

Garbage Collection

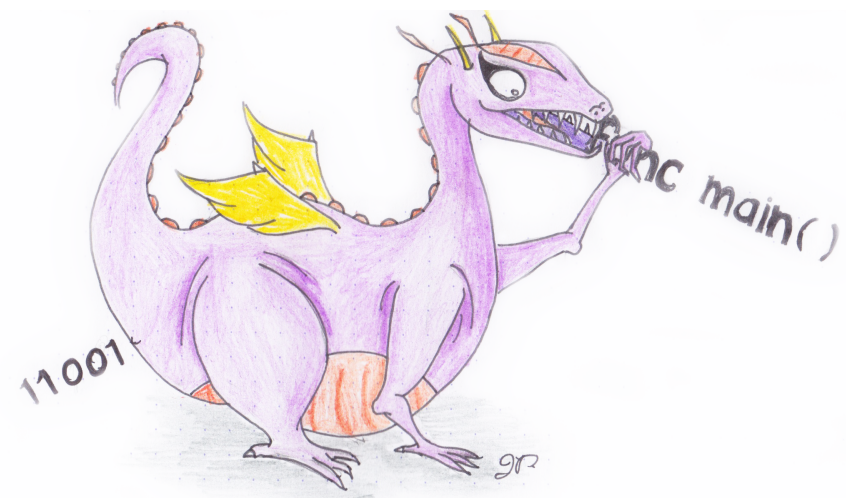
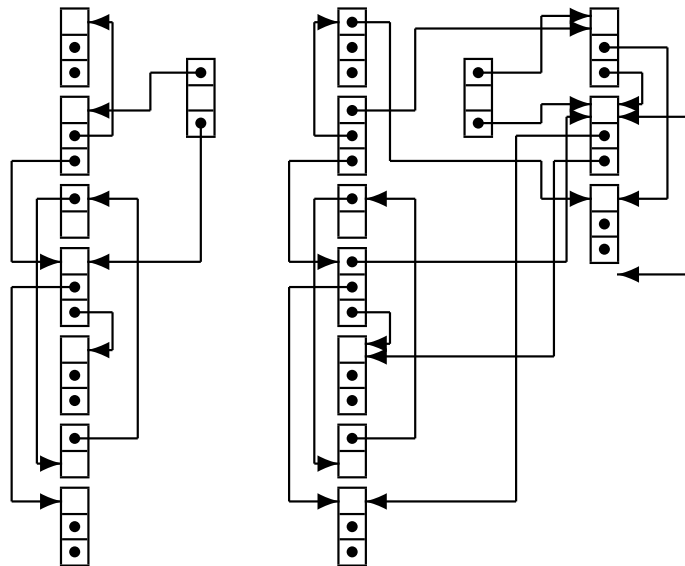
COMP 520: Compiler Design (4 credits)

Alexander Krolik

alexander.krolik@mail.mcgill.ca

MWF 10:30-11:30, TR 1100

<http://www.cs.mcgill.ca/~cs520/2020/>



Dot Gitignore

Announcements (Monday/Wednesday, March 9th/11th)

Milestones

- Milestone 1 graded (programs on myCourses)
- Milestone 2 due: Sunday, March 15th 11:59 PM
- Peephole due: Friday, April 10th 11:59 PM

To come

- Milestone 3 out: Monday, March 16th. Due: Friday, March 27th 11:59 PM
- Milestone 4 out: Monday, March 16th. Due: Friday, April 10th 11:59 PM
- Final report out: Monday, March 16th. Due: Friday, April 10th 11:59 PM
- Group meeting: Week of April 13th
- Final exam: Tuesday, April 21st 2:00 PM

Garbage Collection

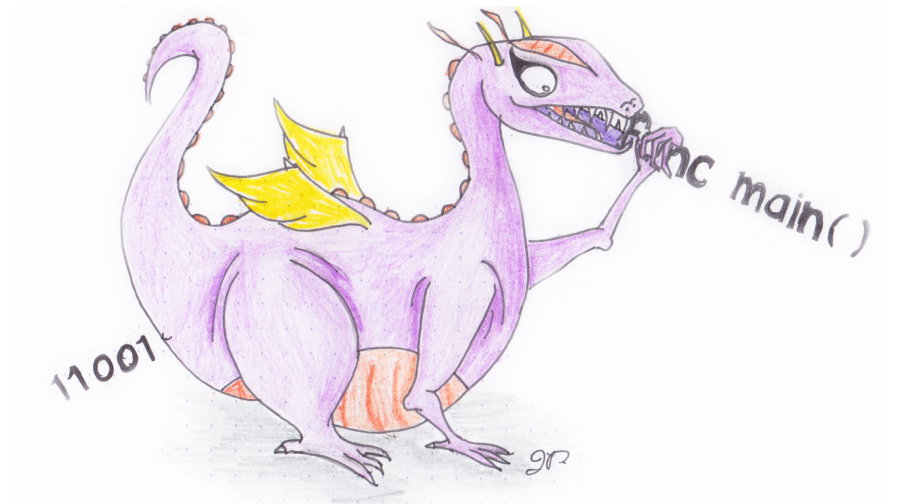
Memory management

Reference counting

Mark-and-sweep

Stop-and-copy

Practical



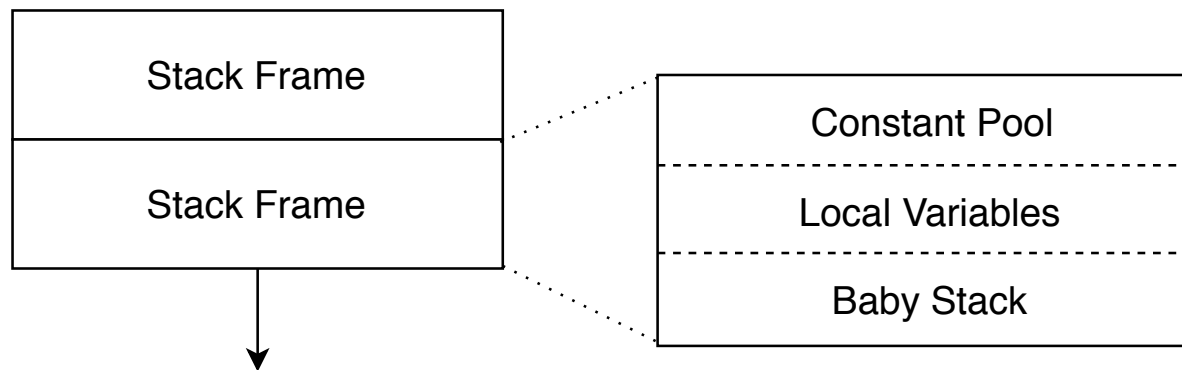
Memory Allocation

Stack Memory Allocation

- Space allocated in the function call stack;
- Is used for function call information, local variables, and return values;
- Typically contains fixed size data; and
- Is allocated and deallocated at the beginning and end of a function.

Information stored in the stack is therefore specific to a particular function invocation (i.e. call).

Example Stack



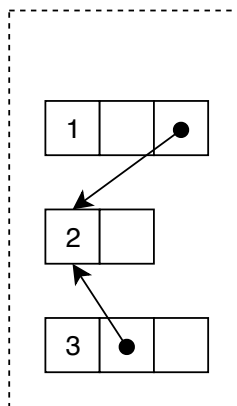
Memory Allocation

Heap Memory Allocation

- Space allocated in the program heap;
- Is very dynamic in nature:
 - Unknown size; and
 - Unknown time;
- Requires additional runtime support for managing the heap space.

Information stored in the heap is therefore not necessarily tied to any particular function invocation.

Example Heap



Heap variables may be referred to by other objects in the heap, or from the stack.

Heap Memory Allocation

Data stored in the heap is controlled by a heap allocator (i.e. `malloc`).

- Manages the memory in the heap space;
- Takes as input an integer representing the size needed for the allocation;
- Finds unallocated space in the heap large enough to accommodate the request; and
- Returns a pointer to the newly allocated space.

You will find more details in an operating systems course.

Heap Memory Deallocation

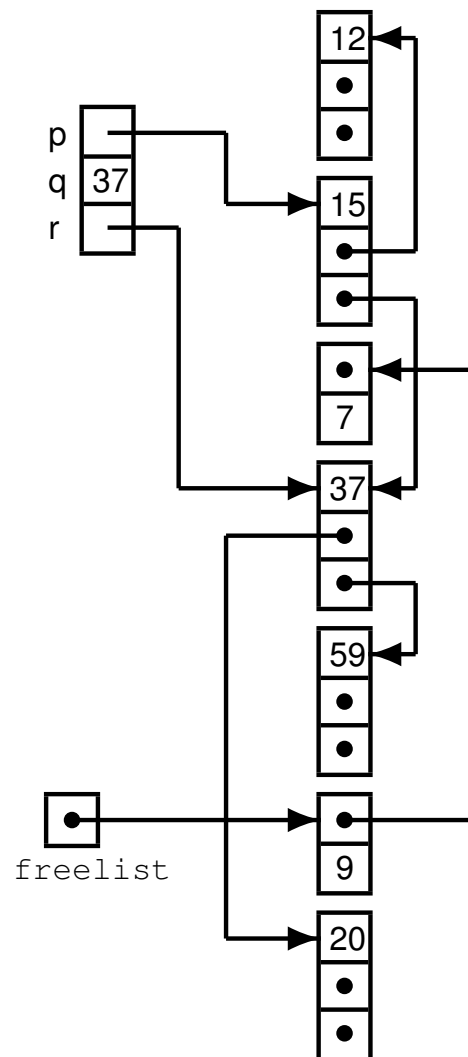
Memory allocated on the heap are freed when they are no longer live (i.e. `free`). This can be:

- Manual: User code making the necessary decisions on what is live;
- Continuous: Runtime code determining on the spot which objects are live; or
- Periodic: Runtime code determining at specific times which objects are live.

Without runtime support it is up to the *program* to return the memory when it is no longer needed.

Heap Memory Deallocation

For this class, we will assume that the freed heap blocks are stored on a `freelist` (a linked list of heap blocks). Freeing an object *prepends* the heap block onto the list.



Manual Deallocation Mechanisms

Heap memory can be freed manually at any point in the program.

- Leaves programmers to determine when an object is no longer live; and
- Requires calls to a deallocator (i.e. `free`).

Consider the following code

```
int *a = malloc(sizeof(int));
```

```
free(a);
```

```
[...]
```

```
*a = 5; // what happens?
```

Manual Deallocation Mechanisms

Advantages

- Reduces runtime complexity;
- Gives the programmer full control on what is live; and
- Can be more efficient in some circumstances.

Disadvantages

- Requires extensive effort from the programmer;
- Gives the programmer full control on what is live;
- Error-prone; and
- Can be less efficient in some circumstances.

Manual Deallocation Mechanisms

Sometimes manual deallocation is slower than automatic methods. Consider the following example code, which allocates 100 integers and then deallocates them one-by-one.

```
for (int i = 0; i < 100; ++i)
{
    a[i] = malloc(sizeof(int));
}
```

[...]

```
for (int i = 0; i < 100; ++i)
{
    free(a[i]);
}
```

This is potentially inefficient. Why?

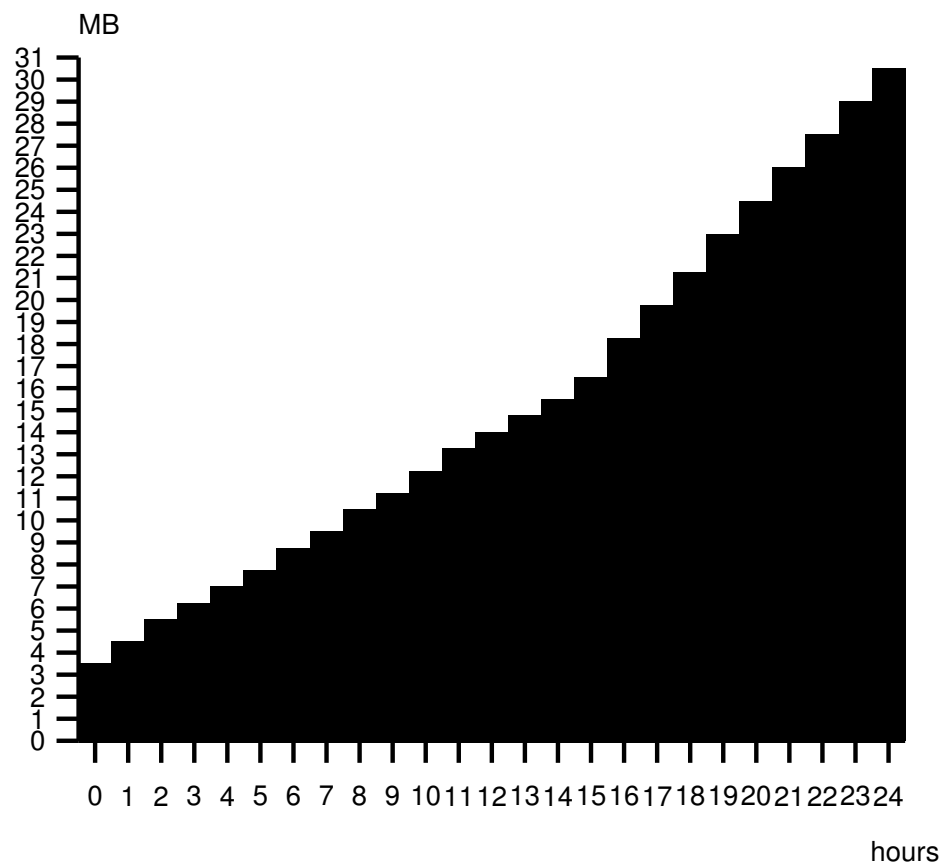
The allocations are potentially contiguous, and could therefore be reclaimed as a *block* instead of one-by-one.

Manual Deallocation Mechanisms

Life Without Garbage Collection

- Dead records must be explicitly deallocated;
- “Superior” if done correctly; but
- It is easy to miss some records;
- It is “dangerous” to handle pointers; and
- May be less efficient in some cases.

Memory leaks in real life (`ical v.2.1`)



Runtime Deallocation Mechanisms

A correct runtime deallocation mechanism must answer the fundamental question:

Which records are *dead*, i.e. no longer in use?

The more precise the answer, the better the deallocation mechanism.

Ideally

- Records that will never be accessed in the future execution of the program; but
- This is undecidable.

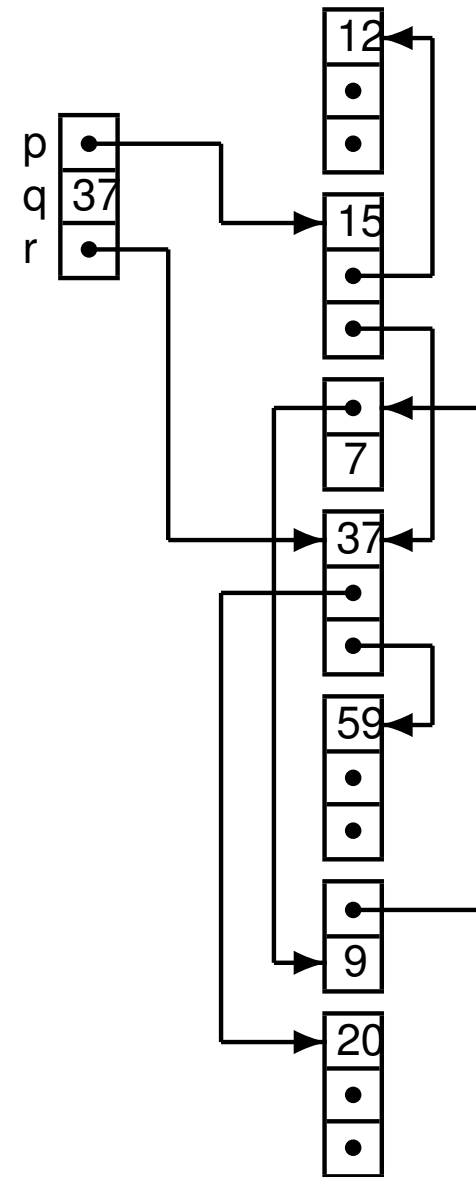
Basic conservative assumption

- A record is *live* if it is reachable from a stack-based program variable (or global variable), otherwise dead.

Example Heap

Consider the following example heap, with stack variables: p , q and r .

- Which records are live?
- Which records are dead?



Runtime Deallocation Mechanisms

A garbage collector

- Is part of the runtime system; and
- Automatically reclaims heap-allocated records that are no longer used.

A garbage collector should

- Reclaim *all* unused records;
- Spend very little time per record;
- Not cause significant delays; and
- Allow all of memory to be used.

These are difficult and often conflicting requirements.

Garbage Collection

In this class we will study 3 types of garbage collection:

- Reference counting;
- Mark-and-sweep; and
- Stop-and-copy.

For each algorithm we will discuss the implementation, an example, and the associated advantages/disadvantages.

Garbage Collection

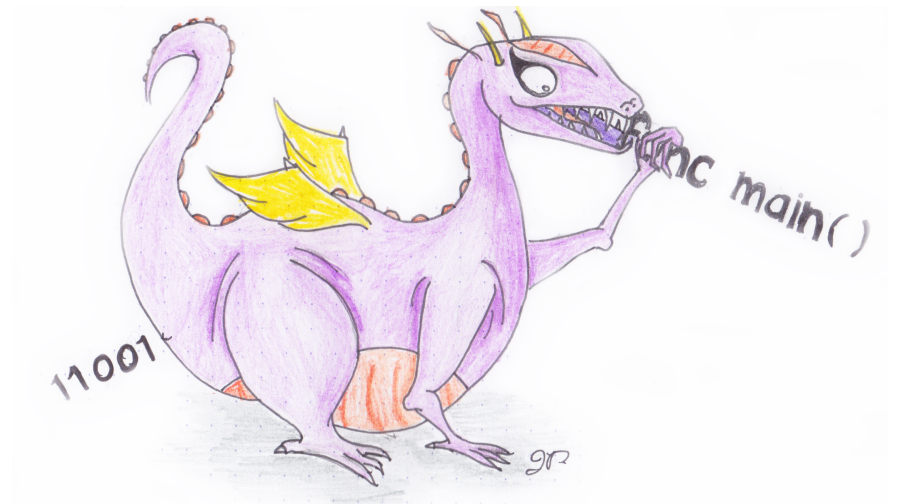
Memory management

Reference counting

Mark-and-sweep

Stop-and-copy

Practical



Reference Counting

- Is a type of continuous (or incremental) garbage collection;
- Uses a field on each object (the reference count) to track incoming pointers; and
- Determines an object is dead when its reference count reaches zero.

The reference count is updated

- Whenever a reference is changed;
 - Created
e.g. `int *a = b; // b refcount++`
 - Destroyed
e.g. `a = c; // b refcount--`
- Whenever a local variable goes out of scope;
- Whenever an object is deallocated (all objects it points to have their reference counts decremented).

Reference Counting

Reference counting inserts calls to `Increase` and `Decrease` in the source program as needed. When the object is no longer needed, the call to `Free` is made.

Pseudo code for reference counting

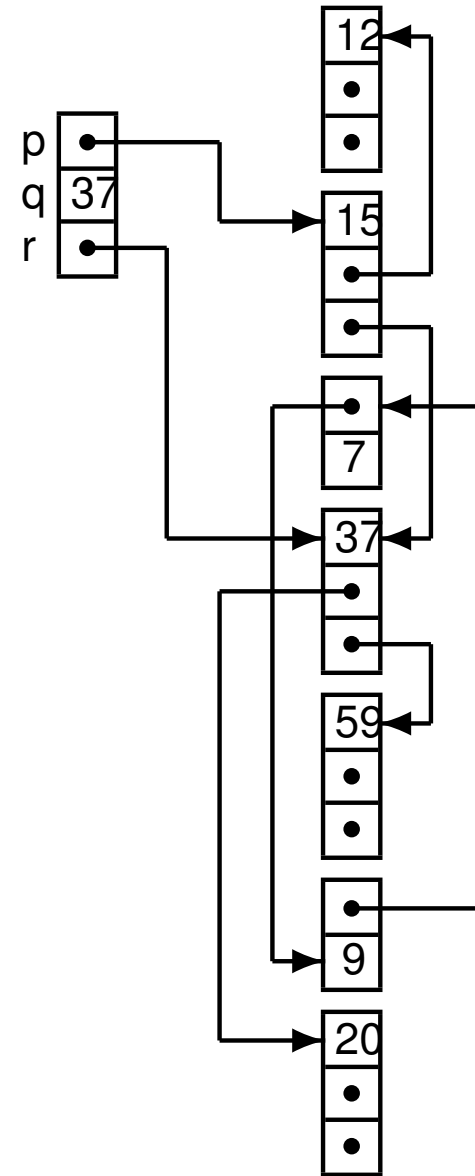
```
function Increase(x)  
    x.count := x.count+1
```

```
function Decrease(x)  
    x.count := x.count-1  
    if x.count = 0 then  
        Free(x)
```

```
function Free(x)  
    for i:=1 to |x| do  
        Decrease(x.fi)  
    x.f1 := freelist  
    freelist := x
```

Reference Counting

Reference counting has one large problem:
What about objects 7 and 9?



Reference Counting

Advantages

- Is incremental, distributing the cost over a long period;
- Does not require long pauses to handle deallocations;
- Catches dead objects immediately; and
- Requires no effort from the user.

Disadvantages

- Is incremental, slowing down the program continuously and unnecessarily;
- Requires a more complex runtime system; and
- Cannot handle circular data structures.

Aside: Automatic Reference Counting (ARC)

Initially for Objective-C (now also for Swift), *automatic reference counting* (ARC) is a reference counting implementation designed by Apple and integrated into Clang.

- Inserts calls to `retain` (increment) and `release` at compile time;
- **Optimizes away unnecessary updates**; and
- Is preferred to garbage collection.

Previously, developers inserted calls to the memory management methods.

Garbage Collection

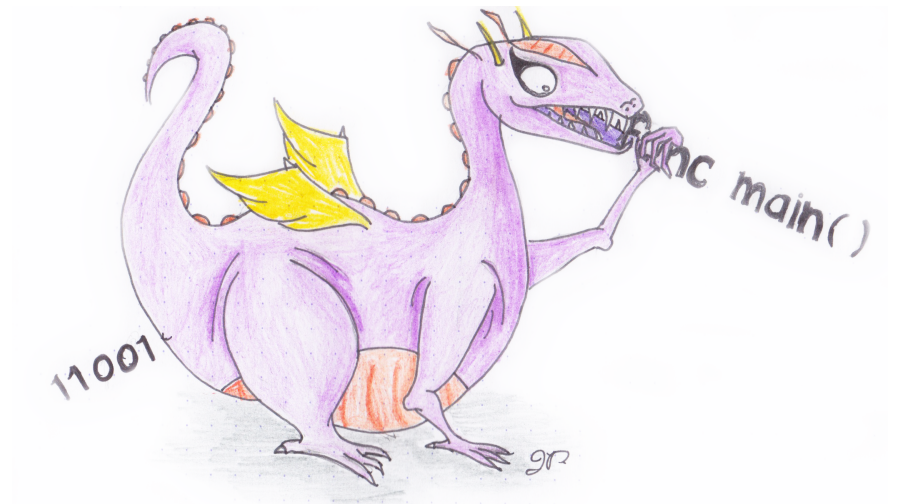
Memory management

Reference counting

Mark-and-sweep

Stop-and-copy

Practical



Mark-and-Sweep

The mark-and-sweep algorithm is a periodic approach to garbage collection that has 3 main steps:

1. Explore pointers starting from the program (stack) variables, and *mark* all records encountered;
2. *Sweep* through all records in the heap and reclaim the unmarked ones; and
3. Finish by unmarking all marked records.

Assumptions

- We know which fields are pointers;
- We know the size of each record; and
- Reclaimed records are kept in a `freelist`.

Mark-and-Sweep

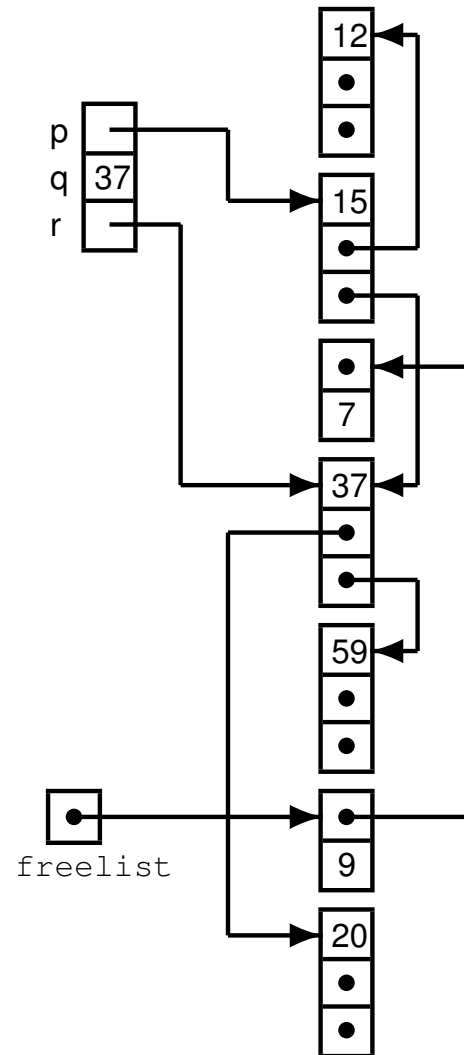
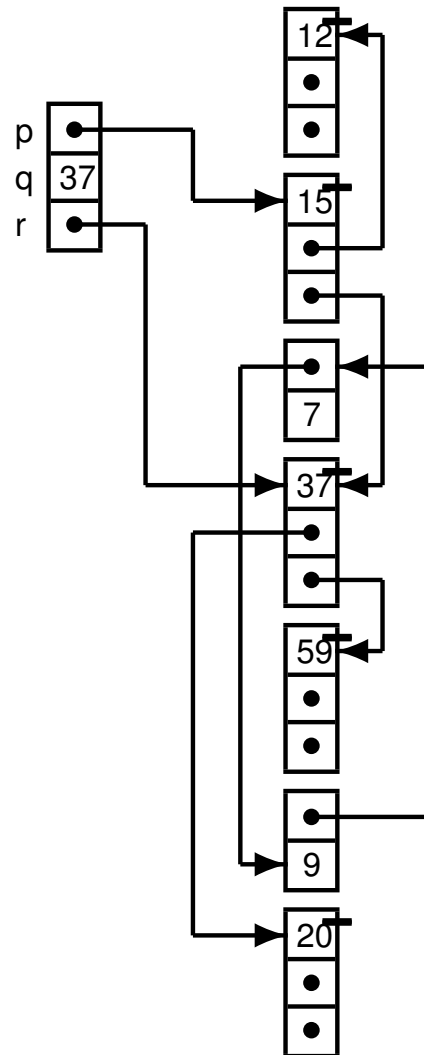
The 3 steps of the mark-and-sweep algorithm are shown below (steps 2 and 3 are merged).

Pseudo code for mark-and-sweep

```
function Mark()  
  for each program variable  $v$  do  
    DFS( $v$ )  
  
function DFS( $x$ )  
  if  $x$  is a pointer into the heap then  
    if record  $x$  is not marked then  
      mark record  $x$   
      for  $i:=1$  to  $|x|$  do  
        DFS( $x.f_i$ )
```

```
function Sweep()  
   $p :=$  first address in heap  
  while  $p <$  last address in heap do  
    if record  $p$  is marked then  
      unmark record  $p$   
    else  
       $p.f_1 :=$  freelist  
      freelist  $:= p$   
   $p := p + \text{sizeof}(\text{record } p)$ 
```

Mark-and-Sweep



Analysis of Mark-and-Sweep

- Assume the heap has size H words; and
- Assume that R words are reachable.

The cost of garbage collection

$$c_1 R + c_2 H$$

Realistic values

$$10R + 3H$$

The cost per reclaimed word

$$\frac{c_1 R + c_2 H}{H - R}$$

- If R is close to H , then this is expensive;
- The lower bound is c_2 ;
- Increase the heap when $R > 0.5H$; then
- The cost per word is $c_1 + 2c_2 \approx 16$.

Mark-and-Sweep

Advantages

- Is periodic, so does not slow down each operation in your program;
- Can be run in parallel to your program;
- Mark and sweep steps can be parallelized too;
- Requires no effort from the user.

Disadvantages

- Scanning the heap can be expensive;
- The heap may become *fragmented*: containing many small free records but none that are large enough for the next allocation.

Heap Fragmentation

To deal with fragmented heaps we can use *compaction*.

- Once mark-and-sweep has finished, collect all live objects at the beginning of the heap;
- Adjust pointers pointing to all moved objects;
- The adjustment depends on the amount of space freed before the object;
- This removes fragmentation and improves locality.

This is not possible in all programming languages as garbage collection is conservative.

⇒ How do we know which fields are pointers?

DFS Recursion Stack

For mark-and-sweep, the recursion stack could have size H (and has at least size $\log H$). However, given that we have (potentially) run out of memory, it may be impossible to allocate!

Pointer reversal

Pointer reversal is a clever method for traversing a tree that embeds the stack in the nodes themselves.

Idea: Instead of recursively calling DFS, temporarily update the field in the record to point to the *parent* node in the traversal.

<http://www.cs.cmu.edu/afs/cs/academic/class/15745-s06/web/handouts/garbage.pdf>

Garbage Collection

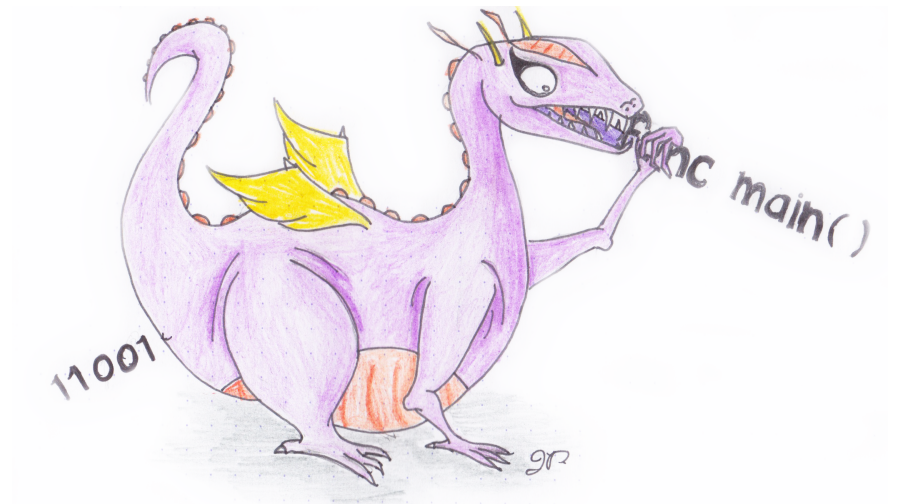
Memory management

Reference counting

Mark-and-sweep

Stop-and-copy

Practical



Stop-and-Copy

The stop-and-copy algorithm is a periodic approach to garbage collection that:

- Divides the heap into two parts;
- Only uses one part at a time;

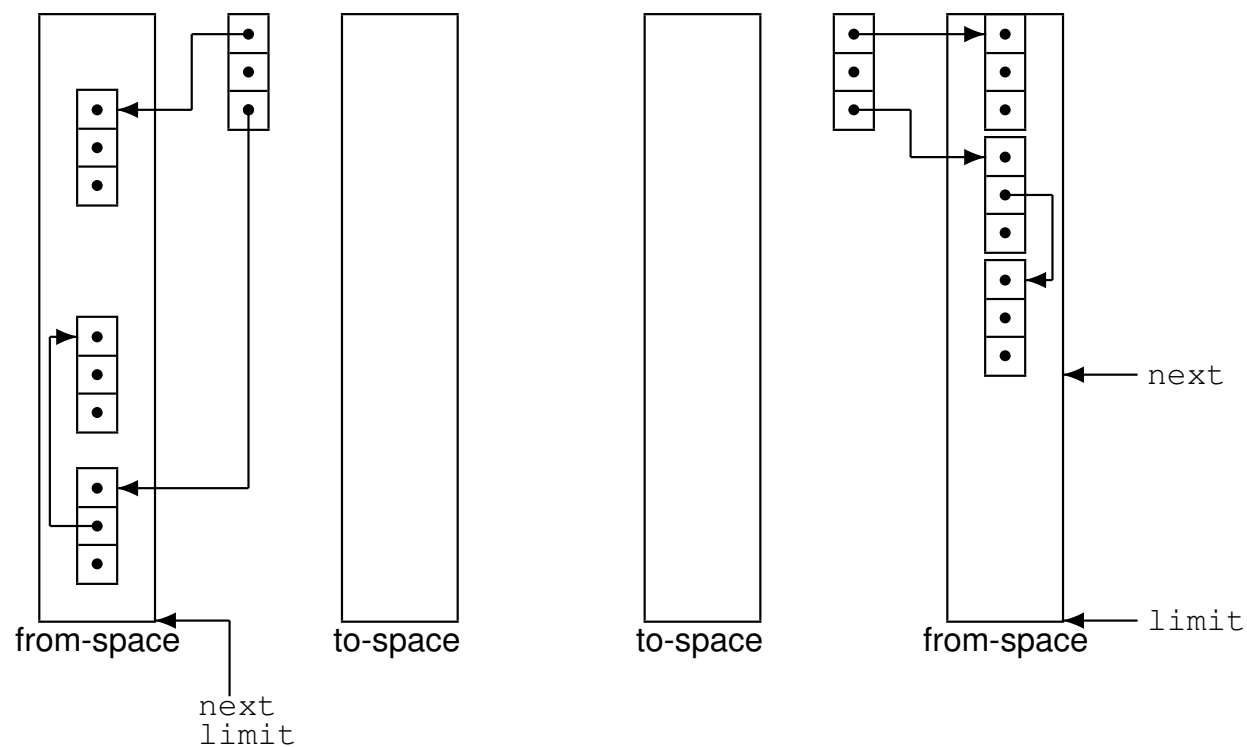
Conceptually this results in a simple high-level algorithm:

1. Use the active half of the heap for all allocations;
2. When it runs full, copy live records to the other part; and
3. Switch the roles of the two parts.

Stop-and-Copy

Consider the following snapshots of stop-and-copy before/after execution.

- `next` and `limit` indicate the available heap space; and
- Copied records are contiguous in memory.



Stop-and-Copy

The stop-and-copy algorithm internals are much more complicated. Intuitively, it *forwards* each record on the heap in a breadth-first manner (starting from the stack).

Pseudo code for stop-and-copy

```

function Copy()
  scan := next := start of to-space
  for each program variable  $v$  do
     $v := \text{Forward}(v)$ 
  while scan < next do
    for  $i:=1$  to |scan| do
      scan. $f_i := \text{Forward}(\text{scan}.f_i)$ 
    scan := scan + sizeof(record scan)

```

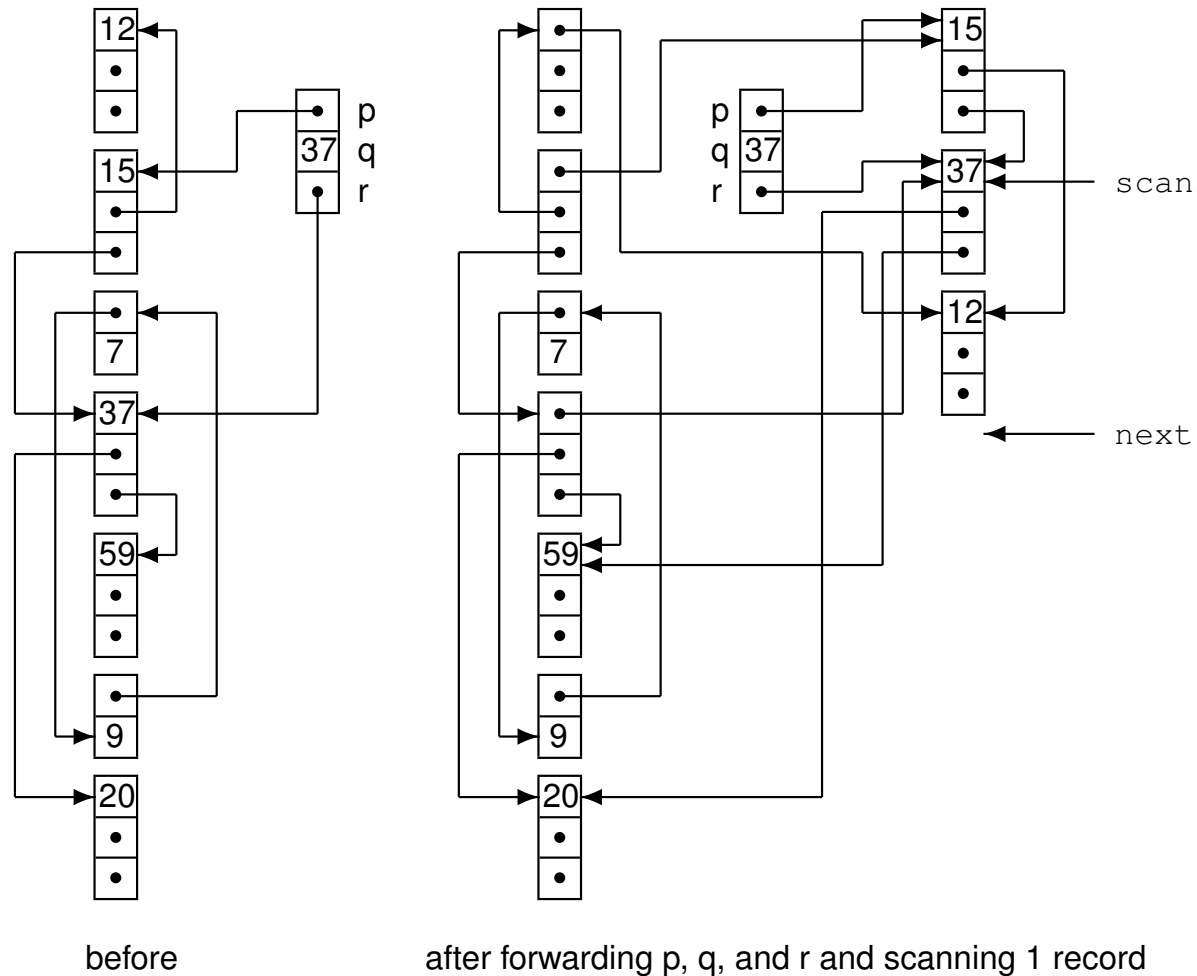
```

function Forward( $p$ )
  if  $p \in$  from-space then
    if  $p.f_1 \in$  to-space then
      return  $p.f_1$ 
    else
      for  $i:=1$  to | $p$ | do
        next. $f_i := p.f_i$ 
       $p.f_1 := \text{next}$ 
      next := next + sizeof(record  $p$ )
      return  $p.f_1$ 
  else return  $p$ 

```

Stop-and-Copy

The follow are snapshots of stop-and-copy before executing and after forwarding the top-level and scanning 1 record.



Analysis of Stop-and-Copy

- Assume the heap has size H words; and
- Assume that R words are reachable.

The cost of garbage collection

$$c_3 R$$

A realistic value

$$10R$$

The cost per reclaimed word

$$\frac{c_3 R}{\frac{H}{2} - R}$$

- This has no lower bound as H grows;
- If $H = 4R$ then the cost is $c_3 \approx 10$.

Stop-and-Copy

Advantages

- Allows fast allocation (no `freelist`);
- Avoids fragmentation;
- Collects in time proportional to R ; and
- Avoids stack and pointer reversal.

Disadvantage

- Wastes half your memory; and
- Stops the program to execute.

Garbage Collection

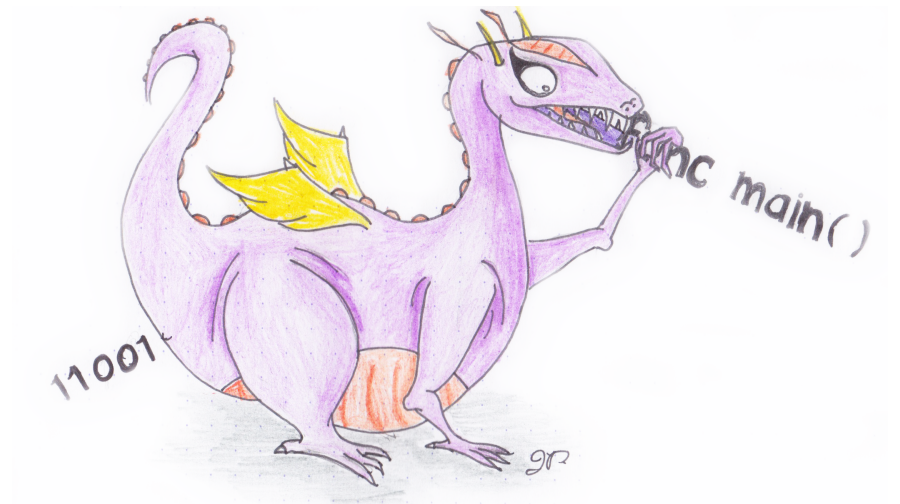
Memory management

Reference counting

Mark-and-sweep

Stop-and-copy

Practical



Practical Considerations

In practice, we use either mark-and-sweep or stop-and-copy (and in some systems ref. counting).

This can lead to better memory management, but garbage collection is still expensive: ≈ 100 instructions for a small object!

Each algorithm can be further extended by

- Generational collection (to make it run faster); and
- Incremental (or concurrent) collection (to make it run smoother).

Generational Collection

Observation: the young die quickly

Given this assumption, the garbage collector should:

- Focus on young records;
- Divide the heap into generations: G_0, G_1, G_2, \dots ;
- All records in G_i are younger than records in G_{i+1} ;
- Collect G_0 often, G_1 less often, and so on; and
- Promote a record from G_i to G_{i+1} when it survives several collections.

Generational Collection

How to collect the G_0 generation

- It might be very expensive to find those pointers;
- Fortunately, they are rare; so
- We can try to remember them.

Ways to remember

- Maintain a list of all updated records (use marks to make this a set); or
- Mark pages of memory that contain updated records (in hardware or software).

Other Optimizations

Incremental collection

- Garbage collection may cause long pauses;
- This is undesirable for interactive or real-time programs; so
- Try to interleave the garbage collection with the program execution.

Two players access the heap

- The *mutator*: creates records and moves pointers around; and
- The *collector*: tries to collect garbage.

Some invariants are clearly required to make this work. The mutator will suffer some slowdown to maintain these invariants.

Earlier Assumptions

The presented garbage collection algorithms assumed that:

- We know the size of each record; and
- We know which fields are pointers.

For object-oriented languages, each record already contains a pointer to a class descriptor, so garbage collection is straightforward to implement.

For general languages, we must sacrifice a few bytes per record to indicate its size, and its organization.