

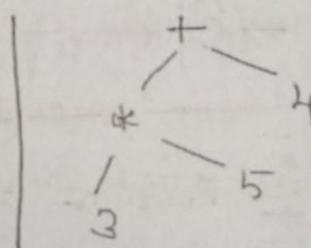
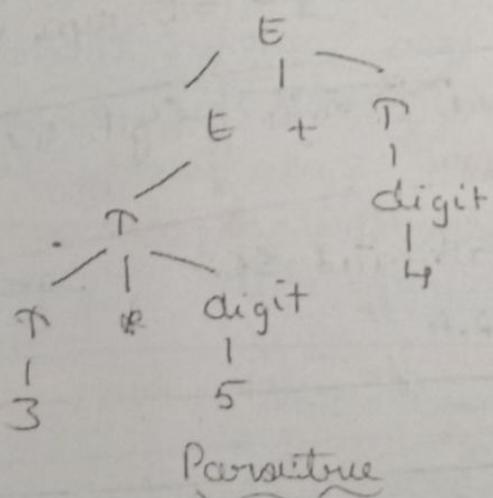
3) $L \rightarrow L, B$ 1) $L.lhs = L.lhs + (2^L.lhs_exponent * B.val)$ 2) $L.rhs = L.rhs + (2^{L.rhs_exponent} * B.val)$ 3) $L.lhs_exponent = L.lhs_exponent + 1$ 4) $L.rhs_exponent = L.rhs_exponent + 1$ 4) $L \rightarrow B$ 1) $L.lhs = R^L.lhs_exponent * B.val$ 2) $L.rhs = R^L.rhs_exponent * B.val$ 3) $L.lhs_exponent = 0$ 4) $L.rhs_exponent = -1$ 5) $B \rightarrow 0/1$ $B.val = \text{digit.lexval}$

Application Of Syntax-Directed Translation

1. Construction Of Syntax Tree

Syntax tree's are useful for construction of syntax tree.

A syntax tree is condensed form of parse tree



Syntax tree

* Syntax trees are useful for representing programming language constructs like expressions & statements.

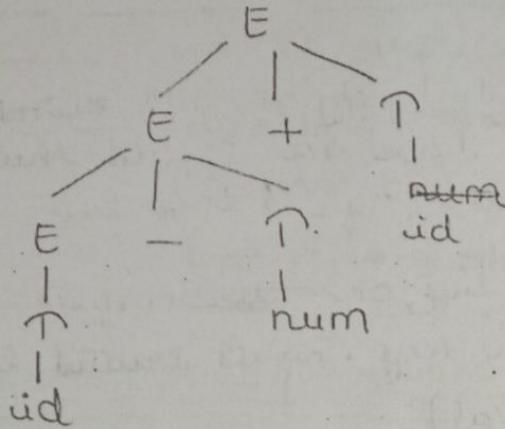
- * They help compiler design by decoupling parsing from translation.
- * Each node of a syntax tree represent a construct; the children of the node represent the meaningful components of a construct
 - Eg: A syntax tree node representing an expression $E_1 + E_2$ has label + & 2 children representing the sub expressions E_1 & E_2
- * Each node is implemented by objects with suitable no of fields; each object will have an opfield that is the label of node with additional fields as follow:
 - If the node is a leaf, an addition field holds the lexical value for the leaf. This is created by function Leaf(op, val).
 - If the node is an interior node, there are as many fields as the node has children in Syntax tree. This is created by function Node(op, c₁, c₂, ..., c_k)

Example: The S-attributed definition in fig below constructs Syntax trees for a simple expr grammar involving only binary operators + & -. As usual these operators are at the same precedence level & are jointly left associative. All nonterminals have one synthesized attr node, which represents a node of the syntax tree.

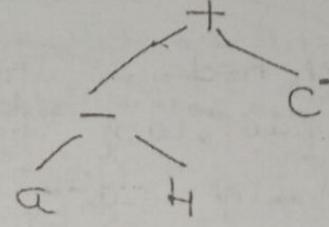
	PRODUCTION	SEMANTIC RULES	Q-H+C
1.	$E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node}('+', E_1.\text{node}, T.\text{node})$	
2.	$E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node}(' - ', E_1.\text{node}, T.\text{node})$	
3.	$E \rightarrow T$	$E.\text{node} = T.\text{node}$	

$T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
$T \rightarrow \text{id}$	$T.\text{node} = \text{newLeaf}(\text{id}, \text{id.entry})$
$T \rightarrow \text{num}$	$T.\text{node} = \text{newLeaf}(\text{num}, \text{num.val})$

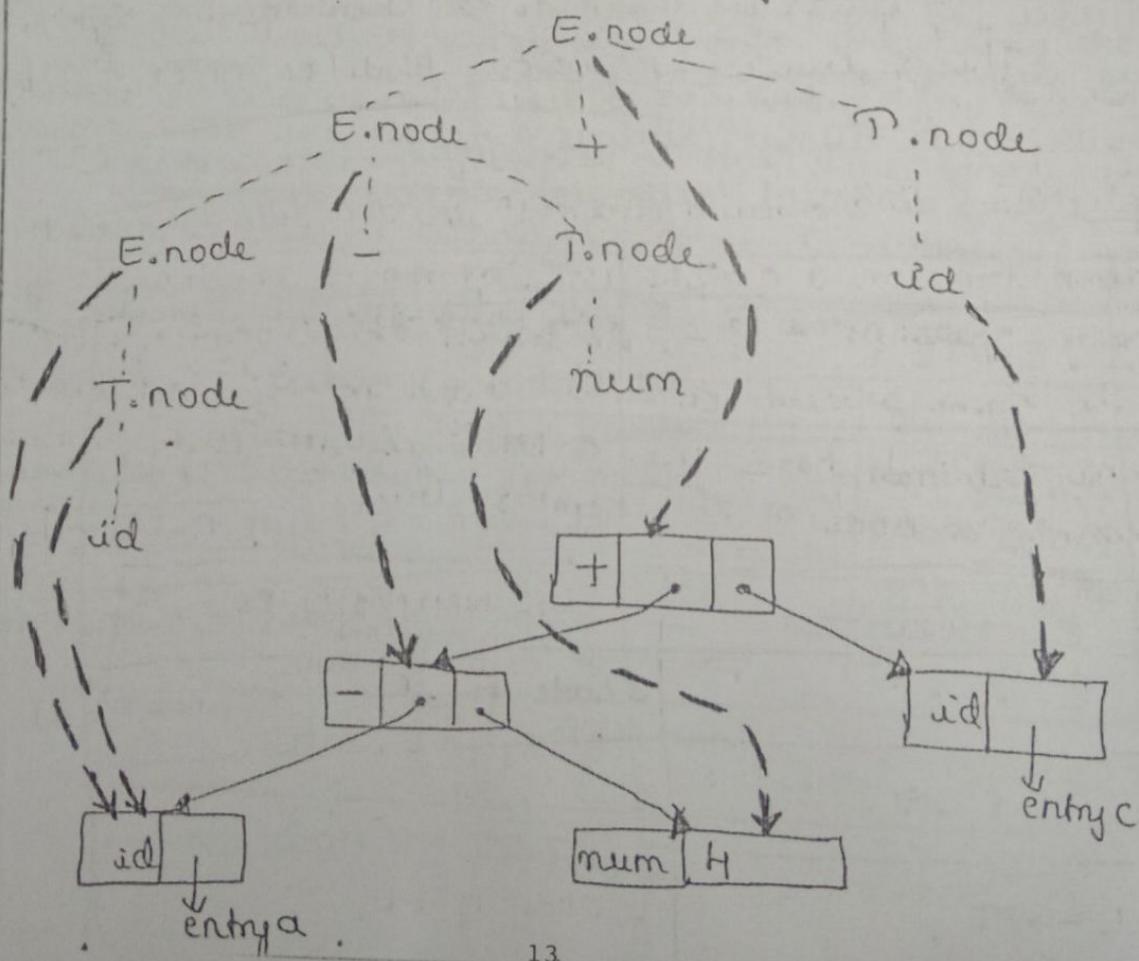
Step 2 : Parse tree



Syntax tree



Step 3 : Syntax tree for a - b + c using above SAIA.



Steps in construction of the syntax tree for $a - b + c$.

If the rules are evaluated during a post order traversal of the parse tree, or with reduceⁿ during a bottom up parse, then the sequence of steps shown below ends with P_5 pointing to the root of the constructed syntax tree.

- 1) $P_1 = \text{new Leaf } (\text{id}, \text{entry} - a)$
- 2) $P_2 = \text{new Leaf } (\text{num}, 4)$
- 3) $P_3 = \text{new Node } ('-', P_1, P_2)$
- 4) $P_4 = \text{new Leaf } (\text{id}, \text{entry} - c)$
- 5) $P_5 = \text{new Node } ('+', P_3, P_4)$

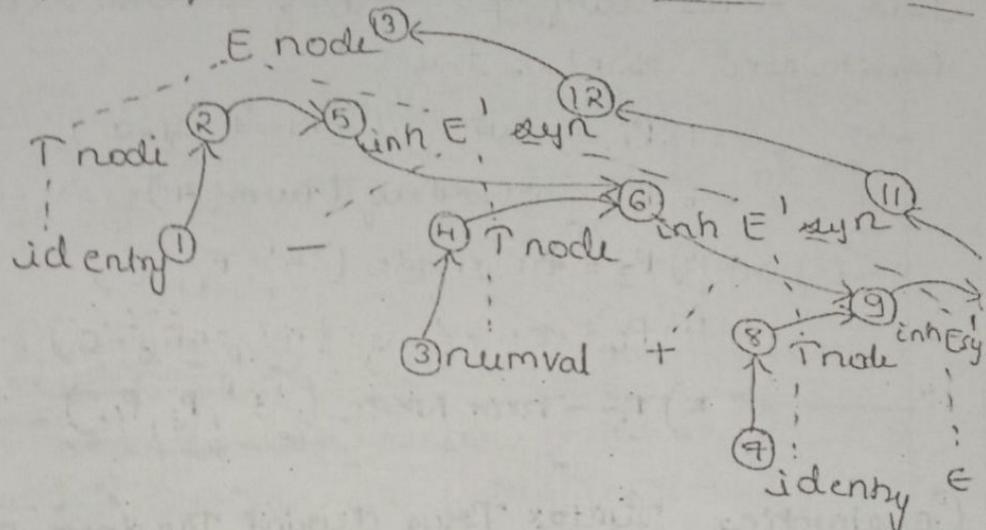
Constructing Syntax Tree during Top down parsing

With a grammar designed for top-down parsing, the syntax trees are constructed, using the same sequence of steps, even though the structure of the parse trees differs significantly from that of syntax trees. The L-attributed definⁿ below performs the same translatⁿ as the S-attributed definition shown before.

1) $E \rightarrow TE'$	$E.\text{node} = E'.\text{syn}$ $E'.\text{inh} = T.\text{node}$
2) $E' \rightarrow +TE'_1$	$E'_1.\text{inh} = E'.\text{inh} + T$ new Node(+, E'.inh, T.node) $E'_1.\text{syn} = E'_1.\text{syn}$
3) $E' \rightarrow -TE'_1$	$E'_1.\text{inh} = \text{new Node} ('-', E'.\text{inh}, T.\text{node})$ $E'_1.\text{syn} = E'_1.\text{syn}$
4) $E' \rightarrow E$	$E'_1.\text{syn} = E'_1.\text{inh}$
5) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$

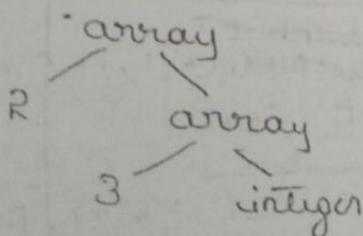
6) $T \rightarrow id$	$T.\text{node} = \text{new leaf}(\text{id}, \text{id}.\text{entry})$
7) $T \rightarrow num$	$T.\text{node} = \text{new leaf}(\text{num}, \text{num}.\text{val})$

Dependency Graph for a-H+C with Lattributed SAD



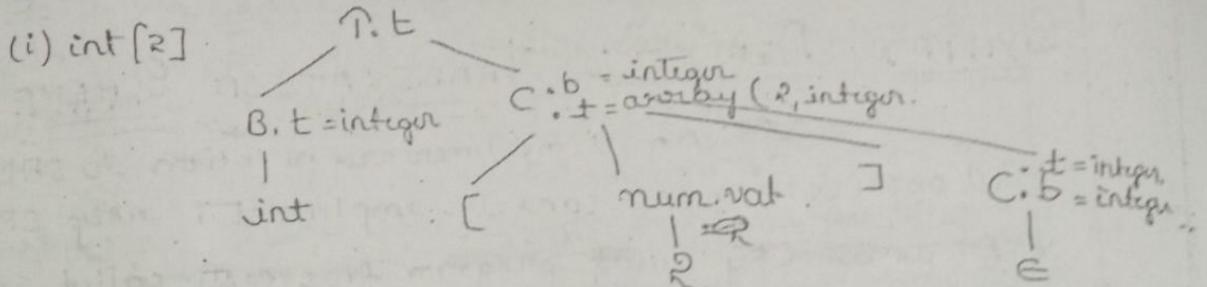
STRUCTURE OF A TYPE

This is an example of how inherited attribute can be used to carry info from one part of the parse tree to another. In C the type $\text{int}[2][3]$ can be used as "array of 2 arrays of 3 int". The corresponding type expression $\text{array}(2, \text{array}(3, \text{integer}))$ is represented by the tree as shown below.

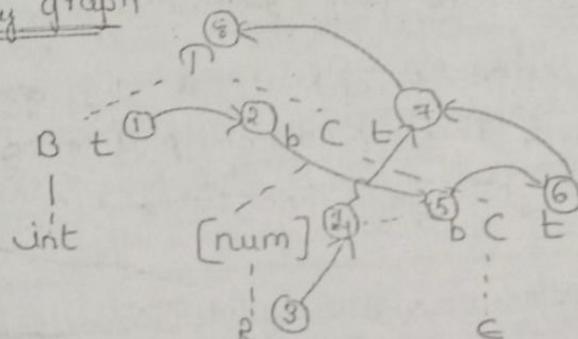


product	semantic rules
1) $T \rightarrow BC$	$T.t = C.t$ $C.b = B.t$
2) $B \rightarrow \text{int}$	$B.t = \text{integer}$
3) $B \rightarrow \text{float}$	$B.t = \text{float}$
4) $C \rightarrow [\text{num}]C_1$	$C.t = \text{array}(\text{num}.val, C_1)$ $C_1.b = C.b$
5) $C \rightarrow \epsilon$	$C.t = C_1.b$

- The non terminals B & T have a synthesized attribute representing a type
- The non terminal C has R attributes: an inherited attr (b) & a synthesized attr (t).
- The inherited attribute, b passes a basic type down the tree
- The synthesized attribute, t accumulates the result.
- In annotated parse tree for i/p: int [2]

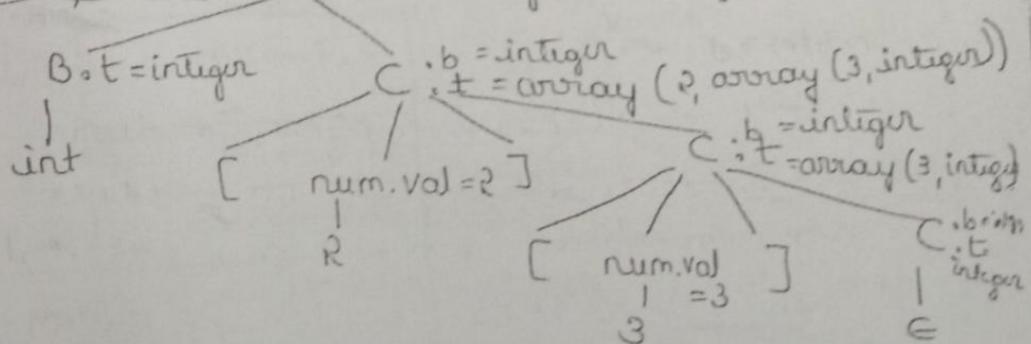


Dependency graph

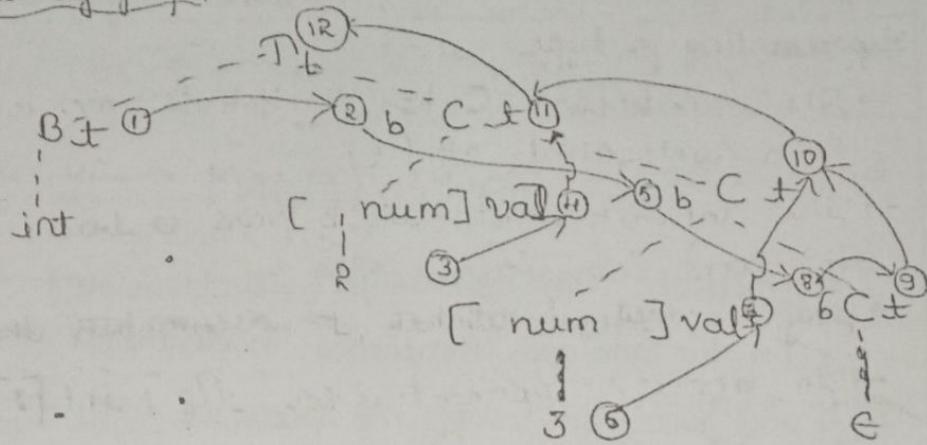


(ii) int [2][3]

Annotated parse tree: T.t = array (2, array (3, integer))



Dependency graph



SYNTAX DIRECTED TRANSLATION SCHEME :

SAT is a complementary notation to SDD.

- * All applicatⁿ of SDD can be implemented using SAT.
- * SAT is a CFG with program fragments called semantic actions embedded with production bodies.
- * Any SAT can be implemented by first building a parse tree & then performing the actions in a left to right, depth first order i-e, during pre-order traversal.
- * Typically SAT's are implemented during parsing without building parse tree. During parsing, an action in a production body is executed as soon as all the grammar symbols to the left of action have been matched.
- * SAT's that can be implemented during parsing can be characterized by introducing distinct marker non terminals in place of each embedded action.
- * Each marker M has only one production $M \rightarrow E$.
- * If grammar with marker non terminals can be parsed by a given method, then SAT can be implemented.

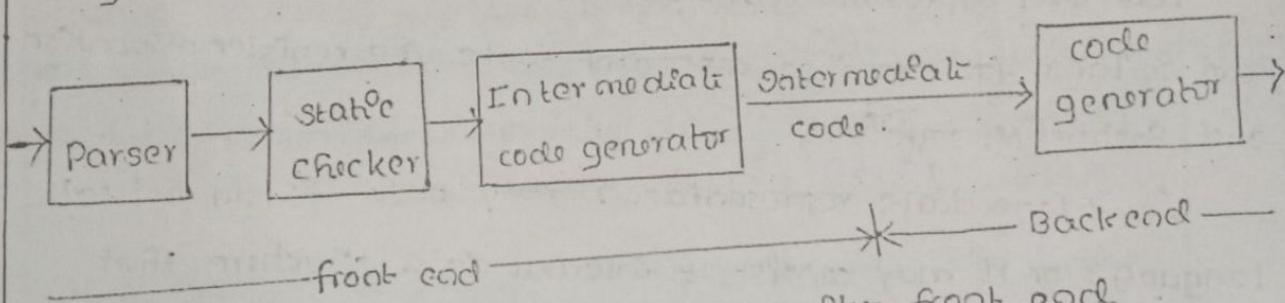
UNIT-6

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.



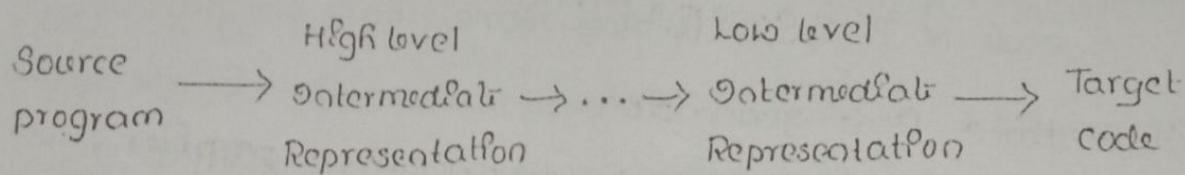
Logical structure of a compiler front end.

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 assures 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 external 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.

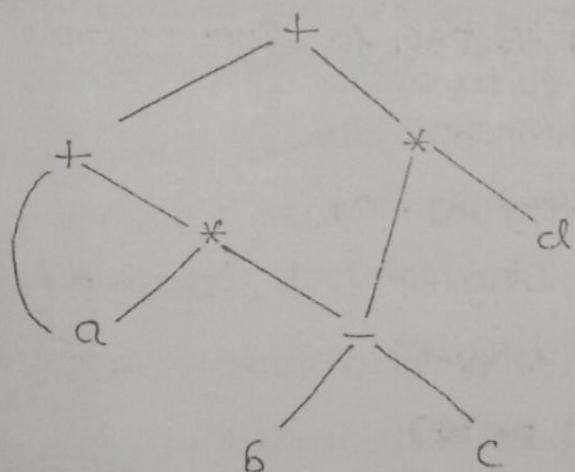
Directed Acyclic Graphs for Expressions.

On 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 expression.

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
(?) $\mathcal{E} \rightarrow \mathcal{E}_1 + \mathcal{T}$	$\mathcal{E}.\text{node} = \text{new Node}('+' , \mathcal{E}_1.\text{node}, \mathcal{T}.\text{node})$
(?) $\mathcal{E} \rightarrow \mathcal{E}_1 - \mathcal{T}$	$\mathcal{E}.\text{node} = \text{new Node}('-' , \mathcal{E}_1.\text{node}, \mathcal{T}.\text{node})$
(?) $\mathcal{E} \rightarrow \mathcal{T}$	$\mathcal{E}.\text{node} = \mathcal{T}.\text{node}$
(?) $\mathcal{T} \rightarrow (\mathcal{E})$	$\mathcal{T}.\text{node} = \mathcal{E}.\text{node}$
(?) $\mathcal{T} \rightarrow \mathcal{E}_d$	$\mathcal{T}.\text{node} = \text{new Leaf}(\mathcal{E}_d, \mathcal{E}_d.\text{entry})$
(?) $\mathcal{T} \rightarrow \text{num}$	$\mathcal{T}.\text{node} = \text{new Leaf}(\text{num}, \text{num}.val)$

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 node exists, the existing node is returned.

Steps for constructing the DAG for above example.

(?) $P_1 = \text{Leaf}(\mathcal{E}_d, \text{entry}-a)$

(?) $P_2 = \text{Leaf}(\mathcal{E}_d, \text{entry}-a) = P_1$

(?) $P_3 = \text{Leaf}(\mathcal{E}_d, \text{entry}-b)$

(?) $P_4 = \text{Leaf}(\mathcal{E}_d, \text{entry}-c)$

(?) $P_5 = \text{Node}('+' , P_3, P_4)$

(?) $P_6 = \text{Node}('x' , P_1, P_5)$

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

(?) $P_8 = \text{Leaf}(\mathcal{E}_d, \text{entry}-b) = P_3$

(?) $P_9 = \text{Leaf}(\mathcal{E}_d, \text{entry}-c) = P_4$

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

(?) $P_{11} = \text{Leaf}(\mathcal{E}_d, \text{entry}-d)$

(exle) $P_{12} = \text{Node}('x^1', P_2, P_{11})$

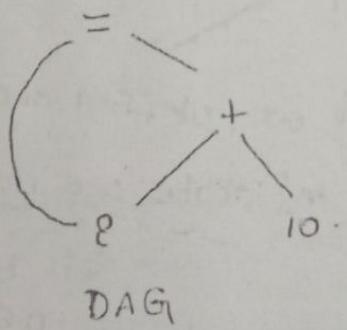
(exle) $P_{13} = \text{Node}('+', P_2, P_{12})$

When the call to Leaf (8ct, entry-a) is repeated at step 2, 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, 8ct/labeling 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

1	Pd		
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" for the node or for the expression represented by the node.
- * 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 $\langle op, l, r \rangle$ where op is the label, l is its left child's value number & r is 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 $\langle op, l, r \rangle$

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 N with label op , left child l and right child r & return its value number.

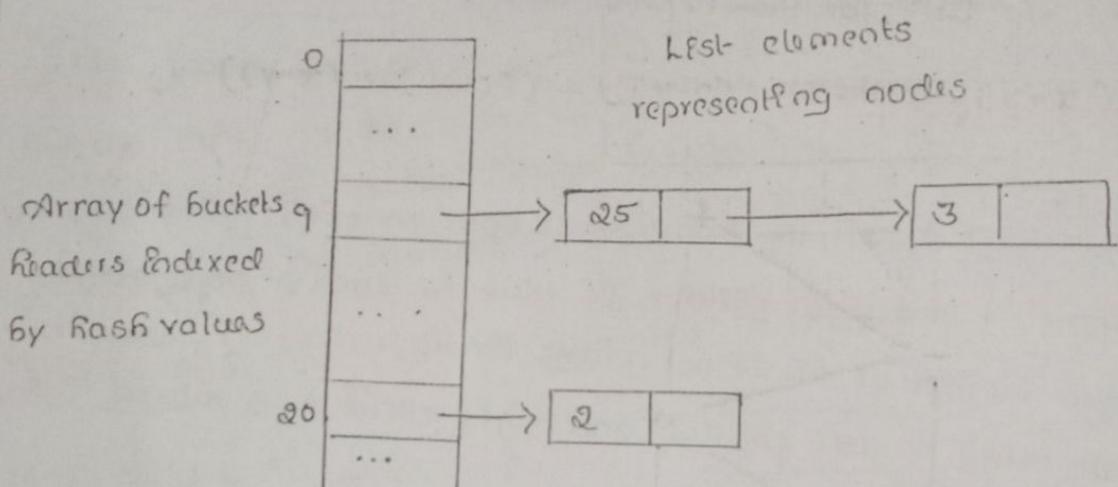
Above algorithm yields the desired output, but searching the entire array every time we are asked to locate one node is expensive.

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

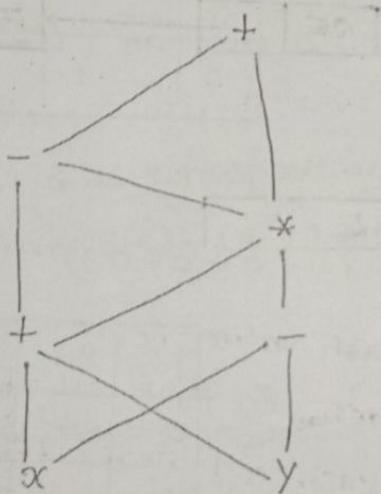
For each value number 'v' found in a cell, we must

check whether the signature $\langle \text{op}, \text{lir} \rangle$ of the Poput node matches the node with value number v in the list of the cell. 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 index $R(\text{op}, \text{lir})$ & return the value number in that new cell.

Problems.

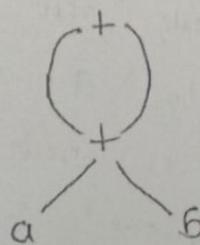
Construct the DAG for the expression.

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



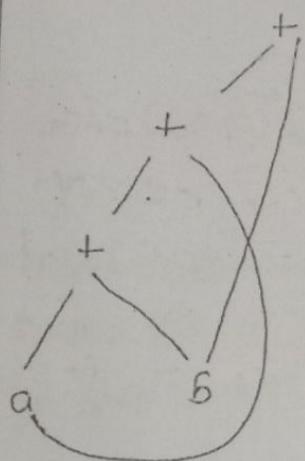
Construct the DAG & identify the value number for the sub expressions of the following expressions, assuming + associativity from the left.

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



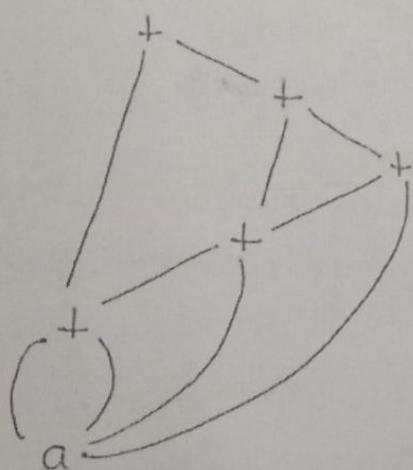
1	8d	a	
2	8d	b	
3	+	1	2
4	+	3	3

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



1	8d	a		
2	8d	b		
3	+	1	2	
4	+	3	1	
5	+	4	2	

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



1	8d	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 hand side of an instruction, i.e., 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

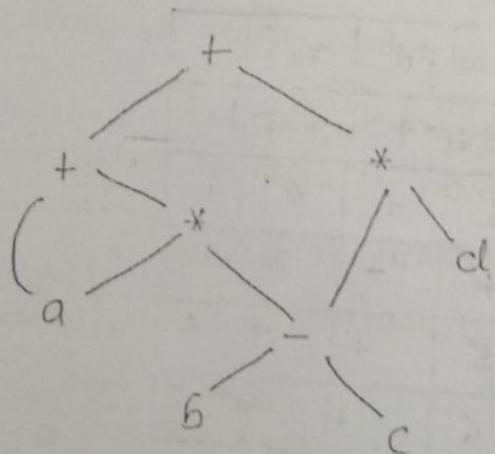
$$t_1 = y + z$$

$$t_2 = x + t_1$$

where t_1 & t_2 are compiler generated names.

3 address code is a weighted 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

$$\begin{aligned}t_1 &= b - c \\t_2 &= a * t_1 \\t_3 &= a + t_2 \\t_4 &= t_1 * d \\t_5 &= t_3 + t_4\end{aligned}$$

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.

(EE) Assignments of the form $x = \text{op } y$, where op is a unary operation

(EEE) Copy instructions of the form $x = y$, where x is assigned the value of y.

- (Pv) 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 ($<$, $=$, $>$ 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.
- (CvE) Procedures calls & returns are implemented using the following instructions: param x_1 for parameters; call p,n & $y = \text{call } p, n$ 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.

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

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

Ex % Consider the statement

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

2 possible translations of this statement are

L : $t_1 = t + 1$

100 : $t_1 = t + 1$

$t = t_1$

101 : $t = t_1$

$t_2 = t * 8$

102 : $t_2 = t * 8$

$t_3 = a[t_2]$

103 : $t_3 = a[t_2]$

if $t_3 < v$ goto L . . .

104 : if $t_3 < v$ goto 100

(a) symbolic labels

(b) position number.

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 $t * 8$ is appropriate for an array of elements that each take 8 units of space.

Quadruples

A quadruples has 4 fields, which we call op, arg₁, arg₂, & result. The op field contains an internal code for the operator.

Ex % $x = y + z$ is represented by placing + in op, y in arg₁, z in arg₂ & x in result.

The following are some exceptions to this rule.

(a) Instructions with unary operators like $x = \text{minus } y$ or $x = y$ do not use arg2. Note that for a copy statement like $x = y$, op is =, while for most other operations, the assignment operator is implied.

(b) Operators like param use neither arg2 nor result.

(c) Conditional & unconditional jumps put the target label in result.

Ex: write 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) 3-address code

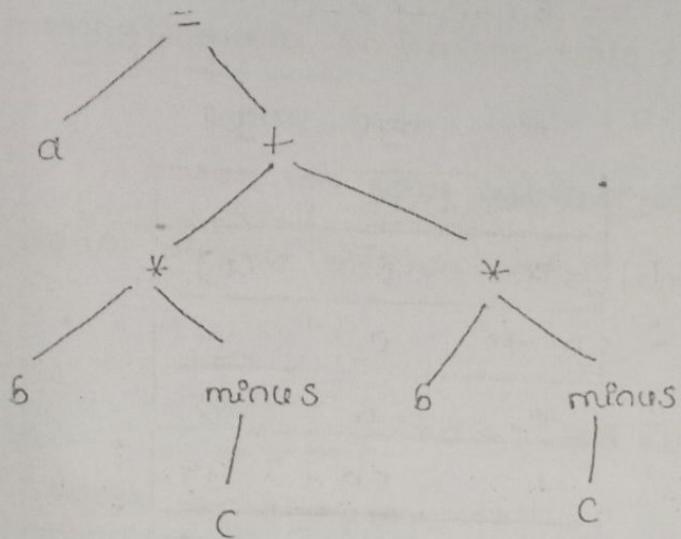
(b) Quadruples.

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

Triples

A triple has only 3 fields, op, arg1 & arg2. 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 $a = 6x - c + 6x - c$



	op	arg 1	arg 2
0	minus	' c	'
1	*	' 6	(0)
2	minus	' c	'
3	*	' 6	(2)
4	+	(1)	(8)
5	=	' a	(4)

(6) triples.

(a) Syntax tree

A ternary operator like $x[i] = y$ requires 2 entries in the triple structure, for example, we can put x & y in one triple & i 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	op	arg1	arg2
35 (0)	0 minus	c	
36 (1)	1 *	6	(0)
37 (2)	2 minus	c	
38 (3)	3 *	6	(2)
39 (4)	4 +	(1)	(3)
40 (5)	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 distinguishes SSA from 3-address code

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

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

$$p_1 = a + b$$

$$q_1 = p - c$$

$$q_{11} = p_1 - c$$

$$p = q_1 * d$$

$$p_2 = q_1 * d$$

$$p = e - p$$

$$p_3 = e - p_2$$

$$q_2 = p + q_1$$

$$q_{12} = p_3 + q_{11}$$

3-address code

Static single assignment form

The same variable may be defined in 2 different control flow paths in a program. For example, the source program

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

 $y = \text{oc}x_1;$

If we use different names for x in the true part & false then conflict arises which name should use in $y = \text{oc}x_1$. (EE) SSA uses a notational convention called ϕ -function to combine the 2 definitions of x .

```
if(flag) x1=-1; else x0=1;
```

 $x_3 = \phi(x_1, x_0);$

Here $\phi(x_1, x_0)$ has the value x_1 if the control flow passes through the true part of the conditional & the value x_0 if the control flow passes through the false part.

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

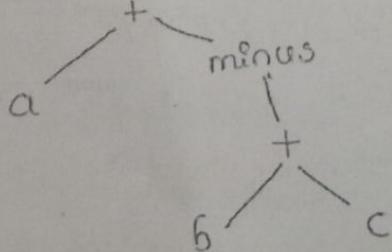
(a) A syntax tree

(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
10	+	b	c	t_1
11	minus	t_1		t_2
12	*	a	t_2	t_3

(c) triples

	op	arg1	arg2
0	+	b	$b \cdot c$
1	minus	(0)	
2	*	a	(1)

(d) Indirect triples.

	Instructions
35	(0)
36	(1)
37	(2)

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

Translation of Expressions.

An expression with more than one operator, like $a + b * c$, will translate into instructions with almost 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 S using attribute 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

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 $a = (6 * -c) + (6 * -c)$

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

$$t_2 = 6 * t_1$$

$$t_3 = \text{minus } c$$

$$t_4 = b * t_3$$

$$t_5 = t_2 + t_4$$

$$a = t_5$$

	op	arg1	arg2	result
0	minus	c		t_1
1	*	6	t_1	t_2
2	minus	c		t_3
3	*	6	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	*	6	(0)
2	minus	c	
3	*	6	(2)
4	+	(1)	(3)
5	=	a	(4)

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

Ex:

$$P = a + b$$

$$P_1 = a + b$$

$$Q = P - C$$

$$Q_1 = P_1 - C$$

$$P = Q \times D$$

$$P_2 = Q_1 \times D$$

$$P = C - P$$

$$P_3 = C - P_2$$

$$Q = P + Q$$

$$Q_2 = P_3 + Q_1$$

(a) 3-address code

(b) Static single assignment form.

2. 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) $E \rightarrow E_1 + T$	$E.\text{node} = \text{new Node} ('+', E_1.\text{node}, T.\text{node})$
(E') $E \rightarrow E_1 - T$	$E.\text{node} = \text{new Node} ('-', E_1.\text{node}, T.\text{node})$
(E'') $E \rightarrow T$	$E.\text{node} = T.\text{node}$
(Pv) $T \rightarrow (E)$	$T.\text{node} = E.\text{node}$
(Pd) $T \rightarrow Pd$	$T.\text{node} = \text{new Leaf} (pd, pd.\text{entry})$
(V) $T \rightarrow \text{num}$	$T.\text{node} = \text{new Leaf} (\text{num}, \text{num}.val)$

Steps for constructing the DAG

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

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

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

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

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

(EV) $P_6 = \text{Node}(^*\&, P_1, P_5)$

(EVPP) $P_7 = \text{Node}(^+, P_1, P_6)$

(EPX) (VEPE) $P_8 = \text{Leaf}(\&d, \text{entry}-b) = P_3$

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

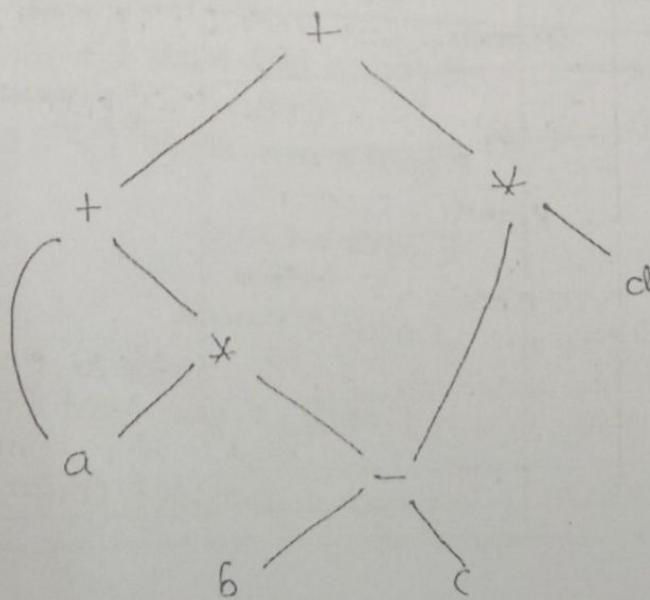
(X) $P_{10} = \text{Node}(^-, P_3, P_4) = P_5'$

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

(XPEP) $P_{12} = \text{Node}(^*\&, P_5, P_{11})$

(XPPEP) $P_{13} = \text{Node}(^+, P_7, P_{12})$

DAG

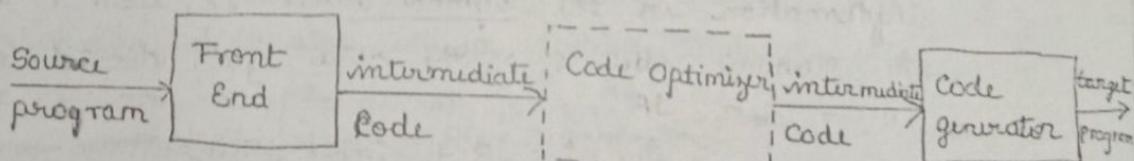


UNIT - 8

CODE GENERATION

INTRODUCTION :

- * Code generation is the final phase in the compiler design.
- * The code optimizer accepts intermediate code representation which is generated from the front end of the compiler.
 - ↳ produces another intermediate code representation which is optimized.
- * Code generator takes intermediate representation produced by code optimizer along with supplementary information in symbol table of the source program.
 - ↳ produce as output an equivalent target program.



- * Code generator has 3 main tasks:
 - 1) Instruction selection
 - 2) Register allocation & assignment
 - 3) Instruction Ordering

INSTRUCTION SELECTION :

Choose an appropriate target machine code instructions to implement the IR [intermediate representation] statements

REGISTER ALLOCATION & ASSIGNMENT :

Decide what values to keep in which registers

3) INSTRUCTION ORDERING

Decide in what order to schedule the execution of instructions.

8.1

ISSUES IN THE DESIGN OF CODE GENERATOR:

- 1) Input to the code generator
- 2) The target program
- 3) Instruction selection
- 4) Register allocation
- 5) Evaluation Order

1) Input to the code generator

* Input to the code generator is the intermediate representation of the source program produced by the front end along with information in the symbol table i.e., used to determine the runtime address of the data objects denoted by the names in IR.

< Input = IR + Symbol table >

* IR has several choices

- (a) 3-address representation : quadruples, triples, indirect triple
- (b) Virtual machine representation : byte codes of stack machine codes
- (c) Linear representation such as postfix notation
- (d) Graphical representation such as syntax tree of DAG

* Assumptions made are

(i) Front end produces low-level IR, i.e., values of names in it can be directly manipulated by the machine instruction.

(ii) Syntactic & semantic errors have been already detected

2) The Target Program:

- * The output of code generator is target program.
- * The instruction set architecture of the target machine has a significant impact on the design of code generator.
- * Most common architectures are:
 - (a) CISC: It has few registers, has maximum of 2 operands & variety of addressing mode, variable length instructions & instruction with side effects.
 - (b) RISC: It has many registers, has maximum of 3 operands with simple addressing modes, & relatively simple instruction set architecture.
- * Output may take variety of forms.
 - a) Absolute machine language [Executable code]
 - b) Relocatable machine language [object files for linker]
 - c) Assembly language [facilitates debugging]
- a) Absolute machine language has advantage that it can be placed in a fixed location in memory & immediately executed.
- b) Relocatable machine language program allows subprograms to be compiled separately.
- c) Producing Assembly language program as output makes the process of code generation somewhat easier.

3) Instruction Selection

The code generator must map the IR program into a code sequence that can be executed by the target machine.

* The complexity of performing this mapping is determined by the factors such as:

- (i) the level of the IR
- (ii) the nature of the instruction set architectures.
- (iii) the desired quality of the generated code.

(i) the levels of the IR:

- > If the [IR is high level], use code templates to translate each IR statements into a sequence of machine instruction.
- > produces poor code, needs further optimization.
- > If the [IR is low level] , use ^{low level} (code, this) information to generate more efficient code sequence.

(ii) the nature of the instruction set architectures has strong effect on difficulty of instruction select?

> Uniformity & completeness of the instruction set are imp factors.

> If we do not care about the efficiency of the target program, instruction selector is straightforward.

> For eg:

$x = y + z \Rightarrow LD R_0, y$
ADD R_0, R_0, z
ST x, R_0

∴ produces redundant LD & store

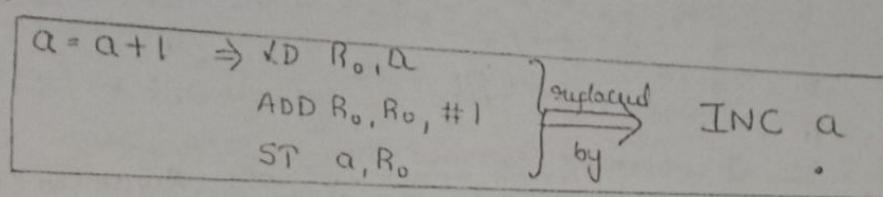
eg2:

$a = b + c \Rightarrow LD R_0, b$
$d = a + b \quad ADD R_0, R_0, c$
ST a, R_0
LD R_0, a
ADD R_0, R_0, c
ST d, R_0

→ REDUNDANT

(iii) the quality of the generated code is determined by its speed & size.

> For eg:



4) Register Allocation:

* Instrucⁿ involving register operands are usually shorter & faster than those involving operands in memory.

* 2 subproblems:

(i) Register allocation: Select the set of variables that will reside in registers at each point in the program.

(ii) Register assignment: Select specific register that a variable will reside in.

* Complications imposed by the hardware architecture
Eg: Register pairs for multiplication & division.

* Multiplication instrⁿ is of the form

$$M \boxed{x, y}$$

where $x \rightarrow$ multiplicand, is the odd register of an even/odd register pair.

$y \rightarrow$ multiplier, is the even register of a single register.
 \Rightarrow Product \rightarrow occupies the entire even/odd register pair.

* Division instrⁿ is of the form

$$D \boxed{x, y}$$

where $x \rightarrow$ dividend, occupies even register

$y \rightarrow$ divisor, occupies odd/even register

\Rightarrow Quotient \rightarrow stored in odd register
remainder \rightarrow stored in even register

Eg: two 3-address code sequences

$$t = a + b$$

$$t = t * c$$

$$t = t / d$$

$$t = a + b$$

$$t = t + c$$

$$t = t / d$$

Optimal machine-Code Sequences

L R1, a

A R1, b

M R0, c

D R0, d

ST R1, t

L R0, a

A R0, b

*A R0, c

SRDA R0, 3R

D R0, d

ST R1, t

5) Evaluation Order:

* The order in which computations are performed can effect the efficiency of the target code.

* when instrucⁿ are independent their evaluation order can be changed.

* Some computatⁿ orders require fewer registers to hold intermediate results than others.

* However picking a best order in the general case is a difficult NP-complete problem.

ADDITIONAL INFORMATION: Eg

$$t1 = a + b$$

$$a + b - (c + d) \rightarrow t2 = c + d$$

$$t3 = e * t2$$

$$t4 = t1 - t3$$

Reorder↓

$$t2 = c + d$$

$$t3 = e * t2$$

$$t1 = a + b$$

$$t4 = t1 - t3$$

MOV R0, a
ADD R0, b
MOV R1, R0
MOV R1, c
ADD R1, d
MOV R0, e
MUL R0, R1
MOV R1, t1
SUB R1, R0
MOV R1, t4, R1

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

TAE

8.2 THE TARGET LANGUAGE:

For designing a good code generator, we need to have familiarity with target machine & its instruction set. Instead of generating code on a specific target machine, a general machine consisting of many registers are considered.

A SIMPLE TARGET MACHINE MODEL:

The characteristics of target machine mode with instruction format & instruction set are shown below:

* Our hypothetical machine:

- (i) It is a 3-address machine with the following format

[OP destination, Source1, Source2]

NOTE:

A 3 address instructⁿ can have 2 operands or 1 operand also but it can have max of 3 operands.

- (ii) The target machine is byte addressable i.e., it can access 8 bit of info from specific address

- (iii) It has n no of registers denoted by $R_0, R_1, R_2, \dots, R_{n-1}$

* Various types of instructⁿ that are used by target m/c

(i) Load Instructⁿ

(ii) Store Instructⁿ

(iii) Computational Instructⁿ

(iv) Unconditional Instructⁿ

(v) Conditional Instructⁿ

- (i) Load Instructⁿ: Used to copy the data into distinctⁿ operand which must be a register.

SYNTAX: LD₄₅ dat,addr

where addr operand \rightarrow register or memory loc.

(ii) Store instruction: used to copy the data into memory location specified in the destinatⁿ operand.

SYNTAX: ST dst, or

where dst \rightarrow destination if it is a mem loc.

or \rightarrow register.

Computational operation.

(iii) Arithmetic instruction: They are performed using these instructions.

SYNTAX: OP dst, Src1, Src2

where 1st operand, dst \rightarrow destination

2nd or 3rd operand \rightarrow Operands where R values fetched for operatⁿ to be

Eg1: ADD R0, R1, R2 // $R0 = R1 + R2$

Eg2: SUB R0, R0, R1 // $R0 = R0 - R1$

Eg3: MUL R2, R0, R1 // $R2 = R0 * R1$

(iv) Unconditional Jumps: The branch instructⁿ without any conditⁿ are called unconditional jumps.

SYNTAX: BR label

where BR \rightarrow Branch instruct

(v) Conditional Jumps: Based on the value stored in a register i.e., whether it is true or zero or -ve, if branching takes place, then the branch instrⁿ are called Conditional jumps.

SYNTAX: Bcond or, label

where B stands from Branch,

Cond can be LT, GT, LTX, GTY

46 less than, greater than, less than or equal, greater than or equal

$R \rightarrow$ register, contains value such as 0, +ve or -ve.

Eg: BL RO, TI

// Branch to TI, if RO contains +ve value

Eg: BLTR RI, TR

// Branch to TR, if RI contains either 0 or -ve value

* Different addressing modes supported by generalized target machine:

1) Direct addressing mode

2) Indexed —————

3) Integer Indexed —————

4) Indirect —————

5) Immediate —————

(i) Direct A/M:

Address of the data to be accessed is directly present in the instruction, i.e., location is identified by a variable name x .

Eg: KDP LD RI, x

// Load value stored in memory locatⁿ x into RI

(ii) Indexed A/M: The data can be accessed from a memory locatⁿ using index. This addressing mode is useful for accessing arrays, where a is the base address of the array & register holds the index value

Eg: LD RI, $a(RR)$

// Accesses the data stored in
 $R_1 = \text{contents}(a + \text{contents}(RR))$

(iii) Indexed A/M where memory locatⁿ is integer

It is same as previous one except that a memory locatⁿ is identified as integer.

Eg: LD RI, 100(RR)

// $RI = \text{contents}(100 + \text{contents}(RR))$

(iv) Indirect A/M: Contents of the data can be accessed by
derefencing using * operators as shown below:

LD R1, *(R2)

// R2 contains memory location
the data stored in that
memory location is copied in
register R1.

LD R1, *100(R2)

// R1 = contents(contents(100)+
contents(R2))

(v) Immediate A/M: The data to be manipulated is
directly present in the instruction & preceded by

LD R1, #100

// R1 $\leftarrow 100$

EXERCISE:

1. Generate 3 address statement for $x = y - z$

LD R1, y

// R1 = y

LD RR, z

// RR = z

ADD R1, RR, RR

// R1 = R1 + RR

ST x, R1

// $x \leftarrow R1$

2. Generate 3 address statement $x = \alpha \oplus p$

LD R1, p

// $R1 \leftarrow p$

LD RR, O(R1)

// RR = contents(O + contents(R1))

ST x, RR

// $x = RR$

(iv) Indirect A/M: Contents of the data can be accessed by dereferencing using * operators as shown below:

LD R1, *(R2)

// R2 = contains memory locaⁿ
the data stored in that
memory locatⁿ is copied in
register R1.

LD R1, *100(R2)

// R1 = contents(contents(100+
contents(R

(v) Immediate A/M: The data to be manipulated is directly present in the instruction & preceded by

LD R1, #100

// R1 $\leftarrow 100$

EXERCISE :

1. Generate 3 address statement for $x = y - z$

LD R1, y

// R1 = y

LD RR, z

// RR = z

ADD R1, R1, RR

// R1 = R1 + RR

ST x, R1

// $x = R1$

2. Generate 3 address statement for $x = *p$

LD R1, p

// $R1 \leftarrow p$

LD RR, 0(R1)

// RR = contents(0 + conte

ST x, RR

// $x = RR$

3. Generate code for 3 address statement $*p = y$

LD R1, p // R1 = p
LD RR, y // RR = y
ST O(R1), RR // content_z(o + content_z(R1)) = RR.

4. Generate m/c code for 3 address statement $b = a[i]$

LD R1, i // R1 = i
MUL R1, R1, 8 // R1 = R1 * 8
LD RR, a[R1] // RR = content_z(a + content_z(R1))
ST b, RR // b = RR.

5. Generate m/c code for 3 address statement $a[j] = c$

LD R1, j // R1 = j
LD RR, C. // RR = C
MUL R1, R1, 8 // R1 = R1 * C
ST a[R1], RR // content_z(a + content_z(R1)) = RR

6. Generate m/c code for 3 address statement

if $x < y$ goto L

LD R1, x // R1 = x
LD RR, y // RR = y
SUB R1, R1, RR // R1 = R1 - RR
BZ R1, M // if R1 < 0 jump to M

Program & Instruction Cost

* For simplicity we take the cost of an instruction to be one plus the costs associated with the addressing modes of the operands.

- * A/M involves register have zero additional cost.
- * A/M involving memory locatⁿ or constant have additional cost of 1.

* For example:

- a) LD A0, R1 \rightarrow cost = 1
- b) RD R0, M \rightarrow cost = 2.
- c) LD RI, *100(RR) \rightarrow cost = 3

Cost of Addressing mode:

	Mode	Form	Address	Added Cost
①	Absolute direct A/M	M	M	1
②	Register direct A/M	R	R	0
③	Indexed A/M	C(R)	C+contents(R)	1
④	Indirect register A/M	*R	contents(R)	0
⑤	Indirect with indexed A/M	*C(R)	contents(C+contents(R))	1
⑥	Immediate A/M	#C	N/A	1

NOTE: Cost of each statement = 1 + cost (addressing mode)

EXERCISES (8.2)

1. Determine the cost of the following instruction sequence

LD R0, y	Cost = 1 + cost(AM)
LD RI, y	Cost = 1 + 1 = 2
ADD R0, R0, RI	Cost = 1 + 1 = 2
ST x, R0	Cost = 1 + 0 = 1
	Cost = 1 + 1 = 2
	<u>Total Cost = 7</u>

2. LD R0, i

MUL R0, R0, 8

LD RI, a(R0)

ST b, RI

LD R0, i	Cost = Cost (AN) + 1.
MUL R0, R0, 8	Cost = 1 + 1 = 2
LD RI, a(R0)	Cost = 1 + 1 = 2
ST b, RI	Cost = 1 + 1 = 2
	<u>Total Cost = 8</u>

3. LD R0, C

LD RI, i

MUL RI, RI, 8

ST a(RI), R0

LD R0, C	Cost = cost(AM) + 1
LD RI, i	1 + 1 = 2
MUL RI, RI, 8	1 + 1 = 2
ST a(RI), R0	1 + 1 = 2
	<u>Total Cost = 8</u>

4. LD R0, P
LD RI, O(R0)
ST X, RI

		$\text{Cost} = \text{Cost(A.M)} + 1$
LD	R0, P	$1+1=R$
LD	RI, O(R0)	$1+1=R$
ST	X, RI	$1+1=R$
TOTAL		$\text{Cost} = 6$

5. LD R0, P
LD RI, X
ST O(R0), RI

		$\text{Cost} = \text{Cost(A.M)} + 1$
LD	R0, P	$1+1=R$
LD	RI, X	$1+1=R$
ST	O(RI), RI	$1+1=R$
TOTAL		$\text{Cost} = 6$

6. LD R0, X
LD RI, Y
SUB R0, R0, RI
BLTR #R3, R0

		$\text{Cost} = 1 + \text{Cost(A.M)}$
LD	R0, X	$1+1=R$
LD	RI, Y	$1+1=R$
SUB	R0, R0, RI	$1+0=1$
BLTR	#R3, R0	$1+1=R$ $\leftarrow \because \text{indirect A}$
TOTAL		$\text{Cost} = 7$