

6 Procedures

Chapter Contents

1. Procedure Calls
2. Name Spaces
 - Name Spaces of Algol-like Languages
 - Runtime Structures to Support Algol-like Languages
3. Communicating Values Between Procedures
 - Passing Parameters
 - Returning Values
 - Establishing Addressability
4. Standardized Linkages

The Procedure Abstraction

The procedure is one of the **central abstractions** that underlie most modern programming languages

- important way for programmers to structure their software
- basic unit of work for most compilers

Procedures provide **three** critical abstractions

1. Procedure Call
2. Name Space
3. External Interface

Together, these abstractions allow the construction of **complex** and **nontrivial** programs.

The Procedure Abstraction

The **procedure call** abstraction is a standard mechanism to...

- invoke a procedure
- map parameters from the caller's name space to the callee's name space
- return control to the caller and resume execution at the point after the call
- return one or more values to the caller

Caller and Callee

In a procedure call, we refer to the procedure that is invoked as the **callee** and to the calling procedure as the **caller**.

The Procedure Abstraction

Each procedure creates a new and protected **name space**

- enable declaration of new names without concern for the surrounding context
- caller can map values and variables in its name space into the callee's name space
- procedure functions correctly and consistently when called from different contexts

Formal and Actual Parameters

A name declared as a parameter of some procedure p is a **formal parameter** of p .
A value or variable passed as a parameter to p at a specific call site is an **actual parameter** of the call.

The Procedure Abstraction

Procedures define the critical **interfaces** among the parts of large software systems

The **linkage convention** is an agreement between the compiler and operating system that defines rules

- that map names to values and locations,
- that preserve the caller's runtime environment, and
- that create the callee's environment

It creates a context in which a programmer can safely invoke code written by others.

Without a linkage convention, both the programmer and the compiler would need **detailed knowledge** about the implementation of the callee at each procedure call.

Procedure Calls

In the following, we focus on Algol-like languages that have a simple and clear call/return discipline

- a procedure call transfers control from the call site in the caller to the **start** of the callee
- on exit from the callee, control returns to the point **immediately following** the call site

If the callee invokes **other** procedures, they are activated in the same way.

Activation

A call to a procedure activates it. Thus, we call an instance of its execution an activation.

Procedure Calls

```

1  program Main(input, output);
2      var x, y, z: integer;
3
4      procedure Fee;
5          var x: integer;
6      begin { Fee }
7          x := 1;
8          y := x * 2 + 1
9      end;
10
11     procedure Fie;
12         var y: real;
13
14         procedure Foe;
15             var z: real;
16
17             procedure Fum;
18                 var y: real;
19             begin { Fum }

```

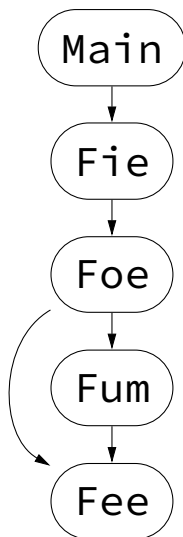
```

20         x := 1.25 * z;
21         Fee;
22         writeln('x_ = ', x)
23     end;
24
25     begin { Foe }
26         z := 1;
27         Fee;
28         Fum
29     end;
30
31     begin { Fie }
32         Foe;
33         writeln('x_ = ', x)
34     end;
35
36     begin { Main }
37         x := 0;
38         Fie
39     end.

```

Procedure Calls

The **call graph** shows the set of potential calls among the procedures. Executing Main can result in two calls to Fee: one from Foe and another from Fum.



1. Main calls Fie
2. Fie calls Foe
3. Foe calls Fum
4. Fum returns to Foe
5. Foe calls Fum
6. Fum calls Fee
7. Fee returns to Fum
8. Fum returns to Foe
9. Foe returns to Fie
10. Fie returns to Main

The **execution history** shows that both calls occur at runtime.

Procedure Calls

When the compiler generates code for calls and returns, that code must **preserve** enough information so that calls and returns operate correctly.

The call and return behavior of Algol-like languages can be modelled with a **stack**

- when a caller calls a callee, it **pushes** the return address in the caller onto the stack
- when the callee returns, it **pops** that address off the stack and jumps to the address

Return Address

When p calls q , the address in p where execution should continue after p 's return is called its **return address**.

Name Space

In most procedural languages, a complete program will contain multiple **name spaces**.

A name space...

- maps a set of names to a set of values and procedures over a range of statements
- may inherit some names from other name spaces
- can create names that are inaccessible outside the name space

Name space rules give the programmer control over access to information.

Scope

In an Algol-like language, **scope** refers to a name space. The term is often used in discussions of the visibility of names.

Name Spaces of Algol-like Languages

Most programming languages inherit many of the conventions that were defined for Algol 60. This is particularly true of the rules that govern the visibility of names.

While Algol 58 introduced the fundamental notion of the **compound statement**, Algol 60 was the first language implementing **nested function** definitions with **lexical scope**.

Lexical Scope

Scopes that nest in the order that they are encountered in the program are often called **lexical scopes**.

Nested Lexical Scopes

Most Algol-like languages allow the programmer to **nest** scopes inside one another.

The limits of a scope are marked by specific terminal symbols

- Pascal demarcated scopes with a begin at the start and an end at the finish
- C uses curly braces, { and }, to begin and end a block, its notion of a scope

Note Each procedure also defines a new scope that covers its entire definition. The programmer can declare new variables and, in some languages, procedures in this scope.

With (nested) lexical scopes a compiler must decide which use of name refers to which declaration. To do so, it uses the **scoping rules** defined by the programming language.

Nested Lexical Scopes

Lexical Scoping

The most common scoping discipline is called **lexical scoping**. It follows a simple general principle.

In a given scope, each name refers to its lexically closest declaration.

Example If s is used in the current scope, it refers to the s declared in the current scope. If none exists, it refers to the declaration of s that occurs in the **closest enclosing** scope.

Note The outermost scope contains **global variables**.

Nested Lexical Scopes

```

1 program Main0(input, output);
2   var x1, y1, z1: integer;
3
4   procedure Fee1;
5     var x2: integer;
6   begin { Fee1 }
7     x2 := 1;
8     y1 := x2 * 2 + 1
9   end;
10
11  procedure Fie1;
12    var y2: real;
13
14    procedure Foe2;
15      var z3: real;
16
17      procedure Fum3
18        var y4: real;
19      begin { Fum3 }

```

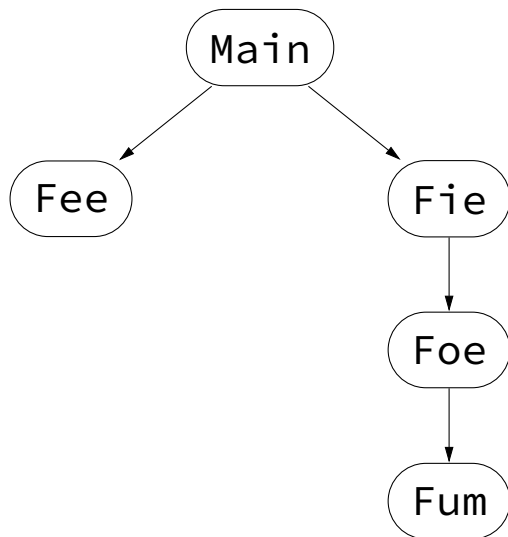
```

20       x1 := 1.25 * z3;
21       Fee1;
22       writeln('x1=', x1)
23     end;
24
25   begin { Foe2 }
26     z3 := 1;
27     Fee1;
28     Fum3
29   end;
30
31  begin { Fie1 }
32    Foe2;
33    writeln('x1=', x1)
34  end;
35
36 begin { Main0 }
37   x1 := 0;
38   Fie1
39 end.

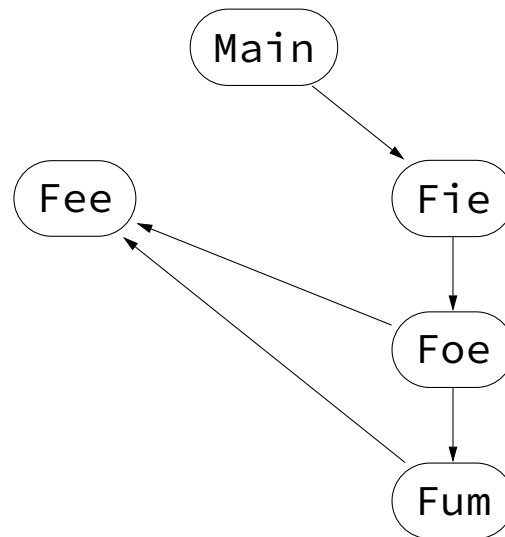
```

Nested Lexical Scopes

Nesting Relationships



Call Graph



Nested Lexical Scopes

The listing on Slide 280 shows each name with a subscript that indicates its **level** number. Names declared in a procedure are **one** level deeper than the level of the procedure name.

To represent names in a lexically scoped language, the compiler can use the **static coordinate** for each name.

Static Coordinate

The static coordinate is a pair $\langle l, o \rangle$, where l is the name's lexical nesting level and o is the its offset in the data area for level l .

To obtain l , the front end uses a **lexically scoped** symbol table. The offset o should be stored with the name and its level in the symbol table.

Nested Lexical Scopes

Example The table below shows the static coordinate for each variable name in each procedure of the code on Slide 280.

Scope	x	y	z
Main	$\langle 1, 0 \rangle$	$\langle 1, 4 \rangle$	$\langle 1, 8 \rangle$
Fee	$\langle 2, 0 \rangle$	$\langle 1, 4 \rangle$	$\langle 1, 8 \rangle$
Fie	$\langle 1, 0 \rangle$	$\langle 2, 0 \rangle$	$\langle 1, 8 \rangle$
Foe	$\langle 1, 0 \rangle$	$\langle 2, 0 \rangle$	$\langle 3, 0 \rangle$
Fum	$\langle 1, 0 \rangle$	$\langle 4, 0 \rangle$	$\langle 3, 0 \rangle$

Scope Rules across Various Languages

Scoping rules vary idiosyncratically from programming language to programming language.

FORTRAN

- single global scope that groups variables into a “common block”
- each procedure or function has a local scope for parameters, variables, and labels
- save statement to give a local variable the lifetime of a global variable

ANSI C

- global scope for procedure names and global variables
- each procedure has a local scope for parameters, variables, and labels
- blocks also define a separate local scope and they can be nested
- file-level scope includes `static` names that are not enclosed in a procedure

Runtime Structures to Support Algol-like Languages

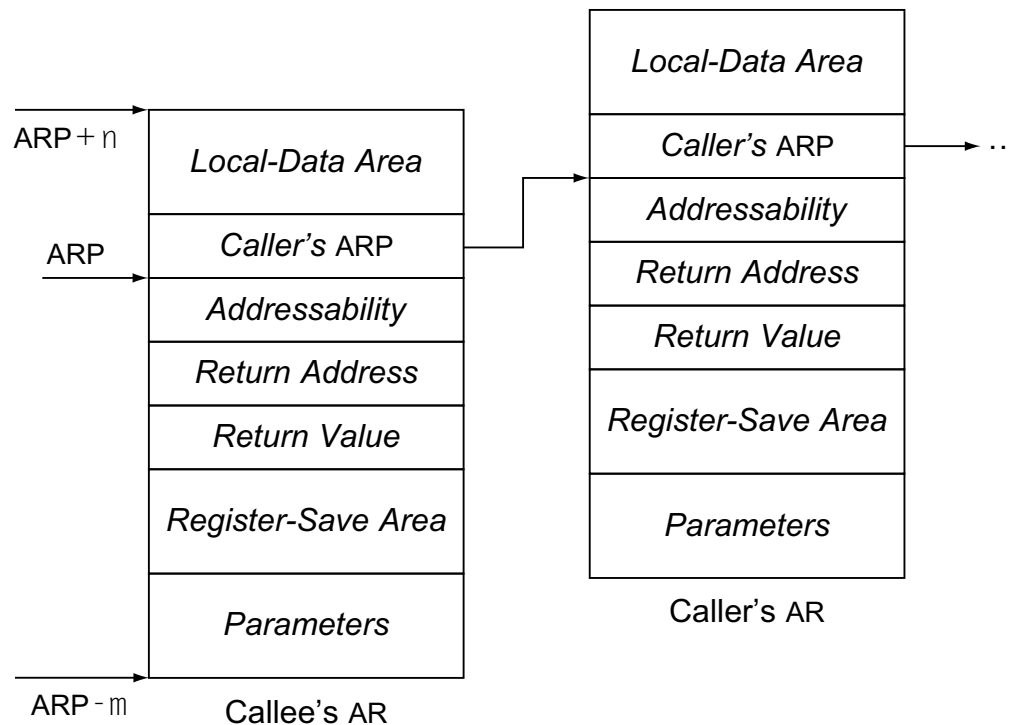
To realize **procedure calls** and **scoped name spaces**, runtime structures are required. The key data structure to implement both abstractions is the **activation record** (AR).

In principle, **each procedure call** gives rise to a new AR

- stores the **return address** where the callee can find it
- maps the **actual parameters** to the formal parameter names known by the callee
- reserves storage space for **local variables** declared in the callee's scope
- provides **other information** needed to connect the callee to the surrounding program

If **multiple instances** of a procedure are active, *e.g.*, recursive calls, each has its own AR.

Runtime Structures to Support Algol-like Languages



Runtime Structures to Support Algol-like Languages

The entire AR is addressed through an **activation record pointer** (ARP), with various fields in the AR found at positive and negative offsets from the ARP.

The ARP always points to a designated location in the AR

- **fixed-length fields** are stored in the center, where compiled code can access those items at fixed offsets from the ARP
- **variable-length fields** whose sizes may change from one invocation to another, e.g., parameters and local data, are stored at either end

Question

Is this a reasonable layout for the activation record?

Runtime Structures to Support Algol-like Languages

Because procedures typically access their AR frequently, most compilers **dedicate** a hardware register to hold the ARP of the current procedure.

Examples

- SPARC Assembly uses registers %fp (frame pointer) and %sp (stack pointer) to demarcate the activation record (or stack frame)
- the corresponding registers in X86-64 Assembly are called %rbp (base pointer) and %rsp (stack pointer)
- in ILOC, we refer to this dedicated register as r_{arp}

Local Storage

The compiler needs to **reserve space for local data** items in the AR

- assign each such item an appropriately sized area
- record its offset from the ARP in the symbol table

Together with its lexical level the item's offset becomes the its static coordinate.

If the compiler **cannot know the size** of a local variable at compile time

1. leave space for a pointer to the actual data or to a descriptor for an array
2. arranges to allocate the actual storage elsewhere at runtime
3. fill the reserved slot with the address of the dynamically allocated memory

Local Storage

If the source language allows the program to specify an initial value for a variable, the compiler must arrange for that **initialization** to occur.

Static and Global Variables

- lifetime independent of any procedure
- data can be inserted directly into the appropriate locations by the loader
- for example, use a **data section** (`.data`) in assembly code

Local Variables

- must be initialized at runtime as a procedure may be invoked multiple times
- generate instructions that store the initial values to the appropriate locations

Local Storage

When p calls q , one of them must **save the register values** that p needs after the call

- it may be necessary to save all the register values or a subset may suffice
- on return to p , these saved values must be restored
- saved registers are stored in the AR of either p or q , or both

If the **callee saves a register**, its value is stored in register save area of the callee's AR. Similarly, if the **caller saves a register**, its value is stored in the caller's register save area.

Calling Convention

Who saves what is defined by the **calling convention** used by the operating system. Solaris, Linux, FreeBSD, and macOS follow the System V AMD64 ABI calling convention. Windows follows the Microsoft x64 calling convention.

Local Storage

Example The table below shows which registers are saved by the caller and the callee, respectively, when following the System V AMD64 ABI calling convention.

Register	Purpose	Saved by
%rax	return value	–
%rbx	scratch	callee
%rcx	4 th parameter	–
%rdx	3 rd parameter	–
%rsi	2 nd parameter	–
%rdi	1 st parameter	–
%rbp	base pointer	callee
%rsp	stack pointer	callee

%r8	5 th parameter	–
%r9	6 th parameter	–
%r10	scratch	caller
%r11	scratch	caller
%r12	scratch	callee
%r13	scratch	callee
%r14	scratch	callee
%r15	scratch	callee

Allocating Activation Records

When p calls q at runtime, the code that implements the call must allocate an AR for q and initialize it with the appropriate values.

Fields of AR are stored in **memory**

- caller needs to have access to callee's AR to store parameters, return address, *etc.*
- this forces allocation of callee's AR into caller, who might not know the local data size

Values of AR passed by **registers**

- allocation of the ar can be performed in the callee, who knows the local data size
- callee may store into its ar some of the values passed in registers after allocation

The compiler can allocate ARs on the **stack** or on the **heap**. In certain cases, it can allocate a single **static** AR for a procedure or **coalesce** ARs of a set of procedures.

Allocating Activation Records

If procedure calls and returns are balanced, they follow a last-in, first-out (LIFO) discipline. Since the ARs also follow this LIFO ordering, they can be allocated on the **stack**.

Keeping activation records on the stack has several advantages

- inexpensive allocation and deallocation: simply move the top-of-stack pointer
- separation of concerns between caller and callee
 - caller allocates space for address parameters, return address, *etc.*
 - callee can extend AR to include local data area and variable-sized objects
- debugger can walk the stack to produce a snapshot of the currently active procedures

Example Pascal, C, and Java are typically implemented with stack-allocated ARs

Allocating Activation Records

Certain language features prevent ARs from being allocated on the stack

- if a procedure can outlive its caller
- if a procedure can return an object that includes references to its local variables

In these situations, ARs can be kept on the **heap**. With heap-allocated ARs, variable-size objects can be allocated as separate objects on the heap.

- explicit deallocation: procedure return code frees the AR and variable-size objects
- implicit deallocation: garbage collector frees them when they are no longer useful

Example Implementations of Scheme and ML typically use heap-allocated ARs

Allocating Activation Records

A **leaf procedure** is a procedure that calls no other procedures. The compiler can allocate ARs for leaf procedures **statically**. This eliminates the runtime costs of AR allocation.

Example If the calling convention requires the caller to save its own registers, then the AR of a leaf procedure needs no register save area.

Under certain circumstances the compiler can do even better! Assuming that callees cannot outlive callers, **only one** leaf procedure can be active at any point during execution.

- allocate a **single** static AR for use by **all** leaf procedures
- AR must be large enough to accommodate any of the program's leaf procedures
- static variables of all leaf procedures can be laid out together in that single AR

Using a single static AR for leaf procedures reduces the space overhead of separate static ARs for each leaf procedure.

Allocating Activation Records

If the compiler discovers a set of procedures that are always invoked in a fixed sequence, it may be able to **coalesce** their ARs and save allocation costs.

Example If a call from p to q always results in calls to r and s, the compiler may find it profitable to allocate the ARs for q, r, and s at the same time.

In practice, this optimization is **limited** by the following factors

- separate compilation
- function-valued parameters

Communicating Values Between Procedures

A procedure **encapsulates** common operations relative to a small set of names. To make this encapsulation work, the compiler needs to **bind these names to values** correctly.

Names in a procedure can refer to the following values

- parameter values
- return value
- values of global variables
- values of local variables in an enclosing procedures

Passing Parameters

Parameter binding plays a critical role in our ability to write abstract, modular code

- procedure can be written without knowledge of the contexts in which it will be called
- procedure can be called from many contexts without exposing its internal details

Most modern programming languages use one of two conventions for mapping actual parameters to formal parameters: **call-by-value** and **call-by-reference**.

Call by Value

Call by Value

Call-by-value is a convention where the caller evaluates the actual parameters and passes their values to the callee.

Any modification of a value parameter in the callee is not visible in the caller.

Call-by-value parameter passing

- caller **copies** the value of an actual parameter into the **appropriate location** for the corresponding formal parameter, *i.e.*, a register or a parameter slot in the callee's AR
- **only one name** refers to that value: the name of the formal parameter
- its value is an **initial condition** determined by evaluating the actual parameter
- if the callee **changes the value**, the change is visible inside the callee, but not outside

Call by Value

```

1  int fee(int x, int y) {
2      x = 2 * x;
3      y = x + y;
4      return y;
5  }
6
7  c = fee(2,3);
8  a = 2;
9  b = 3;
10 c = fee(a,b);
11 a = 2;
12 b = 3;
13 c = fee(a,a);

```

The three invocations produce the following results when invoked using call-by-value parameter binding.

Call by Value	a		b		Return Value
	in	out	in	out	
fee(2,3)	–	–	–	–	7
fee(a,b)	2	2	3	3	7
fee(a,a)	2	2	3	3	6

With call by value, the binding is simple and intuitive.

Call by Reference

Call by Reference

Call-by-reference is a convention where the caller passes an address for the formal parameter to the callee.

If the actual parameter is a variable (rather than an expression), then changing the formal's value also changes the actual's value.

Call-by-reference parameter passing

- caller stores a **pointer** in the register or AR slot for each parameter
- in the callee, each call-by-reference formal parameter needs a **level of indirection**

Call by Reference

Dealing with different kinds of actual parameters

- **variables**: pass the variable's address to the callee
- **expressions**
 1. caller evaluates the expression
 2. stores the result in the local data area of its own AR
 3. passes a pointer to that result to the callee
- **constants**: should be treated as expressions to avoid being changed by the callee

Note Some languages forbid passing expressions and constants as actual parameters to call-by-reference formal parameters.

Call by Reference

Call by reference differs from call by value in **two critical** ways

1. any redefinition of a reference formal parameter is reflected in the corresponding actual parameter
2. any reference formal parameter might be bound to a variable that is accessible by another name inside the callee

Aliasing

When two names can refer to the same location, they are said to be **aliases**. Aliasing can create counterintuitive behavior.

Call by Reference

```

1  int fee(int* x, int* y) {
2      *x = 2 * *x;
3      *y = *x + *y;
4      return *y;
5  }
6
7  c = fee(2,3);
8  a = 2;
9  b = 3;
10 c = fee(&a,&b);
11 a = 2;
12 b = 3;
13 c = fee(&a,&a);

```

With (simulated) call-by-reference parameter passing, the example produces different results.

Call by Value	a		b		Return Value
	in	out	in	out	
fee(2,3)	<i>does not compile in C</i>				
fee(a,b)	2	4	3	7	7
fee(a,a)	2	8	3	3	8

The third call causes x and y to refer to the same location, and thus, the same value.

Space for Parameters

The size of the representation for a parameter has an impact on the cost of procedure calls.

Scalar Values (e.g., variables or pointers)

- stored in registers or in the parameter area of the callee's AR
- cost per parameter is small for both call-by-value and call-by-reference parameters

Large Values (e.g., arrays, records, or structures)

- copying call-by-value parameters will add significant cost to the procedure call
- some languages allow the implementation to pass such parameters by reference
- others include constructs to specify that passing a particular parameter by reference is acceptable, e.g., `const` in C

Returning Values

Because the return value, by definition, is used after the callee terminates, it needs storage **outside** the callee's AR.

If the return value is **small** and has a **fixed size**

- store return value in a **designated register**
- store return value in the **caller's AR**
 1. caller allocates space for return value in its own AR
 2. caller stores a pointer to that space in the return slot of its own AR
 3. callee uses ARP to caller's AR load the pointer

If the caller **cannot know the size** of the returned value

1. callee allocates space for return value, possibly on the heap
2. callee stores the pointer in the return-value slot of the caller's AR
3. caller must free the space allocated by the callee

Establishing Addressability

Apart from binding parameters and passing the return value, the compiler must ensure that each procedure can **generate an address** for each variable that it needs to reference.

In general, the address calculation consists of **two** portions

- finding the **base address** of the data area of the scope that contains the value
- finding the correct **offset** within that data area for the value

Data Area and Base Address

The memory region that holds the data for a specific scope is called its **data area**.

The address of the start of a data area is often called a **base address**. The base address can either be static, *i.e.*, known at compile time, or dynamic, *i.e.*, only known at runtime.

Variables with Static Base Addresses

Compilers typically arrange for global and static data areas to have **static base addresses**.

The strategy to generate an address for such a variable is simple

- compute the data area's base address into a register
- add its offset to the base address

Compiler attaches a **symbolic, assembly-level label** for the variable to the data area.

Depending on the instruction set, the label might be used in a **load immediate** operation or it might evaluate to an address which can be loaded into a register with a **move** operation.

The label becomes a **relocatable symbol** for the assembler and the loader, which convert it into a runtime virtual address.

Variables with Static Base Addresses

The compiler constructs the label for a base address by mangling the name

- add a prefix, a suffix, or both to the original name
- use characters that are legal in the assembly code but not in the source language

Name Mangling

The process of constructing a unique string from a source-language name is called **name mangling**.

Variables with Static Base Addresses

```

1      .data                                # begin of the global data area
2      _msg:
3      .ascii  "Hello, \world!\n"
4      .text                                # begin of the program text
5      .globl _main                        # export main
6      _main:
7      pushq %rbp                          # save current base pointer
8      movq  %rsp, %rbp                    # move base pointer to stack pointer
9      leaq  _msg(%rip), %rdi              # load value from static base address
10     xorq  %rax, %rax                     # zero out %rax (varargs = 0)
11     call  _printf                        # call printf
12     popq  %rbp                          # restore base pointer
13     retq                                # return

```

Note Global constant can also be loaded with `movq _msg@GOTPCREL(%rip), %rdi`.

Variables with Dynamic Base Addresses

Local variables declared within a procedure are typically stored in the procedure's AR

- they have **dynamic base addresses**
- compiler needs a mechanism to find the addresses of various ARs

Fortunately, lexical scoping rules limit the set of ARs that can be accessed from any point in the code to the **current AR** and the **ARs of lexically enclosing procedures**.

Local Variable of the Current Procedure

Accessing a **local variable of the current procedure** is trivial: its base address is simply the address of the current AR, which is stored in the ARP.

- add offset to the ARP and use the result as the value's address
- most processors provide efficient support for these common operations

In some cases, a value is not stored at a **constant** offset from the ARP.

- value might reside in a register, *i.e.*, loads and stores are not needed
- variable might have an unpredictable or changing size
 1. store it in an area reserved for variable-size objects in AR or on the heap
 2. reserve space in AR for a pointer the variable's actual location
 3. generate one additional load to access the variable

Local Variables of Other Procedures

To access a **local variable of some enclosing lexical scope**, the compiler must build runtime data structures that can map a **static coordinate** into a **runtime address**.

Access Links or Static Links

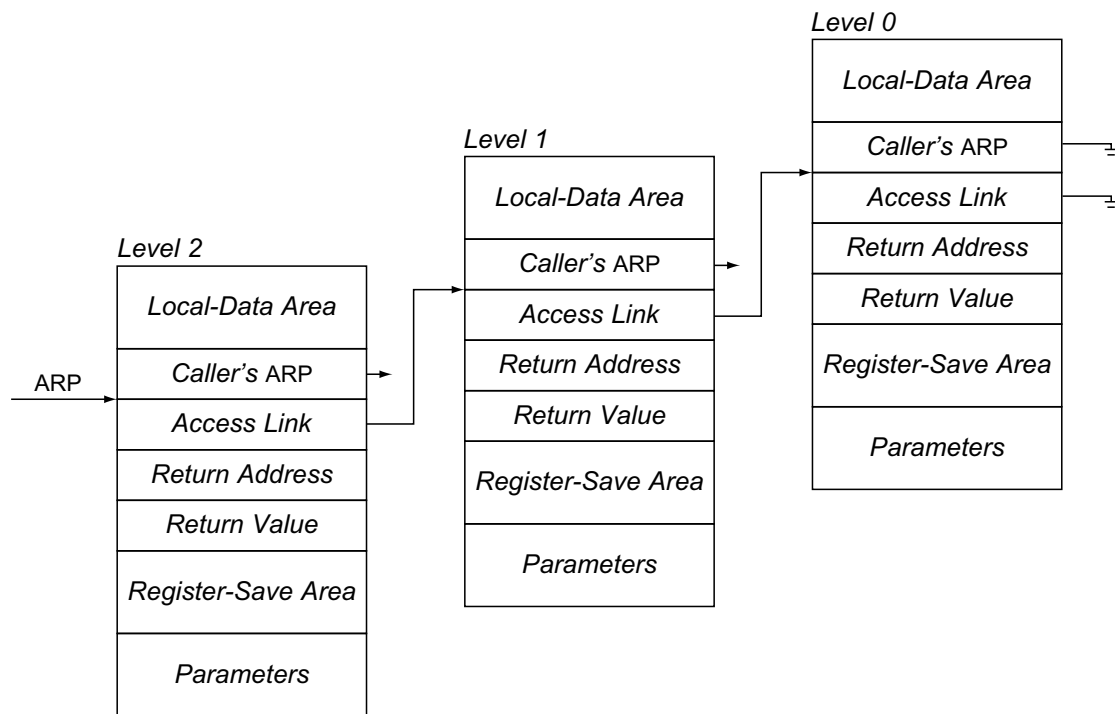
- extend each AR with a pointer to the AR of its **immediate** lexical ancestor
- emit code to walk the chain of links and to find the appropriate ARP
- load value with its offset (from static coordinate) relative to this ARP

Global Display

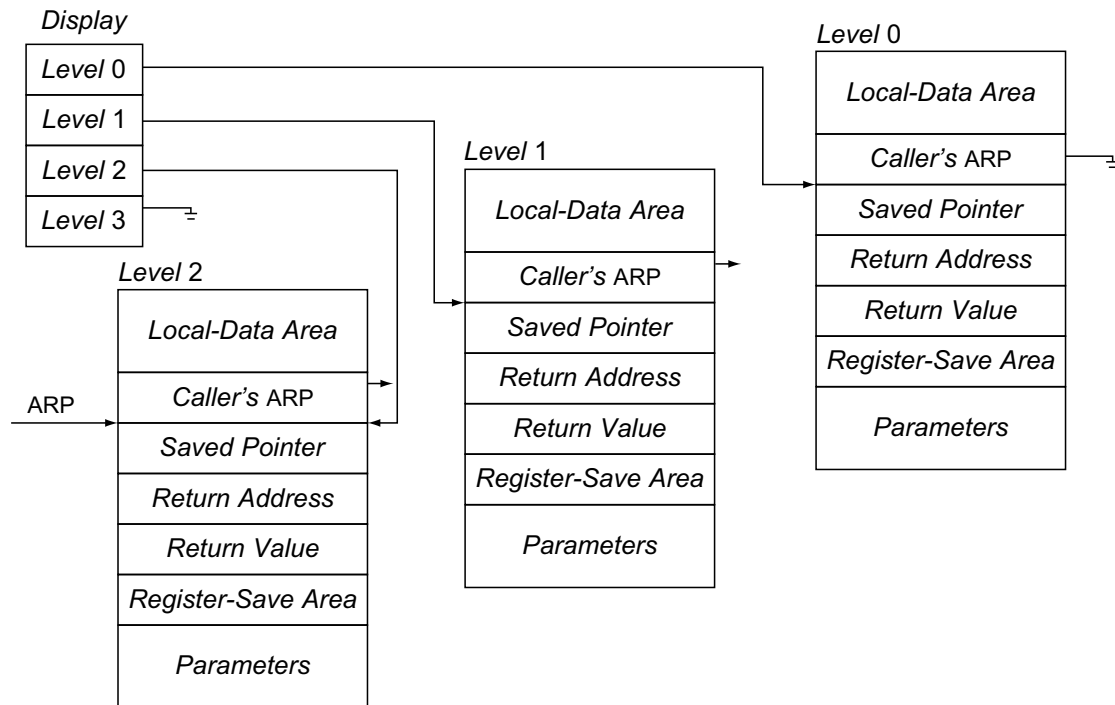
- a global array stores ARP of the **most recent** procedure activation at each lexical level
- local variables of other procedures are referenced indirectly through the display

While the cost of nonlocal access is fixed with a global display, the compiler must insert code where needed to maintain the values in the display.

Access Links



Global Display



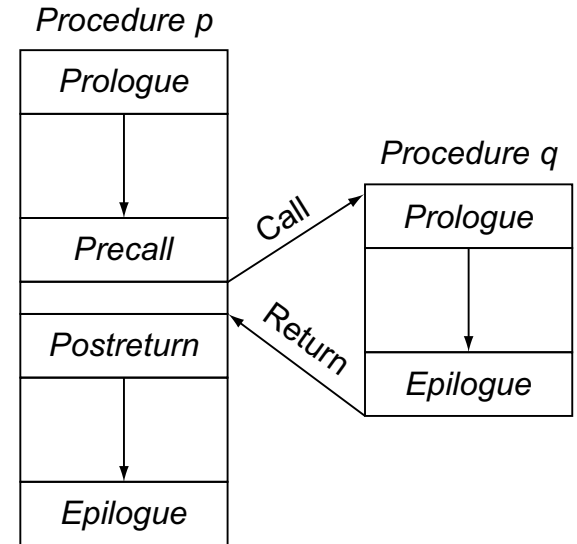
Standardized Linkages

The **procedure linkage** is a contract between the compiler, the operating system, and the target machine

- assigns responsibility for naming, resource allocation, addressability, and protection
- ensures interoperability of code emitted by one compiler with code from other sources

Pieces of a standard procedure linkage

- each procedure has a **prologue sequence** and an **epilogue sequence**
- each call site includes both a **precall sequence** and a **postreturn sequence**



System V AMD64 ABI

Mapping Actual to Formal Parameters

- integer arguments (including pointers) are placed in the registers `%rdi`, `%rsi`, `%rdx`, `%rcx`, `%r8`, and `%r9`, in that order
- floating point arguments are placed in the registers `%xmm0`-`%xmm7`, in that order.
- parameter values in excess of the available registers are pushed onto the stack
- if the function takes a variable number of arguments (like `printf`) then the `%eax` register must be set to the number of floating point arguments

Saving Registers

- the callee may use any registers, but it must restore the values of the registers `%rbx`, `%rbp`, `%rsp`, and `%r12`-`%r15`, if it changes them.

Returning Values

- the return value of a call is placed in `%rax`

7 Code Shape

Chapter Contents

1. Assigning Storage Locations
2. Arithmetic Operators
3. Boolean and Relational Operators
4. Storing and Accessing Arrays
5. Character Strings
6. Structure References
7. Control-Flow Constructs
8. Procedure Calls

Code Shape

Typically, there are **many different ways** how a compiler can map a source-language construct into a sequence of operations in the target instruction set of a given processor.

These variations use different operations and different approaches

- some of these implementations are faster than others
- some use less memory
- some use fewer registers
- some might consume less energy during execution

The concept of **code shape** encapsulates all of the decisions, large and small, that the compiler writer makes on how to represent the computation in both IR and assembly code.

Code shape has a **strong impact** both on the behavior of the compiled code and on the ability of the optimizer and back end to improve it.

Code Shape

Example Consider the way that a C compiler might implement a `switch` statement that switched on a single-byte character value.

```
1 char ch = ...;
2 switch(ch) {
3     key1: ... break;
4     ...
5     key256: ... break;
6     default: ...;
7 }
```

Alternatives

1. use a cascaded series of `if-then-else` statements to implement the `switch` statement
2. use tests that perform a binary search
3. construct a table of 256 labels and interpret the character by loading the corresponding table entry and jumping to it

Which implementation is best for a particular `switch` statement depends on many factors

- number of cases and their relative execution frequencies
- cost structure for branching on the processor

Assigning Storage Locations

The compiler must **assign a storage location** to each value produced by the code, taking the following factors into account.

- the value's type, its size, its visibility, and its lifetime
- runtime layout of memory
- source-language constraints on the layout of data areas and data structures
- target-processor constraints on placement or use of data

We can distinguish **two types** of values

- **named value**: lifetime is defined by source-language rules and actual use in the code
- **unnamed values**: must be handled consistently with the meaning of the program, but the compiler has great leeway in determining where these values reside and how long to retain them.

Assigning Storage Locations

For each value, the compiler must also decide whether to keep it in a **register** or to keep it in **memory**. The compiler's **memory model** guides its choice of locations for values.

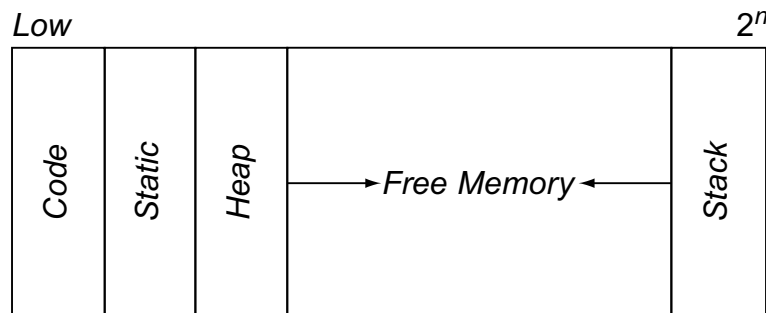
Memory-to-Memory Model

- compiler assumes that all values reside in memory
- values are loaded into registers as needed and stored to memory after each definition
- IR typically uses **physical register** names

Register-to-Register Model

- compiler assumes that it has enough registers to express the computation
- use a distinct **virtual register** for each value that can legally reside in a register
- store a virtual register's value to memory only when absolutely necessary

Placing Runtime Data Structures

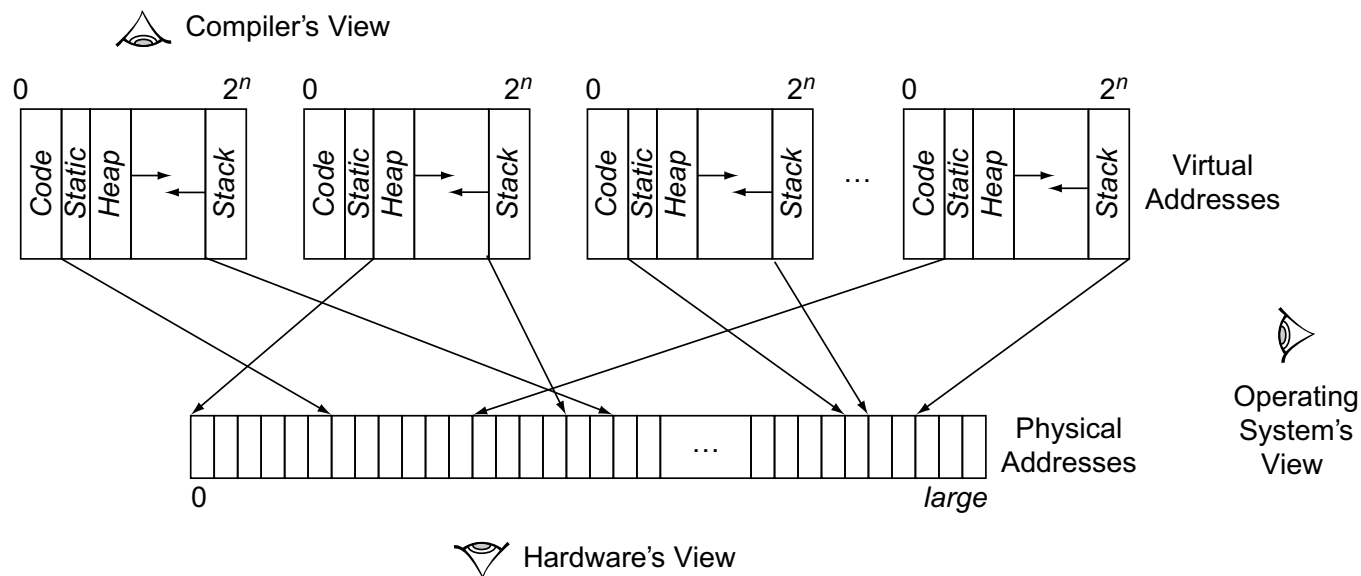


Typical layout for the address space used by a single compiled program

- **code** sits at the bottom of the address space
- **static area** holds both static and global data areas, along with any fixed-size data created by the compiler, e.g., jump tables, debug information
- **heap** contains dynamically allocated data structures
- runtime **stack** is used to stack-allocate activation records, if possible

Placing Runtime Data Structures

The operating system maps multiple **logical address spaces** into the single **physical address space** supported by the processor.



Layout for Data Areas

For convenience, the compiler groups together the storage for values with the **same lifetimes** and **visibility**. It creates distinct **data areas** for them.

Typical placement rules for Algol-like languages

- if x is declared locally in procedure p
 - if its value is not preserved across distinct invocations of p
→ assign it to procedure-local storage
 - if its value is preserved across invocations of p
→ assign it to procedure-local static storage
- if x is declared as globally visible
→ assign it to global storage
- if x is allocated under program control
→ assign it to the runtime heap

Assigning Offsets

In the case of local, static, and global data areas, the compiler must assign each name an **offset** inside the data area. Target ISAs constrain the placement of data items in memory.

Typical alignment rules

- 32-bit integers and floating-point numbers begin on word (32-bit) boundaries
- 64-bit integer and floating-point data begin on doubleword (64-bit) boundaries
- string data begin on halfword (16-bit) boundaries

To **minimize wasted space**, the compiler should order the variables into groups, from those with the most restrictive alignment rules to those with the least.

Note Doubleword alignment is more restrictive than word alignment.

Keeping Values in Registers

The register-to-register memory model has **three** principal advantages

- it is simple
- it can improve the results of analysis and optimization
- it enhances portability by postponing processor-specific constraints until optimization

However, only **unambiguous values** can be kept in registers. **Ambiguous values** cannot be kept in a register across either a definition or a use of another ambiguous value.

Ambiguous and Unambiguous Values

Any value that can be accessed by multiple names is **ambiguous**. In contrast, a value that can be accessed with just one name is **unambiguous**.

Keeping Values in Registers

Ambiguity arises in several ways

- values stored in pointer-based variables are often ambiguous
- call-by-reference formal parameters and name scoping rules can make the formal parameters ambiguous
- array-element values can be ambiguous values since two references could refer to the same location

$$\begin{aligned} a &\leftarrow m + n \\ b &\leftarrow 13 \\ c &\leftarrow a + b \end{aligned}$$

With careful analysis, the compiler can **disambiguate** some of these cases.

- if a and b refer to the same location, c gets the value 26, otherwise it gets $m + n + 13$
- to keep a in a register across assignment to another ambiguous value, the compiler needs to prove that the sets of locations to which a and b can refer are **disjoint**

Since pairwise ambiguity analysis is expensive, compilers typically relegate ambiguous values to memory, with a load before each use and a store after each definition.

Arithmetic Operators

Modern processors provide a full complement of operations to evaluate expressions

- arithmetic operators, e.g., add, sub, imul, and idiv
- shift and rotate operators, e.g., shl, shr, rol, and ror
- boolean operators, e.g., and, or, xor, and not

To generate code for a **trivial expression**, such as $a + b$, the compiler emits code to load the values of a and b into registers, followed by an instruction to perform the addition.

```
loadI  @a          ⇒ r1
loadA0 rarp, r1 ⇒ ra
loadI  @b          ⇒ r2
loadA0 rarp, r2 ⇒ rb
add     ra, rb    ⇒ rt
```

If the value of a or b is already in a register, the compiler can avoid the load instructions and simply use that register in place of r_a or r_b , respectively.

Arithmetic Operators

If the expression is represented in a tree-like IR, this process fits into a **postorder** tree walk.

base

- returns the name of a register holding the base address for an identifier
- if needed, it emits code to get that address into a register

offset

- returns the name of a register holding the identifier's offset
- offset is relative to the address returned by base

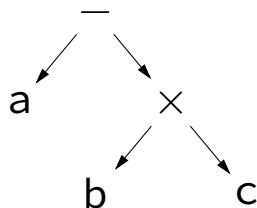
```

1 procedure expr(n)
2   if  $n \in \{+, -, \times, \div\}$  then
3      $t_1 \leftarrow \text{expr}(n.\text{left})$ 
4      $t_2 \leftarrow \text{expr}(n.\text{right})$ 
5      $r \leftarrow \text{NextRegister}()$ 
6     emit(n,  $t_1$ ,  $t_2$ , r)
7   else if  $n = \text{ident}$  then
8      $t_1 \leftarrow \text{base}(n)$ 
9      $t_2 \leftarrow \text{offset}(n)$ 
10     $r \leftarrow \text{NextRegister}()$ 
11    emit(loadAO,  $t_1$ ,  $t_2$ , r)
12  else if  $n = \text{num}$  then
13     $r \leftarrow \text{NextRegister}()$ 
14    emit(loadI, n,  $\_$ , r)
15  return r

```

Arithmetic Operators

Example Invoking the routine `expr` on the AST for $a - b \times c$ shown on the left produces the results shown on the right.



```

loadI  @a      ⇒ r1
loadAO rarp, r1 ⇒ r2
loadI  @b      ⇒ r3
loadAO rarp, r3 ⇒ r4
loadI  @c      ⇒ r5
loadAO rarp, r5 ⇒ r6
mult   r4, r6 ⇒ r7
add    r2, r7 ⇒ r8
  
```

The example assumes that a , b , and c are not already in registers and that each resides in the current AR.

Reducing Demand for Registers

The choice of storage locations has a direct impact on the **quality** of the generated code

- if `a` was in a global data area, we need another `loadI` and register to load it
- if `a` was already in a register, the two instructions to load it could be omitted

Code-shape decisions encoded into the treewalk code generator have an effect on demand for registers

- the naive code uses **eight** registers plus r_{arp}
- a register allocator could rewrite the code as shown on the right
- the rewritten code uses **three** registers plus r_{arp}

<code>loadI</code>	<code>@a</code>	\Rightarrow	r_1
<code>loadAO</code>	r_{arp}, r_1	\Rightarrow	r_1
<code>loadI</code>	<code>@b</code>	\Rightarrow	r_2
<code>loadAO</code>	r_{arp}, r_2	\Rightarrow	r_2
<code>loadI</code>	<code>@c</code>	\Rightarrow	r_3
<code>loadAO</code>	r_{arp}, r_3	\Rightarrow	r_3
<code>mult</code>	r_2, r_3	\Rightarrow	r_2
<code>add</code>	r_1, r_2	\Rightarrow	r_2

Reducing Demand for Registers

A different code shape can reduce the demand for registers. Using a **right-to-left** tree walk instead of a **left-to-right** tree walk produces the code shown below on the left.

```
loadI  @c      ⇒ r1
loadAO rarp, r1 ⇒ r2
loadI  @b      ⇒ r3
loadAO rarp, r3 ⇒ r4
mult   r2, r4 ⇒ r5
loadI  @a      ⇒ r6
loadAO rarp, r6 ⇒ r7
add    r7, r5 ⇒ r2
```

```
loadI  @c      ⇒ r1
loadAO rarp, r1 ⇒ r1
loadI  @b      ⇒ r2
loadAO rarp, r2 ⇒ r2
mult   r1, r2 ⇒ r1
loadI  @a      ⇒ r2
loadAO rarp, r2 ⇒ r2
add    r2, r1 ⇒ r1
```

After register allocation, the code shown above on the right only uses **two** registers plus r_{arp} .

Reducing Demand for Registers

Of course, right-to-left evaluation is not a general solution. To generate an evaluation order that reduces demand for registers, the compiler must choose between right and left children.

Rule

The compiler can **minimize** register use by evaluating first, at each node, the subtree that needs the most registers.

This approach requires **two** passes over the AST

- the first pass computes the demand for registers
- the second pass emits the actual code

Accessing Parameter Values

Formal parameters need may a different treatment

- a **call-by-value** parameter passed in the AR can be treated as if it was a local variable
- a **call-by-reference** parameter passed in the AR requires one additional indirection

Example For the call-by-reference parameter d , the compiler might generate the following instructions to obtain d 's value.

```
loadI    @d           $\Rightarrow$   $r_1$   
loadAO    $r_{arp}$ ,  $r_1$   $\Rightarrow$   $r_2$   
load      $r_2$           $\Rightarrow$   $r_3$ 
```

The first two instructions move the address of d into r_2 , the third instruction loads its value.

Accessing Parameter Values

Many linkage conventions pass the first few parameters in registers. So far, our code generation algorithm cannot handle a value that is **permanently** kept in a register.

The necessary extensions are easy to implement

- **Call-by-Value Parameters**

The ident case must check if the value is already in a register

- if so, it just assigns the register number to `r`
- otherwise, it uses the standard mechanisms to load the value from memory

- **Call-by-Reference Parameters**

- if the address resides in a register, simply load the value into a register
- if the address resides in the AR, must load the address before loading the value

The Small Print Note that the compiler cannot keep the value of a call-by-reference parameter in a register across an assignment, unless the compiler can prove that the reference is unambiguous, across all calls to the procedure.

Function Calls in an Expression

Apart from variables, constants, and temporary values produced by other subexpressions, **function calls** also occur as operands in expressions.

To evaluate a function call, the compiler performs the following steps

- generate the calling sequence needed to invoke the function
- emit the code necessary to move the returned value to a register

Recall The linkage convention limits the callee's impact on the caller.

The presence of a function call may restrict the compiler's ability to **change** the evaluation order of an expression

- function may have **side effects** that modify values of variables used in the expression
- without further analysis, compiler must emit code that is correct in the **worst case**

Other Arithmetic Operators

To handle other arithmetic operations, we can extend the treewalk model.

The basic scheme remains the same

- get the operands into registers
- perform the operation
- store the result

Note Operator precedence, encoded into the expression grammar, ensures the correct evaluation order.

Some operators require complex **multioperation sequences** for their implementation, e.g., exponentiation and trigonometric functions

- expand operation sequence inline
- emit call to a library routine

Mixed-Type Expressions

Many programming languages support operations with operands of **different types**. The compiler must recognize this situation and insert the **conversion code**.

Built-in Types

- definition of programming language specifies a formula for each conversion
- some processors provide explicit conversion operators
- others expect the compiler to generate complex, machine-dependent code

Programmer-defined Types

- compiler has no conversion tables that define each specific case
- source language still defines the meaning of the expression
- compiler must implement this meaning and reject illegal expressions

Assignment as an Operator

Most Algol-like languages implement assignment with the following simple rules

1. evaluate the right-hand side of the assignment to a **value**
2. evaluate the left-hand side of the assignment to a **location**
3. store the right-hand side value into the left-hand side location

Distinguishing between these modes of evaluation

- **rvalue**: result (value) of evaluation on the right-hand side of an assignment
- **lvalue**: result (address) of evaluation on the left-hand side of an assignment

In an assignment, the type of the lvalue can differ from the type of the rvalue, which may require either a compiler-inserted **conversion** or an **error message**.

Note The typical rule for such a conversion is to evaluate the rvalue to its natural type and then convert the result to the type of the lvalue.

Boolean and Relational Operators

Most programming languages can operate on the results of **Boolean** and **relational operators**, both of which produce Boolean values.

To support such values, a compiler writer must...

- augment standard expression grammar with Boolean and relational operators
- define and enforce typing and inference rules for these operators
- decide how to represent Boolean values and how to compute them

Most architectures provide a **rich set** of Boolean operations, but support for relational operators **varies widely** from one architecture to another.

The compiler writer must find an evaluation strategy that matches the needs of the language to the available instruction set.

Boolean and Relational Operators

0	$Expr \rightarrow Expr \vee AndTerm$
1	$AndTerm$
2	$AndTerm \rightarrow AndTerm \wedge RelExpr$
3	$RelExpr$
4	$RelExpr \rightarrow RelExpr < NumExpr$
5	$RelExpr \leq NumExpr$
6	$RelExpr = NumExpr$
7	$RelExpr \neq NumExpr$
8	$RelExpr \geq NumExpr$
9	$RelExpr > NumExpr$
10	$NumExpr$

11	$NumExpr \rightarrow NumExpr + Term$
12	$NumExpr - Term$
13	$Term$
14	$Term \rightarrow Term \times Value$
15	$Term \div Value$
16	$Factor$
17	$Value \rightarrow \neg Factor$
18	$Factor$
19	$Factor \rightarrow (Expr)$
20	num
21	$ident$

Representations

Traditionally, two representations have been proposed for boolean values.

Numerical Encoding

- assign specific values to `true` and `false`
- manipulate them using the target machine's arithmetic and logical operations

Positional Encoding

- encode the value of the expression as a position in the executable code
- use comparisons and conditional branches to evaluate the expression
- different control-flow paths represent the result of evaluation

→ Each approach works well for some examples, but not for others...

Numerical Encoding

When the program **stores** the result of a Boolean or relational operation into a variable, the compiler must assign **numerical values** to true and false.

Typical Values that work with hardware operations such as and, or, and not

- false: zero
- true: one, word of ones, or \neg false

Example If b, c, and d are all in registers, the compiler might produce the following code for the expression $b \vee c \wedge \neg d$.

```
not  rd       $\Rightarrow$  r1
and  rc, r1  $\Rightarrow$  r2
or   rb, r2  $\Rightarrow$  r3
```

Numerical Encoding

For a **comparison**, such as $a < b$, the compiler must generate code that compares a and b and assigns the appropriate value to the result.

1. the target machine provides a comparison operation that **returns a Boolean**

$$\text{cmp_LT } r_a, r_b \Rightarrow r_1$$

2. the comparison defines a **condition code** that must be read with a branch

```

      comp      r_a, r_b  ⇒  cc1
      cbr_LT    cc1      →  L1, L2
L1: loadI     true      ⇒  r1
      jumpI           →  L3
L2: loadI     false     ⇒  r1
      jumpI           →  L3
L3: nop

```


Condition Codes on Intel x86-64

Nearly all arithmetic instructions set condition codes based on their result.

ZF result was **z**ero

CF result caused **c**arry out of most significant bit

SF result was negative (**s**ign bit was set)

OF result caused (signed) **o**verflow

Based on these condition codes, standard condition suffixes *cc* are defined (*cf.* Slide 349) that modify the behaviour of three different kinds of instructions.

jcc jumps to the specified label if *cc* holds

setcc sets the given (single-byte) register to 1 or 0 depending on whether *cc* holds or not

cmovcc performs the specified move only if *cc* holds

The table on the right shows the **standard condition suffixes** defined by the Intel x86-64 instruction set.

Example The expression $a < b$ could be implemented as follows.

```
1 cmp    %rdi, %rsi
2 setl   %al
```

Note Since `setl` needs a single-byte register, the example uses `%al` instead of `%rax`.

cc	Condition Tested	Meaning
e	ZF	equal to zero
ne	\neg ZF	not equal to zero
s	SF	negative
ns	\neg SF	not negative
g	$\neg(\text{SF} \oplus \text{OF}) \wedge \neg \text{ZF}$	greater (signed $>$)
ge	$\neg(\text{SF} \oplus \text{OF})$	greater or equal (signed \geq)
l	$\text{SF} \oplus \text{OF}$	less (signed $<$)
le	$(\text{SF} \oplus \text{OF}) \vee \text{ZF}$	less or equal (signed \leq)
a	$\neg \text{CF} \wedge \neg \text{ZF}$	above (unsigned $>$)
ae	$\neg \text{CF}$	above or equal (unsigned \geq)
b	CF	below (unsigned $<$)
be	$\text{CF} \vee \text{ZF}$	below or equal (unsigned \leq)

Positional Encoding

Positional encoding makes sense if an expression's result is **never** stored

- result of subexpression evaluations
- expressions to determine control flow

Short-Circuit Evaluation

In **short-circuit evaluation**, expressions are only evaluated until their final value is determined. Short-circuit evaluation relies on two Boolean identities.

$$\forall x \text{ false} \wedge x = \text{false}$$

$$\forall x \text{ true} \vee x = \text{true}$$

Some programming languages, e.g., C and Java, **require** the compiler to use short-circuit evaluation.

Positional Encoding

Example Consider the following code for the expression $a < b \vee c < d \wedge e < f$ that a naive code generator would emit.

comp $r_a, r_b \Rightarrow cc_1$	L ₅ : loadI false $\Rightarrow r_2$
cbr_LT $cc_1 \rightarrow L_1, L_2$	jumpI $\rightarrow L_6$
L ₁ : loadI true $\Rightarrow r_1$	L ₆ : comp $r_e, r_f \Rightarrow cc_3$
jumpI $\rightarrow L_3$	cbr_LT $cc_3 \rightarrow L_7, L_8$
L ₂ : loadI false $\Rightarrow r_1$	L ₇ : loadI true $\Rightarrow r_3$
jumpI $\rightarrow L_3$	jumpI $\rightarrow L_9$
L ₃ : comp $r_c, r_d \Rightarrow cc_2$	L ₈ : loadI false $\Rightarrow r_3$
cbr_LT $cc_2 \rightarrow L_4, L_5$	jumpI $\rightarrow L_9$
L ₄ : loadI true $\Rightarrow r_2$	L ₉ : and $r_2, r_3 \Rightarrow r_4$
jumpI $\rightarrow L_6$	or $r_1, r_4 \Rightarrow r_5$

Note Every path takes **eleven** operations, including **three** branches and **three** jumps!

Positional Encoding

If the compiler avoids storing intermediate results of subexpressions and uses short-circuit evaluation, it can emit **much more efficient** code.

```

      comp    ra, rb ⇒ cc1
      cbr_LT  cc1   → L3, L1
L1:  comp    rc, rd ⇒ cc2
      cbr_LT  cc2   → L2, L4
L2:  comp    re, rf ⇒ cc3
      cbr_LT  cc3   → L3, L4
L3:  loadI   true   ⇒ r5
      jumpI           → L5
L4:  loadI   false  ⇒ r5
      jumpI           → L5
L5:  nop

```

Positional Encoding

When the code uses the result of an expression to determine control flow, positional encoding often avoids extraneous operations.

```

1 if (a < b)
2 {
3     statement1;
4 }
5 else
6 {
7     statement2;
8 }

```

```

      comp    ra, rb ⇒ cc1
      cbr_LT  cc1    → L1, L2
L1: code for statement1
      jumpI           → L3
L2: code for statement2
      jumpI           → L3
L3: nop

```

Note The code combines the evaluation of $a < b$ with the selection between *statement*₁ and *statement*₂. The result of $a < b$ is represented as a position, either L₁ or L₂.

Hardware Support for Relational Operations

Specific, low-level details in the target machine's instruction set strongly influence the choice of a representation for relational values.

We will consider four schemes for supporting relational expressions

- straight condition codes
- condition codes augmented with a conditional move operation
- boolean-valued comparisons
- predicated operations

For each of these, we will examine the implementation of an `if-then-else` statement and an assignment of a Boolean value.

Note The first two schemes are supported by the Intel x84-64 instruction set.

Source Code	<pre> if (x < y) then a ← c + d else a ← e + f </pre>	
ILOC Code	<pre> comp r_x, r_y ⇒ cc₁ cbr_LT cc₁ → L₁, L₂ L₁: add r_c, r_d ⇒ r_a jumpI → L_{out} L₂: add r_e, r_f ⇒ r_a jumpI → L_{out} L_{out}: nop Straight Condition Codes </pre>	<pre> cmp_LT r_x, r_y ⇒ r₁ cbr r₁ → L₁, L₂ L₁: add r_c, r_d ⇒ r_a jumpI → L_{out} L₂: add r_e, r_f ⇒ r_a jumpI → L_{out} L_{out}: nop Boolean Compare </pre>
	<pre> comp r_x, r_y ⇒ cc₁ add r_c, r_d ⇒ r₁ add r_e, r_f ⇒ r₂ i2i_LT cc₁, r₁, r₂ ⇒ r_a Conditional Move </pre>	<pre> cmp_LT r_x, r_y ⇒ r₁ not r₁ ⇒ r₂ (r₁)? add r_c, r_d ⇒ r_a (r₂)? add r_e, r_f ⇒ r_a Predicated Execution </pre>

Source Code	$x \leftarrow a < b \wedge c < d$	
ILOC Code	<pre> comp r_a, r_b ⇒ cc₁ cbr_LT cc₁ → L₁, L₂ L₁: comp r_c, r_d ⇒ cc₂ cbr_LT cc₂ → L₃, L₂ L₂: loadI false ⇒ r_x jumpI → L_{out} L₃: loadI true ⇒ r_x jumpI → L_{out} L_{out}: nop Straight Condition Codes </pre>	<pre> comp r_a, r_b ⇒ cc₁ i2i_LT cc₁, r_T, r_F ⇒ r₁ comp r_c, r_d ⇒ cc₂ i2i_LT cc₂, r_T, r_F ⇒ r₂ and r₁, r₂ ⇒ r_x Conditional Move </pre>
		<pre> cmp_LT r_a, r_b ⇒ r₁ cmp_LT r_c, r_d ⇒ r₂ and r₁, r₂ ⇒ r_x Boolean Compare </pre>
		<pre> cmp_LT r_a, r_b ⇒ r₁ cmp_LT r_c, r_d ⇒ r₂ and r₁, r₂ ⇒ r_x Predicated Execution </pre>

Hardware Support for Relational Operations

Condition Codes

- code has at least one conditional branch per relational operator
- comparison operation may be omitted if condition codes are set by default

Conditional Move

- leads to faster code by avoiding branches
- safe as long as neither operation can raise an exception

Boolean Compare

- works without a branch and without converting comparison results to Boolean values
- a weakness of this model is that it requires explicit comparisons

Predicated Execution

- code is simple and concise
- predication can lead to the same code as the boolean-comparison scheme

Relational Operations on Intel x86-64

The Intel x86-64 instruction set supports both the **straight condition codes** and **conditional move** implementation scheme.

Straight condition codes are typically used to implement control-flow constructs, whereas conditional moves are used to evaluate and assign expressions.

A compiler for Intel x86-64 might implement the two examples from before as follows.

```
1      cmpq    %rdi, %rsi
2      jl     _L1
3      leaq    (%rdx, %rcx), %rax
4      jmp     _L2
5 _L1:  leaq    (%r8, %r9), %rax
6 _L2:  nop
```

```
1      cmpq    %rdi, %rsi
2      setl    %al
3      cmpq    %rdx, %rcx
4      setl    %bl
5      and     %al, %bl
6      movzbq  %bl, %rax
```

Storing and Accessing Arrays

So far, we have assumed that variables stored in memory contain scalar values

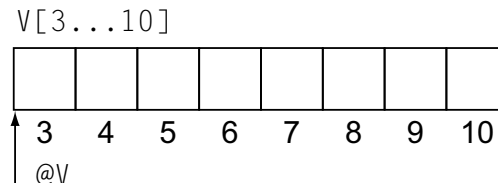
- many programs need arrays or similar structures
- locating and referencing an element of an array can be **surprisingly complex**

We will begin by looking at **one-dimensional** arrays, *i.e.*, vectors, and then generalize the approach to **multi-dimensional** arrays.

Referencing a Vector Element

One-dimensional arrays, or vectors, are typically stored in **contiguous memory**, so that the i^{th} element immediately precedes the $i + 1^{\text{st}}$ element.

A vector $V[3 \dots 10]$ leads to the following memory layout, where the number below a cell indicates its index in the vector.



When the compiler encounters a reference, it must use the **index** into the vector, along with **facts** available from the declaration of V , to generate an offset for the reference.

The actual address is then computed as the **sum** of the offset and a pointer to the start of V , which we write as $@V$.

Referencing a Vector Element

Assume that a vector V has been declared as $V[low...high]$, where *low* and *high* are the vector's lower and upper bounds.

To translate the reference $V[i]$, the compiler needs both a pointer to the start of storage for V and the offset of element i within V .

The offset is simply $(i - low) \times w$, where w is the length of a single element of V .

Example If w is 4, the offset of $V[6]$ in $V[3...10]$ is: $(6 - 3) \times 4 = 12$.

Referencing a Vector Element

Assuming that r_i holds the value of i , the following code fragment computes the address of $V[i]$ into r_3 and loads its value into r_v .

loadI	@V	\Rightarrow	$r_{@v}$...get V's address
subI	$r_i, 3$	\Rightarrow	r_1	...offset — lower bound
multI	$r_1, 4$	\Rightarrow	r_2	... \times element length (4)
add	$r_{@v}, r_2$	\Rightarrow	r_3	...address of $V[i]$
load	r_3	\Rightarrow	r_v	...value of $V[i]$

Notice that the simple reference $V[i]$ introduces **three** arithmetic operations. The compiler can improve this sequence.

Referencing a Vector Element

If w is a power of two, the multiply can be replaced with an **arithmetic shift**. Many base types in real programming languages have this property.

Even though adding the address and offset is unavoidable, most processors include an **addressing mode** specifically designed for this use case.

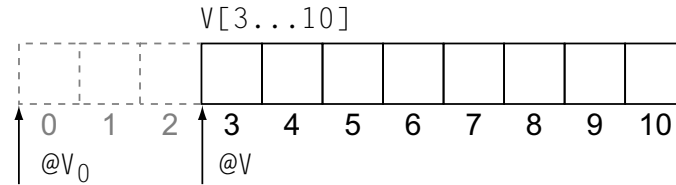
loadI	@V	\Rightarrow	$r_{@v}$...get V's address
subI	$r_i, 3$	\Rightarrow	r_1	...offset — lower bound
lshiftI	$r_1, 2$	\Rightarrow	r_2	... \times element length (4)
loadAO	$r_{@v}, r_2$	\Rightarrow	r_v	...value of $V[i]$

Example In the Intel x86-64 ISA, `move` has an **offset-scaled-base-relative** addressing mode. The instruction `movq -16(%rbx, %rcx, 8), %rax` loads the value at the address $-16 + \%rbx + \%rcx \times 8$ into register `%rax`.

Referencing a Vector Element

Using a lower bound of zero eliminates the subtraction

- if the compiler knows the lower bound of V , it can fold the subtraction into $@V$
- instead of using $@V$ as the base address for V , it can also use $V_0 = @V - low \times w$



False Zero

The **false zero** of a vector V is the address where $V[0]$ would be. In multiple dimensions, it is the location of a zero in each dimension.

Referencing a Vector Element

Using $@V_0$ and assuming that i is in r_i , the code for accessing $V[i]$ becomes shorter and, presumably, faster.

loadI	$@V_0$	$\Rightarrow r_{@V_0}$...adjusted address of V
lshiftI	$r_i, 2$	$\Rightarrow r_1$... \times element length (4)
loadAO	$r_{@V_0}, r_1$	$\Rightarrow r_v$...value of $V[i]$

Note In a compiler, the longer sequence may produce better results by exposing details such as the multiply and add instruction to optimization.

An alternative strategy, employed in languages like C or Java, forces the use of zero as a lower bound, which ensures that $@V_0 = @V$ and simplifies all array-address calculations.

Array Storage Layout

Accessing an element of a multi-dimensional array requires more work. The code that the compiler needs to generate depends on the mapping of array indices to memory locations.

Most implementations use one of three schemes

- row-major order
- column-major order
- indirection vectors

The source-language definition usually specifies one of these mappings.

Array Storage Layout

Example Consider the array $A[1 \dots 2, 1 \dots 4]$. Conceptually, it looks as follows.

A	1,1	1,2	1,3	1,4
	2,1	2,2	2,3	2,4

In **row-major order**, the elements of a are mapped onto consecutive memory locations so that adjacent elements of a single row occupy consecutive memory locations.

1,1	1,2	1,3	1,4	2,1	2,2	2,3	2,4
-----	-----	-----	-----	-----	-----	-----	-----

Array Storage Layout

The array storage layout has an impact on **memory access patterns** and, thus, run-time performance.

Example Consider the following nested loops.

```
for i ← 1 to 2
  for j ← 1 to 4
    A[i,j] ← A[i,j] + 1
```

In row-major order, the assignment statement steps through memory in **sequential order**, beginning with $A[1,1]$, $A[1,2]$, $A[1,3]$, and on through $A[2,4]$.

Switching the two loops produces an access pattern that jumps between rows. For arrays larger than the cache, lack of sequential access can lead to poor run-time performance.

Array Storage Layout

The obvious alternative to row-major order is **column-major order**. It keeps the columns of a in contiguous locations, producing the following layout.

1,1	2,1	1,2	2,2	1,3	2,3	1,4	2,4
-----	-----	-----	-----	-----	-----	-----	-----

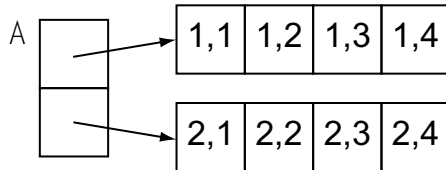
While row-major order leads to sequential access if the **rightmost** subscript varies fastest, column-major order leads to sequential access if the **leftmost** subscript varies fastest.

Row-major vs. column-major order

For languages that store arrays in contiguous storage, row-major order has been the typical choice. The one notable exception is FORTRAN, which uses column-major order.

Array Storage Layout

A third alternative are **indirection vectors**. This scheme reduces all multi-dimensional arrays to a set of (one-dimensional) vectors.



Indirection vectors appear simple, but they introduce their own complexity

- need more storage than either of the contiguous storage schemes
- require that application initializes all of the indirection pointers at runtime

An advantage of indirection vectors is that they can easily represent **ragged arrays**.

Note Java supports indirection vectors.

Referencing an Array Element

Programs that use arrays typically contain **references** to individual array elements.

As with vectors, the compiler must translate an array reference into a base address for the array's storage and an offset where the element is located relative to the starting address.

We will look at the following cases

- row-major order
- indirection vectors
- array-valued parameters

Note The calculations for column-major order follow the same basic scheme as those for row-major order, with the dimensions reversed.

Row-Major Order

Before deriving the corresponding formula, we introduce the following notation

- low_i and $high_i$ denote *low* and *high* of the i -th dimension, respectively
- $len_i = high_i - low_i + 1$ denotes the length of the i -th dimension

In row-major order, the address calculation must find the start of the row and then generate an offset within the row as if it were a vector.

To access element $A[i, j]$ of a two-dimensional array the compiler therefore must emit code that computes the...

- address of row i : $(i - low_1) \times len_2 \times w$
- offset of element j : $(j - low_2) \times w$

Putting all parts together, the resulting address computation is as follows.

$$@A + (i - low_1) \times len_2 \times w + (j - low_2) \times w$$

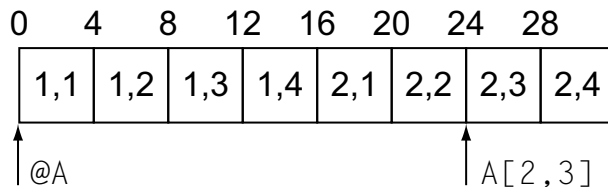
Row-Major Order

Example In the array $A[1 \dots 2, 1 \dots 4]$, element $A[2, 3]$ lies at offset

$$(2-1) \times (4-1+1) \times 4 + (3-1) \times 4 = 24$$

from `A[1,1]`, assuming that `@A` points at `A[1,1]` at offset 0.

Looking at A in memory, we find that the address of $A[1,1] + 24$ is, in fact, the address of $A[2,3]$.



Row-Major Order

As in the vector case, we can simplify the address calculation if upper and lower bounds are known at compile time by introducing a false zero.

$$@A + (i - low_1) \times len_2 \times w + (j - low_2) \times w$$

$$@A + (i \times len_2 \times w) - (low_1 \times len_2 \times w) + (j \times w) - (low_2 \times w)$$

$$@A + (i \times len_2 \times w) + (j \times w) - (low_1 \times len_2 \times w + low_2 \times w)$$

The term $(low_1 \times len_2 \times w + low_2 \times w)$ is independent of i and j . Thus, it can be factored directly into the base address: $@A_0 = @A - (low_1 \times len_2 \times w + low_2 \times w) = @A - 20$.

Now, the array reference is simply $@A_0 + i \times len_2 \times w + j \times w$. Finally, we can refactor and move w outside, saving an extraneous multiplication.

$$@A_0 + (i \times len_2 + j) \times w$$

Row-Major Order

Assuming that i and j are in r_i and r_j , and that len_2 is a constant, this form of the polynomial leads to the following code sequence.

loadI	@A ₀	\Rightarrow	$r_{@A_0}$...adjusted base address of A
multI	r_i, len_2	\Rightarrow	r_1	... $i \times len_2$
add	r_1, r_j	\Rightarrow	r_2	...+ j
multI	$r_2, 4$	\Rightarrow	r_3	... \times element length (4)
loadA0	$r_{@A_0}, r_3$	\Rightarrow	r_a	...value of $A[i, j]$

Note The computation takes two multiplications and two additions (one in the loadA0). The second multiply can be rewritten as a shift.

Row-Major Order

If the compiler does not have access to the array bounds

- compute the false zero at runtime
- use the more complex polynomial with subtractions to adjust for lower bounds

The first option only makes sense if the elements of the array are accessed multiple times in a procedure.

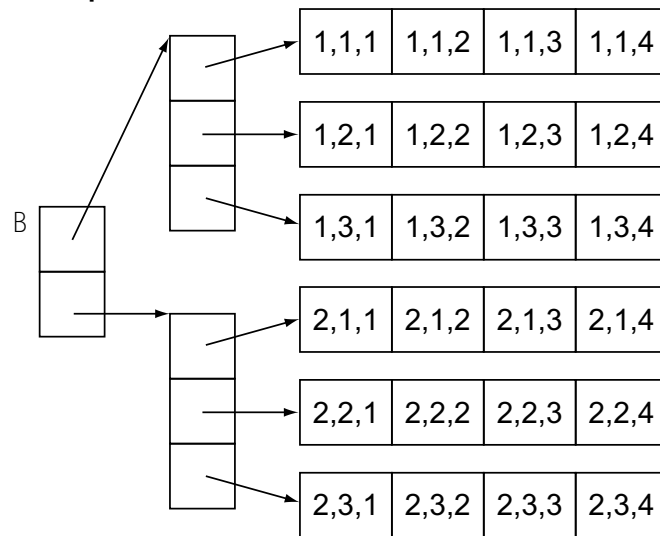
Note These ideas behind address computation for two-dimensional arrays generalize to **arrays of higher dimension**. The address polynomial for an array stored in **column-major order** can be derived in a similar fashion.

Indirection Vectors

Using indirection vectors simplifies the code generated to access an individual element. For an n -dimensional array, the first $n - 1$ steps navigate the indirection vectors, while the last step performs a regular (one-dimensional) vector lookup.

Example Accessing element $B[i, j, k]$ in the array $B[1 \dots 2, 1 \dots 3, 1 \dots 4]$

1. use $@B0$, i , and the length of a pointer, to find the vector for the subarray $B[i, *, *]$
2. use result, along with j and the length of a pointer to find the vector for the subarray $B[i, j, *]$
3. use base address in the vector-address computation with k and element length w to find the address of $B[i, j, k]$



Indirection Vectors

If the values for i , j , and k exist in registers r_i , r_j , and r_k , respectively, and $@B_0$ is the zero-adjusted address of the first dimension, then $B[i, j, k]$ can be loaded as follows.

loadI	$@B_0$	\Rightarrow	$r_{@B_0}$...adjusted base address of B
multI	$r_i, 4$	\Rightarrow	r_1	...pointer size (4)
loadAO	$r_{@B_0}, r_1$	\Rightarrow	r_2	...load $@B[i, *, *]$
multI	$r_j, 4$	\Rightarrow	r_3	...pointer size (4)
loadAO	r_2, r_3	\Rightarrow	r_4	...load $@B[i, j, *]$
multI	$r_k, 4$	\Rightarrow	r_5	...element size (4)
loadAO	r_4, r_5	\Rightarrow	r_b	...value of $B[i, j, k]$

Note This code assumes that the pointers in the indirection structure have already been adjusted to account for non-zero lower bounds.

Accessing Array-Valued Parameters

Note Even in languages that use call by value for all other parameters, arrays are usually passed by reference, for the sake of runtime performance.

To generate array references in the callee, we need information about the dimensions of the array that is bound to the parameter.

- some languages, e.g., FORTRAN, Oberon-0, give the **programmer** responsibility for passing the information to the callee needed to address a parameter array correctly
- other languages leave the task of collecting, organizing, and passing the necessary information to the **compiler**

Accessing Array-Valued Parameters

Dope Vector

A descriptor for an actual parameter array that contains both a pointer to the **start of the array** and the **necessary information for each dimension**. Dope vectors may also be used for arrays whose bounds are determined at runtime.

Caller

- builds the dope vector and stores it in callee's AR
- value passed in the array's parameter slot is a pointer to the dope vector

Callee

- loads values out of the dope vector as needed
- generates array reference using the same polynomial as for a local array

Range Checking

A program that references an out-of-bounds array element is not well formed

- some languages, e.g., Ada, Java, require such accesses to be detected and reported
- compilers for other languages may include optional mechanisms for range checking

The simplest implementation of **range checking** inserts a check before each array access

- introduces additional overhead, which may be significant in some programs
- may need more information in dope vector, e.g., explicit bounds rather than length

An optimizing compiler can improve on this naive implementation in several ways

- combine checks
- move checks out of loops
- prove checks redundant

Character Strings

The operations that programming languages provide for **character data** are different from those provided for numerical data.

Programming language support for character strings varies

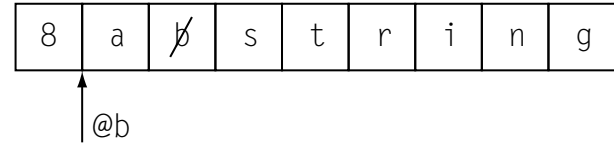
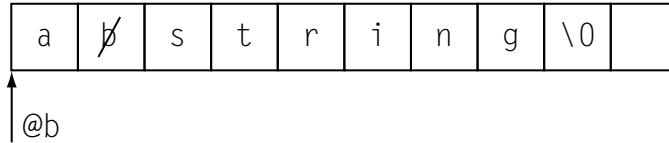
- in C, most manipulation takes the form of calls to library routines
- PL/I provides first-class mechanisms to assign individual characters, specify arbitrary substrings, and concatenate strings

CISC architectures provide extensive support for string manipulation, whereas RISC machines rely on the compiler to code these operations using a set of simpler instructions.

- string assignment
- string concatenation
- string-length computation

String Representation

The string representation chosen by the compiler writer has a strong impact on the cost of string operations.



- C traditionally uses **null termination** (left) with a designated character (`\0`) serving as a terminator
- many other language implementations use an **explicit length field** (right) that stores the length of the string (8) alongside its contents

If the length field takes more space than the null terminator, then storing the length will marginally increase the size of the string in memory. However, storing the length simplifies several operations on strings.

String Assignment

String assignment is conceptually simple. The complexity of its implementation depends on whether or not the target machine supports **character-sized memory operations**.

Example In C, an assignment from the third character of `b` to the second character of `a` can be written as `a[1] = b[2];`.

With support for character-sized memory operations, this assignment translates into the following simple code.

```
loadI    @b      ⇒ r@b
cloadAI  r@b, 2  ⇒ r1
loadI    r@a     ⇒ r@a
cstoreAI r1     ⇒ r@a, 1
```

Without hardware support for character-sized memory operations, the compiler must generate more complex code that **masks** the characters explicitly.

String Assignment

Assumptions Strings a and b begin on word boundaries, a character occupies 1 byte, and that a word is 4 bytes.

loadI	0x0000FF00	\Rightarrow	r_{C2}	...mask for the character 2
loadI	0xFF00FFFF	\Rightarrow	r_{C124}	...mask for characters 1, 2, and 4
loadI	@b	\Rightarrow	$r_{@b}$...load address of b
load	$r_{@b}$	\Rightarrow	r_1	...load first word of b
and	r_1, r_{C2}	\Rightarrow	r_2	...mask away other characters
lshiftI	$r_2, 8$	\Rightarrow	r_3	...move character to first byte
loadI	@a	\Rightarrow	$r_{@a}$...load address of a
load	$r_{@a}$	\Rightarrow	r_4	...load first word of a
and	r_4, r_{C124}	\Rightarrow	r_5	...mask away second character
or	r_3, r_5	\Rightarrow	r_6	...insert new second character
store	r_6	\Rightarrow	$r_{@a}$...store result

String Assignment

The code is similar for longer strings. The following code for `a = b`; assumes explicit length representation and hardware support for character-sized memory operations.

	loadI	@b	⇒	r _{@b}	...load address of b
	loadAI	r _{@b} , -4	⇒	r ₁	...load b's length
	loadI	@a	⇒	r _{@a}	...load address of a
	loadAI	r _{@a} , -4	⇒	r ₂	...load a's length
	cmp_LT	r ₂ , r ₁	⇒	r ₃	...does b fit into a?
	cbr	r ₃	→	L _{SOV} , L ₁	...raise string overflow
L ₁ :	loadI	0	⇒	r ₄	...initialize counter
	cmp_LT	r ₄ , r ₁	⇒	r ₅	...more to copy?
	cbr	r ₅	→	L ₂ , L ₃	...jump
L ₂ :	cloadA0	r _{@b} , r ₄	⇒	r ₆	...load character from b
	cstoreA0	r ₆	⇒	r _{@a} , r ₄	...store character in a
	addI	r ₄ , 1	⇒	r ₄	...increment counter
	cmp_LT	r ₄ , r ₁	⇒	r ₇	...more to copy?
	cbr	r ₇	→	L ₂ , L ₃	...jump
L ₃ :	storeAI	r ₁	⇒	r _{@a} , -4	...set length of a

String Assignment

In C, which uses null termination for strings, the same assignment would be written as a character-copying loop.

```

1 t1 = a;
2 t2 = b;
3 do {
4     *t1++ = *t2++;
5 } while (*t2 != '\0')

```

	loadI	@b	\Rightarrow	$r_{@b}$...load address of b
	loadI	@a	\Rightarrow	$r_{@a}$...load address of a
	loadI	NULL	\Rightarrow	r_1	...load terminator
L_1 :	cload	$r_{@b}$	\Rightarrow	r_2	...load character
	cstore	r_2	\Rightarrow	$r_{@a}$...store character
	addI	$r_{@b}, 1$	\Rightarrow	$r_{@b}$...move pointer
	addI	$r_{@a}, 1$	\Rightarrow	$r_{@a}$...move pointer
	cmp_NE	r_1, r_2	\Rightarrow	r_3	...more to copy?
	cbr	r_3	\rightarrow	L_1, L_2	...jump
L_2 :	nop				...next statement

String Assignment

Further Possibilities

- if the target machine supports **auto-increment** on load and store operations, the two adds in the loop can be performed in the cload and cstore operations
- the compiler could generate **whole-word** loads and stores, followed by a character-oriented loop to handle any leftover characters at the end of the string
- the compiler might also call a carefully optimized **library routine** to implement the nontrivial cases

String Concatenation

Concatenation is simply a shorthand for a sequence of one or more assignments. It comes in **two basic forms**.

1. Appending string b to string a

- compiler emits code to determine the length of a
- space permitting, it assigns b to the space that immediately follows a

2. Creating a new string that contains a followed immediately by b

- requires copying each character in a and each character in b
- can be treated as a pair of assignments in code generation

The compiler should ensure that **enough space** is allocated to hold the result. If the lengths of a and b are unknown at compile time, runtime code needs to be generated

- to compute the lengths of the strings
- to perform corresponding test and branch

String Length

Programs that manipulate strings often need to compute a character string's length.

Different string representations lead to different costs for the length computation

- a **null-terminated string** requires time proportional to the length of the string
- with an **explicit length field** the cost is constant and small

Tradeoff Null termination saves a small amount of space, but requires more code and more time for the length computation. An explicit length field costs one more word per string, but makes the length computation take constant time.

Example With explicit length fields a statement `length(a + b)` can be optimized to two loads and one add.

String Length on Intel x86-64

The Intel x86-64 ISA provides several **dedicated operations** for string manipulation.

- movs, lods, stos: move, load, and store strings
- cmps: compare strings
- scas: scan strings

These instructions are combined with **repetition prefixes**: rep, repe, repne, repz, and repnz).

```
1 movq    -1, %rcx      # "clear" counting register
2 movq    %rdi, %rsi    # backup %rdi
3 mov     0, %al        # look for \0
4 repne scasb          # actually do the search
5 subq    %rsi, %rdi     # save the string length
6 decq    %rdi          # do not count the \0 in the string length
7 movq    %rdi, %rax     # save the return value
```

Structure References

Most programming languages provide a mechanism to aggregate data together into a **structure**.

Example In C, we could use the following structure to create lists of integers.

```
1 struct node {  
2     int value;  
3     struct node *next;  
4 };  
5 struct node NILNode = { 0, (struct node*) 0 };  
6 struct node *NIL = &NILNode;
```

The introduction of structures and pointers creates two distinct problems for the compiler: **anonymous values** and **structure layout**.

Understanding Structure Layouts

To emit code for structure references, the compiler needs to know

- **starting address** of the structure instance
- **offset** and **length** of each structure element

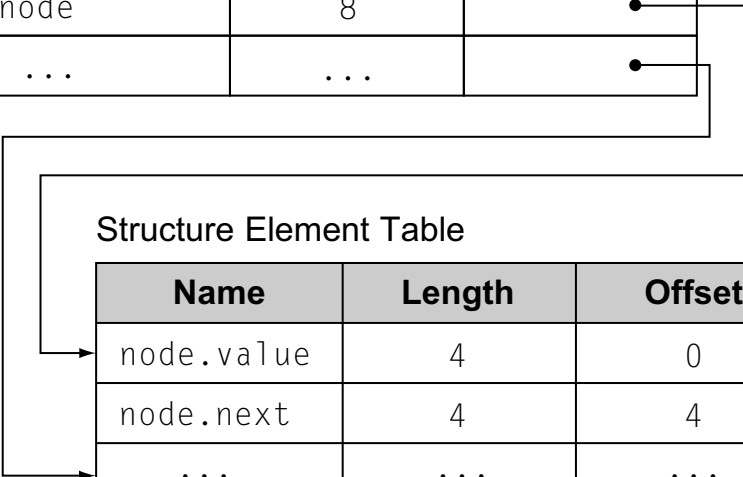
The compiler can build a separate **table of structure layouts** to maintain this information.

- textual name for each structure element
- its offset within the structure
- its source-language data type

Understanding Structure Layouts

Structure Layout Table

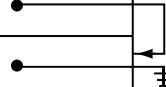
Name	Length	1 st Element
node	8	•
...	...	•



The diagram shows two arrows originating from the '1st Element' column of the Structure Layout Table. One arrow points to the 'node.value' row in the Structure Element Table, and the other points to the 'node.next' row. A third arrow points from the 'node.next' row to a ground symbol (three horizontal lines of decreasing width).

Structure Element Table

Name	Length	Offset	Type	Next
node.value	4	0	int	•
node.next	4	4	struct node *	•
...



The diagram shows two arrows originating from the 'Next' column of the Structure Element Table. One arrow points to the 'node.value' row, and the other points to the 'node.next' row. A third arrow points from the 'node.next' row to a ground symbol (three horizontal lines of decreasing width).

Understanding Structure Layouts

Note Entries in the **structure element table** (shown on the previous slide) use fully qualified names to avoid conflicts due to reuse of a name in several distinct structures.

With this information, the compiler can easily generate code for structure references.

Example The compiler might translate the reference `p1->next`, for a pointer to node `p1`, into the following ILOC code.

```
loadI    4            $\Rightarrow$   $r_1$     ...offset of next  
loadAO    $r_{p1}$ ,  $r_1$   $\Rightarrow$   $r_2$     ...value of p1->next
```

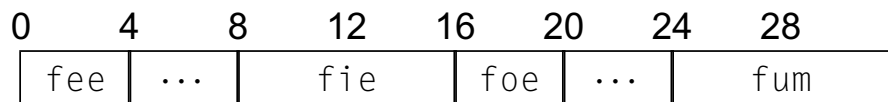

Understanding Structure Layouts

In laying out a structure and assigning offsets to its elements, the compiler must obey the **alignment rules** of the target architecture.

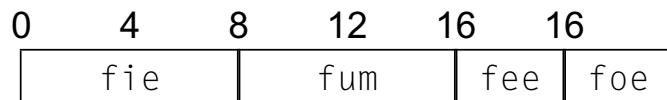
```

1 struct example {
2     int fee;
3     double fie;
4     int foe;
5     double fum;
6 };

```



Elements in Declaration Order



Elements Ordered by Alignment

Note Whether the compiler can reorder structure elements in memory arbitrarily depends whether the source language definition exposes the layout of a structure to the user.

Array of Structures

Many programming languages allow the user to declare an **array of structures**. The compiler has two options how to layout such data in memory.

1. a structure-valued array with multiple copies of the structure layout
2. a structure composed of elements that are arrays

Note The second option is only possible if the programmer cannot take the address of a structure-valued element of an array.

Depending on how the surrounding code accesses the data, these two strategies may have strikingly **different performance** on a system with cache memory.

Unions and Runtime Tags

Many languages support structures with **multiple, data-dependent** interpretations.

Unions and variants present one additional complication

- the possibility exists that element names are **not unique**
- compiler must resolve each reference to a unique offset and type in the runtime object

```
1 struct n1 {  
2     int kind;  
3     int value;  
4 };
```

```
1 struct n2 {  
2     int kind;  
3     float value;  
4 };
```

```
1 union one {  
2     struct n1 inode;  
3     struct n2 fnode;  
4 } u1;
```

This problem has a **linguistic solution**: the programming language can force the programmer to make the reference unambiguous.

Unions and Runtime Tags

Example To reference an integer value, the programmer specifies `u1.inode.value`, whereas to reference a floating-point value, the programmer specifies `u1.fnode.value`.

The **fully qualified name** resolves any ambiguity.

The same mechanism is used to disambiguate references to implicitly defined structure-valued union elements.

Example Union `two` has the same properties as `one`.

Again, the programmer has to specify a fully qualified name such as `u2.inode.value` to reference an integer value and `u2.fnode.value` to reference a floating-point value.

```
1 union two {  
2     struct {  
3         int kind;  
4         int value;  
5     } inode;  
6     struct {  
7         int kind;  
8         float value;  
9     } fnode;  
10 } u2;
```

Unions and Runtime Tags

As an alternative to linguistic solutions, some systems rely on **runtime discrimination**.

- each variant in the union has a field (tag) that distinguishes it from all other variants
- compiler emits code to check the value of the tag field and handle object accordingly

```
1 type shapeKind = (square, rectangle, circle);
2   shape = record
3       centerx : integer;
4       centery : integer;
5       case kind : shapeKind of
6           square : (side : integer);
7           rectangle : (length, height : integer);
8           circle : (radius : integer);
9   end;
```

Pointers and Anonymous Values

A C program creates an instance of a structure in one of two ways

- declare a structure instance, e.g., NILNode in the earlier example
- allocate a structure instance explicitly

Example For a variable declared as a pointer to node, the allocation looks as follows.

```
1 fee = (struct node*) malloc(sizeof(node));
```

The only access to this new node is through the pointer fee. Thus, we think of it as an **anonymous value**, since it has no permanent name.

Because the only name for an anonymous value is a pointer, the compiler cannot easily determine if two pointer references specify the **same** memory location.

Pointers and Anonymous Values

Example Consider the following code fragment.

```
1 p1 = (node*) malloc(sizeof(node));
2 p2 = (node*) malloc(sizeof(node));
3 if (...) {
4     p3 = p1;
5 } else {
6     p3 = p2;
7 }
8 p1->value = ...;
9 p3->value = ...;
10 ... = p1->value;
```

The first two lines create anonymous nodes.

Line 8 writes through p1, while line 9 writes through p3.

On line 9, p3 can refer to either the node allocated in line 1 or in line 2.

Finally, line 10 references p1->value. The compiler cannot easily determine at this point which value is used.

The uncertainty introduced by pointers **prevents** the compiler from keeping values used in pointer-based references in registers.

Pointers and Anonymous Values

Anonymous objects further complicate the problem because they introduce an **unbounded set** of objects to track.

A similar effect occurs with arrays: the compiler would need to perform in-depth analysis of array subscripts to determine whether two array references overlap.

While challenging, the problem of disambiguating array references is easier than the problem of disambiguating pointer references.

Analysis to disambiguate pointer references and array references is a **major** source of potential improvement in program performance.

- interprocedural data-flow analysis to discover all objects each pointer can point to
- data-dependence analysis to understand the patterns of array references

Control-Flow Constructs

As the compiler generates code, it can build up **basic blocks** by simply aggregating consecutive, unlabeled, non-control-flow operations.

To tie a set of blocks together so that they form a procedure, the compiler must insert code that implements the **control-flow operations** of the source program.

- build a **control-flow graph** and use it for analysis, optimization, and code generation
- code for control-flow constructs resides **at or near the end** of each basic block

While many different syntactic conventions have been used to express control flow, the number of underlying concepts is **small**.

Conditional Execution

Most programming languages provide some version of an if-then-else construct.

```
if expr
  then statement1
  else statement2
statement3
```

The compiler needs to generate code that

- evaluates *expr* and branches to *statement*₁ or *statement*₂, accordingly
- implements the two statements and at the end jumps to *statement*₃

As we already saw on Slide 355, the compiler has many options for implementing if-then-else constructs.

Conditional Execution

With trivial `then` and `else` parts the primary consideration for the compiler is matching the expression evaluation to the underlying hardware.

As the `then` and `else` parts grow, the importance of efficient execution inside the `then` and `else` parts begins to **outweigh** the cost of executing the controlling expression.

Example On a machine that supports **predicated execution**, using predicates for large blocks in the `then` and `else` parts can waste execution cycles.

With large blocks of code under both the `then` and `else` parts, the cost of unexecuted instructions may outweigh the overhead of using a conditional branch.

7 Code Shape - 7.7 Control-Flow Constructs - Conditional Execution

Unit 1	Unit 2
<i>comparison</i> $\Rightarrow r_1$	
(r_1) op ₁	$(\neg r_1)$ op ₁₁
(r_1) op ₂	$(\neg r_1)$ op ₁₂
(r_1) op ₃	$(\neg r_1)$ op ₁₃
(r_1) op ₄	$(\neg r_1)$ op ₁₄
(r_1) op ₅	$(\neg r_1)$ op ₁₅
(r_1) op ₆	$(\neg r_1)$ op ₁₆
(r_1) op ₇	$(\neg r_1)$ op ₁₇
(r_1) op ₈	$(\neg r_1)$ op ₁₈
(r_1) op ₉	$(\neg r_1)$ op ₁₉
(r_1) op ₁₀	$(\neg r_1)$ op ₂₀

Unit 1	Unit 2
<i>compare and branch</i>	
L ₁ : op ₁	op ₂
op ₃	op ₄
op ₅	op ₆
op ₇	op ₈
op ₉	op ₁₀
jumpI \rightarrow L ₃	
L ₂ : op ₁₁	op ₁₂
op ₁₃	op ₁₄
op ₁₅	op ₁₆
op ₁₇	op ₁₈
op ₁₉	op ₂₀
jumpI \rightarrow L ₃	
L ₃ : nop	

Conditional Execution

Choosing between branching and predication to implement an `if-then-else` requires some care. Several issues need to be considered.

1. **Expected frequency of execution:** if one path executes significantly more often, techniques that speed up its execution (e.g., branch prediction, speculative execution, and instruction reordering) may produce faster code
2. **Uneven amounts of code:** if one path contains many more instructions, this may weigh against predication or for a combination of predication and branching
3. **Control flow inside the construct:** if either path contains nontrivial control flow, in particular nested `if`-statements, then predication may be a poor choice

Conditional Execution

Predication in the Real World

Nvidia's CUDA parallel computing platform 32 threads are executed together as a **warp**. All threads in a warp execute the **same** instruction at the **same** time.

In order to support **different** threads doing **different** things, CUDA has predicated instructions that are executed only if a logical flag is true.

```
if (x < 0.0)                p = (x < 0.0);  
    r = -1.0;                p:  r = -1.0;  
else                        !p:  r = sqrt(x);  
    r = sqrt(x);
```

This is called **warp divergence**. The performance of the code generated by the nvcc compiler is impacted by the same issues as outlined on the previous slide.

Loops and Iteration

Most programming languages include loop constructs to perform iteration.

The first FORTRAN compiler introduced the do loop to perform iteration. Today, loops are found in many forms. For the most part, they have a similar structure.

Example Consider the C for loop, which has three controlling expressions.

e_1 performs initialization

e_2 evaluates to a Boolean and governs execution of the loop

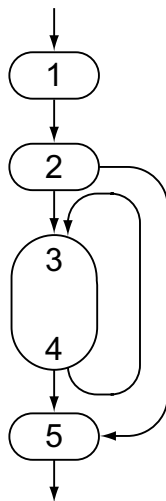
e_3 executes at the end of each iteration and, potentially,
updates the values used in e_2

```
for ( $e_1$ ;  $e_2$ ;  $e_3$ ) {  
    loop body  
}
```

Loops and Iteration

The following schema shows how the compiler might lay out the code.

```
for ( $e_1$ ;  $e_2$ ;  $e_3$ ) {  
    loop body  
}
```



Step	Action
1	evaluate e_1
2	if ($\neg e_2$) then goto 5
3	loop body
4	evaluate e_3 if (e_2) then goto 5
5	code following loop

We will use this basic schema to explain the implementation of several kinds of loops.

Loops and Iteration

If the loop body consists of a single basic block, *i.e.*, contains no other control flow, then the loop that results from this schema has **an initial branch** plus **one branch per iteration**.

The compiler can hide the latency of this branch in several ways

- **branch hint prefix**: mark the branch in Step 4 as likely to be taken
- **delay slot**: move instructions from loop body into the branch delay slot(s)
- **loop unrolling**: if number of iterations is predictable, combine multiple iterations

Delay Slot

Some computer architectures feature **delay slots** to mask the latency of time-consuming instructions, such as loads and branches. A delay slot is an instruction that gets executed without the effects of a preceding instruction.

Loops and Iteration

Branch Prediction in Intel x86-64

In the past, Intel's ISA supported static branch prediction prefixes (2EH and 3EH), a feature that has been abandoned since.

Modern Intel CPUs use a Branch Target Buffer (BTB) to manage branching history dynamically. If this information is not available, the CPU uses a simple heuristic:

- conditional forward branches are predicted as **not being taken**
- conditional backward branches are predicted as **being taken**

This heuristic fits well with the basic schema that we will use to implement loops.

For Loops

To map a for loop into code, the compiler follows the general schema from before.

```

for (i = 1; i <= 100; i++) {
    loop body
}
next statement

```

```

loadI    1      ⇒ ri      ...Step 1
loadI    100    ⇒ r1      ...Step 2
cmp_GT   ri, r1 ⇒ r2
cbr      r2    → L1, L2

L1: loop body      ...Step 3
addI     ri, 1  ⇒ ri      ...Step 4
cmp_LE   ri, r1 ⇒ r3
cbr      r3    → L1, L2

L2: next statement ...Step 5

```

The compiler can also shape the loop to have only **one copy** of the test. In this form, Step 4 evaluates e_3 and then jumps to Step 2, *i.e.*, replace `cmp_LE` and `cbr` with `jumpI`.

For Loops

Overall, this form of the loop is **one operation smaller** than the two-test form.

However, it creates a **two-block loop** for even the simplest loops, and it **lengthens the path** through the loop by at least one operation.

	loadI	1	\Rightarrow	r_i	...Step 1
	loadI	100	\Rightarrow	r_1	...Step 2
L_1 :	cmp_GT	r_i, r_1	\Rightarrow	r_2	
	cbr	r_2	\rightarrow	L_1, L_2	
	<i>loop body</i>				...Step 3
	addI	$r_i, 1$	\Rightarrow	r_i	...Step 4
	jumpI		\rightarrow	L_1	
L_2 :	<i>next statement</i>				...Step 5

This more compact loop form is only appropriate if **code size** is a serious consideration.

While Loops

A `while` loop can also be implemented with the loop schema. Since it has no initialization, the code is even more compact.

```
while (x < y) {
    loop body
}
next statement
```

	<code>cmp_LT</code>	$r_x, r_y \Rightarrow r_1$...Step 2
	<code>cbr</code>	$r_1 \rightarrow L_1, L_2$	
L_1 :	<i>loop body</i>		...Step 3
	<code>cmp_LT</code>	$r_x, r_y \Rightarrow r_2$...Step 4
	<code>cbr</code>	$r_2 \rightarrow L_1, L_2$	
L_2 :	<i>next statement</i>		...Step 5

Replicating the test in Step 4 creates the possibility of a loop with a **single basic block**.

Until Loops

An `until` loop iterates as long as the controlling expression is `false`.

It checks the controlling expression **after** each iteration. Thus, it always enters the loop and performs at least one iteration and produces a particularly simple loop structure.

```
{
    loop body
} until (x < y)
next statement
```

```
L1: loop body                ...Step 3
      cmp_LT rx, ry ⇒ r1      ...Step 4
      cbr    r1    → L1, L2
L2: next statement          ...Step 5
```

Note The `do` loop known from C, C++, and Java is similar to the `until` loop with the difference that it iterates as long as the controlling expression is `true`.

Break and Exit

Several languages provide a **structured way** to exit a control-flow construct.

In a loop, `break` transfers control to the first statement following the loop. For nested loops, a `break` typically exits the innermost loop.

Some languages, e.g., Ada and Java, use labels to control which enclosing constructs will be exited by the `break` statement.

C also uses `break` in `switch` statements to transfer control to the statement that follows the `switch` statement.

All of these actions have simple implementations

- each loop and each case statement ends with a label for the statement that follows it
- a `break` can be implemented as an immediate jump to that label

Skip and Continue

Some languages include a `skip` or `continue` statement that jumps to the next iteration of a loop.

This construct can be implemented in two ways

- **immediate jump** to the code that reevaluates the controlling expression, tests its value, and branches accordingly
- **insert a copy** of the evaluation, test, and branch at the point where the skip occurs

Case Statements

Many programming languages include some variant of a case statement.

The basic strategy is straightforward

1. evaluate the controlling expression
2. branch to the selected case
3. execute the code for that case

Steps 1 and 3 are well understood. The complex part of case-statement implementation lies in choosing an efficient method to locate the designated case.

No single method works well for all case statements. We examine three strategies: a **linear search**, a **binary search**, and a **computed address**.

Linear Search

The simplest way to locate the appropriate case is to treat the case statement as the specification for a nested set of if-then-else statements.

```
switch(e) {  
  case 0:  block0;  
           break;  
  case 1:  block1;  
           break;  
  case 3:  block3;  
           break;  
  default: blockd;  
           break;  
}
```

```
t1 ← e  
if (t1 = 0)  
  then block0  
else if (t1 = 1)  
  then block1  
else if (t1 = 3)  
  then block3  
  else blockd
```

The compiler should order cases according to estimated execution frequency.

Binary Search

As the number of cases rises, the efficiency of linear search becomes a problem.

If the compiler can impose an order on the case labels, it can use **binary search** to obtain a logarithmic search rather than a linear one.

- build a compact ordered table of case labels and corresponding branch labels
- use binary search to discover a matching case label, or the absence of a match
- finally, either branch to the corresponding label or to the default case.

Binary Search

```

switch(e) {
  case 0:    block0;
             break;
  case 15:   block15;
             break;
  case 23:   block23;
             break;
  ...
  case 99:   block99;
             break;
  default:   blockd;
             break;
}

```

```

t1 ← e
down ← 0      ...lower bound
up ← 10       ...upper bound + 1
while (down + 1 < up) {
  middle ← (up + down) ÷ 2
  if (Value[middle] ≤ t1)
    then down ← middle
    else up ← middle
}
if (Value[down] = t1)
  then jump to Label[down]
  else jump to LBd

```

Value	Label
0	LB ₀
15	LB ₁₅
23	LB ₂₃
37	LB ₃₇
41	LB ₄₁
50	LB ₅₀
68	LB ₆₈
72	LB ₇₂
83	LB ₈₃
99	LB ₉₉

Computing the Address Directly

If the case labels form a compact set, the compiler can do better than binary search.

In this case, the compiler can build a compact vector, or **jump table**, that contains the block labels, and find the appropriate label by index into the table.

- for a dense label set, this scheme generates compact and efficient code
- the cost is small and constant: a brief calculation, a memory reference, and a jump
- if a few holes exist in the label set, they can be filled with the label for the default case

Computing the Address Directly

```

switch(e) {
  case 0:  block0;
           break;
  case 1:  block1;
           break;
  case 2:  block2;
           break;
  ...
  case 9:  block9;
           break;
  default: blockd;
           break;
}

```

```

 $t_1 \leftarrow e$ 
if ( $t_1 < 0$  or  $t_1 > 9$ )
  then jump to  $LB_d$ 
  else
     $t_2 \leftarrow @Table + t_1 \times 4$ 
     $t_3 \leftarrow \text{memory}(t_2)$ 
    jump to  $t_3$ 

```

Label

LB_0
LB_1
LB_2
LB_3
LB_4
LB_5
LB_6
LB_7
LB_8
LB_9

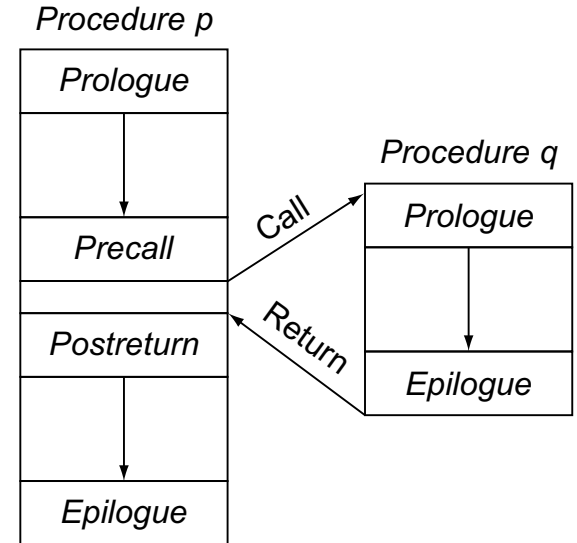
Procedure Calls

The implementation of procedure calls is, for the most part, straightforward.

Recall: a procedure call consists of

- a **precall** and a **postreturn** sequence in the caller
- a **prologue** and an **epilogue** sequence in the callee

In the following, we focus on issues that affect the compiler's ability to generate **efficient**, **compact**, and **consistent** code for procedure calls.



As a general rule, moving operations from the precall and postreturn sequences into the prologue and epilogue sequences should **reduce the overall size** of the final code.

Evaluating Actual Parameters

When it builds the precall sequence, the compiler must emit code to evaluate the actual parameters to the call. The compiler treats each actual parameter as an expression.

- **call-by-value parameters:** evaluate the expression and store its value in location designated for that parameter, *i.e.*, a register or callee's AR
- **call-by-reference parameter:** evaluate the parameter to an address and store it in the location designated for that parameter

Note If a call-by-reference parameter has no storage location, then the compiler may need to allocate space for that value so that it has an address to pass to the callee.

The compiler should use a consistent **evaluation order** for parameters, *i.e.*, either left to right or right to left, as evaluating parameters might have side effects.

Saving and Restoring Registers

As both the cost of memory operations and the number of registers have risen, the cost of **saving and restoring registers** has increased to the point that it needs careful attention.

1. **Using multi-register memory operations:** many ISAs support doubleword and quadword load and store operations for adjacent registers.
2. **Using a library routine:** replace the sequence of individual memory operations with a call to a compiler-supplied save or restore routine
3. **Combining responsibilities:** caller and callee pass a value back and forth that specifies which registers each must save