13BECS603

PRINCIPLES OF COMPILER DESIGN

COURSE OBJECTIVES:

- •To understand and list the different stages in the process of compilation.
- To understand and design Lexical analyzers and parsers
- •To Develop algorithms to generate code for a target machine
- To learn and develop techniques for optimization of code.

COURSE OUTCOMES:

Upon completing the course the students will be able to

- Understand the complete process of compilation from source code to target code
- develop the lexical analyzer and parsers
- Develop algorithms to generate code for a target machine
- Optimize the generated code

UNIT 1:Introduction

Introduction - What is a Compiler? - Cousins of a Compiler- Assembler - Interpreter - Phases of compilation and overview, Lexical Analysis, Syntax Analysis, Semantic Analysis, Intermediate code Generation, Code Optimization, Code Generation - Specification of Tokens.

UNIT 2:Lexical Analysis(scanner)

Regular l a n g u a g e s, finite automata, regular expressions, from regular expressions to finite automata, scanner generator (lex, flex).

SyntaxAnalysis(Parser):Context-freelanguagesandgrammars,push-downautomata,LL(1)gram-marsandtop-downparsing,operatorgrammars,LR(O),SLR(1),LR(1),LALR(1)grammarsandbottom-upparsing, ambiguityandLRparsing,LALR(1)parser generator(yacc, bison)bottom-upbottom-upbottom-upbottom-upbottom-up

UNIT 3:Semantic Analysis

Attribute grammars, syntax directed definition, evaluation and flow of attribute in a syntax tree.

Symbol Table:Its structure, symbol attributes and management. Run-time environment: Procedure activation, parameter passing, value return, memory allocation, and scope.Intermediate Code Generation:Translation of different language features, different types of intermediate forms.

UNIT 4: Code Improvement(optimization)

Analysis: control-flow, data-flow dependence etc.; Code improvement local optimization, global optimization, loop optimization, peep-hole optimization etc.Architecture dependent code improvement: instruction scheduling(for pipeline), loop optimization (for cache memory) etc. Register allocation and target code generation

UNIT 5: Advanced topics

Type systems, data abstraction, compilation of Object Oriented features and non-imperative programming languages.

Total Hours: 45

TEXT BOOKS:

1. Alfred Aho, Ravi Sethi, Jeffrey D Ullman, Compilers Principles, Techniques and Tools, Pearson Education Asia, 2nd Edition, 2017.

(9)

(9)

(9)

(9)

(9)

2. Allen I Holub, Compiler Design in C, Prentice Hall of india, 2016.

REFERENCES:

- 1. Keith Cooper and lindaTorczon, Engineering a compiler, 2nd edition, 2016.
- Bennet.J.P, Introduction to Compiler Techniques, Tata McGraw-Hill, 2015.
- 3. R.Levine, Tony Mason, Doug Brown John, Lex &Yacc, 2nd Edition (October 2012) O'Reilly & Associates.
- 4. Kenneth c.Louden, Compiler Construction: Principles and Pratice, Thomson Learning, 2018.

WEBSITES:

- 1. http://www.tenouk.com/ModuleW.html/
- 2. http://www.mactech.com/articles/mactech/Vol.06/06.04/Lexical Analysis/index.html



KARPAGAM ACADEMY OF HIGHER EDUCATION

Coimbatore-21. FACULTY OF ENGINEERING DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

13BECS603- PRINCIPLES OF COMPILER DESIGN

LECTURE PLAN

S.NO	DESCRIPTION OF PORTION TO BE COVERED	HOURS	Reference Book & Page Nos. Used for teaching	TEACHING AIDS
1	Discussion on the Fundamentals of Compilers	1	R[1] Page no 17-33	PPT
2	Introduction to Types of Compilers-uses of compilers	1	R[1] Page no 12-14	PPT
	UNIT- I I	INTRODU	CTION TO COMPILING	
3	Compilers Analysis of the source program	1	R[2]-Page no 1.1-1.3 R[1]-Page no1-3	BB
4	Phases of compiler	1	R[1]-Page no 4-11 R[2]-Page no1.6-1.10	BB
5	Tutorial: Compilers Phases of compiler	1	R[2]-Page no 1.1-1.3 R[1]-Page no1-3	PPT
6	Cousins of the compiler Grouping of phases	1	R[2]-Page no 1.14-1.17	BB
7	Compiler construction tools Lexical Analysis	1	R[2]-Page no 1.19	PPT
8	The role of the lexical analyzer	1	R[2]-Page no 1.25-1.28	PPT

9	Tutorial: Lexical Analysis The role of the lexical analyzer	1	R[2]-Page no 1.25-1.28	РРТ
10	Input buffering-Tokens Specification	1	R[2]-Page no 1.28-1.29	PPT
,	FOTAL HOURS FOR UNIT-I		10	
	UNIT	- II S	YNTAX ANALYSIS	
11	The role of the parser writing a grammar	1	R[1]-Page no 191-195 R[2]-Page no 2.9-216	BB
12	Context-free grammars	1	R[2]-Page no 197-206	РРТ
13	Tutorial hour – Context-free grammars	1	R[2]-Page no 197-206	PPT
14	Top-down parsing Recursive-descent parser Predictive parser	1	R[1]-Page no 217-231 R[2]-Page no 2.16-2.19	BB
15	Constructing an SLR(1) parsing table	1	R[1]-Page no 252-248 R[2]-Page no 2.43	BB
16	Bottom-up Parsing Shift reduce parsing Operator-precedence parsing	1	R[2]-Page no 2.26-2.37 R[1] Page no 233-240	PPT
17	Tutorial : Bottom-up Parsing Shift reduce parsing	1	R[2]-Page no 2.26-2.37 R[1] Page no 233-240	BB
18	LR Parsers SLR Parser	1	R[1]-Page no 241-248 R[2]-Page no 2.39-2.42	BB
19	Canonical LR Parser	1	R[1]-Page no 259-261	РРТ

20	LALR Parser	1	R[1]-Page no 266-277 R[2]-Page no 2.62	РРТ	
Г	TOTAL HOURS FOR UNIT-II 10				
	UNIT -III INTERMEDIATE CODE GENERATION				
21	Intermediate languages	1	R[2]-Page no 3.1-3.6	BB	
22	Declarations Assignment statements	1	R[1]-Page no 370-379 R[2] Page no 3.13-3.14	BB	
23	Boolean expressions	1	R[1]-Page no 399-409 R[2] Page no 3.24-3.26	BB	
24	Tutorial: Boolean expressions	1	R[1]-Page no 399-409 R[2] Page no 3.24-3.26	BB	
25	Case statements	1	R[1]-Page no 418-421 R[2] Page no 33.31-3.32	РРТ	
26	Backpatching	1	R[2] Page no 3.40	BB	
27	Tutorial: Case statements	1	R[1]-Page no 418-421 R[2] Page no 33.31-3.32	PPT	
28	Procedure calls	1	R[2] Page no 3.41-3.45	РРТ	
29	Symbol table	1	R[2] Page no 3.41-3.45	РРТ	
TOTAL HOURS FOR UNIT-III9					
	UNIT- IV CODE GENERATION				
30	Issues in the design of a code generator	1	R[1]-Page no 501-505 R[2] Page no 4.2-4.3	BB	
31	The target machine	1	R[1]-Page no 512-516 R[2] Page no 4.6	BB	

32	Run-time storage management	1	R[1]-Page no 427-440 R[2] Page no 4.6-4.10	BB
33	Tutorial : Run-time storage management	1	R[1]-Page no 427-440 R[2] Page no 4.6-4.10	BB
34	Basic blocks and flow graphs Next use information, a simple code generator	1	R[2] Page no 4.10-4.14	BB
35	The dag representation of basic blocks	1	R[2] Page no 4.10-4.14	PPT
36	Tutorial: The dag representation of basic blocks	1	R[2] Page no 4.10-4.14	PPT
37	Peephole optimization	1	R[1]-Page no 549-553 R[2]-Page no 4.22-4.24	BB
Т	OTAL HOURS FOR UNIT-IV		8	
	UNIT- V CODE OP	TIMIZATIO	ON AND RUN TIME ENVIRONMEN	rs
	Introduction to code			
38	optimization	1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2	BB
38 39	optimization The principle sources of optimization	1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2 R[1]-Page no 583-596 R[2]-Page no 5.3-5.11	BB BB
38 39 40	optimization The principle sources of optimization Optimization of basic blocks Global data flow analysis	1 1 1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2 R[1]-Page no 583-596 R[2]-Page no 5.3-5.11 R[2]-Page no 5.3-5.11	BB BB BB
38 39 40 41	 optimization The principle sources of optimization Optimization of basic blocks Global data flow analysis Tutorial Code optimization 	1 1 1 1 1 1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2 R[1]-Page no 583-596 R[2]-Page no 5.3-5.11 R[2]-Page no 5.3-5.11 R[2]-Page no 5.3-5.11	BB BB BB BB
38 39 40 41 42	Infoduction to codeoptimizationThe principle sources of optimizationOptimization of basic blocks Global data flow analysisTutorial Code optimizationRun time environment	1 1 1 1 1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2 R[1]-Page no 583-596 R[2]-Page no 5.3-5.11	BB BB BB BB PPT
38 39 40 41 42 43	 Infoduction to code optimization The principle sources of optimization Optimization of basic blocks Global data flow analysis Tutorial Code optimization Run time environment Source Language issues 	1 1 1 1 1 1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2 R[1]-Page no 583-596 R[2]-Page no 5.3-5.11	BB BB BB BB PPT BB
38 39 40 41 42 43 44	 Infoduction to code optimization The principle sources of optimization Optimization of basic blocks Global data flow analysis Tutorial Code optimization Run time environment Source Language issues Storage Organization Storage Allocation strategies 	1 1 1 1 1 1 1	R[1]-Page no 583-596 R[2]-Page no 5.1-5.2 R[1]-Page no 583-596 R[2]-Page no 5.3-5.11 R[2]-Page no 5.20 R[2]-Page no 5.20 R[2]-Page no 5.20 R[2]-Page no 5.27-5.31	BB BB BB BB PPT BB BB BB

46	Access to non-local names, parameter missing	R[2]-page no 5.32	BB
47	Discussion on Past Five Year	End Semester Question Paper 1	
T	OTAL HOURS FOR UNIT-V	10	
	TOTAL LECTURE HOURS	37	
TOTAL TUTORIAL HOURS		10	
TOTAL HOURS		47	

REFERENCES

1	Alfred Aho, Ravi Sethi, Jeffrey D Ullman, 2006, Compilers Principles, Techniques and Tools, 4 th Edition, Pearson Education Asia
2	Author: P.Kalaiselvi, Principles Of Compiler Design A.A.R.Senthikumaar, 2008, 3 rd edition, charulatha publication, India
3	Allen I. Holub, 2003, Compiler Design in C, 4 th Edition, Prentice Hall of India.
4	Fischer.C.N and R.J.LeBlanc, 2003, Crafting a compiler with C, 3 rd Edition, Benjamin Cummings.
5	Bennet.J.P, 2003, Introduction to Compiler Techniques, 2 nd Edition, Tata McGraw-Hill

Faculty In charge

HOD

LECTURE NOTES

<u>CHAPTER I- LEXICAL ANALYSIS</u> <u>1.1 INRODUCTION TO COMPILING</u>

Translator:

It is a program that translates one language to another.



Types of Translator:

1.Interpreter

2.Compiler

3.Assembler

J.Assembler

1.Interpreter:

It is one of the translators that translate high level language to low level language.



Figure 1.2: Interpreter

During execution, it checks line by line for errors. Example: Basic, Lower version of Pascal.

2.Assembler:



Figure 1.3:Assembler

Example: Microprocessor 8085, 8086.

3.Compiler:

It is a program that translates one language(source code) to another language (target code).



Figure 1.4:Compiler

It executes the whole program and then displays the errors. Example: C, C++, COBOL, higher version of Pascal.

	Difference	between	compiler	and	inter	preter:
--	-------------------	---------	----------	-----	-------	---------

Compiler	Interpreter
It is a translator that translates high level to low level language	It is a translator that translates high level to low level language
It displays the errors after the whole program is	It checks line by line for errors.
executed.	
Examples: Basic, lower version of Pascal.	Examples: C, C++, Cobol, higher version of
	Pascal.

<u>1.1.1</u> PARTS OF COMPILATION

There are 2 parts to compilation:

- 1. Analysis
- 2. Synthesis

Analysis part breaks down the source program into constituent pieces and creates an intermediate representation of the source program.

Synthesis part constructs the desired target program from the intermediate representation.



Figure 1.5:Parts of Compilation

Software tools used in Analysis part:

1) Structure editor:

Takes as input a sequence of commands to build a source program. The structure editor not only performs the text-creation and modification functions of an ordinary text editor, but it also analyzes the program text, putting an appropriate hierarchical structure on the source program. For example, it can supply key words automatically - while do and begin..... end.

2) Pretty printers :

A pretty printer analyzes a program and prints it in such a way that the structure of the program becomes clearly visible. For example, comments may appear in a special font.

3) Static checkers :

A static checker reads a program, analyzes it, and attempts to discover potential bugs without running the program. For example, a static checker may detect that parts of the source program can never be executed.

4) Interpreters :

Translates from high level language (BASIC, FORTRAN, etc..) into machine language. An interpreter might build a syntax tree and then carry out the operations at the nodes as it walks the tree. Interpreters are frequently used to execute command language since each operator executed in a command language is usually an invocation of a complex routine such as an editor or complier.

1.2 ANALYSIS OF THE SOURCE PROGRAM

Analysis consists of 3 phases:

Linear/Lexical Analysis :

It is also called scanning. It is the process of reading the characters from left to right and grouping into tokens having a collective meaning. For example, in the assignment statement a=b+c*2, the characters would be grouped into the following tokens:

- i) The identifier1 'a'
- ii) The assignment symbol (=)
- iii) The identifier2 'b'
- iv) The plus sign (+)
- v) The identifier3 'c'
- vi) The multiplication sign (*)
- vii) The constant '2'

Syntax Analysis :

It is called parsing or hierarchical analysis. It involves grouping the tokens of the source program into grammatical phrases that are used by the compiler to synthesize output. They are represented using a syntax tree as shown below:



Figure 1.6:Syntax Analysis

A **syntax tree** is the tree generated as a result of syntax analysis in which the interior nodes are the operators and the exterior nodes are the operands. This analysis shows an error when the syntax is incorrect.

Semantic Analysis :

It checks the source programs for semantic errors and gathers type information for the subsequent code generation phase. It uses the syntax tree to identify the operators and operands of statements. An important component of semantic analysis is **type checking**. Here the compiler checks that each operator has operands that are permitted by the source language specification.

1.3 PHASES OF COMPILER

A Compiler operates in phases, each of which transforms the source program from one representation into another. The following are the phases of the compiler:

Main phases:

Lexical analysis
 Syntax analysis
 Semantic analysis
 Intermediate code generation
 Code optimization
 Code generation

Sub-Phases:

Symbol table management
 Error handling



Figure 1.7: Phases of Compiler

LEXICAL ANALYSIS:

It is the first phase of the compiler. It gets input from the source program and produces tokens as output. It reads the characters one by one, starting from left to right and forms the tokens.

Token : It represents a logically cohesive sequence of characters such as keywords, operators, identifiers, special symbols etc.

Example: a + b = 20

Here, a,b,+,=,20 are all separate tokens.

Group of characters forming a token is called the Lexeme.

The lexical analyser not only generates a token but also enters the lexeme into the symbol table if it is not already there.

SYNTAX ANALYSIS:

It is the second phase of the compiler. It is also known as parser. It gets the token stream as input from the lexical analyser of the compiler and generates syntax tree as the output.

Syntax tree: It is a tree in which interior nodes are operators and exterior nodes are operands. Example: For a=b+c*2, syntax tree is



Figure 1.8: Syntax Tree

SEMANTIC ANALYSIS:

It is the third phase of the compiler. It gets input from the syntax analysis as parse tree and checks whether the given syntax is correct or not. It performs type conversion of all the data types into real data types.

INTERMEDIATE CODE GENERATION:

It is the fourth phase of the compiler. It gets input from the semantic analysis and converts the input into output as intermediate code such as three-address code. The three-address code consists of a sequence of instructions, each of which has atmost three operands.

Example: t1=t2+t3

CODE OPTIMIZATION:

It is the fifth phase of the compiler. It gets the intermediate code as input and produces optimized intermediate code as output. This phase reduces the redundant code and attempts to improve the intermediate code so that faster-running machine code will result. During the code optimization, the result of the program is not affected. To improve the code generation, the optimization involves,

- deduction and removal of dead code (unreachable code).
- calculation of constants in expressions and terms.
- collapsing of repeated expression into temporary string.
- loop unrolling.
- moving code outside the loop.
- removal of unwanted temporary variables.

CODE GENERATION:

It is the final phase of the compiler. It gets input from code optimization phase and produces the target code or object code as result.Intermediate instructions are translated into a sequence of machine instructions that perform the same task. The code generation involves

- allocation of register and memory
- generation of correct references
- generation of correct data types
- generation of missing code

SYMBOL TABLE MANAGEMENT:

Symbol table is used to store all the information about identifiers used in the program. It is a data structure containing a record for each identifier, with fields for the attributes of the identifier. It allows to find the record for each identifier quickly and to store or retrieve data from that record. Whenever an identifier is detected in any of the phases, it is stored in the symbol table.

ERROR HANDLING:

Each phase can encounter errors. After detecting an error, a phase must handle the error so that compilation can proceed. In lexical analysis, errors occur in separation of tokens. In syntax analysis, errors occur during construction of syntax tree. In semantic analysis, errors occur when the compiler detects constructs with right syntactic structure but no meaning and duringtype conversion. In code optimization, errors occur when the result is affected by the optimization. In code generation, it shows error when code is missing etc.

To illustrate the translation of source code through each phase, consider the statement a=b+c*2. The figure shows the representation of this statement after each phase:



<u>1.4 COUSINS OF COMPILER</u>

- 1. Preprocessor
- 2. Assembler
- 3. Loader and Link-editor

PREPROCESSOR

A preprocessor is a program that processes its input data to produce output that is used as input to another program. The output is said to be a preprocessed form of the input data, which is often used by some subsequent programs like compilers.

They may perform the following functions :

- 1. Macro processing
- 2. File Inclusion
- 3. Rational Preprocessors
- 4. Language extension

1. Macro processing:

A macro is a rule or pattern that specifies how a certain input sequence should be mapped to an output sequence according to a defined procedure. The mapping process that instantiates a macro into a specific output sequence is known as macro expansion.

2. File Inclusion:

Preprocessor includes header files into the program text. When the preprocessor finds an #include directive it replaces it by the entire content of the specified file.

3. Rational Preprocessors:

These processors change older languages with more modern flow-of-control and datastructuring facilities.

4. Language extension :

These processors attempt to add capabilities to the language by what amounts to built-in macros. For example, the language Equel is a database query language embedded in C.

ASSEMBLER

Assembler creates object code by translating assembly instruction mnemonics into machine code. There are two types of assemblers:

- One-pass assemblers go through the source code once and assume that all symbols will be defined before any instruction that references them.
- Two-pass assemblers create a table with all symbols and their values in the first pass, and then use the table in a second pass to generate code.

LINKER AND LOADER

A **linker** or **link editor** is a program that takes one or more objects generated by a compiler and combines them into a single executable program.

Three tasks of the linker are :

- 1. Searches the program to find library routines used by program, e.g. printf(), math routines.
- 2. Determines the memory locations that code from each module will occupy and relocates its instructions by adjusting absolute references
- 3. Resolves references among files.

A **loader** is the part of an operating system that is responsible for loading programs in memory, one of the essential stages in the process of starting a program.

1.5 GROUPING OF THE PHASES

Compiler can be grouped into front and back ends:

- Front end: analysis (machine independent)

These normally include lexical and syntactic analysis, the creation of the symbol table, semantic analysis and the generation of intermediate code. It also includes error handling that goes along with each of these phases.

- Back end: synthesis (machine dependent)

It includes code optimization phase and code generation along with the necessary error handling and symbol table operations.

Compiler passes

A collection of phases is done only once (single pass) or multiple times (multi pass)

- Single pass: usually requires everything to be defined before being used in source program.
- Multi pass: compiler may have to keep entire program representation in memory.

Several phases can be grouped into one single pass and the activities of these phases are interleaved during the pass. For example, lexical analysis, syntax analysis, semantic analysis and intermediate code generation might be grouped into one pass.

1.6 COMPILER CONSTRUCTION TOOLS

These are specialized tools that have been developed for helping implement various phases of a compiler. The following are the compiler construction tools:

1) Parser Generators:

-These produce syntax analyzers, normally from input that is based on a context-free grammar.

-It consumes a large fraction of the running time of a compiler. -

Example-YACC (Yet Another Compiler-Compiler).

2) Scanner Generator:

-These generate lexical analyzers, normally from a specification based on regular expressions. -The basic organization of lexical analyzers is based on finite automation.

3) Syntax-Directed Translation:

-These produce routines that walk the parse tree and as a result generate intermediate code. -Each translation is defined in terms of translations at its neighbor nodes in the tree.

4) Automatic Code Generators:

-It takes a collection of rules to translate intermediate language into machine language. The rules must include sufficient details to handle different possible access methods for data.

5) Data-Flow Engines:

-It does code optimization using data-flow analysis, that is, the gathering of information about how values are transmitted from one part of a program to each other part.

1.7 LEXICAL ANALYSIS

Lexical analysis is the process of converting a sequence of characters into a sequence of tokens. A program or function which performs lexical analysis is called a lexical analyzer or scanner. A lexer often exists as a single function which is called by a parser or another function.

1.7.1 THE ROLE OF THE LEXICAL ANALYZER

The lexical analyzer is the first phase of a compiler. Its main task is to read the input characters and produce as output a sequence of tokens that the parser uses for syntax analysis.



Figure 1.10:Role of Lexical Analyzer

Upon receiving a "get next token" command from the parser, the lexical analyzer reads input characters until it can identify the next token.

1.7.2 ISSUES OF LEXICAL ANALYZER

There are three issues in lexical analysis:

- 1. To make the design simpler.
- 2. To improve the efficiency of the compiler.
- 3. To enhance the computer portability.

1.7.3 TOKENS

A token is a string of characters, categorized according to the rules as a symbol (e.g., IDENTIFIER, NUMBER, COMMA). The process of forming tokens from an input stream of characters is called **tokenization**.

A token can look like anything that is useful for processing an input text stream or text file. Consider this expression in the C programming 10language: sum=3+2;

Lexeme	Token type
sum	Identifier
=	Assignment operator
3	Number
+	Addition operator
2	Number
	End of statement

Table 1.1:Tokens

LEXEME:

Collection or group of characters forming tokens is called Lexeme.

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.

1.7.4 Attributes for Tokens

Some tokens have attributes that can be passed back to the parser. The lexical analyzer collects information about tokens into their associated attributes. The attributes influence the translation of tokens.

i) Constant : value of the constant

ii) Identifiers: pointer to the corresponding symbol table entry.

1.7.5 ERROR RECOVERY STRATEGIES IN LEXICAL ANALYSIS:

The following are the error-recovery actions in lexical analysis:

1)Deleting an extraneous character.

2) Inserting a missing character.

3)Replacing an incorrect character by a correct character.

4)Transforming two adjacent characters.

5)Panic mode recovery: Deletion of successive characters from the token until

error is resolved.

<u>1.8 INPUT BUFFERING</u>

We often have to look one or more characters beyond the next lexeme before we can be sure we have the right lexeme. As characters are read from left to right, each character is stored in the buffer to form a meaningful token as shown below:



We introduce a two-buffer scheme that handles large look aheads safely. We then consider an improvement involving "sentinels" that saves time checking for the ends of buffers.

1.8.1 BUFFER PAIRS

A buffer is divided into two N-character halves, as shown below

$$:: E :: = :: M : * C : * :: * : 2 : eof$$

$$lexeme_beginning$$
forward

Figure 1.12: Buffer Pair

Each buffer is of the same size N, and N is usually the number of characters on one disk block. E.g., 1024 or 4096 bytes. Using one system read command we can read N characters into a buffer. If fewer than N characters remain in the input file, then a special character, represented by **eof**, marks the end of the source file. Two pointers to the input are maintained:

- 1. Pointer lexeme_beginning, marks the beginning of the current lexeme,
 - whose extent we are attempting to determine.
- 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.

The string of characters between the two pointers is the current lexeme. After the lexeme is recorded as an attribute value of a token returned to the parser, lexeme_beginning is set to the character immediately after the lexeme just found.

Advancing forward pointer:

Advancing forward pointer requires that we first test whether we have reached the end of one of the buffers, and if so, we must reload the other buffer from the input, and move forward to the beginning of the newly loaded buffer. If the end of second buffer is reached, we must again reload the first buffer with input and the pointer wraps to the beginning of the buffer.

Code to advance forward pointer:

```
if forward at end of first half then begin
    reload second half;
    forward := forward + 1
end
else if forward at end of second half then
    begin reload second half;
    move forward to beginning of first half
end
else forward := forward + 1;
```

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.

The sentinel arrangement is as shown below:



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.

Code to advance forward pointer:

```
forward : = forward + 1; if
forward ↑ = eof then begin
if forward at end of first half then begin
reload second half;
forward := forward + 1
end
else if forward at end of second half then
begin reload first half;
move forward to beginning of first
half end
else /* eof within a buffer signifying end of input
*/ terminate lexical analysis
end
```

1.9 SPECIFICATION OF TOKENS

There are 3 specifications of tokens:1) Strings2) Language3)Regular expression

Strings and Languages

An **alphabet** or character class is a finite set of symbols.

A string over an alphabet is a finite sequence of symbols drawn from that alphabet.

A language is any countable set of strings over some fixed alphabet.

In language theory, the terms "sentence" and "word" are often used as synonyms for "string." The length of a string s, usually written |s|, is the number of occurrences of symbols in s. For example, banana is a string of length six. The empty string, denoted ε , is the string of length zero.

Operations on strings

The following string-related terms are commonly used:

- 2. A **suffix** of string s is any string obtained by removing zero or more symbols from the beginning of s. For example, nana is a suffix of banana.
- 3. A **substring** of s is obtained by deleting any prefix and any suffix from s. For example, nan is a substring of banana.
- 4. The **proper prefixes**, suffixes, and substrings of a string s are those prefixes, suffixes, and substrings, respectively of s that are not ε or not equal to s itself.
- 5. A subsequence of s is any string formed by deleting zero or more not necessarily consecutive positions of s. For example, baan is a subsequence of banana.

Operations on languages:

The following are the operations that can be applied to languages:

1.Union
 2.Concatenation
 3.Kleene closure
 4.Positive closure

The following example shows the operations on strings: Let $L=\{0,1\}$ and $S=\{a,b,c\}$

- 1. Union : $L U S = \{0, 1, a, b, c\}$
- 2. Concatenation : L.S= $\{0a, 1a, 0b, 1b, 0c, 1c\}$
- 3. Kleene closure : L ={ ϵ , 0, 1, 00....}
- 4. Positive closure : $L^{+}=\{0,1,00....\}$

Regular Expressions

Each regular expression r denotes a language L(r). Here are the rules that define the regular expressions over some alphabet Σ and the languages that those expressions denote:

- 1. ε is a regular expression, and L(ε) is { ε }, that is, the language whose sole member is the empty string.
- 2. If 'a' is a symbol in Σ , then 'a' is a regular expression, and $L(a) = \{a\}$, that is, the language with one string, of length one, with 'a' in its one position.
- 3. Suppose r and s are regular expressions denoting the languages L(r) and L(s). Then,
 - a) (r)|(s) is a regular expression denoting the language L(r) U L(s).
 - b) (r)(s) is a regular expression denoting the language L(r)L(s).
 - c) $(r)^*$ is a regular expression denoting $(L(r))^*$.
 - d) (r) is a regular expression denoting L(r).
 - 4. The unary operator * has highest precedence and is left associative.
 - 5. Concatenation has second highest precedence and is left associative.
 - 6. It has lowest precedence and is left associative.

Regular set

A language that can be defined by a regular expression is called a regular set. If two regular expressions r and s denote the same regular set, we say they are equivalent and write r = s. There are a number of algebraic laws for regular expressions that can be used to manipulate into equivalent forms.

For instance, r|s = s|r is commutative; r|(s|t)=(r|s)|t is associative.

Regular Definitions

Giving names to regular expressions is referred to as a Regular definition. If Σ is an alphabet of basic symbols, then a regular definition is a sequence of definitions of the form

```
\begin{array}{l} d_{l} \longrightarrow r \ _{1} \ d_{2} \\ \longrightarrow r_{2} \\ \dots \\ d_{n} \longrightarrow r_{n} \end{array}
```

- 1. Each di is a distinct name.
- 2. Each ri is a regular expression over the alphabet $\Sigma \cup \{d_1, d_2, \dots, d_{i-1}\}$.

Example: Identifiers is the set of strings of letters and digits beginning with a letter. Regular definition for this set:

 $\begin{array}{l} \text{letter} \rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid \\ \text{digit} \rightarrow 0 \mid 1 \mid \dots \mid 9 \\ \text{id} \rightarrow \text{letter} (\text{letter} \mid \text{digit}) * \end{array}$

Shorthands

Certain constructs occur so frequently in regular expressions that it is convenient to introduce notational shorthands for them.

1. One or more instances (+):

- The unary postfix operator + means " one or more instances of".
- If r is a regular expression that denotes the language L(r), then $(r)^+$ is a regular expression that denotes the language $(L(r))^+$
- Thus the regular expression a^+ denotes the set of all strings of one or more a's.
- The operator + has the same precedence and associativity as the operator $\hat{}$.

2. Zero or one instance (?):

- The unary postfix operator ? means "zero or one instance of".
- The notation r? is a shorthand for $r \mid \epsilon.$
- If 'r' is a regular expression, then (r)? is a regular expression that denotes the language L(r) U { ϵ }.

3. Character Classes:

- The notation [abc] where a, b and c are alphabet symbols denotes the regular expression $a \mid b \mid c$.
- Character class such as [a z] denotes the regular expression $a | b | c | d | \dots | z$.
- We can describe identifiers as being strings generated by the regular expression, [A–Za–z][A–Za–z0–9]*

Non-regular Set

A language which cannot be described by any regular expression is a non-regular set. Example: The set of all strings of balanced parentheses and repeating strings cannot be described by a regular expression. This set can be specified by a context-free grammar.

RECOGNITION OF TOKENS

Consider the following grammar fragment:

```
stmt \rightarrow if expr then stmt

|if expr then stmt else stmt

|\epsilon

expr \rightarrow term relop term

|term

term \rightarrow id
```

```
num
```

where the terminals if , then, else, relop, id and num generate sets of strings given by the following regular definitions:

If \rightarrow if Then \rightarrow then Else \rightarrow else Relop $\rightarrow <|<=|=|<>|>|>=$ Id \rightarrow letter(letter|digit) Num \rightarrow digit⁺(.digit⁺)?(E(+|-)?digit⁺)?

For this language fragment the lexical analyzer will recognize the keywords if, then, else, as well as the lexemes denoted by relop, id, and num. To simplify matters, we assume keywords are reserved; that is, they cannot be used as identifiers.

Transition diagrams

It is a diagrammatic representation to depict the action that will take place when a lexical analyzer is called by the parser to get the next token. It is used to keep track of information about the characters that are seen as the forward pointer scans the input.

Transition diagram for relational operators



Figure 1.14:Transition Diagram

1.10 A LANGUAGE FOR SPECIFYING LEXICAL ANALYZER

There is a wide range of tools for constructing lexical analyzers.

- ✤ Lex
- ✤ YACC

LEX

Lex is a computer program that generates lexical analyzers. Lex is commonly used with the yacc parser generator.

Creating a lexical analyzer

First, a specification of a lexical analyzer is prepared by creating a program lex.l in the Lex language. Then, lex.l is run through the Lex compiler to produce a C program lex.yy.c. Finally, lex.yy.c is run through the C compiler to produce an object program a.out, which is the lexical analyzer that transforms an input stream into a sequence of tokens.





Figure 1.15: Creating Lexical Analyzer

Lex Specification

A Lex program consists of three parts:

{ definitions }
%%
{ rules }
%%
{ user subroutines }

> **Definitions** include declarations of variables, constants, and regular definitions

- > **Rules** are statements of the form
 - p1 {action1}
 - p2 {action2} ...
 - pn {actionn}

where p_i is regular expression and action_i describes what action the lexical analyzer should take when pattern p_i matches a lexeme. Actions are written in C code.

User subroutines are auxiliary procedures needed by the actions. These can be compiled separately and loaded with the lexical analyzer.

YACC- YET ANOTHER COMPILER-COMPILER

Yacc provides a general tool for describing the input to a computer program. The Yacc user specifies the structures of his input, together with code to be invoked as each such structure is recognized. Yacc turns such a specification into a subroutine that handles the input process; frequently, it is convenient and appropriate to have most of the flow of control in the user's application handled by this subroutine.

1.11 FINITE AUTOMATA

Finite Automata is one of the mathematical models that consist of a number of states and edges. It is a transition diagram that recognizes a regular expression or grammar.

Types of Finite Automata

There are tow types of Finite Automata :

- Non-deterministic Finite Automata (NFA)
- Deterministic Finite Automata (DFA)

1.11.1 Non-deterministic Finite Automata

NFA is a mathematical model that consists of five tuples denoted by

 $M = \{Q_n, \Sigma, \delta, q_0, f_n\}$

- Q_n finite set of states
- Σ finite set of input symbols
- δ transition function that maps state-symbol pairs to set of states

qo – starting state

 f_n – final state

1.11.2 Deterministic Finite Automata

DFA is a special case of a NFA in which

- i) no state has an ε -transition.
- ii) there is at most one transition from each state on any input.

DFA has five tuples denoted by

 $\mathbf{M} = \{\mathbf{Q}_{d}, \boldsymbol{\Sigma}, \, \boldsymbol{\delta}, \, \mathbf{q}_{0}, \, \mathbf{f}_{d}\}$

Qd – finite set of states

- Σ finite set of input symbols
- δ transition function that maps state-symbol pairs to set of states
- qo starting state
- fd final state

1.11.3 Construction of DFA from regular expression

The following steps are involved in the construction of DFA from regular expression:

- i) Convert RE to NFA using Thomson's rules
- ii) Convert NFA to DFA
- iii) Construct minimized DFA

UNIT-II

CHAPTER-II

SYNTAX ANALYSIS AND RUNTIME ENVIRONMENT

2.1 SYNTAX ANALYSIS

Syntax analysis is the second phase of the compiler. It gets the input from the tokens and generates a syntax tree or parse tree.

Advantages of grammar for syntactic specification :

- 1. A grammar gives a precise and easy-to-understand syntactic specification of a programming language.
- 2. An efficient parser can be constructed automatically from a properly designed grammar.
- 3. A grammar imparts a structure to a source program that is useful for its translation into object code and for the detection of errors.
- 4. New constructs can be added to a language more easily when there is a grammatical description of the language.

2.1.1 THE ROLE OF PARSER

The parser or syntactic analyzer obtains a string of tokens from the lexical analyzer and verifies that the string can be generated by the grammar for the source language. It reports any syntax errors in the program. It also recovers from commonly occurring errors so that it can continue processing its input.

Position of parser in compiler model



Figure 2.1:Role of Parser

Functions of the parser:

5) It verifies the structure generated by the tokens based on the grammar.

6) It constructs the parse tree.

7) It reports the errors.

8) It performs error recovery. **Issues :**

Parser cannot detect errors such as:

- 1. Variable re-declaration
- 2. Variable initialization before use.
- 3. Data type mismatch for an operation.

The above issues are handled by Semantic Analysis phase.

Syntax error handling :

Programs can contain errors at many different levels. For example :

- 1. Lexical, such as misspelling a keyword.
- 2. Syntactic, such as an arithmetic expression with unbalanced parentheses.
- 3. Semantic, such as an operator applied to an incompatible operand.
- 4. Logical, such as an infinitely recursive call.

Functions of error handler:

- 1. It should report the presence of errors clearly and accurately.
- 2. It should recover from each error quickly enough to be able to detect subsequent errors.
- 3. It should not significantly slow down the processing of correct programs.

2.1.2 Error recovery strategies:

The different strategies that a parse uses to recover from a syntactic error are:

- 1. Panic mode
- 2. Phrase level
- 3. Error productions
- 4. Global correction

Panic mode recovery:

On discovering an error, the parser discards input symbols one at a time until a synchronizing token is found. The synchronizing tokens are usually delimiters, such as semicolon or **end**. It has the advantage of simplicity and does not go into an infinite loop. When multiple errors in the same statement are rare, this method is quite useful.

Phrase level recovery:

On discovering an error, the parser performs local correction on the remaining input that allows it to continue. Example: Insert a missing semicolon or delete an extraneous semicolon etc.

Error productions:

The parser is constructed using augmented grammar with error productions. If an error production is used by the parser, appropriate error diagnostics can be generated to indicate the erroneous constructs recognized by the input.

Global correction:

Given an incorrect input string x and grammar G, certain algorithms can be used to find a parse tree for a string y, such that the number of insertions, deletions and changes of tokens is as small as possible. However, these methods are in general too costly in terms of time and space.

2.2 CONTEXT-FREE GRAMMARS

A Context-Free Grammar is a quadruple that consists of **terminals**, **non-terminals**, **start symbol** and **productions**.

Terminals : These are the basic symbols from which strings are formed.

Non-Terminals : These are the syntactic variables that denote a set of strings. These help to define the language generated by the grammar.

Start Symbol : One non-terminal in the grammar is denoted as the "Start-symbol" and the set of strings it denotes is the language defined by the grammar.

Productions : It specifies the manner in which terminals and non-terminals can be combined to form strings. Each production consists of a non-terminal, followed by an arrow, followed by a string of non-terminals and terminals.

Example of context-free grammar: The following grammar defines simple arithmetic expressions:

```
expr \rightarrow expr op expr
expr \rightarrow (expr)
expr \rightarrow - expr
expr \rightarrow id op
\rightarrow + op \rightarrow -
op \rightarrow *
op \rightarrow /
op \rightarrow \uparrow
```

In this grammar,

- $\mathbf{id} + \mathbf{-} * / \uparrow ()$ are terminals.
- *expr*, *op* are non-terminals.
- *expr* is the start symbol.
- Each line is a production.

2.2.1 Derivations:

Two basic requirements for a grammar are :

- \Box To generate a valid string.
- \Box To recognize a valid string.

Derivation is a process that generates a valid string with the help of grammar by replacing the non-terminals on the left with the string on the right side of the production.

Example : Consider the following grammar for arithmetic expressions :

 $E \rightarrow E+E |E*E|(E)| - E | id$

- $\Box \quad E \to E$
- $\Box \quad E \rightarrow (E)$
- $\Box \quad E \rightarrow (E+E)$
- \Box $E \rightarrow (id+E)$
- \Box E \rightarrow (id+id)

In the above derivation,

- ⁽²⁾ E is the start symbol.
- ^(b) (id+id) is the required sentence (only terminals).
- O Strings such as E, -E, -(E), . . . are called sentinel forms.

Types of derivations:

The two types of derivation are:

Left most derivation Right most derivation.

- In leftmost derivations, the leftmost non-terminal in each sentinel is always chosen first for replacement.
- In rightmost derivations, the rightmost non-terminal in each sentinel is always chosen first for replacement.

Example:

Given grammar G : $E \rightarrow E+E |E*E|(E)| - E |id$

Sentence to be derived : -(id+id)

LEFTMOST DERIVATION RIGHTMOST DERIVATION

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

> String that appear in leftmost derivation are called **left sentinel forms.**

> String that appear in rightmost derivation are called **right sentinel forms.**

Sentinels:

Given a grammar G with start symbol S, if $S \rightarrow \alpha$, where α may contain non-terminals or terminals, then α is called the sentinel form of G.

Yield or frontier of tree:

Each interior node of a parse tree is a non-terminal. The children of node can be a terminal or non-terminal of the sentinel forms that are read from left to right. The sentinel form in the parse tree is called **yield** or **frontier** of the tree.

Ambiguity:

A grammar that produces more than one parse for some sentence is said to be **ambiguous** grammar.

Example : Given grammar G : $E \rightarrow E+E | E*E | (E) | -E | id$

The sentence id+id*id has the following two distinct leftmost derivations:

$E \rightarrow E + E$	$E \rightarrow E^* E$
$E \rightarrow id + E$	$E \rightarrow E + E * E$
$E \rightarrow id + E * E$	$E \rightarrow id + E * E$
$E \rightarrow id + id * E$	$E \rightarrow id + id * E$
$E \rightarrow id + id * id$	$E \rightarrow id + id * id$

The two corresponding parse trees are :



2.3 WRITING A GRAMMAR

There are four categories in writing a grammar :

- 1. Regular Expression Vs Context Free Grammar
- □ Eliminating ambiguous grammar.
- □ Eliminating left-recursion
- \Box Left-factoring.

Each parsing method can handle grammars only of a certain form hence, the initial grammar may have to be rewritten to make it parsable.

2.3.1 Regular Expressions vs. Context-Free Grammars:

REGULAR EXPRESSION	CONTEXT-FREE GRAMMAR
It is used to describe the tokens of programming languages.	It consists of a quadruple where $S \rightarrow$ start symbol, $P \rightarrow$ production, $T \rightarrow$ terminal, $V \rightarrow$ variable or non- terminal.
It is used to check whether the given input is valid or not using transition diagram.	It is used to check whether the given input is valid or not using derivation .
The transition diagram has set of states and edges.	The context-free grammar has set of productions.
It has no start symbol.	It has start symbol.
It is useful for describing the structure of lexical constructs such as identifiers, constants, keywords, and so forth.	It is useful in describing nested structures such as balanced parentheses, matching begin-end's and so on.

- > The lexical rules of a language are simple and RE is used to describe them.
- Regular expressions provide a more concise and easier to understand notation for tokens than grammars.
- > Efficient lexical analyzers can be constructed automatically from RE than from grammars.
- Separating the syntactic structure of a language into lexical and nonlexical parts provides a convenient way of modularizing the front end into two manageable-sized components.

Eliminating ambiguity:

Ambiguity of the grammar that produces more than one parse tree for leftmost or rightmost derivation can be eliminated by re-writing the grammar.

Consider this example, G: $stmt \rightarrow if expr$ then $stmt \mid if expr$ then stmt else $stmt \mid other$

This grammar is ambiguous since the string **if E1 then if E2 then S1 else S2** has the following two parse trees for leftmost derivation :



To eliminate ambiguity, the following grammar may be used:

stmt → *matched_stmt* | *unmatched_stmt*

 $matched_stmt \rightarrow if expr then matched_stmt else matched_stmt | other$

 $unmatched_stmt \rightarrow if expr then stmt | if expr then matched_stmt else unmatched_stmt$

2.3.2 Eliminating Left Recursion:

A grammar is said to be *left recursive* if it has a non-terminal A such that there is a derivation $A=>A\alpha$ for some string α . Top-down parsing methods cannot handle left-recursive grammars. Hence, left recursion can be eliminated as follows:

A → βA'A' → αA' | ε without

If there is a production $A \rightarrow A\alpha \mid \beta$ it can be replaced with a sequence of two productions

changing the set of strings derivable from A.

Example : Consider the following grammar for arithmetic expressions:

 $E \rightarrow E+T | T$

 $T \rightarrow T^*F | F$

 $F \rightarrow (E) | id$

First eliminate the left recursion for E

as $E \rightarrow TE'$

 $E' \rightarrow +TE' | \epsilon$

Then eliminate for T

as $T \rightarrow FT'$

 $T' \rightarrow *FT' \mid \epsilon$

Thus the obtained grammar after eliminating left recursion

is $E \rightarrow TE'$

 $\mathrm{E}^{\prime} \rightarrow +\mathrm{TE}^{\prime} \mid \epsilon$

 $T \rightarrow FT'$

 $T' \rightarrow *FT' | \epsilon$

 $F \rightarrow (E) | id$

Algorithm to eliminate left recursion:

- 1. Arrange the non-terminals in some order $A_1, A_2 \dots A_n$.
- 2. for i := 1 to n do begin
 - for j := 1 to i-1 do begin

```
replace each production of the form A_i \rightarrow A_j \gamma by
the productions A_i \rightarrow \delta_1 \gamma | \delta_2 \gamma | \dots | \delta_k \gamma
where A_j \rightarrow \delta_1 | \delta_2 | \dots | \delta_k are all the current A_j-productions;
```

end

eliminate the immediate left recursion among the Ai-productions

end

2.3.3 Left factoring:

Left factoring is a grammar transformation that is useful for producing a grammar suitable for predictive parsing. When it is not clear which of two alternative productions to use to expand a non-terminal A, we can rewrite the A-productions to defer the decision until we have seen enough of the input to make the right choice.

If there is any production $A \rightarrow \alpha \beta_1 \mid \alpha \beta_2$, it can be rewritten as

 $\begin{array}{l} \mathbf{A} \rightarrow \boldsymbol{\alpha} \mathbf{A}^{\prime} \\ \mathbf{A}^{\prime} \rightarrow \boldsymbol{\beta} \mathbf{1} \mid \boldsymbol{\beta} \mathbf{2} \end{array}$

Consider the grammar , $G:S \rightarrow iEtS \mid iEtSeS \mid a$ $E \rightarrow b$

Left factored, this grammar becomes

 $S \rightarrow iEtSS' \mid a$ S' $\rightarrow eS \mid \epsilon$ E $\rightarrow b$

PARSING

It is the process of analyzing a continuous stream of input in order to determine its grammatical structure with respect to a given formal grammar.

Parse tree:

Graphical representation of a derivation or deduction is called a parse tree. Each interior node of the parse tree is a non-terminal; the children of the node can be terminals or non-terminals.

Types of parsing:

- 1. Top down parsing
- 2. Bottom up parsing
- Top-down parsing : A parser can start with the start symbol and try to transform it to the input string. Example : LL Parsers.
- Bottom-up parsing : A parser can start with input and attempt to rewrite it into the start symbol. Example : LR Parsers.

2.4 TOP-DOWN PARSING

It can be viewed as an attempt to find a left-most derivation for an input string or an attempt to construct a parse tree for the input starting from the root to the leaves.
Types of top-down parsing :

- 1. Recursive descent parsing
- 2. Predictive parsing

2.4.1 RECURSIVE DESCENT PARSING

- Recursive descent parsing is one of the top-down parsing techniques that uses a set of recursive procedures to scan its input.
- This parsing method may involve **backtracking**, that is, making repeated scans of the input.

Example for backtracking :

Consider the grammar $G : S \rightarrow cAd$ $A \rightarrow ab | a$ and the input string w=cad.

The parse tree can be constructed using the followingtop-down approach :

Step1:

Initially create a tree with single node labeled S. An input pointer points to 'c', the first symbol of w. Expand the tree with the production of S.



Step2:

The leftmost leaf 'c' matches the first symbol of w, so advance the input pointer to the second symbol of w 'a' and consider the next leaf 'A'. Expand A using the first alternative.



Step3:

The second symbol 'a' of w also matches with second leaf of tree. So advance the input pointer to third symbol of w 'd'. But the third leaf of tree is b which does not match with the input symbol **d**.

Hence discard the chosen production and reset the pointer to second position. This is called **backtracking**.

Step4:

Now try the second alternative for A.



Now we can halt and announce the successful completion of parsing.

Example for recursive decent parsing:

A left-recursive grammar can cause a recursive-descent parser to go into an infinite loop. Hence, **elimination of left-recursion** must be done before parsing.

Consider the grammar for arithmetic expressions

 $E \rightarrow E+T | T$

 $T \to T^*F \mid F$

 $F \rightarrow (E) | id$

After eliminating the left-recursion the grammar

becomes, $E \rightarrow TE'$

 $E' \rightarrow +TE' \mid \epsilon$

 $T \rightarrow FT'$

 $T' \rightarrow *FT' \mid \epsilon$

 $F \rightarrow (E) | id$

Now we can write the procedure for grammar as follows:

Recursive procedure:

```
Procedure E()
begin
T();
EPRIME();
```

Procedure EPRIME() begin If input symbol='+' then ADVANCE(); T(); EPRIME(); end Procedure T() begin F(); TPRIME(); end Procedure TPRIME() begin If input symbol='*' then ADVANCE(); F(); TPRIME(); end Procedure F() begin If input-symbol='id' then ADVANCE(); else if input-symbol='(' then ADVANCE(); E(); else if input-symbol=')' then ADVANCE(); end else ERROR(); **Stack implementation:**

To recognize input **id+id*id** :

Table 2.1: Stack implementation

PROCEDURE	INPUT STRING
Ε()	<u>id</u> +id*id
Τ()	<u>id</u> +id*id
F()	<u>id</u> +id*id
ADVANCE()	id <u>+</u> id*id

	id <u>+</u> id*id
TPRIME()	id <u>+</u> id*id
EPRIME()	id+ <u>id</u> *id
ADVANCE()	id+ <u>id</u> *id
Τ()	id+ <u>id</u> *id
F()	id+id <u>*</u> id
ADVANCE()	id+id <u>*</u> id
TPRIME()	id+id* <u>id</u>
ADVANCE()	id+id* <u>id</u>
F()	id+id* <u>id</u>
ADVANCE()	id+id* <u>id</u>
TPRIME()	i

2.4.2 PREDICTIVE PARSING

.

- 2.5 Predictive parsing is a special case of recursive descent parsing where no backtracking is required.
- 2.6 The key problem of predictive parsing is to determine the production to be applied for a non-terminal in case of alternatives.

Non-recursive predictive parser



The table-driven predictive parser has an input buffer, stack, a parsing table and an output stream.

Input buffer:

It consists of strings to be parsed, followed by \$ to indicate the end of the input string.

Stack:

It contains a sequence of grammar symbols preceded by \$ to indicate the bottom of the stack. Initially, the stack contains the start symbol on top of \$.

Parsing table:

It is a two-dimensional array M[A, a], where 'A' is anon-terminal and 'a' is aterminal.

Predictive parsing program:

The parser is controlled by a program that considers X, the symbol on top of stack, and a, the current input symbol. These two symbols determine the parser action. There are three possibilities:

- 1. If X = a =\$, the parser halts and announces successful completion of parsing.
- 2. If $X = a \neq \$$, the parser pops X off the stack and advances the input pointer to the next input symbol.
- 3. If X is a non-terminal, the program consults entry M[X, a] of the parsing table M. This entry will either be an X-production of the grammar or an error entry. If M[X, a] = {X → UVW}, the parser replaces X on top of the stack by WVU. If M[X, a] = error, the parser calls an error recovery routine.

Algorithm for nonrecursive predictive parsing:

Input : A string *w* and a parsing table *M* for grammar *G*.

Output : If w is in L(G), a leftmost derivation of w; otherwise, an error indication.

Method : Initially, the parser has S on the stack with *S*, the start symbol of *G* on top, and w in the input buffer. The program that utilizes the predictive parsing table *M* to produce a parse for the input is as follows:

set *ip* to point to the first symbol of w\$; **repeat** let X be the top stack symbol and a the symbol pointed to by *ip*; **if** X is a terminal or \$ **then if** X = a **then** pop X from the stack and advance *ip* **else** error() **else** /* X is a non-terminal */ **if** $M[X, a] = X \rightarrow Y_1Y_2 \dots Y_k$ **then begin** pop X from the stack; push Y_k , Y_{k-1} , ..., Y_1 onto the stack, with Y_1 on top; output the production $X \rightarrow Y_1 Y_2 \dots Y_k$ end else *error*() until X = /* stack is empty */

Predictive parsing table construction:

The construction of a predictive parser is aided by two functions associated with a grammar G :

1. FIRST

2. FOLLOW **Rules for first():**

- 1. If X is terminal, then FIRST(X) is $\{X\}$.
- 2. If $X \rightarrow \varepsilon$ is a production, then add ε to FIRST(*X*).
- 3. If *X* is non-terminal and $X \rightarrow a\alpha$ is a production then add *a* to FIRST(X).
- 4. If X is non-terminal and X → Y1 Y2...Yk is a production, then place a in FIRST(X) if for some i, a is in FIRST(Yi), and ε is in all of FIRST(Y1),...,FIRST(Yi-1); that is, Y1,....Yi-1 => ε. If ε is in FIRST(Yj) for all j=1,2,...,k, then add ε to FIRST(X).

Rules for follow():

- 1. If *S* is a start symbol, then FOLLOW(*S*) contains \$.
- 2. If there is a production $A \rightarrow \alpha B\beta$, then everything in FIRST(β) except ε is placed in follow(*B*).
- 3. If there is a production $A \rightarrow \alpha B$, or a production $A \rightarrow \alpha B\beta$ where FIRST(β) contains ε , then everything in FOLLOW(*A*) is in FOLLOW(*B*).

Algorithm for construction of predictive parsing table:

Input : Grammar *G*

Output : Parsing table *M*

Method :

- 1. For each production $A \rightarrow \alpha$ of the grammar, do steps 2 and 3.
- 2. For each terminal *a* in FIRST(α), add $A \rightarrow \alpha$ to *M*[*A*, *a*].
- 3. If ε is in FIRST(α), add A $\rightarrow \alpha$ to M[A, b] for each terminal b in FOLLOW(A). If ε is in FIRST(α) and \$\$ is in FOLLOW(A), add A $\rightarrow \alpha$ to M[A, \$\$].
- 4. Make each undefined entry of *M* be error.

Example:

Consider the following grammar :

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

After eliminating left-recursion the grammar is

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

First() :

 $FIRST(E) = \{ (, id \}$

FIRST(E') ={+, ε }

 $FIRST(T) = \{ (, id \}$

FIRST(T') = {*, ε }

 $FIRST(F) = \{ (, id \}$

Follow():

FOLLOW(E) ={ \$,) }

FOLLOW(E') = { \$,) }

FOLLOW(T) ={ +, \$,) }

FOLLOW(T') = { +, \$,) }

FOLLOW(F) = {+, *, \$, } }
Table 2.2: Predictive parsEr

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

stack	Input	Output
\$E	id+id*id \$	
\$E'T	id+id*id \$	$E \rightarrow TE'$
\$E'T'F	id+id*id \$	$T \rightarrow FT'$
\$E'T'id	id+id*id \$	$F \rightarrow id$
\$E'T'	+id*id \$	
\$E'	+id*id \$	$T' \rightarrow \epsilon$
\$E'T+	+id*id \$	$E' \rightarrow +TE'$
\$E'T	id*id \$	
\$E'T'F	id*id \$	$T \rightarrow FT'$
\$E'T'id	id*id \$	$F \rightarrow id$
\$E'T'	*id \$	
\$E'T'F*	*id \$	$T' \rightarrow *FT'$
\$E'T'F	id \$	
\$E'T'id	id \$	$F \rightarrow id$
\$E'T'	\$	
\$E'	\$	$T' \rightarrow \epsilon$
\$	\$	$E' \rightarrow \varepsilon$

 Table 2.3:Stack implen

LL(1) grammar:

The parsing table entries are single entries. So each location has not more than one entry. This type of grammar is called LL(1) grammar.

Consider this following grammar:

$$\begin{split} S &\rightarrow i E t S \mid i E t S e S \mid a \\ E &\rightarrow b \end{split}$$

After eliminating left factoring, we have

$$\begin{split} S &\to i EtSS' \mid a \\ S' &\to eS \mid \epsilon \\ E &\to b \end{split}$$

To construct a parsing table, we need FIRST()and FOLLOW() for all the non-terminals.

FIRST(S) ={ i, a }
FIRST(S') = { e, ε }
FIRST(E) ={ b}
FOLLOW(S) ={ \$,e }

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

$FOLLOW(E) = \{t\}$

Table 2.4: Parsing table

NON- TERMINAL	а	b	e	i	t	\$
S	$S \rightarrow a$			$S \rightarrow iEtSS'$		
S'			$S' \rightarrow eS$ $S' \rightarrow \varepsilon$			S'→ ε
Е		$E \rightarrow b$				

Since there are more than one production, the grammar is not LL(1) grammar.

Actions performed in predictive parsing:

- 1. Shift
- 2. Reduce
- 3. Accept
- 4. Error

Implementation of predictive parser:

- 1. Elimination of left recursion, left factoring and ambiguous grammar.
- 2. Construct FIRST() and FOLLOW() for all non-terminals.
- 3. Construct predictive parsing table.
- 4. Parse the given input string using stack and parsing table.

2.5 BOTTOM-UP PARSING

Constructing a parse tree for an input string beginning at the leaves and going towards the root is called bottom-up parsing.

A general type of bottom-up parser is a **shift-reduce parser**.

2.5.1 SHIFT-REDUCE PARSING

Shift-reduce parsing is a type of bottom-up parsing that attempts to construct a parse tree for an input string beginning at the leaves (the bottom) and working up towards the root (the top).

Example:

Consider the grammar: $S \rightarrow aABe$ $A \rightarrow Abc \mid b$ $B \rightarrow d$ The sentence to be recognized is **abbcde.**

REDUCTION (LEFTMOST)

RIGHTMOST DERIVATION

abbcde	$(A \rightarrow b)$	$\mathbf{S} \rightarrow \mathbf{a}\mathbf{A}\mathbf{B}\mathbf{e}$
a Abc de	$(A \rightarrow Abc)$	\rightarrow aAde
aA d e	$(B \rightarrow d)$	\rightarrow aAbcde
aABe	$(S \rightarrow aABe)$	\rightarrow abbcde
S		

The reductions trace out the right-most derivation in reverse.

Handles:

A handle of a string is a substring that matches the right side of a production, and whose reduction to the non-terminal on the left side of the production represents one step along the reverse of a rightmost derivation.

Example:

Consider the grammar:

 $E \rightarrow E+E$ $E \rightarrow E*E$ $E \rightarrow (E)$ $E \rightarrow id$

And the input string id1+id2*id3

The rightmost derivation is :

 $E \rightarrow \underline{E+E}$ $\rightarrow E+\underline{E*E}$ $\rightarrow E+E*\underline{id3}$ $\rightarrow E+\underline{id2}*id3$ $\rightarrow \underline{id1}+id2*id3$

In the above derivation the underlined substrings are called handles.

Handle pruning:

A rightmost derivation in reverse can be obtained by "handle pruning".

(i.e.) if *w* is a sentence or string of the grammar at hand, then $w = \gamma_n$, where γ_n is the *n*th right-sentinel form of some rightmost derivation.

Stack	Input	Action
\$	id1+id2*id3 \$	shift
\$ id1	+id2*id3 \$	reduce by $E \rightarrow id$
\$ E	+id2*id3 \$	shift
\$ E+	id2*id3 \$	shift
\$ E+id2	*id3 \$	reduce by $E \rightarrow id$
\$ E+E	*id3 \$	shift
\$ E+E*	id3 \$	shift
\$ E+E*id3	\$	reduce by $E \rightarrow id$
\$ E+E*E	\$	reduce by $E \rightarrow E * E$
\$ E+E	\$	reduce by $E \rightarrow E + E$
\$ E	\$	accept

Table 2.5: Stack implementation of shift-reduce parsing :

Actions in shift-reduce parser:

- shift The next input symbol is shifted onto the top of the stack.
- reduce The parser replaces the handle within a stack with a non-terminal.
- accept The parser announces successful completion of parsing.
- error The parser discovers that a syntax error has occurred and calls an error recovery routine.

Conflicts in shift-reduce parsing:

There are two conflicts that occur in shift shift-reduce parsing:

- 1. Shift-reduce conflict: The parser cannot decide whether to shift or to reduce.
- 2. Reduce-reduce conflict: The parser cannot decide which of several reductions to make.

1. Shift-reduce conflict:

Example:

Consider the grammar:

 $E \rightarrow E + E | E^*E | id and input id+id^*id$

Stack	Input	Action	Stack	Input	Action
\$ E+E	*id \$	Reduce by $E \rightarrow E + E$	\$E+E	*id \$	Shift
\$ E	*id \$	Shift	\$E+E*	id \$	Shift
\$ E*	id \$	Shift	\$E+E*id	\$	Reduce by $E \rightarrow id$
\$ E*id	\$	Reduce by E→id	\$E+E*E	\$	Reduce by $E \rightarrow E^*E$
\$ E*E	\$	Reduce by $E \rightarrow E^*E$	\$E+E	\$	Reduce by $E \rightarrow E^*E$
\$ E			\$E		

.

2. <u>Reduce-reduce conflict:</u>

Consider the grammar:

 $\label{eq:matrix} \begin{array}{l} M \rightarrow R{+}R \; |R{+}c \; |R \\ R \rightarrow c \\ \mbox{and input } c{+}c \end{array}$

Stack	Input	Action	Stack	Input	Action
\$	c+c \$	Shift	\$	c+c \$	Shift
\$ c	+c \$	Reduce by R→c	\$ c	+c \$	Reduce by R→c
\$ R	+c \$	Shift	\$ R	+c \$	Shift
\$ R+	c \$	Shift	\$ R+	c \$	Shift
\$ R+c	\$	Reduce by R→c	\$ R+c	\$	Reduce by M→R+c
\$ R+R	\$	Reduce by M→R+R	\$ M	\$	
\$ M	\$				

Viable prefixes:

- > α is a viable prefix of the grammar if there is w such that αw is a right sentinel form.
- The set of prefixes of right sentinel forms that can appear on the stack of a shift-reduce parser are called viable prefixes.
- > The set of viable prefixes is a regular language.

2.5.2 OPERATOR-PRECEDENCE PARSING

An efficient way of constructing shift-reduce parser is called operator-precedence parsing.

Operator precedence parser can be constructed from a grammar called Operator-grammar. These grammars have the property that no production on right side is ε or has two adjacent non-terminals.

Example:

Consider the grammar:

$$\begin{split} & E \rightarrow EAE \mid (E) \mid -E \mid id \\ & A \rightarrow + \mid - \mid * \mid / \mid \uparrow \end{split}$$

Since the right side EAE has three consecutive non-terminals, the grammar can be written as follows:

 $E \rightarrow E+E | E-E | E*E | E/E | E\uparrow E | -E | id$ Operator precedence relations:

There are three disjoint precedence relations

namely < ' - less than = - equal to '> - greater than The relations give the followingmeaning:

a < b - a yields precedence to b

a=b – a has the same precedence as b

a' > b - a takes precedence over b

Rules for binary operations:

 \Box If operator θ_1 has higher precedence than operator $\theta_2,$ then

make $\theta_1 \rightarrow \theta_2$ and $\theta_2 < \theta_1$

 \Box If operators θ_1 and θ_2 , are of equal precedence, then make

 $\theta_1 \rightarrow \theta_2$ and $\theta_2 \rightarrow \theta_1$ if operators are left associative

 $\theta_1 < \theta_2$ and $\theta_2 < \theta_1$ if right associative

 \Box Make the following for all operators θ :

Also make

$$(=), (<`(,)`>), (<`id,id`>), $<`id,id`>$, $<`(,)`>$$$

Example:

Operator-precedence relations for the grammar

 $E \rightarrow E+E |E-E|E*E|E/E |E/E|(E)|-E |id is given in the following table assuming$

- 1. \uparrow is of highest precedence and right-associative
- 2. * and / are of next higher precedence and left-associative, and
- 3. + and are of lowest precedence and left-associative

Note that the **blanks** in the table denote error entries.

	+	-	*	/	↑	id	()	\$
+	->	->	<.	<-	<.	<:	<	->	->
-	·>	·>	<`	<`	<.	<:	<.	->	->
*	·>	·>	·>	·>	<i>.</i>		`<	·>	·>
/	·>	·>	·>	->	<.	<:	<:	->	->
\uparrow	·>	·>	·>	·>	<	<	`<	·>	·>
id	·>	->	·>	->	·>			->	->
(<i>.</i>	<i>.</i>	<i><</i> .	<i><</i> .	<i>.</i>	, Ś		=	
)	·>	->	·>	->	·>			->	->
\$	<.	<i>.</i>	<i>.</i>	<.	<i>.</i>	<i>.</i>	<i>.</i>		

TABLE : Operator-precedence relations

Operator precedence parsing algorithm:

Input : An input string *w* and a table of precedence relations.

Output : If *w* is well formed, a *skeletal* parse tree ,with a placeholder non-terminal E labeling all interior nodes; otherwise, an error indication.

Method : Initially the stack contains and the input buffer the string *w* . To parse, we execute the following program :

(1)Set *ip* to point to the first symbol of *w*\$;

2. repeat forever

- 3. **if** \$ is on top of the stack and *ip* points to \$ **then**
- 4. return

else begin

- 4. let *a* be the topmost terminal symbol on the stack and let *b* be the symbol pointed to by *ip*;
- 5. **if** a < b or a = b **then begin**
- 6. push b onto the stack;
- 7. advance *ip* to the next input symbol;

end;

(9) else if a > b then /*reduce*/

- (10) repeat
- (11) pop the stack
- (12) **until** the top stack terminal is related by <⁻ to the terminal most recently popped
- (13) else error()

end

Stack implementation of operator precedence parsing:

Operator precedence parsing uses a stack and precedence relation table for its implementation of above algorithm. It is a shift-reduce parsing containing all four actions shift, reduce, accept and error.

The initial configuration of an operator precedence parsing is STACK\$

where w is the input string to be parsed. **Example:**

Consider the grammar $E \rightarrow E+E | E-E | E*E | E/E | E\uparrow E | (E) | id.$ Input string is **id+id*id**. The implementation is as follows:

STACK	INPUT	COMMENT
\$	<· id+id*id \$	shift id
\$ id	·> +id*id \$	pop the top of the stack id
\$	<- +id*id \$	shift +
\$ +	<· id*id \$	shift id
\$ +id	·> *id \$	pop id
\$+	<· *id \$	shift *
\$+*	<· id \$	shift id
\$ + * id	·> \$	pop id
\$ + *	·> \$	pop *
\$ +	·> \$	pop +
\$	\$	accept

Advantages of operator precedence parsing:

- 3. It is easy to implement.
- 4. Once an operator precedence relation is made between all pairs of terminals of a grammar , the grammar can be ignored. The grammar is not referred anymore during implementation.

Disadvantages of operator precedence parsing:

- 2. It is hard to handle tokens like the minus sign (-) which has two different precedence.
- 3. Only a small class of grammar can be parsed using operator-precedence parser.

INPUT w\$

2.6 LR PARSERS

An efficient bottom-up syntax analysis technique that can be used to parse a large class of CFG is called LR(k) parsing. The 'L' is for left-to-right scanning of the input, the 'R' for constructing a rightmost derivation in reverse, and the 'k' for the number of input symbols. When 'k' is omitted, it is assumed to be 1.

Advantages of LR parsing:

- ✓ It recognizes virtually all programming language constructs for which CFG can be written.
- ✓ It is an efficient non-backtracking shift-reduce parsing method.
- ✓ A grammar that can be parsed using LR method is a proper superset of a grammar that can be parsed with predictive parser.
- ✓ It detects asyntactic error as soon as possible.

Drawbacks of LR method:

It is too much of work to construct a LR parser by hand for a programming language grammar. A specialized tool, called a LR parser generator, is needed. Example: YACC.

Types of LR parsing method:

- 1. SLR- Simple LR
 - Easiest to implement, least powerful.
- 2. CLR- Canonical LR
 - Most powerful, most expensive.
- 3. LALR- Look-Ahead LR
 - Intermediate in size and cost between the other two methods.

The LR parsing algorithm:

The schematic form of an LR parser is as follows:



It consists of : an input, an output, a stack, a driver program, and a parsing table that has two parts (*action* and *goto*).

- The driver program is the same for all LR parser.
- The parsing program reads characters from an input buffer one at a time.
- The program uses a stack to store a string of the form soX1s1X2s2...Xmsm, where sm is on top. Each Xi is a grammar symbol and each si is a state.
- The parsing table consists of two parts : *action* and *goto* functions.

Action : The parsing program determines s_m , the state currently on top of stack, and a_i , the current input symbol. It then consults *action*[s_m , a_i] in the action table which can have one of four values :

□ shift s, where s is a state,
 □ reduce by a grammar production A → β,
 □ accept, and

 \Box error.

Goto : The function goto takes a state and grammar symbol as arguments and produces a state.

LR Parsing algorithm:

Input: An input string w and an LR parsing table with functions *action* and *goto* for grammar G.

Output: If *w* is in L(G), a bottom-up-parse for *w*; otherwise, an error indication.

Method: Initially, the parser has so on its stack, where so is the initial state, and w\$ in the input buffer. The parser then executes the following program :

```
set ip to point to the first input symbol of
w$; repeat forever begin
      let s be the state on top of the stack
          and a the symbol pointed to by ip;
      if action[s, a] =shift s' then begin push
         a then s' on top of the stack;
         advance ip to the next input symbol
      end
      else if action[s, a]=reduce A \rightarrow \beta then begin
           pop 2^* |\beta| symbols off the stack;
           let s' be the state now on top of the stack;
           push A then goto[s', A] on top of the
           stack; output the production A \rightarrow \beta
      end
      else if action[s, a]=accept then
           return
      else error()
end
```

2.7 CONSTRUCTING SLR(1) PARSING TABLE:

To perform SLR parsing, take grammar as input and do the following:

- \Box Find LR(0) items.
- \Box Completing the closure.
- \Box Compute *goto*(I,X), where, I is set of items and X is grammar 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. For example, production A \rightarrow XYZ yields the four items :

 $A \rightarrow . XYZ$ $A \rightarrow X \cdot YZ$ $A \rightarrow XY \cdot Z$ $A \rightarrow XYZ \cdot Z$

Closure operation:

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

- \{ Initially, every item in I is added to closure(I).
- \{ If $A \to \alpha$. B β is in closure(I) and $B \to \gamma$ is a production, then add the item $B \to . \gamma$ to I, if it is not already there. We apply this rule until no more new items can be added to closure(I).

Goto operation:

Goto(I, X) is defined to be the closure of the set of all items $[A \rightarrow \alpha X \cdot \beta]$ such

that $[A \rightarrow \alpha . X\beta]$ is in I.

Steps to construct SLR parsing table for grammar G are:

- 1. Augment G and produce G'
- 2. Construct the canonical collection of set of items C for G'
- 3. Construct the parsing action function *action* and *goto* using the following algorithm that requires FOLLOW(A) for each non-terminal of grammar.

Algorithmfor construction of SLR parsing table:

Input : An augmented grammar G'

Output : The SLR parsing table functions action and goto for G'

Method :

- 1. Construct $C = \{I_0, I_1, \dots, I_n\}$, the collection of sets of LR(0) items for G'.
- 2. State *i* is constructed from I*i*. The parsing functions for state *i* are determined as follows:
 - (a) If $[A \rightarrow \alpha \cdot a\beta]$ is in Ii and goto(Ii,*a*) = Ij, then set *action*[*i*,*a*] to "shift j". Here *a* must be terminal.
 - (b) If $[A \rightarrow \alpha^{-}]$ is in I_i, then set *action*[*i*,*a*] to "reduce $A \rightarrow \alpha$ " for all *a* in FOLLOW(A).
 - (c) If $[S' \rightarrow S.]$ is in I_i, then set *action*[i,\$] to "accept".

If any conflicting actions are generated by the above rules, we say grammar is not SLR(1).

3. The *goto* transitions for state *i* are constructed for all non-terminals A using the rule: If $goto(Ii,A) = I_j$, then goto[i,A] = j.

- Σ All entries not defined by rules (2) and (3) are made "error"
- Σ The initial state of the parser is the one constructed from the set of items containing $[S' \rightarrow .S]$.

Example for SLR parsing:

Construct SLR parsing for the following grammar :

 $G: E \to E + T \mid T$ $T \to T * F \mid F$ $F \to (E) \mid id$

The given grammar is :

$G: E \rightarrow E + T$	(1)
$E \rightarrow T$	(2)
$T \to T * F$	(3)
$T \rightarrow F$	(4)
$F \rightarrow (E)$	(5)
$F \rightarrow id$	(6)

Step 1 : Convert given grammar into augmented grammar.

Augmented grammar :

 $E' \rightarrow E$ $E \rightarrow E + T$ $E \rightarrow T$ $T \rightarrow T * F$ $T \rightarrow F$ $F \rightarrow (E)$ $F \rightarrow id$ **Step 2 :** Find LR (0) items.

 $I_{0}: E' \rightarrow \cdot E$ $\delta \rightarrow \cdot E + T$ $iv) \rightarrow \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot id$ $GOTO(I_{0}, E)$ $I_{1}: E' \rightarrow E \cdot$ $E \rightarrow E \cdot + T$

 $\frac{\text{GOTO (I_4, id)}}{\text{I_5}: F \rightarrow \text{id}}$

$\frac{\text{GOTO (I0, T)}}{\text{I2}: E \to T \cdot T \cdot T \to T \cdot F}$	
$\frac{\text{GOTO (I_0, F)}}{\text{I}_3: T \rightarrow F}$	
$\frac{\text{GOTO}(I_0, ())}{\text{I4}: F \rightarrow (\cdot E)}$ $E \rightarrow \cdot E + T$ $E \rightarrow \cdot T$ $T \rightarrow \cdot T * F$ $T \rightarrow \cdot F$ $F \rightarrow \cdot (E)$ $F \rightarrow \cdot \text{id}$	
$\frac{\text{GOTO (I_0, id)}}{\text{I5}: F \rightarrow \text{id}}$	
$\frac{\text{GOTO}(I_1, +)}{\text{I}_6: E \to E + \cdot T}$ $T \to \cdot T * F$ $T \to \cdot F$ $F \to \cdot (E)$ $F \to \cdot \text{id}$	
$\frac{\text{GOTO}(I_2, *)}{I_7: T \to T * \cdot F}$ $F \to \cdot (E)$ $F \to \cdot \text{id}$	
$\frac{\text{GOTO (I4, E)}}{\text{Is}: F \rightarrow (E \cdot) E}$ $\rightarrow E \cdot + T$	
$\frac{\text{GOTO (I4, T)}}{\text{I2}: E \to T \cdot T \cdot F}$	
$\frac{\text{GOTO (I4, F)}}{\text{I3}: T \rightarrow F}$	

 $\underline{\text{GOTO}(I_6,T)}$ I₉: $E \rightarrow E + T$. $T \rightarrow T \cdot * F$ <u>GOTO (I6</u>, F) I₃ : T \rightarrow F. $\underline{\text{GOTO}(I_6},()$ I4: $F \rightarrow (\cdot, E)$ GOTO (I6, id) Is : $F \rightarrow id$. $\underline{\text{GOTO}(I_7,F)}$ In : $T \rightarrow T * F$. <u>GOTO (I</u>₇, () I4: $F \rightarrow (\cdot, E)$ $E \rightarrow . E + T$ $E \rightarrow . T$ $T \to \centerdot T * F$ $T \rightarrow . F$ $F \rightarrow . (E)$ $F \rightarrow .id$ GOTO (I7, id) Is : $F \rightarrow id$. $\underline{\text{GOTO}(I_8,)})$ In: $F \rightarrow (E)$. $\underline{\text{GOTO}(I_8, +)}$ I6: $E \rightarrow E + . T$ $T \rightarrow . T * F$ $T \rightarrow . F$ $F \rightarrow . (E)$ $F \rightarrow . id$ <u>GOTO (I9</u>, *) I7 : T \rightarrow T * . F $F \rightarrow . (E)$ $F \rightarrow . id$

$$\frac{\text{GOTO (I_4, ())}}{\text{I}_4 : F \to (\cdot E)}$$

$$E \to \cdot E + T$$

$$E \to \cdot T$$

$$T \to \cdot T * F$$

$$T \to \cdot F$$

$$F \to \cdot (E)$$

$$F \to \text{id}$$

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

SLR parsing table:

	ACTION			GOTO					
ĺ	id	+	*	()	\$	Е	Т	F
Io	s5			s4			1	2	3
I1		s6				ACC			
I2		r2	s7		r2	r2			
I3		r4	r4		r4	r4			
I4	s5			s4			8	2	3
I5		r6	r6		rб	r6			
I6	s5			s4				9	3
I 7	s5			s4					10
Is		s6			s11				
I 9		r1	s7		r1	r1			
I 10		r3	r3		r3	r3			
I 11		r5	r5		r5	r5			

Blank entries are error entries.

Stack implementation:

Check whether the input **id** + **id** * **id** is valid or not

INFUI	ACTION
id + id * id \$	GOTO (Io, id) = s5; shift
+ id * id \$	GOTO (I ₅ , +) = r6; reduce by $F \rightarrow id$
+ id * id \$	$GOTO (I_0, F) = 3$
	$GOTO(13, +) = 14; \text{ reduce by } 1 \rightarrow F$
+ id * id \$	$\begin{array}{c} \text{GOTO} (\text{Io}, \text{T}) = 2 \\ \text{GOTO} (\text{I}, \text{T}) = 2 \\ \text{GOTO} (\text{I}$
	GOTO (12, +) = r2; reduce by $E \rightarrow 1$
+ id * id \$	GOTO (Io , E) = 1
	GOTO $(I_1, +) = s6$; shift
id * id \$	GOTO (I ₆ , id) = s5 ; shift
* id \$	GOTO (I ₅ , *) = r6; reduce by $F \rightarrow id$
* id \$	GOTO (I ₆ , F) = 3
	GOTO (I ₃ , *) = r4; reduce by $T \rightarrow F$
* id \$	GOTO (I ₆ , T) = 9
	GOTO (I_9 , *) = s7; shift
id \$	GOTO (I ₇ , id) = s5 ; shift
\$	GOTO (I ₅ , \$) = r6; reduce by $F \rightarrow id$
\$	GOTO (I7, F) = 10
	GOTO (I ₁₀ , \$) = r3; reduce by $T \rightarrow T * F$
\$	GOTO (I ₆ , T) = 9
	GOTO (I ₉ , \$) = r1; reduce by $E \rightarrow E + T$
\$	GOTO (Io , E) = 1
	GOTO (I1 , \$) = accept
	id + id * id \$ + id * id \$ + id * id \$ + id * id \$ + id * id \$ id * id \$ * id \$ * id \$ * id \$ * id \$ \$ \$ \$

2.8 TYPE CHECKING

A compiler must check that the source program follows both syntactic and semantic conventions of the source language.

This checking, called *static checking*, detects and reports programming errors.

Some examples of static checks:

1. **Type checks** – A compiler should report an error if an operator is applied to an incompatible operand. Example: If an array variable and function variable are added together.

2. Flow-of-control checks – Statements that cause flow of control to leave a construct must have some place to which to transfer the flow of control. Example: An error occurs when an enclosing statement, such as break, does not exist in switch statement.

Position of type checker



- A *type checker* verifies that the type of a construct matches that expected by its context. For example : arithmetic operator *mod* in Pascal requires integer operands, so a type checker verifies that the operands of *mod* have type integer.
- Type information gathered by a type checker may be needed when code is generated.

2.9 TYPE SYSTEMS

The design of a type checker for a language is based on information about the syntactic constructs in the language, the notion of types, and the rules for assigning types to language constructs.

For example : " if both operands of the arithmetic operators of +,- and * are of type integer, then the result is of type integer "

Type Expressions

- The type of a language construct will be denoted by a "type expression."
- A type expression is either a basic type or is formed by applying an operator called a *type constructor* to other type expressions.
- The sets of basic types and constructors depend on the language to be checked.

The following are the definitions of type expressions:

1. Basic types such as *boolean, char, integer, real* are type expressions.

A special basic type, *type_error*, will signal an error during type checking; *void* denoting "the absence of a value" allows statements to be checked.

- 2. Since type expressions may be named, a type name is a type expression.
- 3. A type constructor applied to type expressions is a type expression. Constructors include:

Arrays : If T is a type expression then *array* (I,T) is a type expression denoting the type of an array with elements of type T and index set I.

Products : If T_1 and T_2 are type expressions, then their Cartesian product $T_1 X T_2$ is a type expression.

Records : The difference between a record and a product is that the fields of a record have names. The *record* type constructor will be applied to a tuple formed from field names and field types.

For example:

```
type row = record

address: integer;

lexeme: array[1..15] of char

end;

var table: array[1...101] of row;
```

declares the type name *row* representing the type expression *record*((*address X integer*) *X* (*lexeme X array*(*1..15,char*))) and the variable *table* to be an array of records of this type.

Pointers : If T is a type expression, then pointer(T) is a type expression denoting the type "pointer to an object of type T". For example, *var p:* \uparrow *row* declares variable p to have type *pointer*(row).

Functions : A function in programming languages maps a *domain type D* to a *range type R*. The type of such function is denoted by the type expression $D \rightarrow R$

4. Type expressions may contain variables whose values are type expressions.

Tree representation for char x char \rightarrow *pointer* (integer)



Type systems

- A type system is a collection of rules for assigning type expressions to the various parts of a program.
- ▶ A type checker implements a type system. It is specified in a syntax-directed manner.
- Different type systems may be used by different compilers or processors of the same language.

Static and Dynamic Checking of Types

- Checkingdone by a compiler is said to be static, while checking done when the target program runs is termed dynamic.
- Any check can be done dynamically, if the target code carries the type of an element along with the value of that element.

Sound type system

A *sound* type system eliminates the need for dynamic checking for type errors because it allows us to determine statically that these errors cannot occur when the target program runs. That is, if a sound type system assigns a type other than *type_error* to a program part, then type errors cannot occur when the target code for the program part is run.

Strongly typed language

A language is strongly typed if its compiler can guarantee that the programs it accepts will execute without type errors.

Error Recovery

- Since type checking has the potential for catching errors in program, it is desirable for type checker to recover from errors, so it can check the rest of the input.
- Error handling has to be designed into the type system right from the start; the type checking rules must be prepared to cope with errors.

2.10 SPECIFICATION OF A SIMPLE TYPE CHECKER

Here, we specify a type checker for a simple language in which the type of each identifier must be declared before the identifier is used. The type checker is a translation scheme that synthesizes the type of each expression from the types of its subexpressions. The type checker can handle arrays, pointers, statements and functions.

<u>A Simple Language</u>

Consider the following grammar: $P \rightarrow D$; E $D \rightarrow D$; D | id : T $T \rightarrow char | integer | array [num] of T | \uparrow T$ $E \rightarrow literal | num | id | E mod E | E [E] | E \uparrow$

Translation scheme:

$\mathbf{P} \rightarrow \mathbf{D}$; E	
$D \rightarrow D$; D	
$D \rightarrow id : T$	{ <i>addtype</i> (id. <i>entry</i> , T. <i>type</i>)}
$T \rightarrow char$	{ T. <i>type</i> : = char }
$T \rightarrow integer$	{ T. <i>type</i> : = integer }
$T \rightarrow \uparrow T1$	{ T. <i>type</i> : = pointer(T1. <i>type</i>) }
$T \rightarrow array [num] of T1$	$\{ T.type := array (1 num.val, T1.type) \}$

In the above language,

- \rightarrow There are two basic types : char and integer ;
- \rightarrow *type_error* is used to signal errors;
- → the prefix operator ↑ builds a pointer type. Example , ↑ integer leads to the type expression pointer (integer).

Type checking of expressions

In the following rules, the attribute *type* forE gives the type expression assigned to the expression generated by E.

- 1. $E \rightarrow literal$ { E.type := char } $E \rightarrow num$ { E.type := integer }Here, constants represented by the tokens literal and num have type *char* and *integer*.
- 2. $E \rightarrow id$ { E.type := lookup (id.entry) } lookup (e) is used to fetch the type saved in the symbol table entry pointed to by e.
- 3. $E \rightarrow E_1 \mod E_2$ { *E.type* := **if** *E1. type* = *integer* **and** *E2. type* = *integer* **then**

integer else type_error }

The expression formed by applying the mod operator to two subexpressions of type integer has type integer; otherwise, its type is *type_error*.

4. $E \rightarrow E_1 [E_2] \{ E.type := if E_2.type = integer and$ $E_1.type = array(s,t) then t else type_error \}$

In an array reference $E_1 [E_2]$, the index expression E_2 must have type integer. The result is the element type *t* obtained from the type array(s,t) of E_1 .

5. $E \rightarrow E_1 \uparrow \{ E.type := \text{if } E_1.type = pointer (t) \text{ then } t \text{ else } type_error \}$

The postfix operator \uparrow yields the object pointed to by its operand. The type of E \uparrow is the type *t* of the object pointed to by the pointer E.

Type checking of statements

Statements do not have values; hence the basic type *void* can be assigned to them. If an error is detected within a statement, then *type_error* is assigned.

Translation scheme for checking the type of statements:

1. Assignment statement:

 $S \rightarrow id := E \{ S.type := if id.type = E.type then void else type_error \}$

2. Conditional statement:

 $S \rightarrow if E$ then $S_1 \{ S.type := if E.type = boolean$ then $S_1.type$

3. While statement:

 $S \rightarrow$ while E do S₁ { S.type : = if E.type = boolean then S₁.type

else type_error }

4. Sequence of statements:

 $S \rightarrow S_1$; $S_2 \{ S.type := if S_1.type = void and S_1.type = void then void else type_error \}$

Type checking of functions

The rule for checking the type of a function application is : $E \rightarrow E_1 (E_2) \{ E.type := if E_2.type = s and$ $E_1.type = s \rightarrow t$ then t else $type_error \}$

2.11 SOURCE LANGUAGE ISSUES

Procedures:

A *procedure definition* is a declaration that associates an identifier with a statement. The identifier is the *procedure name*, and the statement is the *procedure body*. For example, the following is the definition of procedure named *readarray* :

> procedure readarray; var i : integer; begin for i : = 1 to 9 do read(a[i]) end;

When a procedure name appears within an executable statement, the procedure is said to be *called* at that point.

Activation trees:

An *activation tree* is used to depict the way control enters and leaves activations. In an activation tree,

- 1. Each node represents an activation of a procedure.
- 2. The root represents the activation of the main program.
- 3. The node for a is the parent of the node for b if and only if control flows from activation a to b.
- 4. The node for *a* is to the left of the node for *b* if and only if the lifetime of *a* occurs before the lifetime of *b*.

Control stack:

- A *control stack* is used to keep track of live procedure activations. The idea is to push the node for an activation onto the control stack as the activation begins and to pop the node when the activation ends.
- The contents of the control stack are related to paths to the root of the activation tree. When node *n* is at the top of control stack, the stack contains the nodes along the path from *n* to the root.

The Scope of a Declaration:

A declaration is a syntactic construct that associates information with a name. Declarations may be explicit, such as:

var i : integer ;

or they may be implicit. Example, any variable name starting with I is assumed to denote an integer.

The portion of the program to which a declaration applies is called the *scope* of that declaration.

Binding of names:

Even if each name is declared once in a program, the same name may denote different data objects at run time. "Data object" corresponds to a storage location that holds values.

The term *environment* refers to a function that maps a name to a storage location. The term *state* refers to a function that maps a storage location to the value held there.



When an *environment* associates storage location *s* with a name *x*, we say that *x* is *bound* to *s*. This association is referred to as a *binding* of *x*.

2.12 STORAGE ORGANISATION

- The executing target program runs in its own logical address space in which each program value has a location.
- The management and organization of this logical address space is shared between the complier, operating system and target machine. The operating system maps the logical address into physical addresses, which are usually spread throughout memory.

Typical subdivision of run-time memory:



- Run-time storage comes in blocks, where a byte is the smallest unit of addressable memory. Four bytes form a machine word. Multibyte objects are stored in consecutive bytes and given the address of first byte.
- The storage layout for data objects is strongly influenced by the addressing constraints of the target machine.
- A character array of length 10 needs only enough bytes to hold 10 characters, a compiler may allocate 12 bytes to get alignment, leaving 2 bytes unused.
- This unused space due to alignment considerations is referred to as padding.
- The size of some program objects may be known at run time and may be placed in an area called static.
- The dynamic areas used to maximize the utilization of space at run time are stack and heap.

Activation records:

- Procedure calls and returns are usually managed by a run time stack called the *control stack*.
- Each live activation has an activation record on the control stack, with the root of the activation tree at the bottom, the latter activation has its record at the top of the stack.
- The contents of the activation record vary with the language being implemented. The diagram below shows the contents of activation record.



- Temporary values such as those arising from the evaluation of expressions.
- Local data belonging to the procedure whose activation record this is.
- A saved machine status, with information about the state of the machine just before the call to procedures.
- An access link may be needed to locate data needed by the called procedure but found elsewhere.
- A control link pointing to the activation record of the caller.

- Space for the return value of the called functions, if any. Again, not all called procedures return a value, and if one does, we may prefer to place that value in a register for efficiency.
- The actual parameters used by the calling procedure. These are not placed in activation record but rather in registers, when possible, for greater efficiency.

2.13 STORAGE ALLOCATION STRATEGIES

The different storage allocation strategies are :

- 1. Static allocation lays out storage for all data objects at compile time
- 2. Stack allocation manages the run-time storage as a stack.
- 3. **Heap allocation** allocates and deallocates storage as needed at run time from a data area known as heap.

2.13.1 STATIC ALLOCATION

- In static allocation, names are bound to storage as the program is compiled, so there is no need for a run-time support package.
- Since the bindings do not change at run-time, everytime a procedure is activated, its names are bound to the same storage locations.
- Therefore values of local names are *retained* across activations of a procedure. That is, when control returns to a procedure the values of the locals are the same as they were when control left the last time.
- From the type of a name, the compiler decides the amount of storage for the name and decides where the activation records go. At compile time, we can fill in the addresses at which the target code can find the data it operates on.

2.13.2 STACK ALLOCATION OF SPACE

- All compilers for languages that use procedures, functions or methods as units of userdefined actions manage at least part of their run-time memory as a stack.
- Each time a procedure is called , space for its local variables is pushed onto a stack, and when the procedure terminates, that space is popped off the stack.

Calling sequences:

- Procedures called are implemented in what is called as calling sequence, which consists of code that allocates an activation record on the stack and enters information into its fields.
- A return sequence is similar to code to restore the state of machine so the calling procedure can continue its execution after the call.
- The code in calling sequence is often divided between the calling procedure (caller) and the procedure it calls (callee).
- When designing calling sequences and the layout of activation records, the following principles are helpful:
 - Values communicated between caller and callee are generally placed at the beginning of the callee's activation record, so they are as close as possible to the caller's activation record.

- ➢ Fixed length items are generally placed in the middle. Such items typically include the control link, the access link, and the machine status fields.
- Items whose size may not be known early enough are placed at the end of the activation record. The most common example is dynamically sized array, where the value of one of the callee's parameters determines the length of the array.
- We must locate the top-of-stack pointer judiciously. A common approach is to have it point to the end of fixed-length fields in the activation record. Fixed-length data can then be accessed by fixed offsets, known to the intermediate-code generator, relative to the top-of-stack pointer.



Division of tasks between caller and callee

- The calling sequence and its division between caller and callee are as follows.
 - > The caller evaluates the actual parameters.
 - The caller stores a return address and the old value of *top_sp* into the callee's activation record. The caller then increments the *top_sp* to the respective positions.
 - > The callee saves the register values and other status information.
 - > The callee initializes its local data and begins execution.
- A suitable, corresponding return sequence is:

 \succ The callee places the return value next to the parameters.

- ➤ Using the information in the machine-status field, the callee restores top_sp and other registers, and then branches to the return address that the caller placed in the status field.
- Although top_sp has been decremented, the caller knows where the return value is, relative to the current value of top_sp; the caller therefore may use that value.

Variable length data on stack:

- The run-time memory management system must deal frequently with the allocation of space for objects, the sizes of which are not known at the compile time, but which are local to a procedure and thus may be allocated on the stack.
- The reason to prefer placing objects on the stack is that we avoid the expense of garbage collecting their space.
- The same scheme works for objects of any type if they are local to the procedure called and have a size that depends on the parameters of the call.



Access to dynamically allocated arrays

- Procedure p has three local arrays, whose sizes cannot be determined at compile time. The storage for these arrays is not part of the activation record for p.
- Access to the data is through two pointers, *top* and *top-sp*. Here the *top* marks the actual top of stack; it points the position at which the next activation record will begin.
- The second *top-sp* is used to find local, fixed-length fields of the top activation record.
- The code to reposition *top* and *top-sp* can be generated at compile time, in terms of sizes that will become known at run time.

2.13.3 HEAP ALLOCATION

Stack allocation strategy cannot be used if either of the following is possible :

- 1. The values of local names must be retained when an activation ends.
- 2. A called activation outlives the caller.
 - Heap allocation parcels out pieces of contiguous storage, as needed for activation records or other objects.
 - Pieces may be deallocated in any order, so over the time the heap will consist of alternate areas that are free and in use.



- The record for an activation of procedure r is retained when the activation ends.
- Therefore, the record for the new activation q(1, 9) cannot follow that for s physically.
- If the retained activation record for r is deallocated, there will be free space in the heap between the activation records for s and q.

3.1 INTRODUCTION

The front end translates a source program into an intermediate representation from which the back end generates target code.

Benefits of using a machine-independent intermediate form are:

- 5. Retargeting is facilitated. That is, a compiler for a different machine can be created by attaching a back end for the new machine to an existing front end.
- 6. A machine-independent code optimizer can be applied to the intermediate representation.

Position of intermediate code generator



Figure 3.1:Intermediate code generator

3.2 INTERMEDIATE LANGUAGES

Three ways of intermediate representation:

- 1. Syntax tree
- 2. Postfix notation
- 3. Three address code

The semantic rules for generating three-address code from common programming language constructs are similar to those for constructing syntax trees or for generating postfix notation.

3.2.1 Graphical Representations:

Syntax tree:

A syntax tree depicts the natural hierarchical structure of a source program. A **dag** (**Directed Acyclic Graph**) gives the same information but in a more compact way because common subexpressions are identified. A syntax tree and dag for the assignment statement $\mathbf{a} := \mathbf{b}^* - \mathbf{c} + \mathbf{b}^* - \mathbf{c}$ are as follows:



Figure 3.2:Syntax Tree and DAG

Postfix notation:

Postfix notation is a linearized representation of a syntax tree; it is a list of the nodes of the tree in which a node appears immediately after its children. The postfix notation for the syntax tree given above is

a b c uminus * b c uminus * + assign

Syntax-directed definition:

Syntax trees for assignment statements are produced by the syntax-directed definition. Non-terminal S generates an assignment statement. The two binary operators + and * are examples of the full operator set in a typical language. Operator associativities and precedences are the usual ones, even though they have not been put into the grammar. This definition constructs the tree from the input $a := b * - c + b^* - c$.

Table 3.1:Syntax	directed	definition
------------------	----------	------------

PRODUCTION	SEMANTIC RULE
$S \rightarrow id := E$	S.nptr : = mknode('assign',mkleaf(id, id.place), E.nptr)
$E \rightarrow E_1 + E_2$	E.nptr : = mknode('+', E1.nptr, E2.nptr)
$E \rightarrow E_1 * E_2$	E.nptr : = mknode('*', E1.nptr, E2.nptr)
$E \rightarrow - E_1$	E.nptr : = mknode('uminus', E1.nptr)
$E \rightarrow (E_1)$	$E.nptr := E_1.nptr$
$E \rightarrow id$	E.nptr : = mkleaf(id, id.place)

Syntax-directed definition to produce syntax trees for assignment statement

The token **id** has an attribute *place* that points to the symbol-table entry for the identifier. A symbol-table entry can be found from an attribute **id**.*name*, representing the lexeme associated with that occurrence of **id**. If the lexical analyzer holds all lexemes in a single array of characters, then attribute *name* might be the index of the first character of the lexeme.

Two representations of the syntax tree are as follows. In (a) each node is represented as a record with a field for its operator and additional fields for pointers to its children. In (b), nodes are allocated from an array of records and the index or position of the node serves as the pointer to the node. All the nodes in the syntax tree can be visited by following pointers, starting from the root at position 10.



Two representations of the syntax tree

Figure 3.3:Two representations of syntax tree
3.2.2 Three-Address Code:

Three-address code is a sequence of statements of the general form

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

where x, y and z are names, constants, or compiler-generated temporaries; *op* stands for any operator, such as a fixed- or floating-point arithmetic operator, or a logical operator on boolean-valued data. Thus a source language expression like $x + y^*z$ might be translated into asequence

$$t_1 := y * z t_2 :$$

= x + t₁

wheret1 and t2 are compiler-generated temporary names.

Advantages of three-address code:

- The unraveling of complicated arithmetic expressions and of nested flow-ofcontrol statements makes three-address code desirable for target code generation and optimization.
- The use of names for the intermediate values computed by a program allows threeaddress code to be easily rearranged – unlike postfix notation.

Three-address code is a linearized representation of a syntax tree or a dag in which explicit names correspond to the interior nodes of the graph. The syntax tree and dag are represented by the three-address code sequences. Variable names can appear directly in three-address statements.

Three-address code corresponding to the syntax tree and dag given above

$t_1 := -c$	$t_1 := -c$
$t_2 := b * t_1$	$t_2 := b * t_1$
$t_3 := -c$	$t_5 := t_2 + t_2$
$t_4 := b * t_3$	a : = t5
$t_5 := t_2 + t_4$	
$a := t_5$	

(a) Code for the syntax tree (b) Code for the dag

The reason for the term "three-address code" is that each statement usually contains three addresses, two for the operands and one for the result.

Types of Three-Address Statements:

The common three-address statements are:

- 5. Assignment statements of the form $\mathbf{x} := \mathbf{y} op \mathbf{z}$, where op is a binary arithmetic or logical operation.
- 6. Assignment instructions of the form $\mathbf{x} := op \mathbf{y}$, where op is a unary operation. Essential unary operations include unary minus, logical negation, shift operators, and conversion operators that, for example, convert a fixed-point number to a floating-point number.
- 7. Copy statements of the form $\mathbf{x} := \mathbf{y}$ where the value of y is assigned to x.
- 8. The unconditional jump goto L. The three-address statement with label L is the next to be executed.
- 9. Conditional jumps such as if x relop y goto L. This instruction applies a relational operator (
 <, =, >=, etc.) to x and y, and executes the statement with label L next if x stands in relation relop to y. If not, the three-address statement following if x relop y goto L is executed next, as in the usual sequence.
- 4. *param x* and *call p, n* for procedure calls and *return y*, where y representing a returned value is optional. For example,

param x1 param x2 \dots param xn call p,n generated as part of a call of the procedure p(x1, x2,, xn).

- 5. Indexed assignments of the form x := y[i] and x[i] := y.
- 6. Address and pointer assignments of the form x := &y, x := *y, and *x := y.

Syntax-Directed Translation into Three-Address Code:

When three-address code is generated, temporary names are made up for the interior nodes of a syntax tree. For example, id := E consists of code to evaluate E into some temporary t, followed by the assignment id.place := t.

Given input a := b * - c + b * - c, the three-address code is as shown above. The synthesized attribute *S.code* represents the three-address code for the assignment *S*. The nonterminal *E* has two attributes :

- 5. *E.place*, the name that will hold the value of E, and
- 6. *E.code*, the sequence of three-address statements evaluating *E*.

PRODUCTION	SEMANTIC RULES
$S \rightarrow id := E$	S.code : = E.code gen(id.place ':=' E.place)
$E \rightarrow E_1 + E_2$	E.place := newtemp; E.code := E1.code E2.code gen(E.place ':=' E1.place '+' E2.place)
$E \rightarrow E_1 * E_2$	E.place := newtemp; E.code := E1.code E2.code gen(E.place ':=' E1.place '*' E2.place)
$E \rightarrow - E_1$	E.place := newtemp; E.code := E1.code gen(E.place ':=' 'uminus' E1.place)
$E \rightarrow (E_1)$	E.place : = E1.place; E.code : = E1.code
E ⇒id	E.place : = id.place; E.code : = ' '

Table 3.2: Syntax-directed definition to produce three-address code for assignments



Semantic rules generating code for a while statement

PRODUCTION

 $S \rightarrow$ while *E* do *S*₁

S.begin := newlabel; S.after := newlabel; S.code := gen(S.begin ':') || E.code || gen ('if' E.place '=' '0' 'goto' S.after)|| S1.code || gen ('goto' S.begin) || gen (S.after ':')

SEMANTIC RULES

Figure 3.4:semantic rule generating code for while statement

- The function *newtemp* returns a sequence of distinct names t₁,t₂,.... in response to successive calls.
- (*) Notation gen(x := y + z) is used to represent three-address statement x := y + z. Expressions appearing instead of variables like x, y and z are evaluated when passed to *gen*, and quoted operators or operand, like '+' are taken literally.
- The Flow-of-control statements can be added to the language of assignments. The code for $S \rightarrow$ while E do S_1 is generated using new attributes *S.begin* and *S.after* to mark the first statement in the code for E and the statement following the code for S, respectively.
- ⑦ The function *newlabel* returns a new label every time it is called.
- O We assume that a non-zero expression represents true; that is when the value of E becomes zero, control leaves the while statement.

Implementation of Three-Address Statements:

A three-address statement is an abstract form of intermediate code. In a compiler, these statements can be implemented as records with fields for the operator and the operands. Three such representations are:Quadruples

- (b) Triples
- ⑦ Indirect triples

Quadruples:

- ② A quadruple is a record structure with four fields, which are, *op*, *arg1*, *arg2* and *result*.
- The *op* field contains an internal code for the operator. The three-address statement $\mathbf{x} := \mathbf{y}$ op \mathbf{z} is represented by placing y in *arg1*, z in *arg2* and x in *result*.
- The contents of fields arg1, arg2 and result are normally pointers to the symbol-table entries for the names represented by these fields. If so, temporary names must be entered into the symbol table as they are created.

Triples:

- To avoid entering temporary names into the symbol table, we might refer to a temporary value by the position of the statement that computes it.
- If we do so, three-address statements can be represented by records with only three fields: *op*, *arg1* and *arg2*.
- The fields *arg1* and *arg2*, for the arguments of *op*, are either pointers to the symbol table or pointers into the triple structure (for temporary values).
- ② Since three fields are used, this intermediate code format is known as *triples*.

	on	aro l	aro?	result		op	arg1	
	op	ur 81	u182	resuit	$\langle 0 \rangle$			
(0)	uminus	с		t1	(0)	uminus	С	
					(1)	*	b	(0)
(1)	*	b	t 1	t2	(-)		-	
					(2)	uminus	с	
(2)	uminus	С		t3				
(3)	*	h	to.	ŧ.	(3)	*	b	(2)
(\mathbf{J})		U	13	L4	(\mathbf{A})		(1)	(2)
(4)	+	t2	t4	t5	(4)	+	(1)	(3)
					(5)	assign	а	(4)
(5)	:=	t3		а	(0)	455151	u	

(a) Quadruples

(b) Triples

Figure 3.4: Quadruple and triple representation of three-address statements given above

A ternary operation like x[i] := y requires two entries in the triple structure as shown as below while x := y[i] is naturally represented as two operations.

	ор	arg1	arg2	-		op	arg1	arg2
(0)	[]=	Х	i		(0)	=[]	у	i
(1)	assign	(0)	У		(1)	assign	Х	(0)

(a) **x**[i] : = **y**

(b) x := y[i]



Indirect Triples:

- ② Another implementation of three-address code is that of listing pointers to triples, rather than listing the triples themselves. This implementation is called indirect triples.
- ③ For example, let us use an array statement to list pointers to triples in the desired order. Then the triples shown above might be represented as follows:

	statement		ор	argl	arg2
(0) (1)	(14)	(14) (15)	uminus *	c b	(14)
(1) (2)	(16)	(15)	uminus	c	(11)
(3) (4)	(17) (18)	(17) (18)	*	b (15)	(16) (17)
(5)	(19)	(19)	assign	a	(18)

Figure 3.6:Indirect triples representation of three-address statements

3.3 DECLARATIONS

As the sequence of declarations in a procedure or block is examined, we can lay out storage for names local to the procedure. For each local name, we create a symbol-table entry with information like the type and the relative address of the storage for the name. The relative address consists of an offset from the base of the static data area or the field for local data in an activation record.

3.3.1 Declarations in a Procedure:

The syntax of languages such as C, Pascal and Fortran, allows all the declarations in a single procedure to be processed as a group. In this case, a global variable, say *offset*, can keep track of the next available relative address.

In the translation scheme shown below:

- ② Nonterminal P generates a sequence of declarations of the form **id** : *T*.
- ③ Before the first declaration is considered, *offset* is set to 0. As each new name is seen, that name is entered in the symbol table with offset equal to the current value of *offset*, and *offset* is incremented by the width of the data object denoted by that name.
- The procedure *enter(name, type, offset*) creates a symbol-table entry for *name*, gives its type *type* and relative address *offset* in its data area.
- Attribute *type* represents a type expression constructed from the basic types *integer* and *real* by applying the type constructors *pointer* and *array*. If type expressions are represented by graphs, then attribute *type* might be a pointer to the node representing a type expression.
- The width of an array is obtained by multiplying the width of each element by the number of elements in the array. The width of each pointer is assumed to be 4.

Computing the types and relative addresses of declared names

$P \rightarrow D$	{ offset : = 0 }
$D \rightarrow D; D$	
$D \rightarrow id: T$	{ enter(id.name, T.type, offset); offset : = offset + T.width }
$T \rightarrow integer$	{ T.type : = integer; T.width : = 4 }
$T \rightarrow real$	{ T.type : = real; T.width : = 8 }
$T \rightarrow array [num] of T_1$	{ T.type : = array(num.val, T1.type); T.width : = num.val X T1.width }
$T \rightarrow \uparrow T_1$	{ T.type : = pointer (T1.type); T.width : = 4 }

3.3.2 Keeping Track of Scope Information:

When a nested procedure is seen, processing of declarations in the enclosing procedure is temporarily suspended. This approach will be illustrated by adding semantic rules to the following language:

 $P \rightarrow D$ $D \rightarrow D; D / id : T / proc id; D; S$

One possible implementation of a symbol table is a linked list of entries for names.

A new symbol table is created when a procedure declaration $D \rightarrow proc id D_{1}S$ is seen, and entries for the declarations in D₁ are created in the new table. The new table points back to the symbol table of the enclosing procedure; the name represented by id itself is local to the enclosing procedure. The only change from the treatment of variable declarations is that the procedure *enter* is told which symbol table to make an entry in.

For example, consider the symbol tables for procedures *readarray, exchange*, and *quicksort* pointing back to that for the containing procedure *sort*, consisting of the entire program. Since *partition* is declared within *quicksort*, its table points to that of *quicksort*.



Figure 3.7: Symbol tables for nested procedures

The semantic rules are defined in terms of the following operations:

- 1. *mktable(previous)* creates a new symbol table and returns a pointer to the new table. The argument *previous* points to a previously created symbol table, presumably that for the enclosing procedure.
- 2. *enter(table, name, type, offset)* creates a new entry for name *name* in the symbol table pointed to by *table*. Again, *enter* places type *type* and relative address *offset* in fields within the entry.
- 3. *addwidth(table, width)* records the cumulative width of all the entries in table in the header associated with this symbol table.
- 4. *enterproc(table, name, newtable)* creates a new entry for procedure *name* in the symbol table pointed to by *table*. The argument *newtable* points to the symbol table for this procedure *name*.

Syntax directed translation scheme for nested procedures

$P \rightarrow M D$	<pre>{ addwidth (top(tblptr) , top (offset)); pop (tblptr); pop (offset) }</pre>
$M \rightarrow \varepsilon$	{ t : = mktable (nil); push (t,tblptr); push (0,offset) }
$D \rightarrow D_1; D_2$	
$D \rightarrow proc id; ND_1; S$	<pre>{ t : = top (tblptr); addwidth (t, top (offset)); pop (tblptr); pop (offset); enterproc (top (tblptr), id.name, t) }</pre>
$D \rightarrow id: T$	{ enter (top (tblptr), id.name, T.type, top (offset)); top (offset) := top (offset) + T.width }
$N \rightarrow \varepsilon$	{ t := mktable (top (tblptr)); push (t, tblptr); push (0,offset) }

- > The stack *tblptr* is used to contain pointers to the tables for **sort**, **quicksort**, and **partition** when the declarations in **partition** are considered.
- The top element of stack *offset* is the next available relative address for a local of the current procedure.
- > All semantic actions in the subtrees for B and C in

 $A \rightarrow BC \{actionA\}$

are done before *actionA* at the end of the production occurs. Hence, the action associated with the marker M is the first to be done.

- > The action for nonterminal M initializes stack *tblptr* with a symbol table for the outermost scope, created by operation mktable(nil). The action also pushes relative address 0 onto stack offset.
- Similarly, the nonterminal N uses the operation *mktable(top(tblptr))* to create a new symbol table. The argument *top(tblptr)* gives the enclosing scope for the new table.
- For each variable declaration id: T, an entry is created for id in the current symbol table.
 The top of stack offset is incremented by T.width.
- > When the action on the right side of D → proc id; ND1; S occurs, the width of all declarations generated by D1 is on the top of stack offset; it is recorded using addwidth. Stacks tblptr and offset are then popped.

At this point, the name of the enclosed procedure is entered into the symbol table of its enclosing procedure.

3.4 ASSIGNMENT STATEMENTS

Suppose that the context in which an assignment appears is given by the following grammar.

 $P \rightarrow M D$ $M \rightarrow \varepsilon$ $D \rightarrow D; D | id : T | proc id; N D; S$ $N \rightarrow \varepsilon$

Nonterminal P becomes the new start symbol when these productions are added to those in the translation scheme shown below.

Translation scheme to produce three-address code for assignments

$S \rightarrow id := E \{ p := lookup (id.name);$
if p ≠nil then
emit(p ' : =' E.place)
else error }
$E \rightarrow E_1 + E_2 \{ \text{ E.place : = newtemp;} \\ \text{emit(E.place ': =' E1.place ' + ' E2.place) } \}$
$E \rightarrow E_1 * E_2 $ { E.place : = newtemp; emit(E.place ': =' E1.place ' * ' E2.place) }
$E \rightarrow - E_1$ { E.place : = newtemp; emit (E.place ': =' 'uminus' E1.place) }
$E \rightarrow (E_1) \{ E.place := E_1.place \}$

 $E \rightarrow id \{ p := lookup (id.name); \}$

if p ≠nil **then** E.place : = p **else** error }

Reusing Temporary Names

- The temporaries used to hold intermediate values in expression calculations tend to clutter up the symbol table, and space has to be allocated to hold their values.
- > Temporaries can be reused by changing *newtemp*. The code generated by the rules for E
 - \rightarrow E₁ + E₂ has the general form:

evaluate E_1 into t_1 evaluate E_2 into t_2 $t := t_1 + t_2$

- > The lifetimes of these temporaries are nested like matching pairs of balanced parentheses.
- Keep a count c , initialized to zero. Whenever a temporary name is used as an operand, decrement c by 1. Whenever a new temporary name is generated, use \$c and increase c by 1.
- For example, consider the assignment x := a * b + c * d e * f

statement	value of c
	0
\$0 := a * b	1
1 := c * d	2
0 := 0 + 1	1
\$1 := e * f	2
\$0 := \$0 - \$1	1
x := \$0	0

Table 3.3:Three-address code with stack temporaries

Addressing Array Elements:

Elements of an array can be accessed quickly if the elements are stored in a block of consecutive locations. If the width of each array element is *w*, then the *i*th element of array A begins in location

$$base + (i - low) \ge w$$

where *low* is the lower bound on the subscript and *base* is the relative address of the storage allocated for the array. That is, *base* is the relative address of A*[low]*.

The expression can be partially evaluated at compile time if it is rewritten as

$$i \mathbf{x} w + (base - low \mathbf{x} w)$$

The subexpression $c = base - low \ge w$ can be evaluated when the declaration of the array is seen. We assume that c is saved in the symbol table entry for A, so the relative address of A[i] is obtained by simply adding *i* x *w* to *c*.

Address calculation of multi-dimensional arrays:

A two-dimensional array is stored in of the two forms :

- Row-major (row-by-row)
- Column-major (column-by-column)



Layouts for a 2 x 3 array

(b) COLUMN-MAJOR

Figure 3.8:Address calculation of multi-dimensional arrays

In the case of row-major form, the relative address of A[i1,i2] can be calculated by the formula

$$base + ((i_1 - low_1) \ge n_2 + i_2 - low_2) \ge w$$

(a) **ROW-MAJOR**

where, low_1 and low_2 are the lower bounds on the values of i_1 and i_2 and n_2 is the number of values that i_2 can take. That is, if *high*₂ is the upper bound on the value of i_2 , then $n_2 = high_2 - low_2 + 1$.

Assuming that i1 and i2 are the only values that are known at compile time, we can rewrite the above expression as

$$((i_1 \ge n_2) + i_2) \ge w + (base - ((low_1 \ge n_2) + low_2) \ge w)$$

Generalized formula:

The expression generalizes to the following expression for the relative address of A[*i1,i2,...,ik*]

 $((\ldots ((i_{1n2} + i_{2})n_{3} + i_{3})\ldots)n_{k} + i_{k}) \ge w + base - ((\ldots ((low_{1n2} + low_{2})n_{3} + low_{3})\ldots)n_{k} + low_{k}) \ge w$

for all j, $n_j = high_j - low_j + 1$

The Translation Scheme for Addressing Array Elements :

Semantic actions will be added to the grammar :

(1) $S \rightarrow L := E$ (2) $E \rightarrow E + E$ (3) $E \rightarrow (E)$ (4) $E \rightarrow L$ (5) $L \rightarrow Elist$] (6) $L \rightarrow id$ (7) $Elist \rightarrow Elist, E$ (8) $Elist \rightarrow id [E]$

We generate a normal assignment if L is a simple name, and an indexed assignment into the location denoted by L otherwise :

(1)
$$S \rightarrow L := E$$
 { if L.offset = null then / * L is a simple id */
emit (L.place ': =' E.place);
else
emit (L.place '[' L.offset ']' ': ='E.place) }
(2) $E \rightarrow E_1 + E_2$ { E.place : = newtemp;
emit (E.place ': =' E1.place ' +' E2.place) }
(3) $E \rightarrow (E_1)$ { E.place : = E1.place }

When an array reference L is reduced to E, we want the *r*-value of L. Therefore we use indexing to obtain the contents of the location *L.place* [*L.offset*]:

- (4) E → L
 { if L.offset = null then /* L is a simple id*/
 E.place : = L.place
 else begin
 E.place : = newtemp;
 emit (E.place ': =' L.place ' [' L.offset ']')
 end }
- (5) L → Elist] { L.place := newtemp; L.offset := newtemp; emit (L.place ': =' c(Elist.array)); emit (L.offset ': =' Elist.place '*' width (Elist.array)) }

(6)
$$L \rightarrow id$$
 { L.place := id.place;
L.offset := null }
(7) Elist \rightarrow Elist1, E { $t := newtemp;$
 $m := Elist1.ndim + 1;$
 $emit (t ': =' Elist 1.place '*' limit$
(Elist1.array,m)); $emit (t ': =' t '+' E.place);$
Elist.array := Elist1.array;
Elist.place := t;
Elist.ndim := m }

(8) *Elist* \rightarrow **id** [*E* { *Elist.array* : = **id**.*place*;

Elist.place : = E.place; Elist.ndim : = 1 }

Type conversion within Assignments :

Consider the grammar for assignment statements as above, but suppose there are two types – real and integer, with integers converted to reals when necessary. We have another attribute *E.type*, whose value is either *real* or *integer*. The semantic rule for *E.type* associated with the production $E \rightarrow E + E$ is :

$$E \rightarrow E + E$$
 { $E.type :=$
if $E_{1.type} = integer$ and
 $E_{2.type} = integer$ then integer
else real }

The entire semantic rule for $E \rightarrow E + E$ and most of the other productions must be modified to generate, when necessary, three-address statements of the form x := inttoreal y, whose effect is to convert integer y to areal of equal value, called x.

```
Semantic action for E \rightarrow E_1 + E_2
```

```
E.place := newtemp;

if E1.type = integer and E2.type = integer then

begin emit( E.place ': =' E1.place 'int +'

E2.place); E.type : = integer

end

else if E1.type = real and E2.type = real then begin

emit( E.place ': =' E1.place 'real +'

E2.place); E.type : = real

end
```

```
else if E1 .type = integer and E2.type = real then
    begin u := newtemp;
    emit( u ': =' 'inttoreal' E1.place); emit(
      E.place ': =' u ' real +' E2.place);
      E.type := real
```

end

```
else if E1.type = real and E2.type =integer then
    begin u : = newtemp;
    emit( u ': = ' 'inttoreal' E2.place); emit(
        E.place ': = ' E1.place ' real + ' u);
        E.type : = real
```

end else

E.type : = *type_error*;

For example, for the input x := y + i * jassuming *x* and *y* have type *real*, and i and j have type *integer*, the output would look like

t1 := i int* j t3 : = inttoreal t1 t2 : = y real+ t3 x := t2

3.5 BOOLEAN EXPRESSIONS

Boolean expressions have two primary purposes. They are used to compute logical values, but more often they are used as conditional expressions in statements that alter the flow of control, such as if-then-else, or while-do statements.

Boolean expressions are composed of the boolean operators (**and**, **or**, and **not**) applied to elements that are boolean variables or relational expressions. Relational expressions are of the form E_1 relop E_2 , where E_1 and E_2 are arithmetic expressions.

Here we consider boolean expressions generated by the following grammar :

$E \rightarrow E$ or E | E and E | not E | (E) | id relop id | true | false

Methods of Translating Boolean Expressions:

There are two principal methods of representing the value of a boolean expression. They are :

- To encode true and false *numerically* and to evaluate a boolean expression analogously to an arithmetic expression. Often, 1 is used to denote true and 0 to denote false.
- To implement boolean expressions by *flow of control*, that is, representing the value of a boolean expression by a position reached in a program. This method is particularly convenient in implementing the boolean expressions in flow-of-control statements, such as the if-then and while-do statements.

3.5.1 Numerical Representation

Here, 1 denotes true and 0 denotes false. Expressions will be evaluated completely from left to right, in a manner similar to arithmetic expressions. For example :

The translation for

 a or b and not c is
 the three-address sequence
 t1 := not c t2
 := b and t1 t3
 := a or t2

 \blacktriangleright A relational expression such as a < b is equivalent to the conditional

statement if a < b then 1 else 0

which can be translated into the three-address code sequence (again, we arbitrarily start statement numbers at 100) :

100 : if a < b goto 103 101 : t := 0 102 : goto 104 103 : t := 1 104 :

Translation scheme using a numerical representation for booleans

$E \rightarrow E_1 \text{ or } E_2$	{ E.place : = newtemp;
	emit(E.place ': =' E1.place 'or 'E2.place)}
$E \rightarrow E_1$ and E_2	{ E.place : = newtemp;
	emit(E.place ': =' E1.place 'and 'E2.place)}
$E \rightarrow \mathbf{not} \ E_l$	{ E.place : = newtemp;
	emit(E.place ': =' ' not ' E1.place)}
$E \rightarrow (E_1)$	{ E.place : = E1.place }
$E \rightarrow id_1 relop id_2$	{ E.place : = newtemp;
	<pre>emit('if' id1.place relop.op id2.place 'goto' nextstat + 3);</pre>
	emit(E.place ': =' ' 0 ');
	emit('goto' nextstat +2);
	emit(E.place ': =' ' 1 ') }
$E \rightarrow$ true	{ E.place : = newtemp;
	emit(E.place ': =' ' 1 ') }
E →false	{ E.place : = newtemp;
	emit(E.place ': =' '0') }

Short-Circuit Code:

We can also translate a boolean expression into three-address code without generating code for any of the boolean operators and without having the code necessarily evaluate the entire expression. This style of evaluation is sometimes called "**short-circuit**" or "**jumping**" code. It is possible to evaluate boolean expressions without generating code for the boolean operators **and**, **or**, and **not** if we represent the value of an expression by a position in the code sequence.

Translation of a < b or c < d and e < f

100	: if	a < b goto	103	$107: t_2:=1$
101	$: t_1 := 0$	0		108 : if e < f goto 111
102	: goto 1	04		$109:t_3:=0$
103	: t 1	: = 1		110 : goto 112
104	: if	c < d goto	107	$111: t_3: = 1$
105	: t2	:=0		$112: t_4: = t_2 and t_3$
106	: goto 1	08		113 : t5 : = t1 or t4

Flows-of-Control Statements

We now consider the translation of boolean expressions into three-address code in the context of if-then, if-then-else, and while-do statements such as those generated by the following grammar:

- $S \rightarrow if E then S_1$
 - | if E then S1 else S2
 - | while E do S1

In each of these productions, E is the Boolean expression to be translated. In the translation, we assume that a three-address statement can be symbolically labeled, and that the function *newlabel* returns a new symbolic label each time it is called.

- E.true is the label to which control flows if E is true, and E.false is the label to which control flows if E is false.
- The semantic rules for translating a flow-of-control statement S allow control to flow from the translation S.code to the three-address instruction immediately following S.code.
- S.next is a label that is attached to the first three-address instruction to be executed after the code for S.

Code for if-then , if-then-else, and while-do statements





(c) while-do Figure 3.9: code for if-then,if-then-else,while-do statements

PRODUCTION SEMANTIC RULES $S \rightarrow if E then S_1$ *E.true* : = *newlabel*; *E.false* : = *S.next*; $S_{1.next}$:= $S_{.next}$; $S.code := E.code || gen(E.true ':') || S_1.code$ $S \rightarrow if E then S_1 else S_2$ *E.true* : = *newlabel*; *E.false* : = *newlabel*; $S_{1.next}$:= $S_{.next}$; $S_{2.next}$:= $S_{.next}$; *S.code* : = *E.code* || *gen*(*E.true* ':') || *S1.code* || gen('goto' S.next) || gen(E.false ':') || S2.code $S \rightarrow$ while *E* do *S*₁ S.begin : = newlabel; *E.true* : = newlabel; *E.false* : = *S.next*; $S_{1.next}$:= $S_{.begin}$; S.code := gen(S.begin ':') || E.code ||gen(E.true ':') || S1.code // gen('goto' S.begin)

Table 3.4:Syntax-directed definition for flow-of-control statements

3.5.2 Control-Flow Translation of Boolean Expressions:

PRODUCTION	SEMANTIC RULES
$E \rightarrow E_1 \text{ or } E_2$	E1.true : = E.true; E1.false : = newlabel; E2.true : = E.true; E2.false : = E.false; E.code : = E1.code // gen(E1.false ':') E2.code
$E \rightarrow E_1$ and E_2	E.true : = newlabel; E1.false : = E.false; E2.true : = E.true; E2.false : = E.false; E.code : = E1.code // gen(E1.true ':') E2.code
$E \rightarrow \mathbf{not} \ E_1$	$E_{1.true} := E.false;$ $E_{1.false} := E.true;$ $E.code := E_{1.code}$
$E \rightarrow (E1)$	$E_{1.true} := E_{.true};$ $E.code := E_{1.code}$
$E \rightarrow id_1 relop id_2$	E.code : = gen('if' id1.place relop.op id2.place 'goto' E.true) // gen('goto' E.false)
$E \rightarrow true$	E.code := gen(`goto' E.true)
$E \rightarrow \mathbf{false}$	E.code : = gen('goto' E.false)

Table 3.5:Syntax-directed definition to produce three-address code for booleans

3.6 CASE STATEMENTS

The "switch" or "case" statement is available in a variety of languages. The switch-statement syntax is as shown below :

```
Switch-statement syntax

switch expression

begin

case value : statement case

value : statement

....

case value : statement

default : statement

end
```

There is a selector expression, which is to be evaluated, followed by n constant values that the expression might take, including a default "value" which always matches the expression if no other value does. The intended translation of a switch is code to:

- 4. Evaluate the expression.
- 5. Find which value in the list of cases is the same as the value of the expression.
- 6. Execute the statement associated with the value found.

Step (2) can be implemented in one of several ways :

- By a sequence of conditional **goto** statements, if the number of cases is small.
- By creating a table of pairs, with each pair consisting of a value and a label for the code of the corresponding statement. Compiler generates a loop to compare the value of the expression with each value in the table. If no match is found, the default (last) entry is sure to match.
- If the number of cases s large, it is efficient to construct a hash table.
- There is a common special case in which an efficient implementation of the n-way branch exists. If the values all lie in some small range, say imin to imax, and the number of different values is a reasonable fraction of imax imin, then we can construct an array of labels, with the label of the statement for value j in the entry of the table with offset j imin and the label for the default in entries not filled otherwise. To perform switch,

evaluate the expression to obtain the value of j, check the value is within range and transfer to the table entry at offset j-imin.

Syntax-Directed Translation of Case Statements:

Consider the following switch statement:

```
switch E begin

case V_1 : S_1 case

V_2 : S_2

...

case V_{n-1} : S_{n-1}

default : S_n
```

```
end
```

This case statement is translated into intermediate code that has the following form :

Translation of a case statement

	code to evaluate E into t
	goto test
L1 :	code for <i>S</i> ₁
	goto next
L2:	code for S ₂
	goto next
	•••
L_{n-1} :	code for S_{n-1}

```
goto next

Ln:

code for S_n

goto next

test:

if t = V_1 goto L1

if t = V_2 goto L2

...

if t = V_{n-1} goto Ln-1

goto Ln

next:
```

To translate into above form :

 $\blacksquare \blacksquare \blacksquare$ When keyword **switch** is seen, two new labels **test** and **next**, and a new temporary **t** are generated.

As expression *E* is parsed, the code to evaluate *E* into **t** is generated. After processing *E*, the jump **goto test** is generated.

- Solution As each **case** keyword occurs, a new label L_i is created and entered into the symbol table. A pointer to this symbol-table entry and the value V_i of case constant are placed on a stack (used only to store cases).
- Each statement case V_i : S_i is processed by emitting the newly created label L_i , followed by the code for S_i , followed by the jump goto next.
- Then when the keyword **end** terminating the body of the switch is found, the code can be generated for the n-way branch. Reading the pointer-value pairs on the case stack from the bottom to the top, we can generate a sequence of three-address statements of the form

```
case V1 L1 case
V2 L2
...
case Vn-1 Ln-1
case t Ln label
next
```

where t is the name holding the value of the selector expression E, and L_n is the label for the default statement.

3.7 BACKPATCHING

The easiest way to implement the syntax-directed definitions for boolean expressions is to use two passes. First, construct a syntax tree for the input, and then walk the tree in depth-first order, computing the translations. The main problem with generating code for boolean expressions and flowof-control statements in a single pass is that during one single pass we may not know the labels that control must go to at the time the jump statements are generated. Hence, a series of branching statements with the targets of the jumps left unspecified is generated. Each statement will be put on a list of goto statements whose labels will be filled in when the proper label can be determined. We call this subsequent filling in of labels *backpatching*. To manipulate lists of labels, we use three functions :

- 5. *makelist*(*i*) creates a new list containing only *i*, an index into the array of quadruples; *makelist* returns a pointer to the list it has made.
- 6. $merge(p_1,p_2)$ concatenates the lists pointed to by p_1 and p_2 , and returns a pointer to the concatenated list.
- 7. backpatch(p,i) inserts i as the target label for each of the statements on the list pointed to by p.

Boolean Expressions:

We now construct a translation scheme suitable for producing quadruples for boolean expressions during bottom-up parsing. The grammar we use is the following:

- 5. $E \rightarrow E_1$ or $M E_2$
- 6. | *E*¹ and *M E*²
- 7. | not E_1
- 8. $|(E_l)|$
- 9. | **id**1 **relop id**2
- 10. | **true**
- 11. | **false**
- 12. $M \rightarrow \varepsilon$

Synthesized attributes *truelist* and *falselist* of nonterminal *E* are used to generate jumping code for boolean expressions. Incomplete jumps with unfilled labels are placed on lists pointed to by *E.truelist* and *E.falselist*.

Consider production $E \rightarrow E_1$ and $M E_2$. If E_1 is false, then E is also false, so the statements on E_1 .falselist become part of *E*.falselist. If E_1 is true, then we must next test E_2 , so the target for the statements E_1 .truelist must be the beginning of the code generated for E_2 . This target is obtained using marker nonterminal M.

Attribute *M.quad* records the number of the first statement of *E2.code*. With the production $M \rightarrow \varepsilon$ we associate the semantic action

 $\{ M.quad := nextquad \}$

The variable *nextquad* holds the index of the next quadruple to follow. This value will be backpatched onto the *E1.truelist* when we have seen the remainder of the production $E \rightarrow E1$ and *ME2*. The translation scheme is as follows:

(1) $E \rightarrow E_1$ or $M E_2$	{ backpatch (E1.falselist, M.quad);		
	<i>E.truelist</i> : = merge(<i>E1.truelist</i> , <i>E2.truelist</i>);		
	$E.falselist := E_2.falselist $		
(2) $E \rightarrow E_1$ and $M E_2$	{ backpatch (E1.truelist, M.quad);		

	E.truelist : = E2.truelist; E.falselist : = merge(E1.falselist, E2.falselist)}
(3) $E \rightarrow \operatorname{not} E_1$	{ E.truelist : = E1.falselist; E.falselist : = E1.truelist; }
$(4) E \rightarrow (E_I)$	{ E.truelist := E1.truelist; E.falselist := E1.falselist; }
(5) $E \rightarrow id_1 relop id_2$	<pre>{ E.truelist := makelist (nextquad); E.falselist := makelist(nextquad + 1); emit('if' id1.place relop.op id2.place 'goto_') emit('goto_') }</pre>
(6) $E \rightarrow$ true	<pre>{ E.truelist := makelist(nextquad); emit('goto_') }</pre>
(7) $E \rightarrow \mathbf{false}$	<pre>{ E.falselist : = makelist(nextquad); emit('goto_') }</pre>
(8) $M \rightarrow \varepsilon$	{ <i>M.quad</i> : = <i>nextquad</i> }

Flow-of-Control Statements:

A translation scheme is developed for statements generated by the following grammar :

5.	$S \rightarrow \text{if } E \text{ then } S$
6.	if E then S else S
7.	while E do S
8.	begin L end
9.	A
10.	$L \rightarrow L$; S
11.	$\mid S$

Here *S* denotes a statement, *L* a statement list, *A* an assignment statement, and E a boolean expression. We make the tacit assumption that the code that follows a given statement in execution also follows it physically in the quadruple array. Else, an explicit jump must be provided.

Scheme to implement the Translation:

The nonterminal E has two attributes E.truelist and E.falselist. L and S also need a list of

unfilled quadruples that must eventually be completed by backpatching. These lists are pointed to by the attributes *L..nextlist* and *S.nextlist*. *S.nextlist* is a pointer to a list of all conditional and unconditional jumps to the quadruple following the statement S in execution order, and *L.nextlist* is defined similarly.

The semantic rules for the revised grammar are as follows:

F

S → if E then M₁ S₁ N else M₂ S₂ { backpatch (E.truelist, M₁.quad); backpatch (E.falselist, M₂.quad); S.nextlist : = merge (S₁.nextlist, merge (N.nextlist, S₂.nextlist)) }

We backpatch the jumps when E is true to the quadruple $M_1.quad$, which is the beginning of the code for S₁. Similarly, we backpatch jumps when E is false to go to the beginning of the code for S₂. The list *S.nextlist* includes all jumps out of S₁ and S₂, as well as the jump generated by N.

(2)	$N \rightarrow \varepsilon$	<pre>{ N.nextlist : = makelist(nextquad); emit('goto _') }</pre>
(3)	$M \rightarrow \varepsilon$	{ <i>M.quad</i> : = <i>nextquad</i> }
(4)	$S \rightarrow \mathbf{if} \ E \mathbf{then} \ MS_1$	{ backpatch(E.truelist, M.quad); S.nextlist : = merge(E.falselist, S1.nextlist) }
(5)	$S \rightarrow$ while $M_1 E$ do $M_2 S_1$	<pre>{ backpatch(S1.nextlist, M1.quad); backpatch(E.truelist, M2.quad); S.nextlist := E.falselist emit('goto' M1.quad) }</pre>
(6)	$S \rightarrow begin L end$	{ S.nextlist : = L.nextlist }
(7)	$S \rightarrow A$	{ S.nextlist : = nil }

The assignment *S.nextlist* : = **nil** initializes *S.nextlist* to an empty list.

(8) $L \rightarrow L1$; MS { backpatch(L1.nextlist, M.quad); L.nextlist := S.nextlist }

The statement following *L*¹ in order of execution is the beginning of *S*. Thus the *L*1.nextlist list is backpatched to the beginning of the code for *S*, which is given by *M.quad*.

(9) $L \rightarrow S$ { *L.nextlist* := *S.nextlist* }

3.8 PROCEDURE CALLS

The procedure is such an important and frequently used programming construct that it is imperative for a compiler to generate good code for procedure calls and returns. The run-time routines that handle procedure argument passing, calls and returns are part of the run-time support package.

Let us consider a grammar for a simple procedure call statement

- $\Box \qquad S \rightarrow \textbf{call id} (Elist)$
- $\Box \quad Elist \rightarrow Elist, E$
- $\Box \quad Elist \rightarrow E$

Calling Sequences:

The translation for a call includes a calling sequence, a sequence of actions taken on entry to and exit from each procedure. The falling are the actions that take place in a calling sequence :

- When a procedure call occurs, space must be allocated for the activation record of the called procedure.
- The arguments of the called procedure must be evaluated and made available to the called procedure in a known place.
- Environment pointers must be established to enable the called procedure to access data in enclosing blocks.
- The state of the calling procedure must be saved so it can resume execution after the call. \mathbb{E}
- Also saved in a known place is the return address, the location to which the called routine must transfer after it is finished.
- Finally a jump to the beginning of the code for the called procedure must be generated. For example, consider the following syntax-directed translation

 $S \rightarrow$ call id (*Elist*) { for each item p on queue do emit (' param' p); emit ('call' id.place) }

3. *Elist* \rightarrow *Elist* , *E*

{ append *E.place* to the end of *queue* }

4. *Elist* \rightarrow *E*

{ initialize queue to contain only *E.place* }

- ② Here, the code for S is the code for *Elist*, which evaluates the arguments, followed by a **param** p statement for each argument, followed by a **call** statement.
- ② queue is emptied and then gets a single pointer to the symbol table location for the name that denotes the value of E.

CHAPTER IV - CODE GENERATION

The final phase in compiler model is the code generator. It takes as input an intermediate representation of the source program and produces as output an equivalent target program. The code generation techniques presented below can be used whether or not an optimizing phase occurs before code generation.

Position of code generator



Figure 4.1:Code Generator

4.1 ISSUES IN THE DESIGN OF A CODE GENERATOR

The following issues arise during the code generation phase :

- 7. Input to code generator
- 8. Target program
- 9. Memory management
- 10. Instruction selection
- 11. Register allocation
- 12. Evaluation order

9) Input to code generator:

The input to the code generation consists of the intermediate representation of the source program produced by front end, together with information in the symbol table to determine runtime addresses of the data objects denoted by the names in the intermediate representation. Intermediate representation can be :

> Linear representation such as postfix notation Three address representation such as quadruples Virtual machine representation such as stack machine code Graphical representations such as syntax trees and dags.

Prior to code generation, the front end must be scanned, parsed and translated into intermediate representation along with necessary type checking. Therefore, input to code generation is assumed to be error-free.

10) Target program:

The output of the code generator is the target program. The output may be :

Absolute machine language

It can be placed in a fixed memory location and can be executed immediately.

10. Relocatable machine language

It allows subprograms to be compiled separately.

- 11. Assembly language
 - Code generation is made easier.

3. Memory management:

Names in the source program are mapped to addresses of data objects in run-time memory by the front end and code generator. It makes use of symbol table, that is, a name in a three-address statement refers to a symbol-table entry for the name. Labels in three-address statements have to be converted to addresses of instructions. For example,

- *j* : **goto** *i* generates jump instruction as follows :
- ➢ if *i* < *j*, a backward jump instruction with target address equal to location of code for quadruple *i* is generated.
- if i > j, the jump is forward. We must store on a list for quadruple i the location of the first machine instruction generated for quadruple j. When i is processed, the machine locations for all instructions that forward jumps to i are filled.

4. Instruction selection:

The instructions of target machine should be complete and uniform. Instruction speeds and machine idioms are important factors when efficiency of target program is considered. The quality of the generated code is determined by its speed and size. The former statement can be translated into the latter statement as shown below:



Figure 4.2:Instruction selection

5. Register allocation

Instructions involving register operands are shorter and faster than those involving operands in memory. The use of registers is subdivided into two subproblems :

- Register allocation the set of variables that will reside in registers at a point in the program is selected
- Register assignment the specific register that a variable will reside in is picked.

Certain machine requires even-odd *register pairs* for some operands and results. For example, consider the division instruction of the form :

D x, y

where, x - dividend even register in even/odd register pair y - divisor even register holds the remainder odd register holds the quotient

6. Evaluation order

The order in which the computations are performed can affect the efficiency of the target code. Some computation orders require fewer registers to hold intermediate results than others.

4.2 TARGET MACHINE

Familiarity with the target machine and its instruction set is a prerequisite for designing a good code generator. The target computer is a byte-addressable machine with 4 bytes to a word. It has *n* general-purpose registers, $R_0, R_1, \ldots, R_{n-1}$.

It has two-address instructions of the form:

op source, destination where, *op* is an op-code, and *source* and *destination* aredata fields.

It has the following op-codes :

MOV (move *source* to *destination*) ADD (add *source* to *destination*) SUB (subtract *source* from *destination*)

The *source* and *destination* of an instruction are specified by combining registers and memory locations with address modes.

MODE	FORM	ADDRESS	ADDED COST
Absolute	М	М	1
Register	R	R	0
Indexed	$c(\mathbf{R})$	c+contents(R)	1
indirect register	*R	contents (R)	0
indirect indexed	* <i>c</i> (R)	<pre>contents(c+ contents(R))</pre>	1
Literal	# <i>c</i>	С	1

Fable 4.1:Address modes	s with their	assembly-	language forms
-------------------------	--------------	-----------	----------------

For example : MOV R₀, M stores contents of Register R₀ into memory location M; MOV 4(R₀), M stores the value *contents*(4+*contents*(R₀)) into M.

Instruction costs :

Instruction cost = 1 + cost for source and destination address modes. This cost corresponds to the length of the instruction. Address modes involving registers have cost zero. Address modes involving memory location or literal have cost one. Instruction length should be minimized if space is important. Doing so also minimizes the time taken to fetch and perform the instruction. For example : MOV R0, R1 copies the contents of register R0 into R1. It has cost one, since it occupies only one word of memory.

The three-address statement $\mathbf{a} := \mathbf{b} + \mathbf{c}$ can be implemented by many different instruction sequences :

i) MOV b, Ro	
ADD c, Ro	$\cos t = 6$
MOV R ₀ , a	
ii) MOV b, a	
ADD c, a	$\cos t = 6$
ii) Assuming Ro, R	1 and R2 contain

ii the addresses of a, b, and c :

= 2

MOV *R1, *R0	
ADD *R2, *R0	cost

In order to generate good code for target machine, we must utilize its addressing capabilities efficiently.

4.3 RUN-TIME STORAGE MANAGEMENT

Information needed during an execution of a procedure is kept in a block of storage called an activation record, which includes storage for names local to the procedure. The two standard storage allocation strategies are:

Static allocation Stack allocation

In static allocation, the position of an activation record in memory is fixed at compile time. In stack allocation, a new activation record is pushed onto the stack for each execution of a procedure. The record is popped when the activation ends.

The following three-address statements are associated with the run-time allocation and deallocation of activation records:

Call. Return. Halt, and Action, a placeholder for other statements. We assume that the run-time memory is divided into areas for: Code Static data Stack

4.3.1 Static allocation

Implementation of call statement:

The codes needed to implement static allocation are as follows:

MOV #here +20, callee.static_area /*It saves return address*/

GOTO callee.code_area /*It transfers control to the target code for the called procedure */

where,

callee.static_area – Address of the activation record *callee.code_area* – Address of the first instruction for called procedure *#here* +20 – Literal return address which is the address of the instruction following GOTO.

Implementation of return statement:

A return from procedure *callee* is implemented by :

GOTO **callee.static_area*

This transfers control to the address saved at the beginning of the activation record.

Implementation of action statement:

The instruction ACTION is used to implement action statement.

Implementation of halt statement:

The statement HALT is the final instruction that returns control to the operating system.

4.3.2 Stack allocation

Static allocation can become stack allocation by using relative addresses for storage in activation records. In stack allocation, the position of activation record is stored in register so words in activation records can be accessed as offsets from the value in this register.

The codes needed to implement stack allocation are as follows: **Initialization of stack:**

MOV #stackstart, SP /* initializes stack */

Code for the first procedure

HALT /* terminate execution */ Implementation of Call statement:

ADD #caller.recordsize, SP /* increment stack pointer */

MOV #here +16, *SP /*Save return address */

GOTO *callee.code_area*

where, *caller.recordsize* – size of the activation record *#here* +16 – address of the instruction following the **GOTO**

Implementation of Return statement:

GOTO *0 (SP) /*return to the caller */

SUB #caller.recordsize, SP /* decrement SP and restore to previous value */

4.4 BASIC BLOCKS AND FLOW GRAPHS

4.4.1 Basic Blocks

A *basic block* is a sequence of consecutive statements in which flow of control enters at the beginning and leaves at the end without any halt or possibility of branching except at the end. The following sequence of three-address statements forms a basic block: $t_1:=a^*a$

 $\begin{array}{l} t_2 := a \, * \, b \\ t_3 := 2 \, * \, t_2 \\ t_4 := t_1 + t_3 \\ t_5 := b \, * \, b \\ t_6 := t_4 + t_5 \end{array}$

Basic Block Construction:

Algorithm: Partition into basic blocks
Input: A sequence of three-address statements
Output: A list of basic blocks with each three-address statement in exactly one block
Method:

(9) We first determine the set of *leaders*, the first statements of basic blocks. The rules we use are of the following:
The first statement is a leader.
Any statement that is the target of a conditional or unconditional goto is a leader.
Any statement that immediately follows a goto or conditional goto statement is a leader.

(10) For each leader, its basic block consists of the leader and all statements up to but not including the next leader or the end of the program.

begin
 prod :=0;
 i:=1; do
 begin
 prod :=prod+ a[i]* b[i];
 i :=i+1;
 end
 while i <= 20
end</pre>

Consider the following source code for dot product of two vectors a and b of length 20

Figure 4.4:Source code for dot product

The three-address code for the above source program is given as :

(1)	prod := 0
(2)	i := 1
(3)	t ₁ := 4* i
(4)	t ₂ := a[t ₁] /*compute a[i] */
(5)	t ₃ := 4*i
(6)	t ₄ := b[t ₃] /*compute b[i] */
(7)	$t_5 := t_2^* t_4$
(8)	$t_6 := prod+t_5$
(9)	prod := t ₆
(10)	t ₇ := i+1
(11)	i := t ₇
(12)	if i<=20 goto (3)

Figure 4.5:Three address code for above fig4.4 Basic block 1: Statement (1) to (2)

Basic block 2: Statement (3) to (12)

4.4.2 Transformations on Basic Blocks:

A number of transformations can be applied to a basic block without changing the set of expressions computed by the block. Two important classes of transformation are :

- 1. Structure-preserving transformations
- 2. Algebraic transformations

1. Structure preserving transformations:

a. Common subexpression elimination:

a:=b+c	$\mathbf{a} := \mathbf{b} + \mathbf{c}$
$\mathbf{b}:=\mathbf{a}-\mathbf{d}$	$\mathbf{b}:=\mathbf{a}-\mathbf{d}$
c := b + c	c := b + c
d:=a-d	d := b

Since the second and fourth expressions compute the same expression, the basic block can be transformed as above.

b) Dead-code elimination:

Suppose x is dead, that is, never subsequently used, at the point where the statement x : = y + z appears in a basic block. Then this statement may be safely removed without changing the value of the basic block.

c) Renaming temporary variables:

A statement $\mathbf{t} := \mathbf{b} + \mathbf{c}$ (t is a temporary) can be changed to $\mathbf{u} := \mathbf{b} + \mathbf{c}$ (u is a new temporary) and all uses of this instance of t can be changed to u without changing the value of the basic block.

Such a block is called a *normal-form block*.

d) Interchange of statements:

Suppose a block has the following two adjacent statements:

t1 := b + ct2 := x + y

We can interchange the two statements without affecting the value of the block if and only if neither \mathbf{x} nor \mathbf{y} is \mathbf{t}_1 and neither \mathbf{b} nor \mathbf{c} is \mathbf{t}_2 .

8. Algebraic transformations:

Algebraic transformations can be used to change the set of expressions computed by a basic block into an algebraically equivalent set.

Examples:

i) x := x + 0 or x := x * 1 can be eliminated from a basic block without changing the set of expressions it computes.

ii) The exponential statement x := y * * 2 can be replaced by x := y * y.

4.4.2 Flow Graphs

Flow graph is a directed graph containing the flow-of-control information for the set of basic blocks making up a program. The nodes of the flow graph are basic blocks. It has a distinguished initial node. E.g.: Flow graph for the vector dot product is given as follows:



Figure 4.6:Flow graph for vector dot product

B₁ is the *initial* node. B₂ immediately follows B₁, so there is an edge from B₁ to B₂. The target of jump from last statement of B₁ is the first statement B₂, so there is an edge from B₁ (last statement) to B₂ (first statement). B₁ is the *predecessor* of B₂, and B₂ is a *successor* of B₁.

4.4.4 Loops

A loop is a collection of nodes in a flow graph such that All nodes in the collection are *strongly connected*. The collection of nodes has a unique *entry*. A loop that contains no other loops is called an inner loop.

<u>4.5</u> NEXT-USE INFORMATION

If the name in a register is no longer needed, then we remove the name from the register and the register can be used to store some other names. Input: Basic block B of three-address statements

Output: At each statement i: x= y op z, we attach to i the liveliness and next-uses of x, y and z.

Method: We start at the last statement of B and scan backwards.

- 1. Attach to statement i the information currently found in the symbol table regarding the next-use and liveliness of x, y and z.
- 2. In the symbol table, set x to "not live" and "no next use".
- 3. In the symbol table, set y and z to "live", and next-uses of y and z to i.

Figure 4.7:Next-Use Information

 Table 4.2: Symbol Table:

Names	Liveliness	Next-use
х	not live	no next-use
У	Live	i
Z	Live	i

<u>4.6</u> A SIMPLE CODE GENERATOR

A code generator generates target code for a sequence of three- address statements and effectively uses registers to store operands of the statements. For example: consider the three-address statement $\mathbf{a} := \mathbf{b} + \mathbf{c}$ It can have the following sequence of codes:

ADD Rj, Ri		$Cost = 1 // if R_i \text{ contains } b \text{ and } R_j \text{ contains } c$	
	(or)		
ADD c, Ri		Cost = 2	// if c is in a memory location
	(or)		
MOV c, R _j		Cost = 3	// move c from memory to Rj and add
ADD Rj, Ri			

Register and Address Descriptors:

A register descriptor is used to keep track of what is currently in each registers. The register descriptors show that initially all the registers are empty. An address descriptor stores the location where the current value of the name can be found at run time.
A code-generation algorithm:

The algorithm takes as input a sequence of three -address statements constituting a basic block. For each three-address statement of the form x := y op z, perform the following actions:

- \Box Invoke a function *getreg* to determine the location L where the result of the computation y op z should be stored.
- \Box Consult the address descriptor for y to determine y', the current location of y. Prefer the register for y' if the value of y is currently both in memory and a register. If the value of y is not already in L, generate the instruction **MOV** y', L to place a copy of y in L.
- \Box Generate the instruction **OP z'**, **L** where z' is a current location of z. Prefer a register to a memory location if z is in both. Update the address descriptor of x to indicate that x is in location L. If x is in L, update its descriptor and remove x from all other descriptors.
- \Box If the current values of y or z have no next uses, are not live on exit from the block, and are in registers, alter the register descriptor to indicate that, after execution of x : = y op z , those registers will no longer contain y or z.

Generating Code for Assignment Statements:

- = The assignment d := (a-b) + (a-c) + (a-c) might be translated into the following threeaddress code sequence:
 - t := a bu := a cv := t + ud := v + u

with d live at the end.

Table 4.3: Code sequence for the example:

Statements	Code Generated	Register descriptor	Address descriptor
		Register empty	
t : = a - b	MOV a, Ro SUB b, RO	Ro contains t	t in Ro
u : = a - c	MOV a , R1 SUB c , R1	Ro contains t R1 contains u	t in Ro u in R1
v : =t + u	ADD R1, R0	Ro contains v R1 contains u	u in R1 v in R0
$\mathbf{d}:=\mathbf{v}+\mathbf{u}$	ADD R1, R0 MOV R0, d	Ro contains d	d in R0 d in R0 and memory

Generating Code for Indexed Assignments

The table shows the code sequences generated for the indexed assignment

statements **a** : = **b** [**i**] and **a** [**i**] : = **b**

Table 4.4: Indexed assignment

Statements	Code Generated	Cost
a := b[i]	MOV b(Ri), R	2
a[i] : = b	MOV b, a(Ri)	3

Generating Code for Pointer Assignments

The table shows the code sequences generated for the pointer assignments

 $\mathbf{a} := \mathbf{p}$ and $\mathbf{p} := \mathbf{a}$

 Table 4.5:Pointer assignment

Statements	Code Generated	Cost
a : = *p	MOV *R _p , a	2
*p : = a	MOV a, *R _p	2

 Table 4.6:Generating Code for Conditional Statements

Statement	Code
if x < y goto z	CMP x, y CJ< z /* jump to z if condition code is negative */
x : = y +z if x <0 goto z	MOV y, R0 ADD z, R0 MOV R0,x CJ< z

4.7 THE DAG REPRESENTATION FOR BASIC BLOCKS

A DAG for a basic block is a **directed acyclic graph** with the following labels on nodes:
 Leaves are labeled by unique identifiers, either variable names or constants.
 Interior nodes are labeled by an operator symbol.
 Nodes are also optionally given a sequence of identifiers for labels to store the computed values.

DAGs are useful data structures for implementing transformations on basic blocks.

- It gives a picture of how the value computed by a statement is used in subsequent statements.
- > It provides a good way of determining common sub expressions.

Algorithm for construction of DAG

```
Input: A basic block
Output: A DAG for the basic block containing the following information:
            1. A label for each node. For leaves, the label is an identifier. For interior nodes,
                an operator symbol.
            2. For each node a list of attached identifiers to hold the computed values.
  Case (i)x := y OP z
  Case (ii)x := OP y
  Case (iii)x := y
Method:
Step 1: If y is undefined then create node(y).
        If z is undefined, create node(z) for case(i).
Step 2: For the case(i), create a node(OP) whose left child is node(y) and right child is
         node(z). (Checkingfor common sub expression). Let n be this node.
        For case(ii), determine whether there is node(OP) with one child node(y). If not create
        such a node.
        For case(iii), node n will be node(y).
Step 3: Delete x from the list of identifiers for node(x). Append x to the list of attached
        identifiers for the noden found in step 2 and set node(x) to n.
```

Figure 4.8: Algorithm for construction of DAG

Example: Consider the block of three- address statements:

4. $t_1 := 4^* i$ 5. $t_2 := a[t_1]$ 6. $t_3 := 4^* i$ 7. $t_4 := b[t_3]$ 8. $t_5 := t_2^* t_4$ 9. $t_6 := prod+t_5$ 10. $prod := t_6$ 11. $t_7 := i+1$ 12. $i := t_7$ 13. if i<=20 goto (1)

Stages in DAG Construction







Application of DAGs:

- 5. We can automatically detect common sub expressions.
- 6. We can determine which identifiers have their values used in the block.
- 7. We can determine which statements compute values that could be used outside the block.

GENERATING CODE FROM DAGs

The advantage of generating code for a basic block from its dag representation is that, from a dag we can easily see how to rearrange the order of the final computation sequence than we can starting from a linear sequence of three-address statements or quadruples.

Rearranging the order

The order in which computations are done can affect the cost of resulting object code.

For example, consider the following basic block:

 $t_1 := a + b$ $t_2 := c + d$ $t_3 := e - t_2$ $t_4 := t_1 - t_3$

Generated code sequence for basic block:

MOV a, R_0 ADD b, R_0 MOV c, R_1 ADD d, R_1 MOV R_0 , t_1 MOV e, R_0 SUB R_1 , R_0 MOV t_1 , R_1 SUB R_0 , R_1 MOV R_1 , t_4

Rearranged basic block:

Now t1 occurs immediately before t4.

 $t_{2} := c + d$ $t_{3} := e - t_{2}$ $t_{1} := a + b$ $t_{4} := t_{1} - t_{3}$

Revised code sequence:

MOV c , R0 ADD d , R0 MOV a , R0 SUB R0 , R1 MOV a , R0 ADD b , R0 SUB R1 , R0 MOV R0 , t4

In this order, two instructions MOV Ro, t1 and MOV t1, R1 have been saved.

A Heuristic ordering for Dags

The heuristic ordering algorithm attempts to make the evaluation of a node immediately follow the evaluation of its leftmost argument.

The algorithm shown below produces the ordering in reverse.

Algorithm:

- 8. while unlisted interior nodes remain do begin
- 9. select an unlisted node n, all of whose parents have been listed;
- 10. list n;
- while the leftmost child m of n has no unlisted parents and is not a leaf do begin
- 12. list m;
- 13. n : = m

end

end

Example: Consider the DAG shown below:



Figure 4.9:Example DAG

Initially, the only node with no unlisted parents is 1 so set n=1 at line (2) and list 1 at line (3). Now, the left argument of 1, which is 2, has its parents listed, so we list 2 and set n=2 at line (6). Now, at line (4) we find the leftmost child of 2, which is 6, has an unlisted parent 5. Thus we select anew n at line (2), and node 3 is the only candidate. We list 3 and proceed down its left chain, listing 4, 5 and 6. This leaves only 8 among the interior nodes so we list that.

The resulting list is 1234568 and the order of evaluation is 8654321.

Code sequence:

 $ts := d + e t_{6} :$ = a + b t_{5} := t_{6} - c t_{4} := t_{5} * t_{8} t_{3} := t_{4} - e t_{2} : = t_{6} + t_{4} t_{1} := t_{2} * t_{3}

This will yield an optimal code for the DAG on machine whatever be the number of registers.

CHAPTER V - CODE OPTIMIZATION

5.1 INTRODUCTION

The code produced by the straight forward compiling algorithms can often be made to run faster or take less space, or both. This improvement is achieved by program transformations that are traditionally called optimizations. Compilers that apply code-improving transformations are called optimizing compilers.

Optimizations are classified into two categories. They are Machine independent optimizations: Machine dependant optimizations:

Machine independent optimizations:

Machine independent optimizations are program transformations that improve the target code without taking into consideration any properties of the target machine.

Machine dependant optimizations:

Machine dependant optimizations are based on register allocation and utilization of special machine-instruction sequences.

The criteria for code improvement transformations:

- ✓ Simply stated, the best program transformations are those that yield the most benefit for the least effort.
- ✓ The transformation must preserve the meaning of programs. That is, the optimization must not change the output produced by a program for a given input, or cause an error such as division by zero, that was not present in the original source program. At all times we take the "safe" approach of missing an opportunity to apply a transformation rather than risk changing what the program does.

- ✓ A transformation must, on the average, speed up programs by a measurable amount. We are also interested in reducing the size of the compiled code although the size of the code has less importance than it once had. Not every transformation succeeds in improving every program, occasionally an "optimization" may slow down a program slightly.
- ✓ The transformation must be worth the effort. It does not make sense for a compiler writer to expend the intellectual effort to implement a code improving transformation and to have the compiler expend the additional time compiling source programs if this effort is not repaid when the target programs are executed. "Peephole" transformations of this kind are simple enough and beneficial enough to be included in any compiler.

Organization for an Optimizing Compiler:



Figure 5.1:Organization for an Optimizing Compiler

Flow analysis is a fundamental prerequisite for many important types of code improvement. Generally control flow analysis precedes data flow analysis.Control flow analysis (CFA) represents flow of control usually in form of graphs, CFA constructs such as

- 1. control flow graph
- 2. Call graph

Data flow analysis (DFA) is the process of ascerting and collecting information prior to program execution about the possible modification, preservation, and use of certain entities (such as values or attributes of variables) in a computer program.

5.2 PRINCIPAL SOURCES OF OPTIMISATION

A transformation of a program is called local if it can be performed by looking only at the statements in a basic block; otherwise, it is called global. Many transformations can be perormed at both the local and global levels. Local transformations are usually performed first.

5.2.1 Function-Preserving Transformations

There are a number of ways in which a compiler can improve a program without changing the function it computes. The transformations

- ✓ Common sub expression elimination,
- \checkmark Copy propagation,
- ✓ Dead-code elimination, and
- ✓ Constant folding

are common examples of such function-preserving transformations. The other transformations come up primarily when global optimizations are performed.

Frequently, a program will include several calculations of the same value, such as an offset in an array. Some of the duplicate calculations cannot be avoided by the programmer because they lie below the level of detail accessible within the source language.

1.Common Sub expressions elimination:

An occurrence of an expression E is called a common sub-expression if E was previously computed, and the values of variables in E have not changed since the previous computation. We can avoid recomputing the expression if we can use the previously computed value. For example

```
t1: =4*i
t2: =a [t1]
t3: =4*j
t4:=4*i
t5: =n
t6: =b [t4] +t5
```

The above code can be optimized using the common sub-expression elimination as

```
t1: =4*i
t2: =a [t1]
t3: =4*j
t5: =n
t6: =b [t1] +t5
```

The common sub expression t 4: =4*i is eliminated as its computation is already in t1. And value of i is not been changed from definition to use.

2.Copy Propagation:

Assignments of the form f := g called copy statements, or copies for short. The idea behind the copy-propagation transformation is to use g for f, whenever possible after the copy statement f := g. Copy propagation means use of one variable instead of another. This may not appear to be an improvement, but as we shall see it gives us an opportunity to eliminate x. For example

x=Pi;

A=x*r*r;The optimization using copy propagation can be done as follows: A=Pi*r*r;Here the variable x is eliminated

3.Dead-Code Eliminations:

A variable is live at a point in a program if its value can be used subsequently; otherwise, it is dead at that point. A related idea is dead or useless code, statements that compute the values that never get used. While the programmer is unlikely to introduce any dead code intentionally, it may appear as the result of previous transformations. An optimization can be done by eliminating dead code.

Example:

i=0; if(i=1) { a=b+5; }

Here, 'if'statement is dead code because this condition will never get satisfied.

4.Constant folding:

We can eliminate both the test and printing from the object code. More generally, deducing at compile time that the value of an expression is a constant and using the constant instead is known as constant folding. One advantage of copy propagation is that it often turns the copy statement into dead code. For example,

a=3.14157/2 can be replaced by a=1.570 there by eliminating a division operation.

5.2.2 Loop Optimizations:

We now give a brief introduction to a very important place for optimizations, namely loops, especially the inner loops where programs tend to spend the bulk of their time. The running time of a program may be improved if we decrease the number of instructions in an inner loop, even if we increase the amount of code outside that loop. Three techniques are important for loop optimization:

code motion, which moves code outside a loop;Induction-variable elimination, which we apply to replace variables from inner loop.Reduction in strength, which replaces and expensive operation by a cheaper one, such as a multiplication by an addition.

1.Code Motion:

An important modification that decreases the amount of code in a loop is code motion. This transformation takes an expression that yields the same result independent of the number of times a loop is executed (a loop-invariant computation) and places the expression before the loop. Note that the notion "before the loop" assumes the existence of an entry for the loop. For example, evaluation of limit-2 is a loop-invariant computation in the following while-statement:

while (i <= limit-2) /* statement does not change limit*/

Code motion will result in the equivalent of

t= limit-2; while (i<=t) /* statement does not change limit or t */

2.Induction Variables :

Loops are usually processed inside out. For example consider the loop around B3. Note that the values of j and t4 remain in lock-step; every time the value of j decreases by 1, that of t4 decreases by 4 because 4*j is assigned to t4. Such identifiers are called induction variables. When there are two or more induction variables in a loop, it may be possible to get rid of all but one, by the process of induction-variable elimination. For the inner loop around B3 in Fig. we cannot get rid of either j or t4 completely; t4 is used in B3 and j in B4.

However, we can illustrate reduction in strength and illustrate a part of the process of induction-variable elimination. Eventually j will be eliminated when the outer loop of B2 - B5 is considered. Example: As the relationship t4:=4*j surely holds after such an assignment to t4 in Fig. and t4 is not changed elsewhere in the inner loop around B3, it follows that just after the statement j:=j -1 the relationship t4:= 4*j-4 must hold. We may therefore replace the assignment t 4:= 4*j by t4:= t4-4. The only problem is that t 4 does not have a value when we enter block B3 for the first time. Since we must maintain the relationship t4=4*j on entry to the block B3, we place an initializations of t4 at the end of the block where j itself is



Figure 5.2:Induction variable example

initialized, shown by the dashed addition to block B1 in second Fig. The replacement of a multiplication by a subtraction will speed up the object code if multiplication takes more time than addition or subtraction, as is the case on many machines.

3.Reduction In Strength:

Reduction in strength replaces expensive operations by equivalent cheaper ones on the target machine. Certain machine instructions are considerably cheaper than others and can often be used as special cases of more expensive operators. For example, x^2 is invariably cheaper to implement as x^*x than as a call to an exponentiation routine. Fixed-point multiplication or division by a power of two is cheaper to implement as a shift. Floating-point division by a constant can be implemented as multiplication by a constant, which may be cheaper.

5.2.3 OPTIMIZATION OF BASIC BLOCKS

There are two types of basic block optimizations. They are :

- ⁽²⁾ Structure-Preserving Transformations
- ② Algebraic Transformations

Structure-Preserving Transformations:

The primary Structure-Preserving Transformation on basic blocks are:

- Common sub-expressionelimination
- \Im Dead code elimination
- $\mathfrak{M}\mathfrak{O}$ Renaming of temporary variables
- $\mathfrak{L} \mathfrak{D}$ Interchange of two independent adjacent statements.

1.Common sub-expression elimination:

Common sub expressions need not be computed over and over again. Instead they can be computed once and kept in store from where it's referenced when encountered again – of course providing the variable values in the expression still remain constant.

Example:

 $\Box =b+c$ $\Box =a-d$ $\Box =b+c$ $\Box =a-d$

The 2^{nd} and 4^{th} statements compute the same expression: b+c and a-d

Basic block can be transformed to

 $\Box =b+c$ $\Box =a-d$ $\Box =a$ $\Box =b$

2. Dead code elimination:

It's possible that a large amount of dead (useless) code may exist in the program. This might be especially caused when introducing variables and procedures as part of construction or error -correction of a program – once declared and defined, one forgets to remove them in case they serve no purpose. Eliminating these will definitely optimize the code.

3.Renaming of temporary variables:

A statement t:=b+c where t is a temporary name can be changed to u:=b+c where u is another temporary name, and change all uses of t to u. In this we can transform a basic block to its equivalent block called normal-form block.

4.Interchange of two independent adjacent statements:

Two statements

 $t_1:=b+c$

 $t_2:=x+y$

can be interchanged or reordered in its computation in the basic block when value of t1 does not affect the value of t2.

Algebraic Transformations:

Algebraic identities represent another important class of optimizations on basic blocks. This includes simplifying expressions or replacing expensive operation by cheaper ones i.e. reduction in strength. Another class of related optimizations is constant folding. Here we evaluate constant expressions at compile time and replace the constant expressions by their values. Thus the expression 2*3.14 would be replaced by 6.28.

The relational operators $\langle =, \rangle =, \langle, \rangle, +$ and = sometimes generate unexpected common sub expressions. Associative laws may also be applied to expose common sub expressions. For example, if the source code has the assignments

a :=b+c e :=c+d+b

the following intermediate code may be generated:

a :=b+ct :=c+d:=t+b

Example:

x:=x+0 can be removed

x:=y**2 can be replaced by a cheaper statement x:=y*y

The compiler writer should examine the language carefully to determine what rearrangements of computations are permitted, since computer arithmetic does not always obey the algebraic identities of mathematics. Thus, a compiler may evaluate x*y-x*z as x*(y-z) but it may not evaluate a+(b-c) as (a+b)-c.

5.3 LOOPS IN FLOW GRAPH

A graph representation of three-address statements, called a **flow graph**, is useful for understanding code-generation algorithms, even if the graph is not explicitly constructed by a code-generation algorithm. Nodes in the flow graph represent computations, and the edges represent the flow of control.

Dominators:

In a flow graph, a node d dominates node n, if every path from initial node of the flow graph to n goes through d. This will be denoted by *d dom n*. Every initial node dominates all the remaining nodes in the flow graph and the entry of a loop dominates all nodes in the loop. Similarly every node dominates itself.

Example:

*In the flow graph below,

*Initial node, node1 dominates every

node. *node 2 dominates itself

*node 3 dominates all but 1 and 2.

*node 4 dominates all but 1,2 and 3.

*node 5 and 6 dominates only themselves, since flow of control can skip around either by goin through the other.

*node 7 dominates 7,8,9 and 10.

*node 8 dominates 8,9 and 10.

*node 9 and 10 dominates only themselves.



Figure 5.3:Flow graph

The way of presenting dominator information is in a tree, called the dominator tree in which the initial node is the root. The parent of each other node is its immediate dominator. Each node d dominates only its descendents in the tree. The existence of dominator tree follows from a property of dominators; each node has a unique immediate dominator in that is the last dominator of n on any path from the initial node to n. In terms of the dom relation, the immediate dominator m has the property is d=!n and d dom n, then d dom m.



Figure 5.4:Dominator Tree

- $D(1) = \{1\}$
- $D(2) = \{1, 2\}$
- $D(3) = \{1,3\}$
- $D(4) = \{1,3,4\}$
- $D(5) = \{1,3,4,5\}$
- $D(6) = \{1,3,4,6\}$
- $D(7) = \{1,3,4,7\}$
- $D(8) = \{1,3,4,7,8\}$
- $D(9) = \{1,3,4,7,8,9\}$
- D(10)={1,3,4,7,8,10}

Natural Loop:

One application of dominator information is in determining the loops of a flow graph suitable for improvement. The properties of loops are

A loop must have a single entry point, called the header. This entry point-dominates all nodes in the loop, or it would not be the sole entry to the loop.

There must be at least one wayto iterate the loop(i.e.)at least one path back to the header.

One way to find all the loops in a flow graph is to search for edges in the flow graph whose heads dominate their tails. If $a \rightarrow b$ is an edge, b is the head and a is the tail. These types of edges are called as back edges. Example:

In the above graph,

 $7 \rightarrow 4 \qquad 4 \text{ DOM } 7$ $10 \rightarrow 7 \qquad 7 \text{ DOM } 10$ $4 \rightarrow 3$ $8 \rightarrow 3$ $9 \rightarrow 1$

The above edges will form loop in flow graph. Given a back edge $n \rightarrow d$, we define the natural loop of the edge to be d plus the set of nodes that can reach n without going through d. Node d is the header of the loop.

Algorithm: Constructing the natural loop of a back edge.

Input: A flow graph G and a back edge $n \rightarrow d$.

Output: The set loop consisting of all nodes in the natural loop $n \rightarrow d$.

Method: Beginning with node n, we consider each node m*d that we know is in loop, to make sure that m's predecessors are also placed in loop. Each node in loop, except for d, is placed once on stack, so its predecessors will be examined. Note that because d is put in the loop initially, we never examine its predecessors, and thus find only those nodes that reach n without going through d.

```
Procedure insert(m);
if m is not in loop then
    begin loop := loop U
    {m}; push m onto stack
```

end;

```
stack : =empty;
```

loop:

If we use the natural loops as "the loops", then we have the useful property that unless two loops have the same header, they are either disjointed or one is entirely contained in the other. Thus, neglecting loops with the same header for the moment, we have a natural notion of inner loop: one that contains no other loop.

When two natural loops have the same header, but neither is nested within the other, they are combined and treated as a single loop.

Pre-Headers:

Several transformations require us to move statements "before the header". Therefore begin treatment of a loop L by creating a new block, called the preheater. The pre-header has only the header as successor, and all edges which formerly entered the header of Lfrom outside L instead enter the pre-header. Edges from inside loop L to the header are not changed. Initially the pre-header is empty, but transformations on L may place statements in it.



Figure 5.5:Pre-Header

Reducible flow graphs:

Reducible flow graphs are special flow graphs, for which several code optimization transformations are especially easy to perform, loops are unambiguously defined, dominators can be easily calculated, data flow analysis problems can also be solved efficiently. Exclusive use of structured flow-of-control statements such as if-then-else, while-do, continue, and break statements produces programs whose flow graphs are always reducible. The most important properties of reducible flow graphs are that there are no jumps into the middle of loops from outside; the only entry to a loop is through its header.

Definition:

A flow graph G is reducible if and only if we can partition the edges into two disjoint groups, *forward* edges and *back* edges, with the following properties.

- ✓ The forward edges from an acyclic graph in which every node can be reached from initial node of G.
- \checkmark The back edges consist only of edges where heads dominate theirs tails.
- \checkmark Example: The above flow graph is reducible.

If we know the relation DOM for a flow graph, we can find and remove all the back edges. The remaining edges are forward edges. If the forward edges form an acyclic graph, then we can say the flow graph reducible. In the above example remove the five back edges $4\rightarrow 3$, $7\rightarrow 4$, $8\rightarrow 3$, $9\rightarrow 1$ and $10\rightarrow 7$ whose heads dominate their tails, the remaining graph is acyclic. The key property of reducible flow graphs for loop analysis is that in such flow graphs every set of nodes that we would informally regard as a loop must contain a back edge.

5.4 PEEPHOLE OPTIMIZATION

A statement-by-statement code-generations strategy often produce target code that contains redundant instructions and suboptimal constructs .The quality of such target code can be improved by applying "optimizing" transformations to the target program. A simple but effective technique for improving the target code is peephole optimization, a method for trying to improving the performance of the target program by examining a short sequence of target instructions (called the peephole) and replacing these instructions by a shorter or faster sequence, whenever possible. The peephole is a small, moving window on the target program. The code in the peephole need not contiguous, although some implementations do require this.it is characteristic of peephole optimization that each improvement may spawn opportunities for additional improvements.

We shall give the following examples of program transformations that are characteristic of peephole optimizations:

Redundant-instructions elimination Flow-of-control optimizations Algebraic simplifications Use of machine idioms Unreachable Code

Redundant Loads And Stores:

If we see the instructions sequence

- (1) MOV R0,a
- (2) MOV a,R0

we can delete instructions (2) because whenever (2) is executed. (1) will ensure that the value of **a** is already in register R₀.If (2) had a label we could not be sure that (1) was always executed immediately before (2) and so we could not remove (2). Unreachable Code:

Another opportunity for peephole optimizations is the removal of unreachable instructions. An unlabeled instruction immediately following an unconditional jump may be removed. This operation can be repeated to eliminate a sequence of instructions. For example, for debugging purposes, a large program may have within it certain segments that are executed only if a variable **debug** is 1. In C, the source code might look like:

#define debug
0
If (debug) {
Print debugging information
}
termediate representations the

In the intermediate representations the if-statement may be translated as: If

```
debug =1 goto L2
goto L2
L1: print debugging information
L2: .....(a)
```

One obvious peephole optimization is to eliminate jumps over jumps .Thus no matter what the value of **debug**; (a) can be replaced by:

If debug ≠1 goto L2 Print debugging information

(b)

As the argument of the statement of (b) evaluates to a constant true it can be replaced by

If debug $\neq 0$ goto L2

Print debugging information

L2:

As the argument of the first statement of (c) evaluates to a constant true, it can be replaced by goto L2. Then all the statement that print debugging aids are manifestly unreachable and can be eliminated one at a time.

.....(c)

Flows-Of-Control Optimizations:

The unnecessary jumps can be eliminated in either the intermediate code or the target code by the following types of peephole optimizations. We can replace the jump sequence

goto L1 L1: gotoL2 by the sequence goto L2 L1: goto L2

If there are now no jumps to L1, then it may be possible to eliminate the statement L1:goto L2 provided it is preceded by an unconditional jump .Similarly, the sequence

if a < b goto L1 L1: goto L2 can be replaced by Ifa < b goto L2

L1: goto L2

Finally, suppose there is only one jump to L1 and L1 is preceded by an unconditional goto. Then the sequence

goto L1

.

L1: if a <b goto<="" th=""><th>L2</th>	L2
L3:	(1)
Maybe replaced by	
Ifa b goto L2	
goto L3	

L3:(2)

While the number of instructions in (1) and (2) is the same, we sometimes skip the unconditional jump in (2), but never in (1). Thus (2) is superior to (1) in execution time

Algebraic Simplification:

There is no end to the amount of algebraic simplification that can be attempted through peephole optimization. Only a few algebraic identities occur frequently enough that it is worth considering implementing them .For example, statements such as

(d) :=
$$x+0$$

Or
 $x := x * 1$

Are often produced by straightforward intermediate code-generation algorithms, and they can be eliminated easily through peephole optimization.

Reduction in Strength:

Reduction in strength replaces expensive operations by equivalent cheaper ones on the target machine. Certain machine instructions are considerably cheaper than others and can often be used as special cases of more expensive operators.

For example, x^2 is invariably cheaper to implement as x^*x than as a call to an exponentiation routine. Fixed-point multiplication or division by a power of two is cheaper to implement as a shift. Floating-point division by a constant can be implemented as multiplication by a constant, which may be cheaper.

$$X^2 \rightarrow X^*X$$

Useof Machine Idioms:

The target machine may have hardware instructions to implement certain specific operations efficiently. For example, some machines have auto-increment and auto-decrement addressing modes. These add or subtract one from an operand before or after using its value. The use of these modes greatly improves the quality of code when pushing or popping a stack, as in parameter passing. These modes can also be used in code for statements like i :=i+1.

 $i:=i+1 \rightarrow i++$

 $i:=i-1 \rightarrow i--$

5.5 INTRODUCTION TO GLOBAL DATAFLOW ANALYSIS

In order to do code optimization and a good job of code generation, compiler needs to collect information about the program as a whole and to distribute this information to each block in the flow graph. A compiler could take advantage of "reaching definitions", such as knowing where a variable like *debug* was last defined before reaching a given block, in order to perform transformations are just a few examples of data-flow information that an optimizing compiler collects by a process known as data-flow analysis.

Data-flow information can be collected by setting up and solving systems of equations of the form :

This equation can be read as " the information at the end of a statement is either generated within the statement, or enters at the beginning and is not killed as control flows through the statement."

The details of how data-flow equations are set and solved depend on three factors.

- The notions of generating and killing depend on the desired information, i.e., on the data flow analysis problem to be solved. Moreover, for some problems, instead of proceeding along with flow of control and defining out[s] in terms of in[s], we need to proceed backwards and define in[s] in terms of out[s].
- Since data flows along control paths, data-flow analysis is affected by the constructs in a program. In fact, when we write out[s] we implicitly assume that there is unique end point where control leaves the statement; in general, equations are set up at the level of basic blocks rather than statements, because blocks do have unique end points.
- There are subtleties that go along with such statements as procedure calls, assignments through pointer variables, and even assignments to array variables.

Points and Paths:

Within a basic block, we talk of the point between two adjacent statements, as well as the point before the first statement and after the last. Thus, block B1 has four points: one before any of the assignments and one after each of the three assignments.



Figure 5.6:Points and paths

Now let us take a global view and consider all the points in all the blocks. A path from p_1 to p_n is a sequence of points p_1 , p_2 ,..., p_n such that for each i between 1 and n-1, either

- Pi is the point immediately preceding a statement and pi+1 is the point immediately following that statement in the same block, or
- \oint P_i is the end of some block and p_{i+1} is the beginning of a successor block.

Reaching definitions:

A definition of variable x is a statement that assigns, or may assign, a value to x. The most common forms of definition are assignments to x and statements that read a value from an i/o device and store it in x. These statements certainly define a value for x, and they are referred to as **unambiguous** definitions of x. There are certain kinds of statements that may define a value for x; they are called **ambiguous** definitions. The most usual forms of **ambiguous** definitions of x are:

- ✤ A call of a procedure with x as a parameter or a procedure that can access x because x is in the scope of the procedure.
- An assignment through a pointer that could refer to x. For example, the assignment *q: = y is a definition of x if it is possible that q points to x. we must assume that an assignment through a pointer is a definition of every variable.

We say a definition d reaches a point p if there is a path from the point immediately following d to p, such that d is not "killed" along that path. Thus a point can be reached by an unambiguous definition and an ambiguous definition of the same variable appearing later along one path.

Data-flow analysis of structured programs:

Flow graphs for control flow constructs such as do-while statements have a useful property: there is a single beginning point at which control enters and a single end point that control leaves from when execution of the statement is over. We exploit this property when we talk of the definitions reaching the beginning and the end of statements with the following syntax.

S \longrightarrow id: = E| S; S | if E then S else S | do S while

 $E E \longrightarrow id + id | id$

Expressions in this language are similar to those in the intermediate code, but the flow graphs for statements have restricted forms.



IF E then S1 else S2 do S1 while E

Figure 5.7:Data flow analysis of structured programs

We define a portion of a flow graph called a *region* to be a set of nodes N that includes a header, which dominates all other nodes in the region. All edges between nodes in N are in the region, except for some that enter the header. The portion of flow graph corresponding to a statement S is a region that obeys the further restriction that control can flow to just one outside block when it leaves the region.

we say that the beginning points of the dummy blocks at the entry and exit of a statement's region are the beginning and end points, respectively, of the statement. The equations are inductive, or syntax-directed, definition of the sets in[S], out[S], gen[S], and kill[S] for all statements S.

gen[S] is the set of definitions "generated" by S while kill[S] is the set of definitions that never reach the end of S.

Consider the following data-flow equations for reaching definitions :



Observe the rules for a single assignment of variable a. Surely that assignment is a definition of a, say d. Thus

Gen[S]={d} On the other hand, d "kills" all other definitions of a, so we write

$$Kill[S] = D_a - \{d\}$$

Where, D_a is the set of all definitions in the program for variable a.





Under what circumstances is definition d generated by $S=S_1$; S_2 ? First of all, if it is generated by S_2 , then it is surely generated by S_1 if d is generated by S_1 , it will reach the end of S provided it is not killed by S_2 . Thus, we write

gen[S]=gen[S2] U (gen[S1]-kill[S2])

Similar reasoning applies to the killing of a definition, so we have

 $Kill[S] = kill[S_2] U (kill[S_1] - gen[S_2])$

Conservative estimation of data-flow information:

There is a subtle miscalculation in the rules for gen and kill. We have made the assumption that the conditional expression E in the if and do statements are "uninterpreted"; that is, there exists inputs to the program that make their branches go either way. We assume that any graph-theoretic path in the flow graph is also an execution path, i.e., a path that is executed when the program is run with least one possible input.

When we compare the computed gen with the "true" gen we discover that the true gen is always a subset of the computed gen. on the other hand, the true kill is always a superset of the computed kill. These containments hold even after we consider the other rules. It is natural to wonder whether these differences between the true and computed gen and kill sets present a serious obstacle to data-flow analysis. The answer lies in the use intended for these data.

Overestimating the set of definitions reaching a point does not seem serious; it merely stops us from doing an optimization that we could legitimately do. On the other hand, underestimating the set of definitions is a fatal error; it could lead us into making a change in the program that changes what the program computes. For the case of reaching definitions, then, we call a set of definitions safe or conservative if the estimate is a superset of the true set of reaching definitions. We call the estimate unsafe, if it is not necessarily a superset of the truth. Returning now to the implications of safety on the estimation of gen and kill for reaching definitions, note that our discrepancies, supersets for gen and subsets for kill are both in the safe direction. Intuitively, increasing gen adds to the set of definitions that can reach a point, and cannot prevent a definition from reaching a place that it truly reached. Decreasing kill can only increase the set of definitions reaching any given point.

Computation of in and out:

any data-flow problems can be solved by synthesized translations similar to those used to compute gen and kill. It can be used, for example, to determine loop-invariant computations. However, there are other kinds of data-flow information, such as the reaching-definition problem. It turns out that in is an inherited attribute, and out is a synthesized attribute depending on in. we intend that in[S] be the set of definitions reaching the beginning of S, taking into account the flow of control throughout the entire program, including statements outside of S or within which S is nested.

The set out[S] is defined similarly for the end of s. it is important to note the distinction between out[S] and gen[S]. The latter is the set of definitions that reach the end of S without following paths outside S. Assuming we know in[S] we compute out by equation, that is

Out[S] = gen[S] U (in[S] - kill[S])

Considering cascade of two statements S_1 ; S_2 , as in the second case. We start by observingin[S_1]=in[S]. Then, we recursively compute out[S_1], which gives us in[S_2], since a definition reaches the beginning of S_2 if and only if it reaches the end of S_1 . Now we can compute out[S_2], and this set is equal to out[S_1].

Considering if-statement we have conservatively assumed that control can follow either branch, a definition reaches the beginning of S₁ or S₂ exactly when it reaches the beginning of S.

 $In[S_1] = in[S_2] = in[S]$

If a definition reaches the end of S if and only if it reaches the end of one or both sub statements; i.e,

Out[S]=out[S1] U out[S2] **Representation of sets:**

Sets of definitions, such as gen[S] and kill[S], can be represented compactly using bit vectors. We assign a number to each definition of interest in the flow graph. Then bit vector representing a set of definitions will have 1 in position I if and only if the definition numbered I is in the set. The number of definition statement can be taken as the index of statement in an array holding pointers to statements. However, not all definitions may be of interest during global data-flow analysis. Therefore the number of definitions of interest will typically be recorded in a separate table.

A bit vector representation for sets also allows set operations to be implemented efficiently. The union and intersection of two sets can be implemented by logical or and logical and, respectively, basic operations in most systems-oriented programminglanguages. The difference A-B of sets A and B can be implemented by taking the complement of B and then using logical and to compute A.

Local reaching definitions:

Space for data-flow information can be traded for time, by saving information only at certain points and, as needed, recomputing information at intervening points. Basic blocks are usually treated as a unit during global flow analysis, with attention restricted to only those points that are the beginnings of blocks. Since there are usually many more points than blocks, restricting our effort to blocks is a significant savings. When needed, the reaching definitions for all points in a block can be calculated from the reaching definitions for the beginning of a block.

Use-definition chains:

It is often convenient to store the reaching definition information as" use-definition chains" or "ud-chains", which are lists, for each use of a variable, of all the definitions that reaches that use. If a use of variable a in block B is preceded by no unambiguous definition of a, then ud-chain for that use of a is the set of definitions in in[B] that are definitions of a.in addition, if there are ambiguous definitions of a ,then all of these for which no unambiguous definition of a lies between it and the use of a are on the ud-chain for this use of a.

Evaluation order:

The techniques for conserving space during attribute evaluation, also apply to the computation of data-flow information using specifications. Specifically, the only constraint on the evaluation order for the gen, kill, in and out sets for statements is that imposed by dependencies between these sets. Having chosen an evaluation order, we are free to release the space for a set after all uses of it have occurred. Earlier circular dependencies between attributes were not allowed, but we have seen that data-flow equations may have circular dependencies.

General control flow:

Data-flow analysis must take all control paths into account. If the control paths are evident from the syntax, then data-flow equations can be set up and solved in a syntax-directed manner. When programs can contain goto statements or even the more disciplined break and continue statements, the approach we have taken must be modified to take the actual control paths into account. Several approaches may be taken. The iterative method works arbitrary flow graphs. Since the flow graphs obtained in the presence of break and continue statements are reducible, such constraints can be handled systematically using the interval-based methods. However, the syntax-directed approach need not be abandoned when break and continue statements are allowed.

5.6 CODE IMPROVIG TRANSFORMATIONS

Algorithms for performing the code improving transformations rely on data-flow information. Here we consider common sub-expression elimination, copy propagation and transformations for moving loop invariant computations out of loops and for eliminating induction variables. Global transformations are not substitute for local transformations; both must be performed.

Elimination of global common sub expressions:

The available expressions data-flow problem discussed in the last section allows us to determine if an expression at point p in a flow graph is a common sub-expression. The following algorithm formalizes the intuitive ideas presented for eliminating common sub-expressions.

✤ ALGORITHM: Global common sub expression elimination.

INPUT: A flow graph with available expression information.

OUTPUT: A revised flow graph.

METHOD: For every statement s of the form $x := y+z^6$ such that y+z is available at the beginning of block and neither y nor r z is defined prior to statement s in that block, do the following.

- **To discover the evaluations of y+z that reach s's block, we follow flow graph edges, searching backward from s's block. However, we do not go through any block that evaluates y+z. Thelast evaluation of y+z in each block encountered is an evaluation of y+z that reaches s.**
- **Create new variable u.**

Replace each statement w: =y+z found in (1) by u : = y + z w : = u

Replace statement s by x:=u.

Some remarks about this algorithm are in order.

The search in step(1) of the algorithm for the evaluations of y+z that reach statement s can also be formulated as a data-flow analysis problem. However, it does not make sense to solve it for all expressions y+z and all statements or blocks because too much irrelevant information is gathered.

- \checkmark Not all changes made by algorithm are improvements. We might wish to limit the number of different evaluations reaching s found in step (1), probably to one.
- \checkmark Algorithm will miss the fact that a*z and c*z must have the same value in

$$a :=x+y$$
 $c :=x+y$
vs
 $b :=a^*z$ $d :=c^*z$

✓ Because this simple approach to common sub expressions considers only the literal expressions themselves, rather than the values computed by expressions.

Copy propagation:

Various algorithms introduce copy statements such as x :=copies may also be generated directly by the intermediate code generator, although most of these involve temporaries local to one block and can be removed by the dag construction. We may substitute y for x in all these places, provided the following conditions are met every such use u of x. Statement s must be the only definition of x reaching u. On every path from s to including paths that go through u several times, there are no assignments to y.

Condition (1) can be checked using ud-changing information. We shall set up a new dataflow analysis problem in which in[B] is the set of copies s: x:=y such that every path from initial node to the beginning of B contains the statement s, and subsequent to the last occurrence of s, there are no assignments to y.

✤ ALGORITHM: Copy propagation.

INPUT: a flow graph G, with ud-chains giving the definitions reaching block B, and with c_in[B] representing the solution to equations that is the set of copies x:=y that reach block B along every path, with no assignment to x or y following the last occurrence of x:=y on the path. We also need ud-chains giving the uses of each definition.

OUTPUT: A revised flow graph.

METHOD: For each copy s : x:=y do the following:

- \checkmark Determine those uses of x that are reached by this definition of namely, s: x: =y.
- ✓ Determine whether for every use of x found in (1), s is in c_in[B], where B is the block of this particular use, and moreover, no definitions of x or y occur prior to this use of x within B. Recall that if s is in c_in[B]then s is the only definition of x that reaches B.
- ✓ If s meets the conditions of (2), then remove s and replace all uses of x found in (1) by y.

Detection of loop-invariant computations:

Ud-chains can be used to detect those computations in a loop that are loop-invariant, that is, whose value does not change as long as control stays within the loop. Loop is a region consisting of set of blocks with a header that dominates all the other blocks, so the only way to enter the loop is through the header.

If an assignment x := y+z is at a position in the loop where all possible definitions of y and z are outside the loop, then y+z is loop-invariant because its value will be the same each time x:=y+z is encountered. Having recognized that value of x will not change, consider v := x+w, where w could only have been defined outside the loop, then x+w is also loop-invariant.

♦ ALGORITHM: Detection of loop-invariant computations.

INPUT: A loop L consisting of a set of basic blocks, each block containing sequence of three-address statements. We assume ud-chains are available for the individual statements.

OUTPUT: the set of three-address statements that compute the same value each time executed, from the time control enters the loop L until control next leaves L.

METHOD: we shall give a rather informal specification of the algorithm, trusting that the principles will be clear.

- ✓ Mark "invariant" those statements whose operands are all either constant or have all their reaching definitions outside L.
- ✓ Repeat step (3) until at some repetition no new statements are marked "invariant".
- ✓ Mark "invariant" all those statements not previously so marked all of whose operands either are constant, have all their reaching definitions outside L, or have exactly one reaching definition, and that definition is a statement in L marked invariant.

Performing code motion:

Having found the invariant statements within a loop, we can apply to some of them an optimization known as code motion, in which the statements are moved to pre-header of the loop. The following three conditions ensure that code motion does not change what the program computes. Consider s: x: =y+z.

- ✓ The block containing s dominates all exit nodes of the loop, where an exit of a loop is a node with a successor not in the loop.
- ✓ There is no other statement in the loop that assigns to x. Again, if x is a temporary assigned only once, this condition is surely satisfied and need not be changed.No use of x in the loop is reached by any definition of x other than s. This condition too will be satisfied, normally, if x is temporary.

✤ ALGORITHM: Code motion.

INPUT: A loop L with ud-chaining information and dominator information.

OUTPUT: A revised version of the loop with a pre-header and some statements moved to the pre-header.

METHOD:

✓ Use loop-invariant computation algorithm to find loop-invariant statements.

✓ For each statement s defining x found in step(1), check:

- i) That it is in a block that dominates all exits of L,
- ii) That x is not defined elsewhere in L, and
- iii) That all uses in L of x can only be reached by the definition of x in statement s.
- ✓ Move, in the order found by loop-invariant algorithm, each statement s found in (1) and meeting conditions (2i), (2ii), (2iii), to a newly created pre-header, provided any operands of s that are defined in loop L have previously had their definition statements moved to the pre-header.

To understand why no change to what the program computes can occur, condition (2i) and (2ii) of this algorithm assure that the value of x computed at s must be the value of x after any exit block of L. When we move s to a pre-header, s will still be the definition of x that reaches the end of any exit block of L. Condition (2iii) assures that any uses of x within L did, and will continue to, use the value of x computed by s.

Alternative code motion strategies:

The condition (1) can be relaxed if we are willing to take the risk that we may actually increase the running time of the program a bit; of course, we never change what the program computes. The relaxed version of code motion condition (1) is that we may move a statement s assigning x only if:

1'. The block containing s either dominates all exists of the loop, or x is not used outside the loop. For example, if x is a temporary variable, we can be sure that the value will be used only in its own block.

If code motion algorithm is modified to use condition (1'), occasionally the running time will increase, but we can expect to do reasonably well on the average. The modified algorithm may move to pre-header certain computations that may not be executed in the loop. Not only does this risk slowing down the program significantly, it may also cause an error in certain circumstances.

Even if none of the conditions of (2i), (2ii), (2iii) of code motion algorithm are met by an assignment x: =y+z, we can still take the computation y+z outside a loop. Create a new temporary t, and set t: =y+z in the pre-header. Then replace x: =y+z by x: =t in the loop. In many cases we can propagate out the copy statement x: = t.
Maintaining data-flow information after code motion:

The transformations of code motion algorithm do not change ud-chaining information, since by condition (2i), (2ii), and (2iii), all uses of the variable assigned by a moved statement s that were reached by s are still reached by s from its new position. Definitions of variables used by s are either outside L, in which case they reach the pre-header, or they are inside L, in which case by step (3) they were moved to pre-header ahead of s. If the ud-chains are represented by lists of pointers to pointers to statements, we can maintain ud-chains when we move statement s by simply changing the pointer to s when we move it. That is, we create for each statement s pointer ps, which always points to s. We put the pointer on each ud-chain containing s. Then, no matter where we move s, we have only to change ps , regardless of how many ud-chains s is on.

The dominator information is changed slightly by code motion. The pre-header is now the immediate dominator of the header, and the immediate dominator of the pre-header is the node that formerly was the immediate dominator of the header. That is, the pre-header is inserted into the dominator tree as the parent of the header.

Elimination of induction variable:

A variable x is called an induction variable of a loop L if every time the variable x changes values, it is incremented or decremented by some constant. Often, an induction variable is incremented by the same constant each time around the loop, as in a loop headed by for i := 1 to 10. However, our methods deal with variables that are incremented or decremented zero, one, two, or more times as we go around a loop. The number of changes to an induction variable may even differ at different iterations.

A common situation is one in which an induction variable, say i, indexes an array, and some other induction variable, say t, whose value is a linear function of i, is the actual offset used to access the array. Often, the only use made of i is in the test for loop termination. We can then get rid of i by replacing its test by one on t. We shall look for basic induction variables, which are those variables i whose only assignments within loop L are of the form i := i+c or i-c, where c is a constant.

♦ ALGORITHM: Elimination of induction variable

INPUT: A loop L with reaching definition information, loop-invariant computation information and live variable information.

OUTPUT: A revised loop.

METHOD:

✓ Consider each basic induction variable i whose only uses are to compute other induction variables in its family and in conditional branches. Take some j in i's family, preferably one such that c and d in its triple are as simple as possible and modify each test that i appears in to use j instead. We assume in the following tat c is positive. A test of the form 'if i relop x goto B', where x is not an induction variable, is replaced by

r := c*x /* r := x if c is 1. */ r := r+d /* omit if d is 0 */

if j relop r goto B

where, r is a new temporary. The case 'if x relop i goto B'is handled analogously. If there are two induction variables i1 and i2 in the test if i1 relop i2 goto B, then we check if both i1 and i2 can be replaced. The easy case is when we have j1 with triple and j2 with triple, and $c_1=c_2$ and $d_1=d_2$. Then, i1 relop i2 is equivalent to j1 relop j2.

Now, consider each induction variable j for which a statement j: =s was introduced. First check that there can be no assignment to s between the introduced statement j :=s and any use of j. In the usual situation, j is used in the block in which it is defined, simplifying this check; otherwise, reaching definitions information, plus some graph analysis is needed to implement the check. Then replace all uses of j by uses of s and delete statement j:

ONLINE QUESTIONS

UNIT-I

				•			
Questions	opt1	opt2	opt3	opt4	opt 5	opt 6	answer
translates							
level							
to an							
machine level	Compiler	Assemble	Loader	Preprocesso			Assembler
			Lodder	1			Assembler
translates high level							
language in to an							
equivalent low level		Assemble		Preprocesso			
language	Compiler	r	Loader	r			Compiler
File inclusion		Assemble		Dressreeses			Drenzesses
by	Compiler	r	Loader	r			r
performs type checking	Lexical analysis	Semantic analysis	Linear analysis	Syntax analysis			Semantic analysis
Grouping of							
called	String	Stream	Token	Record			Token
	oung	Olican	TOKCH				TOKCH
groups tokens in to							
grammatical phrases	Parser	Scanner	Analyzer	Processor			Parser
Example for Token	Svntax	Character	Symbol	Keyword			Kevword

	1	1	1	1	 1 I
The Idem potent law in regular	R *** – r *	R ** _ r ***	R ** - r *	R *** - r **	R ** - r *
		K =1	K =1	K =1	K =1
breaks up the source program into pieces & creates intermediate code representatio n	Linear phase	Analysis phase	Syntax phase	Synthesis phase	Analysis phase
constructs the target program from intermediate code representatio n	Linear phase	Analysis phase	Syntax phase	Synthesis phase	Synthesis phase
Grouping of tokens into syntactic structure is performed by	Linear analysis	Parser	Scanner	Code optimization	Parser
transforming parse tree in to intermediate language representatio n	Three address code	Code generatio n	Intermediat e code generation	Post fix notation	Intermedia code generation
converts intermediate code in to low level language	Intermediate code generation	Code generatio n	Assembler	Loader	Code

the input of structure editors	Sequence of commands	Sequence of characters	Sequence of tokens	String	Sequence of commands
Pretty printers performs	Printing only	Analyzing and printing	Debugging and printing	Debugging only	Analyzing and printing
Static checker work is	Debugging and printing	Analyzing and printing	Analyzing and debugging	Debugging only	Analyzing and debugging
translates high level language in to an equivalent low level Language	Interpreter	Assemble r	Loader	Preprocesso r	Interpreter
the input of text formatters	Sequence of commands	Stream of characters	Stream of tokens	Lexeme	Stream of characters
Query interpreters translates a predicate contains	Boolean operators only	Relational operators only	Relational and Boolean operators	Arithmetic, Relational and Boolean operators	Relational and Boolean operators
Loader performs	Loading program in to cache memory	Loading program in to main memory	Loading program in to secondary memory	Inking	Loading program in to main memory
is the linking-editor job	Linking preprocessor directives	Linking library functions	Linking machine code	Linking object modules	Linking object modules
Parser generators produce	Scanners	Lexical analyzers	Syntax analyzers	Little languages	Syntax analyzers
uses Scanner generators	Semantic Analyzers	Lexical analyzers	Syntax analyzers	Little languages	Lexical analyzers

Intermediate code generation using tool	Syntax direction engine	Syntax directed trlation engine	Syntax trlation scheme	Syntax directed scheme	Syntax directed trlation engine
Code Optimization phase using tool	Data flow engine	Automatic code generator	Code generator	Code optimizer	Data flow engine
Code Generation phase using tool	Data flow engine	Automatic code generator	Code generator	Code optimizer	Automatic code generator
Lexeme is a	Sequence of characters	Sequence of command s	Set of strings	Pattern	Sequence of characters
Relational operators is a	Lexeme	Pattern	Token	Character	Token
Deleting an extraneous character is a action of					
phase	Lexical	Syntax	Semantic	Synthesis	Lexical
Trposing two adjacent characters is a action of					
phase	Semantic	Syntax	Lexical	Synthesis	Lexical
an error recovery action in Lexical analysis	Inserting a missing character	Function call return	Semicolon missing	Misspelled keyword	Inserting a missing character
Sentinel is an	Foe	Foef	Feof	Eof	Eof

1	1	1	I	1 1	1	1 1
is an one way of implementing lexical analyzer	Using Lex	Using Lexeme	Using Yacc	Using Operating System		Using Lex
The pointer used in buffer pair scheme is	Backward	Forward	Lexeme- End	Lexeme- Start		Forward
is an example of Computer Alphabets	ASC	EBDCIC	ASCI	EBCDIC		EBCDIC
Finite sequence of symbols called	String	Character	Sequence	Group		String
Any set of strings over some fixed alphabet is a	Abstract	Alphabets	Language	Sequence		Language
Set of letters and digits is represented by	LD	LUD	(LU*	(L*		LUD
Set of all four letter strings is represented in a language as	L4	LLLL	L*	(LLLL)*		L4

Set of strings including empty string is represented is a language as	L+	L*	D+	D*	L*
is an representatio n of one or more digits	L+	L*	D+	D*	D+
If X = class , Y = room then XY is	Class	Room	Class Room	Classroom	Classroom
Prefix of Banana is	Ban	Ana	Na	Banana	Ban
is the subsequence of banana	Can	Вааа	Nand	Nanan	Baaa
Suffix of Banana is	Ban	Baaa	Nana	Banana	Nana
Substring of banana is	Nan { sl s is in L or s	Baaa { s s is in L and s is	Aaa { s s is in L nor s is in	Bnn { s s is in L nand s is in	Nan { sl s is in L
LUM is is a notation for	is in M }	in M }	M }	M }	orsisinM}
Regular Expression	Letter(letterdigit) +	(letterdigit) +	Digit (letter digit) *	Letter(letter digit) *	Letter(letter digit) *
Definition of LM is	{ st s is in L or s is in M }	{ st s is in L and t is in M }	{ st s is in L or t is in M }	{ st s is in L and s is in M }	{ st s is in L and t is in M }
L* is an representatio n of	Negative closure	Positive closure	Kleene closure	Line closure	Kleene closure

i i	1	1	1	1 1	· · ·	1
Positive closure of L is written as	L+	L*	D+	D*		L+
(r) (s) is a regular expression denoting	$L(\mathbf{r}) \mid L(\mathbf{s})$	L(r) L(s)	L(r)* L(s)	L(r) U L(s)		L(r) UL(s)
(r) (s) is a regular expression denoting	L(r) L(s)	L(r)L(s)	L(r)* L(s)	L(r) U L(s)		L(r)L(s)
(r) * is a regular expression denoting	(L(r))*	Lr*	L(r)	L*(r)		(L(r))*
The regular expression a* denotes	{a}	{?,a}	{ a , aa , aaa,}	{?,a,aa, aaa,}		{?,a,aa, aaa,}
Identifier is represented in character class as	[A-Z] [A-Z0- 9]*	[A-Za-z] [A-Za-z0- 9]*	[A-Za-z] [A-Za-z]*	[A-Z] [A- Za-z0-9]*		[A-Za-z] [A- Za-z0-9]*
cannot be described by a regular expression	???{ wcw w is a string of a's and b's }	{ w w is a string of a's and b's }	??{ w* w is a string of a's and b's }	??{ w+ w is a string of a's and b's }		???{ wcw w is a string of a's and b's }
an associative property of a regular expression	R(sit) = (ris)t	R (s t) =	R (s t) = (r s) t	(sit)r= t(ris)		R (s t) = (r s) t

A						
replacement						
according to						
a production						
is						
called		Productio				
_	Reduction	n	Derivation	Parse tree		Derivation

UNIT-II

Questions	opt1	opt2	opt3	opt4	opt 5	opt 6	answer
Process of replacing a string by an NT according to a production	Reduction	Productio	Derivation	Parse tree			Reduction
Which of the following is not a true statement as a derivation tree?	all the leaf nodes are terminals	root node is start symbol	interior node is terminal	interior node is non - terminal			interior node is terminal
Method of converting regular expression to a recognizer is	Lexical analyzer	Finite automata	Lex	Үасс			Finite automata
Demerits of using transition table is	tough to implement	slower	takes up less space	takes up lot of space			takes up lot of space
In which Finite Automata no states has an ? - transition	NFA	DFA	NDFA	DNFA			DFA

has					
labeled "a"					
leaving state S		DFA	NDFA	DINFA	DFA
Tool used to design the lexical analysis is	Yacc	Lexeme	Lex	Yaccer	Lex
language is used	0	Lev	Veee		
program	C	Lex	Yacc	Linux	Lex
is a c program produced by Lex compiler	Lex.yy.c	Lex.c	a.out	tokens	Lex.yy.c
is the output produced by C compiler in	Lex vy c	Lex c	aout	tokens	a out
in the	201.99.0				
input taken by C	Lex vv c	lexic	a out	tokens	Lex vy c
is one of the field in the Lex specification	Transition	Translatio n rules	Definition	main function	Translation
is not a part of Lex specification	manifest constants	auxiliary procedure s	declaration s	main function	main function

Í	1	I	I	1	1 1	I
Build parse trees from root to leaves is called 	Top down parsing	Bottom up parsing	LR parsing	Root leaf parsing		Top down parsing
scanned from	root to right	loft to right	top to loft	right to loft		loft to right
A ? Aa is called as	Left factoring	Ambiguou s	Left Recursion	Left refactoring		Left Recursion
Elimination of left recursion for the A ? Aa / ß are	A?A'a,A' ?ߑa/?	A ? ßA' , A' ? aA' / ?	A ? Aa , A' ? A'a / ?	A ? Aa / ß , A ? ߑa / ?		A ? ßA' , A' ? aA' / ?
Elimination of left recursion for the E ? E+T / T is	E? TE' , E' ? +TE' / ?	E? T'E , E' ? +E / ?	E? T'E' , E' ? +T'E/ ?	E ? TE , E' ? +T'E' / ?:		E? TE' , E' ? +TE' / ?
Elimination of left recursion for the T?T*F/F is	T ? FT , T' ? *F'T' / ?	T? F'T , T' ? *E / ?	T? F'T' , T' ? *F'T/ ?	T? FT' , T' ? *FT' / ?\		T? FT' , T' ? *FT' / ?\
Process of factoring out the common prefixes is	Left	Ambiguou s	Left Recursion	Left		Left factoring
Elimination of left factoring for the A ? a ß1 / a ß2 are	A ? aßA' , A' ? ß1 / ß2	A ? a'A , A' ? aß2 / ß1	A ? aA' , A' ? ß1 / ß2	A ? a'A' , A' ? a ß1 / a ß2		A?aA', A' ?ß1/ß2

is the left factored grammar for S? iEtS iEtSeS a	S? iEtSS' a , S'? eS ?	S? iEtS' a , S'? iEtSS' ?	S? iEtSS' a , S'? iEtSeS ?	S? iEtSeS a , S'? eS ?	S? iEtSS' a , S'? eS ?
Top down parsing is creating the tree in order	post	pre	in	reverse polish	pre
Top down parsing is used to find	post order	Right most derivation	in order	Left most derivation	Left most derivation
left recursive grammar can cause top down parser to go	elimination	error condition	infinite loop	finite loop	infinite loop
Difficulty of top down parsing is	Forward loop	For tracking	Backward loop	Back tracking	Back tracking
is the Top down parsing technique Demerits of	Recursive descent	Recursion descent	Predicate logic	periodic parsing	Recursive descent
recursive descent parsing is	Backtrackin g	Left factoring	Recursive	Recursion	Recursive
Predictive parsing consists of	input , parsing program, parsing table , output	input , stack, parsing table , output	input , parsing program, stack , output	input , parsing program,stac k, parsing table	input , stack, parsing table , output

In predictive parsing program let X be the top stack symbol and a be the next input symbol, if X is terminal and		push a on				
if X= a then	push X on	to the	pop a from	pop X from		pop X from
In predictive parsing program let X be the top stack symbol and a be the next input symbol, if X is Non - terminal and if M[X,a] =X?Y1, Y2, Yk then	push X from the stack , pop Yk , Yk-1 , Y1 on to the stack	push X from the stack , pop Y1 , Y2 , Yk on to the stack	pop X from the stack , push Yk , Yk-1 , Y1 on to the stack	pop X from the stack , push Y1 , Y2 , Yk on to the stack		pop X from the stack , push Yk , Yk-1 , Y1 on to the stack
In predictive parsing if X is a terminal , then FIRST(X) is	{?}	{X}	terminal	nonterminal		{X}
Two functions used in predictive parsing is	FIRST , LAST	FIRST , FOLLOW	FOLLOW , LAST	FIRST, PREDICT		FIRST , FOLLOW

In predictive parsing if X ? ?, then FIRST(X) is	{?}	{ X }	terminal	nonterminal	{?}
In predictive parsing A? aBß , then everything in is in	FIRST(a) , FOLLOW (B)	FIRST(ß) , FOLLOW (ß)	FIRST(ß) , FOLLOW (B)	FIRST(a) , FOLLOW (ß)	FIRST(ß) , FOLLOW (B)
In predictive parsing A? aB , then everything in is in	FOLLOW (a), FOLLOW (B)	FOLLOW (A) , FOLLOW (B)	FOLLOW (a) , FOLLOW (A)	FOLLOW (A) , FOLLOW (a)	FOLLOW (A) , FOLLOW (B)
In predictive parsing A? aBß, where FIRST(ß) contains ?, then everything in is in	FOLLOW (A) , FOLLOW (B)	FOLLOW (a), FOLLOW (B)	FOLLOW (a) , FOLLOW (A)	FOLLOW (A) , FOLLOW (a)	FOLLOW (A) , FOLLOW (B)
is included in FOLLOW function set of predictive parsing	%	#	\$	@	\$
If a grammar S ? (L) a,then FIRST(S) is	{(,a}	{?}	{L,a}	{\$}	{(,a}
Elimination of left recursion for the L ? L,S S is	L? S,L' , L'? LL' / ?	L? S',L , L'? S,L / ?	L? SL' , L'? SL' / ?	L? S',L' , L'? LL / ?	L? SL' , L'? SL' / ?

LL(1) grammar has the property of	Ambiguity	No ambiguity	No recursion	No backtracking	No ambiguity
LL(1) grammar has property	Ambiguity	Left recursive	No recursion	No Left recursive	No Left recursive
In LL(1) grammar the first L stands for	left to right scanning	left most derivation	left to right derivation	left subtree	left to right scanning
In LL(1) grammar the second L stands for	left to right scanning	left most derivation	left to right derivation	left subtree	left most derivation
In LL(1) grammar the 1 stands for	one output symbol	one time scanning	one sub tree	one input symbol	one input symbol
In predictive parsing , If X = a = \$	error report	successful completio n	advances pointer	pop X off the stack	successful completion
In predictive parsing , If X = a ? \$	error report	successful completio n	advances pointer & pop X off the stack	pop X off the stack & parser halts	advances pointer & pop X off the stack
In predictive parsing , If X = \$	Stack is empty	Stack is full	error report	pop X from Stack	Stack is empty
Ambiguity means	produces more than two parse tree	produces null parse tree	produces more than one parse tree	produces finite parse tree	produces more than one parse tree
Ambiguous means	produces more than two derivation	produces only left most derivation	produces more than one derivation	produces only right most derivation	produces more than one derivation
The output of Lex compiler is	Transition table	Transition diagram	Lex specificatio n	action	Transition table

The input of Lex compiler is	Transition table	Transition diagram	Lex specificatio n	action	Lex specificatio n
The lexical analyzer design has	Lex compiler	Transition table	Lex specificatio n	Transition diagram	Transition table
Recognizer is a	tool	program	string	grammar	program
NFA representation in a directed graph is called	NFA graph	NFA direction graph	Transition edge	Transition graph	Transition graph
is an easiest implementation of NFA in a computer.	Transition diagram	Transition	Transition graph	Finite automata	Transition table
is an input for the syntax analysis	token	source program	parse tree	syntax	token
is an output of the parser	expression	token	parse tree	intermediate representatio n	parse tree
Parser is a	Lexical analyzer	Front end tool	Scanner	Back end tool	Front end tool
is syntax error	misspelled identifier	misspelle d keyword	misspelled operator	unbalanced parenthesis	unbalanced parenthesis

UNIT-III

					opt	opt	
Questions	opt1	opt2	opt3	opt4	5	6	answer

is one of the method of parsing	left to right	Top down parsing	Bottom down parsing	Top up parsing	Top down parsing
constructs parse tree from leaf node to root	Bottom up	Top down parsing	Bottom down parsing	Top up parsing	Bottom up
one of the goal of error handler in parser is	avoid common error	slow down the process	report the presence of error	avoid specific error	report the presence of error
In which error recovery strategy parser discards one symbol at a time ?	Panic mode	phrase level	error production s	global corrections	Panic mode
In which error recovery strategy parser makes local corrections 2	Panic	phrase	error production	global	phrase
In which error recovery strategy parser generate error diagnostics ?	Panic mode	phrase level	error production s	global corrections	error productions
In which error recovery strategy parser using algorithmic approach ?	Panic mode	phrase level	error production s	global	global

The syntactic structure of a programming language is described by	Context free grammar	CNF	CCNF	Regular language	Context free grammar
Context free grammar consists of	T , NT, \$, P	T, NT, S , Production	Terminal , token , Non terminal	Terminal , Token , production	T, NT, S , Production
The keyword else is a	Terminal	Non terminal	Start symbol	Production	Terminal
are syntactic variable	Production	Non terminal	Terminal	string	Non terminal
Non terminal denotes	sets of characters	sets of grammar	sets of production	sets of strings	sets of strings
Start symbol is a	Production	Non terminal	Terminal	string	Non terminal
Upper case letters		Non	—		Non
are	Production	terminal	Terminal	string	terminal
Lower case letters are	Production	Non terminal	Terminal	string	Terminal
Punctuation symbols are	Production	Non terminal	Terminal	string	Terminal
Boldface strings are	Terminal	Non terminal	Start symbol	Production	Terminal
operator symbols are	Terminal	Non terminal	Start symbol	Production	Terminal

Lower case italic names are	Production	Non terminal	Terminal	string	Non terminal
The left side of the first production is called	String	Terminal	Non terminal	Start symbol	Start symbol
Lower case Greek letters represents	Grammar symbols	Terminals	Non terminals	Start symbol	Grammar symbols
sequence of replacements is called	Reduction	Derivation	parse tree	sentence	Derivation
Graphical representation of a derivation is a	parse tree	syntactic tree	syntax graph	pattern	parse tree
substring that matches the right side of a production is called	Handle pruning	Pattern	Handle	parsing	Handle
Right most derivation in reverse is called	Handle	Pattern	Handle	parsing	Handle

1	1	1		1	 i
In action the next input symbol is shifted on to the top of the stack	accept	reduce	shift	error	shift
In action the parser announces the successful completion of	accent	reduce	chift	error	accent
parsing	accept	reduce	Shiit	error	accept
In action parser discovers the syntax error	accept	reduce	shift	error	error
In action parser replacing the handle with the non terminal	accept	reduce	shift	error	reduce
The set of prefixes that appear on the stack are called	prefixes	viable prefixes	reduce conflict	viable suffixes	viable prefixes
Grammar that has no production right side is e is called	operator grammar	shift reduce grammar	operator precedenc e grammar	precedenc e grammar	operator grammar

	I	I	I	1		
Grammar that has no production right side is two adjacent non terminals is called	operator precedenc e grammar	shift reduce grammar	operator grammar	precedenc e grammar		operator grammar
In precedence relation a < b means	a has different precedenc e to b	a has same precedence as b	a takes precedenc e over b	a yields precedenc e to b		a yields precedence to b
In precedence relation a > b means	a has different precedenc e to b	a has same precedence as b	a takes precedenc e over b	a yields precedenc e to b		a takes precedence over b
In precedence relation a = b means	a has different precedenc e to b	a has same precedence as b	a takes precedenc e over b	a yields precedenc e to b		a has same precedence as b
In precedence relation * and + has the precedence of						
In precedence relation \$ and id	<u>`>+</u>	+>*	+ = +	<u>`</u> = *		^>+
of	\$ > id	\$ =id	\$ <id< td=""><td>\$? id</td><td></td><td>\$ <id< td=""></id<></td></id<>	\$? id		\$ <id< td=""></id<>

		i .	1		
In operator precedence parsing if a < b then	pop a from the stack	push a on to the stack	pop b from the stack	push b on to the stack	push b on to the stack
In operator precedence parsing if a = b then	pop a from the stack	push b on to the stack	pop b from the stack	push a on to the stack	push b on to the stack
In operator precedence parsing if a > b then	pop the stack	push the stack	pop b from the stack	push a on to the stack	pop the stack
In LR(k) parsing , the L stands for	left to right scanning	left most derivation	left to right derivation	left subtree	left to right scanning
In LR(k) parsing , the R stands for	left to right scanning	Right most derivation	left to right derivation	Right most derivation in reverse	Right most derivation in reverse
In LR(k) parsing , the k stands for	parsing symbol	number of input symbols	number of characters	look ahead symbol	number of input symbols
LR technique is easy to implement	SLR	canonical LR	LALR	Look ahead LR	SLR
LR technique is most powerful	SLR	canonical LR	LALR	Look ahead LR	canonical LR
LR technique is most expensive	SLR	canonical LR	LALR	Look ahead LR	canonical LR

Ì	l.	1	1	1	1 1	Í.
LR technique is least powerful	SLR	canonical LR	LALR	Look ahead LR		SLR
S' ? .S is included in items	closure	non kernel	kernel	non closure		kernel
The functions performed in LR parsing are	action and shift	action and goto	action and error	goto and shift		action and goto
Action function involved with	Terminals	Non terminals	start symbol	production		Terminals
Goto function involved with	Terminals	Non terminals	Start symbol	Production		Non terminals
LALR is abbreviated from	Left and Right LR	Look ahead Simple LR	Look ahead LR	Left to Right simple LR		Look ahead LR
is a combination of terminals and non terminals	?productio n	token	regular expression	regular definition		?production
is an one of the bottom up parsing technique	Operator parsing	Shift reduce	Recursive descent parsing	Predictive		Shift reduce

translates intermediate representation in to an equivalent low level language	Analyzer	Front end	Back end	Synthesize r	Back end
The input for the intermediate code generator is	Optimized code	Intermediat e Code	Token	Meaningful expression	Meaningful expression
The output for the intermediate code generator is	Optimized code	Intermediat e Code	Token	Meaningful expression	Intermediat e Code
In tree the operators represented in the interior node	Postfix	Parse Tree	Syntax Tree	Prefix	Syntax Tree
is a linearized representation of a syntax tree	Postfix Notation	Parse Tree	Two address code	Prefix notation	Postfix Notation
Postfix notation for the statement a = b * - c is	* uminus a b c	uminus * a b c	a b c * uminus	a b c uminus *	a b c uminus *

UNIT-IV

Questions	opt1	opt2	opt3	opt4	op t5	op t6	answer
Semantic rule to produce syntax tree for the production E ? id is	E.nptr = mkleaf (id.place)	E.nptr = mknode (id.place)	E.nptr = mkleaf (id , id.place)	E.nptr = mknode (id , id.place)			E.nptr = mkleaf (id , id.place)
Semantic rule to produce syntax tree for the production E ? - E1 is	E.nptr = mknode ('uminus' , E1.nptr)	E.nptr = mkunode (ʻuminus' , E1.nptr)	E.nptr = mkleaf (ʻuminus' , E1.nptr)	E.nptr = mkuleaf (ʻuminus' , E1.nptr)			E.nptr = mkunode ('uminus' , E1.nptr)
Semantic rule to produce syntax tree for the production E ? E1 + E2 is	E.nptr = mknode ('+' . E1.nptr)	E.nptr = mkpnode ('+' . E1.nptr)	E.nptr = mkpnode ('+' , E1.nptr , E2.nptr)	E.nptr = mknode ('+' , E1.nptr , E2.nptr)			E.nptr = mknode ('+', E1.nptr, E2.nptr)
is an general form of three address code representation	x = y op z	x = z op	op x = z op	op x = op y z			x = y op z
is an one type of three address code statements	if x y goto L	if x relop y goto L	goto L if x y	goto L if x relop y			if x relop y goto L

		_			
The name that will hold the value of E is called	E.place	E.code	place.E	code.E	E.place
The sequence of three address statements evaluating E is called	E.place	E.code	place.E	code.E	E.code
Record structure with four fields is called	Three address code	Quadruples	Triples	Indirect triples	Quadruples
Three fields of Record structure is called	Three address code	Quadruples	Triples	Indirect triples	Triples
emit is used to	emit 3address statements to an output file	emit assignments to an output file	emit terminals to an output file	emit non terminals to an output file	emit 3address statements to an output file
Translation of E?(E1) is	E1.place = E.place	E.place=E1.pl ace	E.place = (E1.place)	(E1.place)= E.place	E.place=E1. place
Translation of E? id is	E1.place = E.place	E.place=E1.pl ace	E.place = id.place	E1.place= id.place	E.place = id.place
the 3 fields of triples	arg1,arg2,a rg3	arg1,arg2,res ult	arg1,op,resul t	arg1,arg2,o p	arg1,arg2,op
are the 4 fields of Quadruples	arg1,arg2,a rg3 , arg4	arg1,arg2,arg 3,result	arg1,,arg2,o p,result	arg1,arg2,a rg3,op	arg1,,arg2,o p,result

Listing pointers to triples is called	Indirect triples	triples	Quadruples	Indirect ruples	Indirect triples
offset represents	relative address	three address	location	address	relative address
E1.type=integer ,E2.type = integer , E.type	integer	real	inttoreal	float	integer
	lineger				linegoi
E1.type=integer ,E2.type = real , E.type is	integer	real	inttoreal	float	real
E1.type=real,E2 .type = real , E.type is	integer	real	inttoreal	float	real
E1.type=real,E2 .type = integer, E.type is	integer	real	inttoreal	float	real
E?E1 or E2 represents	E1.place = E1.place * E2.place	E.code= E1.place or E2.place	E1.code= E1.place and E2.place	E.place= E1.place or E2.place	E.place= E1.place or E2.place
E2 true					
represents	E code-0	E place – 0	E place – 1	E code-1	E place – 1
	L.0006-0			L.0006-1	
E? false					
	E.code=0	E.place = 0	E.place = 1	E.code=1	E.place = 0
E?not E1 represents	E1.place = E1.place not E2.place	E.code= E1.place not E2.place	E1.code= not E1.place	E.place= not E1.place	E.place= not E1.place

		1			 1
is a one semantic rule of S?if E then S1	E.false=ne wlabel	E.true=newla bel	S.next=S1.n ext	S.code=E.p lace B	E.true=newl abel
is a one of the semantic rule of S?if E then S1	E.false=ne wlabel	S1.next=S.ne xt	S.next=S1.n ext	S.code=E.p lace	S1.next=S.n ext
is a one of the semantic rule of S? while E do S1	E.false=ne wlabel	S1.next=S.ne xt	S.next=S1.n ext	E.false=S.n ext	E.false=S.ne xt
is a one of the semantic rule of S? while E do S1	E.false=ne wlabel	S1.next=S.ne xt	E.true=newl abel	E.false=S1. next	E.true=newl abel
is a one of the semantic rule of S?if E then S1 else S2	E.false=S1.	S1.next=S.ne	S.next=S1.n	S.code=E.c	S.code=E.co
is a one of the semantic rule of S?if E then S1	E.false=E.tr	S1.next=S.ne	S.next=S1.n	S.code=E.p	S1.next=S.n
else S2 is a one of the semantic rule of S? while E do S1	ue E.false=S.n ext	S1.next=S2.n	ext E.true=E.fals e	E.false=S1.	E.false=S.ne

is a one of the semantic rule of E?E1 or E2	E1.true=E.t rue	E1.false=E.fal	E2.true=E.fal se	E2.fasle=E 1.true	E1.true=E.tr ue
is a one of the semantic rule of E?E1 or E2	E1.true=E.f alse	E1.false=E.fal	E2.true=E.tr ue	E2.fasle=E 1.true	E2.true=E.tr ue
one of the semantic rule of E?E1 or E2	E1.true=E.f	E1.false=E.fal	E2.true=E1.t	E2.fasle=E. false	E2.fasle=E.f alse
is a one of the semantic rule of F2F1_and F2	E1.true=E.f	E1.false=E.fal	E2.true=E1.t	E2.fasle=E	E1.false=E.f
is a one of the semantic rule of	E1.true=E.f	E1.false=E.tr	E2.true=E1.t	E2.fasle=E.	E2.fasle=E.f
is a one of the semantic rule of	E1.true=E.f	E1.false=E2.f	E2.true=E.tr	E2.fasle=E	E2.true=E.tr
is a one of the semantic rule of E?E1 and E2	E1.true=ne wlabel	E1.false=E2.f	E2.true=E1.t	E2.fasle=E. true	E1.true=new label

is a one of the semantic rule of E? not E1	E.true= E1.false	E1.true = E.false	E.false= E2.true	E.code=E.p lace	E1.true = E.false
is a one of the semantic rule of E? not E1	E1.true= E.false	E1.true = E2.false	E.false= E2.true	E.code=E.p lace	E1.true= E.false
The use of makelist(i) is	creates a new list containing quadruples	creates a new list containing only i	creates a new list by inserting i	creates a new list pointed to p1	creates a new list containing only i
is the use of merge(p1,p2)	concatenat es the lists pointed by p1 and p2	merge the list containg only i	merge the list pointed by p1	merge the list containing quadruples	concatenate s the lists pointed by p1 and p2
The use of backpatch(p,i) is	concatenat es the lists pointed by p1 and p2	merge the list containing only i	inserts I as the target label	merge the list containing quadruples	inserts I as the target label
is a one of the semantic rule of E? not E1	E.truelist= E1.falselist	E1.truelist = E2.falselist	E.falselist= E2.truelist	E.code=E.p lace	E.truelist= E1.falselist
is a one of the semantic rule of	E1.truelist=	E1.truelist =	E.falselist=	E.code=E.p	E.falselist=
one of the semantic rule of E? (E1)	E.truelist= E1.truelist	E1.truelist = E2.falselist	E.falselist= E2.truelist	E.code=E.p lace	E.truelist= E1.truelist

	1	1	1	I I	1	I
is a one of the semantic rule of E? (E1)	E2.truelist= E1.truelist	E1.truelist = E2.falselist	E.falselist= E1.falselist	E.code=E.p lace		E.falselist= E1.falselist
translation of s?begin L end is	S.list = L.nextlist	S.nextlist=L.n extlist	L.List = S.Listnext	L.nextlist=S .list		S.nextlist=L. nextlist
The translation of Elist?Elist,E	append E.code to the end of the queue	append E.place to the beginning of the queue	append E.code to the beginning of the queue	append E.place to the end of queue		append E.place to the end of queue
The translation of Elist? E is	Initialize E.code to the end of the queue	Initialize E.place to the beginning of the queue	Initialize E.code to the beginning of the queue	Initialize queue to contain only E.place		Initialize queue to contain only E.place
The translation of M? ? with quadruple is	M.quad = nextQuad	M.nextquad = nextQuad	M.next = M.Quad	M.quad = M.nextQua d		M.quad = nextQuad
The translation of N? ? with quadruple is	N.nextlist = makelist(ne xtqua	N.nextlist = list(qua	N.nextlist = makelist(M.q ua	N.nextlist = make (nextqua		N.nextlist = makelist(nex tqua
is an input of code generation phase	optimized code	target code	source program	object code		optimized code
The output for the code generator is	Optimized code	Intermediate Code	Token	Assembly language		Assembly language

The transformation performed only within a basic block is						
called	local	global	preserve	optimization		local
Eliminating the same sub expressions is called	common elimination	common sub expression elimination	common expression deletion	common sub expression deletion		common sub expression elimination
is a transformations of copy statements	copy propagatio n	copy transformatio n	copy for long	copy elimination		copy propagation
Useless code transformation is called	usecode elimination	dead elimniation	useless code elimination	Deadcode elimination		Deadcode elimination
Using the constant and deducing during compile time is called	dead code elimination	copy propagation	constant folding	constant propagation		constant folding
Optimizing inner loops named as	Loop transformat ion	Loop	Deadcode elimination	copy propagation		Loop optimization

UNIT-V

Questions	opt1	opt2	opt3	opt4	opt 5	opt 6	answer
Decreasing				Code			
the amount	Induction	Reduction	Loop motion	motion			Code motion

of code in a inner loop is called as					
is an one way of loop optimization	Induction variable elimination	Copy propagati on	Deadcode elimination	constant folding	Induction variable elimination
an loop optimization technique	Reduction variable elimination	Reduction in strength	Deadcode elimination	constant folding	Reduction in strength
Expansion for DAG is	Directed Acyclic Graph	Directed Action Graph	Direction Asymmetric Graph	Direction Action Graph	Directed Acyclic Graph
The Algebraic transformati on includes	Algebraic Deduction	Algebraic Identities	constant folding	reduction in strength	Algebraic Identities
The output for the code generation phase is	Optimized code	Intermedi ate Code	Machine level language	Token	Machine level language
the input of code generation phase	Optimized code	Intermedi ate Code	Machine level language	Token	Intermediate Code
The use of symbol table is	to determine the run time addresses of the data objects	to determine the run time value of the data	to determine the compile time value of the data	to determine the compile time addresse s of the data objects	to determine the run time addresses of the data objects
is a linear representati ons of intermediate code	Prefix notation	Infix notation	Postfix notation	RP notation	Postfix notation
an representati on of three address code	Quadruples	Indirect Quadrupl es_	Postfix notation	Linear	Quadruples
the virtual machine representati on	Sequence of commands	Stack machine code	Machine code	Stack code	Stack machine code
is a Graphical	DEG tree	Parsing	Syntax trees	Linear tree	Syntax trees

representati on of three address code					
is a Graphical representati on of intermediate code	DAG	Parsing	Semantic tree	Linear tree	DAG
is the output form of a target program	Intermediate code	linking library functions	linking machine code	Absolute machine code	Absolute machine code
is an one of the output form of a target program	Intermediate code	linking library functions	Re locatable machine code	Absolute intermedi ate code	Re locatable machine code
Semantic checking done in	Intermediate code generator	Lexical analyzers	Syntax analyzers	Code generator	Code generator
Mapping names to addresses of data objects is done by	intermediate code generator	code generator	code optimizer	Lexical analysis	code generator
Deducing the number of jumping labels is done by	Backpatching	Quadrupl	Triple	Indirect Triple	Backpatching
The speed is increased based on instruction selection by using	Assignment	Machine	Structure	Register	Machine
During Register Allocation	we select variables reside in the register	we pick specific register that a variable reside in	choose register pairs	Allocate constants to register	we select variables reside in the register
During Register assignment	We select variables reside in the register	We pick specific register that a variable reside in	Choose registers pairs	Allocate constants to register	We pick specific register that a variable reside in

SRDA stands for	Shift Right Double Arithmetic	Shift Round Direct Arithmetic	Scan Right Double Arithmetic	Scan Round Direct Arithmetic	Shift Right Double Arithmetic	
used to improve the efficiency	Choice of Run	Syntax	Choice of Evaluation order	Semantic	Choice of Evaluation order	
an two address instruction form	op source,destina tion	source op destinatio n	source,destina tion op	destinatio n source,op	op source,destina tion	
ADD is an	ADD to register	ADD destinatio n to memory	ADD to memory	ADD source to destinatio n	ADD source to destination	
SUB is an	SUB to register	SUB destinatio n to memory	SUB to memory	SUB source from destinatio n	SUB source from destination	
For absolute mode the added cost is	1	0	2	3	1	
For Register mode the added cost is	1	0	2	3	0	
For Indexed mode the added cost is	1	0	2	3	1	
For Indirect indexed mode the added cost	1	0	2	3	1	
The form for absolute mode is		*R	R	M	M	
The form for Register mode is						
The form for Indexed mode is	<u>c(R)</u>	*R	R	M	R	
The form for Indirect Register	c(R)	*R	R	M	c(R)	
mode is	c(R)	*R	R	М	*R	
The form for Indirect	*c(R)	*R	R	М	*c(R)	
Indexed mode is						
---------------------------------------------------	--------------------------------------	--------------------------------------------------------------------	----------------------------------	---------------------------------------------------	--	--------------------------------------------------------------
The address of Register mode is	o(D)	*D	P	D.4		D
The address		ĸ	ĸ	IVI		ĸ
of Indexed						
mode is	c + contents of					c + contents of
	(R)	R	contents of R	М		(R)
The address of Indirect Register mode is	c + contents of	R	contents of (R)	М		contents of (R)
The address						
of Indirect Indexed mode is	contents (c + contents of (R))	R	contents of R	М		contents (c + contents of (R))
The cost of						
MOV R0,R1	1	0	2	2		1
The cost of	1	0	2	5		I
MOV R5.M						
is	1	0	2	3		2
The cost of ADD #1,R3						
is	1	0	2	3		2
The cost of						
500 4(R0), *12(R1) is						
12(1(1)13	1	0	2	3		3
The cost of						
MOV b,a						
and ADD c,a						
is	1	0	6	3		6
The cost of ADD R2,R1 and MOV R1.a is						
	1	0	2	3		3
				to		
The getreg		to		determine		
denotes	to determine	determine	to determine	the		to determine
	the location L	the value	the Register	memory		the location L
The function of register descriptor is	to keep track of the location	to keep track of what is currently in each register	to keep track of the register	to keep track of the descriptor value		to keep track of what is currently in each register

The function of register descriptor is	to keep track of the location	to keep track of what is currently in each register	to keep track of the register	to keep track of the descriptor value	to keep track of the location
For the statement t = $a - b$, the value of Address Descriptor is	R0 contains t	u in R0	t in R0	R0 contains u	t in R0