whether a given source program satisfies the rules implied by a context-free grammar or not.
 - If it satisfies, the parser creates the parse tree of that program.
 - Otherwise the parser gives the error messages.
- A context-free grammar
 - gives a precise syntactic specification of a programming language.
 - the design of the grammar is an initial phase of the design of a compiler.
 - a grammar can be directly converted into a parser by some tools.

Parser

- Parser works on a stream of tokens.
- The smallest item is a token.



- We categorize the parsers into two groups:
- Top-Down Parser
 - The parse tree is created top to bottom, starting from the root.
- Bottom-Up Parser
 - The parse is created bottom to top; starting from the leaves
- Both top-down and bottom-up parsers scan the input from left to right (one symbol at a time).
- Efficient top-down and bottom-up parsers can be implemented only for sub-classes of context-free grammars.
 - LL for top-down parsing
 - LR for bottom-up parsing

Lecture #9 Context Free Grammars

Inherently recursive structures of a programming language are defined by a context-free grammar.

In a context-free grammar, we have:

- A finite set of terminals (in our case, this will be the set of tokens)
- A finite set of non-terminals (syntactic-variables)
- A finite set of productions rules in the following form

$A \rightarrow \alpha$ where A is a non-terminal and

α is a string of terminals and non-terminals (including the empty string)

- A start symbol (one of the non-terminal symbol)
- $L(G)$ is *the language of G* (the language generated by G) which is a set of sentences.
- A *sentence of $L(G)$* is a string of terminal symbols of G .
- If S is the start symbol of G then
 - (a) ω is a sentence of $L(G)$ iff $S \Rightarrow \omega$ where ω is a string of terminals of G .
- If G is a context-free grammar, $L(G)$ is a *context-free language*.
- Two grammars are *equivalent* if they produce the same language.
- $S \Rightarrow \alpha$
 - If α contains non-terminals, it is called as a *sentential form* of G .
 - If α does not contain non-terminals, it is called as a *sentence* of G .

Derivations

Example:

(a) $E \rightarrow E + E \mid E - E \mid E * E \mid E / E \mid - E$

(b) $E \rightarrow (E)$

(c) $E \rightarrow id$

- $E \Rightarrow E+E$ means that $E+E$ derives from E
 - we can replace E by $E+E$
 - to able to do this, we have to have a production rule $E \rightarrow E+E$ in our grammar.
- $E \Rightarrow E+E \Rightarrow id+E \Rightarrow id+id$ means that a sequence of replacements of non-terminal symbols
- In general a derivation step is
 $\alpha A \beta \Rightarrow \alpha \gamma \beta$ if there is a production rule $A \rightarrow \gamma$ in our grammar

Where α and β are arbitrary strings of terminal and non-terminal Symbol

$$\alpha_1 \Rightarrow \alpha_2 \Rightarrow \dots \Rightarrow \alpha_n \quad (\alpha_n \text{ derives from } \alpha_1 \text{ or } \alpha_1 \text{ derives } \alpha_n)$$

- At each derivation step, we can choose any of the non-terminal in the sentential form of G for the replacement.
- If we always choose the left-most non-terminal in each derivation step, this derivation is called as left-most derivation.

Example:

$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(id+E) \Rightarrow -(id+id)$$

- If we always choose the right-most non-terminal in each derivation step, this derivation is called as right-most derivation.

Example:

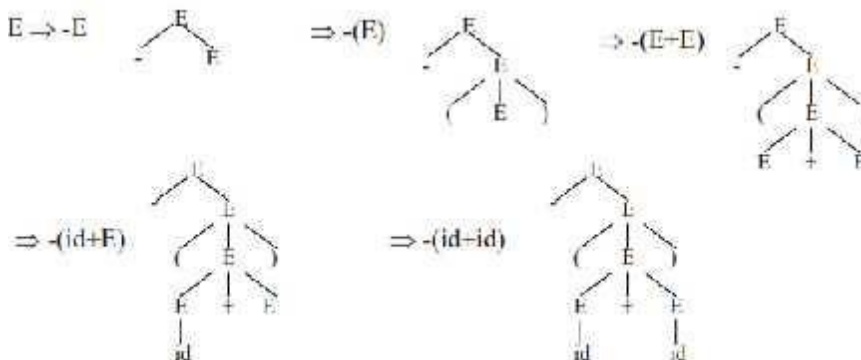
$$E \Rightarrow -E \Rightarrow -(E) \Rightarrow -(E+E) \Rightarrow -(E+id) \Rightarrow -(id+id)$$

- We will see that the top-down parsers try to find the left-most derivation of the given source program.
- We will see that the bottom-up parsers try to find the right-most derivation of the given source program in the reverse order.

Parse Tree

- Inner nodes of a parse tree are non-terminal symbols.
- The leaves of a parse tree are terminal symbols.
- A parse tree can be seen as a graphical representation of a derivation.

Example:



Ambiguity

- A grammar produces more than one parse tree for a sentence is called as an *ambiguous* grammar.
- For the most parsers, the grammar must be unambiguous.
- Unambiguous grammar
- Unique selection of the parse tree for a sentence
- We should eliminate the ambiguity in the grammar during the design phase of the compiler.
- An unambiguous grammar should be written to eliminate the ambiguity.
- We have to prefer one of the parse trees of a sentence (generated by an ambiguous grammar) to disambiguate that grammar to restrict to this choice.
- Ambiguous grammars (because of ambiguous operators) can be disambiguated according to the precedence and associativity rules.

Example:

To disambiguate the grammar

$$E \rightarrow E+E \mid E^*E \mid E^{\wedge}E \mid \text{id} \mid (E),$$

we can use precedence of operators as follows:

\wedge (right to left)

$*$ (left to right)

$+$ (left to right)

We get the following unambiguous grammar: $E \rightarrow E+T \mid T$

$$T \rightarrow T^*F \mid F$$

$$F \rightarrow G^{\wedge}F \mid G$$

$$G \rightarrow \text{id} \mid (E)$$

Lecture #10

Left Recursion

- A grammar is *left recursive* if it has a non-terminal A such that there is a derivation: $A \Rightarrow A\alpha$ for some string α
- Top-down parsing techniques cannot handle left-recursive grammars.
- So, we have to convert our left-recursive grammar into an equivalent grammar which is not left-recursive.
- The left-recursion may appear in a single step of the derivation (*immediate left-recursion*), or may appear in more than one step of the derivation.

Immediate Left-Recursion

$A \rightarrow A\alpha \mid \beta$ where β does not start with A
↓ Eliminate immediate left recursion
 $A \rightarrow \beta A'$
 $A' \rightarrow \alpha A' \mid \epsilon$ an equivalent grammar

In general,
 $A \rightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$ where $\beta_1 \dots \beta_n$ do not start with A
↓ Eliminate immediate left recursion
 $A \rightarrow \beta_1 A' \mid \dots \mid \beta_n A'$
 $A' \rightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \epsilon$ an equivalent grammar

Example:

$E \rightarrow E+T \mid T \quad T \rightarrow T^*F \mid F \quad F \rightarrow \text{id} \mid (E)$
↓ Eliminate immediate left recursion
 $E \rightarrow TE'$
 $E' \rightarrow +TE' \mid \epsilon$
 $T \rightarrow FT'$
 $T' \rightarrow *FT' \mid \epsilon$
 $F \rightarrow \text{id} \mid (E)$

•A grammar cannot be immediately left-recursive, but it still can be left-recursive.

•By just eliminating the immediate left-recursion, we may not get a grammar which is not left-recursive.

Example:

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Sc \mid d$$

This grammar is not immediately left-recursive, but it is still left-recursive.

$$\underline{S} \Rightarrow Aa \Rightarrow \underline{S}ca$$

Or

$$\underline{A} \Rightarrow Sc \Rightarrow \underline{A}ac$$

causes to a left-recursion

- So, we have to eliminate all left-recursions from our grammar.

Elimination

Arrange non-terminals in some order: $A_1 \dots A_n$

```
for i from 1 to n do {
  for j from 1 to i-1 do {
    replace each production
       $A_i \rightarrow A_j \gamma$ 
    by
       $A_i \rightarrow \alpha_1 \gamma \mid \dots \mid \alpha_k \gamma$ 
      where  $A_j \rightarrow \alpha_1 \mid \dots \mid \alpha_k$ 
  }
  eliminate immediate left-recursions among  $A_i$  productions
}
```

Example:

$$S \rightarrow Aa \mid b$$

$$A \rightarrow Ac \mid Sd \mid f$$

Case 1: Order of non-terminals: S, A

for S:

- we do not enter the inner loop.
- there is no immediate left recursion in S.

for A:

- Replace $A \rightarrow Sd$ with $A \rightarrow Aad \mid bd$

So, we will have $A \rightarrow Ac \mid Aad \mid bd \mid f$

- Eliminate the immediate left-recursion in

$A \rightarrow bdA' \mid fA'$

$A' \rightarrow cA' \mid adA' \mid \epsilon$

So, the resulting equivalent grammar which is not left-recursive is:

$S \rightarrow Aa \mid b$

$A \rightarrow bdA' \mid fA'$

$A' \rightarrow cA' \mid adA' \mid \epsilon$

Case 2: Order of non-terminals: A, S

for A:

- we do not enter the inner loop.
- Eliminate the immediate left-recursion in A

$A \rightarrow SdA' \mid fA'$

$A' \rightarrow cA' \mid \epsilon$

for S:

- Replace $S \rightarrow Aa$ with $S \rightarrow SdA'a \mid fA'a$

So, we will have $S \rightarrow SdA'a \mid fA'a \mid b$

- Eliminate the immediate left-recursion in S

$S \rightarrow fA'aS' \mid bS'$

$$S' \rightarrow dA'aS' \mid \epsilon$$

So, the resulting equivalent grammar which is not left-recursive is: $S \rightarrow fA'aS' \mid bSS' \rightarrow dA'aS' \mid \epsilon$

$$A \rightarrow SdA' \mid fA'$$

$$A' \rightarrow cA' \mid \epsilon$$

Left Factoring

- A predictive parser (a top-down parser without backtracking) insists that the grammar must be *left-factored*.

Grammar → a new equivalent grammar suitable for predictive parsing

stmt → if expr then stmt else stmt | if expr then stmt

- when we see if, we cannot now which production rule to choose to re-write *stmt* in the derivation
- In general,

$$A \rightarrow \beta\alpha_1 \mid \beta\alpha_2$$

where α is non-empty and the first symbols of β_1 and β_2 (if they have one) are different.

- when processing α we cannot know whether expand

$$A \text{ to } \beta\alpha_1 \text{ or}$$

$$A \text{ to } \beta\alpha_2$$

- But, if we re-write the grammar as follows

$$A \rightarrow \alpha A'$$

$$A' \rightarrow \beta_1 \mid \beta_2 \text{ so, we can immediately expand } A \text{ to } \alpha A'$$

10.1 Algorithm

- For each non-terminal A with two or more alternatives (production rules) with a common non-empty prefix, let say

$$A \rightarrow \beta\alpha_1 \mid \dots \mid \beta\alpha_n \mid \gamma_1 \mid \dots \mid \gamma_m \text{ convert it into}$$

$$A \rightarrow \alpha A' \mid \gamma_1 \mid \dots \mid \gamma_m$$

$$A' \rightarrow \beta_1 \mid \dots \mid \beta_n$$

Example:

$A \rightarrow \underline{a}bB \mid \underline{a}B \mid cdg \mid cdeB \mid cdfB$

↓

$A \rightarrow aA' \mid \underline{cdg} \mid \underline{cde}B \mid \underline{cdf}B$

$A' \rightarrow bB \mid B$

↓

$A \rightarrow aA' \mid cdA''$

$A' \rightarrow bB \mid B$

$A'' \rightarrow g \mid eB \mid fB$

Example:

$A \rightarrow ad \mid a \mid ab \mid abc \mid b$

↓

$A \rightarrow aA' \mid b$

$A' \rightarrow d \mid \varepsilon \mid b \mid bc$

↓

$A \rightarrow aA' \mid b$

$A' \rightarrow d \mid \varepsilon \mid bA''$

$A'' \rightarrow \varepsilon \mid c$

Lecture #11

YACC

YACC generates C code for a syntax analyzer, or parser. YACC uses grammar rules that allow it to analyze tokens from LEX and create a syntax tree. A syntax tree imposes a hierarchical structure on tokens. For example, operator precedence and associativity are apparent in the syntax tree. The next step, code generation, does a depth-first walk of the syntax tree to generate code. Some compilers produce machine code, while others output assembly.

YACC takes a default action when there is a conflict. For shift-reduce conflicts, YACC will shift. For reduce-reduce conflicts, it will use the first rule in the listing. It also issues a warning message whenever a conflict exists. The warnings may be suppressed by making the grammar unambiguous.

```
... definitions ...  
%%  
... rules ...  
%%  
... subroutines ...
```

Input to YACC is divided into three sections. The definitions section consists of token declarations, and C code bracketed by “%{“ and “%}”. The BNF grammar is placed in the rules section, and user subroutines are added in the subroutines section.

Lecture #12

TOP-DOWN PARSING

- The parse tree is created top to bottom.
- Top-down parser
 - Recursive-Descent Parsing
- Backtracking is needed (If a choice of a production rule does not work, we backtrack to try other alternatives.)
- It is a general parsing technique, but not widely used.
- Not efficient
 - Predictive Parsing
- No backtracking
- Efficient
- Needs a special form of grammars i.e. LL (1) grammars.
- Recursive Predictive Parsing is a special form of Recursive Descent parsing without backtracking.
- Non-Recursive (Table Driven) Predictive Parser is also known as LL (1) parser.

Recursive-Descent Parsing (uses Backtracking)

- Backtracking is needed.
- It tries to find the left-most derivation.

Example:

If the grammar is $S \rightarrow aBc$; $B \rightarrow bc \mid b$ and the input is abc:



Predictive Parser

| | | | |
|---------|-----------------------------|----------------|--|
| Grammar | -----> | -----> | |
| | Eliminate left recursion | Left Factor | a grammar suitable for predictive parsing (a LL(1) grammar) no %100 guarantee. |

- When re-writing a non-terminal in a derivation step, a predictive parser can uniquely choose a production rule by just looking the current symbol in the input string.

Example:

```

stmt → if ..... |
      while ..... |
      begin ..... |
      for .....
  
```

When we are trying to write the non-terminal *stmt*, we have to choose first production rule.

When we are trying to write the non-terminal *stmt*, we can uniquely choose the production rule by just looking the current token.

We eliminate the left recursion in the grammar, and left factor it. But it may not be suitable for predictive parsing (not LL (1) grammar).


```
    }  
proc C {
```

```
}
```

```
    match the current token with f and move to the next token;
```

```
proc B {
```

```
    case of the current token {
```

```
    b:      - match the current token with b and move to the next token;
```

```
    - call B
```

```
    e,d:   - do nothing                //Follow Set of B
```

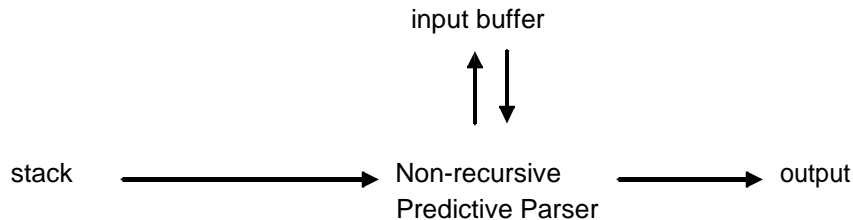
```
    }
```

```
    }
```

Lecture #14

Non-Recursive Predictive Parsing – LL (1) Parser

- Non-Recursive predictive parsing is a table-driven parser.
- It is a top-down parser.
- It is also known as LL(1) Parser.



Parsing Table input buffer

- our string to be parsed. We will assume that its end is marked with a special symbol \$.
- output
- a production rule representing a step of the derivation sequence (left-most derivation) of the string in the input buffer.

stack

- contains the grammar symbols
- at the bottom of the stack, there is a special end marker symbol \$.
- initially the stack contains only the symbol \$ and the starting symbol S. (\$S<-initial stack)
- when the stack is emptied (i.e. only \$ left in the stack), the parsing is completed.

parsing table

- a two-dimensional array $M[A,a]$
- each row is a non-terminal symbol
- each column is a terminal symbol or the special symbol \$
- each entry holds a production rule.

Parser Actions

The symbol at the top of the stack (say X) and the current symbol in the input string (say a) determine the parser action. There are four possible parser actions.

- If X and a are \$- > parser halts (successful completion)
- If X and a are the same terminal symbol (different from \$)
->parser pops X from the stack, and moves the next symbol in the input buffer.
- If X is a non-terminal

- parser looks at the parsing table entry $M[X,a]$. If $M[X,a]$ holds a production rule $X \rightarrow Y_1 Y_2 \dots Y_k$, it pops X from the stack and pushes Y_k, Y_{k-1}, \dots, Y_1 into the stack. The parser also outputs the production rule $X \rightarrow Y_1 Y_2 \dots Y_k$ to represent a step of the derivation.
- None of the above- > error
 - All empty entries in the parsing table are errors.
 - If X is a terminal symbol different from a , this is also an error case.

Example:

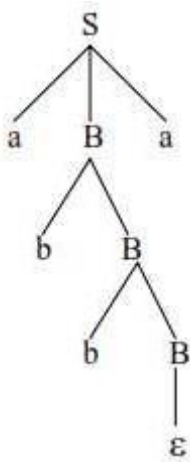
For the Grammar is $S \rightarrow aBa$; $B \rightarrow bB \mid \epsilon$ and the following LL(1) parsing table:

| | A | B | \$ |
|---|--------------------------|--------------------|----|
| S | $S \rightarrow aBa$ | | |
| B | $B \rightarrow \epsilon$ | $B \rightarrow bB$ | |

| <i>stack</i> | <i>input</i> | <i>output</i> |
|--------------|--------------|------------------------------------|
| $\$S$ | abba\$ | $S \rightarrow aBa$ |
| $\$aBa$ | abba\$ | |
| $\$aB$ | bba\$ | $B \rightarrow bB$ |
| $\$aBb$ | bba\$ | |
| $\$aB$ | ba\$ | $B \rightarrow bB$ |
| $\$aBb$ | ba\$ | |
| $\$aB$ | a\$ | $B \rightarrow \epsilon$ |
| $\$a$ | | a\$ |
| $\$$ | | $\$$ accept, successful completion |

Outputs: $S \rightarrow aBa$ $B \rightarrow bB$ $B \rightarrow bB$ $B \rightarrow \epsilon$

Derivation (left-most): $S \Rightarrow aBa \Rightarrow abBa \Rightarrow abbBa \Rightarrow abba$



Constructing LL(1) parsing tables

- Two functions are used in the construction of LL(1) parsing tables -FIRST & FOLLOW
- $FIRST(\alpha)$ is a set of the terminal symbols which occur as first symbols in strings derived from α where α is any string of grammar symbols.
- if α derives to ϵ , then ϵ is also in $FIRST(\alpha)$.
- $FOLLOW(A)$ is the set of the terminals which occur immediately after (follow) the *non-terminal* A in the strings derived from the starting symbol.
 - A terminal a is in $FOLLOW(A)$ if $S \Rightarrow \alpha A a \beta$
 - $\$$ is in $FOLLOW(A)$ if $S \Rightarrow \alpha A$

To Compute FIRST for Any String X:

- If X is a terminal symbol $\rightarrow FIRST(X) = \{X\}$
- If X is a non-terminal symbol and $X \rightarrow \epsilon$ is a production rule $\rightarrow \epsilon$ is in $FIRST(X)$.
- If X is a non-terminal symbol and $X \rightarrow Y_1 Y_2 \dots Y_n$ is a production rule
 - \rightarrow if a terminal a in $FIRST(Y_i)$ and ϵ is in all $FIRST(Y_j)$ for $j=1, \dots, i-1$ then a is in $FIRST(X)$.
 - \rightarrow if ϵ is in all $FIRST(Y_j)$ for $j=1, \dots, n$ then ϵ is in $FIRST(X)$.
- If X is $\epsilon \rightarrow FIRST(X) = \{\epsilon\}$
- If X is $Y_1 Y_2 \dots Y_n$
 - \rightarrow if a terminal a in $FIRST(Y_i)$ and ϵ is in all $FIRST(Y_j)$ for $j=1, \dots, i-1$ then a is in $FIRST(X)$.
 - \rightarrow if ϵ is in all $FIRST(Y_j)$ for $j=1, \dots, n$ then ϵ is in $FIRST(X)$.

To Compute FOLLOW (for non-terminals):

- If S is the start symbol $\rightarrow \$$ is in $FOLLOW(S)$
- If $A \rightarrow \alpha B \beta$ is a production rule \rightarrow everything in $FIRST(\beta)$ is $FOLLOW(B)$ except ϵ
- If ($A \rightarrow \alpha B$ is a production rule) or ($A \rightarrow \alpha B \beta$ is a production rule and ϵ is in $FIRST(\beta)$)
 - \rightarrow everything in $FOLLOW(A)$ is in $FOLLOW(B)$.
- Apply these rules until nothing more can be added to any follow set.

Algorithm for Constructing LL(1) Parsing Table:

- for each production rule $A \rightarrow \alpha$ of a grammar G
 - for each terminal a in $FIRST(\alpha)$ \rightarrow add $A \rightarrow \alpha$ to $M[A, a]$
 - If ϵ in $FIRST(\alpha)$ \rightarrow for each terminal a in $FOLLOW(A)$ add $A \rightarrow \alpha$ to $M[A, a]$
 - If ϵ in $FIRST(\alpha)$ and $\$$ in $FOLLOW(A)$ \rightarrow add $A \rightarrow \alpha$ to $M[A, \$]$
- All other undefined entries of the parsing table are error entries.

Example:

$$\begin{aligned} E &\rightarrow TE' \\ E' &\rightarrow +TE' \mid \epsilon \\ T &\rightarrow FT' \\ T' &\rightarrow *FT' \mid \epsilon \\ F &\rightarrow (E) \mid id \end{aligned}$$

$FIRST(F) = \{ (, id \}$
 $FIRST(T') = \{ *, \epsilon \}$
 $FIRST(T) = \{ (, id \}$
 $FIRST(E') = \{ +, \epsilon \}$
 $FIRST(E) = \{ (, id \}$
 $FIRST(TE') = \{ (, id \}$
 $FIRST(+TE') = \{ + \}$
 $FIRST(\epsilon) = \{ \epsilon \}$
 $FIRST(FT') = \{ (, id \}$
 $FIRST(*FT') = \{ * \}$
 $FIRST((E)) = \{ \}$
 $FIRST(id) = \{ id \}$
 $FOLLOW(E) = \{ \$,) \}$
 $FOLLOW(E') = \{ \$,) \}$
 $FOLLOW(T) = \{ +,), \$ \}$
 $FOLLOW(T') = \{ +,), \$ \}$
 $FOLLOW(F) = \{ +, *,), \$ \}$

LL(1) Parsing Table

| | | |
|---------------------------|---|---|
| $E \rightarrow TE'$ | $FIRST(TE') = \{ (, id \}$ | -> $E \rightarrow TE'$ into $M[E, (]$ and $M[E, id]$ |
| $E' \rightarrow +TE'$ | $FIRST(+TE') = \{ + \}$ | -> $E' \rightarrow +TE'$ into $M[E', +]$ |
| $E' \rightarrow \epsilon$ | $FIRST(\epsilon) = \{ \epsilon \}$ | -> none |
| | but since ϵ in $FIRST(\epsilon)$ and $FOLLOW(E') = \{ \$,) \}$ | -> $E' \rightarrow \epsilon$ into $M[E', \$]$ and $M[E',)]$ |
| $T \rightarrow FT'$ | $FIRST(FT') = \{ (, id \}$ | -> $T \rightarrow FT'$ into $M[T, (]$ and $M[T, id]$ |
| $T' \rightarrow *FT'$ | $FIRST(*FT') = \{ * \}$ | -> $T' \rightarrow *FT'$ into $M[T', *]$ |
| $T' \rightarrow \epsilon$ | $FIRST(\epsilon) = \{ \epsilon \}$ | -> none |
| | but since ϵ in $FIRST(\epsilon)$ and $FOLLOW(T') = \{ \$, + \}$ | -> $T' \rightarrow \epsilon$ into $M[T', \$]$, $M[T', +]$ and $M[T',)]$ |
| $F \rightarrow (E)$ | $FIRST((E)) = \{ \}$ | -> $F \rightarrow (E)$ into $M[F, (]$ |
| $F \rightarrow id$ | $FIRST(id) = \{ id \}$ | -> $F \rightarrow id$ into $M[F, id]$ |

| | id | + | * | (|) | \$ |
|----|---------------------|---------------------------|-----------------------|---------------------|---------------------------|---------------------------|
| E | $E \rightarrow TE'$ | | | $E \rightarrow TE'$ | | |
| E' | | $E' \rightarrow +TE'$ | | | $E' \rightarrow \epsilon$ | $E' \rightarrow \epsilon$ |
| T | $T \rightarrow FT'$ | | | $T \rightarrow FT'$ | | |
| T' | | $T' \rightarrow \epsilon$ | $T' \rightarrow *FT'$ | | $T' \rightarrow \epsilon$ | $T' \rightarrow \epsilon$ |
| F | $F \rightarrow id$ | | | $F \rightarrow (E)$ | | |

Lecture #15

LL(1) Grammars

- A grammar whose parsing table has no multiply-defined entries is said to be LL(1) grammar.
- The parsing table of a grammar may contain more than one production rule. In this case, we say that it is not a LL(1) grammar.
- A grammar G is LL(1) if and only if the following conditions hold for two distinctive production rules $A \rightarrow \alpha$ and $A \rightarrow \beta$:
 1. Both α and β cannot derive strings starting with same terminals.
 2. At most one of α and β can derive to ϵ .
 3. If β can derive to ϵ , then α cannot derive to any string starting with a terminal in $\text{FOLLOW}(A)$.

Non- LL(1) Grammars

Example:

$S \rightarrow iCtSE \mid a$

$E \rightarrow eS \mid \epsilon$

$C \rightarrow b$

$\text{FOLLOW}(S) = \{ \$, e \}$

$\text{FOLLOW}(E) = \{ \$, e \}$

$\text{FOLLOW}(C) = \{ t \}$

$\text{FIRST}(iCtSE) = \{ i \}$

$\text{FIRST}(a) = \{ a \}$

$\text{FIRST}(eS) = \{ e \}$

$\text{FIRST}(\epsilon) = \{ \epsilon \}$

$\text{FIRST}(b) = \{ b \}$

| | a | b | e | i | t | \$ |
|---|-------------------|-------------------|--|-----------------------|---|--------------------------|
| S | $S \rightarrow a$ | | | $S \rightarrow iCtSE$ | | |
| E | | | $E \rightarrow eS$ $E \rightarrow \epsilon$ | | | $E \rightarrow \epsilon$ |
| C | | $C \rightarrow b$ | | | | |

two production rules for $M[E,e]$

The Problem with multiple entries here is that of Ambiguity.

- What do we have to do if the resulting parsing table contains multiply defined entries?
 - If we didn't eliminate left recursion, eliminate the left recursion in the grammar.
 - If the grammar is not left factored, we have to left factor the grammar.
 - If its (new grammar's) parsing table still contains multiply defined entries, that grammar is ambiguous or it is inherently not a LL(1) grammar.
- A left recursive grammar cannot be a LL(1) grammar.
 - $A \rightarrow A\alpha \mid \beta$
 - >any terminal that appears in $\text{FIRST}(\beta)$ also appears $\text{FIRST}(A\alpha)$ because $A\alpha \Rightarrow \beta\alpha$.
 - >If β is ϵ , any terminal that appears in $\text{FIRST}(\alpha)$ also appears in $\text{FIRST}(A\alpha)$ and

FOLLOW(A).

- A grammar is not left factored, it cannot be a LL(1) grammar
 - $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$
 - >any terminal that appears in FIRST($\alpha\beta_1$) also appears in FIRST($\alpha\beta_2$).
- An ambiguous grammar cannot be a LL(1) grammar.

Lecture #16

BASIC BOTTOM-UP PARSING TECHNIQUES

- A bottom-up parser creates the parse tree of the given input starting from leaves towards the root.
- A bottom-up parser tries to find the right-most derivation of the given input in the reverse order.
 - (a) $S \Rightarrow \dots \Rightarrow \omega$ (the right-most derivation of ω)
 - (b) \leftarrow (the bottom-up parser finds the right-most derivation in the reverse order)
- Bottom-up parsing is also known as shift-reduce parsing because its two main actions are shift and reduce.
 - At each shift action, the current symbol in the input string is pushed to a stack.
 - At each reduction step, the symbols at the top of the stack (this symbol sequence is the right side of a production) will be replaced by the non-terminal at the left side of that production.
 - There are also two more actions: accept and error.

Shift-Reduce Parsing

- A shift-reduce parser tries to reduce the given input string into the starting symbol.
- At each reduction step, a substring of the input matching to the right side of a production rule is replaced by the non-terminal at the left side of that production rule.
- If the substring is chosen correctly, the right most derivation of that string is created in the reverse order.

Example:

For Grammar $S \rightarrow aABb$; $A \rightarrow aA \mid a$; $B \rightarrow bB \mid b$ and Input string $aaabb$,

```
aaabb
⇒ aaAbb
⇒ aAbb
⇒ aABb
⇒ S
```

The above reduction corresponds to the following rightmost derivation: $S \Rightarrow aABb \Rightarrow aAbb \Rightarrow aaAbb \Rightarrow aaabb$

Handle

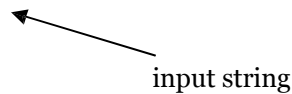
- Informally, a handle of a string is a substring that matches the right side of a production rule.
 - But not every substring that matches the right side of a production rule is handle.
- A handle of a right sentential form $\gamma (\equiv \alpha\beta\omega)$ is a production rule $A \rightarrow \beta$ and a position of γ where the string β may be found and replaced by A to produce the previous right-sentential form in a

rightmost derivation of γ .

$$S \Rightarrow \alpha A \omega \Rightarrow \alpha \beta \omega$$

- If the grammar is unambiguous, then every right-sentential form of the grammar has exactly one handle.
- We will see that ω is a string of terminals.
- A right-most derivation in reverse can be obtained by handle-pruning.

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n = \omega$$



- Start from γ_n , find a handle $A_n \rightarrow \beta_n$ in γ_n , and replace β_n in by A_n to get γ_{n-1} .
- Then find a handle $A_{n-1} \rightarrow \beta_{n-1}$ in γ_{n-1} , and replace β_{n-1} in by A_{n-1} to get γ_{n-2} .
- Repeat this, until we reach S .

Example:

$$\begin{aligned} E &\rightarrow E+T \mid T \\ T &\rightarrow T^*F \mid F \\ F &\rightarrow (E) \mid id \end{aligned}$$

Right-Most Derivation of $id+id*id$ is

$$E \Rightarrow E+T \Rightarrow E+T^*F \Rightarrow E+T^*id \Rightarrow E+F^*id \Rightarrow E+id^*id \Rightarrow T+id^*id \Rightarrow F+id^*id \Rightarrow id+id^*id$$

Right-Most Sentential Form

$$F \rightarrow id$$

$$\underline{F}+id^*id$$

$$E \rightarrow T \underline{E+id^*id}$$

$$T \rightarrow F \underline{E+T^*id}$$

$$T \rightarrow T^*F \underline{E+T}$$

Reducing Production $id+id^*id$

$$T \rightarrow F \underline{T+id^*id}$$

$$F \rightarrow id \underline{E+F^*id}$$

$$F \rightarrow id \underline{E+T^*F}$$

$$E \rightarrow E+T \underline{E}$$

Handles are underlined in the right-sentential forms.

Stack Implementation

- There are four possible actions of a shift-parser action:
- Shift : The next input symbol is shifted onto the top of the stack.
- Reduce: Replace the handle on the top of the stack by the non-terminal.
- Accept: Successful completion of parsing.
- Error: Parser discovers a syntax error, and calls an error recovery routine.
- Initial stack just contains only the end-marker \$.

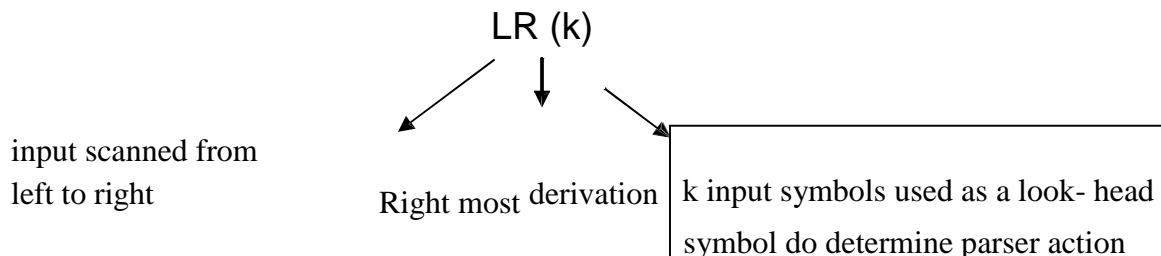
- The end of the input string is marked by the end-marker \$.

Example:

| <u>Stack</u> | <u>Input</u> | <u>Action</u> |
|--------------|--------------|-------------------|
| \$ | id+id*id\$ | shift |
| \$id | +id*id\$ | reduce by F → id |
| \$F | +id*id\$ | reduce by T → F |
| \$T | +id*id\$ | reduce by E → T |
| \$E | +id*id\$ | shift shift |
| \$E+ | id*id\$ | |
| \$E+id | *id\$ | reduce by F → id |
| \$E+F | *id\$ | reduce by T → F |
| \$E+T | *id\$ | shift shift |
| \$E+T* | id\$ | |
| \$E+T*id | \$ | reduce by F → id |
| \$E+T*F | \$ | reduce by T → T*F |
| \$E+T | \$ | reduce by E → E+T |
| \$E | \$ | accept |

Lecture #17
Conflicts during Shift Reduce Parsing

- There are context-free grammars for which shift-reduce parsers cannot be used.
- Stack contents and the next input symbol may not decide action:
 - shift/reduce conflict: Whether make a shift operation or a reduction.
 - reduce/reduce conflict: The parser cannot decide which of several reductions to make.
- If a shift-reduce parser cannot be used for a grammar, that grammar is called as non-LR(k) grammar.



- An ambiguous grammar can never be a LR grammar.

Types of Shift Reduce Parsing

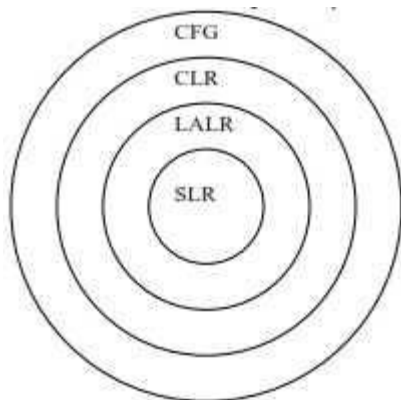
There are two main categories of shift-reduce parsers

1. Operator-Precedence Parser

- simple, but only a small class of grammars.

2. LR-Parsers

- Covers wide range of grammars.
- SLR – Simple LR parser
- CLR – most general LR parser (Canonical LR)
- LALR – intermediate LR parser (Look Ahead LR)
- SLR, CLR and LALR work same, only their parsing tables are different.



Lecture #18

Operator Precedence Parsing

- Operator grammar
 - small, but an important class of grammars
 - we may have an efficient operator precedence parser (a shift-reduce parser) for an operator grammar.
- In an *operator grammar*, no production rule can have:
 - ϵ at the right side
 - two adjacent non-terminals at the right side.

Examples:

| | | |
|----------------------|-----------------------|--------------------------|
| $E \rightarrow AB$ | $E \rightarrow EOE$ | $E \rightarrow E+E \mid$ |
| $A \rightarrow a$ | $E \rightarrow id$ | $E^*E \mid$ |
| $B \rightarrow b$ | $O \rightarrow + * /$ | $E/E \mid id$ |
| not operator grammar | not operator grammar | operator grammar |

Precedence Relations

- In operator-precedence parsing, we define three disjoint precedence relations between certain pairs of terminals.

$a < b$ b has higher precedence than a $a = b$ b has same precedence as a
 $a > b$ b has lower precedence than a

- The determination of correct precedence relations between terminals are based on the traditional notions of associativity and precedence of operators. (Unary minus causes a problem).
- The intention of the precedence relations is to find the handle of a right-sentential form,
 - < with marking the left end,
 - = appearing in the interior of the handle, and
 - > marking the right hand.
- In our input string $\$a_1a_2\dots a_n\$$, we insert the precedence relation between the pairs of terminals (the precedence relation holds between the terminals in that pair).

Example:

$E \rightarrow E+E \mid E-E \mid E^*E \mid E/E \mid E^{\wedge}E \mid (E) \mid -E \mid id$

The partial operator-precedence table for this grammar is as shown.

| | | | | |
|----|----|----|----|----|
| | id | + | * | \$ |
| id | | .> | .> | .> |
| + | <. | .> | <. | .> |
| * | <. | .> | .> | .> |
| \$ | <. | <. | <. | |

Then the input string $id+id*id$ with the precedence relations inserted will be:

$\$ < . id .> + < . id .> * < . id .> \$$

Using Precedence relations to find Handles

- Scan the string from left end until the first \rightarrow is encountered.
- Then scan backwards (to the left) over any $=$ until a $<$ is encountered.
- The handle contains everything to left of the first \rightarrow and to the right of the $<$ is encountered.

The handles thus obtained can be used to shift reduce a given string.

Operator-Precedence Parsing Algorithm

- The input string is $w\$$, the initial stack is $\$$ and a table holds precedence relations between certain terminals

Parsing Algorithm

The input string is $w\$$, the initial stack is $\$$ and a table holds precedence relations between certain terminals.

```
set p to point to the first symbol of w$ ;
repeat forever
if ( $ is on top of the stack and p points to $ ) then return else {
let a be the topmost terminal symbol on the stack and let b be the symbol pointed to by p;
if ( a <. b or a =. b )then { /* SHIFT */
push b onto the stack;
advance p to the next input symbol;
}
else if ( a .> b) then /* REDUCE */
repeat pop stack
until ( the top of stack terminal is related by <. to the terminal most recently popped);
else error();
}
```

Example:

| <i>stack</i> | <i>input</i> | <i>action</i> |
|--------------|--------------|------------------------|
| \$ | id+id*id\$ | \$ <. id shift |
| \$id | +id*id\$ | id >. + reduce E → id |
| \$ | +id*id\$ | shift |
| \$+ | id*id\$ | shift |
| \$+id | *id\$ | id >. * reduce E → id |
| \$+ | *id\$ | shift |
| \$+* | id\$ | shift |
| \$+*id | \$ | id >. \$ reduce E → id |
| \$+* | \$ | * >. \$ reduce E → E*E |

| | | | |
|-----|----|---------|----------------|
| \$+ | \$ | + :> \$ | reduce E → E+E |
| \$ | \$ | | accept |

Creating Operator-Precedence Relations from Associativity and Precedence

- If operator O1 has higher precedence than operator O2, $\rightarrow O1 \cdot > O2$ and $O2 < \cdot O1$
- If operator O1 and operator O2 have equal precedence, they are left-associative $\rightarrow O1 \cdot > O2$ and $O2 \cdot > O1$ they are right-associative $\rightarrow O1 < \cdot O2$ and $O2 < \cdot O1$
- For all operators O,
 $O < \cdot id$, $id \cdot > O$, $O < \cdot ($, $(< \cdot O$, $O \cdot >)$, $) \cdot > O$, $O \cdot > \$$, and $\$ < \cdot O$
- Also, let

| | | | |
|--------------|---------------|---------------|--------------|
| (=) | \$ < \cdot (| id \cdot >) |) \cdot > \$ |
| (< \cdot (| \$ < \cdot id | id \cdot > \$ |) \cdot >) |
| (< \cdot id | | | |

Example:

The complete table for the Grammar $E \rightarrow E+E \mid E-E \mid E * E \mid E / E \mid E ^ E \mid (E) \mid -E \mid id$ is:

| | + | - | * | / | ^ | id | (|) | \$ |
|----|---------|---------|---------|---------|---------|---------|---------|---------|---------|
| + | \cdot > | \cdot > | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | \cdot > | \cdot > |
| - | \cdot > | \cdot > | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | \cdot > | \cdot > |
| * | \cdot > | \cdot > | \cdot > | \cdot > | < \cdot | < \cdot | < \cdot | \cdot > | \cdot > |
| / | \cdot > | \cdot > | \cdot > | \cdot > | < \cdot | < \cdot | < \cdot | \cdot > | \cdot > |
| ^ | \cdot > | \cdot > | \cdot > | \cdot > | < \cdot | < \cdot | < \cdot | \cdot > | \cdot > |
| id | \cdot > | \cdot > | \cdot > | \cdot > | \cdot > | | | \cdot > | \cdot > |
| (| < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | = \cdot | |
|) | \cdot > | \cdot > | \cdot > | \cdot > | \cdot > | | | \cdot > | \cdot > |
| \$ | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | < \cdot | | |

Operator-Precedence Grammars

There is another more general way to compute precedence relations among terminals:

- $a = b$ if there is a right side of a production of the form $\alpha a \beta b \gamma$, where β is either a single non-terminal or ϵ .
- $a < b$ if for some non-terminal A there is a right side of the form $\alpha a A \beta$ and A derives to $\gamma b \delta$

where γ is a single non-terminal or ϵ .

3. $a > b$ if for some non-terminal A there is a right side of the form $\alpha A b \beta$ and A derives to $\gamma a \delta$ where δ is a single non-terminal or ϵ .

Note that the grammar must be unambiguous for this method. Unlike the previous method, it does not take into account any other property and is based purely on grammar productions. An ambiguous grammar will result in multiple entries in the table and thus cannot be used.

Handling Unary Minus

- Operator-Precedence parsing cannot handle the unary minus when we also use the binary minus in our grammar.
- The best approach to solve this problem is to let the lexical analyzer handle this problem.
 - The lexical analyzer will return two different operators for the unary minus and the binary minus.
 - The lexical analyzer will need a look ahead to distinguish the binary minus from the unary minus.
- Then, we make

| | |
|-------------------------|---|
| $O < \cdot$ unary-minus | for any operator |
| unary-minus $\cdot > O$ | if unary-minus has higher precedence than O |
| unary-minus $< \cdot O$ | if unary-minus has lower (or equal) |
| precedence than O | |

Precedence Functions

- Compilers using operator precedence parsers do not need to store the table of precedence relations.
- The table can be encoded by two precedence functions f and g that map terminal symbols to integers.
- For symbols a and b .

$f(a) < g(b)$ whenever $a < \cdot b$

$f(a) = g(b)$ whenever $a = \cdot b$

$f(a) > g(b)$ whenever $a > \cdot b$

Advantages and Disadvantages

- Advantages:
 - simple
 - powerful enough for expressions in programming languages
- Disadvantages:
 - It cannot handle the unary minus (the lexical analyzer should handle the unary minus).
 - Small class of grammars.
 - Difficult to decide which language is recognized by the grammar.

Lecture #19

LR PARSING

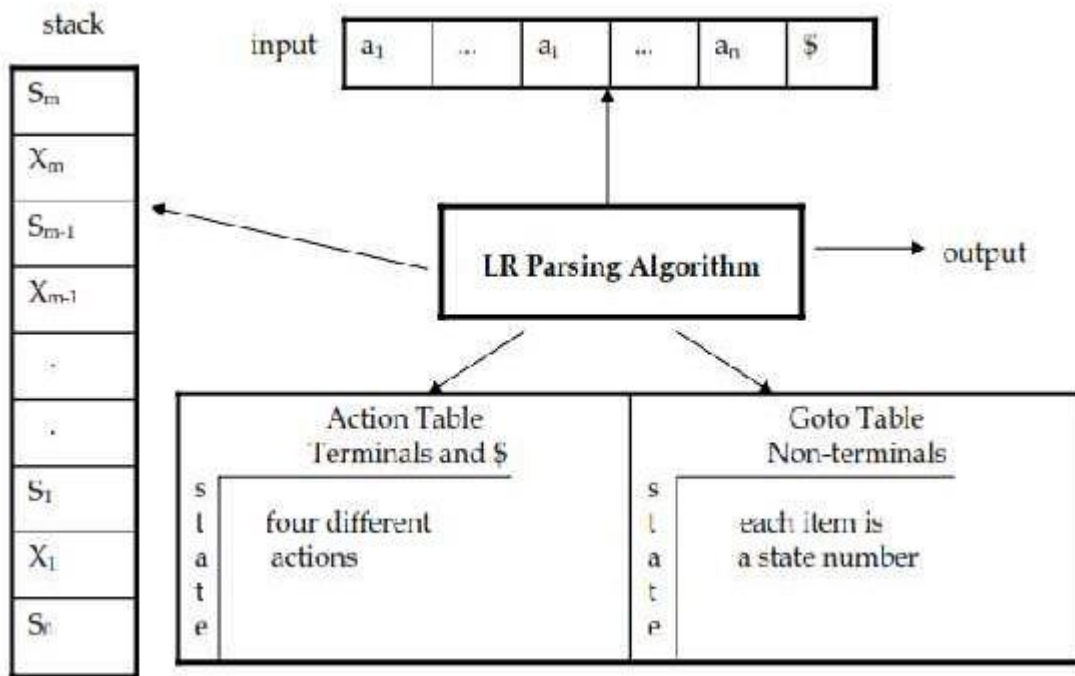
LR parsing is attractive because:

- LR parsing is most general non-backtracking shift-reduce parsing, yet it is still efficient.
- The class of grammars that can be parsed using LR methods is a proper superset of the class of grammars that can be parsed with predictive parsers.

$$LL(1)\text{-Grammars} \subset LR(1)\text{-Grammars}$$

- An LR-parser can detect a syntactic error as soon as it is possible to do so a left-to-right scan of the input.

Parser Configuration



• A configuration of a LR parsing is:

$$(\underbrace{S_0 X_1 S_1 \dots X_m S_m}_{\text{Stack}}, \underbrace{a_i a_{i+1} \dots a_n \$}_{\text{Rest of Input}})$$

• S_m and a_i decides the parser action by consulting the parsing action table. (*Initial Stack* contains just S_0)

• A configuration of a LR parsing represents the right sentential form:

$$X_1 \dots X_m a_i a_{i+1} \dots a_n \$$$

Parser Actions

1. **shift s** -- shifts the next input symbol and the state **s** onto the stack
 (So $X_1 S_1 \dots X_m S_m, a_i a_{i+1} \dots \text{an } \$$) \rightarrow (So $X_1 S_1 \dots X_m S_m a_i s, a_{i+1} \dots \text{an } \$$)

2. **reduce $A \rightarrow \beta$** (or **rn** where n is a production number)

- pop $2|\beta| (=r)$ items from the stack; let us assume that $\beta = Y_1 Y_2 \dots Y_r$
 - then push **A** and **s** where **$s = \text{goto}[s_{m-r}, A]$**

(So $X_1 S_1 \dots X_m S_m, a_i a_{i+1} \dots \text{an } \$$) \rightarrow (So $X_1 S_1 \dots X_{m-r} S_{m-r} A s, a_i \dots \text{an } \$$)

- Output is the reducing production reduce $A \rightarrow \beta$
- In fact, $Y_1 Y_2 \dots Y_r$ is a handle.

$X_1 \dots X_{m-r} A a_i \dots \text{an } \$ \Rightarrow X_1 \dots X_m Y_1 \dots Y_r a_i a_{i+1} \dots \text{an } \$$

3. **Accept** – Parsing successfully completed.

4. **Error** -- Parser detected an error (an empty entry in the action table) Example:
 Let following be the grammar and its LR parsing table.

1) $E \rightarrow E+T$

2) $E \rightarrow T$

3) $T \rightarrow T^*F$

4) $T \rightarrow F$

5) $F \rightarrow (E)$

6) $F \rightarrow \text{id}$

| state | Action | | | | | | Goto | | |
|-------|--------|----|----|----|-----|-----|------|---|----|
| | id | + | * | (|) | \$ | E | T | F |
| 0 | s5 | | | s4 | | | 1 | 2 | 3 |
| 1 | | s6 | | | | acc | | | |
| 2 | | r2 | s7 | | r2 | r2 | | | |
| 3 | | r4 | r4 | | r4 | r4 | | | |
| 4 | s5 | | | s4 | | | 8 | 2 | 3 |
| 5 | | r6 | r6 | | r6 | r6 | | | |
| 6 | s5 | | | s4 | | | | 9 | 3 |
| 7 | s5 | | | s4 | | | | | 10 |
| 8 | | s6 | | | s11 | | | | |
| 9 | | r1 | s7 | | r1 | r1 | | | |
| 10 | | r3 | r3 | | r3 | r3 | | | |
| 11 | | r5 | r5 | | r5 | r5 | | | |

The action of the parser would be as follows:

| <u>stack</u> | <u>input</u> | <u>action</u> | <u>output</u> |
|--------------|--------------|---------------------------------|-----------------------|
| 0 | id*id+id\$ | shift 5 | |
| 0id5 | *id+id\$ | reduce by $F \rightarrow id$ | $F \rightarrow id$ |
| 0F3 | *id+id\$ | reduce by $T \rightarrow F$ | $T \rightarrow F$ |
| 0T2 | *id+id\$ | | |
| 0T2*7 | id+id\$ | shift 7 | |
| 0T2*7id5 | +id\$ | reduce by $F \rightarrow id$ | $F \rightarrow id$ |
| 0T2*7F10 | +id\$ | reduce by $T \rightarrow T * F$ | $T \rightarrow T * F$ |
| 0T2 | +id\$ | reduce by $E \rightarrow T$ | $E \rightarrow T$ |
| 0E1 | +id\$ | | |
| 0E1+6 | id\$ | shift 6 | |
| 0E1+6id5 | \$ | reduce by $F \rightarrow id$ | $F \rightarrow id$ |
| 0E1+6F3 | \$ | reduce by $T \rightarrow F$ | $T \rightarrow F$ |
| 0E1+6T9 | \$ | reduce by $E \rightarrow E + T$ | $E \rightarrow E + T$ |
| 0E1 | \$ | | |

Lecture #20

Constructing SLR Parsing tables

- An LR parser using SLR parsing tables for a grammar G is called as the SLR parser for G .
- If a grammar G has an SLR parsing table, it is called SLR grammar.
- Every SLR grammar is unambiguous, but every unambiguous grammar is not a SLR grammar.
- *Augmented Grammar.* G' is G with a new production rule $S' \rightarrow S$ where S' is the new starting symbol.

LR(0) Items

- An **LR(0) item** of a grammar G is a production of G with a dot at some position of the right side.

Example:

$A \rightarrow aBb$

Possible LR(0) Items (four different possibilities):

$A \rightarrow .aBb$ $A \rightarrow a.Bb$ $A \rightarrow aB.b$ $A \rightarrow aBb.$

- Sets of LR(0) items will be the states of action and goto table of the SLR parser.
- A collection of sets of LR(0) items (**the canonical LR(0) collection**) is the basis for constructing SLR parsers.

Closure Operation

If I is a set of LR(0) items for a grammar G , then **closure(I)** is the set of LR(0) items constructed from I by the two rules:

1. Initially, every LR(0) item in I is added to closure(I).
2. If $A \rightarrow \alpha.B\beta$ is in closure(I) and $B\gamma \rightarrow$ is a production rule of G ; then $B \rightarrow .\gamma$ will be in the closure(I). We will apply this rule until no more new LR(0) items can be added to closure(I).

Example:

$E' \rightarrow E$; $E \rightarrow E+T$;

$E \rightarrow T$;

$T \rightarrow T^*F$; $T \rightarrow F$;

$F \rightarrow (E)$;

$F \rightarrow id$

$\text{closure}(\{E' \rightarrow .E\}) = \{ E' \rightarrow .E, E \rightarrow .E+T, E \rightarrow .T, T \rightarrow .T^*F, T \rightarrow .F, F \rightarrow .(E), F \rightarrow .id \}$

GOTO Operation

If I is a set of LR(0) items and X is a grammar symbol (terminal or non-terminal), then $\text{goto}(I, X)$ is defined as follows:

If $A \rightarrow \alpha.X\beta$ in I then every item in $\text{closure}(\{A \rightarrow \alpha X \beta\})$ will be in $\text{goto}(I, X)$.

Example:

$$I = \{ E' \rightarrow .E, E \rightarrow .E+T, E \rightarrow .T, T \rightarrow .T^*F, T \rightarrow .F, F \rightarrow .(E), F \rightarrow .id \}$$
$$\text{goto}(I, E) = \{ E' \rightarrow E., E \rightarrow E.+T \}$$
$$\text{goto}(I, T) = \{ E \rightarrow T., T \rightarrow T.*F \}$$
$$\text{goto}(I, F) = \{ T \rightarrow F. \}$$
$$\text{goto}(I, () = \{ F \rightarrow (.E), E \rightarrow .E+T, E \rightarrow .T, T \rightarrow .T^*F, T \rightarrow .F, F \rightarrow .(E), F \rightarrow .id \}$$
$$\text{goto}(I, id) = \{ F \rightarrow id. \}$$

Lecture #21

Construction of The Canonical LR(0) Collection

To create the SLR parsing tables for a grammar G , we will create the canonical LR(0) collection of the grammar G .

Algorithm:

C is { closure($\{S' \rightarrow S\}$) }

repeat the followings until no more set of LR(0) items can be added to C .

for each I in C and each grammar symbol X

if goto(I, X) is not empty and not in C

add goto(I, X) to C

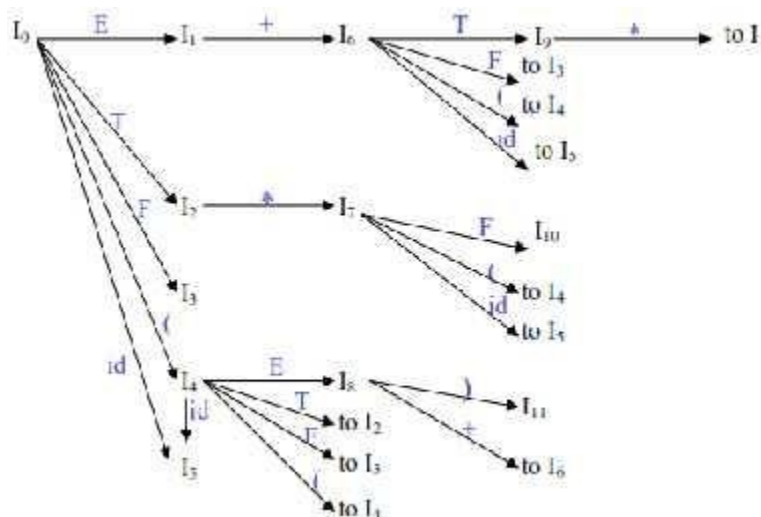
GOTO function is a DFA on the sets in C . Example:

For grammar used above, Canonical LR(0) items are as follows-

| | | | |
|-------------------------|--------------------------|--------------------------|---------------------------|
| I0: $E' \rightarrow .E$ | I1: $E' \rightarrow E.$ | I6: $E \rightarrow E+.T$ | I9: $E \rightarrow E+T.$ |
| $E \rightarrow .E+T$ | $E \rightarrow E.+T$ | $T \rightarrow .T*F$ | $T \rightarrow T.*F$ |
| $E \rightarrow .T$ | | $T \rightarrow .F$ | |
| $T \rightarrow .T*F$ | I2: $E \rightarrow T.$ | $F \rightarrow .(E)$ | I10: $T \rightarrow T*F.$ |
| $T \rightarrow .F$ | $T \rightarrow T.*F$ | $F \rightarrow .id$ | |
| $F \rightarrow .(E)$ | | | |
| $F \rightarrow .id$ | I3: $T \rightarrow F.$ | I7: $T \rightarrow T*.F$ | I11: $F \rightarrow (E).$ |
| | | $F \rightarrow .(E)$ | |
| | I4: $F \rightarrow (.E)$ | $F \rightarrow .id$ | |
| | $E \rightarrow .E+T$ | | |
| | $E \rightarrow .T$ | I8: $F \rightarrow (E.)$ | |
| | $T \rightarrow .T*F$ | $E \rightarrow E.+T$ | |
| | $T \rightarrow .F$ | | |
| | $F \rightarrow .(E)$ | | |
| | $F \rightarrow .id$ | | |

I5: $F \rightarrow id.$

Transition Diagram (DFA) of GOTO Function is as follows-



Parsing Table

1. Construct the canonical collection of sets of LR(0) items for G' .

$$C \leftarrow \{I_0, \dots, I_n\}$$

2. Create the parsing action table as follows

- a. If a is a terminal, $A\alpha \rightarrow a\beta$ in I_i and $\text{goto}(I_i, a) = I_j$ then $\text{action}[i, a]$ is **shift j** .
- b. If $A\alpha \rightarrow \cdot$ is in I_i , then $\text{action}[i, a]$ is **reduce $A\alpha \rightarrow$** for all a in $\text{FOLLOW}(A)$ where

$$A \neq S'$$

- c. If $S' \rightarrow S \cdot$ is in I_i , then $\text{action}[i, \$]$ is **accept**.
 - d. If any conflicting actions generated by these rules, the grammar is not SLR(1).
3. Create the parsing goto table
 - a. for all non-terminals A , if $\text{goto}(I_i, A) = I_j$ then $\text{goto}[i, A] = j$
 4. All entries not defined by (2) and (3) are errors.
 5. Initial state of the parser contains $S' \rightarrow \cdot S$

Example:

For the Grammar used above, SLR Parsing table is as follows:

| state | Action | | | | | | Goto | | |
|-------|--------|----|----|----|-----|-----|------|---|----|
| | id | + | * | (|) | \$ | E | T | F |
| 0 | s5 | | | s4 | | | 1 | 2 | 3 |
| 1 | | s6 | | | | acc | | | |
| 2 | | r2 | s7 | | r2 | r2 | | | |
| 3 | | r4 | r4 | | r4 | r4 | | | |
| 4 | s5 | | | s4 | | | 8 | 2 | 3 |
| 5 | | r6 | r6 | | r6 | r6 | | | |
| 6 | s5 | | | s4 | | | | 9 | 3 |
| 7 | s5 | | | s4 | | | | | 10 |
| 8 | | s6 | | | s11 | | | | |
| 9 | | r1 | s7 | | r1 | r1 | | | |
| 10 | | r3 | r3 | | r3 | r3 | | | |
| 11 | | r5 | r5 | | r5 | r5 | | | |

Lecture #22

Shift/reduce and reduce/reduce conflicts

- If a state does not know whether it will make a shift operation or reduction for a terminal, we say that there is a **shift/reduce conflict**.

Example:

| | | | | |
|-------------------------------|--------------------------------------|-------------------------------------|-------------------------------------|---------------------------|
| S → L=R S → R | I ₀ : S' → .S S → .L=R | I ₁ : S' → S. | I ₆ : S → L=.R R → .L | I ₉ : S → L=R. |
| L → *R L → id | S → .R L → .*R | I ₂ : S → L.=R R → L. | L → .*R L → .id | |
| R → L | L → .id | | | |
| | R → .L | I ₃ : S → R. | | |
| Problem in I ₂ | | I ₄ : L → *.R R → .L | I ₇ : L → *R. | |
| FOLLOW(R)={=,\$} = shift 6 | | L → .*R L → .id | I ₈ : R → L. | |
| & reduce by R → L | | | | |
| shift/reduce conflict | | I ₅ : L → id. | | |

- If a state does not know whether it will make a reduction operation using the production rule i or j for a terminal, we say that there is a **reduce/reduce conflict**.

Example:

| | |
|----------|--------------------------|
| S → AaAb | I ₀ : S' → .S |
| S → BbBa | S → .AaAb |
| A → ε | S → .BbBa |
| B → ε | A → . |
| B → . | |

Problem FOLLOW(A)={a,b} FOLLOW(B)={a,b}

| | | | |
|---|--|---|--|
| a | reduce by $A \rightarrow \epsilon$ reduce by $B \rightarrow \epsilon$ reduce/reduce conflict | b | reduce by $A \rightarrow \epsilon$ reduce by $B \rightarrow \epsilon$ reduce/reduce conflict |
|---|--|---|--|

If the SLR parsing table of a grammar G has a conflict, we say that that grammar is not SLR grammar.

Constructing Canonical LR(1) Parsing tables

- In SLR method, the state i makes a reduction by $A\alpha \rightarrow$ when the current token is a:
 - if the $A\alpha \rightarrow$ in the li and a is FOLLOW(A)
- In some situations, βA cannot be followed by the terminal a in a right-sentential form when $\alpha\beta$ and the state i are on the top stack. This means that making reduction in this case is not correct.

| | | |
|--------------------------|---|---|
| $S \rightarrow AaAb$ | $S \Rightarrow AaAb \Rightarrow Aab \Rightarrow ab$ | $S \Rightarrow BbBa \Rightarrow Bba \Rightarrow ba$ |
| $S \rightarrow BbBa$ | | |
| $A \rightarrow \epsilon$ | $Aab \Rightarrow \epsilon ab$ | $Bba \Rightarrow \epsilon ba$ |
| $B \rightarrow \epsilon$ | $AaAb \Rightarrow Aa \epsilon b$ | $BbBa \Rightarrow Bb \epsilon a$ |

LR(1) Item

- To avoid some of invalid reductions, the states need to carry more information.
- Extra information is put into a state by including a terminal symbol as a second component in an item.
- A LR(1) item is: item $A \rightarrow \alpha.\beta,a$ where a is the look-head of the LR(1) (a is a terminal or end-marker.)
- When β (in the LR(1) item $A \rightarrow \alpha.\beta,a$) is not empty, the look-head does not have any affect.
- When β is empty ($A \rightarrow \alpha.,a$), we do the reduction by $A\alpha \rightarrow$ only if the next input symbol is a (not for any terminal in FOLLOW(A)).
- A state will contain $A \rightarrow \alpha.,a_1$ where $\{a_1, \dots, a_n\} \subseteq \text{FOLLOW}(A)$
- ...
- $A \rightarrow \alpha.,a_n$

Closure and GOTO Operations

closure(I) is: (where I is a set of LR(1) items

every LR(1) item in I is in closure(I)

if $A\alpha \rightarrow \cdot B\beta, a$ in $\text{closure}(I)$ and $B\gamma \rightarrow$ is a production rule of G ; then $B \rightarrow \cdot \gamma, b$ will be in the $\text{closure}(I)$ for each terminal b in $\text{FIRST}(\beta a)$.

If I is a set of LR(1) items and X is a grammar symbol (terminal or non-terminal), then $\text{goto}(I, X)$ is defined as follows:

If $A \rightarrow \alpha \cdot X \beta, a$ in I then every item in $\text{closure}(\{A \rightarrow \alpha X \cdot \beta, a\})$ will be in $\text{goto}(I, X)$.

Lecture #23 Construction of The Canonical LR(1) Collection

Algorithm:

C is $\{ \text{closure}(\{S' \rightarrow \cdot S, \$\}) \}$

repeat the followings until no more set of LR(1) items can be added to C .

for each I in C and each grammar symbol X

if $\text{goto}(I, X)$ is not empty and not in C

add $\text{goto}(I, X)$ to C

GOTO function is a DFA on the sets in C .

A set of LR(1) items containing the following items

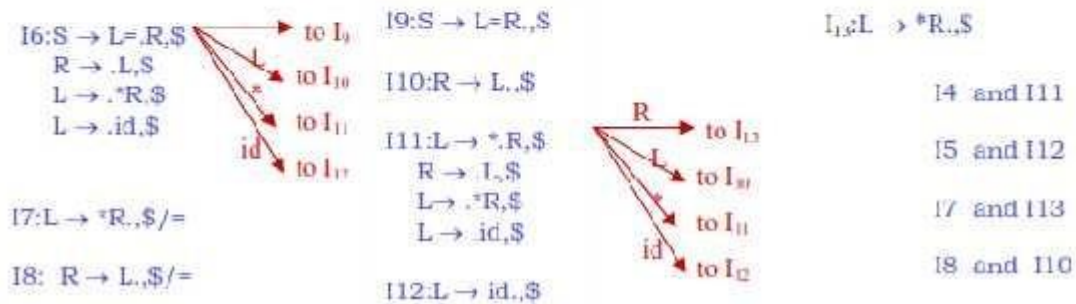
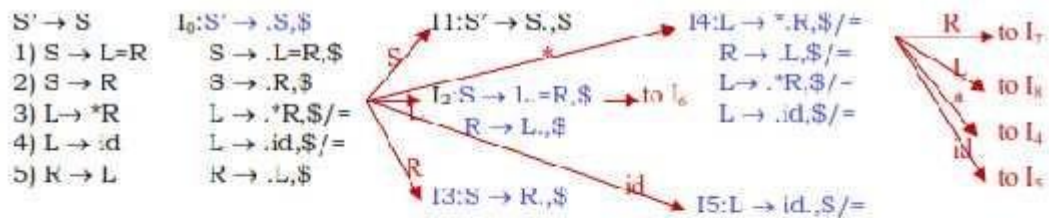
$A \rightarrow \alpha \cdot \beta, a_1$

...

$A \rightarrow \alpha \cdot \beta, a_n$

can be written as $A \rightarrow \alpha \cdot \beta, a_1/a_2/.../a_n$

Example:



Parsing Table

1. Construct the canonical collection of sets of LR(1) items for G' .

$C \leftarrow \{I_0, \dots, I_n\}$

2. Create the parsing action table as follows

- a. If a is a terminal, $A\alpha \rightarrow .a\beta$, b in li and $goto(li,a)=lj$ then $action[i,a]$ is *shift j*.
 - b. If $A\alpha \rightarrow .,a$ is in li , then $action[i,a]$ is *reduce $A\alpha \rightarrow$* where $A \neq S'$.
 - c. If $S' \rightarrow S., \$$ is in li , then $action[i, \$]$ is *accept*.
 - d. If any conflicting actions generated by these rules, the grammar is not LR(1).
3. Create the parsing goto table
 - a. for all non-terminals A , if $goto(li,A)=lj$ then $goto[i,A]=j$
 4. All entries not defined by (2) and (3) are errors.
 5. Initial state of the parser contains $S' \rightarrow .S, \$$

Example:

For the above used Grammar, the parse table is as follows:

| | id | * | = | \$ | S | L | R |
|----|-----|-----|----|-----|---|----|----|
| 0 | S5 | s4 | | | 1 | 2 | 3 |
| 1 | | | | acc | | | |
| 2 | | | s6 | r5 | | | |
| 3 | | | | r2 | | | |
| 4 | S5 | s4 | | | | 8 | 7 |
| 5 | | | r4 | r4 | | | |
| 6 | s12 | s11 | | | | 10 | 9 |
| 7 | | | r3 | r3 | | | |
| 8 | | | r5 | r5 | | | |
| 9 | | | | r1 | | | |
| 10 | | | | r5 | | | |
| 11 | s12 | s11 | | | | 10 | 13 |
| 12 | | | | r4 | | | |
| 13 | | | | r3 | | | |

->cores of $\text{goto}(I1,X), \dots, \text{goto}(I2,X)$ must be same.
 -So, $\text{goto}(J,X)=K$ where K is the union of all sets of items having same cores as $\text{goto}(I1,X)$.

- If no conflict is introduced, the grammar is LALR(1) grammar. (We may only introduce reduce/reduce conflicts; we cannot introduce a shift/reduce conflict)

Shift/Reduce Conflict

- We say that we cannot introduce a shift/reduce conflict during the shrink process for the creation of the states of a LALR parser.
- Assume that we can introduce a shift/reduce conflict. In this case, a state of LALR parser must have:

$A \rightarrow \alpha.,a$ and $B \rightarrow \beta.a\gamma,b$

- This means that a state of the canonical LR(1) parser must have: $A \rightarrow \alpha.,a$ and

$B \rightarrow \beta.a\gamma,c$

But, this state has also a shift/reduce conflict. i.e. The original canonical

LR(1) parser has a conflict.

(Reason for this, the shift operation does not depend on Lookaheads)

Reduce/Reduce Conflict

But, we may introduce a reduce/reduce conflict during the shrink process for the creation of the states of a LALR parser.

$I1 : A \rightarrow \alpha.,a$

$I2 : A \rightarrow \alpha.,b$

$B \rightarrow \beta.,b$

$B \rightarrow \beta.,c$

↓

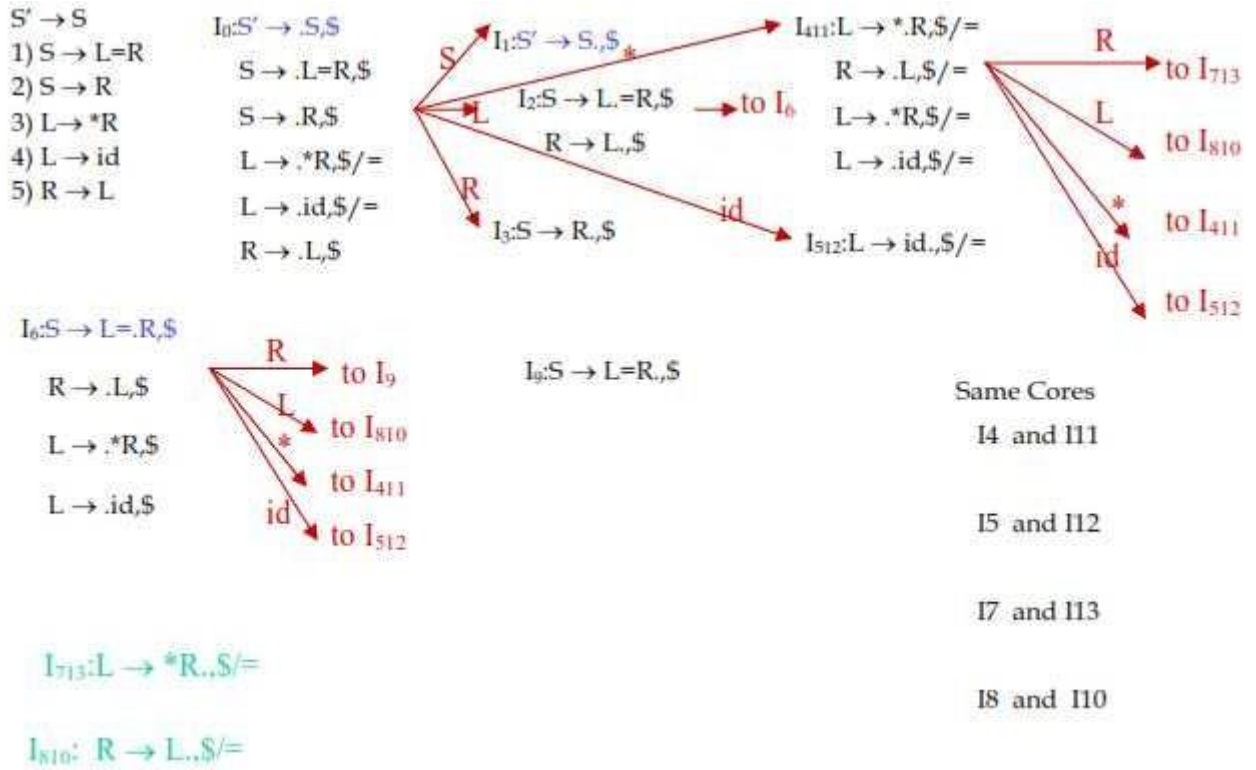
$I12 : A \rightarrow \alpha.,a/b$

- > reduce/reduce conflict

$B \rightarrow \beta.,b/c$

Example:

For the above Canonical LR Parsing table, we can get the following LALR(1) collection



Lecture #25

Using Ambiguous Grammars

- All grammars used in the construction of LR-parsing tables must be un-ambiguous.
- Can we create LR-parsing tables for ambiguous grammars?
 - Yes, but they will have conflicts.
 - We can resolve these conflicts in favor of one of them to disambiguate the grammar.
 - At the end, we will have again an unambiguous grammar.
- Why we want to use an ambiguous grammar?
 - Some of the ambiguous grammars are **much natural**, and a corresponding unambiguous grammar can be very complex.
 - Usage of an ambiguous grammar may **eliminate unnecessary reductions**.

Example:

$E \rightarrow E+T \mid T$

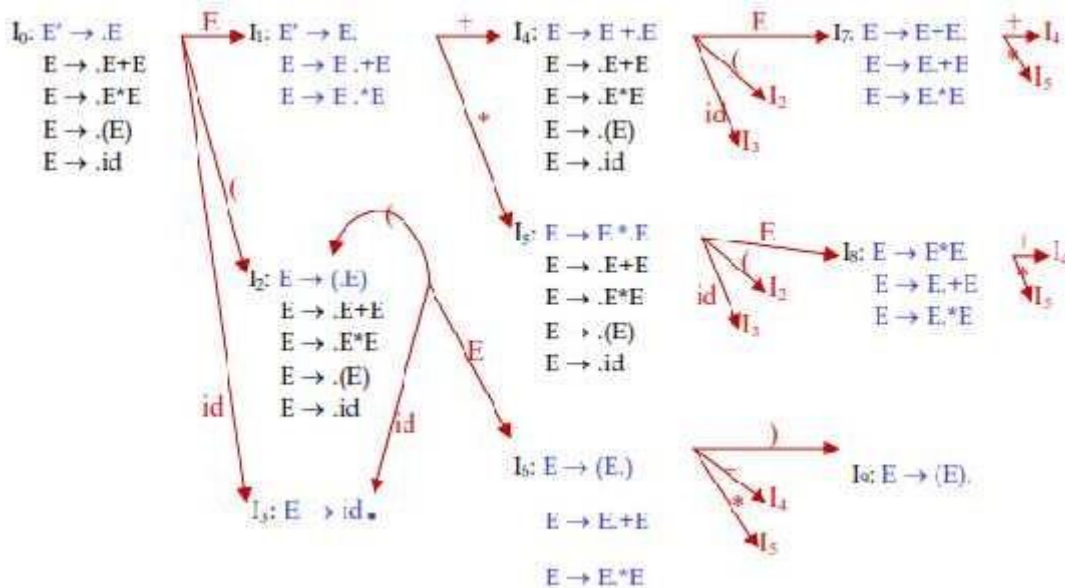
$E \rightarrow E+E \mid E^*E \mid$

$(E) \mid id \mid T \rightarrow T^*F$

$\mid F$

$F \rightarrow (E) \mid id$

Sets of LR(0) Items for Ambiguous Grammar



SLR-Parsing Tables for Ambiguous Grammar

$FOLLOW(E) = \{ \$, +, *,) \}$

State I7 has shift/reduce conflicts for symbols + and *.

when current token is +

shift $\rightarrow +$ is right-associative reduce $\rightarrow +$ is left-associative

when current token is *

shift $\rightarrow *$ has higher

precedence than + reduce->+
has higher precedence than *

State 18 has shift/reduce conflicts for symbols +

and *. when current token is *

shift ->* is right-associative reduce ->* is left-associative

when current token is +

shift -> + has higher

precedence than * reduce->*

has higher precedence than +

| | id | + | * | (|) | \$ | E |
|---|----|----|----|----|----|-----|---|
| 0 | s3 | | | s2 | | | 1 |
| 1 | | s4 | s5 | | | acc | |
| 2 | s3 | | | s2 | | | 6 |
| 3 | | r4 | r4 | | r4 | r4 | |
| 4 | s3 | | | s2 | | | 7 |
| 5 | s3 | | | s2 | | | 8 |
| 6 | | s4 | s5 | | s9 | | |
| 7 | | r1 | s5 | | r1 | r1 | |
| 8 | | r2 | r2 | | r2 | r2 | |
| 9 | | r3 | r3 | | r3 | r3 | |

Module-2

Lecture #26

SYNTAX-DIRECTED TRANSLATION

- Grammar symbols are associated with **attributes** to associate information with the programming language constructs that they represent.
- Values of these attributes are evaluated by the **semantic rules** associated with the production rules.
- Evaluation of these semantic rules:
 - may generate intermediate codes
 - may put information into the symbol table
 - may perform type checking
 - may issue error messages
 - may perform some other activities
 - In fact, they may perform almost any activities.
- An attribute may hold almost any thing.
 - A string, a number, a memory location, a complex record.
- Evaluation of a semantic rule defines the value of an attribute. But a semantic rule may also have some side effects such as printing a value.

Example:

| <u>Production</u> | <u>Semantic Rule</u> | <u>Program Fragment</u> |
|---------------------------------|-------------------------------------|-----------------------------------|
| $L \rightarrow E$ return | print(E.val) | print(val[top-1]) |
| $E \rightarrow E + T$ | E.val = E.val + T.val | val[ntop] = val[top-2] + val[top] |
| $E \rightarrow T$ | E.val = T.val | |
| $T \rightarrow T^1 * F$ | T.val = T ¹ .val * F.val | val[ntop] = val[top-2] * val[top] |
| $T \rightarrow F$ | T.val = F.val | |
| $F \rightarrow (E)$ | F.val = E.val | val[ntop] = val[top-1] |
| $F \rightarrow$ digit | F.val = digit .lexval | val[ntop] = digit.lexval |

- Symbols E, T, and F are associated with an attribute *val*.
- The token **digit** has an attribute *lexval* (it is assumed that it is evaluated by the lexical analyzer).
- The *Program Fragment* above represents the implementation of the semantic rule for a bottom-up parser.
 - At each shift of **digit**, we also push **digit.lexval** into *val-stack*.
 - At all other shifts, we do not put anything into *val-stack* because other terminals do not have attributes (but we increment the stack pointer for *val-stack*).
 - The above model is suited for a desk calculator where the purpose is to evaluate and to generate code.

1. Intermediate Code Generation

- *Intermediate codes* are machine independent codes, but they are close to machine instructions.
- The given program in a source language is converted to an equivalent program in an intermediate language by the intermediate code generator.
- Intermediate language can be many different languages, and the designer of the compiler decides this intermediate language.
 - syntax trees can be used as an intermediate language.
 - postfix notation can be used as an intermediate language.
 - three-address code (Quadruples) can be used as an intermediate language

- we will use quadruples to discuss intermediate code generation
- quadruples are close to machine instructions, but they are not actual machine instructions.

2 Syntax Tree

Syntax Tree is a variant of the Parse tree, where each leaf represents an operand and each interior node an operator.

Example:

Production

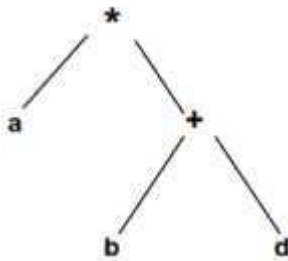
$E \rightarrow E1 \text{ op } E2$
 $E \rightarrow - E1$

Semantic Rule

$E.val = \text{NODE}(\text{op}, E1.val, E2.val) \quad E \rightarrow (E1)$
 $E.val = \text{UNARY}(-, E1.val) \quad E \rightarrow \text{id}$

$E.val = E1.val$
 $E.val = \text{LEAF}(\text{id})$

A sentence $a*(b+d)$ would have the following syntax tree:



Postfix Notation

Postfix Notation is another useful form of intermediate code if the language is mostly expressions.

Example:

Production

$E \rightarrow E1 \text{ op } E2$
 $E \rightarrow (E1) \quad E \rightarrow \text{id}$

Semantic Rule

$E.code = E1.code \parallel E2.code \parallel \text{op}$
 $E.code = E1.code$
 $E.code = \text{id}$

Program Fragment

print op print id

3 Three Address Code

• We use the term “three-address code” because each statement usually contains three addresses (two for operands, one for the result).

• The most general kind of three-address code is:

$x := y \text{ op } z$

where x , y and z are names, constants or compiler-generated temporaries; **op** is any operator.

• But we may also use the following notation for quadruples (much better notation because it looks like a machine code instruction)

op y,z,x

apply operator op to y and z , and store the result in x .

4 Representation of three-address codes

Three-address code can be represented in various forms viz. Quadruples, Triples and Indirect Triples. These forms are demonstrated by way of an example below. Example:
 $A = -B * (C + D)$ Three-Address code is as follows:

$T1 = -B$

$T2 = C + D$ $T3 = T1 * T2$

$A = T3$

Quadruple:

| | Operator | Operand 1 | Operand 2 | Result |
|-----|-----------------|------------------|------------------|---------------|
| (1) | - | B | | T1 |
| (2) | + | C | D | T2 |
| (3) | * | T1 | T2 | T3 |
| (4) | = | A | T3 | |

Triple:

| | Operator | Operand 1 | Operand 2 |
|-----|-----------------|------------------|------------------|
| (1) | - | B | |
| (2) | + | C | D |
| (3) | * | (1) | (2) |
| (4) | = | A | (3) |

Indirect Triple:

| | | Statement | | |
|------|-----------------|------------------|------------------|--|
| | (0) | (56) | | |
| | (1) | (57) | | |
| | (2) | (58) | | |
| | (3) | (59) | | |
| | Operator | Operand 1 | Operand 2 | |
| (56) | - | B | | |
| (57) | + | C | D | |
| (58) | * | (56) | (57) | |
| (59) | = | A | (58) | |

Lecture #27

Translation of Assignment Statements

A statement $A := - B * (C + D)$ has the following three-address translation: $T1 := - B$
 $T2 := C+D$
 $T3 := T1* T2$
 $A := T3$

| <u>Production</u> | <u>Semantic Action</u> |
|-------------------------|---|
| $S \rightarrow id := E$ | $S.code = E.code \parallel gen(id.place = E.place)$ |
| $E \rightarrow E1 + E2$ | $E.place = newtemp();$ $E.code = E1.code \parallel E2.code \parallel gen(E.place = E1.place + E2.place)$ |
| $E \rightarrow E1 * E2$ | $E.place = newtemp();$ $E.code = E1.code \parallel E2.code \parallel gen(E.place = E1.place * E2.place)$ |
| $E \rightarrow - E1$ | $E.place = newtemp();$ $E.code = E1.code \parallel gen(E.place = - E1.place)$ |
| $E \rightarrow (E1)$ | $E.place = E1.place; E.code = E1.code$ |
| $E \rightarrow id$ | $E.place = id.place; E.code = null$ |

1. Translation of Boolean Expressions

Grammar for Boolean Expressions is: $E \rightarrow E$ or E
 $E \rightarrow E$ and E
 $E \rightarrow \text{not } E$ $E \rightarrow (E)$ $E \rightarrow id$
 $E \rightarrow id \text{ relop } id$

There are two representations viz. Numerical and Control-Flow.

Numerical Representation of Boolean

- o TRUE is denoted by 1 and FALSE by 0.
- o Expressions are evaluated from left to right, in a manner similar to arithmetic expressions.

Example:

The translation for **A or B and C** is the three-address sequence:

$T1 := B \text{ and } C$ $T2 := A \text{ or } T1$

Also, the translation of a relational expression such as $A < B$ is the three-address sequence:

(1) if $A < B$ goto (4) (2) $T := 0$
(3) goto (5) (4) $T := 1$ (5)

Therefore, a Boolean expression $A < B$ or C can be translated as: (1) if $A < B$ goto (4)
 (2) $T1 := 0$
 (3) goto (5) (4) $T1 := 1$
 (5) $T2 := T1$ or C

| <u>Production</u> | <u>Semantic Action</u> |
|---------------------------------------|---|
| $E \rightarrow E1$ or $E2$ | $T = \text{newtemp} ();$ $E.\text{place} = T;$ $\text{Gen} (T = E1.\text{place} \text{ or } E2.\text{place})$ |
| $E \rightarrow E1$ and $E2$ | $T = \text{newtemp} (); E.\text{place} = T;$ $\text{Gen} (T = E1.\text{place} \text{ and } E2.\text{place})$ |
| $E \rightarrow \text{not } E1$ | $T = \text{newtemp} (); E.\text{place} = T;$ $\text{Gen} (T = \text{not } E1.\text{place})$ |
| $E \rightarrow (E1)$ | $E.\text{place} = E1.\text{place}; E.\text{code} = E1.\text{code}$ |
| $E \rightarrow \text{id}$ | $E.\text{place} = \text{id}.\text{place}; E.\text{code} = \text{null}$ |
| $E \rightarrow \text{id1 rel op id2}$ | $T = \text{newtemp} (); E.\text{place} = T;$ $\text{Gen} (\text{if id1.place rel op id2.place goto NEXTQUAD}+3) \text{Gen} (T = 0)$ $\text{Gen} (\text{goto NEXTQUAD}+2)$ $\text{Gen} (T = 1)$ |

o Quadruples are being generated and NEXTQUAD indicates the next available entry in the quadruple array.

Control-Flow Representation of Boolean Expressions

- o If we evaluate Boolean expressions by program position, we may be able to avoid evaluating the entire expressions.
- o In A or B , if we determine A to be true, we need not evaluate B and can declare the entire expression to be true.
- o In A and B , if we determine A to be false, we need not evaluate B and can declare the entire expression to be false.
- o A better code can thus be generated using the above properties.

Example:

The statement **if ($A < B$ || $C < D$) $x = y + z$;** can be translated as

(1) if $A < B$ goto (4) (2) if $C < D$ goto (4) (3) goto (6)
 (4) $T = y + z$
 (5) $X = T$ (6)

Here (4) is a true exit and (6) is a false exit of the Boolean expressions.

Lecture #28

Generating 3-address code for Numerical Representation of Boolean expressions

- o Consider a production **E→E1 or E2** that represents the OR Boolean expression. If E1 is true, we know that E is true so we make the location TRUE for E1 be the same as TRUE for E. If E1 is false, then we must evaluate E2, so we make FALSE for E1 be the first statement in the code for E2. The TRUE and FALSE exits can be made the same as the TRUE and FALSE exits of E, respectively.
- o Consider a production **E→E1 and E2** that represents the AND Boolean expression. If E1 is false, we know that E is false so we make the location FALSE for E1 be the same as FALSE for E. If E1 is true, then we must evaluate E2, so we make TRUE for E1 be the first statement in the code for E2. The TRUE and FALSE exits can be made the same as the TRUE and FALSE exits of E, respectively.
- o Consider the production **E→not E** that represents the NOT Boolean expression. We may simply interchange the TRUE and FALSE exits of E1 to get the TRUE and FALSE exits of E.
- o To generate quadruples in the manner suggested above, we use three functions- Makelist, Merge and Backpatch that shall work on the list of quadruples as suggested by their name.
- o If we need to proceed to E2 after evaluating E1, we have an efficient way of doing this by modifying our grammar as follows:

E→E or M E
E→E and M E E→not E
E→(E) E→id
E→id relop id
M→ε

The translation scheme for this grammar would as follows:

| <u>Production</u> | <u>Semantic Action</u> |
|-------------------|--|
| E→E1 or M E2 | BACKPATCH (E1.FALSE, M.QUAD); E.TRUE = MERGE (E1.TRUE, E2.TRUE); E.FALSE = E2.FALSE; |
| E→E1 and M E2 | BACKPATCH (E1.TRUE, M.QUAD); E.TRUE = E2.TRUE; E.FALSE = MERGE (E1.FALSE, E2.FALSE); |
| E→not E1 | E.TRUE = E1.FALSE; E.FALSE = E1.TRUE; |

| | |
|------------------|---|
| E->(E1) | E.TRUE = E1.TRUE; E.FALSE = E1.FALSE; |
| E->id | E.TRUE = MAKELIST (NEXTQUAD); E.FALSE = MAKELIST (NEXTQUAD + 1); GEN (if id.PLACE goto _); GEN (goto _); |
| E->id1 relop id2 | E.TRUE = MAKELIST (NEXTQUAD); E.FALSE = MAKELIST (NEXTQUAD + 1); GEN (if id1.PLACE relop id2.PLACE goto _); GEN (goto _); |
| M-> ϵ | M.QUAD = NEXTQUAD; |

Example:

For the expression $P < Q$ or $R < S$ and T , the parsing steps and corresponding semantic actions are shown below. We assume that NEXTQUAD has an initial value of 100.

Step 1: $P < Q$ gets reduced to E by $E \rightarrow id \text{ relop } id$. The grammatical form is $E1$ or $R < S$ and T .

We have the following code generated (Makelist).

```
100: if P<Q goto _
101: goto _
```

$E1$ is true if goto of 100 is reached and false if goto of 101 is reached.

Step 2: $R < S$ gets reduced to E by $E \rightarrow id \text{ relop } id$. The grammatical form is $E1$ or $E2$ and T .

We have the following code generated (Makelist).

```
102: if R<S goto _
103: goto _
```

$E2$ is true if goto of 102 is reached and false if goto of 103 is reached. Step 3: T gets reduced to E by $E \rightarrow id$. The grammatical form is $E1$ or $E2$ and $E3$.

We have the following code generated (Makelist).

```
104: if T goto _
105: goto _
```

$E3$ is true if goto of 104 is reached and false if goto of 105 is reached.

Step 4: $E2$ and $E3$ gets reduced to E by $E \rightarrow E \text{ and } E$. The grammatical form is $E1$ or $E4$.

We have no new code generated but changes are made in the already generated code (Backpatch).


```
100: if P<Q goto _
101: goto _
102: if R<S goto 104
103: goto _
104: if T goto _
105: goto _
```

E4 is true only if E3.TRUE (goto of 104) is reached. E4 is false if E2.FALSE (goto of 103) or E3.FALSE (goto of 105) is reached (Merge).

Step 5: E1 or E4 gets reduced to E by $E \rightarrow E$ or E. The grammatical form is E.

We have no new code generated but changes are made in the already generated code (Backpatch).

```
100: if P<Q goto _
101: goto 102
102: if R<S goto 104
103: goto _
104: if T goto _
105: goto _
```

E is true only if E1.TRUE (goto of 100) or E2.TRUE (goto of 104) is reached (Merge). E is false if E4.FALSE (goto of 103 or 105) is reached.

1. Mixed Mode Expressions

- o Boolean expressions may in practice contain arithmetic sub expressions e.g. $(A+B)>C$.
- o We can accommodate such sub-expressions by adding the production $E \rightarrow E \text{ op } E$ to our grammar.
- o We will also add a new field MODE for E. If E has been achieved after reduction using the above (arithmetic) production, we make $E.MODE = \text{arith}$, otherwise make $E.MODE = \text{bool}$.
- o If $E.MODE = \text{arith}$, we treat it arithmetically and use E.PLACE. If $E.MODE = \text{bool}$, we treat it as Boolean and use E.FALSE and E.TRUE.

Lecture #29

Statements that Alter Flow of Control

->In order to implement goto statements, we need to define a LABEL for a statement. A production can be added for this purpose:

S -> LABEL : S LABEL- > id

->The semantic action attached with this production is to record the LABEL and its value (NEXTQUAD) in the symbol table. It will also Backpatch any previous references to this LABEL with its current value.

->Following grammar can be used to incorporate structured Flow-of-control constructs: (1) S->if E then S

(2) S->if E then S else S

(3) S->while E do S (4) S->begin L end

(5) S->A

(6) L->L ; S (7) L->S

Here, S denotes a statement, L a statement-list, A an assignment statement and E a Boolean-valued expression.

1. Translation Scheme for statements that alter flow of control

->We introduce a new field NEXT for S and L like TRUE and FALSE for E. S.NEXT and L.NEXT are respectively the pointers to a list of all conditional and unconditional jumps to the quadruple following statement S and statement-list L in execution order.

->We also introduce the marker non-terminal M as in the case of grammar for Boolean expressions. This is put before statement in if-then, before both statements in if-then-else and the statement in while-do as we may need to proceed to them after evaluating E. In case of while-do, we also need to put M before E as we may need to come back to it after executing S.

->In case of if-then-else, if we evaluate E to be true, first S will be executed. After this we should ensure that instead of second S, the code after this if-then-else statement be executed. We thus place another non-terminal marker N after first S i.e. before else.

The grammar now is as follows:

(1) S->if E then M S

(2) S->if E then M S N else M S (3) S->while M E do M S

(4) S->begin L end

(5) S->A

(6) L->L ; M S (7) L->S

(8) M-> ϵ

(9) N-> ϵ

The translation scheme for this grammar would as follows:

ProductionSemantic Action

| | |
|--|---|
| S->if E then M S1 | BACKPATCH (E.TRUE, M.QUAD) S.NEXT = MERGE (E.FALSE, S1.NEXT) |
| S->if E then M1 S1 N else M2 S2 (E.FALSE, S.NEXT = MERGE | BACKPATCH (E.TRUE, M1.QUAD) BACKPATCH M2.QUAD) (S1.NEXT, N.NEXT, S2.NEXT) |
| S->while M1 E do M2 S1 (E.TRUE, | BACKPATCH (S1.NEXT, M1.QUAD) BACKPATCH M2.QUAD) S.NEXT = E.FALSE GEN (goto M1.QUAD) |
| S->begin L end | S.NEXT = L.NEXT |
| S->A | S.NEXT = MAKELIST () |
| L->L1 ; M S | BACKPATCH (L1.NEXT, M.QUAD) L.NEXT = S.NEXT |
| L->S | L.NEXT = S.NEXT |
| M-> ϵ | M.QUAD = NEXTQUAD |
| N-> ϵ | N.NEXT = MAKELIST (NEXTQUAD) GEN (goto _) |

Lecture #30

Postfix Translations

In an production $A \rightarrow \alpha$, the translation rule of A.CODE consists of the concatenation of the CODE translations of the non-terminals in α in the same order as the non-terminals appear in α . Productions can be factored to achieve Postfix form.

1. Postfix translation of while statement

The production

$S \rightarrow \text{while } M1 \text{ E do } M2 \text{ S1}$ can be factored as
 $S \rightarrow C \text{ S1}$
 $C \rightarrow W \text{ E do}$
 $W \rightarrow \text{while}$

A suitable translation scheme would be

| <u>Production</u> | <u>Semantic Action</u> |
|--------------------------------|--|
| $W \rightarrow \text{while}$ | $W.QUAD = \text{NEXTQUAD}$ |
| $C \rightarrow W \text{ E do}$ | $C.QUAD = W.QUAD$ $\text{BACKPATCH (E.TRUE, NEXTQUAD) C.FALSE} = \text{E.FALSE}$ |
| $S \rightarrow C \text{ S1}$ | $\text{BACKPATCH (S1.NEXT, C.QUAD) S.NEXT} = \text{C.FALSE}$ GEN (goto C.QUAD) |

2. Postfix translation of for statement

Consider the following production which stands for the for-statement

$S \rightarrow \text{for } L = E1 \text{ step } E2 \text{ to } E3 \text{ do } S1$

Here L is any expression with l-value, usually a variable, called the index. E1, E2 and E3 are expressions called the initial value, increment and limit, respectively. Semantically, the for statement is equivalent to the following program.

Begin

```
INDEX = addr ( L );  
*INDEX = E1; INCR = E2; LIMIT = E3;  
while *INDEX <= LIMIT do begin
```

```
end
```

```
end
```

```
code for statement S1;  
*INDEX = *INDEX + INCR;
```

The non-terminals L, E1, E2, E3 and S appear in the same order as in the production. The production can be factored as

- (1) F->for L
- (2) T->F = E1 by E2 to E3 do
- (3) S->T S1

A suitable translation scheme would be

Production

Semantic Action

F->for L

F.INDEX = L.INDEX

T->F = E1 by E2 to E3 do

GEN (*F.INDEX = E1.PLACE) INCR = NEWTEMP ()
LIMIT = NEWTEMP () GEN (INCR = E2.PLACE)
GEN (LIMIT = E3.PLACE)
T.QUAD = NEXTQUAD
T.NEXT = MAKELIST (NEXTQUAD) GEN (IF *F.INDEX >
LIMIT goto _) T.INDEX = F.INDEX
T.INCR = INCR

S->T S1

BACKPATCH (S1.NEXT, NEXTQUAD) GEN (*T.INDEX =
*T.INDEX + T.INCR) GEN (goto T.QUAD)

S.NEXT = T.NEXT

Lecture #31

Array references in arithmetic expressions


Elements of arrays can be accessed quickly if the elements are stored in a block of consecutive locations.

For a one-dimensional array A:

Base (A) is the address of the first location of the array A,
width is the width of each array element.
low is the index of the first array element

location of $A[i] = \text{baseA} + (i - \text{low}) * \text{width}$
 can be re-written as

$$i * \text{width} + (\text{baseA} - \text{low} * \text{width})$$



should be computed at run-time can be computed at compile-time

So, the location of $A[i]$ can be computed at the run-time by evaluating the formula $i * \text{width} + c$ where c is $(\text{baseA} - \text{low} * \text{width})$ which is evaluated at compile-time.

Intermediate code generator should produce the code to evaluate this formula $i * \text{width} + c$ (one multiplication and one addition operation).

A two-dimensional array can be stored in either row-major (row-by-row) or column-major (column-by-column).

Most of the programming languages use row-major method.

The location of $A[i_1, i_2]$ is $\text{baseA} + ((i_1 - \text{low}_1) * n_2 + i_2 - \text{low}_2) * \text{width}$

baseA is the location of the array A.
low₁ is the index of the first row
low₂ is the index of the first column
n₂ is the number of elements in each row
width is the width of each array element

Again, this formula can be re-written as

$$((i_1 * n_2) + i_2) * \text{width} + (\text{baseA} - ((\text{low}_1 * n_1) + \text{low}_2) * \text{width})$$



should be computed at run-time can be computed at compile-time

Arrays of any dimension can be dealt in a similar but general manner.

In general, the location of $A[i_1, i_2, \dots, i_k]$ is

$$((\dots ((i_1 * n_2) + i_2) \dots) * n_k + i_k) * \text{width} + (\text{baseA} - ((\dots ((\text{low}_1 * n_1) + \text{low}_2) \dots) * n_k + \text{low}_k) * \text{width})$$

So, the intermediate code generator should produce the codes to evaluate the following formula (to find the location of $A[i_1, i_2, \dots, i_k]$):

$$((\dots ((i_1 * n_2) + i_2) \dots) * n_k + i_k) * \text{width} + c$$

To evaluate the $((\dots ((i_1 * n_2) + i_2) \dots) * n_k + i_k)$ portion of this formula, we can use the recurrence equation:

$$e_1 = i_1$$

$$e_m = e_{m-1} * n_m + i_m$$

1. Grammar and Translation Scheme

The grammar and suitable translation scheme for arithmetic expressions with array references is as given below:

| <u>Production</u> | <u>Semantic Action</u> |
|---|--|
| $S \rightarrow L = E$ | if (L.OFFSET = NULL) then GEN (L.PLACE = E.PLACE) else GEN(L.PLACE [L.OFFSET] = E.PLACE) |
| $E \rightarrow E_1 + E_2$ | E.PLACE = NEWTEMP () GEN (E.PLACE = E1.PLACE + E2.PLACE) $E \rightarrow (E_1)$ E.PLACE = E1.PLACE |
| $E \rightarrow L$ | if (L.OFFSET = NULL) then E.PLACE = L.PLACE else {E.PLACE = NEWTEMP (); GEN (E.PLACE = L.PLACE[L.OFFSET])} |
| $L \rightarrow \text{id}$ | L.PLACE = id.PLACE L.OFFSET = NULL |
| $L \rightarrow \text{ELIST }]$ | L.PLACE = NEWTEMP () L.OFFSET = NEWTEMP () GEN (L.PLACE = ELIST.ARRAY - C) GEN (L.OFFSET = ELIST.PLACE * WIDTH (ELIST.ARRAY)) |
| $\text{ELIST} \rightarrow \text{ELIST1 } , E$ | ELIST.ARRAY = ELIST1.ARRAY ELIST.PLACE = NEWTEMP () ELIST.NDIM = ELIST1.NDIM + 1 GEN (ELIST.PLACE = ELIST1.PLACE * LIMIT (ELIST.ARRAY, ELIST.NDIM)) GEN (ELIST.PLACE = E.PLACE + ELIST.PLACE) |
| $\text{ELIST} \rightarrow \text{id } [E$ | ELIST.ARRAY = id.PLACE ELIST.PLACE = E.PLACE E = 1 |

Here, NDIM denotes the number of dimensions, LIMIT (ARRAY, i) function returns the upper limit along the i th dimension of ARRAY i.e. n_i , WIDTH (ARRAY) returns the number of bytes for one element of ARRAY.

2. Declarations

Following is the grammar and a suitable translation scheme for declaration statements:

| <u>Production</u> | <u>Semantic Action</u> |
|-------------------|--|
| D->integer, id | ENTER (id.PLACE, integer) D.ATTR = integer |
| D->real, id | ENTER (id.PLACE, real) D.ATTR = real |
| D->D1, id | ENTER (id.PLACE, D1.ATTR) D.ATTR = D1.ATTR |

Here, ENTER makes the entry into symbol table while ATTR is used to trace the data type.

3. Procedure Calls

Following is the grammar and a suitable translation scheme for Procedure Calls:

| <u>Production</u> | <u>Semantic Action</u> |
|-------------------------------------|---|
| S->call id (ELIST) GEN (param p) | for each item p on QUEUE do GEN (call id.PLACE) |
| ELIST->ELIST, E | append E.PLACE to the end of QUEUE |
| ELIST->E | initialize QUEUE to contain only E.PLACE QUEUE is used to store the list of parameters in the procedure call. |

4. Case Statements

The case statement has following syntax:

```
switch E
begin
end
case V1: S1 case V2: S2
.
case Vn-1: Sn-1 default: Sn The translation scheme for this shown below:
code to evaluate E into T
goto TEST L1: code for S1
goto NEXT L2: code for S2
goto NEXT
.
.
Ln-1: code for Sn-1 goto NEXT
Ln: code for Sn goto NEXT
TEST: if T = V1 goto L1
If T = V2 goto L2
.
if T = Vn-1 goto Ln-1 goto Ln
```


Lecture #32

SYMBOL

TABLES

- Symbol table is a data structure meant to collect information about names appearing in the source program.
- It keeps track about the scope/binding information about names.
- Each entry in the symbol table has a pair of the form (name and information).
- Information consists of attributes (e.g. type, location) depending on the language.
- Whenever a name is encountered, it is checked in the symbol table to see if already occurs. If not, a new entry is created.
- In some cases, the symbol table record is created by the lexical analyzer as soon as the name is encountered in the input, and the attributes of the name are entered when the declarations are processed.
- If same name can be used to denote different program elements in the same block, the symbol table record is created only when the name's syntactic role is discovered.

Operations on a Symbol Table

- Determine whether a given name is in the table
- Add a new name to the table
- Access information associated to a given name
- Add new information for a given name
- Delete a name (or a group of names) from the table

Implementation

- Each entry in a symbol table can be implemented as a record that consists of several fields.
- The entries in symbol table records are not uniform and depend on the program element identified by the name.
- Some information about the name may be kept outside of the symbol table record and/or some fields of the record may be left vacant for the reason of uniformity. A pointer to this information may be stored in the record.
- The name may be stored in the symbol table record itself, or it can be stored in a separate array of characters and a pointer to it in the symbol table.
- The information about runtime storage location, to be used at the time of code generation, is kept in the symbol table.
- There are various approaches to symbol table organization e.g. Linear List, Search Tree and Hash Table.

Linear List

- It is the simplest approach in symbol table organization.
- The new names are added to the table in the order they arrive.
- A name is searched for its existence linearly.
- The average number of comparisons required are proportional to $0.5*(n+1)$ where n =number of entries in the table.
- It takes less space but more access time.

Search Tree

- It is more efficient than Linear Trees.
- We provide two links- left and right, which point to record in the search tree.

- A new name is added at a proper location in the tree such that it can be accessed alphabetically.
- For any node name1 in the tree, all names accessible by following the left link precede name1 alphabetically.
- Similarly, for any node name1 in the tree, all names accessible by following the right link succeed name1 alphabetically.
- The time for adding/searching a name is proportional to $(m+n) \log_2 n$.

Hash Table

- A hash table is a table of k-pointers from 0 to k-1 that point to the symbol table and record within the symbol table.
- To search a value, we find out the hash value of the name by applying suitable hash function.
- The hash function maps the name into an integer value between 0 and k-1 and uses it as an index in the hash table to search the list of the table records that are built on that hash index.
- To add a non-existent name, we create a record for that name and insert it at the head of the list.

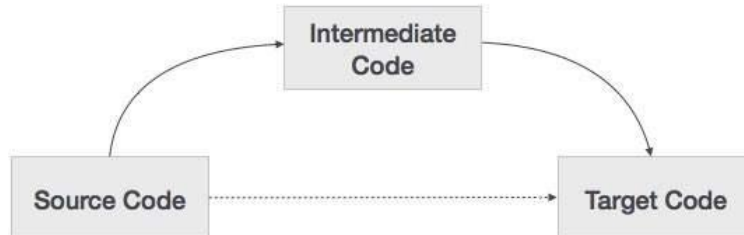
32.3 Scope Information

- Each name possesses a region of validity within the source program called the scope of that name.
- The rules governing the scope of names in a block-structured language are as follows:
 - A name declared within block B is valid only within B.
 - If block B1 is nested within B2, then any name that is valid for B2 is also valid for B1, unless identifier for that name is re-declared in B1.
- These rules require a more complicated symbol table organization than simply a list of associations between names and attributes.
 - One technique is to keep multiple symbol tables for each active block:
 - Each table is list of names and their associated attributes, and the tables are organized on stack.
 - Whenever a new block is entered, a new table is pushed on the stack.
 - When a declaration is compiled, the table on the stack is searched for the name.
 - If name is not found it is inserted.
 - When a reference is translated, it is searched in all tables starting from top.
 - Another technique is to represent scope information in the symbol table.
 - Store the nesting depth of each procedure block in the symbol table.
 - Use the (procedure name, nesting depth) pair as the key to accessing the information from the table.
- The nesting depth of a procedure is a number that is obtained by starting with a value of one for the main and adding one to it every time we go from an enclosing to an enclosed procedure. It counts the number of procedure in the referencing environment of a procedure.

Lecture#33

Intermediate Code Generation

A source code can directly be translated into its target machine code, then why at all we need to translate the source code into an intermediate code which is then translated to its target code? Let us see the reasons why we need an intermediate code.



- If a compiler translates the source language to its target machine language without having the option for generating intermediate code, then for each new machine, a full native compiler is required.
- Intermediate code eliminates the need of a new full compiler for every unique machine by keeping the analysis portion same for all the compilers.
- The second part of compiler, synthesis, is changed according to the target machine.
- It becomes easier to apply the source code modifications to improve code performance by applying code optimization techniques on the intermediate code.

Intermediate Representation

Intermediate codes can be represented in a variety of ways and they have their own benefits.

- **High Level IR** - High-level intermediate code representation is very close to the source language itself. They can be easily generated from the source code and we can easily apply code modifications to enhance performance. But for target machine optimization, it is less preferred.
- **Low Level IR** - This one is close to the target machine, which makes it suitable for register and memory allocation, instruction set selection, etc. It is good for machine-dependent optimizations.

Intermediate code can be either language specific (e.g., Byte Code for Java) or language independent (three-address code).

Three-Address Code

Intermediate code generator receives input from its predecessor phase, semantic analyzer, in the form of an annotated syntax tree. That syntax tree then can be converted into a linear representation, e.g.,

postfix notation. Intermediate code tends to be machine independent code. Therefore, code generator assumes to have unlimited number of memory storage (register) to generate code.

For example:

```
a = b + c * d;
```

The intermediate code generator will try to divide this expression into sub-expressions and then generate the corresponding code.

```
r1 = c * d;
```

```
r2 = b + r1;
```

```
a = r2
```

r being used as registers in the target program.

A three-address code has at most three address locations to calculate the expression. A three-address code can be represented in two forms : quadruples and triples.

Quadruples

Each instruction in quadruples presentation is divided into four fields: operator, arg1, arg2, and result.

The above example is represented below in quadruples format:

| Op | arg ₁ | arg ₂ | result |
|----|------------------|------------------|--------|
| * | c | d | r1 |
| + | b | r1 | r2 |
| + | r2 | r1 | r3 |
| = | r3 | | a |

Triples

Each instruction in triples presentation has three fields : op, arg1, and arg2. The results of respective sub-expressions are denoted by the position of expression. Triples represent similarity with DAG and syntax tree. They are equivalent to DAG while representing expressions.

| Op | arg ₁ | arg ₂ |
|----|------------------|------------------|
| * | c | d |
| + | b | (0) |
| + | (1) | (0) |
| = | (2) | |

Triples face the problem of code immovability while optimization, as the results are positional and changing the order or position of an expression may cause problems.

Indirect Triples

This representation is an enhancement over triples representation. It uses pointers instead of position to store results. This enables the optimizers to freely re-position the sub-expression to produce an optimized code.

Declarations

A variable or procedure has to be declared before it can be used. Declaration involves allocation of space in memory and entry of type and name in the symbol table. A program may be coded and designed keeping the target machine structure in mind, but it may not always be possible to accurately convert a source code to its target language.

Taking the whole program as a collection of procedures and sub-procedures, it becomes possible to declare all the names local to the procedure. Memory allocation is done in a consecutive manner and names are allocated to memory in the sequence they are declared in the program. We use offset variable and set it to zero {offset = 0} that denote the base address.

The source programming language and the target machine architecture may vary in the way names are stored, so relative addressing is used. While the first name is allocated memory starting from the memory location 0 {offset=0}, the next name declared later, should be allocated memory next to the first one.

Example:

We take the example of C programming language where an integer variable is assigned 2 bytes of memory and a float variable is assigned 4 bytes of memory.

```
int a;
float b;
Allocation process:
{offset = 0}
  int a;
  id.type = int
  id.width = 2

offset = offset + id.width
{offset = 2}
  float b;
  id.type = float
  id.width = 4

offset = offset + id.width
{offset = 6}
```

To enter this detail in a symbol table, a procedure *enter* can be used. This method may have the following structure:

```
enter(name, type, offset)
```

This procedure should create an entry in the symbol table, for variable *name*, having its type set to *type* and relative address *offset* in its data area.

Lecture #34

Directed Acyclic Graph

Directed Acyclic Graph (DAG) is a tool that depicts the structure of basic blocks, helps to see the flow of values flowing among the basic blocks, and offers optimization too. DAG provides easy transformation on basic blocks. DAG can be understood here:

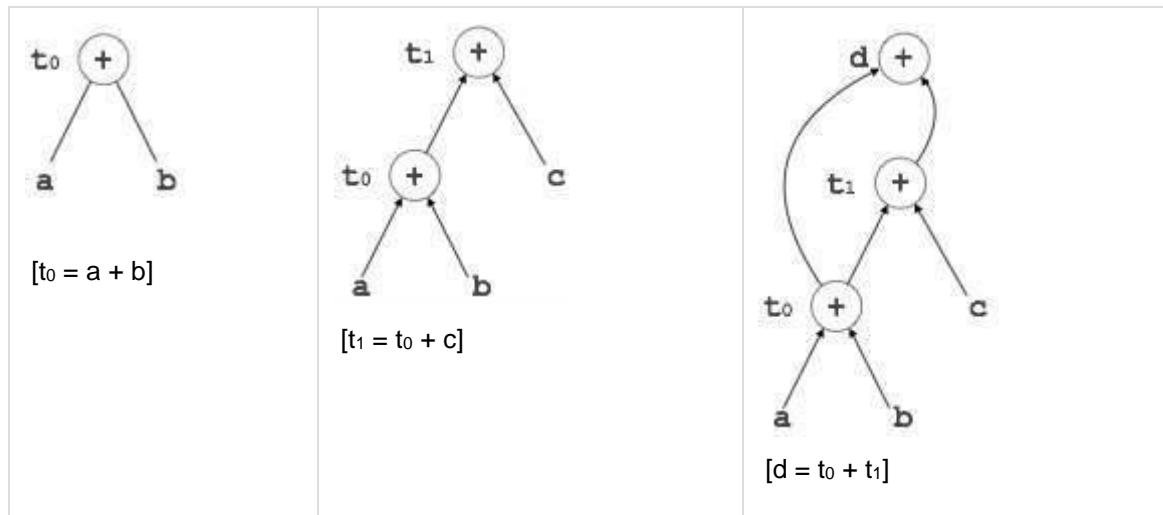
- Leaf nodes represent identifiers, names or constants.
- Interior nodes represent operators.
- Interior nodes also represent the results of expressions or the identifiers/name where the values are to be stored or assigned.

Example:

$t_0 = a + b$

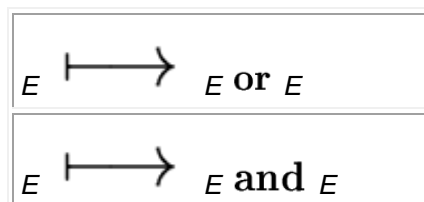
$t_1 = t_0 + c$

$d = t_0 + t_1$



Boolean Expressions and Control Flow

BOOLEAN EXPRESSIONS are constructed using boolean operators. We will consider here the following rules.



| |
|--------------------------------|
| $E \mapsto \text{not } E$ |
| $E \mapsto (E)$ |
| $E \mapsto \text{id relop id}$ |
| $E \mapsto \text{true}$ |
| $E \mapsto \text{false}$ |

- Boolean expressions are used as conditions for statements changing the flow of control.
- Evaluation of boolean expressions can be optimized if it is sufficient to evaluate a part of the expression that determines its value.
- When translating Boolean expressions into three-address code, we can use two different methods.

Numerical method.

Assign numerical values to **true** and **false** and evaluate the expression analogously to an arithmetic expression. This is convenient for boolean expressions which are not involved in flow of control constructs.

Jump method.

Evaluate a boolean expression E as a sequence of conditional and unconditional jumps to location $E.true$ (if E is **true**) or to $E.false$. To be detailed shortly.

WHILE STATEMENTS WITH THE NUMERICAL METHOD.

The following syntax-directed translation generates code for a while statement.

| Production | Semantic Rule |
|---|---|
| $S \mapsto \text{while } E \text{ repeat } S_1$ | $S.begin := \text{newlabel}$ |
| | $S.after := \text{newlabel}$ |
| | $code_1 := \text{generate}(S.begin ':') \parallel E.code \parallel$ |
| | $code_2 := \text{generate}(\text{'if ' } E.place = 0 \text{ 'goto' } S.after) \parallel S_1.code$ |
| | $code_3 := \text{generate}(\text{'goto' } S.begin) \parallel \text{generate}(S.after)$ |
| | $S.code := code_1 \parallel code_2 \parallel code_3$ |

With respect to the previous syntax-directed translation

- we have added two new synthesized attributes $S.begin$ and $S.after$.
- When the value of E becomes zero, control leaves the while statement

Lecture #35

FLOW OF CONTROL STATEMENTS WITH THE JUMP METHOD.

We will consider here the following rules.

| | | |
|---|---|---|
| S | ⟶ | if E then S_1 |
| S | ⟶ | if E then S_1 else S_2 |
| S | ⟶ | while E repeat S_1 |

- In each of these productions, E is the boolean expression to be translated.
- The boolean expression E is associated with two labels (that are **inherited** attributes in the following semantic rules)
 - $E.true$ the label to which control flows if E is **true**,
 - $E.false$ the label to which control flows if E is **false**.
- In each of these productions, S is a flow of control statement associated with two attributes
 - $S.next$ which is a label that is attached to the first 3-address statement to be executed after the code for S , $S.next$ is an **inherited** attribute,
 - $S.code$ is the translation code for S , as usual it is a **synthesized** attribute.

| Production | Semantic Rule |
|--|---|
| $S \longrightarrow \text{if } E \text{ then } S_1$ | $E.true := \text{newlabel}$ |
| | $E.false := S.next$ |
| | $S_1.next := S.next$ |
| | $S.code := E.code \parallel \text{generate}(E.true ':') \parallel S_1.code$ |
| $S \longrightarrow \text{if } E \text{ then } S_1 \text{ else } S_2$ | $E.true := \text{newlabel}$ |
| | $E.false := \text{newlabel}$ |
| | $S_1.next := S.next$ |
| | $S_2.next := S.next$ |
| | $code_1 := E.code \parallel \text{generate}(E.true ':') \parallel S_1.code$ |
| | $code_2 := \text{generate}(\text{goto } S.next) \parallel$ |
| | $code_3 := \text{generate}(E.false ':') \parallel S_2.code$ |
| $S.code := code_1 \parallel code_2 \parallel code_3$ | |
| $S \longrightarrow \text{while } E \text{ repeat } S_1$ | $S.begin := \text{newlabel}$ |
| | $E.true := \text{newlabel}$ |
| | $E.false := S.next$ |
| | $S_1.next := S.begin$ |

| | |
|--|--|
| | $code_1 := generate(S.begin ':') E.code$ |
| | $code_2 := generate(E.true ':') S_1.code$ |
| | $code_3 := generate('goto' S.begin)$ |
| | $S.code := code_1 code_2 code_3$ |

Warning! The above is a syntax-directed definition: It provides formulas for the computation of the attributes $S.code$ (via the computations of the other attributes).

- Since several attributes are inherited and since each action above appears after its associated production, **this is not a translation scheme**.
- However it is an L -attributed definition.
- Then its conversion into a translation scheme is obvious.
- From now on, we may present a translation scheme as a syntax-directed definition if the latter is an L -attributed definition.
- The reason is to make large translation schemes easier to read.

TRANSLATION OF BOOLEAN EXPRESSIONS WITH THE JUMP METHOD. We will consider here the following rules.

| | | |
|-----|-----------|---------------------|
| E | \mapsto | E_1 or E_2 |
| E | \mapsto | E_1 and E_2 |
| E | \mapsto | not E_1 |
| E | \mapsto | (E_1) |
| E | \mapsto | id_1 relop id_2 |
| E | \mapsto | true |
| E | \mapsto | false |

- Here again each symbol E is associated with two inherited attributes
 - $E.true$ the label to which control flows if E is **true**,
 - $E.false$ the label to which control flows if E is **false**.
- The attributes $E.true$ and $E.false$ of E will be defined when the flow of control (where E appears) is translated.

| Production | Semantic Rule |
|--------------------------|-------------------------|
| $E \mapsto E_1$ or E_2 | $E_1.true := E.true$ |
| | $E_1.false := newlabel$ |
| | $E_2.true := E.true$ |
| | $E_2.false := E.false$ |

| | |
|--------------------------------------|--|
| | $E.code := E_1.code \parallel generate(E_1.false:') \parallel E_2.code$ |
| $E \mapsto E_1 \text{ and } E_2$ | $E_1.true := newlabel$ |
| | $E_1.false := E.false$ |
| | $E_2.true := E.true$ |
| | $E_2.false := E.false$ |
| | $E.code := E_1.code \parallel generate(E_1.true:') \parallel E_2.code$ |
| $E \mapsto \text{not } E_1$ | $E_1.true := E.false$ |
| | $E_1.false := E.true$ |
| | $E.code := E_1.code$ |
| $E \mapsto (E_1)$ | $E_1.true := E.true$ |
| | $E_1.false := E.false$ |
| | $E.code := E_1.code$ |
| $E \mapsto id_1 \text{ relop } id_2$ | $code_{e_1} := generate('if' id_1 .place \text{ relop } id_2 .place \text{ goto } E.true)$ |
| | $code_{e_2} := generate(\text{goto } E.false)$ |
| | $E.code := code_{e_1} \parallel code_{e_2}$ |
| $E \mapsto \text{true}$ | $E.code := generate(\text{goto } E.true)$ |
| $E \mapsto \text{false}$ | $E.code := generate(\text{goto } E.false)$ |

Example 5 Let us consider the expression

$a < b$

Assume that

- the attributes *true* and *false* exist for the entire expression
- as labels *Ltrue* and *Lfalse* respectively.

Then the translation is

if $a < b$ goto *Ltrue*

goto *Lfalse*

Observations.

- Of course, this is not optimal and looks funny since the expression to translate is a sentence of the target language!
- But this *jump method* allows us to translate more involved expressions (which are not part of the target language) like those of the following examples.

Example 6 Now consider the expression

$a < b$ or $c < d$

Again assume that the labels *Ltrue* and *Lfalse* have been set for the entire expression. Then the translation is

if $a < b$ goto *Ltrue*

goto *L1*

L1: if $c < d$ goto *Ltrue*

goto *Lfalse*

Example 7 Now consider the expression

$a < b$ or ($c < d$ and $e < f$)

Then the translation is

if $a < b$ goto *Ltrue*

goto *L1*

L1: if $c < d$ goto *L2*

goto *Lfalse*

L2: if $e < f$ goto *Ltrue*

goto *Lfalse*

Of course the generated code is not optimal! Indeed the second statement can be eliminated.

Example 8 Finally consider the expression

while $a < b$ do

if $c < d$ then

$x := y + z$

else

$x := y - z$

Then the translation is

L1: if $a < b$ goto *L2*

goto *Lnext*

L2: if $c < d$ goto *L3*

goto *L4*

L3: $t1 := y + z$

$x := t1$

goto *L1*

L4: $t2 := y - z$

$x := t2$

goto *L1*

Lnext:

Lecture #36

Backpatching

A key problem when generating code for boolean expressions and flow-of-control statements is that of matching a jump instruction with the target of the jump. For example, the translation of the boolean expression B in $\text{if } (B) S$ contains a jump, for when B is false, to the instruction following the code for S . In a one-pass translation, B must be translated before S is examined. Labels can be passed as inherited attributes to where the relevant jump instructions were generated. But a separate pass is then needed to bind labels to addresses. In backpatching, lists of jumps are passed as synthesized attributes. Specifically, when a jump is generated, the target of the jump is temporarily left unspecified. Each such jump is put on a list of jumps whose labels are to be filled in when the proper label can be determined. All of the jumps on a list have the same target label.

One-Pass Code Generation Using Backpatching

Backpatching can be used to generate code for boolean expressions and flow-of-control statements in one pass. Synthesized attributes `truelist` and `falselist` of nonterminal B are used to manage labels in jumping code for boolean expressions. B .`truelist` will be a list of jump or conditional jump instructions into which we must insert the label to which control goes if B is true. B .`falselist` likewise is the list of instructions that eventually get the label to which control goes when B is false. As code is generated for B , jumps to the true and false exits are left incomplete, with the label field unfilled. These incomplete jumps are placed on lists pointed to by B .`truelist` and B .`falselist`, as appropriate. A statement S has a synthesized attribute S .`nextlist`, denoting a list of jumps to the instruction immediately following the code for S .

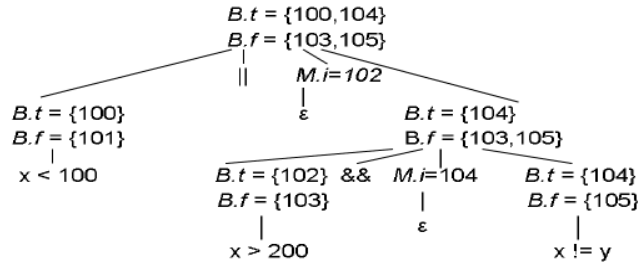
Instructions are generated into an instruction array, and labels will be indices into this array. To manipulate lists of jumps, we use three functions: `makelist(i)`, `merge(p1,p2)`, `backpatch(p,i)`

- `makelist(i)` creates a new list containing only i , an index into the array of instructions; `makelist` returns a pointer to the newly created list.

- `merge(p1,p2)` concatenates the lists pointed to by $p1$ and $p2$, and returns a pointer to the concatenated list.

- `backpatch(p,i)` inserts i as the target label for each of the instructions on the list pointed to by p .

Backpatching for Boolean Expressions



Backpatching for Flow-Of-Control Statements

The grammar is given by the following productions :

$$S \rightarrow \text{if } (B) S \mid \text{if } (B) S \text{ else } S \mid \text{while } (B) S$$

$$\mid \{ L \} \mid A ;$$

$$L \rightarrow L S \mid S$$

- Where S – statement
- L – list of statements
- A- assignment statement
- B – boolean expression

Translation of flow-of-control Statements Using Backpatching

```

S → if (B) M S1
    { backpatch(B.truelist, M.instr);
      S.nextlist = merge(B.falselist, S1.nextlist); }

S → if (B) M1 S1 N else M2 S2
    { backpatch(B.truelist, M1.instr);
      backpatch(B.falselist, M2.instr);
      temp = merge(S1.nextlist, N.nextlist);
      S.nextlist = merge(temp, S2.nextlist); }

S → while M1 (B) M2 S1
    { backpatch(S1.nextlist, M1.instr);
      backpatch(B.truelist, M2.instr);
      S.nextlist = B.falselist;
      emit('goto' M1.instr); }

S → { L }
    { S.nextlist = L.nextlist; }

S → A ;
    { S.nextlist = null; }

M → ε
    { M.instr = nextinstr; }

S → { L }
    { N.nextlist = makelist(nextinstr);
      emit('goto _'); }

L → L1 M S
    { backpatch(L1.nextlist, M.instr)
      L.Nextlist = S.nextlist; }

L → S
    { L.nextlist = S.nextlist; }
  
```

Intermediate Code for Procedures

Assume that parameters are passed by value. Suppose that a is an array of integers, and that f is a function from integers to integers. Then, the assignment $n = f(a[i])$ might translate into the following three-address code:

1. $t_1 = i * 4$
2. $t_2 = a[t_1]$
3. param t_2
4. $t_3 = \text{call } f, 1$
5. $n = t_3$

The grammar below adds functions to the source language

$D \rightarrow \text{define } T \text{ id } (F) \{ S \}$

$F \rightarrow \varepsilon \mid T \text{ id } , F$

$S \rightarrow \text{return } E ;$

$E \rightarrow \text{id } (A)$

$A \rightarrow \varepsilon \mid E , A$

Module-3: Lecture #37

RUN TIME ADMINISTRATION

- The places of the data objects that can be determined at compile time will be *allocated statically*.
- But the places for the some of data objects will be *allocated at run-time*.
The allocation of de-allocation of the data objects is managed by the *run-time support package*.
→Run-time support package is loaded together with the generate target code.
→The structure of the run-time support package depends on the semantics of the programming language (especially the semantics of procedures in that language).

Procedure Activations

- Each execution of a procedure is called as *activation of that procedure*.
- An execution of a procedure starts at the beginning of the procedure body;
- When the procedure is completed, it returns the control to the point immediately after the place where that procedure is called.
- Each execution of a procedure is called as its *activation*.
- *Lifetime* of an activation of a procedure is the sequence of the steps between the first and the last steps in the execution of that procedure (including the other procedures called by that procedure).
- If a and b are procedure activations, then their lifetimes are either non-overlapping or are nested.
- If a procedure is recursive, a new activation can begin before an earlier activation of the same procedure has ended.

Activation Tree

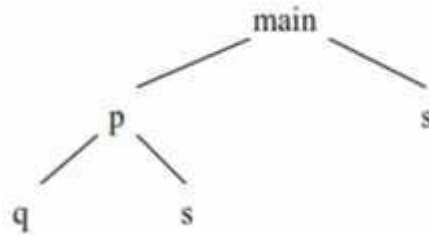
- We can use a tree (called activation tree) to show the way control enters and leaves activations.
- In an activation tree:
 - Each node represents an activation of a procedure.
 - The root represents the activation of the main program.
 - The node a is a parent of the node b iff the control flows from a to b.
 - The node a is left to to the node b iff the lifetime of a occurs before the lifetime of b.

Example:

| | |
|----------------|------------|
| program main; | enter main |
| procedure s; | enter p |
| begin ... end; | enter q |
| procedure p; | exit q |
| procedure q; | enter s |

```
begin ... end;
begin q; s; end;
begin p; s; end;
```

```
exit s
exit p
enter s
exit s
exit main
```



Control Stack

- The flow of the control in a program corresponds to a depth-first traversal of the activation tree that:
 - starts at the root,
 - visits a node before its children, and
 - recursively visits children at each node in a left-to-right order.
- A stack (called **control stack**) can be used to keep track of live procedure activations.
 - An activation record is pushed onto the control stack as the activation starts.
 - That activation record is popped when that activation ends.
- When node n is at the top of the control stack, the stack contains the nodes along the path from n to the root.

Variable Scopes

- The same variable name can be used in the different parts of the program.
- The scope rules of the language determine which declaration of a name applies when the name appears in the program.
- An occurrence of a variable (a name) is:
 - **local**: If that occurrence is in the same procedure in which that name is declared.
 - **non-local**: Otherwise (ie. it is declared outside of that procedure) Example:

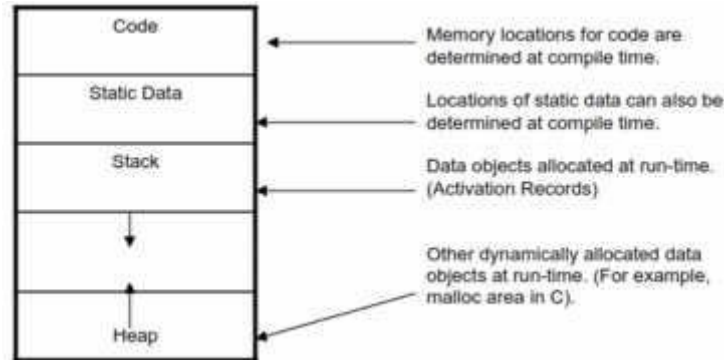
```

procedure p;
var b:real;
procedure p;
var a: integer;
begin a := 1; b := 2; end;
begin ... end;
```

a is local
b is non-local

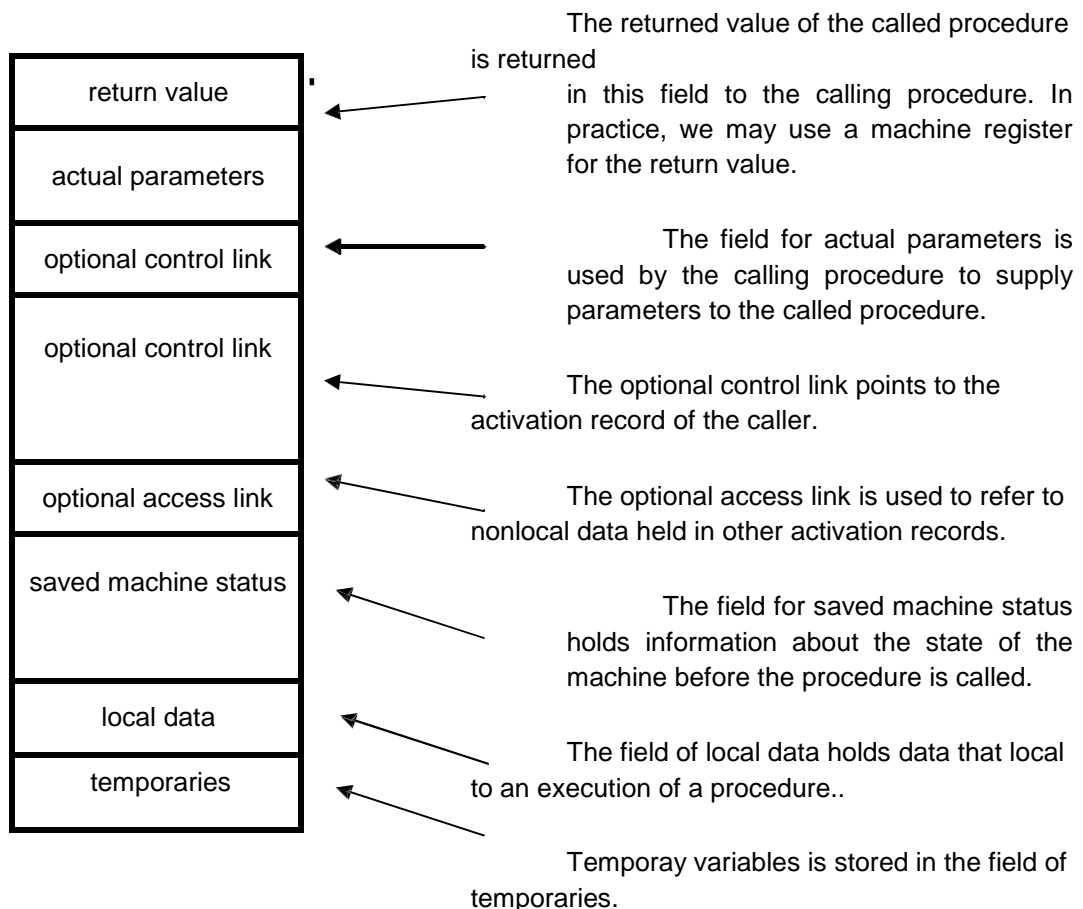
Lecture #38

Storage Organization



Activation Records

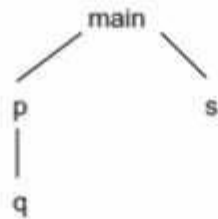
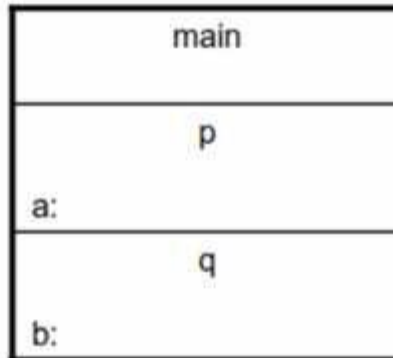
- Information needed by a single execution of a procedure is managed using a contiguous block of storage called **activation record**.
- An activation record is allocated when a procedure is entered, and it is de-allocated when that procedure exited.
- Size of each field can be determined at compile time (Although actual location of the activation record is determined at run-time).
 - Except that if the procedure has a local variable and its size depends on a parameter, its size is determined at the run time.



Example:

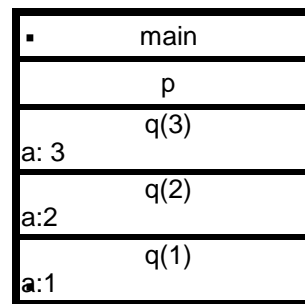
(For a non-recursive procedure)

```
program main;  
  procedure p;  
    var a:real;  
  procedure q;  
    var b:integer;  
    begin ... end;  
  begin q; end;  
  procedure s;  
    var c:integer;  
    begin ... end;  
  begin p; s; end;
```



(For a recursive procedure)

```
program main; procedure p; function  
q(a:integer):integer;  
begin  
if (a=1) then q:=1; else q:=a+q(a-1);  
end;  
begin q(3); end;  
begin p; end;
```



Creation of Activation Records

- Who allocates an activation record of a procedure?
 - Some part of the activation record of a procedure is created by that procedure immediately after that procedure is entered.
 - Some part is created by the caller of that procedure before that procedure is entered.
- Who deallocates?
 - Callee de-allocates the part allocated by Callee.
 - Caller de-allocates the part allocated by Caller.

Displays

- An array of pointers to activation records can be used to access activation records.
- This array is called as displays.
- For each level, there will be an array entry.

| | |
|----|--------------------------------------|
| 1: | Current activation record at level 1 |
| 2: | Current activation record at level 2 |
| 3: | Current activation record at level 3 |
| | |

Lecture #39

ERROR DETECTION AND RECOVERY

- What should the parser do in an error case?
 - The parser should be able to give an error message (as much as possible meaningful error message).
 - It should recover from that error case, and it should be able to continue the parsing with the rest of the input.

Error Recovery Techniques

- Panic-Mode Error Recovery
 - Skipping the input symbols until a synchronizing token is found.
- Phrase-Level Error Recovery
 - Each empty entry in the parsing table is filled with a pointer to a specific error routine to take care of that error case.
- Error-Productions
 - If we have a good idea of the common errors that might be encountered, we can augment the grammar with productions that generate erroneous constructs.
 - When an error production is used by the parser, we can generate appropriate error diagnostics.
 - Since it is almost impossible to know all the errors that can be made by the programmers, this method is not practical.
- Global-Correction
 - Ideally, we would like a compiler to make as few changes as possible in processing incorrect inputs.
 - We have to globally analyze the input to find the error.
 - This is an expensive method, and it is not in practice.

Lecture #40

Error Recovery in Predictive Parsing

- An error may occur in the predictive parsing (LL(1) parsing)
 - if the terminal symbol on the top of stack does not match with the current input symbol.
 - if the top of stack is a non-terminal A, the current input symbol is a, and the parsing table entry $M[A,a]$ is empty.

Panic-Mode Error Recovery in LL(1) Parsing

- In panic-mode error recovery, we skip all the input symbols until a synchronizing token is found.
- What is the synchronizing token?
 - All the terminal-symbols in the follow set of a non-terminal can be used as a synchronizing token set for that non-terminal.
- So, a simple panic-mode error recovery for the LL(1) parsing:
 - All the empty entries are marked as *sync* to indicate that the parser will skip all the input symbols until a symbol in the follow set of the non-terminal A which on the top of the stack. Then the parser will pop that non-terminal A from the stack. The parsing continues from that state.
 - To handle unmatched terminal symbols, the parser pops that unmatched terminal symbol from the stack and it issues an error message saying that that unmatched terminal is inserted.

Example:

$S \rightarrow AbS \mid e \mid \epsilon$
 $A \rightarrow a \mid cAd$
 FOLLOW(S)={ ϵ }
 FOLLOW(A)={b,d}

| | | | | | | |
|---|---------------------|-------------|---------------------|-------------|-------------------|--------------------------|
| | A | b | c | | e | ϵ |
| S | $S \rightarrow AbS$ | <i>sync</i> | $S \rightarrow AbS$ | <i>sync</i> | $S \rightarrow e$ | $S \rightarrow \epsilon$ |
| A | $A \rightarrow a$ | <i>sync</i> | $A \rightarrow cAd$ | <i>sync</i> | <i>sync</i> | <i>sync</i> |

For string aab

| Stack | Input | Output |
|-----------|-------|--|
| \$\$ | aab\$ | $S \rightarrow AbS$ |
| AbS | | |
| \$\$bA | aab\$ | $A \rightarrow a$ |
| cAd | | |
| \$\$ba | aab\$ | |
| \$\$b | ab\$ | Error: missing b, inserted (illegal A) |
| \$\$ | ab\$ | $S \rightarrow AbS$ |
| d, pop A) | | |
| \$\$bA | ab\$ | $A \rightarrow a$ |
| \$\$ba | ab\$ | |
| \$\$b | b\$ | |
| \$\$ | \$ | $S \rightarrow \epsilon$ |
| \$\$ | | accept |

For string ceadb

| Stack | Input | Output |
|---|---------|--------------------------|
| \$\$ | ceadb\$ | $S \rightarrow$ |
| \$\$bA | ceadb\$ | $A \rightarrow$ |
| \$\$bdAc | ceadb\$ | |
| \$\$bdA | eadb\$ | Error: unexpected e |
| (Remove all input tokens until first b or | | |
| \$\$bd | db\$ | |
| \$\$b | b\$ | |
| \$\$ | \$ | $S \rightarrow \epsilon$ |
| \$ | \$ | accept |

Phrase-Level Error Recovery

- Each empty entry in the parsing table is filled with a pointer to a special error routine which will

take care that error case.

- These error routines may:
 - change, insert, or delete input symbols.
 - issue appropriate error messages
 - pop items from the stack.
- We should be careful when we design these error routines, because we may put the parser into an infinite loop.

Error Recovery in Operator-Precedence Parsing

Error Cases:

- No relation holds between the terminal on the top of stack and the next input symbol.
- A handle is found (reduction step), but there is no production with this handle as a right side

Error Recovery:

- Each empty entry is filled with a pointer to an error routine.
- Decides the popped handle “looks like” which right hand side. And tries to recover from that situation.

Error Recovery in LR Parsing

- An LR parser will detect an error when it consults the parsing action table and finds an error entry. All empty entries in the action table are error entries.
- Errors are never detected by consulting the goto table.
- An LR parser will announce error as soon as there is no valid continuation for the scanned portion of the input.
- A canonical LR parser (LR(1) parser) will never make even a single reduction before announcing an error.
- The SLR and LALR parsers may make several reductions before announcing an error.
- But, all LR parsers (LR(1), LALR and SLR parsers) will never shift an erroneous input symbol onto the stack.

Panic Mode Error Recovery in LR Parsing

- Scan down the stack until a state **s** with a goto on a particular nonterminal **A** is found. (Get rid of everything from the stack before this state s).
- Discard zero or more input symbols until a symbol **a** is found that can legitimately follow A.
 - The symbol a is simply in FOLLOW (A), but this may not work for all situations.
- The parser stacks the nonterminal **A** and the state **goto[s,A]**, and it resumes the normal parsing.
- This nonterminal A is normally is a basic programming block (there can be more than one choice for A).
 - stmt, expr, block, ...

Phrase-Level Error Recovery in LR Parsing

- Each empty entry in the action table is marked with a specific error routine.
- An error routine reflects the error that the user most likely will make in that case.
- An error routine inserts the symbols into the stack or the input (or it deletes the symbols from the stack and the input, or it can do both insertion and deletion).
 - missing operand
 - unbalanced right parenthesis

Lecture #41

CODE OPTIMIZATION

- Code optimization is aimed at obtaining a more efficient code.
- Two constraints on the technique used to perform optimizations
 - They must ensure that the transformed program is semantically equivalent to the original program.
 - The improvement of the program efficiency must be achieved without changing the algorithms which are used in the program.
- Optimization may be classified as Machine dependent and Machine independent.
 - Machine dependent optimizations exploit characteristics of the target machine.
 - Machine independent optimizations are based on mathematical properties of a sequence of source statements.

Optimizing Transformations

Common Sub-expression Elimination

- An expression need not be evaluated if it was previously computed and values of variables in this expression have not changed since the earlier computations.

Example:

```
a = d * c;  
. . .  
  
d = b * c + x - y;
```

We can eliminate the second evaluation of $b*c$ from this code if none of the intervening statements has changed its value. The code can be rewritten as given below.

```
T1 = b * c;  
a = T1;  
. . .  
d = T1 + x - y;
```

Compile Time Evaluation

- We can improve the execution efficiency of a program by shifting execution time actions to compile time.
- We can evaluate an expression by a single value (known as *folding*).

Example:

```
A = 2 * (22.0/7.0) * r
```

Here we can perform the computation $2 * (22.0/7.0)$ at compile time itself.

- If a variable is assigned a constant value and is used in an expression without being assigned other value to it, we can evaluate some portion of the expression using the constant value (known as *Constant Propagation*).

Example:

```
x = 12.4
```

```
.
```

```
.
```

```
y = x / 2.3
```

Here we evaluate $x / 2.3$ as $12.4 / 2.3$ at compile time.

Variable Propagation

- If a variable is assigned to another variable, we use one in place of another.
- This will be useful to carry out other optimization that were otherwise not possible.

Example:

```
c = a * b; x = a;
```

```
.
```

```
.
```

```
.
```

```
d = x * b;
```

Here, if we replace x by a then $a * b$ and $x * b$ will be identified as common sub- expressions.

Dead Code Elimination

- If the value contained in a variable at that point is not used anywhere in the program subsequently, the variable is said to be dead at that place.
- If an assignment is made to a dead variable, then that assignment is a dead assignment and it can be safely removed from the program.
- A piece of code is said to be dead if it computes values that are never used anywhere in the program.
- Dead Code can be eliminated safely.
- Variable propagation often leads to making assignment statement into dead code.

Example:

```
c = a * b;
```

```
x = a;
```

```
.
```

```
.  
.   
d = x * b + 4;
```

Variable propagation will lead to following changes. $c = a * b$;

```
x = a;  
.   
.   
.   
d = a * b + 4;
```

This assignment $x = a$ is now useless and can be removed

```
c = a * b;  
.   
.   
.   
d = a * b + 4;
```

Code Motion

- We aim to improve the execution time of the program by reducing the evaluation frequency of expressions.
- Evaluation of expressions is moved from one part of the program to another in such a way that it is evaluated lesser frequently.
- Loops are usually executed several times.
- We can bring the loop-invariant statements out of the loop.

Example:

```
a = 200;  
while (a > 0)  
{  
    b = x + y;  
    if ( a%b == 0)  
        printf ("%d", a);  
}
```

The statement $b = x + y$ is executed every time with the loop. But because it is loop- invariant, we can bring it outside the loop. It will then be executed only once.

```
a = 200;  
b = x + y;  
while (a > 0)  
{  
    if ( a%b == 0)  
        printf ("%d", a);  
}
```

Induction Variables and Strength Reduction

- An induction variable may be defined as an integer scalar variable which is used in loop for the following kind of assignments $i = i + \text{constant}$.
- Strength Reduction means replacing the high strength operator by a low strength operator.
- Strength Reduction used on induction variables to achieve a more efficient code.

Example:

```

i = 1;
while (i < 10)
{
    .
    .
    .
    y = i * 4;
    .
    .
}

```

This code can be replaced by the following code. $i = 1$;

```

t = 4;
while (t < 40)
{
    .
    .
    .
    y = t;
    .
    .
    .
    t = t + 4;
    .
    .
    .
}

```

Use of Algebraic Identities

- Certain computations that look different to the compiler and are not identified as common sub-expressions are actually same.
- An expression $B \text{ op } C$ will usually be treated as being different to $C \text{ op } B$.
- However, for certain operations (like addition and multiplication), they will produce the same result.
- We can achieve further optimization by treating them as common sub-expressions for such operations.

Lecture #42

Local Optimizations

- Target code generated statement by statement generally contains redundant instructions.
- We can improve the quality of such code by applying optimizing transformations locally by examining a short sequence of code instructions and replacing them by faster or shorter sequence, if possible.
- This technique is known as *Peephole Optimization* where the peephole is a small moving window on the program.
- Many of the code optimization techniques can be carried out by a single portion of a program known as *Basic Block*.

Basic Block

- A basic Block is defined as a sequence of consecutive statements with only one entry (at the beginning) and one exit (at the end).
- When a Basic Block of a program is entered, all the statements are executed in sequence without a halt or possibility of branch except at the end.
- In order to determine all the Basic Block in a program, we need to identify the *leaders*, the first statement of each Basic Block.
 - Any statement that satisfies the following conditions is a leader;
 - The first statement is leader.
 - Any statement which is the target of any goto (jump) is a leader.
 - Any statement that immediately follows a goto (jump) is a leader.
 - A basic block is defined as the portion of code from one leader to the statement up to but including the next leader or the end of the program.

Flow Graph

- It is a directed graph that is used to portray basic block and their successor relationships.
- The nodes of a flow graph are the basic blocks.
- The basic block whose leader is the first statement is known as the initial block.
- There is a directed edge from block B1 to B2 if B2 could immediately follow B1 during execution.
- To determine whether there should be directed edge from B1 to B2, following criteria is applied:
 - There is a jump from last statement of B1 to the first statement of B2, OR
 - B2 immediately follows B1 in order of the program and B1 does not end in an unconditional jump.
- B1 is known as the *predecessor* of B2 and B2 is a *successor* of B1.

Loops

- We need to identify all the loops in a flow graph to carry out many optimizations discussed earlier.
- A loop is a collection of nodes that
 - is strongly connected i.e. from any node in the loop to any other, there is a path of length one or more wholly within the loop, and
 - has a unique entry, a node in the loop such that the only way to reach a node in the loop from a node outside the loop is to first go through the entry.

DAG Representation of a Basic Block

- Many optimizing transformations can be implemented using the DAG representation of a basic block.
- DAG stands for Directed Acyclic Graph i.e. a graph with directed edges and no cycles.
- DAG is very much like a tree but differs in that it may contain shared nodes where shared nodes indicate common sub-expressions.
- A DAG has following components;
 - Leaves are labeled by unique identifiers, either variable names or constants.
 - Interior nodes are labeled by an operator symbol.
 - Nodes are optionally given an extra set of identifiers known as *attached identifiers*.

DAG Construction

- We assume there are initially no nodes and NODE () is undefined for all arguments.
- The 3-address statements has one of three cases:
 - (i) $A = B \text{ op } C$
 - (ii) $A = \text{op } B$
 - (iii) $A = B$
- We shall do the following steps (1) through (3) for each 3-address statement of the basic block:
 - (1) If NODE (B) is undefined, create a leaf labeled B, and let NODE (B) be this node. In case (i), if NODE (C) is undefined, create a leaf labeled C and let that leaf be NODE (C); (2) In case (i), determine if there is a node labeled op whose left child is NODE (B) and whose right child is NODE (C). (This is to catch common sub-expressions.) If not create such a node. In case (ii), determine whether there is a node labeled op whose lone child is NODE (B). If not create such a node. Let n be the node found or created in both cases. In case (iii), let n be NODE (B).
 - (3) Append A to the list of attached identifiers for the node n in (2). Delete A from the list of attached identifiers for NODE (A). Finally, set NODE (A) to n.

Applications of DAG

- We automatically detect common sub-expressions while constructing DAG.
- It is also known as to which identifiers have their values used inside the block; they are exactly those for which a leaf is created in Step (1).
- We can also determine which statements compute values which could be used outside the block; they are exactly those statements S whose node n in step (2) still has NODE (A) = n at the end of DAG construction, where A is the identifier assigned by statement S i. e. A is still an attached identifier for n.

Global Data Flow Analysis

- Certain optimizations can be achieved by examining the entire program and not just a portion of the program.
- User-defined chaining is one particular problem of this kind.
- Here we try to find out as to which definition of a variable is applicable in a statement using the value of that variable.