

KESHAV MEMORIAL ENGINEERING COLLEGE

Expt. No. _____

Date _____

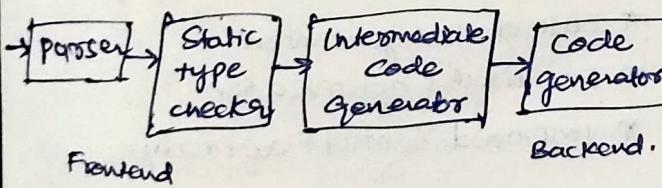
CD4:

Intermediate Code Generation: Variants of Syntax trees, three-Address code, Types and Declarations, Translation of Expressions, Type checking, Control Flow, Backpatching, switch statements, Intermediate code for procedures.

Run-time environment: Storage organization, Stack allocation of space, Access to Non local Data on the stack, parameters passing, heap management and Garbage collection.

Intermediate Code Generation:

- * converting source program into machine code in one pass is not possible always. So compiler generates an easy representation of source language (program) called as Intermediate Language which leads to an efficient code generation.



Forms of Intermediate code

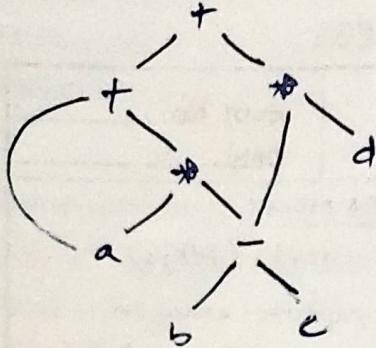
1. Abstract Syntax Tree
2. Polish Notation.
3. Three Address Code

Variants of Syntax Trees:

Directed Acyclic Graphs (DAG's)

- * It efficiently capture common subexpressions and optimize their evaluation.
 - * Initially, expression in the source code are represented as (ASTs)
 - These can contain redundant nodes representing common subexpressions.
 - * The compiler traverses the expression trees and constructs a DAG by identifying common subexpressions and sharing their nodes.
 - * Each node in the DAG represents an operation / a variable.
 - * Edges represent dependencies between the nodes.
 - * Optimizations used: Constant folding, algebraic simplification and common subexpression elimination.
 - * These optimizations help reduce the size of the DAG and improve the efficiency.
 - * The DAG serves as an intermediate representation that simplifies subsequent stages of compilation process while preserving the semantics of the original program.
- It is condensed form of parse trees. Ex: $a + b * c - d$
- leaf \rightarrow operand
node \rightarrow operator.
- the original program.
-

$$\text{Ex:- } a + a^*(b-c) + (b-c)^*d$$



Polish Notation :

- * It is a linearization of AST.
- * Most natural way
- * also called prefix expression.
(operator op₁ op₂)

Ex:-

$$(a+b)^*(c-d) \Rightarrow *(+ab)(-cd)$$

e) $\star +ab-cd$.

Reverse Polish Notation (postfix evaluation)

Ex:-

$$(a+b)^*(c-d) \Rightarrow (ab+)(cd-)^*$$

$\Rightarrow ab+cd-^*$

Three-Address Code :

- * built from two concepts called addresses and instructions.
- * given expression is breakdown into several separate instructions.
- * these instructions can easily translate into assembly language.

General form is $a := b \text{ op } c$

- * In a three-address code there is at most one operator at the right side of an operation.

Ex:-

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

Three-Address code :

- * It is a sequence of statements of the general form.

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

- | | |
|----------------------------------------|--------------------------|
| → names, | → any operator |
| → constants, | → fixed / floating point |
| → comprising generated temporary names | arithmetic operators |
| | → logical operator |

Ex:- $x + y * z$

$t_1 := y * z$

$t_2 = x + t_1$

↳ temporary names.

Advantages :

- * Simplicity * Platform independence
- * Ease of code generation
- * Facilitates optimization
- * Improved Control Flow Analysis
- * Three Address statements can be implemented as records with fields for the operator and the operands.

3-Such representations are :

- * Quadruples
- * Triples
- * Indirect triples.

Expt. No. _____

Date _____

Quadruples:

- * record structure with four fields which are op, arg1, arg2 and result.

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

$x := y * z$

$(*, y, z, x)$:

{ entered in symbol table }

Triples

- * refer to a temporary value by the position of the statement that computes it.

- * record structure with three fields op, arg1 and arg2.

Ex: $x := y * z$

$(*, y, z)$:

Indirect Triples:

- * using pointers to triples, rather than listing the triples themselves.

Ex: $b^* - c + b^* - c$

Three-address code.

$t1 = \text{minus } c$

$t2 = b^* t1$

$t3 = \text{minus } c$

$t4 = b^* t3$

$t5 = t2 + t4$

$a = t5$

- * 3-address code simplifies the task of generating machine code.

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

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

Quadruples	op	arg1	arg2	result	t1	t2	t3	t4	t5	a
op	c	b	c	t1	t1	t1	t3	t4	t5	a
minus	*	*	*	t2	t2	t2	t4	t4	t5	

Types and Declarations:

- * Data type of a name / identifier and type checking of the expression / Statement is to be done before Intermediate code Generation.
- * Type checking uses logical rules to reason about the behaviour of a program at run time.
 - It ensures that the types of the operands match the type expected by an operator.
- * Type information is also needed
 - to calculate the address denoted by an array reference.
 - To insert explicit type conversions
 - To choose the right version of an arithmetic operator.

Basic Types: integer, character, real, Boolean.

Derived types: Array, Structure (record), set and pointer.

Declarations:

- * Declares just one name at a time.
- Generates
- D → T id ; D/E → Sequence of declarations.
- T → BC | record → basic, array or
'D' & record types.
- B → int | float → basic int, float types.
- E → E/[num]C → strings of zero or more integers.

Type expressions

- * Types have structure, represented by using type expressions.
- * A type expression is either a basic type or it is formed by applying an operator called a type constructor to a type expression.
- * Type constructors are derived types (array, structure, pointer function).

Type Expression examples:

Array: int array [20]

→ array (1, 2 ... 19, int)

int [2][3] → array (2, array (3, integer))

Product: t1 * t2 (* is left associative)

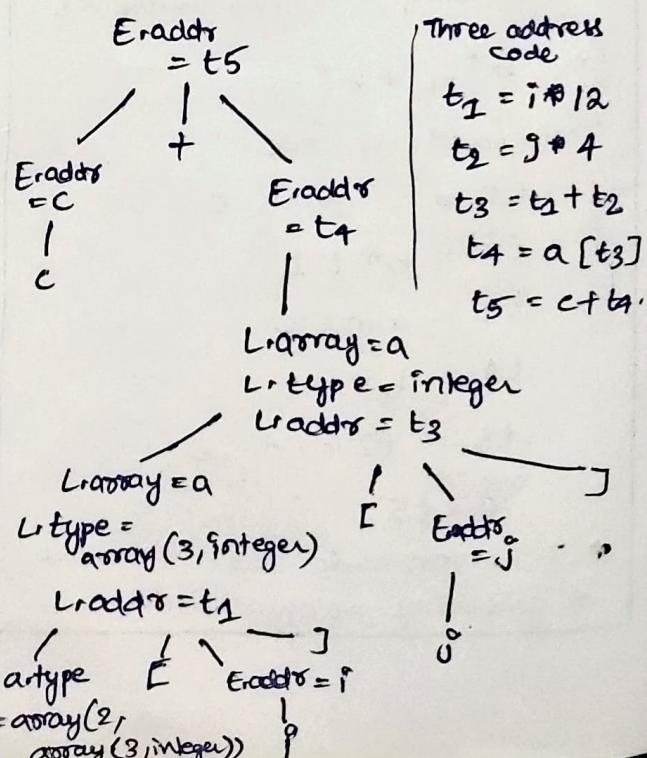
Pointer: float * y; pointer (float).

Type equivalence:

- * If two type expressions are equal then return a certain type else error.
- * When type expressions are represented by graphs, two types are structurally equivalent if and only if one of the following is true:
 - If they are of the same basic type.
 - If they are formed by applying the same constructor to structurally equivalent types.
 - One is a type name that denotes the others.

Translation of array Reference:

Annotated parse tree for c + a[i][j]



Expt. No. _____

Date _____

CONTROL Flow

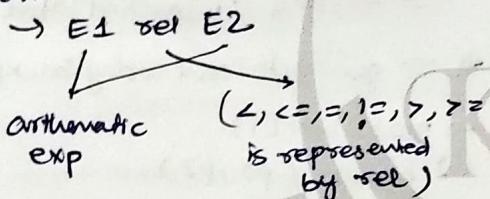
- * The translation of statements such as if-else-statement and while-statements is tied to the translation of boolean expressions.

Boolean expressions are often used

- to compute logical values.
- Alter the flow of control.
- Short circuit evaluation.

Boolean operators (if , II , !)

Relational expressions are of form



Flow of control statements:

- * In the below grammars, non-terminal B represents a boolean expression and non terminal S represents a statement.

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

$\quad | \quad \text{if}(B) S_1 \text{ else } S_2$

$\quad | \quad \text{while}(B) S_1$

$B.\text{true} = \text{newlabel}()$

$\rightarrow B.\text{false} = S.\text{next}()$

$\rightarrow S_1.\text{next} = S.\text{next}$

$S.\text{code} = B.\text{code} // \text{label}(B.\text{true})$

$// S_1.\text{code}$

Type checking

- * A compiler needs to assign a type expression to each component of the source program.
- * Type checking can take one of the two forms:
 - Synthesis → inference.

Type Synthesis: builds the type of an expression from the types of its subexpressions.

- * A typical rule for type synthesis has the form
 - "if f has type $s \rightarrow t$ and x has type s , then expression $f(x)$ has type t ."

f and x denote expressions.

$s \rightarrow t$ denotes a function.

- * The rule can be adapted for $E_1 + E_2$ by viewing it as a function application $\text{add}(E_1, E_2) \Rightarrow E_1.\text{type}$

Type inference: determines the types of a language construct from the way it is used.

- * Some statements do not have values. So, void type is used.

- * A typical rule for type inference has the form.

"if $f(x)$ is an expression, then for some α and β , f has type $\alpha \rightarrow \beta$ and x has type α "

- * Rule:

$S \rightarrow \text{if}(E) S_1 :$

$\quad \text{if } E.\text{type} = \text{boolean type},$

$\quad S_1.\text{type} = S_1.\text{type}$

Type conversion

2 * 3,14

t1 = (float)2 // explicit.

t2 = t1 * 3,14

Two types of conversions happen :

<u>Widening conversion.</u>	<u>Narrowing conversion..</u>
* Intended to preserve information.	* Which can lose information.
* Any type lower in the hierarchy can be widened to a higher type.	* A type 'S' can be narrowed to a type T if there is a path from S to T.
* <pre>double float long int short \ char byte</pre>	* <pre>double ↓ float ↓ long ↓ int ↓ char ↔ short ↔ byte</pre>

Implicit → coercions

Explicit → casts

Backpatching

* Generation of three address code in a single pass creates a problem in identifying addresses of the label statements.

Ex :- Jump statements (goto (Boolean)) flow of control statements.

Def :- Backpatching is the activity of filling up unspecified / missing information of labels using appropriate semantic actions during the code generation phase.

* In Boolean expressions we must insert symbolic labels for jumps.

⇒ So, Boolean expressions needs a separate pass to set them to appropriate addresses for labels.

* Backpatching is needed to avoid the need of two passes.

⇒ Here we save the sequence of instructions into an array and labels will be indices of the array.

* For Non-Terminal, we use two attributes

→ truelist ; it contains the addresses of the true instructions.

⇒ falselist ; it contains the addresses of the false instructions.

* A new non-terminal called M (monkey) $M \rightarrow E$, which gives the address of missing / unspecified label.

* To generate code using backpatching;

→ makeList(i)

→ merge (P1, P2)

→ backpatch (P, i)

Switch statements

Switch (E)

{ case 1 : S1

break;

case 2 : S2

break;

default : S3

three-address code for Switch statements:

evaluate E value

go to test

L1 : S1

goto last

L2 : S2

goto last

L3 : S3

goto last

test : if (E==1) goto L1

Expt. No. _____
Date _____

```
if (E==2) goto L2
goto L3 // default
Last; end.
```

Intermediate code for procedures.
In three-address code, a function call is specified with the below steps.

- Evaluation of parameters in preparation for a call (pass by value).
- Function call itself.

Ex. 1: $n = f(a[i]);$

↓ ↓ → array.
function memory

might translate into

- 1) $t_1 = i * 4$
- 2) $t_2 = a[t_1]$
- 3) param t_2
- 4) $t_3 = \text{call } f, 1$ // 1 indicates no of args of function &
- 5) $n = t_3.$

Ex. 2: int f(int x, int y)
 {
 return x+y+1;
 }

Function call: $f(2+3, 4)$

- 1) $t_1 = x$
- 2) $t_2 = y$
- 3) param t_1
- 4) param t_2
- 5) $t_3 = \text{call } f, 2$

Run-Time Environment:

- * A compiler must accurately implement the abstractions like names, scope, bindings, data types, operators, procedures, and flow-control statements of source language.

The environment deals with:

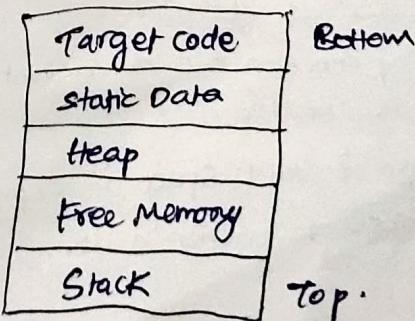
- * Layout and storage allocation for objects.
- * Mechanism of accessing variables.
- * Linkage between procedures.
- * parameter passing mechanism.
- * memory management (stack & heap)
- * Garbage collection.

Storage Organization:

* Compiler demands a block of memory from operating system to run the compiled program. This block is called as run time storage.

* Run-time storage comes in blocks of contiguous bytes divided into parts to hold code and data.

- Generated target code.
- Data objects.
- Information which keeps track of procedure activations.



- * Size of the generated target code is fixed (So, placed in static space).
- * Memory required for data objects is known at compile time.
- * Stack is used to manage active procedures.
- * Heap stores other information.
- * Size of stack and heap is not fixed.

Storage Allocation Strategies:

- ① Static → for all data objects at compile time.
- ② Stack → used to manage the run time storage.
- ③ Heap → used to manage the dynamic memory allocation.

Pros: - easy for the compiler to find the address of the objects in activation record.

Cons: - knowing the size of data object at compile time is mandatory.
⇒ Data structures cannot be created dynamically.
⇒ Recursive procedures are not supported.

- * Also called Control stack.
- * Each time a function is called, space for its local variables is pushed onto a stack.
- * Local variables are stored in activation record and bound to it.
- * Data structures can be created dynamically.

Stack Allocation of Space:

- * for every function call an activation record is created.
- * Allocation of stack space is depends on number of function calls.

fp → (frame pointer)	Returned value
	Actual parameters
	Optional control link
	Optional access link
	Saved machine status (w/optional return addresses)
	Local data
	Temporaries.

Caller's responsibility to finalize

Callee's responsibility to initialize

Activation record: Also known as a stack frame / activation frame, is a data structure used to manage the execution of a subroutine or function within a program.

- ⇒ It typically contains various fields that store information such as:

Returned value: When the subroutine or function finishes, the value stored in this field is passed back to the caller.

Actual parameters: When a subroutine or function is called, the actual parameters passed to it are stored in the activation record.

Optional control link: This link points to the activation record of the callee subroutine or function.

Optional access link: It points to the activation record of the nearest enclosing scope or outer subroutine.

Saved machine status (with optional return address): It typically consists of the values of CPU registers that need to be preserved across subroutine calls.

Expt. No. _____

Date _____

Local data: Space is allocated within the activation record to store local variables declared within the subroutine / function.

Temporaries: are typically used for intermediate calculations in storage within the subroutine / function.

Limitation:

- * Stack allocation would not be feasible if function calls did not nest in time.
- * Memory addressing is slow because of pointers and index registers.

Heap Allocation

- * If nonlocal variables must be retained even after the activation record, then heap allocation is required.
- * Heap allocations allocate continuous blocks of memory when required for storage of activation records / other data obj.
- * It can be deallocated when activation ends, and it can be reused by heap manager.
- * Linked list is suitable for efficient Management.

Access to Nonlocal Data on the Stack.

- ⇒ Storage allocation is done for two types of data : Local and Non-Local data.
- ⇒ If a procedure refers to variables that are not local to it, such variables are called non-local variables.
- ⇒ Global variables are placed in static storage, storage and location is known at compile time.
- ⇒ Accessing these variables is easy, Static scope is used to access non-local names.
- ⇒ In some languages, which support nesting of procedures, In these languages non-local data are handled by using scope information.
- ⇒ Two types of Scope rules :
- * Static scope (used by block structured language).
- * Dynamic scope (used by non-block structured language).

Static Scope / Lexical Scope:

- * Scope is verified by examining the text of the program.
- * If procedures are nested, nesting depth of the procedure is used.

* The lexical scope of nested procedures can be implemented by using Access links and Display.

Display :- Faster to access non-local names.

Dynamic Scope:

* used in non-block structured languages (LISP).

- * By considering the current activation, it determines the scope of declaration of the names at runtime.
- * In this type, the scope is verified at runtime.

Deep Access : The idea is to keep a stack of active variables.

~~Shallow~~ * use control links instead of access links.
* to find a variable, search the stack from top to bottom.

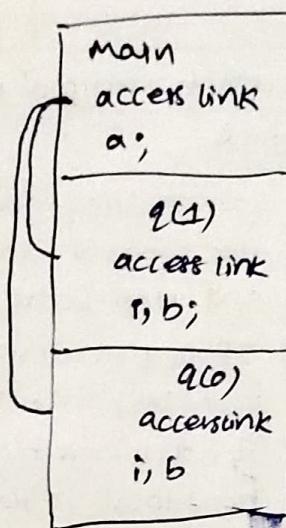
Shallow Access : - The idea to keep a central storage and allot one slot for every variable name.

Parameter passing:

* If one procedure calls another, communication between them is made by using non-local names and parameters.

Actual params : Variables specified in function call.

Formal params : Variables declared in function definition.



L-value → refers to storage.
R-value → refers to value contained in the storage.

Parameter passing Technique,

1. Call by value.
2. Call by reference.
3. Copy restore ← ←
4. Call by name
(Actual params are substituted for the formals).

Heap Management:

- * the heap is the portion of the store that is used for data that lives indefinitely, or until the program explicitly deletes it.
- * local variables typically become inaccessible when their procedure end.
- * Memory Manager is a subsystem that allocates and deallocates space within the heap.
 - ⇒ Interface b/w app prog & OS.
 - ⇒ always keep track of the free space in heap storage.
 - ⇒ It performs:
 1. Allocation
 2. Deallocation.
- * Memory Management would be simpler if
 - (a) all allocation requests were for chunks of the same size, and
 - (b) storage was released predictably, say, first-allocated first deallocated.
- * Desirable properties of memory manager:
 - 1) Space efficiency.
 - 2) Program U.
 - 3) Low overhead.

Memory hierarchy in runtime environment:

Typical sizes

>2GB

256MB - 2GB

128KB - 4MB

16 - 64KB

32 words

