

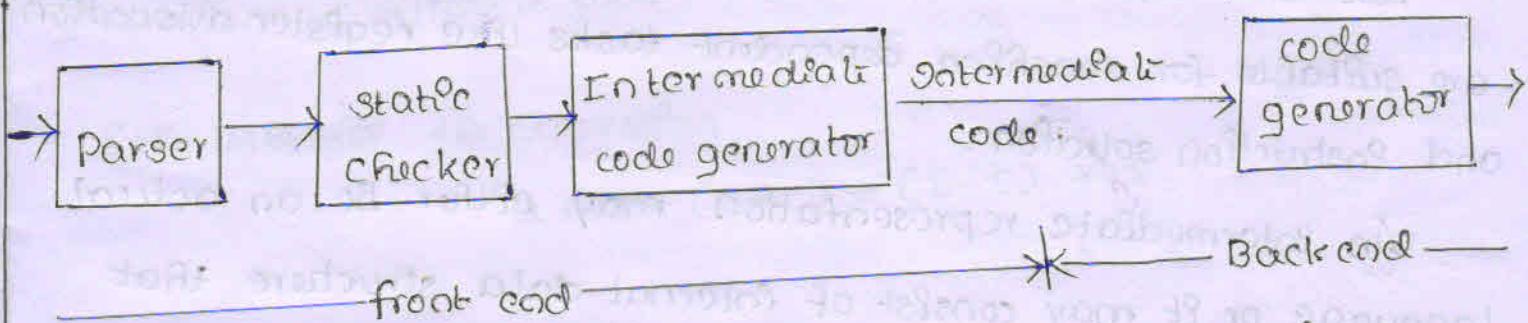
## UNIT-6

# INTERMEDIATE CODE

## GENERATION

Intermediate Code generation.

In the analysis-synthesis model of a compiler, the front end analyzes a source program and creates an intermediate representation, from which the back end generates target code.

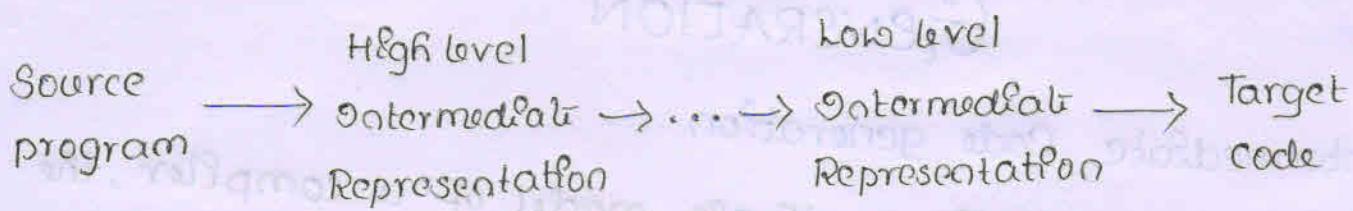


Logical structure of a compiler front end

Parsing, static checking and intermediate code generation are done sequentially; sometimes they can be combined and folded into parsing.

Static checking includes type checking, which ensures that operators are applied to compatible operands. It also includes any syntactic checks that remain after parsing. Ex: It ensures that a break-statement in C is enclosed within a while-, for- or switch-statement; an error is reported if such an enclosing statement does not exist.

In the process of translating a program in a given source language into code for a given target machine, a compiler may construct a sequence of intermediate representation as



High level representations are close to the source language and are well suited to tasks like static type checking.

Ex: Syntax tree

Low level representations are close to the target machine & are suitable for machine dependent tasks like register allocation and instruction selection.

An intermediate representation may either be an actual language or it may consist of internal data structures that are shared by phases of the compiler.

### Variants of Syntax Trees

Nodes in a syntax tree represent constructs in the source program; the children of a node represent the meaningful components of a construct.

A directed acyclic graph (DAG) for an expression identifies the common subexpressions of the expression.

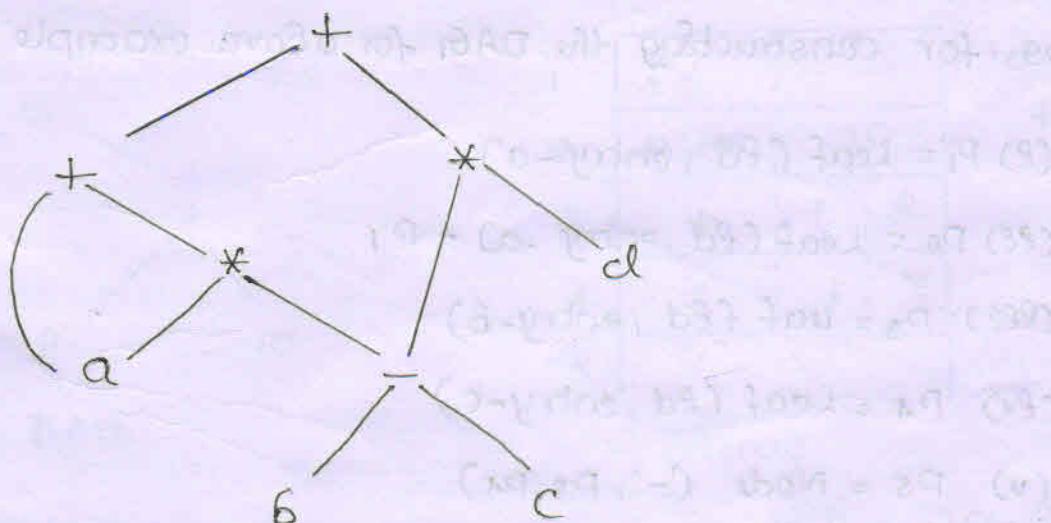
# ① Directed Acyclic Graphs for Expressions.

In DAG leaves represents the atomic operands and interior nodes represents the operators. as in the syntax tree A node  $N$  in a DAG has more than one parent if  $N$  represents a common subexpression ; But in the syntax tree , the tree for the common subexpression would be duplicated as many times as the subexpression appears. In the original expression .

DAG gives the compiler important clues regarding the generation of efficient code to evaluate the expressions .

Ex: DAG for the expression

$$a + a * (b - c) + (b - c) * d$$



↳ The leaf for 'a' has 2 parents, because 'a' appears twice in the expression

↳ The 2 occurrence of the common subexpression  $b - c$  are represented by one node , the node labeled '-'

## SDD to produce DAG.

PRODUCTION	SEMANTIC RULES
(2) $E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node}('+', E_1.\text{node}, T.\text{node})$
(2E) $E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node}('-', E_1.\text{node}, T.\text{node})$
(2EE) $E \rightarrow T$	$E.\text{node} = T.\text{node}$
(2V) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
(V) $T \rightarrow Ed$	$T.\text{node} = \text{new Leaf}(Ed, Ed.\text{entry})$
(V2) $T \rightarrow \text{num}$	$T.\text{node} = \text{new Leaf}(\text{num}, \text{num}.val)$

It will construct a DAG, before creating a new node, these functions first check whether an identical node already exists.

If a previously created identical  $\text{ex}$  node exists, the existing node is returned.

Steps for constructing the DAG for above example .

$$(2) P_1 = \text{Leaf}(Ed, \text{entry}-a)$$

$$(2E) P_2 = \text{Leaf}(Ed, \text{entry}-a) = P_1$$

$$(2EE) P_3 = \text{Leaf}(Ed, \text{entry}-b)$$

$$(2V) P_4 = \text{Leaf}(Ed, \text{entry}-c)$$

$$(V) P_5 = \text{Node}('-', P_3, P_4)$$

$$(V2) P_6 = \text{Node}('* ', P_1, P_5)$$

$$(V2E) P_7 = \text{Node}('+ ', P_1, P_6)$$

$$(2) P_8 = \text{Leaf}(Ed, \text{entry}-b) = P_3$$

$$(2X) P_9 = \text{Leaf}(Ed, \text{entry}-c) = P_4$$

$$(X) P_{10} = \text{Node}('-', P_3, P_4) = P_5$$

$$(X2) P_{11} = \text{Leaf}(Ed, \text{entry}-d)$$

(Ex8)  $P_{12} = \text{Node} ('*', P_5, P_{11})$

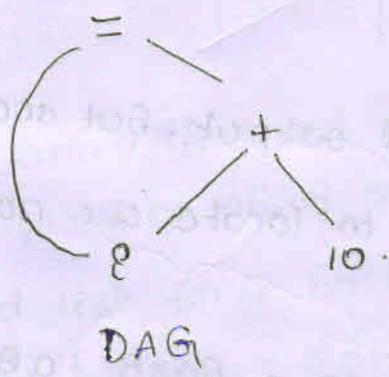
(Ex8)  $P_{13} = \text{Node} ('+', P_7, P_{12})$

When the call to Leaf (8th entry-a) is repeated at step 8, the node created by the previous call is returned, so  $P_2 = P_1$ .

### The Value-Number Method for Constructing DAG's

- \* The nodes of a DAG are stored in an array of records
- \* Each row of array represents one record & therefore one node.
- \* In each record, the first field is an operation code, indicating the label of the node. Leaves have one additional field which holds the lexical value and interior nodes have 2 additional fields indicating the left and right children.

Ex: DAG for  $E = E + 10$  allocated in an Array



DAG

	Op		
1	Op		
2	Num	10	
3	+	1	2
4	=	1	3
5			

Array

to entry  
for E

- \* In this array, we refer to nodes by giving the integer index of the record for that node within the array.
- \* This integer historically has been called the "value number".
- \* The integer historically has been called the "value number".
- \* For above example node labeled + has value number 3 & its left & right children have value numbers 1 & 2 respectively.

Suppose that nodes are stored in an array & each node is referred to by its value numbers. Let the signature of an interior node be the triple  $\langle \text{op}, l, r \rangle$  where op is the label, l is its left child's value number & r its right child's value number. A unary operator may be assumed to have  $r=0$ .

ALGORITHM: To construct the nodes of a DAG using value number method.

INPUT: Label op, node l and node r

OUTPUT: The value number of a node in the array with signature  $\langle \text{op}, l, r \rangle$

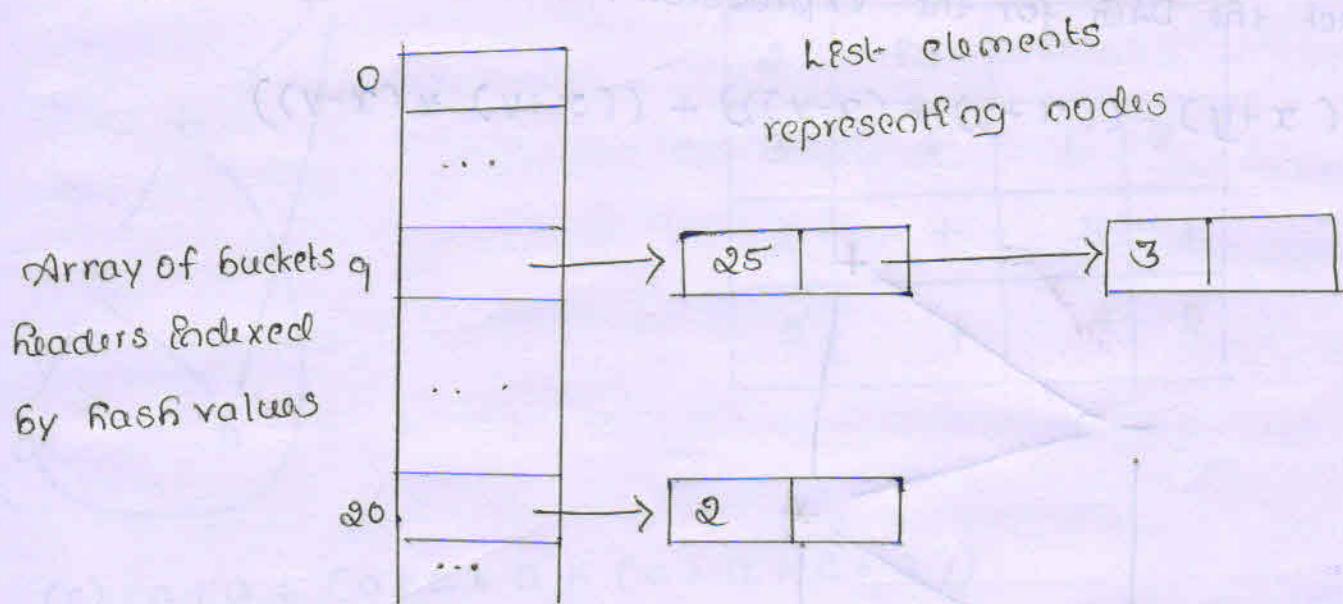
METHOD: Search the array for a node M with label op, left child l and right child r. If there is such a node, return the value number of M. If not, create in the array a new node M with label op, left child l and right child r & return its value number.

Above algorithm yields the desired output, but searching the entire array every time we are asked to locate one node is expensive.

A more efficient approach is to use a hash table, in which the nodes are put into "buckets" each of which typically will have only a few nodes. It supports dictionaries which is an abstract data type that allows us to insert & delete elements of a set & to determine whether a given element is currently in the set.

To construct a hash table for the nodes of a DAG, we need a hash function  $R$  that computes the index of the bucket for a signature  $\langle op, l, r \rangle$ .

The bucket index  $R(\langle op, l, r \rangle)$  is computed deterministically from  $op, l$  &  $r$  so that we may repeat the calculation & always get to the same bucket index for node  $\langle op, l, r \rangle$ .  
The buckets can be implemented as linked list as,



An array indexed by hash value, holds the bucket readers, each of which points to the first cell of a list. Within the linked list for a bucket, each cell holds the value number of one of the nodes that hash to that bucket. That is, node  $\langle op, l, r \rangle$  can be found on the list whose Reader is at index  $R(\langle op, l, r \rangle)$  of the array.

Thus, given the output input node  $op, l$  &  $r$  we compute the bucket index  $R(\langle op, l, r \rangle)$  & search the list of cells in this bucket for the given input node.

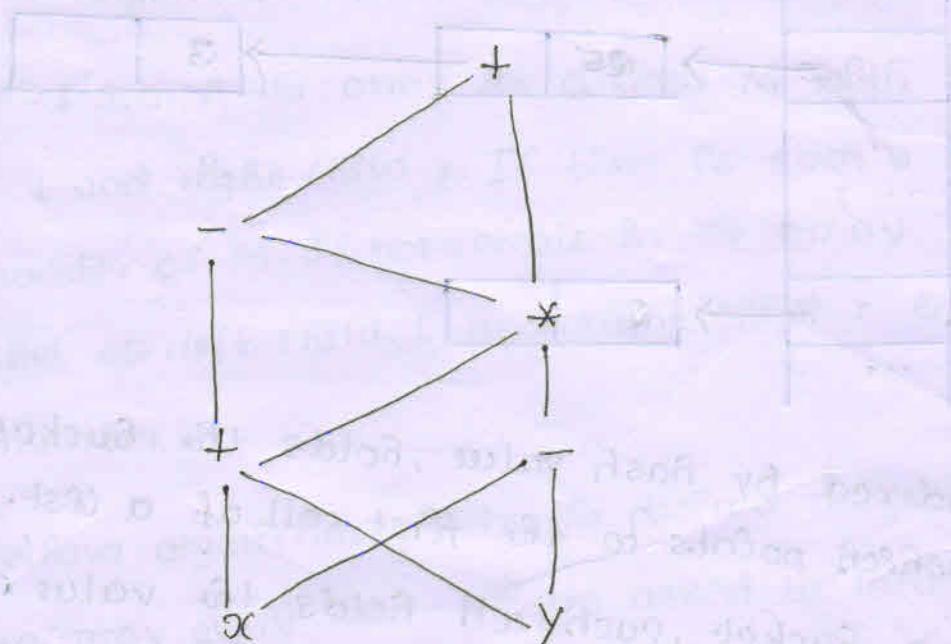
For each value number 'v' found in a cell, we must-

check whether the signature  $\langle op, l, r \rangle$  of the input node matches the node with value number  $v$  in the list of the cells. If we find a match, we return  $v$ . If we find no match, we know no such node can exist in any other bucket, so we create a new cell, add it to the list of cells for bucket- $i$ , index  $R\langle op, l, r \rangle$  & return the value number in that new cell.

Problems.

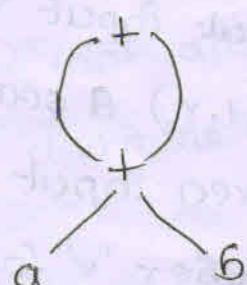
Construct the DAG for the expression.

$$((x+y) - ((x+y) * (x-y))) + ((x+y) * (x-y))$$



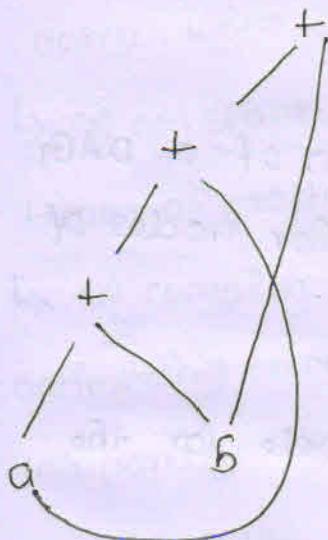
Construct the DAG & identify the value number for the sub expressions of the following expressions, assuming + associates from the left.

(a)  $a+b+(a+b)$



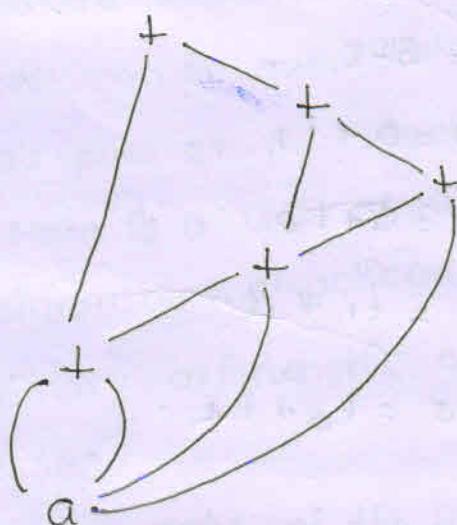
1	$\text{Ed}$	a
2	$\text{Ed}$	b
3	+	1   2
4	+	3   3

(b)  $a + b + a + b$



1	$\text{Ed}$	a	$\text{Ed}$
2	$\text{Ed}$	b	
3	+	1   2	
4	+	3   1	
5	+	4   2	

(c)  $a + a + (a + a + a + (a + a + a + a))$



1	$\text{Ed}$	a	
2	+	1   1	
3	+	2   1	
4	+	3   1	
5	+	3   4	
6	+	2   5	

## Three - Address Code

In 3-address code, there is at most one operator on the right & side of an instruction, e.g., no built up arithmetic expression are permitted.

Thus a source-language expression like  $x+y * z$  might be translated into the sequence of 3 address instructions

$$t_1 = y * z$$

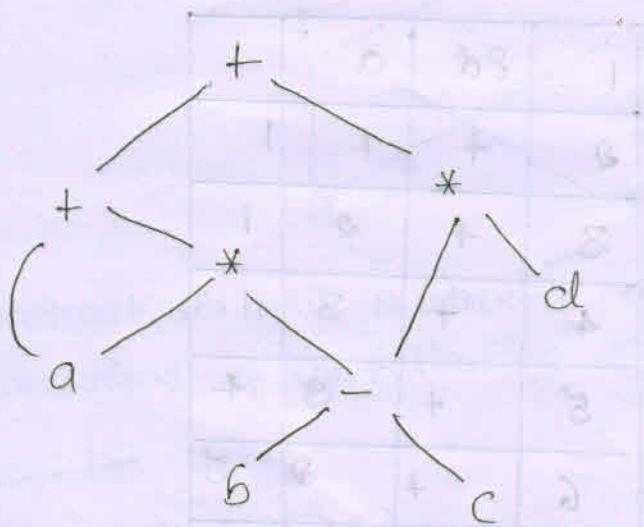
$$t_2 = x + t_1$$

3 ADDRESS (3)

where  $t_1$  &  $t_2$  are compiler generated names.

3 address code is a linearized representation of a DAG in which explicit names correspond to the interior nodes of the graph.

Ex: Write DAG & its corresponding 3 address code for the expression  $a + a * (b - c) + (b - c) * d$ .



DAG

$$t_1 = b - c$$

$$t_2 = a * t_1$$

$$t_3 = a + t_2$$

$$t_4 = t_1 * d$$

$$t_5 = t_3 + t_4$$

3 -address code

## Addresses and Instructions

3-address code is built from 2 concepts : address & instructions.

An address can be one of the following.

- ↳ A name : For convenience, we allow source program names to appear as addresses in 3-address code. In an implementation, a source name is replaced by a pointer to its symbol table entry, where all information about the name is kept.
- ↳ A constant : A compiler must deal with many different types of constants and variables.
- ↳ A compiler generated temporary : It is useful, especially in optimizing compilers, to create a distinct name each time a temporary is needed.

Symbolic labels will be used by instructions that alter the flow of control. A symbolic label represents the index of a 3-address instruction in the sequence of instructions. Actual indexes can be substituted for the labels, either by making a separate pass or by "backpatching".

- Here is a list of the common 3-address instruction forms
- (E) Assignment instructions of the form  $x = y \text{ op } z$  where  $op$  is a binary arithmetic or logical operation &  $x, y$  &  $z$  are addresses.
  - (UE) Assignments of the form  $x = \text{op } y$ , where  $\text{op}$  is a unary operation
  - (PE) Copy instructions of the form  $x = y$ , where  $x$  is assigned the value of  $y$ .

(cv) An unconditional jump goto L. The 3-address instruction with label L is the next to be executed.

(cv) Conditional jumps of the form If  $x$  goto L and If False  $x$  goto L. These instructions execute the instruction with label L next if  $x$  is true and false respectively.

(cvf) Conditional jumps such as If  $x$  relop  $y$  goto L, which apply a relational operator ( $<$ ,  $=$ ,  $\geq$  etc) to  $x$  &  $y$  & execute the instruction with label L next if  $x$  stands in relation relop to  $y$ . If not, the 3 address instruction following If  $x$  relop  $y$  goto L is executed next, in sequence.

(cvff) Procedures calls & returns are implemented using the following instructions: param  $x$  for parameters; call  $p,n$  &  $y=call p,n$  for procedure & function calls respectively & return  $y$ , where  $y$  representing a returned value, is optional.

param  $x_1$

param  $x_2$

...

param  $x_n$

call  $p,n$

The integer  $n$  indicating the no. of actual parameters in call  $p,n$  is not redundant because calls can be nested.

(cvfff) Indexed copy instructions of the form  $x=y[?]$  and  $x[?] = y$ . The instruction  $x=y[?]$  sets  $x$  to the value in the location  $y[?]$ . The instruction  $x[?] = y$  sets the contents of the locations  $y$  units beyond  $x$  to the value of  $y$ .

(Ex) Address & pointer assignments of the form  $x = \&y$ ,  $x = *y$  and  $*x = y$ .

Ex % Consider the statement

do  $\ell = \ell + 1$ ; while ( $a[\ell] < v$ );

2 possible translations of this statement are

L :  $t_1 = \ell + 1$

$\ell = t_1$

$t_2 = \ell * 8$

$t_3 = a[t_2]$

if  $t_3 < v$  goto L.

100 :  $t_1 = \ell + 1$

101 :  $\ell = t_1$

102 :  $t_2 = \ell * 8$

103 :  $t_3 = a[t_2]$

104 : if  $t_3 < v$  goto 100

(b) position number.

(a) symbolic labels

The translation in (a) uses a symbolic label L, attached to the first instruction. The translation in (b) shows position number for the instructions, starting arbitrarily at position 100. In both translations, the last instruction is a conditional jump to the first instruction. The multiplication  $\ell * 8$  is appropriate for an array of elements that each take 8 units of space.

### Quadruples

A quadruple has 4 fields, which we call op, arg<sub>1</sub>, arg<sub>2</sub>, & result. The op field contains an internal code for the operator.

Ex %  $x = y + z$  is represented by placing + in op, y in arg<sub>1</sub>, z in arg<sub>2</sub> & x in result.

The following are some exceptions to this rule.

(e) Instructions with unary operators like  $x = \text{minus } y$  or  $x = y$  do not use  $\text{arg}_2$ . Note that for a copy statement like  $x = y, \text{op } \text{ps} =$ , while for most other operations, the assignment operator  $\text{ps}$  is employed.

operator is implied.  
cpp operators like param use refl for argz nor result.

(C++) Operators like `param` use `label` as a label. Put the target label

(88) Conditional & unconditional jumps put the target label

as a result.

Ex: write quadruples for  $a = b * -c + b * -c$

$$t_1 = \text{minus } C$$

$$t_2 = b * t_1$$

$$\text{tg} = \frac{\text{opposite}}{\text{adjacent}}$$

$$t_4 = 6 * t_3$$

$$t_5 = t_2 + t_4$$

$$a = 15$$

	op	arg1	arg2	result
0	minus	c		t1
1	*	6	t1	t2
2	minus	c		t3
3	*	6	t3	t4
4	+	t2	t4	t5
5	=	t5		a

(a) 3-address code

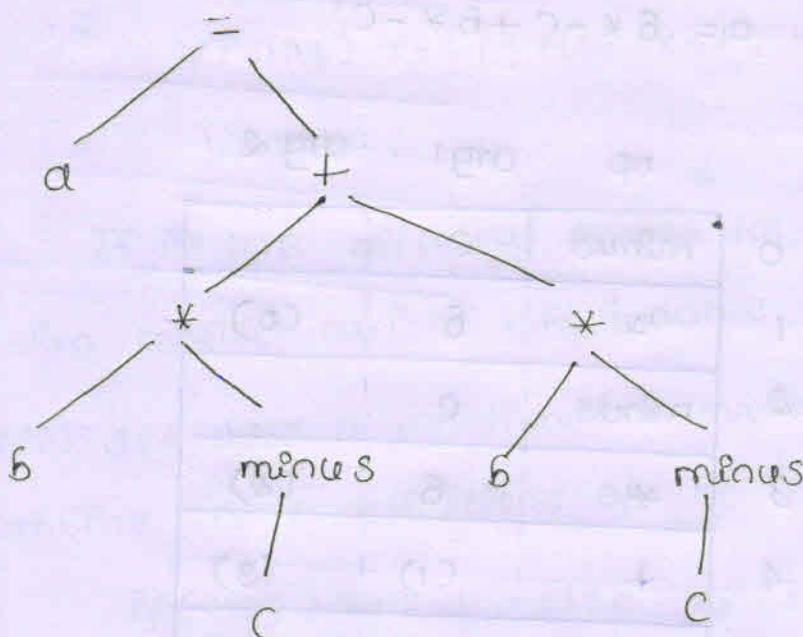
### (6) Quadruples.

The special operator minus is used to distinguish the unary minus operator.

## Triples

Triples  
A triples has only 3 fields , op , arg<sub>1</sub> & arg<sub>2</sub>. Note that the result field in quadruples is used primarily for temporary names . Using triples , we refer to the result of the operation x op y by its position , rather than by an explicit temporary names .

Ex: write triples for  $a = b * c + b * c$



	op	arg1	arg2
0	minus	'c'	' '
1	*	'b'	(0)
2	minus	'c'	' '
3	*	'b'	(2)
4	+	(1)	(3)
5	=	'a'	(4)

(6) triples.

(a) Syntax tree

A ternary operator like  $x[e] = y$  requires 2 entries in the triple structure, for example, we can put  $x$  &  $e$  in one triple &  $y$  in the next.

The benefits of quadruples over triples can be seen in an optimizing compiler, where instructions are often moved around. With quadruples, if we move an instruction that computes a temporary  $t$ , then the instructions that use  $t$  require no change. With triples, the result of an operation is referred to by its position, so moving an instruction may require us to change all references to that result.

### Indirect triples.

Indirect triples consist of a listing of pointers to triples, rather than a listing of triples themselves.

With indirect triples, an optimizing compiler can move an

instruction by reordering the instruction list without affecting the triples themselves.

Ex: write indirect triples for  $a = 6 * -c + 6 * -c$

Instruction	Section
35 (0)	
36 (1)	
37 (2)	
38 (3)	
39 (4)	
40 (5)	
...	

	op	arg1	arg2
0	minus	c	
1	*	6	(0)
2	minus	c	
3	*	6	(2)
4	+	(1)	(3)
5	=	a	(4)
			...

### Static Single Assignment Form

Static single assignment form (SSA) is an intermediate representation that facilitates certain code optimization. 2 distinctive aspects of SSA that distinguish SSA from 3-address code

- (e) All assignments in SSA are to variables with distinct names.

$$\text{Ex: } p = a + b$$

$$p_1 = a + b$$

$$q_1 = p - c$$

$$q_1 = p_1 - c$$

$$p = q_1 * d$$

$$p_2 = q_1 * d$$

$$p = e - p$$

$$p_3 = e - p_2$$

$$q_2 = p + q_1$$

$$q_2 = p_3 + q_1$$

3-address code

static single assignment form

The same variable may be defined in different control flow paths in a program. For example, the source program

```
if(flag) x = -1; else x = 1;
```

```
y = x * a;
```

If we use different names for  $x$  in the true part & false then conflict arises which name should use in  $y = x * a$ .  
(P2) SSA uses a notational convention called  $\phi$ -function to combine the definitions of  $x$ .

```
if(flag) x1 = -1; else x2 = 1;
```

```
x3 =  $\phi(x_1, x_2);$ 
```

Here  $\phi(x_1, x_2)$  has the value  $x_1$  if the control flow passes through the true part of the conditional & the value  $x_2$  if the control flow passes through the false part.

Translate the arithmetic expression  $a + (b+c)$  into

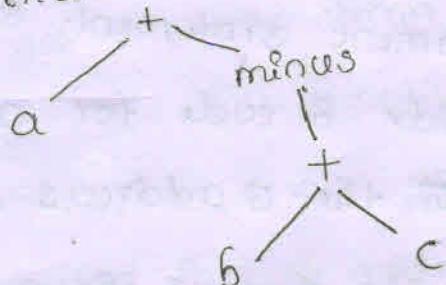
(a) A syntax tree

(b) Quadruples

(c) Triples

(d) Indirect triples.

(a) Syntax tree



$$t_1 = b + c$$

$$t_2 = \text{minus } t_1$$

$$t_3 = a + t_2$$

### (6) Quadruples

	op	arg1	arg2	result
1 0	+	b	c	t <sub>1</sub>
2 1	minus	t <sub>1</sub>		t <sub>2</sub>
3 2	=	a	t <sub>2</sub>	t <sub>3</sub>

### (c) triples

	op	arg1	arg2
0	+	b	c
1	minus	(0)	
2	=	a	(1)

### (d) Indirect triples.

Instructions

	op	arg1	arg2
35	(0)	+	b
36	(1)	minus	(0)
37	(2)	=	a

### Translation of Expressions.

An expression with more than one operator, like  $a+b*c$ , will translate into instructions with at most one operator per instruction. An array reference  $A[i][j]$  will expand into a sequence of 3-address instructions that calculate an address for the reference.

### Operations within Expressions

The following syntax-directed definition builds up the 3-address code for an assignment statement & using attributes code for s & attributes addr & code for an expression e. Attributes s.code & e.code denotes the 3 address code for s & e respectively. Attribute e.addr denotes the address

that will hold the value of E.

PRODUCTION	SEMANTIC RULES
$S \rightarrow Ed = E;$	$E.code = E.code \parallel$ $\text{gen}(\text{top.get}(Ed.\text{lexeme}) = 'E.addr')$
$E \rightarrow E_1 + E_2$	$E.addr = \text{new temp}()$ $E.code = E_1.code \parallel E_2.code \parallel$ $\text{gen}(E.addr = 'E_1.addr + E_2.addr')$
$1 - E_1$	$E.addr = \text{new temp}()$ $E.code = E_1.code \parallel$ $\text{gen}(E.addr = 'minus' E_1.addr)$
$  (E_1)$	$E.addr = E_1.addr$ $E.code = E_1.code$
$  Ed$	$E.addr = \text{top.get}(Ed.\text{lexeme})$ $E.code = ''$

Consider  $E \rightarrow Ed$ , when an expression is a single identifier say x then x itself holds the value of the expression. So semantic rules E.addr point to the symbol table entry. If top denote the current symbol table then function top.get retrieves the entry when it is applied to the string representation Ed.lexeme. E.code is set to empty string.

For  $E \rightarrow E_1 + E_2$ , compute value of E from  $E_1$  &  $E_2$ . Values are computed into newly generated temporary names. If  $E_1$  is computed into  $E_1.addr$  &  $E_2$  into  $E_2.addr$  then  $E_1 + E_2$

translates into  $t = \epsilon_1.\text{addr} + \epsilon_2.\text{addr}$ , where  $t$  is a new temporary name.  $\epsilon.\text{addr}$  is set to  $t$ .

The notation  $\text{gen}(x' = 'y' + 'z')$  is used to represent the 3 address instruction  $x = y + z$ . Variables like  $x, y$  &  $z$  are evaluated & when passed to  $\text{gen}$  & quoted strings like  $'='$  are taken literally.

$\epsilon.\text{code}$  is built up by concatenating  $\epsilon_1.\text{code}$ ,  $\epsilon_2.\text{code}$  & an instruction that adds the values of  $\epsilon_1$  &  $\epsilon_2$ . The construction puts the result of the addition into a new temporary names for  $\epsilon$ , denoted by  $\epsilon.\text{addr}$ .

The production  $\mathcal{S} \rightarrow \mathcal{E}d = \epsilon$ ; generates instructions that assign the value of expression  $\mathcal{E}$  to the identifier  $\mathcal{E}d$ . The semantic rule for this production uses function  $\text{top.get}$  to determine the address of the Identifier represented by  $\mathcal{E}d$ .

$\mathcal{S}.\text{code}$  consists of the instructions to compute the value of  $\mathcal{E}$  into an address given by  $\epsilon.\text{addr}$ , followed by an assignment to the address  $\text{top.get}(\mathcal{E}d, \text{lexeme})$  for this instance of  $\mathcal{E}d$ .

## Incremental Translation.

In Incremental translation, instead of building up  $\epsilon.\text{code}$  we can arrange to generate only the new 3-address instruction. In the incremental approach,  $\text{gen}$  not only constructs a 3 address instruction, it appends the instruction to the sequence of instruction generated so far.

With the incremental approach, the  $\text{code}$  attribute is not used, since there is a single sequence of instructions that is created by successive calls to  $\text{gen}$ .

PRODUCTION	SEMANTIC RULES
$s \rightarrow \text{Pd} = \epsilon$	{gen (top.get (Pd.lexeme) '=' E.addr);}
$\epsilon \rightarrow \epsilon_1 + \epsilon_2$	{E.addr = new Temp(); gen (E.addr !='+' E1.addr + E2.addr);}
$1 \div \epsilon_1$	{E.addr = new Temp(); gen (E.addr != 'minus' E1.addr);}
$  \epsilon_1$	{E.addr = E1.addr;}
$1 \text{ Pd}$	{E.addr = top.get (Pd.lexeme);}

The semantic rule for  $\epsilon \rightarrow \epsilon_1 + \epsilon_2$  simply calls gen to generate an add instruction. The instructions to compute  $\epsilon_1$  into  $\epsilon_1.\text{addr}$  &  $\epsilon_2$  into  $\epsilon_2.\text{addr}$  have already been generated.

The new semantic action for  $\epsilon \rightarrow \epsilon_1 + \epsilon_2$  creates a node by using a constructor, as in

$\epsilon \rightarrow \epsilon_1 + \epsilon_2 \quad \& \quad \epsilon.\text{addr} = \text{new Node} ('+', \epsilon_1.\text{addr}, \epsilon_2.\text{addr})$

Here, attribute addr represents the address of a node rather than a variable or constant.

Addressing array elements.

Array elements can be accessed quickly if they are stored in a block of consecutive locations. If the width of each array element is w, then the  $i^{\text{th}}$  element of array A begins in location

Base is the relative address of  $A[0]$ .

The formula (1) generalizes to 2 or more dimensions. In 2 dimension,  $A[\ell_1][\ell_2]$ , let  $w_1$  be the width of a row & let  $w_2$  be the width of an element in a row. The relative address of  $A[\ell_1][\ell_2]$  can then be calculated by

$$\text{base} + \ell_1 * w_1 + \ell_2 * w_2 - (2)$$

In  $k$  dimensions, the formula is,

$$\text{base} + \ell_1 * w_1 + \ell_2 * w_2 + \dots + \ell_k * w_k - (3)$$

The relative address can be calculated in terms of no. of elements of along dimension of the array and the width of a single element of the array.

$w = w_k$  of a single element of the array, the location of  $A[\ell_1][\ell_2]$  is given by

In 2 dimensions, the location of  $A[\ell_1][\ell_2]$  is given by

$$\text{base} + (\ell_1 * n_2 + \ell_2) * w - (4)$$

In  $k$  dimensions,

$$\text{base} + ((\dots((\ell_1 * n_2 + \ell_2) * n_3 + \ell_3) \dots) * n_k + \ell_k) * w - (5)$$

The array elements need not be numbered starting at 0.

In a one-dimensional array, the array elements are numbered  $\text{low}, \text{low}+1, \dots, \text{high}$  & base is the relative address. (1) can be replaced as

$$\text{base} + (\ell - \text{low}) * w - (6)$$

Both (1) & (6) can be rewritten as

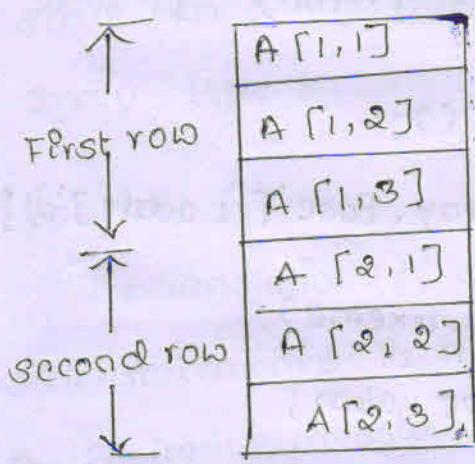
$$\ell * w + c \quad \text{where,}$$

$$c = \underline{\text{base}} * \text{low}$$

$$c = \text{base} - \text{low} * w \quad \text{calculated at complete time}$$

Compile time precalculation can not be used when the array's size is dynamic. If we don't know the values of low & high at compile time, then we can't compute constants such as C.

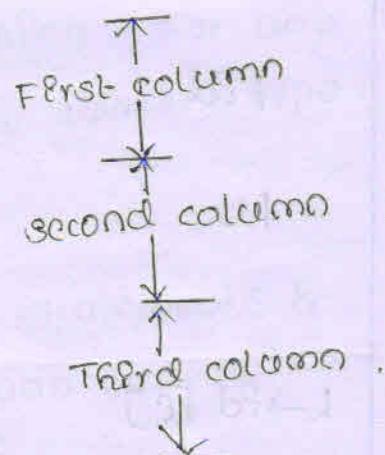
The above address calculations are based on row-major layout for arrays. A 2-dimensional array is normally stored in either row-major or column major forms.



(a) Row major

A [1,1]
A [2,1]
A [1,2]
A [2,2]
A [1,3]
A [2,3]

(b) Column major



We can generalize row or column major form to many dimensions. The generalization of row-major form is to store the elements in such a way that, as we scan down a block of storage, the rightmost subscripts appear to vary fastest, like the numbers on an odometer. Column-major form generalizes to the opposite arrangement, with the leftmost subscripts varying fastest.

### Translation of Array References.

Let nonterminal L generate an array name followed by a sequence of index expressions

$L \rightarrow L[e] \mid Ld[e]$

PRODUCTION	SEMANTIC RULES
$\$ \rightarrow \text{Ed} = E;$	$\{ \text{gen}(\text{top.get}(\text{Ed.lexeme}) = E.\text{addr}); \}$
$L L = E;$	$\{ \text{gen}(L.\text{array.base}[L.\text{addr}] = E.\text{addr}); \}$
$E \rightarrow E_1 + E_2$	$\{ E.\text{addr} = \text{new temp}(); \}$ $\text{gen}(E.\text{addr} = E_1.\text{addr} + E_2.\text{addr}); \}$
$L Ed$	$\{ E.\text{addr} = \text{top.get}(\text{Ed.lexeme}); \}$
$  L$	$\{ E.\text{addr} = \text{new temp}(); \}$ $\text{gen}(E.\text{addr} = L.\text{array.base}[L.\text{addr}]); \}$
$L \rightarrow \text{Ed } [e]$	$\{ L.\text{array} = \text{top.get}(\text{Ed.lexeme}); \}$ $L.\text{type} = L.\text{array.type.elem}; \}$ $L.\text{addr} = \text{new temp}(); \}$ $\text{gen}(L.\text{addr} = E.\text{addr} * L.\text{type.width}); \}$
$  L_i [e]$	$\{ L.\text{array} = L_i.\text{array}; \}$ $L.\text{type} = L_i.\text{type.elem}; \}$ $t = \text{new temp}(); \}$ $L.\text{addr} = \text{new temp}(); \}$ $\text{gen}(t = E.\text{addr} * L.\text{type.width}); \}$ $\text{gen}(L.\text{addr} = L_i.\text{addr} + t); \}$

Non-terminal  $L$  has 3 synthesized attributes

(a)  $L.\text{addr}$  denotes a temporary that is used while computing the offset for the array references by summing the terms

$$l_0 * 10^0$$

(88)  $\text{L.array}$  is a pointer to the symbol table entry name. The base address of the array, say,  $\text{L.array}.base$  is used to determine the actual l-value of an array reference after all the index expressions are analyzed.

(88)  $\text{L.type}$  is the type of the subarray generated by  $\text{L}$ . For any type  $t$ , we assume that its width is given by  $t.width$ . We use types as attributes, rather than widths, since types are needed anyway for type checking. For any array type  $t$ , suppose that  $t.elem$  gives the element type.

## Control Flow

The translation of statements such as if-else-statements & while-statements is tied to the translation of boolean expression. The boolean expression are often used to alter the flow of control.

(8) Alter the flow of control

(88) Compute logical values

## Boolean Expressions.

Boolean expressions are composed of the Boolean operators applied to elements that are Boolean variables or relational expressions.

Relational expressions are of the form  $E_1 \text{rel} E_2$  where  $E_1$  &  $E_2$  are arithmetic expressions &  $\text{rel.op}$  is used to indicate which of the 6 comparison operators  $<, \leq, =, !=, >$  or  $\geq$  is represented by  $\text{rel}$ .

The semantic definition of the programming language determines whether all part of a boolean expression must be evaluated.

Ex: In expression  $B_1 \text{|| } B_2$ , if  $B_1$  is true then the expression is true.

In expression  $B_1 \text{&&} B_2$  if  $B_1$  is false then the entire expression is false.

### Short-Circuit code

In short circuit (or jumping) code, the boolean operators  $\text{&&}$ ,  $\text{||}$ ,  $\text{!}$  translate into jumps. The operators themselves do not appear in the code. Instead, the value of a boolean expression is represented by a position in the code sequence.

Ex: If  $(x < 100 \text{ || } x > 200 \text{ && } x != y)$   $x = 0;$

This statement might be translated into the code

if  $x < 100$  goto L<sub>2</sub>

if false  $x > 200$  goto L<sub>1</sub>

if false  $x != y$  goto L<sub>1</sub>

L<sub>2</sub> :  $x = 0$

L<sub>1</sub> :

The boolean expression is true if control reaches label

### Flow-of-control Statements

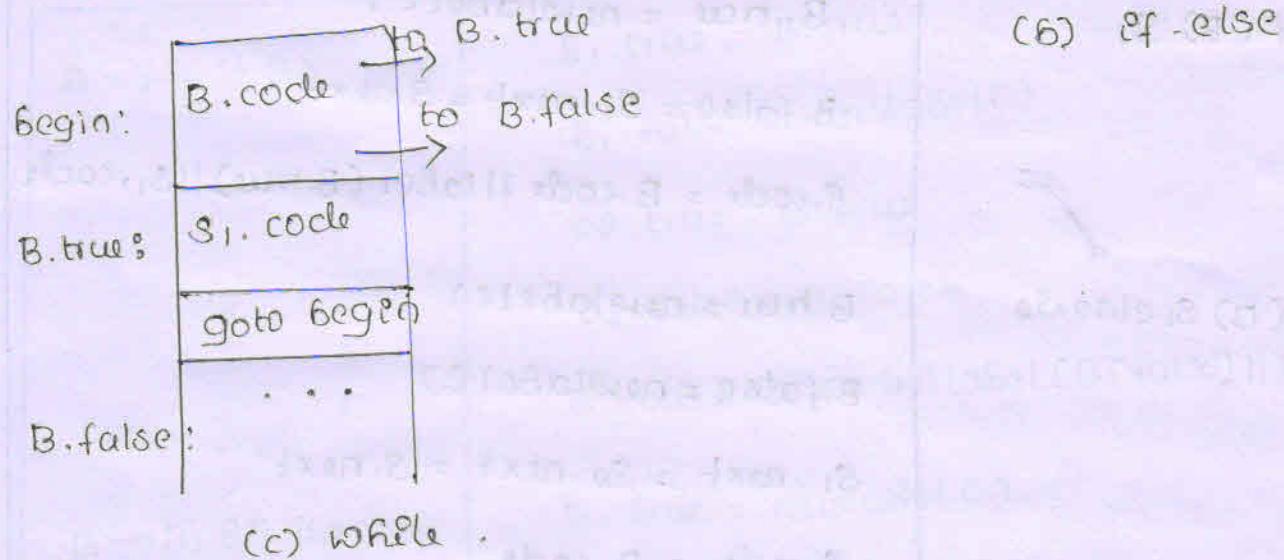
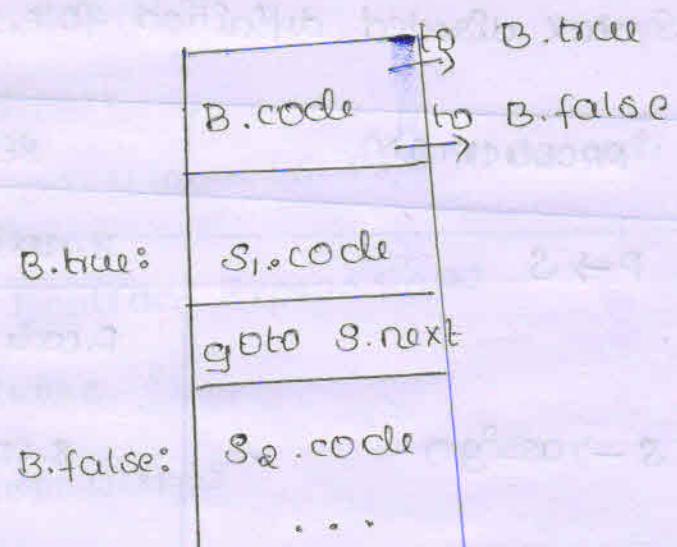
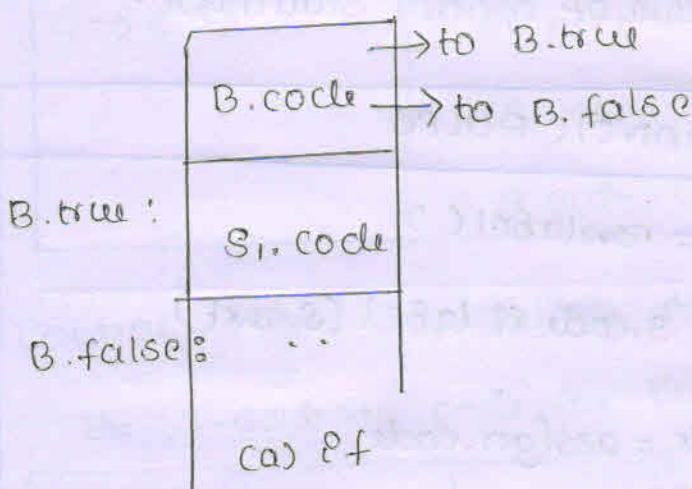
Consider the grammar

$S \rightarrow \text{if } (B) S_1$

$S \rightarrow \text{if } (B) S_1 \text{ else } S_2$

$S \rightarrow \text{while } (B) S_1$

In these productions, nonterminal B represents a boolean expression & non terminal S represents a statement.



The translation of if (B) S<sub>1</sub> consists of B.code followed by

the translation of B<sub>1</sub>.code. Within B<sub>1</sub>.code are jumps based on the value of B<sub>1</sub>.  
If B is true, control flows to the first instruction of S<sub>1</sub>.code  
If B is false, control flows to the instruction immediately following S<sub>1</sub>.code

With boolean expression B, we associate 2 labels

- (P) B.true , the label to which control flows if B is true
  - (P) B.false , the label to which control flows if B is false
- With statement S, we associate an enforced attribute

$s.next$  denoting a label for the instruction immediately after the code for  $s$

Syntax directed definition for flow of control statement.

PRODUCTION	SEMANTIC RULES
$P \rightarrow S$	$s.next = \text{newlabel}()$ $p.code = s.code \sqcup \text{label}(s.next)$
$S \rightarrow \text{assign}$	$s.code = \text{assgn}.code$
$S \rightarrow \text{if } (B) S_1$	$B.\text{true} = \text{newlabel}()$ $B.\text{false} = S_1.\text{next} = S.\text{next}$ $s.code = B.code \sqcup \text{label}(B.\text{true}) \sqcup S_1.\text{code}$
$S \rightarrow \text{if } (B) S_1 \text{ else } S_2$	$B.\text{true} = \text{newlabel}()$ $B.\text{false} = \text{newlabel}()$ $S_1.\text{next} = S_2.\text{next} = S.\text{next}$ $s.code = B.code$ $\sqcup \text{label}(B.\text{true}) \sqcup S_1.\text{code}$ $\sqcup \text{gen('goto', } S.\text{next})$ $\sqcup \text{label}(B.\text{false}) \sqcup S_2.\text{code}$
$S \rightarrow \text{while } (B) S_1$	$\text{begin} = \text{newlabel}()$ $B.\text{true} = \text{newlabel}()$ $B.\text{false} = S.\text{next}$ $S_1.\text{next} = \text{begin}$ $s.code = \text{label } (\text{begin}) \sqcup B.\text{code}$

$S \rightarrow S_1 S_2$

```

|| label(B.true) || S1.code
|| gen('goto' begin)
S1.next = newlabel(C)
S2.next = S.next
S.code = S1.code || label(S1.next) || S2.code

```

### Control - Flow Translation of Boolean Expressions

The 3-address code for booleans is,

PRODUCTION	SEMANTIC RULES
$B \rightarrow B_1 \text{    } B_2$	$B_1.\text{true} = B.\text{true}$ $B_1.\text{false} = \text{newlabel}(C)$ $B_2.\text{true} = B.\text{true}$ $B_2.\text{false} = B.\text{false}$ $B.\text{code} = B_1.\text{code} \text{    } \text{label}(B.\text{false}) \text{    } B_2.\text{code}$
$B \rightarrow B_1 \text{ && } B_2$	$B_1.\text{true} = \text{newlabel}(C)$ $B_1.\text{false} = B.\text{false}$ $B_2.\text{true} = B.\text{true}$ $B_2.\text{false} = B.\text{false}$ $B.\text{code} = B_1.\text{code} \text{    } \text{label}(B_1.\text{true}) \text{    } B_2.\text{code}$
$B \rightarrow ! B_1$	$B_1.\text{true} = B.\text{false}$ $B_1.\text{false} = B.\text{true}$ $B.\text{code} = B_1.\text{code}$
$B \rightarrow E_1 \text{ rel } E_2$	$B.\text{code} = E_1.\text{code} \text{    } E_2.\text{code}$

	if gen('if' E <sub>1</sub> .addr, E <sub>2</sub> .addr, 'goto' B.true)
B → true	if gen('goto' B.true)
B → false	B.code = gen('goto' B.false)

B → E<sub>1</sub> || E<sub>2</sub> is translated directly into a comparison 3-address instruction with jumps to the appropriate place.

Ex: of B of the form a < b

if a < b goto B.true

goto B.false

The remaining production for B translated as

(8) Suppose B is B<sub>1</sub> || B<sub>2</sub>. If B<sub>1</sub> is true, then B itself is true, B<sub>1</sub>.true is same as B.true. If B<sub>1</sub> is false then B<sub>2</sub> must be evaluated, make B<sub>1</sub>.false be the label of the first

instruction in the code for B<sub>2</sub>.

(88) The translation of B<sub>1</sub> && B<sub>2</sub> is similar.

(888) No code is needed for an expression B of the form

!B<sub>1</sub>; just interchange the true & false exits of B to get the true & false exits of B,

(8888) The constants true & false translate into jumps to B.true & B.false respectively.

Ex: consider if (x < 100 || x > 200 && x != y) x = 0;

```

if x < 100 goto L2
      goto L3
L3 : if x > 200 goto L4
      ↳ goto L1
L4 : if x != y goto L2
      goto L1
L2 : x = 0

```

L1 : control flow translation of a simple-if-statement

### Avoiding Redundant Gotos

of the code ps,

```
if x > 200 goto L4
```

```
goto L1
```

L4 :

Instead, consider the instruction

```
ifFalse x > 200 goto L1
```

L4 :

This ifFalse instruction takes advantage of natural flow from one instruction to the next in sequence, so control simply "falls through" to label L4 if  $x > 200$ , thereby avoiding jump. In if-while statements we can use special label fall. The new rules for  $s \rightarrow \text{if}(B) s_1 \text{ else } s_2$  are B.true to fall and B.false =  $s_1$ .next =  $s_2$  next

$s_1.\text{code} = B_1.\text{code} \parallel S_1.\text{code}$

The semantic rules for  $B \rightarrow E_1$  and  $E_2$  PS

$\text{test} = E_1.\text{addr} \text{ rel.op } E_2.\text{addr}$

$s = \text{if } B.\text{true} \neq \text{fall} \text{ and } B.\text{false} \neq \text{fall} \text{ then}$

$\text{gen}(\text{'if' test 'goto' } B.\text{true}) \parallel \text{gen}(\text{'goto' } B.\text{false})$

$\text{else if } B.\text{true} \neq \text{fall} \text{ then } \text{gen}(\text{'if' test 'goto' } B.\text{true})$

$\text{else if } B.\text{false} \neq \text{fall} \text{ then } \text{gen}(\text{'ffFalse' test 'goto' } B.\text{false})$

$\text{else }$

$B.\text{code} = E_1.\text{code} \parallel E_2.\text{code} \parallel s$

If both  $B.\text{true}$  &  $B.\text{false}$  are explicit labels P.e neither

of both  $B.\text{true}$  &  $B.\text{false}$  must be fall

fall.

If  $B.\text{true}$  is an explicit label, then  $B.\text{false}$  must be fall

so they create an if instruction that lets control fall through

if the condition is false.

If  $B.\text{false}$  is an explicit label then they generate

an ffFalse instruction.

The semantic rule for  $B \rightarrow B_1 \parallel B_2 \parallel s$ ,

$B_1.\text{true} = \text{if } B.\text{true} \neq \text{fall} \text{ then } B.\text{true} \text{ else newlabel() }$

$B_1.\text{false} = \text{fall}$

$B_2.\text{false} = B.\text{true}$

$B_2.\text{false} = B.\text{false}$

$B.\text{code} = \text{if } B.\text{true} \neq \text{fall} \text{ then } B_1.\text{code} \parallel B_2.\text{code}$

$\text{else } B_1.\text{code} \parallel B_2.\text{code} \parallel \text{label}(B_1.\text{true})$

Here note that the meaning of label fall for  $B_1$  is different from its meaning for  $B_2$ . Suppose  $B_1$ .true is fall & e control flows falls through  $B_1$ , if  $B_1$  evaluates to true. Although  $B_1$  evaluates to true if  $B_1$  does,  $B_1$ .true must ensure that control jumps over the code for  $B_1$  & to get to the next instruction after  $B_1$ .

On the other hand, if  $B_1$  evaluates to false, the truth value of  $B$  is determined by the value of  $B_2$ .

Ex: If ( $x < 100 \text{ || } x > 200 \text{ && } x \neq y$ )  $x = 0$ ;

translation into:

If  $x < 100$  goto  $L_2$

else  $x > 200$  goto  $L_1$

else  $x \neq y$  goto  $L_1$

$L_2 : x = 0$

$L_1 :$

Boolean values and Jumping code.

A Boolean expression may also be evaluated for its value, as in assignment statement such as  $x = \text{true}$  or  $x = a < b$ ;

A clean way of handling both roles of Boolean expressions is to first build a syntax tree for expression, using either of the following approaches.

(a) use 2 passes. Construct a complete syntax tree for the input & then walk the tree in depth first order, computing

the translations specified by the semantic rules.

(ii) Use one pass for statements, but 2 passes for expressions.  
With this approach, we would translate  $E \rightarrow \text{while}(E) S$ , before  $S$  is examined.

Consider

$$S \rightarrow Pd = E ; | \text{if } (E) S | \text{while } (E) S | SS$$

$$E \rightarrow E11E | E\&E | E\text{rel } E | E+E | (E) | Pd | \text{true} | \text{false}$$

Nonterminal  $E$  governs the flow of control. In  
 $S \rightarrow \text{while}(E) S_1$ , the same  $E$  denotes a value  $S \rightarrow Pd = E$   
and  $E \rightarrow E+E$ .

Suppose  $E.n$  denotes the syntax tree node for an expression  $E$ . Let method jump generate jumping code at an expression node  $E$  & rvalue generate code to compute the value of the node  $E$  onto a temporary.

When  $E$  appears in  $S \rightarrow \text{while}(E) S_1$ , method jump is called at node  $E.n$ . Jumping code is generated by calling  $E.n.jump(t, f)$  where  $t$  is a new label for  $S_1$ . code  $E.f$  is the label  $S_{1ext}$ .

When  $E$  appears in  $S \rightarrow Pd = E ;$  method rvalue is called at node  $E.n$ . If  $E$  has the form  $E+E$  the method call  $E.n.rvalue()$  generates code. If  $E$  has the form  $E.E&E$  we first generate jumping code for  $E$  & then assign true or false to a new temporary  $t$  at the true & false exits respectively from the jumping code.

Ex: the assignment  $sc = a < b \&& c < d$  can be implemented by the code

if False a < b goto L1

if False c < d goto L1

t = true

goto L2

L1 : t = false

L2 : sc = t.

## Backpatching

A key problem when generating code for boolean expressions & flow of control statement is that matching a jump instruction with the target of the jump.

This can be solved by an approach, called backpatching, in which lists of jumps are passed as synthesized attributes. Specifically when a jump is generated, the target of the jump is temporarily left unspecified. Each such jump is put on a list of jumps whose labels are to be filled in when the proper label can be determined. All the jumps on a list have the same target label.

## One-pass code Generation Using Backpatching

Backpatching can be used to generate code for boolean expressions & flow of control statements in one pass.

Synthesized attributes trueList & falseList of nonterminal B are used to manage labels in jumping code for boolean expression.

B.trueList → List of jump or conditional jump instructions into which we must insert the label to which control goes if B is true.

B.falseList → List of instructions that eventually get the label to which control goes when B is false.

As code is generated for B, jumps to the true & false exit are left incomplete, with the label field unfilled.

S has a synthesized attribute S.nextList, denoting a list to jumps to the instruction immediately following the code for S.

To manipulate lists of jumps, we use 3 functions

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

(ii) merge (P<sub>1</sub>, P<sub>2</sub>) concatenates the lists pointed to by P<sub>1</sub>, P<sub>2</sub> & returns a pointer to the concatenated list.

(iii) backpatch (P, E) inserts E as the target label for each of the instruction on the list pointed to by P.

### Backpatching for Boolean Expressions

We can construct a translation scheme suitable for generating code for boolean expression during bottom-up parsing. A marker nonterminal M in the grammar causes a semantic action to pick up, at appropriate time, the index of the next instruction to be generated. The grammar is,

$$B \rightarrow B_1 \mid MB_2 \mid B, \& MB_2 \mid !B_1 \mid (B_1) \mid E \text{ rel } E_2 \mid \text{true} \mid \text{false}$$

$M \rightarrow C$

The translation scheme is,

- (P)  $B \rightarrow B_1 \parallel M B_2$     { Backpatch ( $B_1.\text{falseList}$ ,  $M.\text{instr}$ );  
      $B.\text{trueList} = \text{merge}(B_1.\text{trueList}, B_2.\text{trueList})$ ;  
      $B.\text{falseList} = B_2.\text{falseList}; \}$
- (PP)  $B \rightarrow B_1 \wedge M B_2$     { Backpatch ( $B_1.\text{trueList}$ ,  $M.\text{instr}$ );  
      $B.\text{trueList} = B_2.\text{trueList};$   
      $B.\text{falseList} = \text{merge}(B_1.\text{falseList}, B_2.\text{falseList});$
- (PPP)  $B \rightarrow ! B_1$             {  $B.\text{trueList} = B_1.\text{falseList};$   
      $B.\text{falseList} = B_1.\text{trueList}; \}$
- (IV)  $B \rightarrow (B_1)$             {  $B.\text{trueList} = B_1.\text{trueList};$   
      $B.\text{falseList} = B_1.\text{falseList}; \}$
- (V)  $B \rightarrow E_1 \text{ rel } E_2$     {  $B.\text{trueList} = \text{makelist}(\text{nextInstr});$   
      $B.\text{falseList} = \text{makelist}(\text{nextInstr} + 1);$   
     gen ('ef'  $E_1.\text{addr}$  rel op  $E_2.\text{addr}' \text{ goto } -1$ );  
     gen ('goto -1'); }
- (VI)  $B \rightarrow \text{true}$             {  $B.\text{trueList} = \text{makelist}(\text{nextInstr});$   
     gen ('goto -1'); }
- (VII)  $B \rightarrow \text{false}$             {  $B.\text{falseList} = \text{makelist}(\text{nextInstr});$   
     gen ('goto -1'); }
- (VIII)  $M \rightarrow C$             {  $M.\text{instr} = \text{nextInstr}; \}$

Consider (e)  $B \rightarrow B_1 \parallel M B_2$ . If  $B_1$  is true, then  $B$  is also true, so  $B_1.\text{trueList}$  becomes part of  $B.\text{trueList}$ . If  $B_1$  is false, test  $B_2$ , so the target for the jumps  $B_1.\text{falseList}$  must be beginning of the code generated for  $B_2$ . This target is obtained by  $M$ . It produces synthesized attribute  $M.\text{nextInstr}$ , the index of the next instruction, just before  $B_2$  code starts being generated. To obtain that instruction index, we associate with the production  $M \rightarrow G$

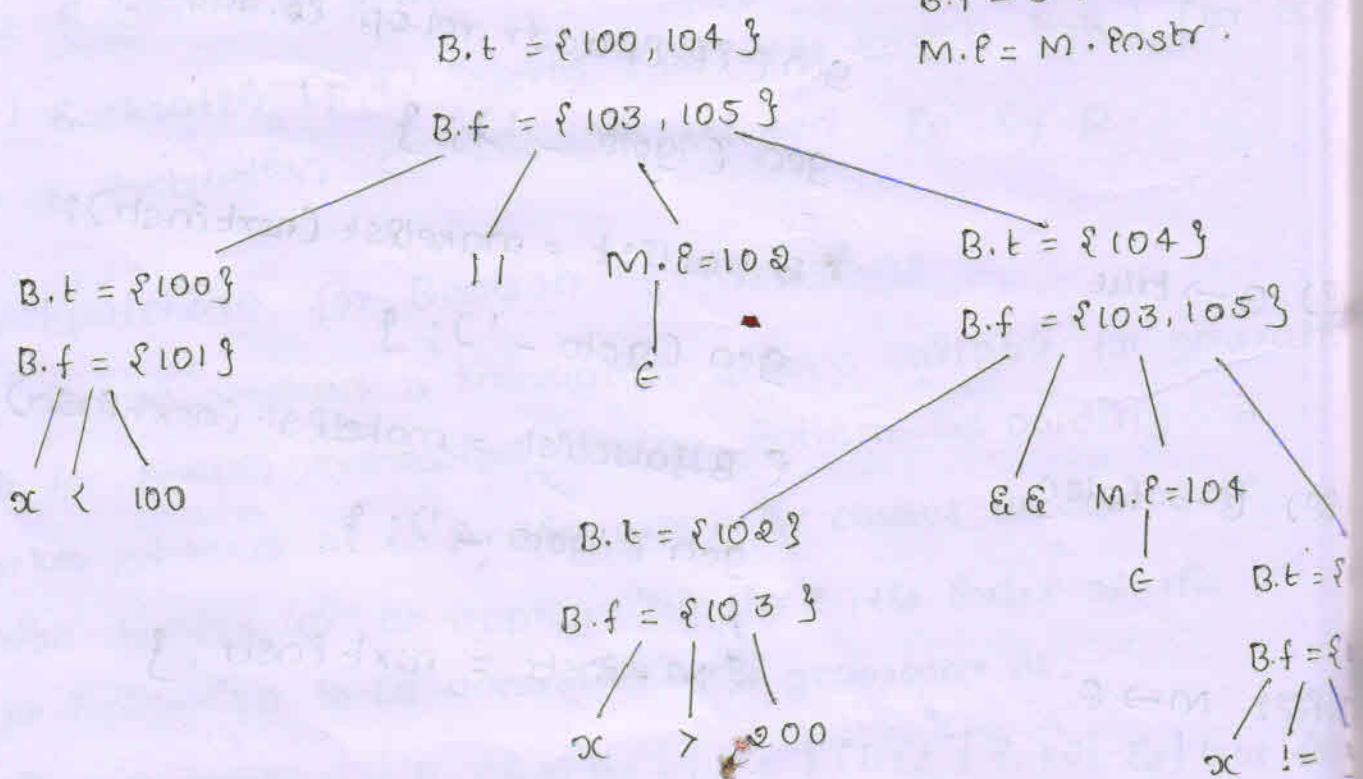
$$\{ M.\text{nextInstr} = \text{nextInstr}; \}$$

The variable  $\text{nextInstr}$  holds the index of the next instruction to follow. This value will be backpatched onto  $B_1.\text{falseList}$  when we have seen the remainder of the production  $B \rightarrow B_1 \parallel M$ .

Example:

Consider the expression  $x < 100 \parallel x > 200 \ \&\& \ x != y$

The annotated parse tree is,



$$\begin{aligned} B.t &= B.\text{trueList} \\ B.f &= B.\text{falseList} \\ M.p &= M.\text{nextInstr} \end{aligned}$$

The reduction of  $x < 100$  to B by production (v), the 2 instructions are generated

100: If  $x < 100$  goto —

101: goto —

The marker nonterminal M in the production

$B \rightarrow B_1 \sqcup M B_2$

records the value of nextInstr, which at this time is 102. The reduction of  $x > 200$  to B by production (v) generates the instructions

102: If  $x > 200$  goto —

103: goto —

The sub expression  $x > 200$  corresponds to  $B_1$  in the production

$B \rightarrow B_1 \& M B_2$

The marker nonterminal M records the current value of nextInstr, which is now 104. Reducing  $x! = y$  into B by production (v) generates

104: If  $x! = y$  goto —

105: goto —

We now reduce by  $B \rightarrow B_1 \& M B_2$ . The corresponding semantic action calls backpatch ( $B_1.\text{trueList}$ , M.instr) to bind the true exit of  $B_1$  to the first instruction of  $B_2$ . Since  $B_1.\text{trueList}$  is {100} & M.instr is 104 this call to backpatch fills pos 104 in instruction 102. The SPX instructions generated so far are

100: If  $x < 100$  goto —

101 : goto —

102 : if  $x > 200$  goto 104

103 : goto —

104 : if  $x != y$  goto —

105 : goto —

The semantic action associated with the final reduction by  $B \rightarrow B_1 \cup B_2$  calls backpatch ({101, 102}) which leaves the instruction as

100 : if  $x < 100$  goto —

101 : goto 102

102 : if  $x > 200$  goto 104

103 : goto —

104 : if  $x != y$  goto —

105 : goto —

The entire expression is true iff goto's of instructions 100 or 104 are reached, & is false iff goto's of instructions 101, 103 or 105 are reached.

### Flow of control statements

We can use backpatching to translate flow of control statements in one pass. Consider

$S \rightarrow \text{if } (B) S_1 \text{ else } S_2 \mid \text{while } (B) S \mid \text{for } L \mid A$

$L \rightarrow LS \mid S$

where  $S \rightarrow \text{statement}$   $L \rightarrow \text{statement list}$

$A \rightarrow \text{assignment statement}$   $B \rightarrow \text{Boolean expression}$

The translation scheme is,

- (8)  $S \rightarrow Pf(B) M S_1$  { Backpatch (B.trueList, M.instr);  
 $S.nextList = \text{merge} (B.falseList, S_1.nextList); }$
- (88)  $S \rightarrow Ef(B) M S_1 N \text{ else } M_2 S_2$   
{ Backpatch (B.trueList, M.instr);  
Backpatch (B.falseList, M\_2.instr);  
 $\text{temp} = \text{merge} (S_1.nextList, N.nextList);$   
 $S.nextList = \text{merge} (\text{temp}, S_2.nextList); }$
- (888)  $S \rightarrow \text{while } M_1 (B) M_2 S_1$   
{ Backpatch (S\_1.nextList, M\_1.instr);  
Backpatch (B.trueList, M\_2.instr);  
 $S.nextList = B.falseList;$   
 $\text{gen} ('goto' M_2.instr); }$   
 $S.nextList = L.nextList; }$
- (8v)  $S \rightarrow \{ L \}$  {  $S.nextList = \text{null}; }$
- (v)  $S \rightarrow A;$  {  $S.nextList = \text{null}; }$
- (v8)  $M \rightarrow E$  {  $M.instr = \text{nextInstr}; }$
- (v88)  $N \rightarrow E$  {  $N.nextList = \text{makelist} (auxInstr);$   
 $\text{gen} ('goto -'); }$
- (v888)  $L \rightarrow L_1 M S$  { Backpatch (L\_1.nextList, M.instr);  
 $L.nextList = S.nextList; }$
- (8x)  $L \rightarrow S$  {  $L.nextList = S.nextList; }$

Consider the production  $S \rightarrow \text{while } M_1 (B) M_2 S_1$ . The occurrences of marker nonterminal  $M$  in the production need the instruction numbers of the beginning of the code for  $B$ , the beginning of the code for  $S_1$ , the production  $M \rightarrow e$  set attribute  $M.e.nstr$  to the no. of next instruction. After  $S_1$  is executed, control flows to the beginning. Therefore, when we reduce to  $S$ , we backpatch  $S_1.\text{nextlist}$  to make all targets on that list be  $M_1.\text{instr}$ . An explicit jump to the beginning of the code for  $B$  is appended after the code for  $S_1$  because control may also "fall out the bottom".  $B.\text{true.list}$  is backpatched to go to the beginning of  $S_1$  by making jumps on  $B.\text{true}$  go to  $M_2.\text{instr}$ .

$S.\text{nextlist}$  &  $L.\text{nextlist}$  comes when code is generated for  $\text{if } (B) S_1 \text{ else } S_2$ . If control "falls out the bottom" of  $S_1$  as when  $S_1$  is an assignment, we must include at the end the code for  $S_1$ , a jump over the code for  $S_2$ . We use another marker nonterminal to generate this jump after  $S_1$ . i.e.  $N$ .

We backpatch the jumps when  $B$  is true to  $M_1.\text{instr}$  the latter is the beginning of the code for  $S_1$ . Similarly we backpatch jumps when  $B$  is false to go to the beginning of the code for  $S_2$ . The list  $S.\text{nextlist}$  includes all jumps out of  $S_1$  as well as the jump generated by  $N$ .

Semantic actions ( $\text{v}(\text{ee})$  &  $(\text{ex})$ ) handles sequence of statements.

$L \rightarrow L_1 M S$

In the instruction following the code for  $L_1$  in order of execution is the beginning of  $S$ . Thus the  $L_1.\text{nextlist}$  list is backpatched to the beginning of the code for  $S$ , which is given by  $M.\text{instr}$ .

On  $L \rightarrow S$ , L.nextList & same as S.nextList.

Break-, Continue - & Goto-Statements.

Goto statements can be implemented by maintaining a list of unfilled jumps for each label and then backpatching the target when it is known.

The break statements sends control out of an enclosing loop & continue-statement triggers the next iteration of an enclosing loop.

Consider

(e) for(;;readch()) {  
(ee) if (peek == '\n') continue;  
(eee) else peek = '\n'; i = i + 1;  
(ev) else break;  
(v) }

ee. Break statement on line (ev) make control to jump to the next statement after the enclosing for loop.

the next statement after the enclosing for loop construct, then a break statement continues on line (ee) makes control to jump to readch()

If S is the enclosing loop construct, then a break statement is a jump to the first instruction after the code for S. We can generate code for the break

(e) keeping track of the enclosing loop statement S

(ee) Generating an unfilled jump for the break statement

(eee) Putting this unfilled jump on S.nextList

In 2 pass front end  $s.\text{nextList}$  can be implemented as a field in the node for  $s$ . We can keep track of  $s$  by using the symbol table to map a special identifier `break` to the  $s$  for the enclosing loop statement  $s$ .

Instead of using the symbol table we can put pointer to  $s.\text{nextList}$  in the symbol table. When a break statement is reached, we generate an unfilled jump, look up `nextList` through the symbol table & add the jump to the list, where it will be backpatched.

Continue statement can be handled in a manner analogous to the break statement. The main difference being the  $\&$  is that the target of the generated jump is different.

### \* Switch - Statements

The switch statement syntax is,

`switch (E) {`

`case V1 : S1`

`case V2 : S2`

`...`

`case Vn-1 : Sn-1`

`default : Sn`

`}`

There is a selector expression  $E$ , which is to be evaluated followed by  $n$  constant values  $V_1, V_2 \dots V_{n-1}$  that the expression might take, including a default "value" which always matches the expression if no other value does.

## Translation of switch-statements.

The intended translation of a switch is code to

(1) evaluate the expression  $E$

(2) find the value  $V_E$  in the list of cases that is the same as the value of the expression.

(3) execute the statement  $S_j$  associated with the value found.

Step (2) is an n-way branch, if the no. of cases is small then it is reasonable to use a sequence of conditional jumps, each tests for an individual value & transfer to the code.

A compact way to implement this is to create a table of

pairs, each consist of value & label for statement's code. The value of the expression itself, paired with label for the default statement. If no match found, then last entry is sure to match.

If no. of values is more, it is more efficient to construct a hash table, with the labels of various statements as entries. If no entry for value possessed by the switch is found, a jump to the default statement is generated.

There is a common special case that can be implemented even more efficiently than by an n-way branch. If the values all lie in some small range say  $\min$  to  $\max$  & no. of diff. values is a reasonable fraction of  $\max - \min$ , then we can construct an array of  $\max - \min$  "buckets", where bucket  $j - \min$  contains the label of the statement with value  $j$ ; any bucket that would otherwise remain unfilled contains the default label.

To perform switch, evaluate the expression to obtain  $t$  & check that  $t$  is in the range min to max & transfer  $t$  directly to the table entry at offset  $f - m[t]$ .

### Syntax-Directed Translation of Switch-Statement.

The translation of a switch statement is,

goto test

L<sub>1</sub> : code for S<sub>1</sub>

goto next

L<sub>2</sub> : code for S<sub>2</sub>

goto next

...

L<sub>n-1</sub> : code for S<sub>n-1</sub>

goto next

L<sub>n</sub> : code for S<sub>n</sub>

goto next

test : if  $t = v_1$  goto L<sub>1</sub>

if  $t = v_2$  goto L<sub>2</sub>

...

if  $t = v_{n-1}$  goto L<sub>n-1</sub>

goto L<sub>n</sub>

Next:

When we see the keyword switch, we generate & now labels test & next, and a new temporary  $t$ . Then, as we parse the expression  $E$ , we generate code to evaluate  $E$  into  $t$ . After processing  $E$ , we generate the jump goto test.

When we see each case keyword, we create a new label  $L_E$  & enter it into the symbol table. We place it on a queue, used only to store cases, a value-label pair consisting of the value  $v_E$  of the case constant &  $L_E$ . We process each statement case  $v_E : S_E$  by emitting the label  $L_E$  attached to the code for  $S_E$  followed by the jump goto next.

When the end of the switch is found, we are ready to generate the code for the n-way branch. Reading the queue of value-label pairs, we can generate a sequence of 3-address statements of the form.

case  $t = v_1 L_1$

case  $t = v_2 L_2$

...

case  $t = v_{n-1} L_{n-1}$

case  $t = v_n L_n$

next:

where  $t$  is the temporary holding the value of the selector expression &  $L_n$  is the label for the default statement.

The case  $t v_E L_E$  instruction is a synonym for  
if  $t = v_E$  goto  $L_E$ .

Intermediate code for Procedures.

In 3-address code, a function call is unraveled into the evaluation of parameters in preparation for a call, followed by the call itself.

Ex: a is an array of integer & that f is a function from integers to integers. Then the assignment

$$n = f(a[8]);$$

might translate into the following 3-address code

$$(e) t_1 = 8 * 4$$

$$(ee) t_2 = a[t_1]$$

(eee) param t\_2

$$(ev) t_3 = \text{call } f, t_2$$

$$(v) n = t_3$$

The first 2 lines compute the value of expression a[8] into t<sub>2</sub>, Line (eee) makes t<sub>2</sub> an actual parameter for the call of f on line (ev). That line also assigns the return value to temporary t<sub>3</sub>. Line (v) assigns the result of f(a[8]) to n.

The below productions allow function definitions & function calls.

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

$$F \rightarrow E | T \text{ fd}, F$$

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

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

$$A \rightarrow E | E, A$$

- \* Non terminal D & T generate declarations and type
- \* D consists of keyword define, a return type, the function name, formal parameters in parentheses & function body

consisting of a bracketed statement.

\* F generates zero or more formal parameters

\* S & E generate statements & expressions.

\* S adds a statement that returns the value of an expression

\* E adds function call with actual parameter generated by A

\* Actual parameter is an expression.

Function types: the type of a function must encode the

return type and types of the formal parameters.

Symbol tables: Let  $s$  be the top symbol table when the function definition is reached. In the production for D, after seeing define & the function name, we push  $s$  & set up a new symbol table.

Env.push (top); top = new Env (top);

Type checking: Within expressions, a function is treated like

any other operator. For ex, if  $f$  is a function with a parameter of type real, then the integer  $a$  is coerced to a real  $\alpha$  in the call  $f(a)$ .

Function calls: When generating 3-address instructions for a function call  $pd(E_1, E_2, \dots, E_n)$  it is sufficient to generate the 3-address instructions for evaluating or reducing the parameters  $E_i$  to addresses, followed by a param instruction for each parameter.

## QUESTIONS

1. Define quadruples, triples and static single assignment form.

A quadruples has 4 fields, op, arg1, arg2 & result. The op field contains an incremental code for the operation.

Ex: the quadruples for  $a = (b * -c) + (b * -c)$

$t_1 = \text{minus } c$

$t_2 = b * t_1$

$t_3 = \text{minus } c$

$t_4 = b * t_3$

$t_5 = t_2 + t_4$

$a = t_5$

	op	arg1	arg2	result
0	minus	c		$t_1$
1	*	b	$t_1$	$t_2$
2	minus	c		$t_3$
3	*	b	$t_3$	$t_4$
4	+	$t_2$	$t_4$	$t_5$
5	=	$t_5$		a

A triples has only 3 fields op, arg1 & arg2

Ex: the triples for  $a = 6 * -c + 6 * -c$

	op	arg1	arg2
0	minus	c	
1	*	b	(0)
2	minus	c	
3	*	b	(2)
4	+	(1)	(3)
5	=	a	(4)

Static single statement assignment form is an intermediate representation that facilitates certain code optimizations

Ex:

$$p = a + b$$

$$q = p - c$$

$$p = q * d$$

$$p = e - p$$

$$q = p + q$$

(a) 3-address code

$$p_1 = a + b$$

$$q_1 = p_1 - c$$

$$p_2 = q_1 * d$$

$$p_3 = e - p_2$$

$$q_2 = p_3 + q_1$$

(b) Static single assignment form.

- Q. Develop SDD to produce directed acyclic graph for an expression show the steps for constructing the DAG for the expression  $a + a * (b - c) + (b - c) * d$ .

Syntax directed definition is,

PRODUCTION	SEMANTIC RULES
(E) $E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node} ('+', E_1.\text{node}, T.\text{node})$
(EE) $E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node} ('-', E_1.\text{node}, T.\text{node})$
(ET) $E \rightarrow T$	$E.\text{node} = T.\text{node}$
(EV) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
(V) $T \rightarrow Pd$	$T.\text{node} = \text{new Leaf} (pd, pd.\text{entry})$
(VI) $T \rightarrow \text{num}$	$T.\text{node} = \text{new Leaf} (\text{num}, \text{num}.\text{val})$

Steps for constructing the DAG

(8)  $P_1 = \text{Leaf} (\&d, \text{entry}-a)$

(88)  $P_2 = \text{Leaf} (\&d, \text{entry}-a) = P_1$

(888)  $P_3 = \text{Leaf} (\&d, \text{entry}-b)$

(8v)  $P_4 = \text{Leaf} (\&d, \text{entry}-c)$

(v)  $P_5 = \text{Node} ('-', P_3, P_4)$

(v8)  $P_6 = \text{Node} ('*', P_1, P_5)$

(v88)  $P_7 = \text{Node} ('+', P_1, P_6)$

(8x) (v888)  $P_8 = \text{Leaf} (\&d, \text{entry}-b) = P_3$

(8x)  $P_9 = \text{Leaf} (\&d, \text{entry}-c) = P_4$

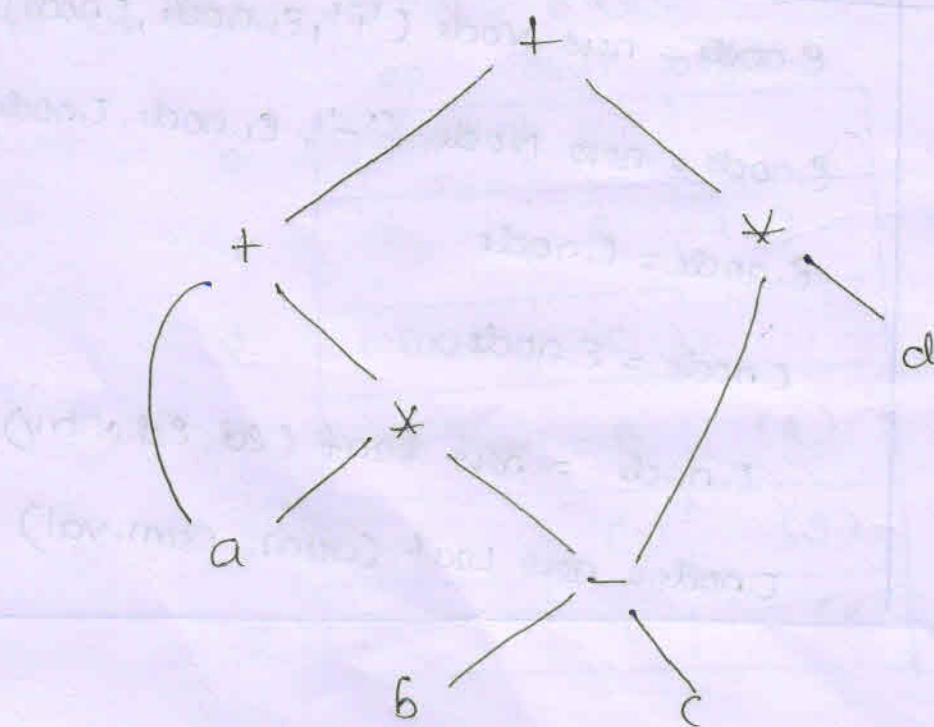
(x)  $P_{10} = \text{Node} ('-', P_3, P_4) = P_5$

(x8)  $P_{11} = \text{Leaf} (\&d, \text{entry}-d)$

(x88)  $P_{12} = \text{Node} ('*', P_5, P_{11})$

(x888)  $P_{13} = \text{Node} ('+', P_7, P_{12})$

DAG



### 3. Write syntax directed definition for flow of control statement

PRODUCTION	SEMANTIC RULES
$P \rightarrow S$	$S.\text{next} = \text{nextlabel}()$ $P.\text{code} = S.\text{code} \parallel \text{label}(S.\text{next})$
$S \rightarrow \text{assgn}$	$S.\text{code} = \text{assgn}.\text{code}$
$S \rightarrow \text{if } (B) S_1$	$B.\text{true} = \text{newlabel}()$ $B.\text{false} = S_1.\text{next} = S.\text{next}$ $S.\text{code} = B.\text{code} \parallel \text{label}(B.\text{true}) \parallel S_1.\text{code}$
$S \rightarrow \text{if } (B) S_1 \text{ else } S_2$	$B.\text{true} = \text{newlabel}()$ $B.\text{false} = \text{newlabel}()$ $S_1.\text{next} = S_2.\text{next} = S.\text{next}$ $S.\text{code} = B.\text{code}$ $\parallel \text{label}(B.\text{true}) \parallel S_1.\text{code}$ $\parallel \text{gen('goto' } S_2.\text{next})$ $\parallel \text{label}(B.\text{false}) \parallel S_2.\text{code}$
$S \rightarrow \text{while } (B) S_1$	$\text{begin} = \text{newlabel}()$ $B.\text{true} = \text{newlabel}()$ $B.\text{false} = S.\text{next}$ $S_1.\text{next} = \text{begin}$ $S.\text{code} = \text{label } (\text{begin}) \parallel B.\text{code}$ $\parallel \text{label } (B.\text{true}) \parallel S_1.\text{code}$ $\parallel \text{gen('goto' } \text{begin})$
$S \rightarrow S_1 S_2$	$S_1.\text{next} = \text{newlabel}()$

$s_2.next = s.next$

$s.code = s_1.code \parallel \text{label } (s_1.next)$

$\parallel s_2.code.$

4. Explain syntax directed translation of switch statement.

Translation of a switch statement

goto test

$L_1 : \text{code for } s_1$

goto next

$L_2 : \text{code for } s_2$

goto next.

$L_{n-1} : \text{code for } s_{n-1}$

goto next

$L_n : \text{code for } s_n$

goto next

test : if  $t = v_1$  goto  $L_1$

if  $t = v_2$  goto  $L_2$

...

if  $t = v_{n-1}$  goto  $L_{n-1}$

goto  $L_n$

next :

F&g(1)

We can write more straightforward sequence as,

if  $t_1 = v_1$  goto  $L_1$

code for  $s_1$

goto next

$L_1 : \text{if } t_1 \neq v_2 \text{ goto } L_2$

code for S<sub>2</sub>

goto next

L<sub>1</sub> : if t != v<sub>2</sub> goto L<sub>2</sub>

code for S<sub>2</sub>

goto next

L<sub>2</sub> :

...

L<sub>n-2</sub> : if t != v<sub>n-1</sub> goto L<sub>n-1</sub>

code for S<sub>n-1</sub>

goto next

L<sub>n-1</sub> : code for S<sub>n</sub>

next t :

Above code require the compiler to do extensive analyses to find the most efficient implementation. It is inconvenient in a one pass compiler.

To translate into the form of fig(1), when we see the keyword switch, we generate new labels test & next, & a new temporary variable. Then, as we parse the expression E, we generate code to evaluate E in to t. After processing E, we generate the jump goto test. When we see each case keyword, we create a new table L<sub>E</sub> & enter it into the symbol table. We place L<sub>E</sub> in a queue used only to store cases, a value-label pair consisting of the value v<sub>E</sub> of the case constant L<sub>E</sub>. We process each statement case v<sub>E</sub>: S<sub>E</sub> by emitting the label L<sub>E</sub> attached to the code for S<sub>E</sub>, followed by the jump goto next.

When the end of the switch is found, we are ready to

generate the code for the n way branch. code is

case  $t = V_1 L_1$

case  $t = V_2 L_2$

...

case  $t = V_{n-1} L_{n-1}$

case  $t = L_n$

next :

where,  $t \rightarrow$  temporary holding the value of the selector expression

$L_n \rightarrow$  label for the default statement.

case  $t = V_i L_i$  is same as if  $t = V_i$  goto  $L_i$ .