

Chapter 7:: Data Types

Programming Language Pragmatics

Michael L. Scott

Data Types

- We all have developed an intuitive notion of what types are; what's behind the intuition?
 - collection of values from a "domain" (the denotational approach)
 - internal structure of a bunch of data, described down to the level of a small set of fundamental types (the structural approach)
 - equivalence class of objects (the implementor's approach)
 - collection of well-defined operations that can be applied to objects of that type (the abstraction approach)

Data Types

- What are types good for?
 - implicit context
 - checking - make sure that certain meaningless operations do not occur
 - type checking cannot prevent all meaningless operations
 - It catches enough of them to be useful
- Polymorphism results when the compiler finds that it doesn't need to know certain things

Data Types

- STRONG TYPING has become a popular buzz-word
 - like *structured programming*
 - informally, it means that the language prevents you from applying an operation to data on which it is not appropriate
- STATIC TYPING means that the compiler can do all the checking at compile time

Type Systems

- Examples
 - Common Lisp is strongly typed, but not statically typed
 - Ada is statically typed
 - Pascal is almost statically typed
 - Java is strongly typed, with a non-trivial mix of things that can be checked statically and things that have to be checked dynamically

Type Systems

- Common terms:
 - discrete types – countable
 - integer
 - boolean
 - char
 - enumeration
 - subrange
 - Scalar types - one-dimensional
 - discrete
 - real

Type Systems

- Composite types:
 - records (unions)
 - arrays
 - strings
 - sets
 - pointers
 - lists
 - files

Type Systems

- ORTHOGONALITY is a useful goal in the design of a language, particularly its type system
 - A collection of features is orthogonal if there are no restrictions on the ways in which the features can be combined (analogy to vectors)

Type Systems

- For example
 - Pascal is more orthogonal than Fortran, (because it allows arrays of anything, for instance), but it does not permit variant records as arbitrary fields of other records (for instance)
- Orthogonality is nice primarily because it makes a language easy to understand, easy to use, and easy to reason about

Type Checking

- A TYPE SYSTEM has rules for
 - type equivalence (when are the types of two values the same?)
 - type compatibility (when can a value of type A be used in a context that expects type B?)
 - type inference (what is the type of an expression, given the types of the operands?)

Type Checking

- Type compatibility / type equivalence
 - Compatibility is the more useful concept, because it tells you what you can DO
 - The terms are often (incorrectly, but we do it too) used interchangeably.

Type Checking

- Certainly format does not matter:

```
struct { int a, b; }
```

is the same as

```
struct {  
    int a, b;  
}
```

We certainly want them to be the same as

```
struct {  
    int a;  
    int b;
```

```
}
```

Type Checking

- Two major approaches: structural equivalence and name equivalence
 - Name equivalence is based on declarations
 - Structural equivalence is based on some notion of meaning behind those declarations
 - Name equivalence is more fashionable these days

Type Checking

- There are at least two common variants on name equivalence
 - The differences between all these approaches boils down to where you draw the line between important and unimportant differences between type descriptions
 - In all three schemes described in the book, we begin by putting every type description in a standard form that takes care of "obviously unimportant" distinctions like those above

Type Checking

- Structural equivalence depends on simple comparison of type descriptions substitute out all names
 - expand all the way to built-in types
- Original types are equivalent if the expanded type descriptions are the same

Type Checking

- Coercion
 - When an expression of one type is used in a context where a different type is expected, one normally gets a type error

- But what about

```
var a : integer; b, c : real;
```

```
...
```

```
c := a + b;
```


Type Checking

- Coercion
 - Many languages allow things like this, and COERCE an expression to be of the proper type
 - Coercion can be based just on types of operands, or can take into account expected type from surrounding context as well
 - Fortran has lots of coercion, all based on operand type

Type Checking

- C has lots of coercion, too, but with simpler rules:
 - all `floats` in expressions become `doubles`
 - short `int` and `char` become `int` in expressions
 - if necessary, precision is removed when assigning into LHS

Type Checking

- In effect, coercion rules are a relaxation of type checking
 - Recent thought is that this is probably a bad idea
 - Languages such as Modula-2 and Ada do not permit coercions
 - C++, however, goes hog-wild with them
 - They're one of the hardest parts of the language to understand

Type Checking

- Make sure you understand the difference between
 - type conversions (explicit)
 - type coercions (implicit)
 - sometimes the word 'cast' is used for conversions (C is guilty here)

Records (Structures) and Variants (Unions)

- Records
 - usually laid out contiguously
 - possible holes for alignment reasons
 - smart compilers may re-arrange fields to minimize holes (C compilers promise not to)
 - implementation problems are caused by records containing dynamic arrays
 - we won't be going into that in any detail

Records (Structures) and Variants (Unions)

- Unions (variant records)
 - overlay space
 - cause problems for type checking
- Lack of tag means you don't know what is there
- Ability to change tag and then access fields hardly better
 - can make fields "uninitialized" when tag is changed (requires extensive run-time support)
 - can require assignment of entire variant, as in Ada

Records (Structures) and Variants (Unions)

- Memory layout and its impact (structures)

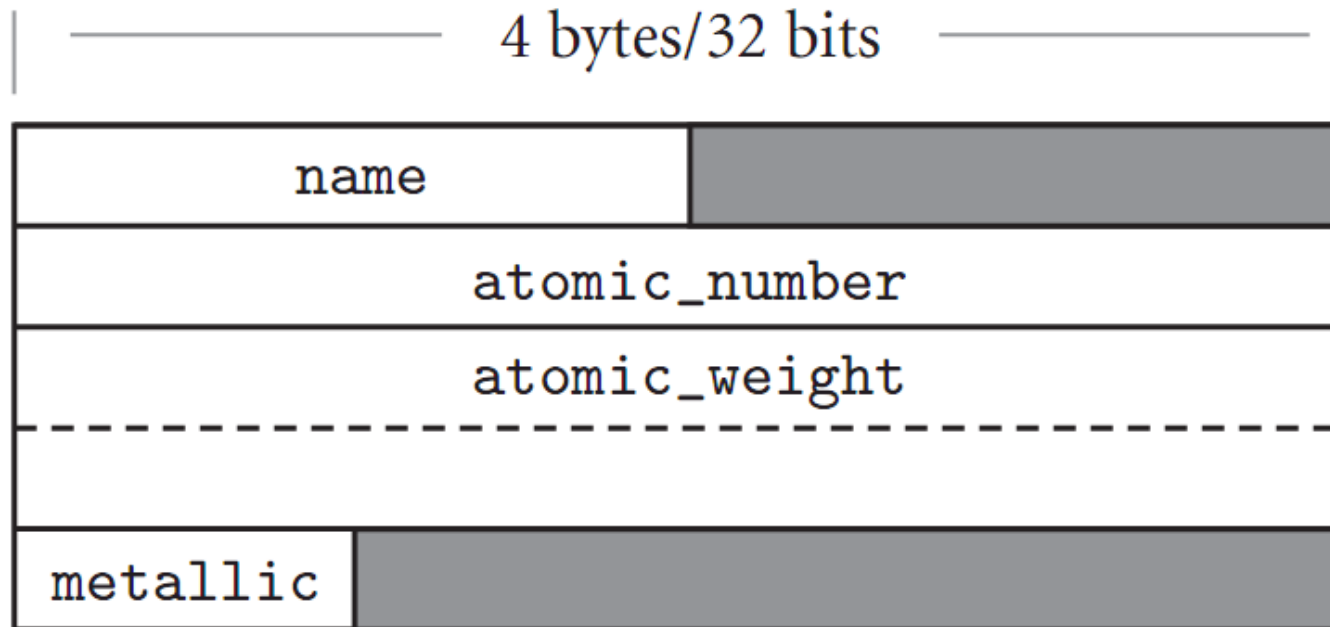


Figure 7.1 Likely layout in memory for objects of type element on a 32-bit machine. Alignment restrictions lead to the shaded “holes.”

Records (Structures) and Variants (Unions)

- Memory layout and its impact (structures)

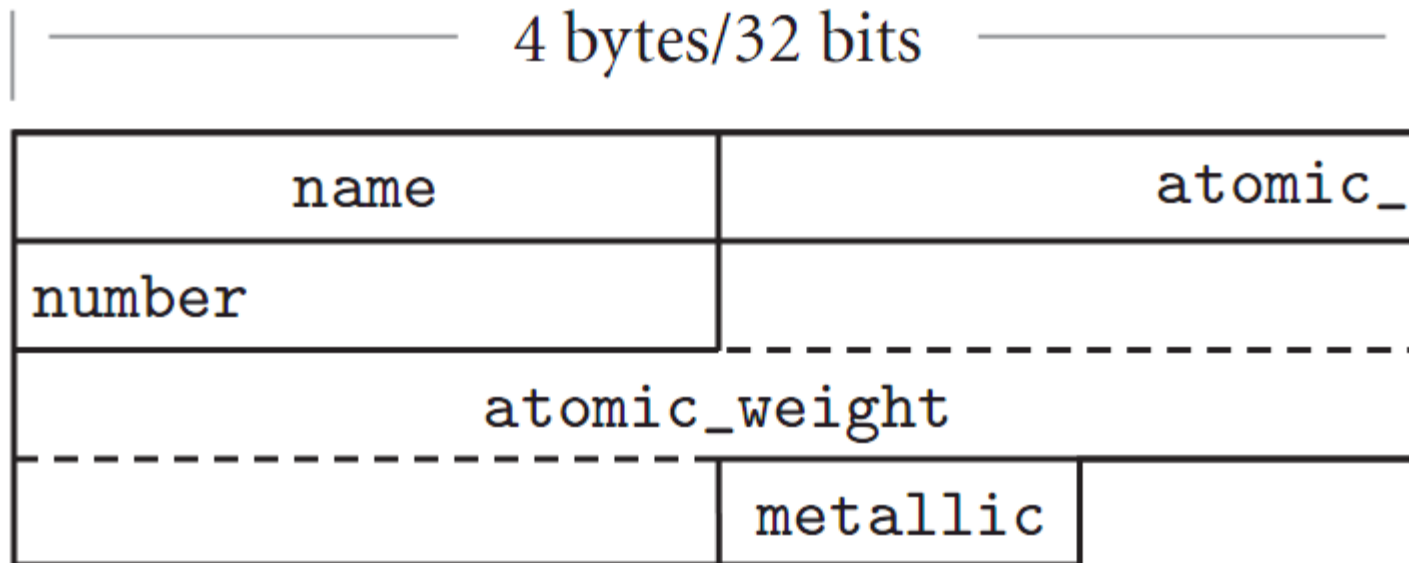


Figure 7.2 Likely memory layout for packed element records. The `atomic_number` and `atomic_weight` fields are nonaligned, and can only be read or written (on most machines) via multi-instruction sequences.

Records (Structures) and Variants (Unions)

- Memory layout and its impact (structures)

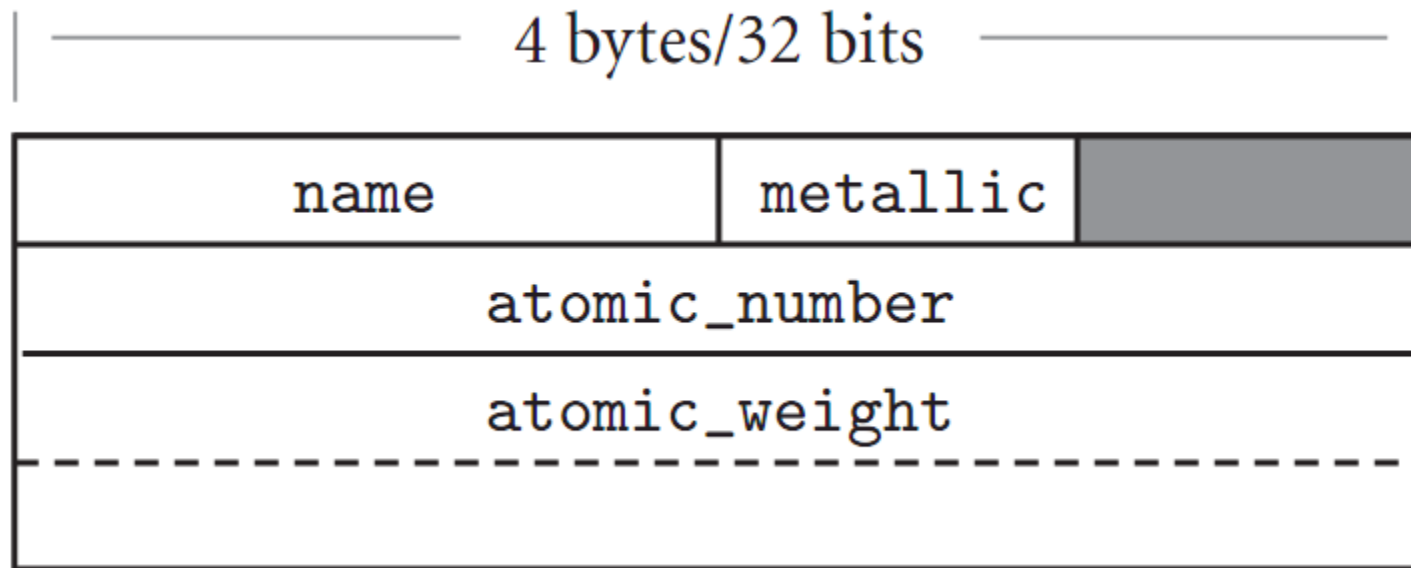


Figure 7.3 Rearranging record fields to minimize holes. By sorting fields according to the size of their alignment constraint, a compiler can minimize the space devoted to holes, while keeping the fields aligned.

Records (Structures) and Variants (Unions)

- Memory layout and its impact (unions)

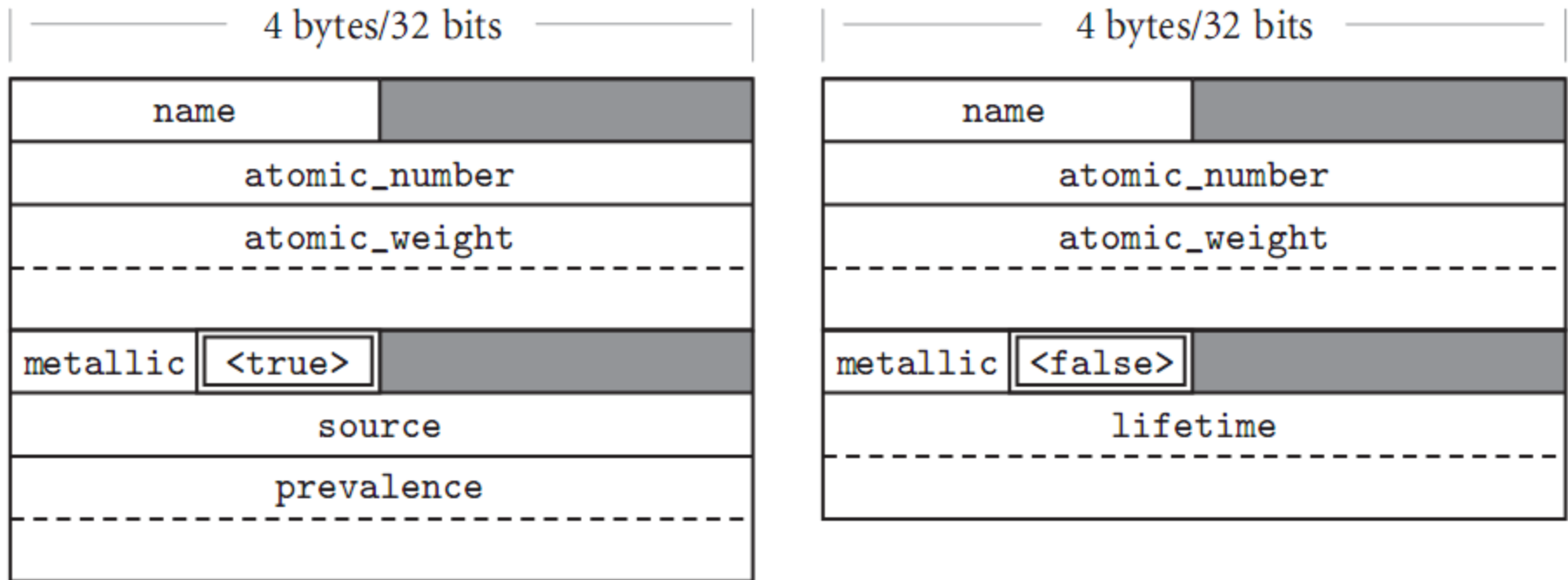
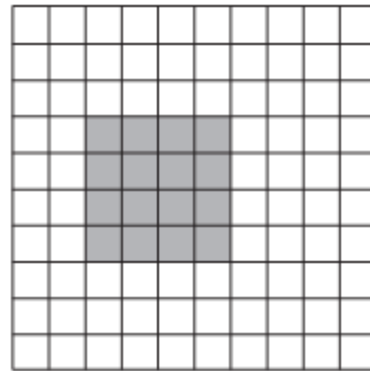


Figure 7.15 (CD) Likely memory layouts for element variants. The value of the naturally occurring field (shown here with a double border) determines which of the interpretations of the remaining space is valid. Type `string_ptr` is assumed to be represented by a (four-byte) pointer to dynamically allocated storage.

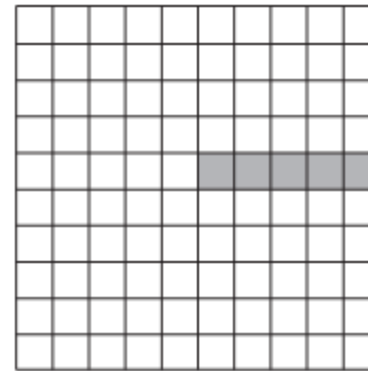
Arrays

- Arrays are the most common and important composite data types
- Unlike records, which group related fields of disparate types, arrays are usually homogeneous
- Semantically, they can be thought of as a mapping from an *index type* to a *component* or *element type*
- A *slice* or *section* is a rectangular portion of an array (See figure 7.4)

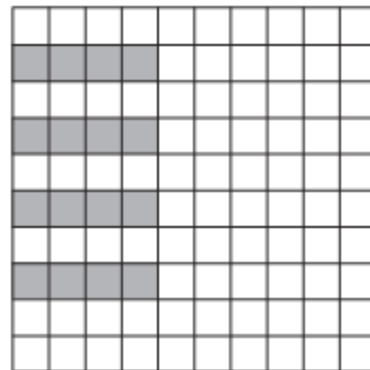
Arrays



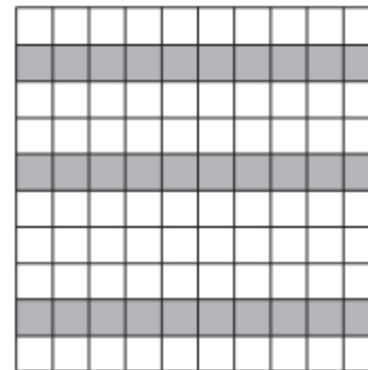
`matrix(3:6, 4:7)`



`matrix(6:, 5)`



`matrix(:, 2:8:2)`



`matrix(:, (/2, 5, 9/))`

Figure 7.4 Array slices(sections) in Fortran90. Much like the values in the header of an enumeration-controlled loop (Section 6.5.1), `a:b:c` in a subscript indicates positions `a`, `a+c`, `a+2c`, ...through `b`. If `a` or `b` is omitted, the corresponding bound of the array is assumed. If `c` is omitted, 1 is assumed. It is even possible to use negative values of `c` in order to select positions in reverse order. The slashes in the second subscript of the lower right example delimit an explicit list of positions.



Arrays

- Dimensions, Bounds, and Allocation
 - ***global lifetime, static shape*** — If the shape of an array is known at compile time, and if the array can exist throughout the execution of the program, then the compiler can allocate space for the array in static global memory
 - ***local lifetime, static shape*** — If the shape of the array is known at compile time, but the array should not exist throughout the execution of the program, then space can be allocated in the subroutine's stack frame at run time.
 - ***local lifetime, shape bound at elaboration time***

Arrays

```
-- Ada:  
procedure foo (size : integer) is  
  M : array (1..size, 1..size) of real;  
  ...  
begin  
  ...  
end foo;
```

```
// C99:  
void foo(int size) {  
    double M[size][size];  
    ...  
}
```

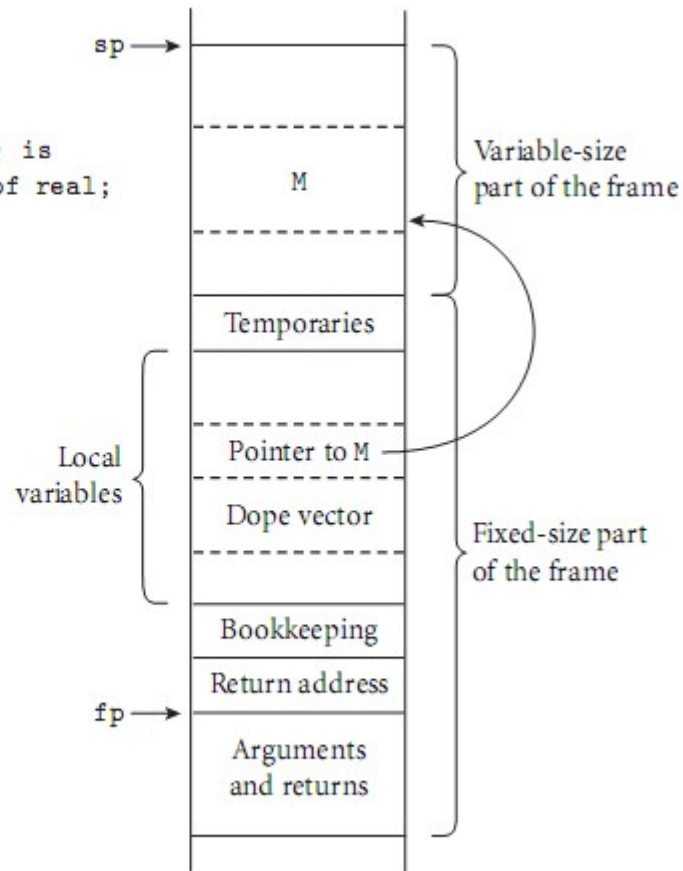


Figure 7.6 Elaboration-time allocation of arrays in Ada or C99.

Arrays

- Contiguous elements (see Figure 7.7)
 - column major - only in Fortran
 - row major
 - used by everybody else
 - makes array [a..b, c..d] the same as array [a..b] of array [c..d]

Arrays

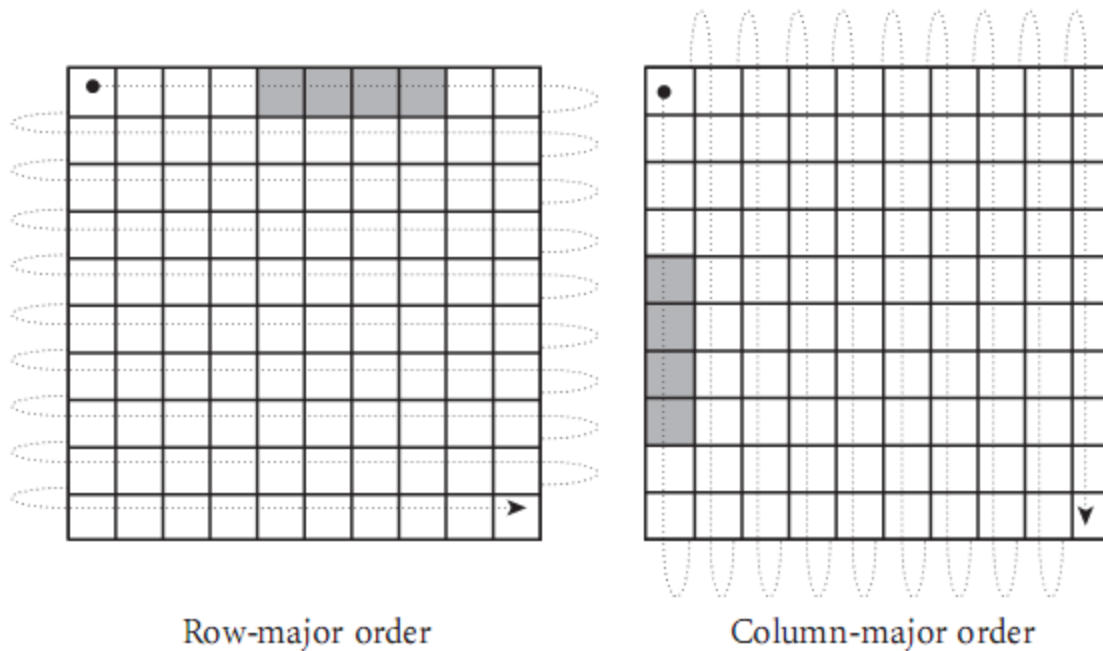


Figure 7.7 Row- and column-major memory layout for two-dimensional arrays. In row-major order, the elements of a row are contiguous in memory; in column-major order, the elements of a column are contiguous. The second cache line of each array is shaded, on the assumption that each element is an eight-byte floating-point number, that cache lines are 32 bytes long (a common size), and that the array begins at a cache line boundary. If the array is indexed from $A[0,0]$ to $A[9,9]$, then in the row-major case elements $A[0,4]$ through $A[0,7]$ share a cache line; in the column-major case elements $A[4,0]$ through $A[7,0]$ share a cache line.

Arrays

- Two layout strategies for arrays (Figure 7.8):
 - Contiguous elements
 - Row pointers
- Row pointers
 - an option in C
 - allows rows to be put anywhere - nice for big arrays on machines with segmentation problems
 - avoids multiplication
 - nice for matrices whose rows are of different lengths
 - e.g. an array of strings
 - requires extra space for the pointers

Arrays

```
char days[][10] = {
    "Sunday", "Monday", "Tuesday",
    "Wednesday", "Thursday",
    "Friday", "Saturday"
};
...
days[2][3] == 's'; /* in Tuesday */
```

S	u	n	d	a	y	/			
M	o	n	d	a	y	/			
T	u	e	s	d	a	y	/		
W	e	d	n	e	s	d	a	y	/
T	h	u	r	s	d	a	y	/	
F	r	i	d	a	y	/			
S	a	t	u	r	d	a	y	/	

```
char *days[] = {
    "Sunday", "Monday", "Tuesday",
    "Wednesday", "Thursday",
    "Friday", "Saturday"
};
...
days[2][3] == 's'; /* in Tuesday */
```

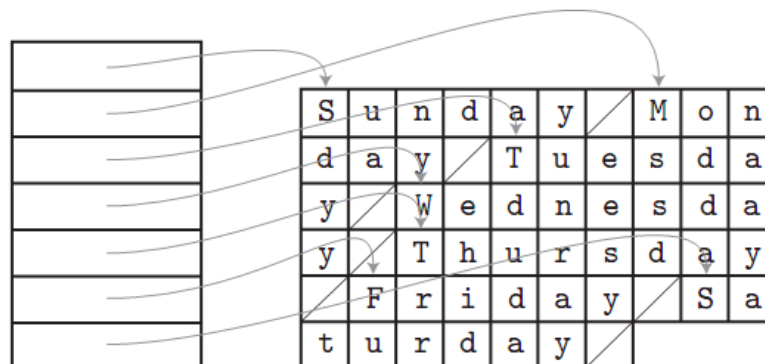


Figure 7.8 Contiguous array allocation v. row pointers in C. The declaration on the left is a true two-dimensional array. The slashed boxes are NUL bytes; the shaded areas are holes. The declaration on the right is a ragged array of pointers to arrays of characters. In both cases, we have omitted bounds in the declaration that can be deduced from the size of the initializer (aggregate). Both data structures permit individual characters to be accessed using double subscripts, but the memory layout (and corresponding address arithmetic) is quite different.

Arrays

- **Example:** Suppose

A : array [L1..U1] of array [L2..U2] of array [L3..U3] of elem;

$$D1 = U1 - L1 + 1$$

$$D2 = U2 - L2 + 1$$

$$D3 = U3 - L3 + 1$$

Let

$$S3 = \text{size of elem}$$

$$S2 = D3 * S3$$

$$S1 = D2 * S2$$

Arrays

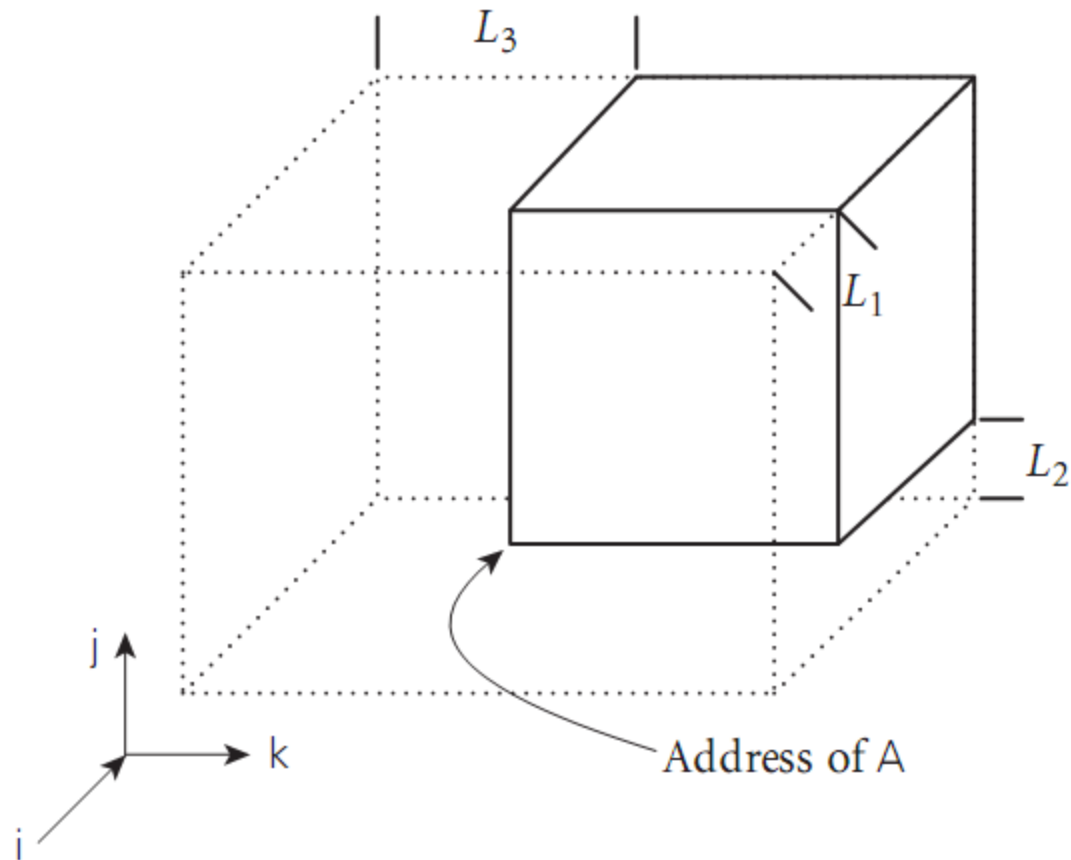


Figure 7.9 Virtual location of an array with nonzero lower bounds. By computing the constant portions of an array index at compile time, we effectively index into an array whose starting address is offset in memory, but whose lower bounds are all zero.

Arrays

- **Example** (continued)

We could compute all that at run time, but we can make do with fewer subtractions:

$$\begin{aligned} &== (i * S1) + (j * S2) + (k * S3) \\ &\quad + \text{address of } A \\ &\quad - [(L1 * S1) + (L2 * S2) + (L3 * S3)] \end{aligned}$$

The stuff in square brackets is compile-time constant that depends only on the type of A

Strings

- Strings are really just arrays of characters
- They are often special-cased, to give them flexibility (like polymorphism or dynamic sizing) that is not available for arrays in general
 - It's easier to provide these things for strings than for arrays in general because strings are one-dimensional and (more important) non-circular

Sets

- We learned about a lot of possible implementations
 - Bitsets are what usually get built into programming languages
 - Things like intersection, union, membership, etc. can be implemented efficiently with bitwise logical instructions
 - Some languages place limits on the sizes of sets to make it easier for the implementor
 - There is really no excuse for this

Pointers And Recursive Types

- Pointers serve two purposes:
 - efficient (and sometimes intuitive) access to elaborated objects (as in C)
 - dynamic creation of linked data structures, in conjunction with a heap storage manager
- Several languages (e.g. Pascal) restrict pointers to accessing things in the heap
- Pointers are used with a value model of variables
 - They aren't needed with a reference model

the 1990s, the number of people in the United States who are 65 years of age or older is projected to increase from 20 million to 35 million, and the number of people 75 years of age or older is projected to increase from 10 million to 15 million (U.S. Census Bureau, 1996).

serve



Pointers And Recursive Types

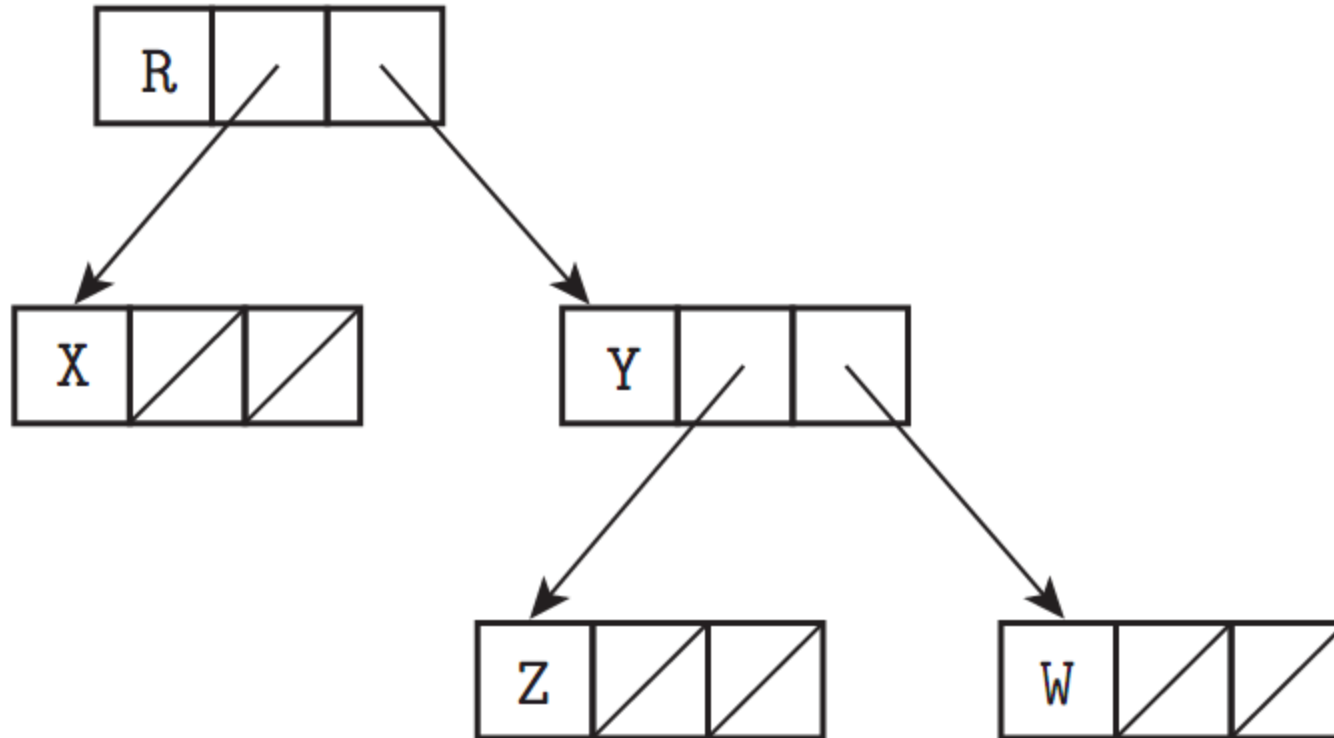


Figure 7.12 Typical implementation of a tree in a language with explicit pointers. As in Figure 7.11, a diagonal slash through a box indicates a null pointer.

Pointers And Recursive Types

- C pointers and arrays

```
int *a == int a[]
```

```
int **a == int *a[]
```

- BUT equivalences don't always hold

- Specifically, a declaration allocates an array if it specifies a size for the first dimension

- otherwise it allocates a pointer

```
int **a, int *a[] pointer to pointer to int
```

```
int *a[n], n-element array of row pointers
```

```
int a[n][m], 2-d array
```

Pointers And Recursive Types

- Compiler has to be able to tell the size of the things to which you point
 - So the following aren't valid:
`int a[][]` bad
`int (*a)[]` bad
 - C declaration rule: read right as far as you can (subject to parentheses), then left, then out a level and repeat
`int *a[n]`, n-element array of pointers to integer
`int (*a)[n]`, pointer to n-element array of integers

Pointers And Recursive Types

- Problems with dangling pointers are due to
 - explicit deallocation of heap objects
 - only in languages that *have* explicit deallocation
 - implicit deallocation of elaborated objects
- Two implementation mechanisms to catch dangling pointers
 - Tombstones
 - Locks and Keys

Pointers And Recursive Types

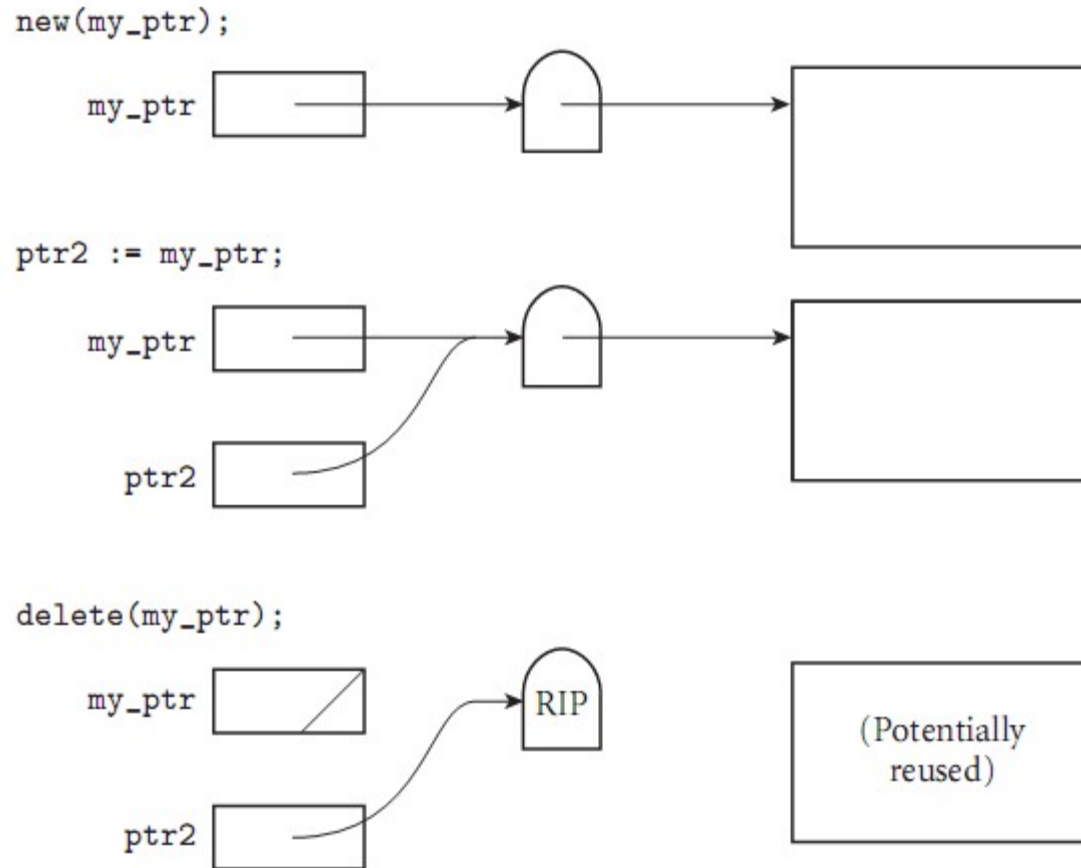


Figure 7.17 (CD) Tombstones. A valid pointer refers to a tombstone that in turn refers to an object. A dangling reference refers to an “expired” tombstone.

Pointers And Recursive Types

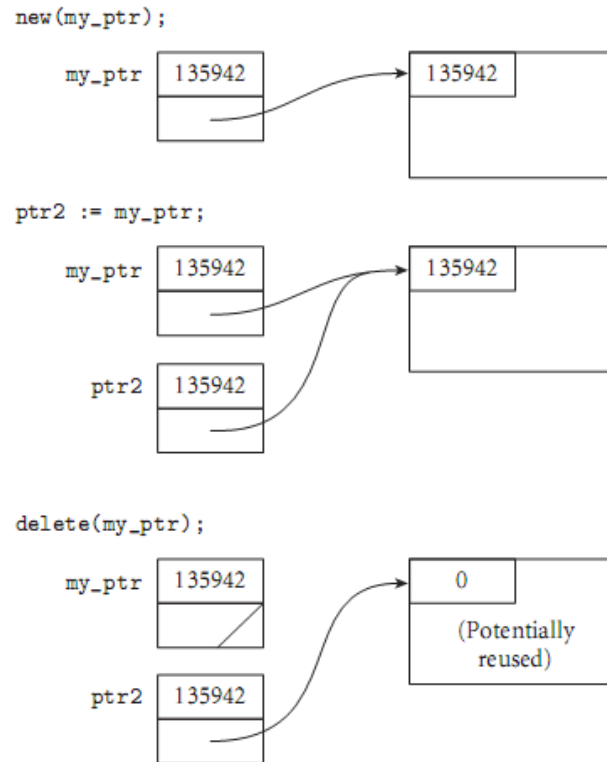


Figure 7.18 (CD) Locks and Keys. A valid pointer contains a key that matches the lock on an object in the heap. A dangling reference is unlikely to match.

Pointers And Recursive Types

- Problems with garbage collection
 - many languages leave it up to the programmer to design without garbage creation - this is VERY hard
 - others arrange for automatic garbage collection
 - reference counting
 - does not work for circular structures
 - works great for strings
 - should also work to collect unneeded tombstones

Pointers And Recursive Types

- Garbage collection with reference counts

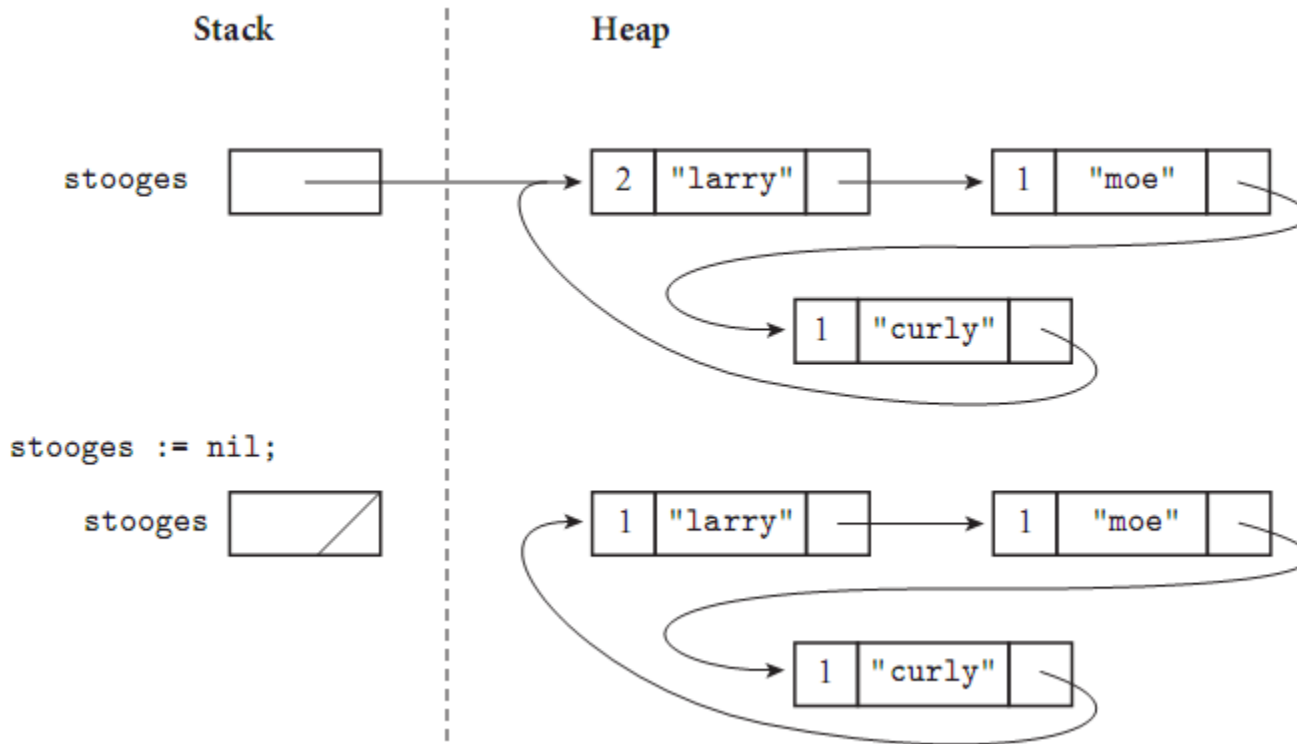


Figure 7.13 Reference counts and circular lists. The list shown here cannot be found via any program variable, but because it is circular, every cell contains a nonzero count.

Pointers And Recursive Types

- Mark-and-sweep
 - commonplace in Lisp dialects
 - complicated in languages with rich type structure, but possible if language is strongly typed
 - achieved successfully in Cedar, Ada, Java, Modula-3, ML
 - complete solution impossible in languages that are not strongly typed
 - conservative approximation possible in almost any language (Xerox Portable Common Runtime approach)

Pointers And Recursive Types

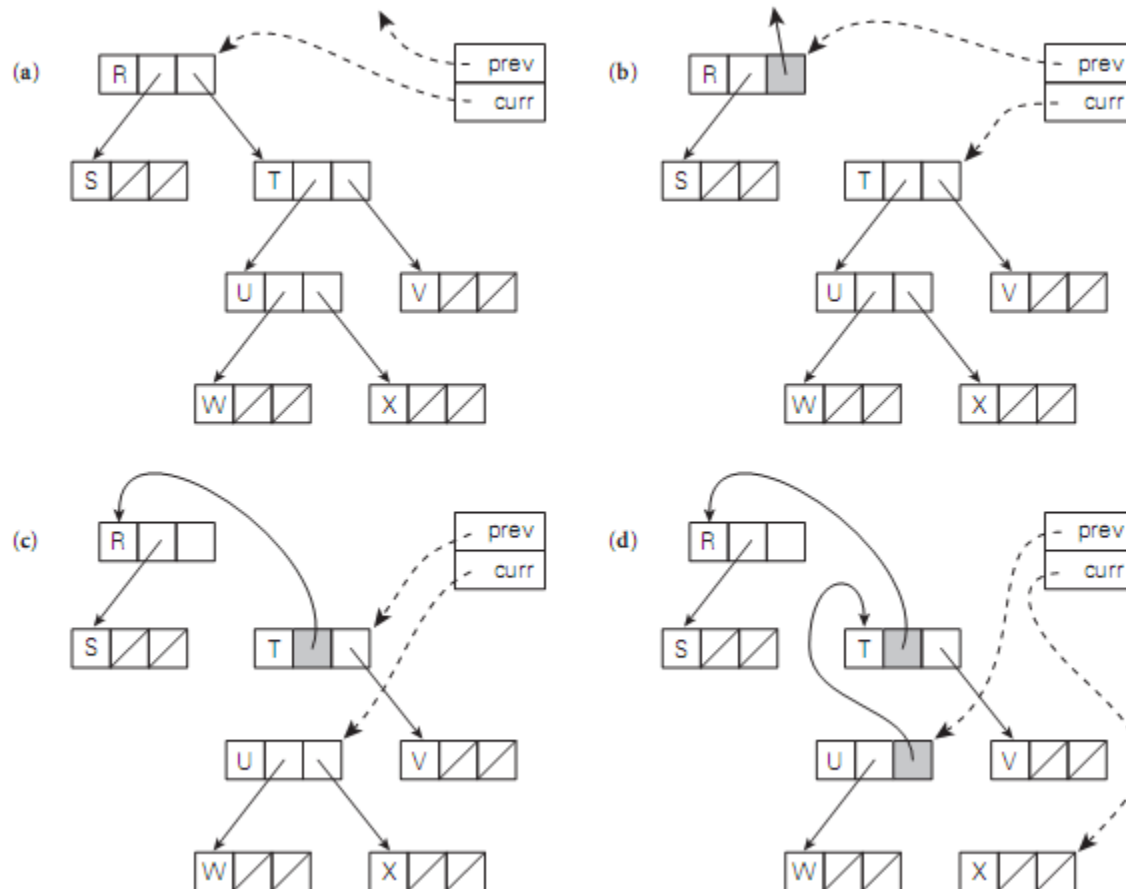


Figure 7.14 Heap exploration via pointer reversal.

Lists

- A list is defined recursively as either the empty list or a pair consisting of an object (which may be either a list or an atom) and another (shorter) list
 - Lists are ideally suited to programming in functional and logic languages
 - In Lisp, in fact, a program *is* a list, and can extend itself at run time by constructing a list and executing it
 - Lists can also be used in imperative programs

Files and Input/Output

- Input/output (I/O) facilities allow a program to communicate with the outside world
 - *interactive* I/O and I/O with files
- Interactive I/O generally implies communication with human users or physical devices
- Files generally refer to off-line storage implemented by the operating system
- Files may be further categorized into
 - *temporary*
 - *persistent*