Martin Kellogg

Course Announcements

- As some of you have noticed, the PA4 leaderboard results are subject to timing variance
 - In particular, sometimes you "get lucky" and your compiler's code appears to be much, much faster than the reference
 - There's not much I can do about this, unfortunately.
 - We will run each final submission several times and take the median Q score (not the best), so if you see a big speedup without doing anything, you should assume it is ephemeral
- PA4c1 is due on Monday
- I will not hold office hours next Wednesday

Agenda

- Why Automatic Memory Management?
- Garbage Collection
- Three Techniques
 - Mark and Sweep
 - \circ Stop and Copy
 - Reference Counting

• **Storage management** is still a hard problem in modern programming

- **Storage management** is still a hard problem in modern programming
- Programs in languages like C and C++ have many storage bugs

- **Storage management** is still a hard problem in modern programming
- Programs in languages like C and C++ have many storage bugs
 - forgetting to free unused memory

- **Storage management** is still a hard problem in modern programming
- Programs in languages like C and C++ have many storage bugs
 - forgetting to free unused memory
 - dereferencing a dangling pointer

- **Storage management** is still a hard problem in modern programming
- Programs in languages like C and C++ have many storage bugs
 - forgetting to free unused memory
 - dereferencing a dangling pointer
 - overwriting parts of a data structure by accident and so on...
 (can be big security problems)

- **Storage management** is still a hard problem in modern programming
- Programs in languages like C and C++ have many storage bugs
 - forgetting to free unused memory
 - dereferencing a dangling pointer
 - overwriting parts of a data structure by accident and so on...
 (can be big security problems)
- Storage bugs are hard to find

- **Storage management** is still a hard problem in modern programming
- Programs in languages like C and C++ have many storage bugs
 - forgetting to free unused memory
 - dereferencing a dangling pointer
 - overwriting parts of a data structure by accident and so on...
 (can be big security problems)
- Storage bugs are hard to find
 - a bug can lead to a visible effect far away in time and program text from the source

• Some storage bugs can be prevented in a strongly-typed language

- Some storage bugs can be prevented in a strongly-typed language
 - e.g., most type systems guarantee no random access into some other object's private data

- Some storage bugs can be prevented in a strongly-typed language
 - e.g., most type systems guarantee no random access into some other object's private data
- Can types prevent errors in programs with manual allocation and deallocation of memory?

- Some storage bugs can be prevented in a strongly-typed language
 - e.g., most type systems guarantee no random access into some other object's private data
- Can types prevent errors in programs with manual allocation and deallocation of memory?
 - Some fancy type systems (linear types) were designed for this purpose, but they complicate programming significantly

- Some storage bugs can be prevented in a strongly-typed language
 - e.g., most type systems guarantee no random access into some other object's private data
- Can types prevent errors in programs with manual allocation and deallocation of memory?
 - Some fancy type systems (linear types) were designed for this purpose, but they complicate programming significantly
 - So, in theory, yes

- Some storage bugs can be prevented in a strongly-typed language
 - e.g., most type systems guarantee no random access into some other object's private data
- Can types prevent errors in programs with manual allocation and deallocation of memory?
 - Some fancy type systems (linear types) were designed for this purpose, but they complicate programming significantly
 - So, in theory, yes
- If you want type safety in practice then you must use *automatic memory management*

- This is an old problem
 - \circ $\,$ studied since the 1950s for Lisp

- This is an old problem
 - \circ studied since the 1950s for Lisp
- There are several well-known techniques for completely automatic memory management

- This is an old problem
 - studied since the 1950s for Lisp
- There are several well-known techniques for completely automatic memory management
 - These are *dynamic analyses*

- This is an old problem
 - studied since the 1950s for Lisp
- There are several well-known techniques for completely automatic memory management
 - These are *dynamic analyses*
 - That is, they involve instrumenting the program so that its behavior at run time is different

- This is an old problem
 - studied since the 1950s for Lisp
- There are several well-known techniques for completely automatic memory management
 - These are *dynamic analyses*
 - That is, they involve instrumenting the program so that its behavior at run time is different

You are probably familiar with some other dynamic analysis techniques like:

- testing
- code coverage

- This is an old problem
 - studied since the 1950s for Lisp
- There are several well-known techniques for completely automatic memory management
 - These are *dynamic analyses*
 - That is, they involve instrumenting the program so that its behavior at run time is different
 - In particular, we want it to automatically free memory :)

- This is an old problem
 - studied since the 1950s for Lisp
- There are several well-known techniques for completely automatic memory management
 - These are *dynamic analyses*
 - That is, they involve instrumenting the program so that its behavior at run time is different
 - In particular, we want it to automatically free memory :)
- Until relatively recently (Java), these techniques were not popular outside the Lisp family of languages
 - Just like static type safety used to be unpopular...

- When an object that takes memory space is created, unused space is automatically allocated
 - In Cool, new objects are created by new X

- When an object that takes memory space is created, unused space is automatically allocated
 - In Cool, new objects are created by new X
- After a while there is no more unused space

- When an object that takes memory space is created, unused space is automatically allocated
 - In Cool, new objects are created by new X
- After a while there is no more unused space
- Some space is occupied by objects that will never be used again (= dead objects?)

- When an object that takes memory space is created, unused space is automatically allocated
 - In Cool, new objects are created by new X
- After a while there is no more unused space
- Some space is occupied by objects that will never be used again (= dead objects?)
- This space can be **freed** to be reused later

• How can we tell whether an object will "never be used again"?

How can we tell whether an object will "never be used again"?
 In general it is impossible (undecidable) to tell (cf. liveness)

- How can we tell whether an object will "never be used again"?
 - In general it is impossible (**undecidable**) to tell (cf. liveness)
 - We will have to use a heuristic to find many (not all) objects that will never be used again

- How can we tell whether an object will "never be used again"?
 - In general it is impossible (undecidable) to tell (cf. liveness)
 - We will have to use a heuristic to find many (not all) objects that will never be used again
- Observation: a program can use only the objects that it can find.

- How can we tell whether an object will "never be used again"?
 - In general it is impossible (undecidable) to tell (cf. liveness)
 - We will have to use a heuristic to find many (not all) objects that will never be used again
- **Observation**: a program can use only the objects **that it can find**.
- For example:

let x : A <- new A in { x <- y; ... }

- How can we tell whether an object will "never be used again"?
 - In general it is impossible (undecidable) to tell (cf. liveness)
 - We will have to use a heuristic to find many (not all) objects that will never be used again
- **Observation**: a program can use only the objects **that it can find**.
- For example:

let x : A <- new A in { x <- y; ... }</pre>

 After x <- y there is no way to access the newly allocated object

Garbage

• **Definition:** An object *x* is *reachable* if and only if:


- **Definition:** An object *x* is *reachable* if and only if:
 - A local variable (or register) contains a pointer to *x*, **or**



- **Definition:** An object *x* is *reachable* if and only if:
 - A local variable (or register) contains a pointer to *x*, **or**
 - Another reachable object y contains a pointer to x



- **Definition:** An object *x* is *reachable* if and only if:
 - A local variable (or register) contains a pointer to *x*, **or**
 - Another reachable object y contains a pointer to x
 - (Note that self is a local variable in Cool)



- **Definition:** An object *x* is *reachable* if and only if:
 - A local variable (or register) contains a pointer to *x*, **or**
 - Another reachable object y contains a pointer to x
 (Note that self is a local variable in Cool)
- You can find all reachable objects by starting from local variables and following all the pointers ("transitive")



- **Definition:** An object *x* is *reachable* if and only if:
 - A local variable (or register) contains a pointer to *x*, **or**
 - Another reachable object y contains a pointer to x
 (Note that self is a local variable in Cool)
- You can find all reachable objects by starting from local variables and following all the pointers ("transitive")
- An unreachable object can never be referred to by the program



- **Definition:** An object *x* is *reachable* if and only if:
 - A local variable (or register) contains a pointer to *x*, **or**
 - Another reachable object y contains a pointer to x
 (Note that self is a local variable in Cool)
- You can find all reachable objects by starting from local variables and following all the pointers ("transitive")
- An unreachable object can never be referred to by the program
 - Such objects are called *garbage*



- Consider the program:
 - x <- new Ant; y <- new Bat; x <- y; if alwaysTrue() then x <- new Cow else x.eat() fi</pre>

- Consider the program:
 - x <- new Ant;
 - y <- new Bat;

- if alwaysTrue() then x <- new Cow else x.eat() fi</pre>
- After x <- y (assuming y becomes dead there):

- Consider the program:
 - x <- new Ant;</pre>
 - y <- new Bat;

- if alwaysTrue() then x <- new Cow else x.eat() fi</pre>
- After x <- y (assuming y becomes dead there):
 - The Ant object is not reachable anymore

- Consider the program:
 - x <- new Ant;</pre>
 - y <- new Bat;

- if alwaysTrue() then x <- new Cow else x.eat() fi</pre>
- After x <- y (assuming y becomes dead there):
 - The Ant object is not reachable anymore
 - The Bat object is reachable (through x)

- Consider the program:
 - x <- new Ant;</pre>
 - y <- new Bat;

- if alwaysTrue() then x <- new Cow else x.eat() fi</pre>
- After x <- y (assuming y becomes dead there):
 - The Ant object is not reachable anymore
 - The Bat object is reachable (through x)
 - Thus the Bat is not garbage and is *not* collected

- Consider the program:
 - x <- new Ant;</pre>
 - y <- new Bat;

- if alwaysTrue() then x <- new Cow else x.eat() fi</pre>
- After x <- y (assuming y becomes dead there):
 - The Ant object is not reachable anymore
 - The Bat object is reachable (through x)
 - Thus the Bat is not garbage and is *not* collected
 - But the Bat object is never going to be used

• Recall that at run-time, a Cool interpreter has two mappings:

Operational Semantics!

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 o for each location I ∈ domain(S)

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location $I \in \text{domain}(S)$
 - let can_reach = false

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location $I \in \text{domain}(S)$
 - let can_reach = false
 - for each $(v, I_2) \in E$

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location I ∈ domain(S)
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location I ∈ domain(S)
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true
 - if not can_reach then reclaim_location(l)

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location $I \in \text{domain}(S)$
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true
 - if not can_reach then reclaim_location(l)

Does this work?

Does That Work?



- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location I ∈ domain(S)
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location I ∈ domain(S)
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true
 - for each $I_3 \in v$

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location I ∈ domain(S)
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true
 - for each $I_3 \in v$ // v is X(..., $a_i = I_i$, ...)

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location $I \in \text{domain}(S)$
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true
 - for each $I_3 \in v$ // v is X(..., $a_i = I_i$, ...) • if $I = I_3$ then can_reach = true

- Recall that at run-time, a Cool interpreter has two mappings:
 - Environment E maps variable identifiers to locations
 - Store **S** maps locations to values
- Proposed Cool garbage collector algorithm:
 - for each location $I \in \text{domain}(S)$
 - let can_reach = false
 - for each $(v, I_2) \in E$
 - if $I = I_2$ then can_reach = true
 - for each $I_3 \in v$ // v is X(..., $a_i = I_i$, ...) • if $I = I_3$ then can_reach = true
 - if not can_reach then reclaim_location(l)

Garbage Analysis

• Could we use this proposed Cooler Garbage Collector in real life?



Garbage Analysis

- Could we use this proposed Cooler Garbage Collector in real life?
 - How long would it take?
 - How much space would it take?
 - Are we forgetting anything?



Garbage Analysis

- Could we use this proposed Cooler Garbage Collector in real life?
 - How long would it take?
 - How much space would it take?
 - Are we forgetting anything?
 - Hint: Yes. It's still wrong



• In Cool, local variables are easy to find:

- In Cool, local variables are easy to find:
 - Use the environment mapping **E**

- In Cool, local variables are easy to find:
 - Use the environment mapping **E**
 - and one object may point to other objects, etc.
- In Cool, local variables are easy to find:
 - Use the environment mapping **E**
 - and one object may point to other objects, etc.
- The **stack** is more complex:

- In Cool, local variables are easy to find:
 - Use the environment mapping **E**
 - and one object may point to other objects, etc.
- The **stack** is more complex:
 - each stack frame (activation record) contains:

- In Cool, local variables are easy to find:
 - Use the environment mapping **E**
 - and one object may point to other objects, etc.
- The **stack** is more complex:
 - $\circ~$ each stack frame (activation record) contains:
 - method parameters! (which are other objects...)

- In Cool, local variables are easy to find:
 - Use the environment mapping E
 - and one object may point to other objects, etc.
- The **stack** is more complex:
 - each stack frame (activation record) contains:
 - method parameters! (which are other objects...)
- If we know the layout of a stack frame we can find the pointers (objects) in it

- In Cool, local variables are easy to find:
 - Use the environment mapping E
 - and one object may point to other objects, etc.
- The stack is more cor
 - \circ each stack frame
 - method para
- If we know the layou (objects) in it
- Many things may look legitimate and reachable but will turn out not to be
- How can we figure this out systematically?

Reachability can be tricky!

...) e pointers

A Simple Example

A Simple Example

local _____ variables

A Simple Example











• Start tracing from local vars and the stack



Start tracing from local vars and the stack
 They are called the *roots*



- Start tracing from local vars and the stack
 - They are called the *roots*
- Note that B and D are not reachable from other local vars or the stack
 - Thus we can reuse their storage when they go out of scope

• Every garbage collection scheme has the following steps

Every garbage collection scheme has the following steps
 Allocate space as needed for new objects

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:
 - Compute what objects might be used again

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:
 - Compute what objects might be used again
 - generally by tracing objects reachable from a set of roots

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:
 - Compute what objects might be used again
 - generally by tracing objects reachable from a set of roots
 - Free space used by objects not found in the previous step

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:
 - Compute what objects might be used again
 - generally by tracing objects reachable from a set of roots
 - Free space used by objects not found in the previous step
- Some strategies perform garbage collection **before** the space actually runs out

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:
 - Compute what objects might be used again
 - generally by tracing objects reachable from a set of roots
 - Free space used by objects not found in the previous step
- Some strategies perform garbage collection before the space actually runs out
 - Why might this be useful?

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:
 - Compute what objects might be used again
 - generally by tracing objects reachable from a set of roots
 - Free space used by objects not found in the previous step
- Some strategies perform garbage collection before the space actually runs out
 - Why might this be useful?
 - Hint: will we need any space to run our garbage collector?

- Every garbage collection scheme has the following steps
 - Allocate space as needed for new objects
 - When space runs out:

Compute where	After trivia, we will see three	
• generally	specific garbage collection	m a set of
roots	algorithms:	
Free space u	 mark and sween 	revious step
Some strategies per	 stop and copy 	he space
actually runs out	reference counting)
• why hight this be userui:		
Hint: will we need any space to run our garbage collector?		

Trivia Break: Music Theory

This musical symbol (examples highlighted in **blue** below) indicates which notes are represented by the lines and spaces on a musical staff. Its position on a staff assigns a particular pitch to one of the five lines or four spaces, which defines the pitches on the remaining lines and spaces. The modern symbols derive from the medieval practice of annotating the reference line of a staff with the name of the note it was intended to bear; over time the shapes of these letters became stylised, leading to their current versions.



Trivia Break: Programming Languages

This general-purpose high-level programming language supports multiple paradigms. Its features include automatic memory management, a strong type system, and good support for internationalization and portability. Its creators originally released it in 2000 as a closed-source language, aligning with their business goals at the time. However, in the decades since, the company responsible for the language has changed its attitude towards open-source, and in 2014 the open-source Roslyn compiler for this language was released; it has been the primary compiler for the language since. The language is famously used in game development (e.g., it is the default scripting language in the Unity game engine).

• Our first garbage collection algorithm

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects
 - the **sweep** phase: collects garbage objects

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects
 - the *sweep* phase: collects garbage objects
- Every object has an extra bit: the *mark bit*

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects
 - the *sweep* phase: collects garbage objects
- Every object has an extra bit: the *mark bit*
 - reserved for memory management
 - initially the mark bit is 0

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects
 - the *sweep* phase: collects garbage objects
- Every object has an extra bit: the *mark bit*
 - reserved for memory management
 - initially the mark bit is 0
 - set to 1 for the **reachable** objects in the mark phase

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects
 - the *sweep* phase: collects garbage objects
- Every object has an extra bit: the *mark bit*
 - reserved for memory management
 - initially the mark bit is 0
 - set to 1 for the **reachable** objects in the mark phase
- In the sweep phase, free all objects whose mark bit is still 0

- Our first garbage collection algorithm
- When memory runs out, GC executes two phases:
 - the *mark* phase: traces reachable objects
 - the *sweep* phase: collects garbage objects
- Every object has an extra bit: the *mark bit*
 - reserved for memory management
 - initially the mark bit is 0
 - set to 1 for the **reachable** objects in the mark phase
- In the sweep phase, free all objects whose mark bit is still 0
 creating a free list of garbage that can be reused

Mark and Sweep: Example
Mark and Sweep: Example



Mark and Sweep: Example







let todo = { all roots } (* worklist *)

let todo = { all roots } (* worklist *)
while todo is non-empty; do

```
let todo = { all roots } (* worklist *)
while todo is non-empty ; do
    pick v \in todo
    todo <- todo - { v }</pre>
```

```
let todo = { all roots } (* worklist *)
while todo is non-empty ; do
    pick v ∈ todo
    todo <- todo - { v }
    if mark(v) = 0 then (* v is unmarked so far *)</pre>
```

```
let todo = { all roots } (* worklist *)
while todo is non-empty ; do
    pick v \in todo
    todo <- todo - { v }
    if mark(v) = 0 then (* v is unmarked so far *)
    mark(v) <- 1</pre>
```

```
let todo = { all roots } (* worklist *)

while todo is non-empty ; do

pick v \in todo

todo <- todo - { v }

if mark(v) = 0 then (* v is unmarked so far *)

mark(v) <- 1

let v_1, ..., v_n be the pointers contained in v
```

```
let todo = { all roots } (* worklist *)
while todo is non-empty; do
    pick v \in todo
    todo <- todo - { v }
    if mark(v) = 0 then (* v is unmarked so far *)
        mark(v) <- 1
        let v_1, ..., v_n be the pointers contained in v
        todo < -todo U \{v_1, ..., v_n\}
```

```
let todo = { all roots } (* worklist *)
while todo is non-empty; do
    pick v \in todo
    todo <- todo - { v }
    if mark(v) = 0 then (* v is unmarked so far *)
        mark(v) <- 1
        let v_1, ..., v_n be the pointers contained in v
        todo <- todo U \{v_1, ..., v_n\}
    fi
```

od



• The sweep phase scans the entire heap looking for objects with mark bit 0



- The sweep phase scans the entire heap looking for objects with mark bit 0
 - these objects have not been visited in the mark phase



- The sweep phase scans the entire heap looking for objects with mark bit 0
 - these objects have not been visited in the mark phase
 - \circ and so they are garbage



- The sweep phase scans the entire heap looking for objects with mark bit 0
 - these objects have not been visited in the mark phase
 - and so they are garbage
- Any such object is added to the free list



- The sweep phase scans the entire heap looking for objects with mark bit 0
 - these objects have not been visited in the mark phase
 - and so they are garbage
- Any such object is added to the free list
- The objects with a mark bit of 1 have their mark bit reset to 0



/* sizeof(p) is size of block starting at p */

/* sizeof(p) is size of block starting at p */
p <- bottom of heap</pre>

/* sizeof(p) is size of block starting at p */
p <- bottom of heap
while p < top of heap do</pre>

/* sizeof(p) is size of block starting at p */
p <- bottom of heap
while p < top of heap do
 if mark(p) = 1 then
 mark(p) <- 0</pre>

```
/* sizeof(p) is size of block starting at p */
p < -bottom of heap
while p < top of heap do
    if mark(p) = 1 then
        mark(p) < -0
    else
        add block p...(p+sizeof(p)-1) to freelist
    fi
```

```
/* sizeof(p) is size of block starting at p */
p <- bottom of heap
while p < top of heap do
    if mark(p) = 1 then
        mark(p) < -0
    else
        add block p...(p+sizeof(p)-1) to freelist
    fi
    p < -p + sizeof(p)
od
```

• While conceptually simple, this algorithm has a number of tricky details

- While conceptually simple, this algorithm has a number of tricky details
 - this is typical of GC algorithms

- While conceptually simple, this algorithm has a number of tricky details
 - this is typical of GC algorithms
- There is a **serious problem** with the mark phase:

- While conceptually simple, this algorithm has a number of tricky details
 - this is typical of GC algorithms
- There is a **serious problem** with the mark phase:
 - \circ $\,$ it is invoked when we are out of space $\,$

- While conceptually simple, this algorithm has a number of tricky details
 - this is typical of GC algorithms
- There is a **serious problem** with the mark phase:
 - \circ $\,$ it is invoked when we are out of space
 - yet it needs space to construct the todo list

- While conceptually simple, this algorithm has a number of tricky details
 - this is typical of GC algorithms
- There is a **serious problem** with the mark phase:
 - it is invoked when we are out of space
 - yet it needs space to construct the todo list
 - the size of the todo list is unbounded, so we cannot reserve space for it a priori

• The todo list is used as an auxiliary data structure to perform the reachability analysis

- The todo list is used as an auxiliary data structure to perform the reachability analysis
- There is a trick that allows the auxiliary data to be stored in the objects themselves

- The todo list is used as an auxiliary data structure to perform the reachability analysis
- There is a trick that allows the auxiliary data to be stored in the objects themselves
 - pointer reversal: when a pointer is followed it is reversed to point to its parent

- The todo list is used as an auxiliary data structure to perform the reachability analysis
- There is a trick that allows the auxiliary data to be stored in the objects themselves
 - pointer reversal: when a pointer is followed it is reversed to point to its parent
- Similarly, the free list is stored in the free objects themselves
• Space for each new object is allocated from the free list

- Space for each new object is allocated from the free list
 - a block large enough is picked

- Space for each new object is allocated from the free list
 - a block large enough is picked
 - an area of the necessary size is allocated from it
 - the leftover is put back in the free list

- Space for each new object is allocated from the free list
 - a block large enough is picked
 - \circ $\,$ an area of the necessary size is allocated from it $\,$
 - the leftover is put back in the free list
- Disadvantage: mark and sweep can **fragment** memory
 - why is this a problem?

- Space for each new object is allocated from the free list
 - a block large enough is picked
 - \circ $\,$ an area of the necessary size is allocated from it $\,$

the leftover is put back in the free list

- Disadvantage: mark and sweep can **fragment** memory
 - why is this a problem?
- Advantage: objects are **not moved** during GC
 - no need to update the pointers to objects
 - works for languages like C and C++ where it's difficult to distinguish pointers from data (more on this later)

• Memory is organized into two areas:

- Memory is organized into two areas:
 - **Old space**: used for allocation

- Memory is organized into two areas:
 - Old space: used for allocation
 - *New space*: used as a reserve for the GC

- Memory is organized into two areas:
 - Old space: used for allocation
 - New space: used as a reserve for the GC



- Memory is organized into two areas:
 - Old space: used for allocation
 - *New space*: used as a reserve for the GC



• The heap pointer points to the next free word in the old space

- Memory is organized into two areas:
 - Old space: used for allocation
 - *New space*: used as a reserve for the GC



The heap pointer points to the next free word in the old space
 Allocation just advances the heap pointer



• Starts when the old space is full



- Starts when the old space is full
- Copies all reachable objects from old space into new space



- Starts when the old space is full
- Copies all reachable objects from old space into new space
 - garbage is left behind



- Starts when the old space is full
- Copies all reachable objects from old space into new space
 - garbage is left behind
 - after the copy phase the new space uses less space than the old one before the collection



- Starts when the old space is full
- Copies all reachable objects from old space into new space
 - garbage is left behind
 - after the copy phase the new space uses less space than the old one before the collection
- After the copy the roles of the old and new spaces are reversed and the program resumes



Stop and Copy: Simple Example

Stop and Copy: Simple Example



Stop and Copy: Simple Example



- We need to find all the reachable objects
 - Just as in mark and sweep

- We need to find all the reachable objects
 - Just as in mark and sweep
- As we find a reachable object we copy it into the new space

- We need to find all the reachable objects
 - Just as in mark and sweep
- As we find a reachable object we copy it into the new space
 - And we have to fix ALL pointers pointing to it!

- We need to find all the reachable objects
 - Just as in mark and sweep
- As we find a reachable object we copy it into the new space
 And we have to fix ALL pointers pointing to it!
- As we copy an object we store in the old copy a *forwarding pointer* to the new copy

- We need to find all the reachable objects
 - Just as in mark and sweep
- As we find a reachable object we copy it into the new space
 And we have to *fix ALL pointers pointing to it*!
- As we copy an object we store in the old copy a *forwarding pointer* to the new copy
 - When we later reach an object with a forwarding pointer, we know it was already copied

- We need to find all the reachable objects
 - Just as in mark and sweep
- As we find a reachable object we copy it into the new space
 And we have to *fix ALL pointers pointing to it*!
- As we copy an object we store in the old copy a *forwarding pointer* to the new copy
 - When we later reach an object with a forwarding pointer, we know it was already copied
 - How can we identify forwarding pointers?

- We need to find all the reachable objects
 - Just as in mark and sweep
- As we find a reachable object we copy it into the new space
 And we have to *fix ALL pointers pointing to it*!
- As we copy an object we store in the old copy a *forwarding pointer* to the new copy
 - When we later reach an object with a forwarding pointer, we know it was already copied
 - How can we identify forwarding pointers?
- We also still have the issue of how to implement the traversal without using extra space

• Both problems can be solved via the following trick:

- Both problems can be solved via the following trick:
 - partition the new space into three contiguous regions:

- Both problems can be solved via the following trick:
 - partition the **new space** into **three** contiguous regions:

- Both problems can be solved via the following trick:
 - partition the new space into three contiguous regions:

start

copied and scanned

copied objects whose pointer fields were followed and fixed

- Both problems can be solved via the following trick:
 - partition the **new space** into **three** contiguous regions:

start ↓	scan
copied and scanned	copied
copied objects whose pointer fields were followed and fixed	copied objects whose pointer fields were NOT followed

- Both problems can be solved via the following trick:
 - partition the new space into three contiguous regions:

start	scan a	alloc
copied and scanned	copied	empty
copied objects whose pointer fields were followed and fixed	copied objects whose pointer fields were NOT followed	_
• Before garbage collection:



• Step 1: Copy the objects pointed by roots and set forwarding pointers (dotted arrow)



- Step 2: Follow the pointer in the next unscanned object (A)
 - copy the pointed objects (just C in this case)
 - fix the pointer in A
 - set forwarding pointer
 alloc
 start
 root
 A B C D E F A C

Step 3: Follow the pointer in the next unscanned object (C)
 o copy the pointed objects (F in this case)



- Step 4: Follow the pointer in the next unscanned object (F)
 - the pointed object (A) was already copied. Set the pointer same as the forwarding pointer



- Step 5: Since scan caught up with alloc, we are done
- Now, we swap the role of the two spaces and then resume the program



Stop and Copy: Algorithm

Stop and Copy: Algorithm

while scan != alloc do let O be the object at the scan pointer for each pointer p contained in O do find O' that p points to if O' is without a forwarding pointer copy O' to new space (update alloc pointer) set 1st word of old O' to point to the new copy change p to point to the new copy of O'else set p in O equal to the forwarding pointer fi rof increment scan pointer to the next object

od

- As with mark and sweep, we must be able to tell how large an object is when we scan it
 - Is this a problem in Cool?

- As with mark and sweep, we must be able to tell how large an object is when we scan it
 - Is this a problem in Cool?
- We also need to know where the pointers are inside the object
 How hard is this?

- As with mark and sweep, we must be able to tell how large an object is when we scan it
 - Is this a problem in Cool?
- We also need to know where the pointers are inside the object
 O How hard is this?
- We must also copy any objects pointed to by the stack and update pointers in the stack

- As with mark and sweep, we must be able to tell how large an object is when we scan it
 - Is this a problem in Cool?
- We also need to know where the pointers are inside the object
 O How hard is this?
- We must also copy any objects pointed to by the stack and update pointers in the stack
 - This can be an **expensive** operation in practice...

• Stop and copy is generally believed to be the fastest GC technique

- Stop and copy is generally believed to be the **fastest** GC technique
- Allocation is very cheap
 - Just increment the heap pointer
- Collection is relatively cheap
 - Especially if there is a lot of garbage
 - Only touch reachable objects
- But some languages do not allow copying
 - C, C++, ...

- Stop and copy is generally believed to be the **fastest** GC technique
- Allocation is very cheap
 - Just increment the heap pointer

- Stop and copy is generally believed to be the **fastest** GC technique
- Allocation is very cheap
 - Just increment the heap pointer
- Collection is relatively cheap
 - Especially if there is a lot of garbage
 - Only touch reachable objects

- Stop and copy is generally believed to be the **fastest** GC technique
- Allocation is very cheap
 - Just increment the heap pointer
- Collection is relatively cheap
 - Especially if there is a lot of garbage
 - Only touch reachable objects
- But some languages do not allow copying
 - C, C++, ...

Garbage collection relies on being able to find all reachable objects
 And it needs to find all pointers in an object

- Garbage collection relies on being able to find all reachable objects
 And it needs to find all pointers in an object
- In languages like C or C++ it is **impossible** to identify the contents of objects in memory

- Garbage collection relies on being able to find all reachable objects
 And it needs to find all pointers in an object
- In languages like C or C++ it is **impossible** to identify the contents of objects in memory
 - e.g., how can you tell that a sequence of two memory words is a list cell (with data and next fields) or a binary tree node (with a left and right fields)?

- Garbage collection relies on being able to find all reachable objects
 And it needs to find all pointers in an object
- In languages like C or C++ it is **impossible** to identify the contents of objects in memory
 - e.g., how can you tell that a sequence of two memory words is a list cell (with data and next fields) or a binary tree node (with a left and right fields)?
 - Thus we cannot tell where all the pointers are

• But we can be **conservative**:

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer
 - e.g., if it is aligned (what does this mean?)

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer
 - e.g., if it is aligned (what does this mean?)
 - it must point to a valid address in the data segment

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer
 - e.g., if it is aligned (what does this mean?)
 - it must point to a valid address in the data segment
 - All such pointers are followed and we overestimate the reachable objects

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer
 - e.g., if it is aligned (what does this mean?)
 - it must point to a valid address in the data segment
 - All such pointers are followed and we overestimate the reachable objects
- But we still cannot move objects because we cannot update pointers to them

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer
 - e.g., if it is aligned (what does this mean?)
 - it must point to a valid address in the data segment
 - All such pointers are followed and we overestimate the reachable objects
- But we still cannot move objects because we cannot update pointers to them
 - What if what we thought to be a pointer is actually an account number?

- But we can be **conservative**:
 - If a memory word "looks like" a pointer, it is considered to be a pointer
 - e.g., if it is aligned (what does this mean?)
 - it must point to a valid address in the data segment
 - All such pointers are followed and we overestimate the reachable objects
- But we still cannot move opointers to them
 - What if what we thoug number?

Thus we **cannot** use stop-and-copy GC for languages with raw pointer fields that are indistinguishable from data at run time

• Rather that wait for memory to be exhausted, try to collect an object when there are no more pointers to it

- Rather that wait for memory to be exhausted, try to collect an object when there are no more pointers to it
- Store in each object the number of pointers to that object

- Rather that wait for memory to be exhausted, try to collect an object when there are no more pointers to it
- Store in each object the number of pointers to that object
 This is the *reference count*
Reference Counting

- Rather that wait for memory to be exhausted, try to collect an object when there are no more pointers to it
- Store in each object the number of pointers to that object
 This is the *reference count*
- Each assignment operation has to manipulate the reference count

Reference Counting

- Rather that wait for memory to be exhausted, try to collect an object when there are no more pointers to it
- Store in each object the number of pointers to that object
 This is the *reference count*
- Each assignment operation has to manipulate the reference count
 - How expensive is that?

• new returns an object with a reference count of 1

- new returns an object with a reference count of 1
- If x points to an object then let rc(x) refer to the object's reference count

- new returns an object with a reference count of 1
- If x points to an object then let rc(x) refer to the object's reference count
- Code generated for every assignment x <- y must be changed:

- new returns an object with a reference count of 1
- If x points to an object then let rc(x) refer to the object's reference count
- Code generated for every assignment x <- y must be changed:

```
rc(y) <- rc(y) + 1
rc(x) <- rc(x) - 1
if (rc(x) == 0) then mark x as free
x <- y</pre>
```

- new returns an object with a reference count of 1
- If x points to an object then let rc(x) refer to the object's reference count
- Code generated for every assignment x <- y must be changed:

• (That's basically it!)

• Advantages:

- Advantages:
 - Easy to implement

- Advantages:
 - Easy to implement
 - Collects garbage incrementally without large pauses in the execution
 - Why would we care about that?

- Advantages:
 - Easy to implement
 - Collects garbage incrementally without large pauses in the execution
 - Why would we care about that?
- Disadvantages:

- Advantages:
 - Easy to implement
 - Collects garbage incrementally without large pauses in the execution
 - Why would we care about that?
- Disadvantages:
 - Manipulating reference counts at each assignment is very slow

- Advantages:
 - Easy to implement
 - Collects garbage incrementally without large pauses in the execution
 - Why would we care about that?
- Disadvantages:
 - Manipulating reference counts at each assignment is very slow
 - Cannot collect circular structures

• Automatic memory management avoids some serious storage bugs

- Automatic memory management avoids some serious storage bugs
- But it **takes away control** from the programmer

- Automatic memory management avoids some serious storage bugs
- But it **takes away control** from the programmer
 - e.g., layout of data in memory

- Automatic memory management avoids some serious storage bugs
- But it **takes away control** from the programmer
 - e.g., layout of data in memory
 - e.g., when is memory deallocated

- Automatic memory management avoids some serious storage bugs
- But it **takes away control** from the programmer
 - e.g., layout of data in memory
 - e.g., when is memory deallocated
- Most garbage collection implementations stop the execution during collection
 - not acceptable in real-time applications

• **Concurrent**: allow the program to run while the collection is happening

- **Concurrent**: allow the program to run while the collection is happening
- Generational: do not scan long-lived objects at every collection (infant mortality)
 - This approach, in particular, has seen wide adoption in practice (Java)

- **Concurrent**: allow the program to run while the collection is happening
- Generational: do not scan long-lived objects at every collection (infant mortality)
 - This approach, in particular, has seen wide adoption in practice (Java)
- Parallel: several collectors working in parallel

- **Concurrent**: allow the program to run while the collection is happening
- Generational: do not scan long-lived objects at every collection (infant mortality)
 - This approach, in particular, has seen wide adoption in practice (Java)
- Parallel: several collectors working in parallel
- Real-Time / Incremental: no long pauses

- Python uses Reference Counting
 - Because of "extension modules", they deem it too difficult to determine the root set
 - Has a special separate cycle detector
- **Perl** does **Reference Counting** + cycles

- Python uses Reference Counting
 - Because of "extension modules", they deem it too difficult to determine the root set
 - Has a special separate cycle detector
- **Perl** does **Reference Counting** + cycles
- Ruby does Mark and Sweep

- Python uses Reference Counting
 - Because of "extension modules", they deem it too difficult to determine the root set
 - Has a special separate cycle detector
- Perl does Reference Counting + cycles
- Ruby does Mark and Sweep
- OCaml does (generational) Stop and Copy
- Java does (generational) Stop and Copy
- Node.js does (generational) Stop and Copy

• An **automatic memory management** system deallocates objects when they are no longer used and reclaims their storage space.

- An **automatic memory management** system deallocates objects when they are no longer used and reclaims their storage space.
- We must be **conservative** and only free objects that won't be used later.

- An **automatic memory management** system deallocates objects when they are no longer used and reclaims their storage space.
- We must be **conservative** and only free objects that won't be used later.
- Garbage collection scans the heap from a set of roots to find reachable objects. We saw three algorithms:

- An **automatic memory management** system deallocates objects when they are no longer used and reclaims their storage space.
- We must be **conservative** and only free objects that won't be used later.
- Garbage collection scans the heap from a set of roots to find reachable objects. We saw three algorithms:
 - Mark and Sweep uses a per-object mark bit to track which objects are reachable.

- An **automatic memory management** system deallocates objects when they are no longer used and reclaims their storage space.
- We must be **conservative** and only free objects that won't be used later.
- Garbage collection scans the heap from a set of roots to find reachable objects. We saw three algorithms:
 - Mark and Sweep uses a per-object mark bit to track which objects are reachable.
 - Stop and Copy has low overhead but stalls the program and requires language support

- An **automatic memory management** system deallocates objects when they are no longer used and reclaims their storage space.
- We must be **conservative** and only free objects that won't be used later.
- Garbage collection scans the heap from a set of roots to find reachable objects. We saw three algorithms:
 - Mark and Sweep uses a per-object mark bit to track which objects are reachable.
 - Stop and Copy has low overhead but stalls the program and requires language support
 - **Reference Counting** stores the number of pointers to an object with that object and frees it when that count reaches zero
Course Announcements

- As some of you have noticed, the PA4 leaderboard results are subject to timing variance
 - In particular, sometimes you "get lucky" and your compiler's code appears to be much, much faster than the reference
 - There's not much I can do about this, unfortunately.
 - We will run each final submission several times and take the median Q score (not the best), so if you see a big speedup without doing anything, you should assume it is ephemeral
- PA4c1 is due on Monday
- I will not hold office hours next Wednesday