Lectures 10-11

Memory management

Iztok Savnik, FAMNIT

May, 2024.

Literature

- John Mitchell, Concepts in Programming Languages, Cambridge Univ Press, 2003 (Chapter 7, Section 7.7.3)
- Michael L. Scott, Programming Language Pragmatics (3rd ed.), Elsevier, 2009 (Chapters 3 and 8)
- Brian W. Kernighan, Dennis M. Ritchie, The C Programming Language, Second Edition, Prentice Hall, Inc., 1988.

Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Introduction

- The topics of this lecture were discussed before, but not in an organized way and in detail.
- How and where are memory for variables allocated?
 - Static and dynamic variables
 - Stack allocation and heap allocation
- Implementation view of scope and lifetime of objects
 - Activation records, stack allocation,
 - Local and global variables
 - How to access global variables?

Introduction

- How the chains of function calls are implemented?
 - Stack activation records
 - Structures formed by the activation records
- Manual and automatic allocation of heap storage
 - Why using heap storage? How memory allocation and deallocation works?
 - Memory for objects, records, arrays, lists, ...
 - Memory leaks, dangling references, possible bugs
- Heap management strategies and algorithms
 - Price for automatic storage allocation
 - Garbage collection

Binding Time

- A <u>binding</u> is an association between two things, such as a name and the thing it names
- <u>Binding time</u> is time at which a binding is created
 - The time at which any implementation decision is made
 - Binding between question and answer
- Important binding times in SW systems
 - <u>Compile time</u>: Mapping of high-level constructs to machine code, including the layout of statically defined data in memory
 - <u>Link time</u>: Name in one module refers to an object in another module, the binding between the two is not finalized until link time
 - <u>Run time</u>: Bindings of values to variables and other run-time decisions

Binding time

- Early binding times are associated with greater <u>efficiency</u>, while later binding times are associated with greater <u>flexibility</u>
 - Compiler-based language implementations tend to be more efficient than interpreter-based implementations because they make earlier decisions
 - Interpreters must analyze the declarations every time the program begins execution
- The terms <u>static</u> and <u>dynamic</u> are generally used to refer to things bound before run time and at run time

Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Lifetime and Storage Management

- It is important to distinguish between <u>names</u> and the <u>objects</u> to which they refer
- Key events in object / binding lifetime:
 - Creation of objects
 - Creation of bindings
 - References to variables, subroutines, types, and so on, all of which use bindings
 - Deactivation and reactivation of bindings that may be temporarily unusable
 - Destruction of bindings
 - Destruction of objects

Lifetime

- The period of time between the creation and the destruction of a name-to-object binding is called the <u>binding's lifetime</u>.
- The time between the creation and destruction of an object is the <u>object's lifetime</u>.
 - Object may retain its value and the potential to be accessed even when a given name can no longer be used to access it
 - When a variable is passed to a subroutine by reference

- Fortran, var in Pascal, or '&' in C, C++

- Lifetime of name-to-object binding is longer than that of object
 - Generally, a sign of a program bug dangling references

Storage allocation mechanisms

- Object lifetimes generally correspond to one of three principal storage allocation mechanisms, used to manage the object's space:
 - 1. <u>Static objects</u> are given an absolute address that is retained throughout the program's execution.
 - 2. <u>Stack objects</u> are allocated and deallocated in last-in, first-out order, usually in conjunction with subroutine calls and returns.
 - 3. <u>Heap objects</u> may be allocated and deallocated at arbitrary times. They require a more general (and expensive) storage management algorithm.

Program memory



Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Static allocation

- <u>Global variables</u> are static objects, but they are not the only ones.
- Other <u>static objects</u>:
 - Instructions that constitute a program's machine language translation
 - Variables that are local to a single subroutine, but retain their values from one invocation to the next
 - Numeric and string-valued constant literals, such as A = B/14.7 or printf("hello, world\n")
 - Tables that are used by run-time support routines for debugging, dynamic-type checking, garbage collection, exception handling, and other purposes

Static allocation

- Statically allocated objects are often allocated in protected, read-only memory
- <u>Static activation records</u>
 - Storing variables in blocks, subrutines
 - One activation record for one subrutine
 - <u>Only one activation</u> of subrutine can live at a given time
 - Location of activation record is determined in compile time
 - Simple and very fast

Static allocation

- <u>Static activation records store</u>:
 - Local variables + Temporary values
 - Subrutine arguments, return address, return value
 - Reference to activation record of the caller
- Problems with static allocation records
 - Recursion can not be used
 - Fortran did not have recursion until very late version
 - Multi-threading can also not be implemented statically

Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Stack allocation

- <u>Activation record</u> or <u>stack frame</u> is allocated when block or subrutine are activated
 - Natural nesting of blocks and subroutine calls makes it easy to allocate space for locals <u>on a stack</u>
- Maintenance of the stack is the responsibility of the subroutine <u>calling-sequence</u>
 - The code executed by the caller immediately before and after the call (prologue/epilogue)
- Access to <u>higher name spaces</u>
 - Chain of activation records following structure of blocks and subrutines

Block-structured languages

• Storage management mechanisms associated with block structures

- outer block $\begin{cases} \{ \text{int } x = 2; \\ \{ \text{ int } y = 3; \\ x = y+2; \end{cases} \text{ block} \\ \end{cases}$
- A block can be: begin-end, procedure, function, let statemnt, ...
- A variable declared in block is said to be <u>local</u> to that block
- A variable declared in an enclosing block is said to be <u>global</u> to the block
- Properties of block-structured languages
 - May define new variables anywhere in block
 - Blocks may be <u>nested</u>, but cannot partially overlap
 - <u>When entering</u>, memory is allocated for variables declared in block
 - <u>When exiting</u>, some or all of the memory allocated to variables declared in that block will be deallocated



Stack-based allocation of space for subroutines



Memory for block variables

- Memory for block variables stores three classes of variables:
- Local variables
 - Stored on the stack in the activation record associated with the block
- Parameters
 - Parameters to subroutine stored in activation record
- Global variables
 - Accessed from an activation record that was placed on the run-time stack before activation of the current block

Example

- { int x=0; int y=x+1; { int z=(x+y)*(x-y); }; }.
- When the outer block is entered, }; an activation record containing space for x and y is pushed onto the stack
- On entry into the inner block, a separate activation record containing space for z will be added to the stack
- After the value of z is set, the activation record containing this value will be popped off the stack
- Finally, on exiting the outer block, the activation record containing space for x and y will be popped off the stack

intermedate results

```
{ int z = (x+y)*(x-y);
}
```

Space for z
Space for x+y
Space for x-y

| Space for global variables |
|----------------------------|----------------------------|----------------------------|----------------------------|
| | Space for x,y | Space for x,y | Space for x,y |
| | | Space for z | |

Control link

- Different activation records have different sizes
 - Operations that push and pop activation records from the run-time stack store a pointer to the top of the preceding activation record
- Control link (dynamic link)
 - The pointer to the top of the previous activation record
- When a new act. record is added
 - Control link of the new activation record is set to the previous value of the environment pointer
 - Environment pointer is updated



Example



Activation Records for Functions

- <u>Control link</u>, pointing to the previous activation record on the stack,
- <u>Access link</u> (static link), pointer to structurally subsuming block
- <u>Return address</u>, giving the address of the first instruction to execute when the function terminates,
- <u>Return-result address</u>, the location in which to store the function return value,
- Actual parameters of the function,
- Local variables declared within the function,
- <u>Temporary storage</u> for intermediate results computed with the function executes



Example fun fact(n) = if n <= 1 then 1 else n * fact(n-1);

 Activation records are added and removed from the run-time stack when tracing the execution of the familiar factorial function



Global Variables

- Identifier x appears in the body of a function, but x is not declared inside the function
- <u>Access to a global x involves finding an appropriate</u> <u>activation record elsewhere on the stack</u>
- There are two main rules for finding the declaration of a global identifier
 - <u>Static Scope</u>:
 - A global identifier refers to the identifier with the name that is declared in the closest enclosing scope of the program text
 - <u>Dynamic Scope</u>:
 - A global identifier refers to the identifier associated with the most recent activation record

Static and dynamic scope

- Difference between static and dynamic scope:
 - Finding declaration under static scope uses the <u>static</u> (<u>unchanging</u>) relationship between blocks in program text.
 - Actual sequence of calls that are executed in the <u>dynamic (changing) execution</u> of the program.
- Example:

int x=1;	outer block	Х	1
function $g(z) = x+z;$			
function $f(y) = \{$	f(3)	У	3
int $x = y+1;$		х	4
return g(y*x)			
};	g(12)	7	12
f(3);	5(12)	L	12

Access link

- How to implement Static scope?
- Access link (<u>static link</u>) of an activation record points to the activation record of the closest enclosing block in the program
- <u>In-line blocks</u> do not need an access link, as the closest enclosing block will be the most recently entered block
- Access link will generally point to a different activation record than the control link



En vironment Pointer

Example

```
int x=1;
function q(z) = x+z;
function f(y) = \{
     int x = y+1;
     return q(y^*x)
};
```

```
f(3);
```

- Declaration of g occurs inside the scope of declaration of x
 - Access link for declaration of g points to activation record for declaration of x
- Declaration of f is similarly inside the scope of the declaration of g
 - Access link for declaration of f points to activation record for the declaration of g
- Calls f(3) and g(12) cause activation record to be allocated for scope associated with body of f and body of g, respectively



Control and access links

• <u>To summarize:</u>

- Control link is a link to the activation record of the previous (calling) block
- Access link is a link to the activation record of the closest enclosing block in program text
- Control link depends on the <u>dynamic behavior</u> of program
- Access link depends on only the <u>static form</u> of the program text
- Access links are used to find the location of global variables in statically scoped languages with nested blocks at run time

Example: quicksort

```
# let rec quicksort = function
    [] -> []
   | pivot::rest ->
      let rec split = function
        [] -> ([],[])
       | x::tail ->
          let (below, above) = split tail
          in
           if x<pivot then (x::below, above)
           else (below, x::above)
      in let (below, above) = split rest
         in quicksort below @ [pivot] @ quicksort above;;
val quicksort : 'a list -> 'a list = <fun>
```

Example: quicksort



Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Heap allocation

- A heap is a region of storage in which subblocks can be <u>allocated</u> and <u>deallocated</u> at arbitrary times
- Heaps are required for the dynamically allocated pieces of linked data structures
 - Character strings, lists, and sets, whose size may change on update
 - Arrays, records, objects, recursive data structures, ...
- <u>Strategies</u> to manage space in a heap
 - Tradeoffs between speed and space
 - Internal and external fragmentation



Storage-management algorithms

- <u>Single linked list</u>—the free list
 - Heap blocks not currently in use
 - Initially list consists of a single block comprising the entire heap
 - Allocation request searches list for a block of appropriate size
 - <u>First fit</u> algorithm
 - <u>Best fit</u> algorithm
 - Unneeded portion below some min threshold is left in block as internal fragmentation, or, inserted back to list

Single linked list

- One would expect the best-fit algorithm to do a better job.
- Time complexity:
 - Best-fit has higher allocation cost than first-fit algorithm.
 - Always goes through all candidates
 - In recent applications, we may have a <u>huge</u> number of blocks!
 - The concept of <u>"current" block</u> in first-fit algorithm
 - Travels in a round-robin fashion
- Any algorithm is <u>linear</u> in the number of free blocks
 - In worst case the algorithm has to inspect all blocks

Single linked list

- Space complexity
 - First-fit tends to behave better then best-fit
 - Best-fit results in a larger number of very <u>small "left-</u> <u>over" blocks</u>
 - Internal as well as external fragmentation?
 - Score depends on the distribution of size requests
 - Distribution depends on the application type

Single linked list

- How to reduce the cost of the algorithm?
 - Maintain separate free lists for <u>blocks of different sizes</u>
 - If block is not found in the appropriate list then the list with <u>larger blocks is searched</u>
 - The leftover is stored in a list with smaller blocks
 - Cost can be reduced to a constant
- We first consider a heap in the C prog. language
 - Then we will go through some solutions that use more lists

- C originally implemented heap based on <u>linked</u> <u>list of free blocks.</u>
 - Calls to malloc() and free() may occur in any order
 - malloc() calls upon the operating system to obtain more memory when needed
 - Space that malloc() manages <u>may not be contiguous</u>
 - Free storage is kept as a list of free blocks
 - Each block contains <u>a size</u>, <u>a pointer to the next block</u>, and <u>the space itself</u>.
 - Blocks are kept in order of <u>increasing storage address</u>
 - Last block (highest address) points to the first

- When a request is made, the free list is scanned until a <u>big-enough block</u> is found, i.e. "first fit".
- If the block is too big, it is split, and the proper amount is returned to the user while the residue remains on the free list.



- <u>Proper alignment</u> for the objects stored
 - The most restrictive type can be stored at a particular address, then all other types may be also
 - On some machines, the most restrictive type is a double; on others, int or long suffices

```
typedef long Align; /* for alignment to long boundary */
union header { /* block header */
struct {
    union header *ptr; /* next block if on free list */
    unsigned size; /* size of this block */
    } s;
    Align x; /* force alignment of blocks */
};
typedef union header Header;
```



A block returned by **malloc**

- Requested size in characters is rounded up to the proper number of <u>header-sized units</u>
 - Block that will be allocated contains <u>one more unit</u>, for the header itself
- Search for a free block of adequate size begins where the <u>last block was found</u> (at freep)
 - This strategy helps keep the list homogeneous
 - If a too-big block is found, the tail end is returned to user
 - Pointer returned to the user points to free space within the block, which begins one unit beyond the header

malloc()

```
static Header base; /* empty list to get started */
static Header *freep = NULL; /* start of free list */
/* malloc: general-purpose storage allocator */
void *malloc(unsigned nbytes)
```

```
{
```

}

```
Header *p, *prevp;
Header *morecore(unsigned);
unsigned nunits;
```

```
nunits = (nbytes+sizeof(Header)-1)/sizeof(header) + 1;
if ((prevp = freep) == NULL) {
    /* no free list yet */
    base.s.ptr = freeptr = prevptr = &base;
    base.s.size = 0;
```

```
for (p = prevp > s.ptr; ; prevp = p, p = p > s.ptr) 
  if (p->s.size >= nunits) { /* big enough */
    if (p->s.size == nunits) /* exactly */
       prevp->s.ptr = p->s.ptr;
                            /* allocate tail end */
     else {
      p->s.size -= nunits;
      p += p -> s.size;
      p->s.size = nunits;
     }
    freep = prevp;
     return (void *)(p+1);
  if (p == freep) /* wrapped around free list */
     if ((p = morecore(nunits)) = = NULL)
       return NULL; /* none left */
  }
```

Î

Û

morecore() free()

```
#define NALLOC 1024 /* minimum #units to request */
/* morecore: ask system for more memory */
static Header *morecore(unsigned nu)
{
  char *cp, *sbrk(int);
  Header *up;
  if (nu < NALLOC)
    nu = NALLOC:
  cp = sbrk(nu * sizeof(Header));
  if (cp == (char *) -1) /* no space at all */
    return NULL:
  up = (Header *) cp;
  up->s.size = nu;
  free((void *)(up+1));
  return freep;
}
```

```
/* free: put block ap in free list */
void free(void *ap)
{
  Header *bp, *p;
  bp = (Header *)ap - 1; /* point to block header */
  for (p = freep; !(bp > p && bp < p->s.ptr); p = p->s.ptr)
       if (p \ge p \ge s.ptr \&\& (bp \ge p || bp 
          break; /* freed block at start or end of arena */
  if (bp + bp -> size == p -> s.ptr) \{ /* join to upper nbr */
     bp->s.size += p->s.ptr->s.size;
     bp > s.ptr = p > s.ptr > s.ptr;
  } else
     bp->s.ptr = p->s.ptr;
  if (p + p \rightarrow size == bp) \{ /* join to lower nbr */
     p->s.size += bp->s.size;
     p \rightarrow s.ptr = bp \rightarrow s.ptr;
  } else
     p \rightarrow s.ptr = bp;
  freep = p;
}
```

Dynamic pools

- Heap is divided into "pools," one for each standard size
 - Request is rounded up to the next standard size
 - Division into ranges may be static or dynamic
- Two common mechanisms for dynamic pool:
 - Buddy system
 - Fibonacci heap
- Buddy system
 - Standard block sizes are powers of two
 - Block of size 2^k is needed
 - none ⇒ block of size 2^{k+1} is split in two (one is put into 2^k free list)

Dynamic pools

- When a block is deallocated, it is coalesced with its "buddy"—if that buddy is free
- Fibonacci heap
 - Fibonacci numbers for the standard sizes
 - Slightly more complex
 - Leads to slightly lower internal fragmentation
- Problems with external fragmentation
 - The ability of the heap to satisfy requests may <u>degrade</u> over time
 - It is always possible to devise a sequence of requests that cannot be satisfied (while enough space 3)
 - <u>Compact the heap</u>, by moving already-allocated blocks
 - update all outstanding references

Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Manual/automatic mem. management

- Explicit (manual) memory management
 - Explicit allocation and deallocation of objects
 - Programer has total control over memory management
 - C, C++, Go (unsafe), Rust (new model), …
- <u>Automatic</u> memory management
 - Compiler and run-time system manage memory
 - Garbage collection
 - Java, Scala, Go, Ocaml, Haskell, Erlang, Python, Perl,

Explicit memory management

- Program allocates memory blocks and has full control over them
 - After block is not needed it is reclaimed
- Usual malloc-free cycle known from C

 Functions malloc() and free() implement heap in C
- Problems:
 - If pointer is "lost" then we have memory leak performance decreases and memory is filled-up eventually
 - If object is reclaimed by mistake we have dangling pointer
- Programmer has to be very careful and design all allocations in pairs.

Example

void function_which_allocates() {
 // allocate an array of 50 floats
 float* a = new float[50];
 // additional code making use of 'a'
 // return to caller, having forgotten
 // to delete the memory we allocated
}

int main() {

```
function_which_allocates();
```

```
// the pointer 'a' no longer exists, and therefore cannot be freed,
```

```
// but the memory is still allocated. a <u>leak</u> has occurred.
```

```
int* p = new(1024);
```

int* q = p;

delete p; // q is dangling pointer by now

// main continues

```
*q = 2048; // memory corruption: write to garbage memory
```

delete q; // memory corruption: double free of memory

Explicit memory management

- Many reclaims are automatic
 - On function return, space for local variables and parameters is reclaimed
- <u>Disciplined allocation/deallocation</u> of memory can lead to efficient programs
 - In reality, all fast programs are implemented in languages that allow explicit memory management
 - Performance of PL with GC is comparable to explicit memory management if there is enough memory :-|
 - Language without GC can perform orders of magnitude better than languages with GC
 - If memory is a problem languages with explicit control are always better

Rust memory model

Ownership rules: Each value in Rust has an owner. There can only be one owner at a time. When the owner goes out of scope, the value will be dropped.

- Rust has a new memory model (Mozilla, 2010)
 - Rules with a "story" are used for managing memory
 - System of ownership with a set of rules that compiler checks
 - Ownership works for heap objects! (for references to heap!)
 - Exactly one variable is an <u>owner</u> of a memory object
 - Owner goes out of scope, the object will be dropped
 - Ownership is <u>moved</u> by assignment
 - Object has one owner and many <u>borrowers</u>
 - Borrowers get a reference to the object
 - Object does not disappear while there are references to it
 - Rust uses a form of <u>reference count</u> introduced for automatic mem. management (presented shortly)
 - No garbage collection!

Outline

- 1. Binding time
- 2. Lifetime and storage management
- 3. Static allocation
- 4. Stack allocation
- 5. Heap allocation
- 6. Explicit memory management
- 7. Garbage collection

Garbage collecton

- <u>Allocation</u> of heap-based objects is always triggered by some specific operation in program:
 - Instantiating an object, appending to the end of a list, assigning a long value into previously short string, ...
- <u>Deallocation</u> can be done in two ways:
 - It is explicit in some languages
 - Examples: C, C++, Go and Pascal
 - In many languages objects are deallocated implicitely
 - After they can not be reached from any program variable
 - Such language must then provide a garbage collection mechanism to identify and reclaim unreachable objects

Garbage collection

- Languages with garbage collection
 - Most functional and scripting languages <u>require</u> garbage collection
 - Many more recent imperative languages (including Modula-3, Java, and C) use garbage collectors
- Arguments in favor of <u>explicit deallocation</u>!
 - Implementation simplicity of prog. languages
 - Even naive implementations of automatic garbage collection add significant complexity to the implementation
 - Execution speed
 - Even the most sophisticated garbage collector can consume nontrivial amounts of time in certain programs

Garbage collection

- Argument in favor of <u>automatic garbage collection</u>
 - Manual deallocation errors are among the most common and costly bugs in real-world programs
 - Object is deallocated too soon
 - Program may follow a <u>dangling reference</u>
 - Accessing memory now used by another object
 - Object is not deallocated at the end of its lifetime
 - Program may <u>leak memory</u>, eventually running out of heap space
 - Deallocation errors are notoriously difficult to identify and fix

Garbage collection

- <u>Automatic garbage collection is an essential</u> <u>language feature</u> (Invariant)
 - Conclusion of both language designers and programmers
 - Many times we do not want to spend many days debugging but want solution »at once«
 - The cost of garbage collection is compensated by faster hardware?
- In many cases it is <u>not possible</u> to implement system efficiently without explicit control of memory allocation
 - Most compilers, browsers, DBMSs, OS routines, ... are written in C or C++

Reference counts

- When is an <u>object no longer useful</u>?
 - There are no pointers to object
- Simple solution:
 - Place the counter of pointers referencing the object in object itself
 - Initially this counter is 1
 - When pointer a is assigned to pointer b:
 - 1.dec_rc(object(b))
 - 2.Make assignment b := a
 - 3.inc_rc(object(a))
 - On subrutine return
 - calling-sequence epilogue has to decrement reference counts of all objects referred to by parameters and local variables

Reference counts

- When reference count is 0 object can be reclaimed
 - Recursively, run-time system must decrement counts for any objects referred to by pointers within the reclaimed object
- In order for reference counts to work
 - Language implementation must be able to identify the location of every pointer
 - Which words in object or stack frame represent pointers?
 - Type descriptors (offsets of components) generated by compiler are used
 - In general languages with strong typing can use such algorithms
 - Solutions for languages not strongly typed also exist

Reference counts

- The most important problem is definition of "<u>useful</u> <u>object</u>".
 - Object may be useless despite there are references to it
 - RCs fail to collect circular structures
- Many languages use RC for var-length strings
 - They do not contain refs
- Perl uses RCs for all dynamically allocated data
 - Programmer is warned to break cycles



Mark-and-sweep

- Better definition of a "useful" object
 - Can be reached by following a chain of valid pointers starting from something that has a name
 - Circularly referenced objects do not stay in heap
- Recursively exploring the heap to determine what is useful
 - Starting from external pointers (very expensive...)
- <u>Mark-and-sweep</u> (collector)
 - Classic mechanism to identify useless blocks, under this more accurate definition
 - When amount of free space remaining in heap falls below some minimum threshold <u>collector</u> proceeds in three main steps

Mark-and-sweep

- 1. Collector walks through the heap, tentatively marking every block as "<u>useless</u>."
- 2. Beginning with all pointers outside the heap, collector recursively explores all linked data structures in the program, marking each newly discovered block as "<u>useful</u>."
- 3. The collector again walks through the heap, moving every block that is still marked "<u>useless</u>" to the free list.

Mark-and-sweep

- Problems with algorithm:
 - We must know where every block in-use begins and ends
 - Every block must begin with an indication of its size, and of whether it is currently free
 - Collector must be able in Step 2 to find the pointers contained within each block
 - *Solution*: put pointer to (block) type descriptor near the beginning of block

Improvements of Mark-and-sweep

Pointer reversal

- Recursive exploration of heap requires storage
 - Heap could be used to track the path
- As collector explores the path to a given block, it reverses the pointer to the block
 - Collector keeps track of current block and the block from whence it came
- Search can be implemented without stack



Improvements of Mark-and-sweep

- Stop-and-copy
 - Reduce external fragmentation by performing storage compaction
 - Eliminating Steps 1 and 3
 - Divide the heap into two regions of equal size
 - All allocations happen in first part
 - Each reachable block is copied into second half of the heap, with no external fragmentation
 - Old copy is marked "useful", in this way
 - Pointers to old block are corrected to point to new
 - When collector finishes, all useful blocks are stored in the second part of heap
 - First part of heap is empty!
 - Collector swaps its notion of first and second halves

Generational Collection

- Observation: most dynamically allocated objects are short-lived
- Heap is divided into multiple regions (often two)
 - When space runs low the collector first examines the youngest region (the "nursery")
 - It is likely to have the highest proportion of garbage
 - If it is unable to reclaim sufficient space in this region the collector examines the <u>next-older region</u>
 - In worst case collector has to examine complete heap
- In most cases, the overhead of collection will be proportional to the size of youngest region only

Generational Collection

- Object that survives few collections (often one) is promoted (<u>moved</u>) to the next older region
 - Reminiscent of stop-and-copy
 - You do not need to move around long living objects while dealing with younger regions
 - Gains the speed
 - Promotion requires that pointers reflect new locations
 - At each pointer assignment, the compiler generates code to check whether new value is <u>old-to-new pointer</u>
 - This instrumentation on assignments is known as a write barrier