

UNIT-3

SYNTAX DIRECTED TRANSLATION

In syntax directed translation, along with the grammar we associate some informal notations and these notations are called as semantic rules.

So, we can say that

1. Grammar + semantic rule = SDT (syntax directed translation)

- In syntax directed translation, every non-terminal can get one or more than one attribute or sometimes 0 attribute depending on the type of the attribute. The value of these attributes is evaluated by the semantic rules associated with the production rule.
- In the semantic rule, attribute is VAL and an attribute may hold anything like a string, a number, a memory location and a complex record
- In Syntax directed translation, whenever a construct encounters in the programming language then it is translated according to the semantic rules define in that particular programming language.

Example

Production	Semantic Rules
$E \rightarrow E + T$	$E.val := E.val + T.val$
$E \rightarrow T$	$E.val := T.val$
$T \rightarrow T * F$	$T.val := T.val * F.val$
$T \rightarrow F$	$T.val := F.val$
$F \rightarrow (F)$	$F.val := F.val$
$F \rightarrow num$	$F.val := num.lexval$

E.val is one of the attributes of E.

Syntax directed translation scheme

- The Syntax directed translation scheme is a context -free grammar.
- The syntax directed translation scheme is used to evaluate the order of semantic rules.
- In translation scheme, the semantic rules are embedded within the right side of the productions.
- The position at which an action is to be executed is shown by enclosed between braces. It is written within the right side of the production.

Example

Production	Semantic Rules
$S \rightarrow E \$$	{ printE.VAL }
$E \rightarrow E + E$	{ E.VAL := E.VAL + E.VAL }
$E \rightarrow E * E$	{ E.VAL := E.VAL * E.VAL }
$E \rightarrow (E)$	{ E.VAL := E.VAL }
$E \rightarrow I$	{ E.VAL := I.VAL }
$I \rightarrow I digit$	{ I.VAL := 10 * I.VAL + LEXVAL }
$I \rightarrow digit$	{ I.VAL := LEXVAL }

Implementation of Syntax directed translation

Syntax direct translation is implemented by constructing a parse tree and performing the actions in a left to right depth first order.

SDT is implementing by parse the input and produce a parse tree as a result.

Example

Production	Semantic Rules
$S \rightarrow E \$$	{ printE.VAL }
$E \rightarrow E + E$	{ E.VAL := E.VAL + E.VAL }
$E \rightarrow E * E$	{ E.VAL := E.VAL * E.VAL }
$E \rightarrow (E)$	{ E.VAL := E.VAL }
$E \rightarrow I$	{ E.VAL := I.VAL }
$I \rightarrow I \text{ digit}$	{ I.VAL := 10 * I.VAL + LEXVAL }
$I \rightarrow \text{digit}$	{ I.VAL := LEXVAL }

Parse tree for SDT:

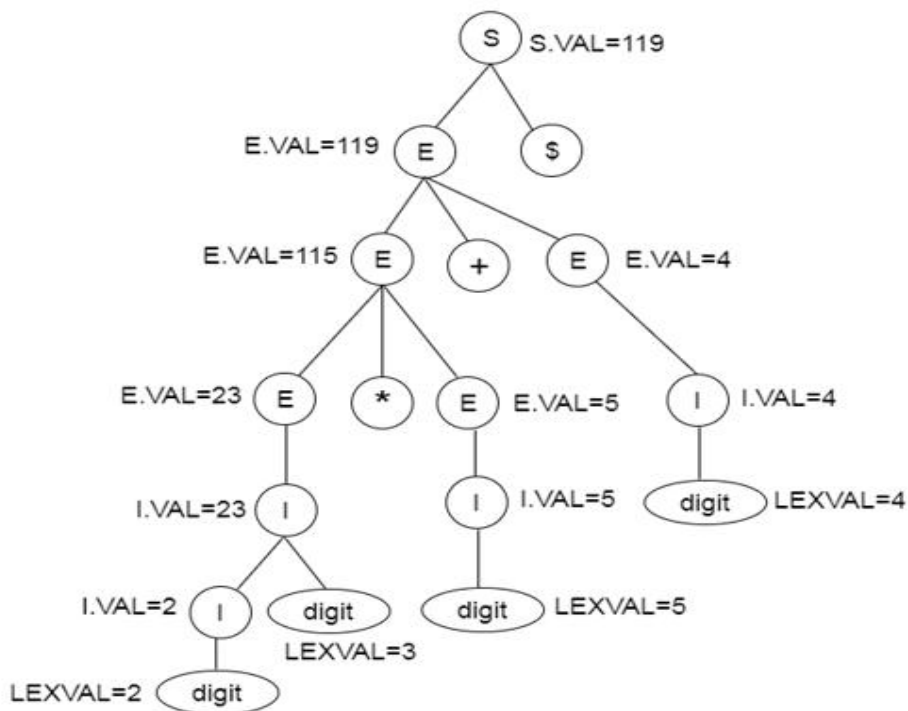


Fig: Parse tree

Intermediate code

Intermediate code is used to translate the source code into the machine code. Intermediate code lies between the high-level language and the machine language.



Fig: Position of intermediate code generator

- If the compiler directly translates source code into the machine code without generating intermediate code then a full native compiler is required for each new machine.
- The intermediate code keeps the analysis portion same for all the compilers that's why it doesn't need a full compiler for every unique machine.
- Intermediate code generator receives input from its predecessor phase and semantic analyzer phase. It takes input in the form of an annotated syntax tree.
- Using the intermediate code, the second phase of the compiler synthesis phase is changed according to the target machine.

Intermediate representation

Intermediate code can be represented in two ways:

1. High Level intermediate code:

High level intermediate code can be represented as source code. To enhance performance of source code, we can easily apply code modification. But to optimize the target machine, it is less preferred.

2. Low Level intermediate code

Low level intermediate code is close to the target machine, which makes it suitable for register and memory allocation etc. it is used for machine-dependent optimizations.

Postfix Notation

- Postfix notation is the useful form of intermediate code if the given language is expressions.
- Postfix notation is also called as 'suffix notation' and 'reverse polish'.
- Postfix notation is a linear representation of a syntax tree.
- In the postfix notation, any expression can be written unambiguously without parentheses.
- The ordinary (infix) way of writing the sum of x and y is with operator in the middle: $x * y$. But in the postfix notation, we place the operator at the right end as $xy *$.
- In postfix notation, the operator follows the operand.

Example

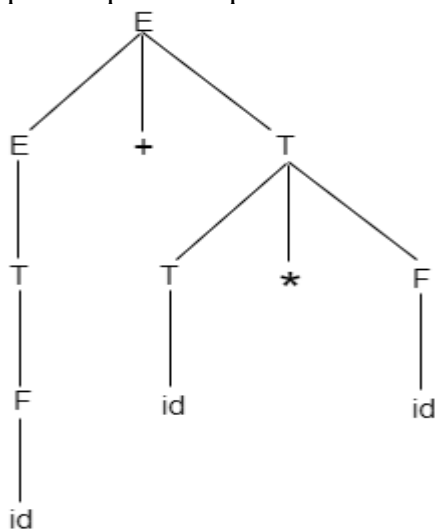
Production

1. $E \rightarrow E1 \text{ op } E2$
2. $E \rightarrow (E1)$
3. $E \rightarrow \text{id}$

Semantic Rule	Program fragment
$E.\text{code} = E1.\text{code} \parallel E2.\text{code} \parallel \text{op}$	print op
$E.\text{code} = E1.\text{code}$	
$E.\text{code} = \text{id}$	print id

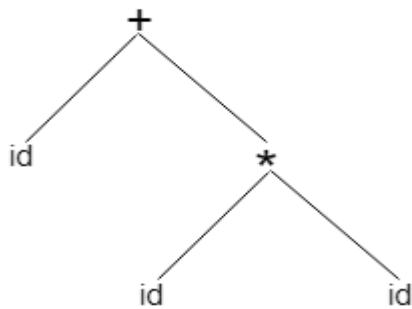
Parse tree and Syntax tree

When you create a parse tree then it contains more details than actually needed. So, it is very difficult to compiler to parse the parse tree. Take the following parse tree as an example:

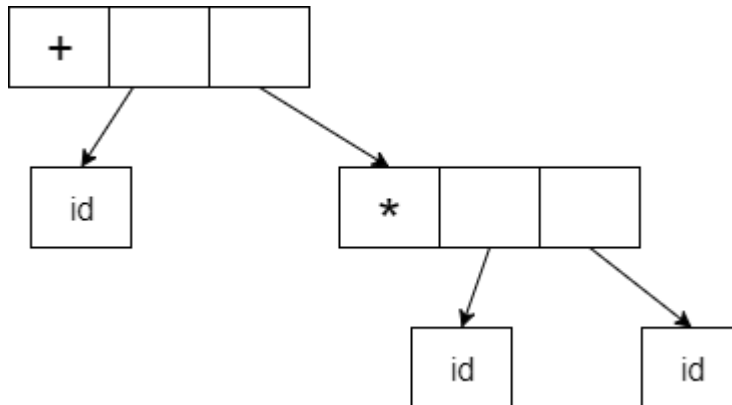


- In the parse tree, most of the leaf nodes are single child to their parent nodes.
- In the syntax tree, we can eliminate this extra information.
- Syntax tree is a variant of parse tree. In the syntax tree, interior nodes are operators and leaves are operands.
- Syntax tree is usually used when represent a program in a tree structure.

A sentence **id + id * id** would have the following syntax tree:



Abstract syntax tree can be represented as:



Abstract syntax trees are important data structures in a compiler. It contains the least unnecessary information.

Abstract syntax trees are more compact than a parse tree and can be easily used by a compiler.

Three address code

- Three-address code is an intermediate code. It is used by the optimizing compilers.
- In three-address code, the given expression is broken down into several separate instructions. These instructions can easily translate into assembly language.
- Each Three address code instruction has at most three operands. It is a combination of assignment and a binary operator.

Example

Given Expression:

1. $a := (-c * b) + (-c * d)$

Three-address code is as follows:

$t_1 := -c \quad t_2 := b * t_1 \quad t_3 := -c \quad t_4 := d * t_3 \quad t_5 := t_2 + t_4 \quad a := t_5$

t is used as registers in the target program.

The three address code can be represented in two forms: **quadruples** and **triples**.

Quadruples

The quadruples have four fields to implement the three-address code. The field of quadruples contains the name of the operator, the first source operand, the second source operand and the result respectively.

Operator
Source 1
Source 2
Destination

Fig: Quadruples field

Example

1. $a := -b * c + d$

Three-address code is as follows:

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

These statements are represented by quadruples as follows:

	Operator	Source 1	Source 2	Destination
(0)	uminus	b	-	t_1
(1)	+	c	d	t_2
(2)	*	t_1	t_2	t_3
(3)	:=	t_3	-	a

Triples

The triples have three fields to implement the three-address code. The field of triples contains the name of the operator, the first source operand and the second source operand.

In triples, the results of respective sub-expressions are denoted by the position of expression. Triple is equivalent to DAG while representing expressions.

Operator
Source 1
Source 2

Fig: Triples field

Example:

1. $a := -b * c + d$

Three address code is as follows:

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

These statements are represented by triples as follows:

	Operator	Source 1	Source 2
(0)	uminus	b	-
(1)	+	c	d
(2)	*	(0)	(1)
(3)	:=	(2)	-

Translation of Assignment Statements

In the syntax directed translation, assignment statement is mainly deals with expressions. The expression can be of type real, integer, array and records.

Consider the grammar

1. $S \rightarrow id := E$
2. $E \rightarrow E1 + E2$
3. $E \rightarrow E1 * E2$
4. $E \rightarrow (E1)$
5. $E \rightarrow id$

The translation scheme of above grammar is given below:

Production rule	Semantic actions
$S \rightarrow id := E$	<pre> {p = look_up(id.name); If p ≠ nil then Emit (p = E.place) Else Error; } </pre>
$E \rightarrow E1 + E2$	<pre> {E.place = newtemp(); Emit (E.place = E1.place '+' E2.place) } </pre>
$E \rightarrow E1 * E2$	<pre> {E.place = newtemp(); Emit (E.place = E1.place '*' E2.place) } </pre>
$E \rightarrow (E1)$	<pre> {E.place = E1.place} </pre>
$E \rightarrow id$	<pre> {p = look_up(id.name); If p ≠ nil then Emit (p = E.place) Else Error; } </pre>

- The p returns the entry for id.name in the symbol table.
- The Emit function is used for appending the three-address code to the output file. Otherwise it will report an error.
- The newtemp() is a function used to generate new temporary variables.
- E.place holds the value of E.

Boolean expressions

Boolean expressions have two primary purposes. They are used for computing the logical values. They are also used as conditional expression using if-then-else or while-do.

Consider the grammar

1. $E \rightarrow E \text{ OR } E$
2. $E \rightarrow E \text{ AND } E$
3. $E \rightarrow \text{NOT } E$
4. $E \rightarrow (E)$
5. $E \rightarrow \text{id relop id}$
6. $E \rightarrow \text{TRUE}$
7. $E \rightarrow \text{FALSE}$

The relop is denoted by $<$, $>$, $<=$, $>=$.

The AND and OR are left associated. NOT has the higher precedence then AND and lastly OR.

Production rule	Semantic actions
$E \rightarrow E1 \text{ OR } E2$	<pre>{E.place = newtemp(); Emit (E.place := E1.place 'OR' E2.place) }</pre>
$E \rightarrow E1 + E2$	<pre>{E.place = newtemp(); Emit (E.place := E1.place 'AND' E2.place) }</pre>
$E \rightarrow \text{NOT } E1$	<pre>{E.place = newtemp(); Emit (E.place := 'NOT' E1.place) }</pre>
$E \rightarrow (E1)$	<pre>{E.place = E1.place }</pre>
$E \rightarrow \text{id relop id2}$	<pre>{E.place = newtemp(); Emit ('if' id1.place relop.op id2.place 'goto' nextstar + 3); EMIT (E.place := '0') EMIT ('goto' nextstat + 2) EMIT (E.place := '1') }</pre>
$E \rightarrow \text{TRUE}$	<pre>{E.place := newtemp(); Emit (E.place := '1') }</pre>
$E \rightarrow \text{FALSE}$	<pre>{E.place := newtemp(); Emit (E.place := '0') }</pre>

The EMIT function is used to generate the three address code and the newtemp() function is used to generate the temporary variables.

The $E \rightarrow \text{id relop id2}$ contains the next_state and it gives the index of next three address statements in the output sequence.

Here is the example which generates the three-address code using the above translation scheme:

1. $p > q$ AND $r < s$ OR $u > r$
2. 100: **if** $p > q$ **goto** 103
3. 101: $t1 := 0$
4. 102: **goto** 104
5. 103: $t1 := 1$
6. 104: **if** $r > s$ **goto** 107
7. 105: $t2 := 0$
8. 106: **goto** 108
9. 107: $t2 := 1$
10. 108: **if** $u > v$ **goto** 111
11. 109: $t3 := 0$
12. 110: **goto** 112
13. 111: $t3 := 1$
14. 112: $t4 := t1$ AND $t2$
15. 113: $t5 := t4$ OR $t3$

Statements that alter the flow of control

The goto statement alters the flow of control. If we implement goto statements then we need to define a LABEL for a statement. A production can be added for this purpose:

1. $S \rightarrow \text{ LABEL } : S$
2. $\text{ LABEL } \rightarrow \text{ id}$

In this production system, semantic action is attached to record the LABEL and its value in the symbol table.

Following grammar used to incorporate structure flow-of-control constructs:

1. $S \rightarrow \text{ if } E \text{ then } S$
2. $S \rightarrow \text{ if } E \text{ then } S \text{ else } S$
3. $S \rightarrow \text{ while } E \text{ do } S$
4. $S \rightarrow \text{ begin } L \text{ end}$
5. $S \rightarrow A$
6. $L \rightarrow L ; S$
7. $L \rightarrow S$

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

Translation scheme for statement that alters flow of control

- We introduce the marker non-terminal M as in case of grammar for Boolean expression.
- This M is put before statement in both if then else. In case of while-do, we need to put M before E as we need to come back to it after executing S.
- In case of if-then-else, if we evaluate E to be true, first S will be executed.
- After this we should ensure that instead of second S, the code after the if-then else will be executed. Then we place another non-terminal marker N after first S.

The grammar is as follows:

1. $S \rightarrow \text{ if } E \text{ then } M S$
2. $S \rightarrow \text{ if } E \text{ then } M S \text{ else } M S$
3. $S \rightarrow \text{ while } M E \text{ do } M S$
4. $S \rightarrow \text{ begin } L \text{ end}$
5. $S \rightarrow A$
6. $L \rightarrow L ; M S$
7. $L \rightarrow S$
8. $M \rightarrow \epsilon$
9. $N \rightarrow \epsilon$

The translation scheme for this grammar is as follows:

Production rule	Semantic actions
$S \rightarrow \text{if } E \text{ then } M \ S1$	BACKPATCH (E.TRUE, M.QUAD) S.NEXT = MERGE (E.FALSE, S1.NEXT)
$S \rightarrow \text{if } E \text{ then } M1 \ S1 \ \text{else } M2 \ S2$	BACKPATCH (E.TRUE, M1.QUAD) BACKPATCH (E.FALSE, M2.QUAD) S.NEXT = MERGE (S1.NEXT, N.NEXT, S2.NEXT)
$S \rightarrow \text{while } M1 \ E \ \text{do } M2 \ S1$	BACKPATCH (S1.NEXT, M1.QUAD) BACKPATCH (E.TRUE, M2.QUAD) S.NEXT = E.FALSE GEN (goto M1.QUAD)
$S \rightarrow \text{begin } L \ \text{end}$	S.NEXT = L.NEXT
$S \rightarrow A$	S.NEXT = MAKELIST ()
$L \rightarrow L ; M \ S$	BACKPATHCH (L1.NEXT, M.QUAD) L.NEXT = S.NEXT
$L \rightarrow S$	L.NEXT = S.NEXT
$M \rightarrow \epsilon$	M.QUAD = NEXTQUAD
$N \rightarrow \epsilon$	N.NEXT = MAKELIST (NEXTQUAD) GEN (goto_)

Postfix Translation

In a production $A \rightarrow \alpha$, the translation rule of A.CODE consists of the concatenation of the CODE translations of the non-terminals in α in the same order as the non-terminals appear in α .

Production can be factored to achieve postfix form.

Postfix translation of while statement

The production

- $S \rightarrow \text{while } M1 \ E \ \text{do } M2 \ S1$

Can be factored as:

- $S \rightarrow C \ S1$
- $C \rightarrow W \ E \ \text{do}$
- $W \rightarrow \text{while}$

A suitable transition scheme would be

Production Rule	Semantic Action
$W \rightarrow \text{while}$	W.QUAD = NEXTQUAD
$C \rightarrow W \ E \ \text{do}$	C W E do

$S \rightarrow C S1$	BACKPATCH ($S1.NEXT, C.QUAD$) $S.NEXT = C.FALSE$ GEN (goto $C.QUAD$)
----------------------	---

Postfix translation of for statement

The production

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

Can be factored as

1. $F \rightarrow \text{for } L$
2. $T \rightarrow F = E1 \text{ by } E2 \text{ to } E3 \text{ do}$
3. $S \rightarrow T S1$

Array references in arithmetic expressions

Elements of arrays can be accessed quickly if the elements are stored in a block of consecutive location. Array can be one dimensional or two dimensional.

For one dimensional array:

1. A: array [low..high] of the ith elements is at:
2. $\text{base} + (i - \text{low}) * \text{width} \rightarrow i * \text{width} + (\text{base} - \text{low} * \text{width})$

Multi-dimensional arrays:

Row major or column major forms

- Row major: $a[1,1], a[1,2], a[1,3], a[2,1], a[2,2], a[2,3]$
- Column major: $a[1,1], a[2,1], a[1, 2], a[2, 2], a[1, 3], a[2,3]$
- In row major form, the address of $a[i1, i2]$ is
- $\text{Base} + ((i1 - \text{low}1) * (\text{high}2 - \text{low}2 + 1) + i2 - \text{low}2) * \text{width}$

Translation scheme for array elements

Limit(array, j) returns $n_j = \text{high}_j - \text{low}_j + 1$

place: the temporary or variables

offset: offset from the base, null if not an array reference

The production:

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

A suitable transition scheme for array elements would be:

Production Rule	Semantic Action
$S \rightarrow L := E$	{ if L.offset = null then emit(L.place := E.place) else EMIT (L.place['L.offset '] := E.place); }
$E \rightarrow E+E$	{ E.place := newtemp; EMIT (E.place := E1.place '+' E2.place); }
$E \rightarrow (E)$	{ E.place := E1.place; }
$E \rightarrow L$	{ if L.offset = null then E.place = L.place else { E.place = newtemp; EMIT (E.place := L.place [' L.offset ']); } }
$L \rightarrow Elist]$	{ L.place = newtemp; L.offset = newtemp; EMIT (L.place := c(Elist.array)); EMIT (L.offset := Elist.place '*' width(Elist.array)); }
$L \rightarrow id$	{ L.place = lookup(id.name); L.offset = null; }
$Elist \rightarrow Elist, E$	{ t := newtemp; m := Elist1.ndim + 1; EMIT (t := Elist1.place '*' limit(Elist1.array, m)); EMIT (t, := t '+' E.place); Elist.array = Elist1.array; Elist.place := t; Elist.ndim := m; }
$Elist \rightarrow id[E$	{ Elist.array := lookup(id.name); Elist.place := E.place Elist.ndim := 1; }

Where:

ndim denotes the number of dimensions.

Limit (array, i) function returns the upper limit along with the dimension of array

width(array) returns the number of bytes for one element of array.

Procedures call

Procedure is an important and frequently used programming construct for a compiler. It is used to generate good code for procedure calls and returns.

Calling sequence:

The translation for a call includes a sequence of actions taken on entry and exit from each procedure. Following actions take place in a calling sequence:

- When a procedure call occurs then space is allocated for activation record.
- Evaluate the argument of the called procedure.
- Establish the environment pointers to enable the called procedure to access data in enclosing blocks.
- Save the state of the calling procedure so that it can resume execution after the call.
- Also save the return address. It is the address of the location to which the called routine must transfer after it is finished.
- Finally generate a jump to the beginning of the code for the called procedure.

Let us consider a grammar for a simple procedure call statement

1. $S \rightarrow \text{call id}(\text{Elist})$
2. $\text{Elist} \rightarrow \text{Elist}, \text{E}$
3. $\text{Elist} \rightarrow \text{E}$

A suitable transition scheme for procedure call would be:

Production Rule	Semantic Action
$S \rightarrow \text{call id}(\text{Elist})$	for each item p on QUEUE do GEN (param p) GEN (call id.PLACE)
$\text{Elist} \rightarrow \text{Elist}, \text{E}$	append E.PLACE to the end of QUEUE
$\text{Elist} \rightarrow \text{E}$	initialize QUEUE to contain only E.PLACE

Queue is used to store the list of parameters in the procedure call.

Declarations

When we encounter declarations, we need to lay out storage for the declared variables. For every local name in a procedure, we create a ST(Symbol Table) entry containing:

1. The type of the name
2. How much storage the name requires

The production:

1. $D \rightarrow \text{integer}, \text{id}$
2. $D \rightarrow \text{real}, \text{id}$
3. $D \rightarrow D1, \text{id}$

A suitable transition scheme for declarations would be:

Production rule	Semantic action
$D \rightarrow \text{integer}, \text{id}$	ENTER (id.PLACE, integer) D.ATTR = integer
$D \rightarrow \text{real}, \text{id}$	ENTER (id.PLACE, real) D.ATTR = real
$D \rightarrow D1, \text{id}$	ENTER (id.PLACE, D1.ATTR) D.ATTR = D1.ATTR

ENTER is used to make the entry into symbol table and **ATTR** is used to trace the data type.

Case Statements

Switch and case statement is available in a variety of languages. The syntax of case statement is as follows:

```
1.  switch E
2.      begin
3.          case V1: S1
4.          case V2: S2
5.          .
6.          .
7.          .
8.  case Vn-1: Sn-1
9.  default: Sn
10.      end
```

The translation scheme for this shown below:

Code to evaluate E into T

```
1.  goto TEST
2.      L1:    code for S1
3.          goto NEXT
4.      L2:    code for S2
5.          goto NEXT
6.          .
7.          .
8.          .
9.      Ln-1:  code for Sn-1
10.         goto NEXT
11.      Ln:   code for Sn
12. goto NEXT
13.      TEST: if T = V1 goto L1
14.         if T = V2 goto L2
15.         .
16.         .
17.         .
18.         if T = Vn-1 goto Ln-1
19.         goto
20. NEXT:
```

- When switch keyword is seen then a new temporary T and two new labels test and next are generated.
- When the case keyword occurs then for each case keyword, a new label Li is created and entered into the symbol table. The value of Vi of each case constant and a pointer to this symbol-table entry are placed on a stack.