

Project Snowflake: Non-blocking Safe Manual Memory Management in .NET

Matthew Parkinson Dimitrios Vytiniotis Kapil Vaswani
Manuel Costa Pantazis Deligiannis
Microsoft Research

Dylan McDermott Aaron Blankstein Jonathan Balkind
University of Cambridge Princeton University

July 26, 2017

Abstract

Garbage collection greatly improves programmer productivity and ensures memory safety. Manual memory management on the other hand often delivers better performance but is typically unsafe and can lead to system crashes or security vulnerabilities. We propose integrating safe manual memory management with garbage collection in the .NET runtime to get the best of both worlds. In our design, programmers can choose between allocating objects in the garbage collected heap or the manual heap. All existing applications run unmodified, and without any performance degradation, using the garbage collected heap. Our programming model for manual memory management is flexible: although objects in the manual heap can have a single owning pointer, we allow deallocation at any program point and concurrent sharing of these objects amongst all the threads in the program. Experimental results from our .NET CoreCLR implementation on real-world applications show substantial performance gains especially in multithreaded scenarios: up to 3x savings in peak working sets and 2x improvements in runtime.

1 Introduction

The importance of garbage collection (GC) in modern software cannot be overstated. GC greatly improves programmer productivity because it frees programmers from the burden of thinking about object lifetimes and freeing memory. Even more importantly, GC prevents temporal memory safety errors, i.e., uses of memory after it has been freed, which often lead to security breaches.

Modern generational collectors, such as the .NET GC, deliver great throughput through a combination of fast thread-local bump allocation and cheap collection of young objects [63, 18, 61]. At the same time several studies show that GC can introduce performance overheads when compared with manual memory management [41, 44, 69]. These overheads are amplified in big data analytics and real time stream processing applications as recent work shows [57, 36, 56, 49], partly due to the need to trace through large heaps. This trend is likely to continue as modern servers make use of ever larger memories – sizes of hundreds of gigabytes, or even terabytes, are already common.

Manual memory management addresses this problem: it avoids tracing the object graph to free objects and instead allows programmers to exploit their knowledge of object lifetimes to free objects at specific program locations. This improves throughput and also achieves better memory usage due to prompt deallocation. The downside is that manual memory management is typically unsafe and can lead to system crashes or security vulnerabilities, because freeing memory may create dangling pointers, i.e., pointers to memory that has been freed, and dereferences of dangling pointers lead to undefined behaviour. Requiring all memory to be manually managed is also unappealing because it negates the productivity benefits of GC.

In this paper, we show how to get the best of both worlds: combining a GC-based system – in our case the Microsoft open-source .NET runtime [3] – with a facility to manage memory manually, without compromising safety or performance. In our design, programmers can choose between allocating objects in the garbage collected heap or the manual heap. All existing applications run entirely unmodified using the garbage collected heap, and enjoy no performance degradation. Our design places no overhead on garbage collections or other operations like write barriers. Programmers who wish to optimize their applications need to incrementally change their code to allocate some objects from the manual heap, and to explicitly deallocate those objects. We allow allocation and deallocation of individual objects at arbitrary program locations, and we guarantee that manually managed objects enjoy full type- and temporal- safety, including in the presence of concurrent accesses. Programmers get dynamic managed-exceptions for use-after-free scenarios, but no crashes or security vulnerabilities.

Our novel programming model for manual memory management builds on the notion of unique *owners* of manual objects: locations in the stack or on the heap that hold the only reference to an object allocated on the manual heap. Our notion of *owners* is unique, compared to similar concepts in C++, Rust [8], and Cyclone [62]: we allow arbitrary client threads to (a) share stack references to owners (but not to the underlying manual objects), (b) create arbitrary stack references to the actual underlying manual objects from these owner references, and (c) freely abandon the owner reference (which will eventually cause deallocation of the underlying manual objects) – while guaranteeing use-after-free exceptions. To allow safe concurrent sharing of manual objects we introduce the notion of *shields*. Accessing a manual object requires getting a reference from a shield, which creates state in thread local storage that prevents deallocation while the object is being used. Shields can only be created from the unique owning reference, thus when the reference is destroyed no more shields can be created and memory can be safely reclaimed once all previously active shields have been disposed.

We implement this model using a novel combination of ideas drawn from hazard pointer literature [50] and epochs for memory reclamation [13, 39, 34] to provide an efficient lock-free manual memory management scheme, without having to scan large portions of the heap. We develop an epoch-based protocol for determining when it is safe to deallocate an object on the manual heap. The protocol accounts for weak memory model effects, but it is non-blocking. That is, it does not require stopping the world or the use of expensive synchronization. We introduce a mechanism that guarantees *liveness* of the epoch protocol by employing virtual memory protection.

We note that our manual memory management scheme and programming model is independent of the integration of manual memory management with garbage collection and could be applicable in a purely manually managed system too, although it would be more difficult to ensure end-to-end memory safety without an additional strong type system.

Our system is implemented as a fork of the Microsoft open-source .NET implementation. We have modified the .NET runtime (CoreCLR) and extended the standard libraries (CoreFX) with APIs that use manual memory. For manual heap allocations we have integrated jemalloc [6], an industrial size-class-based allocator. Experimental results show substantial performance gains with this design: up to 3x savings in peak working set and 2x improvements in run time. In summary, our contributions are:

- A new flexible programming model for manual memory management where objects can be allocated and deallocated at any program point, and can be concurrently and safely shared amongst multiple threads.
- A set of rules at the C# frontend that ensure the safe use of the programming model.
- An efficient implementation of this programming model that does not require stop-the-world synchronization to safely reclaim manual memory.
- A design for safe interoperability with the garbage collected heap that does not adversely impact the write barriers. To keep the latency of Gen0 collections low, we use existing GC mechanisms to scan only the fragments of the manual heap that contain pointers to Gen0 objects, exactly as if those manual objects were in an older generation on the GC heap.

- An implementation and detailed evaluation on industrial use-cases from machine learning, data analytics, caching middleware, and a set of micro-benchmarks.

2 Background and motivation

Consider the code below, taken from `System.Linq.Manual`, a widely used .NET library for LINQ queries on collections that implement the `IEnumerable<T>` interface.

```
IEnumerable<TR> GroupJoinIterator(IEnumerable<TO> outer, IEnumerable<TI> inner, Func<TO, TKey>
    outerKey,
    Func<TI, TKey> innerKey, Func<TO, IEnumerable<TI>, TRes> res) {
    using (var e = outer.GetEnumerator()) {
        if (e.MoveNext()) {
            var lookup = Lookup.CreateForJoin(inner, innerKey);
            do { TOuter item = e.Current;
                yield return res(item, lookup[outerKey(item)]);
            } while (e.MoveNext());
        }
    }
}
```

The code defines an iterator for the results of a join. We iterate through the outer `IEnumerable<TO>`, `outer`. If the outer enumerable is non-empty then we create a `Lookup<TKey, TI>` structure, which is – in effect – a dictionary that maps keys to groupings of elements from the inner enumerable `inner` that share the same key. We then iterate and apply the `res()` function through every `item` of the outer enumerable. The code uses the C# `yield return` construct and as a result will be compiled to a state machine that can return the results of `res()` one by one in successive calls.

The intermediate data structure `lookup` can potentially grow as large as the size of the inner enumerable and we have to hold-on to it throughout the iteration of the outer enumerable. It cannot be stack-allocated because `yield return` compilation has to save the current state of the iteration and pick it up in the next invocation of the iterator. This `lookup` is an object that is likely then to survive many Gen0 collections (which could happen as a result of allocations inside `e.MoveNext()` or `res()`), and possibly end up in the oldest generation before it can be collected.

It is pretty clear though that once the outer enumerable iteration completes, the internal arrays and data structures of the `lookup` dictionary are entirely safe to deallocate. A generational GC may, however, hold on to these objects until the next full heap collection (such as a Gen2 in the .NET garbage collector) which might happen much later, leading to a blowup of the peak working set; confirmed by our evaluation in Section 5.

Instead we would like to enable programmers to allocate ordinary .NET objects like `lookup` in their manually managed heap, and let them deallocate precisely when they wish. Note that `System.Linq` is a *generic* library that can manipulate not only collections of unboxed data (such as `structs` of primitive types) but also of GC-allocated objects. To maintain this genericity and still be able to allocate internal dictionaries like `lookup` on the manual heap we must allow manual heap objects (like the internal arrays associated with `lookup`) to contain references to GC-allocated objects. Hence a key design goal is to maintain full GC interoperability, allowing pointers to and from the two heaps. This full interoperability is also essential to allow gradual pay-as-you-go migration of applications to use manual memory management for certain objects while others remain in the GC discipline, as performance requirements mandate.

2.1 The challenge of safe manual memory

Deallocation of manually managed objects, while preserving memory safety, is a challenging problem. We seek ways to ensure – statically or dynamically – that an object will not be accessed after it has been deallocated. The challenge is that references to the deallocated object could be remaining on the heap or the stack after deallocation. This might lead to access violations, memory corruption, accessing data that belongs to newer objects allocated in the same virtual address range etc. Scanning the roots and the heap

upon any manual object deletion to discover and zero-out such remaining references can be expensive as each individual scan might have similar cost to a Gen0 collection.

Let us demonstrate the challenge of safety with an example from a multi-threaded caching component of the ASP.NET framework [1], a popular framework for web applications. The cache is defined simply as a dictionary of `CacheEntry` objects (for brevity throughout the paper we are omitting C# access modifiers like `public`, `private` etc.):

```

class MemoryCache {
    Dictionary<object, CacheEntry> _entries;
    bool TryGetValue(object key, out object res);
    void RemoveEntry(CacheEntry entry);
    void Set(object key, object value,
             MemoryCacheOptions options);
}

class CacheEntry {
    // cache entry metadata
    ...
    // actual entry
    Object m_Value;
}

```

Each cache entry object contains metadata about this cache entry (such as when it was last accessed) plus the actual payload value of this object `m_Value`. The cache exposes a method `TryGetValue()` that tries to fetch the payload object associated with a key in the cache. Occasionally the elements of the cache are checked for expiration and `RemoveEntry()` is called on those elements that have expired, which removes those entries from the shared dictionary. Client code can use a cache object `_cache` as follows:

```

if (!_cache.TryGetValue(key, out value)) {
    value = ... ; // allocate a new object
    _cache.Set(key, value, _options); // put it in the cache
}
// perform a computation with value

```

The cache entry payload objects, accessed from the `m_Value` field above, are objects that survive into the older generations and may be collected much later than their removal from the dictionary. They also have a very clear lifetime that ends with a call to `RemoveEntry()`. They are, hence, excellent candidates for allocation on the manual heap. But how to deallocate those objects safely, when multiple threads are accessing the cache and one of them decides to remove an entry?

2.2 Key insight: owner references and shields

In the caching example, the lifetime of the cache payload objects and access to them is controlled by the associated `CacheEntry`. The field `m_Value` acts as the only *owner* of those objects. These objects are only temporarily accessed by stack references in client code (while remaining in the cache) but all of those stack references have been first-obtained through the pointer `m_Value`. Consequently, if we are to zero-out `m_Value` in `RemoveEntry()` no *future* references to the payload object can be obtained.

But when can we actually deallocate and reclaim the memory of the underlying payload? In this example, client threads may have *already* obtained stack references to the object in question before `RemoveEntry()` has been called, and they may still be working on those objects after the corresponding cache entries have expired and have been removed from the dictionary. We cannot deallocate these objects while other code is accessing them.

Our solution to this problem is inspired by hazard-pointers [50], a technique originating in the lock-free data structure literature. We introduce a mechanism to publish in thread-local state (TLS) the intention of a thread to access a manual object through one of these owner locations. This registration can be thought of as creating a *shield* that protects the object against deallocation and grants permission to the thread that issued the registration to directly access the manual object e.g. call methods on it or mutate its fields. At the same time no thread (the same or another) is allowed to deallocate the object and reclaim its memory. Once client code no longer needs to access this object, it can dispose the shield, that is remove the reference to this object from its TLS. It is not safe to directly access the object that has been obtained from a shield, after the shield has been disposed because, following this disposal of the shield, the actual deallocation is allowed to proceed (if some thread has asked for it, and provided that this reference does not exist in *any* TLS in the system). If the owner link has been zeroed-out in the meanwhile no new references can be obtained.

```

struct Owner<T> where T : class {
    Shield<T> Defend();
    void Move<S>(ref Owner<S> x)
        where S:class, T;
    void Abandon();
}

class ManualHeap {
    void Create<T>(ref Owner<T> dst) where T:class, new();
    void CreateArray<S>(ref Owner<S[]> dst, int len);
}

struct Shield<T> : IDisposable
    where T:class {
    static Shield<T> Create();
    void Defend(ref Owner<T> u);
    T Value;
    void Dispose();
}

```

Figure 1: Core Snowflake API

This is the key mechanism that we use to ensure manual memory safety. Next, we formally present our programming model and describe the rules that programmers must abide by to preserve memory safety.

3 Programming model

The Snowflake programming model is designed for .NET and we present it here in C# syntax. First, we remind the readers of some .NET features we depend on.

3.1 Preliminaries

C#/.NET introduce constructs that give programmers some control over memory layout. In our programming model we will use **struct** types, which – contrary to classes – can be allocated on the stack or directly inside another struct or class. Struct assignment amounts to memory copying and **struct** arguments are passed by value. In addition, C#/.NET allow to pass arguments (of **struct** or **class** types) by *reference* by explicitly using the **ref** keyword. The address of the corresponding struct will then be passed instead of a copy of the struct, and in the case of classes, the address where the object pointer lives instead of the object pointer.

3.2 Snowflake API

Figure 1 gives the public Snowflake API. To avoid any confusion we emphasize here that – by itself – the API does not guarantee safety. Safe use of the API relies on a set of language frontend checks, described in more detail in Section 3.4.

Owners An `Owner<T>` encapsulates a (private, unmanaged) pointer to a manual object. The runtime implementation of our API relies for safety on the *unique owner condition*: No two `Owner<T>` structs should ever be allowed to refer to the same manual object. This condition is enforced by our C# language frontend (see Section 3.4). Note that `Owner<T>` is defined as a **struct** as opposed to a **class** to avoid an extra GC object allocation per manual object. For reasons explained in Section 3.4 we are only allowed to pass such structs as arguments to functions *by reference*.

Struct `Owner<T>` exposes three methods. The first `Defend()`, returns a `Shield<T>` and prevents deallocation of the manual object associated with this owner (by publishing this manual object pointer in thread-local state.) The second `Abandon()`, zeroes out the pointer to the manual object, so that no new `Shield<T>` can be obtained, and schedules the manual object for deallocation at some safe point in the future, when it is no longer protected by any shield in any thread. The final method `Move(ref Owner<S> x)`, corresponds to transferring ownership from `x` to the receiver struct. The pointer inside the `x` struct is moved to the receiver struct and the `x` struct is zeroed out. If the receiver struct was holding a manual object pointer prior to the call to `Move(ref x)` then that manual object will be scheduled for deallocation at some later safe point, since – by the unique owner condition – the receiver struct was the only owner of that object.

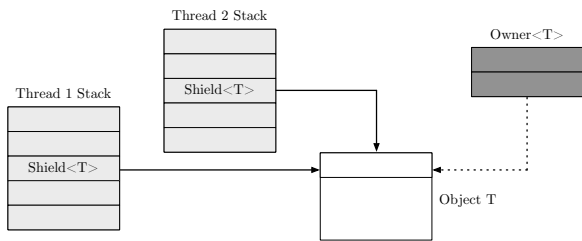


Figure 2: Owners and shields.

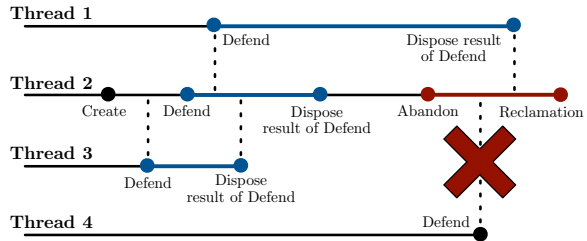


Figure 3: Example of lifetimes of owners and shields. `Defend` refers to calling `Owner.Defend()` to get a shield; `Dispose` refers to disposing that shield.

Shields A `Shield<T>` acts as a stack-only access token to the underlying object. It can be obtained from the `Defend()` method of an `Owner<T>` and encapsulates a reference to thread-local state that records the underlying manual object as one whose memory cannot be reclaimed. It exposes the following members: `Value`, is a *property* that gives access to the underlying manual object; and `Dispose()` un-registers the manual object that this shield protects from thread-local state, making it thus a candidate for deallocation.

The lifetime of a shield is not tied to a specific access of a specific owner. Shields are only references to slots in thread-local state and can be created in *uninitialized form*, and be used to defend multiple objects. For this reason `Shield<T>` exposes two more methods: `Create()` which simply creates a new uninitialized shield that does not yet defend any object against deallocation; and `Defend(ref Owner<T> u)` which defends a *new* owner, and *un-defends* the owner it previously defended, if any. This is done by overwriting the TLS slot that corresponds to this shield with the new object pointer. This method is handy for avoiding frequent creation and disposal of shields on every iteration of a loop that accesses some manual objects. We can instead create a shield before the loop (allocate a TLS slot) and dispose it in the end (deallocate a TLS slot), but continuously re-use it (overwrite the pointer in that slot) to defend each item in each iteration.

Allocating on the manual heap Our API exposes `Create()` and `CreateArray()` methods in Figure 1 for allocating objects and arrays. These methods allocate in the manual heap and transfer ownership of the newly allocated object to the destination owner. In our C# frontend we use syntactic sugar `new Owner<MyClass>(...)` for allocating in the manual heap *and* calling a constructor. Here we just focus on the low-level .NET APIs.

3.3 Putting it all together: lifetimes

Figure 2 demonstrates how the stack and the heap may look during an execution of code that uses our programming model. An owner can live on the stack, the GC heap or manual heap and contains a link to the underlying manual object. Multiple threads may have used this owner to obtain stack-local shields, which internally refer to the manual object directly.

Figure 3 on the other hand describes owner and shield lifetimes in a left to right chronological order. Thread 2 creates an owner object and subsequently creates a shield through `Defend()` that allows it to access the underlying object up to the point of disposal of that shield. At the same time Thread 3 can also access the same object by obtaining another shield. At some point later on, `Abandon()` is called from Thread 2, and the object is scheduled for deallocation. Thread 1 has in the meanwhile *also* defended and obtained a shield, so is keeping the object from being deallocated. Deallocation can proceed once no thread holds an (undisposed) shield to the underlying object. Finally, Thread 4, cannot obtain a shield as it is issuing the `Defend()` call after Thread 2 has abandoned the owner.

```

class Foo {
    Owner<Object> mobj;

    void meth() {
        Owner<Object> mobj1;
        mobj1 = mobj;
        // violated unique owner condition
        mobj.Abandon();
        var sh = mobj1.Defend();
        sh.Value.ToString(); // unsafe!
    }
}

var sh = mobj.Defend();
var sh1 = sh;
// referring to the same TLS slot!
... use sh.Value here ...
sh.Dispose();
// Object can go away now
... use sh1.Value here ... // unsafe!

```

Figure 4: Unsafe owner (left) and shield (right) examples

3.4 Language frontend and safety

As mentioned in the beginning of the section, extra checks are needed at the C# frontend level to ensure safe use of the Snowflake .NET API. Although the emphasis of this paper is on the efficient lock-free runtime, in this section we describe those checks for completeness:

Ensuring the unique owner condition If `Owner<T>` structs are allowed to be duplicated then we may be holding on to references into deallocated manual objects, because we have violated the unique owner condition. As an example, consider the unsafe code fragment in the left part of Figure 4. Similar dangers exist anywhere `Owner<T>` could be subject to copying. Concretely our frontend implements the following rules:

- No assignment or cloning of a heap location or stack variable of type `Owner<T>` is allowed.
- No passing or returning `Owner<T>` by value is allowed. If we were to pass `Owner<T>` by value then the callee and caller would effectively be holding two copies of the struct, compromising safety.
- No instantiation of generic methods with `Owner<T>` types is allowed. Generic code can involve arbitrary assignments where the type of the object we are assigning is a generic parameter, or function calls that pass arguments of generic parameter types, and hence can lead to potential duplication of `Owner<T>` structs if those generic parameters were to be instantiated to `Owner<T>`.

Effectively our frontend guarantees that `Owner<T>` is a *non-copyable* value type.

Ensuring shields are unique Similarly to `Owner<T>`, `Shield<T>` should be a non-copyable value type. Consider the unsafe code fragment in the right part of Figure 4. By duplicating `sh` into `sh1` we kept a reference to the TLS slot that guarded our manual object past the point where that TLS slot stopped guarding the object, resulting in an unsafe access. For this reason, our C# language frontend enforces the same non-copyable value type conditions on `Shield<T>`, with a small improvement: a simple static analysis does allow a method to return by value a *fresh* `Shield<T>`. A fresh `Shield<T>` is one that is created within the scope of a method and is never stored or copied, but is allowed to be returned to the caller (and then freshness propagates). This improvement is useful to allow the line: `var sh = mobj.Defend();` above. That line – superficially – looks like a copy, but is actually innocuous. If we were not to have this improvement, our API would have to provide a `Defend()` method that accepted a `ref Shield<T>` argument.

Ensuring shields are stack-only Shields cannot be stored on the shared heap. The reason is that a shield has a meaning only for the thread that created it – it holds a pointer to a TLS slot. By placing shields on the shared heap, we generate possibilities for races on state that was supposed to be thread-local.

```

internal class CacheEntry : ICacheEntry {
    ...
    Owner<object> m_Value;
}

class MemoryCache : IMemoryCache {
    ...
    object Set(object key, ref Owner<object> value, MemoryCacheEntryOptions options) {...}
    bool TryGetValueShield(object key, ref Shield<object> result) {...}
    bool RemoveEntry(CacheEntry entry) {
        ...
        entry.m_Value.Abandon();
    } }

```

Figure 5: Memory cache API modification for Snowflake

Ensuring manual objects don’t escape past their shields lifetimes Consider the unsafe fragment below:

```

var sh = owner.Defend();
Object mobj = sh.Value;
... use mobj here ... // safe
foo.f = mobj; // unsafe (heap escape)
sh.Dispose();
... use mobj here ... // unsafe (use after shield dispose)

```

Here, we’ve created a shield from `owner`, and we are allowed to access the underlying manual object `mobj`. However, we must ensure that the object does not escape onto the heap because this can result in an access after the local shield has been disposed of and the object has been deallocated. We must also ensure that we do not access the object locally after the shield has been disposed.

Our frontend enforces those restrictions with a conservative dataflow analysis that has proven sufficient for our examples. Whereas the escape past the lifetime of the shield is easier and more local to detect, the heap escape analysis can become more involved as we need to descend through potentially deep hierarchies of method calls. For this reason we also have experimented with an alternative implementation that enforces this restriction uniformly by throwing exceptions in the write barriers if the object we are storing belongs in the virtual address range of the manual heap. We have measured the impact to performance of this check to be negligible because this path is extremely rare and does not involve another memory access.

3.5 Examples of Snowflake in action

We present here some examples of how Snowflake can be used to offload objects to the manual heap, safely.

Lists with manual spines The C# collection library defines a `List<T>` collection for lists of `T` elements. Internally it uses an *array* `T[]` to store those elements, which gets appropriately resized when more space is needed. This array can grow quite large and persist many collections, but is completely internal to the `List<T>` object, and hence it is an ideal candidate for moving it to the manually managed heap. Here is the original (left) and the modified (right) code:

```

class List<T>{ T[] _items; ...} | class List<T>{ Owner<T[]> _items; ...}

```

Multiple methods of `List<T>` use `_items` and each requires modification. Here is the original `Find` method (slightly simplified for space) that finds the first element that matches a predicate `match` or returns a default value. To port this to our new definition for `_items` we have to obtain a shield on the owner struct.


```

T Find(Predicate<T> match) {
    for (int i = 0; i < _size; i++) {
        if (match(_items[i]))
            return _items[i];
    }
    return default(T);
}

T Find(Predicate<T> match) {
    using (Shield<T[]> s_items =
        _items.Defend()){
        for (int i = 0; i < _size; i++) {
            if (match(s_items.Value[i]))
                return s_items.Value[i];
        }
    }
    return default(T);
}

```

The `using` construct automatically inserts the `Dispose()` method on a `IDisposable` object (such as our `Shield<T>`) in ordinary and exceptional return paths and we use it as convenient syntactic sugar. Inside the `using` scope we can access `s_items.Value`, but it must not escape on the heap nor be used past the `using` scope.

Note though, that since the new `List<T>` has a manual spine, when it is no longer needed programmers have to explicitly deallocate it using the following method:

```
void Abandon() { _items.Abandon(); }
```

This is the only new method we added to the `List<T>` API. Finally note that `List<T>` itself can be allocated on the GC heap or on the manual heap. We stress that having owner structs inside GC objects is not problematic for memory safety. Whereas the GC object is shared by many threads, our API ensures there is still a unique owner into the manual spine of the list, the one that is stored inside that GC object. In addition to this distinguished pointer, plus possibly many shield-registered stack references to the manual spine can also exist.

Moving ownership in lists with manual spines Occasionally the internal manually allocated array of a list must be resized to increase its capacity. Here is how we can do that:

```

var new_items = new Owner<T[]>;
ManualHeap.CreateArray(ref new_items, new_size);
using (var s_items = _items.Defend(),
    s_new_items = new_items.Defend()) {
    Array.Copy(s_items.Value, 0,
        s_new_items.Value, 0, _size);
}
_items.Move(ref new_items);

```

We first allocate `new_items`, a new `Owner<T[]>` using our `ManualHeap.CreateArray<T>()` method. Once that is done, we obtain shields to both the old and new items, and copy over the contents of the old items to the new array. Finally, we transfer ownership of the `new_items` to `_items`, which schedules the original manual object for deallocation.

Collections of owner objects `List<T>` is an example where the spine of a data structure can be moved over to the manual heap, but ASP.Net caching actually does store long lived data. For this reason we may port the `m_Value` field to be an owner of a manually allocated object. For convenience we will keep the dictionary spine on the GC heap. We first need to change the `CacheEntry` and the interface to the memory cache, as shown in Figure 5.

Method `Set()` now accepts a reference to a locally created owner value, and will `Move()` it in to the cache if an entry is found with the same key or will create a fresh entry for it. Method `TryGetValueShield()` is passed in a shield reference, and uses `it` to protect an object (if found) against deallocation. The client code can then just access the object through that shield, if `TryGetValueShield()` returns `true`. Finally `RemoveEntry()` abandons the owner, scheduling it for deallocation once it is no longer shielded. The client code can be as follows:

```

using (var res_sh = Shield<object>.Create()) {
    if (!_cache.TryGetValueShield(key, ref res_sh)) {
        Owner<byte[]> tmp;
        ManualHeap.CreateArray(ref tmp, size);
    }
}

```

```

    res_sh.Defend(ref tmp);
    ... // populate tmp, through res_sh
    _cache.Set(key, ref tmp, _options);
}
... // use res_sh to complete the request
}

```

First we create a shield that will protect manual objects against deallocation throughout the request. We create it uninitialized to start with, so it does not protect any object. Subsequently we try to make this new shield defend the object that this key maps to in the cache, if such a cache entry exists. If we do not find it in the cache, we create a new local owner and allocate something, and use the shield to protect it and finally exchange it into the cache. In the rest of the code we can access the manual object (coming from the cache or freshly allocated) through `res_sh`.

A set of simple changes in an application can offload many objects that survive into older generations to the manual heap and result in significant speedups and peak working set savings as we will see in Section 5.

4 Implementation

Next we describe the key parts of our implementation. We have extended CoreCLR with the shield runtime, integrated a modified version of jemalloc to manage the physical memory for the manual objects, and added interoperability with the GC to allow the manual heap to be scanned for roots.

4.1 Shield runtime

Our approach to implementing `Shield<T>` combines ideas from both hazard pointers [50] and epoch-based reclamation [39, 34, 13]. We provide a comparison in the related work, and just explain our implementation here.

The core responsibility of the shield runtime is to safely enable access to (and deletion of) manually managed objects, without requiring expensive synchronisation on accessing an object (a “read barrier”). To achieve this we allow reclaiming manual objects to be delayed.

Each thread holds a thread-local array of slots, each of which protects a manual object (or `0x01` for an unused slot). A `Shield<T>` struct then holds a pointer (`IntPtr` in C#) capturing the address of the thread-local slot (`slot` field) that this shield is referring to. The TLS slot stores the address of the manual object. The same address – for efficient access that avoids TLS indirections – is also cached as a field (`value`) of type `T` inside `Shield<T>` and is what the `value` property returns. Allocation of shields amounts to finding an unused TLS array slot and creating a shield that holds a pointer to that yet uninitialized slot. A call to `v.Defend(ref x.o)` is rewritten to `v.value = SetShield(ref v.slot, ref x.o)` where `SetShield()` simply sets the TLS slot to contain the address of the manual object and returns it.

For abandoning a manual object, we effectively do

```

void Abandon(ref Owner<T> o) {
    var x = InterlockedExchange(o, null);
    if(x != null) AddToDeleteList(x);
}

```

where `AddToDeleteList` adds the object to a thread local list of objects that need to be reclaimed. The call to `InterlockedExchange` is a compiler intrinsic that sets a location to the specified value as an atomic operation and returns the old value.¹

Occasionally, when the objects in the list consume too much space we trigger `Empty()` to reclaim memory:

```

void Empty() {
    for (var i in DeleteList)
        if(IsNotShielded(i)) free(i.ptr);
}

```

¹We use the Microsoft VC++ intrinsics – in gcc this would be `__sync_lock_test_and_set (o, null)`.

`IsNotShielded()` needs to check that *every thread's* TLS shield array does *not contain* the object. However, during that check some thread may be calling `SetShield`. Correctness is predicated on the correct interaction between `IsNotShielded` and `SetShield`.

Naïve approach to synchronisation One way to solve the synchronization problem is to exploit the “stop-the-world” GC synchronization for collecting the roots. GC suspends all threads (which also causes all threads to flush their write buffers so that all memory writes are globally visible). At that point we can iterate through every thread and call `Empty()` on the thread-local `DeleteList`.

To make `Empty()` efficient, we first iterate through all threads and build a Bloom filter [19] that contains all shields of all threads. We use that Bloom filter to over-approximate the `IsNotShielded` check. This avoids the need to repeatedly check with the local shields of each thread.

However, there is a danger if we allow mutator threads to be suspended *during* a `SetShield()` call. In C notation, `SetShield()` executes the following:

```
Object* SetShield(Object** slot, Object** fld) { *slot = *fld; return *slot; }
```

If a thread is suspended for GC right after reading `*fld`; but *before* publishing the object in the shield `slot` and another thread has already issued an `Abandon()` on the same manual object we can get into trouble: a call to `Empty()` from the GC thread will not see the object in the shield TLS array of the first thread and may deallocate it. When threads resume execution, the first thread will publish the address of a (now deallocated) object and continue as normal, leading to safety violation. For this reason we prevent a thread that executes `SetShield()` from being suspended. This is supported in CoreCLR by executing `SetShield()` in “cooperative mode”. This mode sets a flag that forces the runtime to back out of suspending the thread.

We have solved the problem by effectively forbidding calls to `IsNotShielded()` from occurring during `SetShield()`. However, calling `Empty()` this way comes at the cost of suspending the runtime. We want to maintain eager deallocation of manual objects, and avoid suspending the runtime too often, so we only use this approach if a GC is going to occur. Next, we develop another mechanism that does not require suspending any threads, and allows `Empty()` to be called independently by any mutator.

Epochs for concurrent reclamation To enable more concurrency for manual object collection, we use epochs to synchronise the threads’ view of shields, drawing from work on reference counted garbage collectors [13]. Epochs allow thread-local collections without stopping any threads or requiring expensive barriers in the `Defend` code.

We use a 64-bit integer for a global epoch. Each thread has a local epoch that tracks the global. Adding an object to the `DeleteList` may trigger advancement of both this thread’s local and/or the global epoch. If the running thread detects that it is lagging behind the global epoch, it just sets its local epoch to be equal to the global epoch. If the running thread detects that all threads agree on their local epochs – and epoch advancement is heuristically sensible – then it performs a CAS that sets the global epoch to be the agreed upon epoch plus one. It is okay for the CAS to fail, as that means another thread has advanced the global epoch.

This protocol guarantees that the global epoch is *never* more than one ahead of any local epoch. In Figure 6 we illustrate the ordering of such epoch events (arrows denote causal happens-before relations). $L(n)$ signifies a thread writing to its local epoch the value n ; and $G(n)$ signifies the global epoch advancing to n . To advance the global epoch from n to $n + 1$, every thread’s local epoch must be in epoch n .

When we add an object to the `DeleteList` we record which local epoch it was deleted in. Occasionally, mutators will call `Empty()`, and create the Bloom filter of all threads’ shields, without stopping any thread. Unfortunately not all writes to the shield TLS may have hit the main memory so there is a possibility that an object will be residing in the `DeleteList` with its corresponding `SetShield()` write effects still in flight; we cannot actually deallocate those objects. However, all the objects in the `DeleteList` whose recorded epoch count is less than three epochs behind the local epoch, if they are not in the Bloom filter, are actually safe to deallocate.

We illustrate why three epochs is enough in Figure 6. First we define 3 types of *logical event*: *Abandon events*, $A(o)$ denote the exchange of the owner with *null* and scheduling a pending deallocation; *Defend*

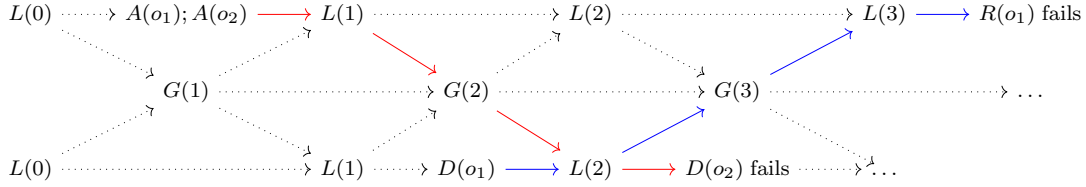


Figure 6: Happens-before ordering of global (G), local (L) epoch events, abandon (A), defend (D), and reclaim (R) operations. Thread 1 in the top line, Thread 2 execution in the bottom line.

events, $D(o)$ denote the combined effects of `setShield()` – that is reading of the object and storing it in the TLS shield array; and *Reclaim events*, $R(o)$ denote successful object deallocation. Assume that in epoch 0, Thread 1 abandons both o_1 and o_2 ($A(o_1); A(o_2)$). We have several possibilities for Thread 2:

- Assume a defend event $D(o_1)$ in epoch 1. The effects of this event (in particular the write to the TLS slot) are guaranteed to be visible by epoch 3 in Thread 1, as the blue solid arrows indicate. This ordering on x86 comes as part of the memory model; for ARM and Power barriers are required in the global and local epoch updates.
- Assume Thread 2 attempts to defend an object o_2 in a later epoch, epoch 2 ($D(o_2)$). However, as the red solid arrows indicate, the abandon event of o_2 ($A(o_2)$) must have become visible by now and hence we will be defending `null` – no violation of safety.
- If Thread 2 defends one of o_1 and o_2 earlier (e.g. in $L(0)$) then trivially less than three epochs suffice (not depicted in the diagram).

We conclude that in all cases we need not wait more than three epochs to safely delete an object if it is not contained in some TLS shield array. Effectively, and taking into account that `isNotShielded()` may also be called with the runtime suspended, our `isNotShielded` code becomes:

```
bool isNotShielded (Node i) {
    return (runtimeIsSuspended || (i.epoch + 3 <= local_epoch) )
    && ... ; //check per thread shields
}
```

To efficiently represent the epoch an object was abandoned in, we use a cyclic ring buffer segmented into four partitions: three for the most recent epochs and one for the spare capacity.

Finally notice that our reasoning relies on $A(\cdot)$ and $D(\cdot)$ events being atomic with respect to local epoch advancement. This holds as each thread is responsible for advancing its own local epoch.

Protocol ejection for liveness A thread being responsible for advancing its own epoch can lead to liveness problems. If a thread blocks, goes into unmanaged code or goes into a tight computation loop, it can hold up the deallocation of objects. To solve this problem we introduce an additional mechanism to *eject* threads from the epoch protocol. If Thread A has too many objects scheduled for deallocation, and Thread B is holding up the global epoch, then Thread A will *eject* Thread B from the epoch consensus protocol and ignore its local epoch; Thread B must rejoin when it next attempts to use `setShield()`.

Each thread has a lock that guards the ejection and rejoining process. When Thread A ejects Thread B, it makes the TLS shield array of Thread B read-only using memory protection (`VirtualProtect` in Windows, `mprotect` in Linux). It then marks Thread B’s epoch as a special value `EJECTED` which allows the global epoch advancing check to ignore Thread B. As multiple threads may be simultaneously trying to eject Thread B, the ejection lock guarantees that only one will succeed (`TryAcquireEjectionLock()`).

```

void Eject(Thread *other) {
    if (other->TryAcquireEjectionLock()) {
        VirtualProtect(other->Shields,
                       READONLY);
        other->local_epoch = EJECTED;
        other->ReleaseEjectionLock();
    }
}

void Thread::Rejoin() {
    this->AcquireEjectionLock();
    VirtualProtect(this->Shields,
                  READWRITE);
    this->local_epoch = global_epoch;
    this->ReleaseEjectionLock();
}

```

Note that we *must* rejoin the protocol if we are to use shields, hence we must wait to acquire the ejection lock (`AcquireEjectionLock()`). We can then un-protect the TLS pages that hold the shield array, and set our local epoch back to a valid value to rejoin the consensus protocol.

```

void AddToDeleteList(Object *o) {
    Epoch curr = local_epoch;
    if (curr == EJECTED) curr = global_epoch;
    DeleteList->push(o, curr);
    ... // possibly call Empty()
}

```

To handle the ejection mechanism, we must adapt `Abandon()` slightly. Recall that `Abandon()` first exchanges `null` in the owner, and then calls `AddToDeleteList()`, with the current local epoch. However, due to ejection, the local epoch may be `EJECTED`. So if the thread is ejected, we use the global epoch. Note that for the argument in Figure 6 to be correct, it actually suffices that the epoch used

to insert the object in the `DeleteList` be at least $(g - 1)$ where g was the global epoch at the point of the atomic exchange of `null` in the owner. If ejection happened between the exchange and `AddToDeleteList` then the new global epoch that we will read is guaranteed to be at least g .

The TLS shield array of an ejected thread will be read-only, hence we guarantee that any call to `SetShield()` from an ejected thread will receive an access violation (AV). We trap this AV, rejoin the protocol with `Rejoin()`, and then replay the `SetShield()` code. By replaying the `SetShield()` code after `Rejoin()` we make it atomic with respect to ejection, and thus local epoch advancement, as required earlier. You can view memory protection as an asynchronous interrupt that is only triggered if the thread resumes using shields.

4.2 GC interoperability and jemalloc

The core changes to the GC and jemalloc are: (1) provide a cheap test if an object is in the manual or the GC heap; (2) extend the GC card table [45] to cover the manual heap; and (3) allow iteration of GC roots from the manual heap.

We modify the OS virtual memory allocation calls both in the GC and jemalloc to use a threshold: 2^{46} . The GC allocates in pages directly above this, and jemalloc directly below. This allows for a cheap test to determine in which heap an object resides without any memory access. As CoreCLR has a generational GC, we need a card table to track pointers from the manual heap into Gen0 and Gen1. By growing both heaps away from a threshold, the used virtual address space is contiguous among the two heaps. Hence, we use a single card table to cover both heaps and we do not modify the write barriers. When jemalloc requests pages from the OS we notify the GC so that it can grow the card table. This requires delicate synchronization to avoid deadlocks.

Finally, we have implemented a trie that covers the manual heap and contains a bit for each 64bit word to represent if it is the start of an object that contains GC heap references. We iterate this trie when we need to find roots in the manual heap during a garbage collection. This trie is used in conjunction with the card table for Gen0 and Gen1 collections, and iterates all objects for a full Gen2 collection.

5 Evaluation

We performed experiments to validate the the impact of our mixed-mode memory management approach to application runtime and peak working sets (PWS), as well as scalability on threads and heap sizes. To measure the PWS of a process we directly query the OS. Specifically, we query the “Working Set Peak”

Config	Mode	%GC	#Gen0	Mean	Max	#Gen1	Mean	Max	#Gen2	Mean	Max
256-byte	GC	26.5%	530	5.2	15.3	197	16.2	31.5	21	27.5	181.8
256-byte	M+GC	18.1%	349	5.4	9.9	185	11.3	25.9	17	17.6	133.6
2048-byte	GC	32.1%	1333	2.0	6.8	698	8.5	74.4	37	40.8	321.6
2048-byte	M+GC	17.6%	370	5.4	14.8	201	10.9	26.1	20	21.6	143.3
4096-byte	GC	34.5%	2481	1.4	7.4	1387	7.0	67.7	46	43.5	409.8
4096-byte	M+GC	14.9%	397	5.3	15.0	199	11.4	27.2	27	17.1	132.8

Table 1: GC Collections and pauses in ASP.NET Caching benchmark (8 threads). Mean and Max times are

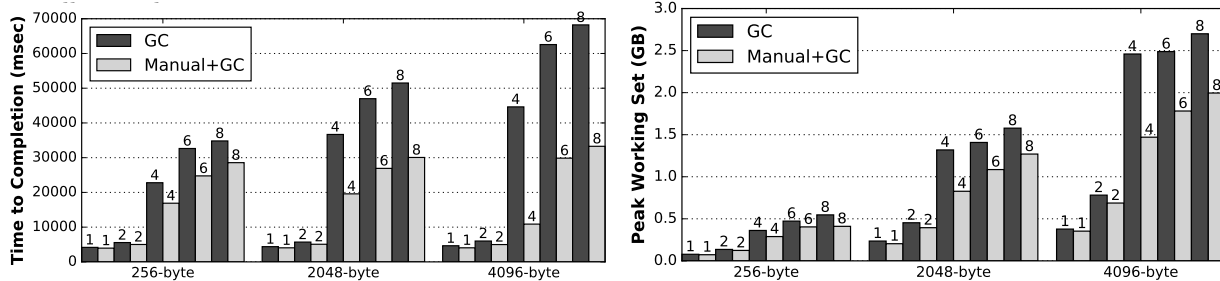


Figure 7: Results for ASP.NET Caching.

Windows Performance Counter for the process that we are testing, which includes both GC heap memory, manual heap memory, runtime metadata etc.

We have ported three industrial benchmarks from data analytics (`System.Linq.Manual` used to implement TPC-H queries), caching (ASP.NET Caching), and machine learning (probabilistic automata operations on top of the Infer.NET framework [51]). We also present a small set of data structure-specific micro-benchmarks on trees and graphs. The benchmarks were written originally in unmodified .NET but we ported some of the allocations to the manual heap based on profiling. For some benchmarks this was as easy as writing a few lines of code (e.g. ASP.NET caching, `System.Linq.Manual`), for some others it was more tedious due to the “by ref” style of owners and shields.

Our results generally demonstrate (i) better scalability of both runtime and PWS than the purely-GC versions, (ii) savings in PWS, though (iii) throughput results can be mixed, depending on the benchmark.

Experimental Setup We performed all experiments on a 3.5GHz Intel Xeon CPU E5-1620 v3 (8 physical cores) with 16GB RAM running Windows 10 Enterprise (64-bit). We used our own branch of CoreCLR, CoreFX and jemalloc. We configured the CoreCLR in “Workstation GC” mode. Note that the CoreCLR GC does not give programmers control over generation/heap size so we do not experiment with space-time tradeoffs in garbage collection. However, CoreCLR additionally provides a “Server GC” mode, which allows for independent thread-local collections. We have only run preliminary results against this configuration and the relative trends are similar to our findings with Workstation GC, but the absolute numbers show that Server GC trades higher PWS for lower times to completion. We refer to the supplementary material for this data.

ASP.NET Caching In Section 3.5 we have presented a modification to a caching middleware component from ASP.NET. We have created a benchmark where a number of threads perform 4 million requests on a *shared* cache. Each request generates a key uniformly at random from a large key space and attempts to get the associated value from the cache. If the entry does not exist then a new one is allocated. We show here experiments with small (256 bytes), medium (2K), or large (4K) payloads and a sliding expiration of 1sec.

Figures 7 show the results. Thread configurations are labeled with a number on each bar. Both time to completion and PWS savings compared to the purely GC version improve substantially, particularly with bulkier allocations and multiple threads. Characteristically, time to completion halves with 8 threads and 4k cache payloads, whereas peak working sets improve up to approximately 25%.

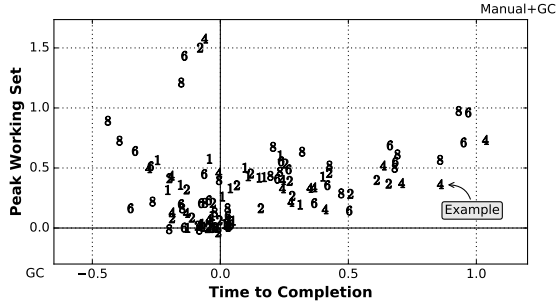


Figure 8: PWS and runtime for TPCB queries.

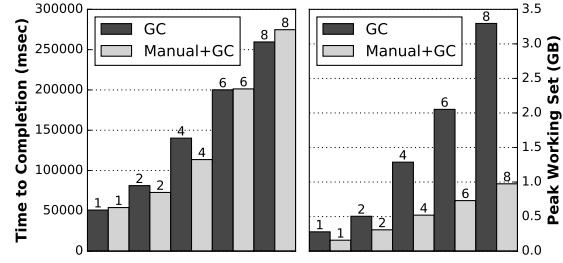


Figure 9: PWS and runtime for the Infer.NET benchmark.

The key reason for the massive throughput speedups for this benchmark is that by off-loading the cache payloads to the manual heap, a lot of write barriers are eliminated entirely, and at the same time fewer GC collections are triggered (Table 1). For the 2048 and 4096 payload sizes, there are approximately half the Gen2 collections and $\frac{1}{6}$ th of the Gen0 and Gen1 collections. There are still quite a few garbage collections remaining with the manual memory modifications, basically because the cache stores GC-based `CacheEntry` objects (which *in turn* store the actual manual payloads), and it is the allocation of those objects that triggers collections. Note also that the cost of collections for Gen0 and Gen1 is higher with manual memory, Gen0 goes from 1.4ms to 5ms with manual memory. This is possibly because our modifications have taken the bulky payload objects out of the ephemeral (Gen0 and Gen1) segments, and hence the ephemeral segments are now full with (many more) small `CacheEntry` objects. Ephemeral collections simply have to collect more objects and take more time. However, since in the purely GC version Gen0 collections collect fewer objects, more pressure is put on Gen2 collections. We can see this as the purely GC code takes 43.5ms for a Gen2 collection, where as with manual heap it is just 17.1ms (for the 4096-byte configuration).

TPCH on `System.Linq.Manual` Here we introduced a drop-in replacement for `System.Linq` (namespace `System.Linq.Manual`), which is a popular C# library for data processing in SQL-like syntax. It is designed to allocate little, but some operators must materialize collections, namely *group*, *join*, and *order-by* variants, as the example of a join iterator from Section 2 explained. We have moved these transient collections to the manual heap, and introduced deallocations when they were no longer needed. For evaluation we used 22 industry-standard TPCB queries [10], ported them in LINQ C#, and run them using `System.Linq` and `System.Linq.Manual`. We have scaled the TPCB dataset to a version where 300MB of tables pre-loaded in memory, and another where 3GB of tables is pre-loaded in memory. We present here the results from the small dataset (the results are similar for the large dataset in the supplementary material). To evaluate scalability in multiple threads we run each query in 10 iterations using 1,2,4,6,8 *independent* threads sharing the same pre-loaded dataset.

For reasons of space Figure 8 presents all 22x5 results as a scatter plot where the vertical (resp. horizontal) axis shows *relative* improvement in the peak working sets (resp. runtime). The labeling indicates number of threads. Positive numbers in both axes mean that we are faster and consume less memory. For example, the number 4 labeled “Example” in the figure indicates that on one query with 4 threads, the purely GC code used $\sim 80\%$ more time, and $\sim 40\%$ more memory. The majority of queries improve in PWS, some up to 1.5, meaning that the purely GC version uses 150% more memory than our modified version. For runtime the results are mixed, though a significant number of queries achieves significantly better runtime. Profiling showed that this is primarily due to the cost of TLS accesses during the creation and defending with shields. GC statistics that can be found in the supplementary material reveal that we generally reduce the number and pauses of collections. For some queries though Gen2 pauses become higher – this may have to do with the more expensive scanning of the trie for GC roots, especially for *background* collections. We discuss this issue further in Section 6.

Machine learning on Infer.NET We used a workload from a research project with known GC overheads based on the Infer.NET framework. The task was to compute the product of two large probabilistic automata for a small number of iterations. We converted various collections of automata states to use manual “spines” but kept the actual automata states in the GC heap to avoid too intrusive re-engineering. To avoid TLS access during iteration of these collections, a significant fragment of the product construction code was rewritten to re-use pre-allocated shields that were passed down the stack as by-ref parameters. This was a more invasive change, than those required in previous benchmarks. As with the TPCB queries, to understand scalability with multiple threads we run the tasks with 1, 2, 4, 6 and 8 independent threads. Figure 9 suggests that our modifications did not improve the runtime, in fact they did have a small negative effect (due to excessive use of shields). However, due to the immediate deallocation of the various intermediate collections we were able to almost half the peak working set in 1 thread, and get approximately 3x improvement when we scaled up this benchmark to 8 threads.

CoreCLR micro-benchmarks Apart from the benchmarks discussed above, we also evaluated our mixed-mode memory management approach on 6 data structure-specific micro-benchmarks adapted from the performance test suite of CoreCLR: BinTree, BinTreeLive, BinTreeGrow, RedBlackTree, DirectedGraph and DirectedGraphReplace. Each of these was configured to run with 3 different input sizes (small, medium and large) and 1, 2, 4, 6 and 8 independent threads. Each benchmark involves typical operations with allocations on an underlying container. Details of each benchmark can be found in the supplementary material. Figure 10 plots all our data points in all configurations, and shows massive gains both in runtime and memory that were obtained by porting selectively parts of these data structures to use the manual heap.

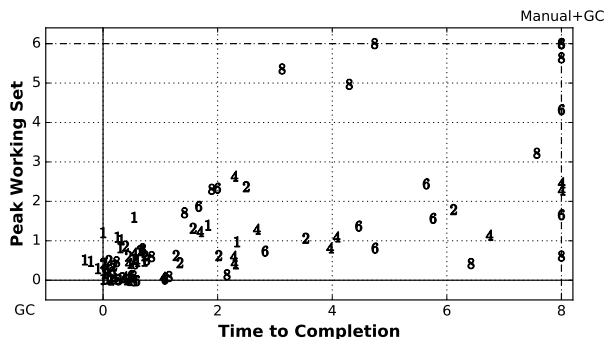


Figure 10: PWS and runtime for micro-benchmarks. We collapse to the border the datapoints that are outside our bounding boxes.

We collapse to the border the datapoints that are outside our bounding boxes.

6 Discussion

6.1 Frontend

We have built a C# language front-end that guarantees the safe use of our API, enforces the correct use of owners and shields, and performs an escape analysis to ensure that manual objects do not escape their shields. It additionally makes programming easier, for example adds implicit `Dispose()` calls for stack-based shields, implicit `Defend()` calls as coercions from owners to shields etc. However the focus of this paper is the runtime and hence we defer any source language modifications for presentation in future work.

6.2 Performance considerations

Cost of shield allocation and setting Shield allocation and setting is a thread-local operation (hence no interlocked barriers are needed); still frequent TLS access can have an adverse affect on performance, hence shield reuse and elimination of redundant shield setting can have beneficial effects. In future work we would like to explore tighter JIT integration for shield allocation, as well as analysis for thread-local objects for which cheaper synchronization can be used.

Allocator characteristics Although jemalloc is a fast and scalable allocator with thread-local caches, it is difficult to match the throughput of the bump-allocator of the .NET generational GC. The situation reverses once surviving objects start to get stored in long-lived and growing collections. In some experiments that

allocate a lot but survive some objects, we have found that it is beneficial to use the GC as nursery and only when we need to store an object *clone* it over to the manual heap. Despite the throughput disadvantage, using a size-class based allocator has the advantage of low fragmentation which means we can avoid compaction (fragmentation and compaction in older generations is a known issue in generational GCs.)

6.3 GC integration

In allocation-intensive workloads in both heaps and lots of pointers from the manual to the GC heap collections will be frequent and have to scan through portions of the manual heap. For ephemeral collections this amounts to a simultaneous iteration over the set cards for the manual heap and the trie, which becomes a hot code path. Furthermore, CoreCLR GC employs concurrent background marking but our implementation at the moment only scans roots synchronously from the manual heap; addressing this is a matter of engineering but explains some of the larger Gen2 pauses in some benchmarks. Finally, when bulky objects are off-loaded to the manual heap, the ephemeral segment budget will fill up with more and smaller objects, potentially increasing the pauses for Gen0 collections. This is often not a problem as the number of Gen0 collections is dramatically reduced but it reveals an opportunity to tune the various parameters of the GC to take the manual heap into account.

Another interesting phenomenon happens when we resize a manual array with GC references to ephemeral objects – if that array was GC then the newly allocated array would be allocated in Gen0 hence no card setting would be required. However in our case, we have references from the manual heap to ephemeral objects and hence card setting is required. Although we have observed this to happen, we have not seen any major performance problems, as the card range is much smaller than the array itself.

6.4 Migrating away from the GC heap

How does a programmer know which objects should be allocated on the manual heap? Our methodology is to use a heap profiling tool (e.g. we have modified `PerfView`, a publicly available tool from Microsoft), to determine if the GC has a significant cost, then identify objects of reasonable size that have survived and have been collected in the older generations, get the stack traces associated with such objects and look for sensible candidates amongst those with clearly defined lifetimes.

6.5 Programming model

By-ref style of owners and shields Owners and shields cannot be duplicated on the stack or returned from arbitrary functions, and hence can only be passed “by-ref”, complicating the porting of an application. For some applications this is not a problem (e.g. `System.Linq.Manual`) but for others it may require redesign of APIs and implementation.

Sharing references on the heap A well known limitation of owner-style objects (and unique references in general) is that sharing of multiple references to the same manual object from the (GC or manual) heap is not allowed, thus making them unsuitable for data structures with a lot of pointer sharing. Snowflake allows shareable pointers to be built around owners. For instance, we have introduced a *class* called `Manual<T>` (as opposed to a struct), that encapsulates an `Owner<T>` and supports similar methods but can be stored on the heap and passed to or returned by other functions like every other GC object. In exchange for greater programming flexibility, object `Manual<T>` incurs extra space cost per manual object, so calls for a careful profile-driven use. Finally, we have also built shareable reference counted objects, `RefCount<T>`, but we are considering API extensions in this space as important future work.

Shields on the heap Shields may accidentally escape on the (shared) heap, resulting in complaints in our C# frontend, predominantly due to the subtle C# *implicit boxing* feature or closure capture. For other examples, we may actually prefer to temporarily store a shield on the heap, as long as we can guarantee that only the very same thread that created that shield will access it and the underlying object (otherwise

the reference to the TLS state is bogus). How to enable this without sacrificing performance is an open challenging problem.

Finalization and owners When owner objects get abandoned our runtime abandons recursively their children owner fields. To do this efficiently we extend the method table with an explicit run-length encoding of the offsets that the owner fields exist in an object and use that to efficiently scan and abandon. An avenue for future work is to modify the layout of objects to group together owner fields to improve this scanning, similarly to what CoreCLR is using for GC pointers. .NET also allows for *finalizers*, which are functions that can be called from a separate thread when objects are no longer live. There is design space to be explored around finalizers for manual objects, e.g. should they be executed by the mutator or scheduled on the GC finalization thread.

Asynchronous tasks In the *async and await* [2] popular C# programming model a thread may produce a value and yield to the scheduler by enqueueing its continuation for later execution. This continuation may be eventually executed *on a different thread*. This decoupling of tasks from threads makes the use of the thread-local shields challenging and it is an open problem of how to allow a task that resumes on a different thread to safely use a shield.

7 Related Work

Several previous studies compared garbage collection with manual memory management [41, 44, 69]. Like us, they concluded that manual memory management can achieve substantial gains in both memory usage and run time, particularly when the heap size is not allowed to grow arbitrarily.

Memory management for managed languages Several systems have proposed optimizing garbage collection for specific tasks – for instance for big data systems [35, 49], taking advantage of idle mutator time [30], specializing to real-time and embedded systems with very low latency constraints [14]. Other work suggests arena-based allocation for moving bulky data out of the GC heap. Scala off-heap [9] provides a mechanism to offload all allocations in a given scope onto the unmanaged heap but no full temporal and thread safety. Stancu et al. [60] introduce a static analysis that can infer a hierarchy of regions and annotate allocation sites so that at runtime a stack-like discipline of regions can be enforced. The analysis ensures that pointers only exist from newer regions to older regions or the GC heap but not vice versa. Broom [36] introduces region-based allocation contexts, but relies on type system for safety. Other work attempts to offload data allocations to the manual heap through program transformations [57]. Recent work proposes hints for safe arena allocations [56], through more expensive write barriers and potential migration of objects when programmer hints are wrong. Our proposal is complementary to these other techniques, and should be seen as yet another tool available to programmers that guarantees memory safety in single- and multithreaded scenarios. For instance, allocating the intermediate dictionaries from our `System.Linq` example in an arena would be cumbersome and less efficient as those allocations are intertwined with other allocations of objects that need to survive, so copying out of the arenas would be required.

Kedia et al. [46] propose a version of .NET with only manually managed memory. Their programming model just exposes a “free object” method, and a dangling pointer exception for accessing reclaimed memory. The runtime reclaims physical pages associated to an object at some point after they are “freed”, and relocates other collocated objects. When an object is accessed an access violation can be triggered due to the object been deallocated or relocated: the runtime determines which and either surfaces the exception, or patches the memory and execution to allow the program to continue. The approach gets surprisingly good results (comparable to the results we report here). The approach has not been integrated with the garbage collector and thus does not provide the pay-for-play cost of our solution. The exceptions for accessing deallocated objects are unpredictable, and dependent on the runtime decisions for when to reclaim memory: some schedule can work by allowing access to a freed but not reclaimed object, while other schedules may not.

Finally, their solution requires the JIT/Compiler to be modified to enable stack walking on every read and write to memory, whereas we do not require any JIT level changes.

Although .NET does not support it, manual or automatic object pre-tenuring [38, 24, 27] would be another way to directly push long lived objects onto a less frequently scanned generation. It is certainly the case that a lot of the performance gains of Snowflake are related to not having to compact and promote between generations, and that is where pre-tenuring would also have a similar effect. However, the significant reduction in the number of Gen2 collections would not be given by pre-tenuring. Using Table 1 for the 4096-byte case we can measure total GC pause times of Gen0 = 3473.4, Gen1 = 9709, Gen2 = 2001, whereas Manual+GC are Gen0 = 2104.0, Gen1 = 2268, and Gen2 = 459. Pre-tenuring could potentially match the Gen 0/1 cost reduction but would not produce savings from not having to scan the older generation.

Safe memory reclamation Our shields are based on the hazard pointers [50] mechanism developed for memory reclamation in lock-free data structures. Hazard pointers require the access to the data structure, in our case defend, to use a heavy write barrier (mfence on x86). We found this was prohibitively expensive for a general-purpose solution. By combining the hazard pointer concept with epochs we can remove the need for this barrier.

Epoch-based reclamation (EBR) has been used for lock-free data structures [39, 34, 40], and has been very elegantly captured in the Rust Crossbeam library [7]. To access a shared data structure you “pin” the global epoch which prevents it advancing (too much), and when you finish accessing the data structure you release the pin. As epochs advance the older deallocated objects can be reclaimed. We could not directly use epochs without hazards, as EBR requires there to be points where you are not accessing the protected data structure. When EBR is protecting a particularly highly-optimised data structure, this is perfectly sensible. We allow pervasive use of manually allocated objects and thus enforcing points where the manual heap is not being accessed is simply impractical. One could view shields as pinning a single object, rather than a global pin as in Crossbeam.

Our use of epochs is actually closer to those used in reference counted garbage collectors such as Recycler [13]. Recycler delays decrements, such that they are guaranteed to be applied after any “concurrent” increments, thus if the reference count reaches zero, then the object can be reclaimed. We are similarly using epochs to ensure that the writes to shields will be propagated to all threads, before any attempt to return it to the underlying memory manager. It is fairly straightforward to extend our mechanism to handle reference counted ownership rather than unique ownership using delayed decrements. Recycler uses stack scanning to deal with stack references into the heap.

Alistarh et al. [12] take the stack scanning approach further and use this instead of hazard pointers for a safe memory reclamation. They use signals to stop each thread and scan its stack. These stack scans are checked before any object can be deallocated. We think our approach may scale better with more threads as we do not have to stop the threads to scan the stacks, but the sequential throughput would be lower for us as we have to perform the shield TLS assignment for what we are accessing.

There are several other schemes that use hazard pointers with another mechanism to remove the requirement for a memory barrier. Balmau et al. [16] extend hazard pointers in a very similar way to our use of epochs. Rather than epochs they track that every thread has performed a context switch, and use this to ensure that the hazard pointers are sufficiently up to date. If the scheduling leads to poor performance, they drop back to a slow path that uses standard hazard pointers with a barrier. Our mechanism for ejection means we do not to have a slow path even when the scheduler starves particular threads, we believe our ejection mechanism could be added to their scheme. Dice et al. [33] also develop a version of hazard pointers without a memory barrier. When the hazards need scanning, virtual memory protection is used to ensure the hazards are read only, which forces the correct ordering of the checking and the write. A mutator threads that attempts to assign a hazard will take an access violation (AV) and then block until the scan has finished. This means any reclamation is going to incur two rounds of inter-processor interrupt (IPIs), and mutators are likely to experience an AV. We only use the write protection in very infrequent cases of threads miss behaving, and rarely have to deal with the AVs.

Morrison and Affek [53] use time as a proxy for an epoch scheme: effectively they wait so many cycles, and

have observed empirically observed that the memory updates will have propagated by this point. This leads to a simple design, no thread can stop time advancing, hence they do not require an Ejection mechanism like we have. However, to achieve this they make assumptions on the hardware that are currently not guaranteed by the vendors about timing.

There are other approaches to reclamation that use data-structure specific fixup [23] or roll-back [29, 28] code to handle cases where deallocation occurs during an operation. This would not be practical for our setting as we are using manual memory management pervasively and do not have nice data structure boundaries that these schemes exploit.

Techniques for unsafe languages Several systems use page-protection mechanism to add temporal safety to existing unsafe languages: [48, 4, 5, 31] these approaches are either probabilistic or suffer from performance problems. Some systems propose weaker guarantees for safe manual memory management. Cling [11] and [32] allow reuse of objects having same type and alignment. DieHard(er) [17, 58] and Archipelago [48] randomize allocations to make the application less vulnerable to memory attacks. Several systems detect accesses to freed objects [55, 47, 68], but do not provide full type safety.

Type systems for manual memory management The Cyclone language [62, 42] is a safe dialect of C with conservative garbage collector [20] and several forms of safe manual memory management, including stack and region-based allocation [64, 37]. Capability types [66] can be used to verify the safety of region-based memory management. Unlike Cyclone, we do not use regions as a basis for safe manual memory management. Furthermore we allow eager deallocation at arbitrary program locations. Several languages have proposed using unique pointers to variables or objects [43, 52, 54] based on linear typing [65, 15]. Our owners are a form of unique pointers, but we allow stack sharing of the references using shields. This is similar to the concept of borrowing references to temporarily use them in specific lexical scopes [67, 22, 25, 62]. Rust [8] incorporates several aspects of the Cyclone design, including integration of manually managed memory with reference counting, unique pointers, and lexically-scoped borrowing. Finally, languages with ownership types [26, 25, 21] and alias types [59] can express complex restrictions on the object graphs that a program can create. Though owners cannot express cycles, we do allow sharing through the GC, and permit cross-heap pointers in both directions.

8 Conclusion

We have presented a design that integrates safe manual memory management with garbage collection in the .NET runtime, based on owners and shields. Our programming model allows for stack-based sharing of owners potentially amongst multiple threads, as well as arbitrary deallocations while guaranteeing safety. We show that this design allows programmers to mix-and-match GC and manually-allocated objects to achieve significant performance gains.

References

- [1] Asp.net/caching: Libraries for in-memory caching and distributed caching. <https://github.com/aspnet/Caching>.
- [2] Asynchronous programming with async and await. <https://msdn.microsoft.com/en-us/library/mt674882.aspx>.
- [3] Coreclr. www.github.com/dotnet/CoreCLR.
- [4] Electric fence malloc debugger. http://elinux.org/Electric_Fence.
- [5] How to use pageheap utility to detect memory errors. <https://support.microsoft.com/en-us/kb/264471>.

- [6] Jemalloc. <http://jemalloc.net>.
- [7] Rust crossbeam library. <http://aturon.github.io/crossbeam-doc/crossbeam/mem/epoch/index.html>.
- [8] Rust programming language. <https://www.rust-lang.org>.
- [9] Scala-offheap: Type-safe off-heap memory for scala. <https://github.com/densh/scala-offheap>.
- [10] Tpch. <http://www.tpch.org/tpch>.
- [11] P. Akritidis. Cling: A memory allocator to mitigate dangling pointers. In *USENIX Security Symposium*, pages 177–192, 2010.
- [12] D. Alistarh, W. M. Leiserson, A. Matveev, and N. Shavit. Threadscan: Automatic and scalable memory reclamation. In *SPAA*, 2015.
- [13] D. F. Bacon, C. R. Attanasio, H. B. Lee, V. T. Rajan, and S. Smith. Java without the coffee breaks: A nonintrusive multiprocessor garbage collector. *PLDI*, 2001.
- [14] D. F. Bacon, P. Cheng, and V. Rajan. The metronome: A simpler approach to garbage collection in real-time systems. In *In Workshop on Java Technologies for Real-Time and Embedded Systems (JTRES), OTM Workshops*, 2003.
- [15] H. G. Baker. Use-once variables and linear objects—storage management, reflection, and multi-threading. *SIGPLAN Notices*, 30(1):45–52, January 1995.
- [16] O. Balmau, R. Guerraoui, M. Herlihy, and I. Zabolotchi. Fast and robust memory reclamation for concurrent data structures. In *SPAA*, 2016.
- [17] E. D. Berger and B. G. Zorn. Diehard: probabilistic memory safety for unsafe languages. In *Acm sigplan notices*, volume 41, pages 158–168. ACM, 2006.
- [18] S. M. Blackburn and K. S. McKinley. Immix: a mark-region garbage collector with space efficiency, fast collection, and mutator performance. In *PLDI*, 208.
- [19] B. H. Bloom. Space/time trade-offs in hash coding with allowable errors. *Commun. ACM*, 13(7), 1970.
- [20] H.-J. Boehm and M. Weiser. Garbage collection in an uncooperative environment. *Software – Practice and Experience*, 18(9):807–820, 1988.
- [21] C. Boyapati, A. Salcianu, W. Beebe, and M. Rinard. Ownership types for safe region-based memory management in real-time Java. In *PLDI*, 2003.
- [22] J. Boyland. Alias burying: Unique variables without destructive reads. *Software – Practice and Experience*, 31(6):533–553, 2001.
- [23] T. A. Brown. Reclaiming memory for lock-free data structures: There has to be a better way. In *PODC*, 2015.
- [24] P. Cheng, R. Harper, and P. Lee. Generational stack collection and profile-driven pretenuing. In *Proceedings of the ACM SIGPLAN 1998 Conference on Programming Language Design and Implementation*, PLDI '98, pages 162–173, New York, NY, USA, 1998. ACM.
- [25] D. Clarke and T. Wrigstad. External uniqueness is unique enough. In *ECOOP*, pages 176–200, July 2003.

- [26] D. G. Clarke, J. M. Potter, and J. Noble. Ownership types for flexible alias protection. In *OOPSLA*, October 1998.
- [27] D. Clifford, H. Payer, M. Stanton, and B. L. Titzer. Memento mori: Dynamic allocation-site-based optimizations. In *Proceedings of the 2015 International Symposium on Memory Management, ISMM '15*, pages 105–117, New York, NY, USA, 2015. ACM.
- [28] N. Cohen and E. Petrank. Automatic memory reclamation for lock-free data structures. In *OOPSLA*, 2015.
- [29] N. Cohen and E. Petrank. Efficient memory management for lock-free data structures with optimistic access. In *SPAA*, 2015.
- [30] U. Degenbaev, J. Eisinger, M. Ernst, R. McIlroy, and H. Payer. Idle time garbage collection scheduling. In *PLDI*, 2016.
- [31] D. Dhurjati and V. Adve. Efficiently detecting all dangling pointer uses in production servers. In *DSN*, June 2006.
- [32] D. Dhurjati, S. Kowshik, V. Adve, and C. Lattner. Memory safety without runtime checks or garbage collection. *ACM SIGPLAN Notices*, 38(7):69–80, 2003.
- [33] D. Dice, M. Herlihy, and A. Kogan. Fast non-intrusive memory reclamation for highly-concurrent data structures. In *ISMM*.
- [34] K. Fraser. Practical lock-freedom. PhD Thesis UCAM-CL-TR-579, Computer Laboratory, University of Cambridge, February 2004.
- [35] L. Gidra, G. Thomas, J. Sopena, M. Shapiro, and N. Nguyen. NumaGiC: a garbage collector for big data on big NUMA machines. In *ASPLOS*, 2015.
- [36] I. Gog, J. Giceva, M. Schwarzkopf, K. Vaswani, D. Vytiniotis, G. Ramalingam, M. Costa, D. G. Murray, S. Hand, and M. Isard. Broom: Sweeping out garbage collection from big data systems. In *HotOS*, 2015.
- [37] D. Grossman, G. Morrisett, and T. Jim. Region-based memory management in Cyclone. In *PLDI*, 2002.
- [38] T. L. Harris. Dynamic adaptive pre-tenuring. In *Proceedings of the 2Nd International Symposium on Memory Management, ISMM '00*, pages 127–136, New York, NY, USA, 2000. ACM.
- [39] T. L. Harris. A pragmatic implementation of non-blocking linked-lists. In *DISC*, 2001.
- [40] T. E. Hart, P. E. McKenney, A. D. Brown, and J. Walpole. Performance of memory reclamation for lockless synchronization. *Journal of Parallel and Distributed Computing*, 67:1270–1285, May 2007.
- [41] M. Hertz and E. D. Berger. Quantifying the performance of garbage collection vs. explicit memory management. In *OOPSLA*, 2005.
- [42] M. Hicks, G. Morrisett, D. Grossman, and T. Jim. Experience with safe manual memory-management in Cyclone. In *ISMM*, 2004.
- [43] J. Hogg. Islands: Aliasing protection in object-oriented languages. In *OOPSLA*, 1991.
- [44] R. Hundt. Loop recognition in C++/Java/Go/Scala. In *Proceedings of Scala Days 2011*, 2011.
- [45] R. Jones, A. Hosking, and E. Moss. *The Garbage Collection Handbook: The Art of Automatic Memory Management*. Chapman & Hall/CRC, 1st edition, 2011.

- [46] P. Kedia, M. Costa, M. Parkinson, K. Vaswani, and D. Vytiniotis. Simple, fast and safe manual memory management. In *PLDI*, 2017.
- [47] B. Lee, C. Song, Y. Jang, and T. Wang. Preventing use-after-free with dangling pointer nullification. In *NDSS*, 2015.
- [48] V. B. Lvin, G. Novark, E. D. Berger, and B. G. Zorn. Archipelago: trading address space for reliability and security. In *ASPLOS*, 2008.
- [49] M. Maas, K. Asanović, T. Harris, and J. Kubiawicz. Taurus: A holistic language runtime system for coordinating distributed managed-language applications. In *ASPLOS*, 2016.
- [50] M. M. Michael. Hazard pointers: Safe memory reclamation for lock-free objects. *IEEE Transactions on Parallel and Distributed Systems*, 15(6):491–504, June 2004.
- [51] T. Minka, J. Winn, J. Guiver, S. Webster, Y. Zaykov, B. Yangel, A. Spengler, and J. Bronskill. Infer.NET 2.6, 2014. Microsoft Research Cambridge. <http://research.microsoft.com/infernet>.
- [52] N. Minsky. Towards alias-free pointers. In *ECOOP*, pages 189–209, July 1996.
- [53] A. Morrison and Y. Afek. Temporally bounding tso for fence-free asymmetric synchronization. In *ASPLOS*, 2015.
- [54] K. Naden, R. Bocchino, J. Aldrich, and K. Bierhoff. A type system for borrowing permissions. In *POPL*, 2012.
- [55] S. Nagarakatte, J. Zhao, M. M. K. Martin, and S. Zdancewic. CETS compiler-enforced temporal safety for c. In *ISMM*, 2010.
- [56] K. Nguyen, L. Fang, G. Xu, B. Demsky, S. Lu, S. Alamian, and O. Mutlu. Yak: A high performance big-data-friendly garbage collector. In *OSDI*, 2016.
- [57] K. Nguyen, K. Wang, Y. Bu, L. Fang, J. Hu, and G. Xu. Facade: A compiler and runtime for (almost) object-bounded big data applications. In *ASPLOS*, 2015.
- [58] G. Novark and E. D. Berger. Dieharder: securing the heap. In *Proceedings of the 17th ACM conference on Computer and communications security*, pages 573–584. ACM, 2010.
- [59] F. Smith, D. Walker, and G. Morrisett. Alias types. In *European Symposium on Programming (ESOP)*, 2000.
- [60] C. Stancu, C. Wimmer, S. Brunthaler, P. Larsen, and M. Franz. Safe and efficient hybrid memory management for java. In *ISMM*, 2015.
- [61] D. Stefanovic, K. S. McKinley, and J. E. B. Moss. Age-based garbage collection. In *OOPSLA*, 1999.
- [62] N. Swamy, M. Hicks, G. Morrisett, D. Grossman, and T. Jim. Safe manual memory-management in Cyclone. *Science of Computer Programming*, 62(2):122–14, October 2006.
- [63] G. Tene, B. Iyengar, and M. Wolk. C4: The continuously concurrent compacting collector. In *ISMM*, 2011.
- [64] M. Tofte and J.-P. Talpin. Region-based memory management. *Information and Computation*, 132(2):109–176, February 1997.
- [65] P. Wadler. Linear types can change the world! In *IFIP TC 2 Working Conference*, 1990.
- [66] D. Walker, K. Crary, and G. Morrisett. Typed memory management in a calculus of capabilities. *ACM Transactions on Programming Languages and Systems*, 24(4):701–771, 2000.

- [67] D. Walker and K. Watkins. On regions and linear types. In *ICFP*, 2001.
- [68] Y. Younan. FreeSentry: protecting against user-after-free vulnerabilities due to dangling pointers. In *NDSS*, 2015.
- [69] B. G. Zorn. The measured cost of conservative garbage collection. *Software – Practice and Experience*, 23(7):733–756, 1993.

A Appendix

We performed experiments to validate the usefulness of our mixed-mode memory management approach. Due to space constraints, we could not fit all experiments and plots in the main paper, and thus we present them here as supplementary material.

Experimental Setup We performed the majority of the experiments on a 3.5GHz Intel Xeon CPU E5-1620 v3 (8 physical cores) with 16GB RAM running Windows 10 Enterprise (64-bit). For the TPCB/LINQ experiments using the big data set, we used a similar machine with 32GB of memory. We used our own branch of CoreCLR, CoreFX and jemalloc.

Workstation versus Server GC For the experiments in the main paper, we configured CoreCLR in “Workstation GC” mode, as this is the configuration that our implementation is currently stable on. As discussed in the main paper, CoreCLR additionally provides a “Server GC” mode, which allows for independent thread-local collection of each thread’s ephemeral heap. We have only run preliminary experiments against this configuration and the relative trends are similar to our findings with Workstation GC, but the absolute numbers show that Server GC trades higher peak working sets for lower times to completion (as seen in the experiments below that support both Workstation and Server GC).

A.1 ASP.NET Caching

The ASP.NET Caching benchmark spawns a number of threads that perform 4 million requests on a *shared* cache. Each request randomly generates a key to an entry in the cache, and attempts to get its value. If the entry does not exist, then a new entry is allocated. The size of the entry payload is randomly generated using a distribution (in the results presented here drawn from a triangular distribution, but we have experimented with uniform sizes and constant sizes with similar results). The cache is also configured with a sliding expiration of 1 second (we experimented with various expiration values, and found similar results). The mean size of each entry is 256, 2048, 4096 bytes. The cache is shared among 1, 2, 4, 6 and 8 threads in our experiments.

Figure 11 shows the results for time to completion, and Figure 12 shows the peak working sets. Both plots show results for both Workstation and Server GC. Thread configurations are denoted by a number on top of each bar.

A.2 TPCB on System.Linq.Manual

For evaluating `System.Linq.Manual` we used the 22 industry-standard TPCB queries,² ported them in LINQ C#, and simply run them using `System.Linq` (for the GC experiments) and `System.Linq.Manual` (for the Manual+GC experiments). We have scaled the TPCB dataset to a version where 300MB of tables is pre-loaded in memory (small dataset), and another one where 3GB of tables is pre-loaded in memory (large dataset). To understand the scalability with multiple threads we run each query in 10 iterations using 1,2,4,6,8 *independent* threads, i.e. we did not parallelize the query, simply had a number of worker threads executing the same query once the data is loaded (which happens just once for all threads). We only run this experiment using Workstation GC, since the integration with Server GC is not yet stable enough to run this benchmark.

For spacing reasons, in the main paper we only presented a scatter plot with the results from the small dataset. Here, we present detailed time-to-completion and peak working set results (see bar plots in Figures 14-25) for each individual TPCB query using both the small and large dataset.³ Figure 13 presents the 21x5 (5 is the number of thread configurations) results from the large dataset as a scatter plot where the vertical (respectively horizontal) axis shows *relative* improvement in the peak working sets

²<http://www.tpch.org/tpch>

³Note that query 16 is missing from the large dataset due to an unknown bug discovered during submission.

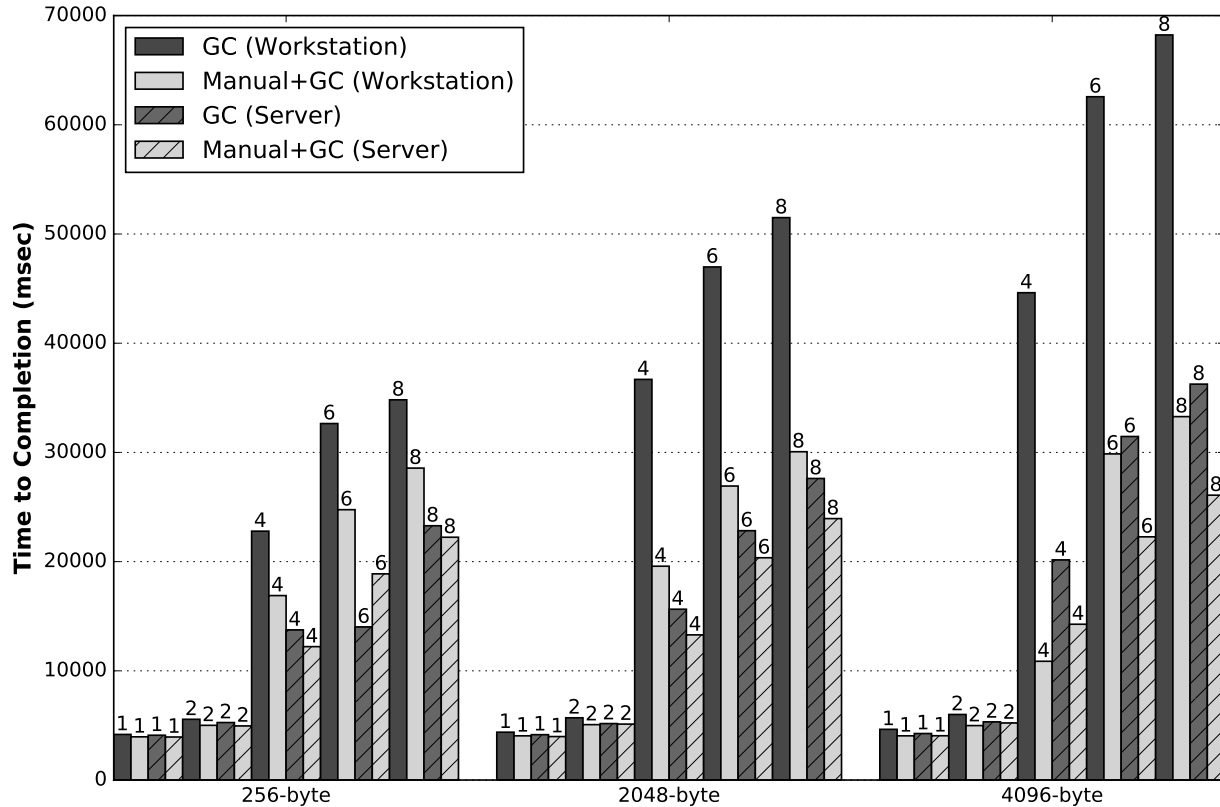


Figure 11: Comparison of time to completion in ASP.NET Caching using both Workstation and Server GC.

(respectively runtime). The labeling indicates number of threads. Positive numbers in both axes mean that we are faster *and* consume less memory.

To calculate the relative improvement, for all queries in the TPCB benchmark, and all thread configurations, if the GC peak working set is less than our peak working set then we plot $(1.0 - \text{Manual_PWS} / \text{GC_PWS})$ else we plot $(1.0 - \text{GC_PWS} / \text{Manual_PWS})$. We plot the relative time to completion using an analogous computation.

Finally, Table 2 presents the full GC collections and pauses from running the 22 TPCB queries using 1 thread and the small dataset.

A.3 Machine learning on Infer.NET

We evaluated Snowflake on a workload from a research project with known GC overheads based on the Infer.NET framework. The task was to compute the product of two large probabilistic automata for a small number of iterations.

Table 3 presents the GC collections and pauses from running this benchmark using 1 thread.

A.4 Micro-benchmarks

Besides the industrial benchmarks discussed above, we also evaluated our mixed-mode memory management approach on 6 data structure-specific micro-benchmarks from the CoreCLR GC performance testing suite: DirectedGraph, DirectedGraphReplace, BinTree, BinTreeLive, BinTreeGrow and RedBlackTree. We also evaluated our approach on HavlakPool, an algorithm for finding loops in an input control-flow graph (CFG)

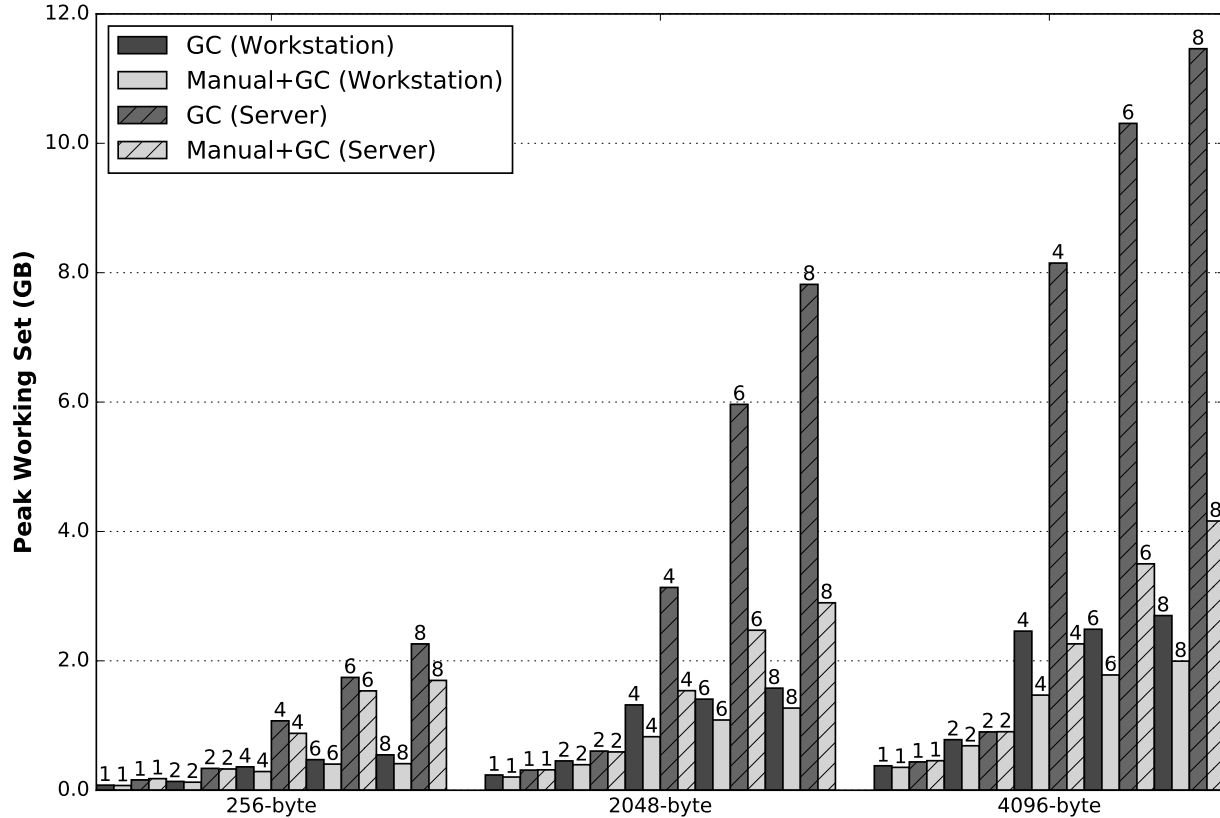


Figure 12: Comparison of peak working set in ASP.NET Caching using both Workstation and Server GC.

not coming from the CoreCLR performance suite. Each of these micro-benchmarks was configured to run with 3 different input sizes (small, medium and large) and 1, 2, 4, 6 and 8 independent threads.

The BinTree benchmark constructs a binary tree (50000 nodes in the large configuration, 25000 in medium and 10000 in the small), then removes a fixed number (20000 in large, 10000 in medium, 5000 in small) of randomly selected nodes (using a uniform distribution) from the tree, and then finally inserts a fixed number (10000 in large, 5000 in medium, 2500 in small) of newly allocated nodes with randomly selected keys (again using a uniform distribution). This process repeats for 100 iterations. Each node contains a reference to its left and right child and an array of bytes. In our port, we moved all the nodes, the children references and the payloads in the manual heap. The BinTreeLive is a variation of the BinTree benchmark: each time a node is visited, its array is resized (the first time its visited it grows from 10 elements to 1000, then the second time it goes back to 10 elements, then back to 1000 elements, repeat). The BinTreeGrow is similar to BinTreeLive, but each time a node is visited, the internal array grows by 100 elements (to not run out of space, we only run 10 iterations of BinTreeGrow).

RedBlackTree constructs a red-black tree (10000 nodes in the large configuration, 5000 in medium and 1000 in the small), and then in each of 1000 iterations it adds and deletes a number of nodes. Each node has an internal array of integers. In this benchmark, we left the actual nodes in the GC heap, but placed the arrays in the manual heap. DirectedGraph constructs a directed graph (100000 nodes in the large configuration, 50000 in medium and 10000 in the small), and then in each iteration it randomly replaces a number of nodes (selected with a uniform distribution). Every graph node contains an array of adjacent nodes, which we ported to the manual heap. The DirectedGraphReplace is a modification of DirectedGraph; instead of deleting a node from the graph, we replace it with a newly allocated node.

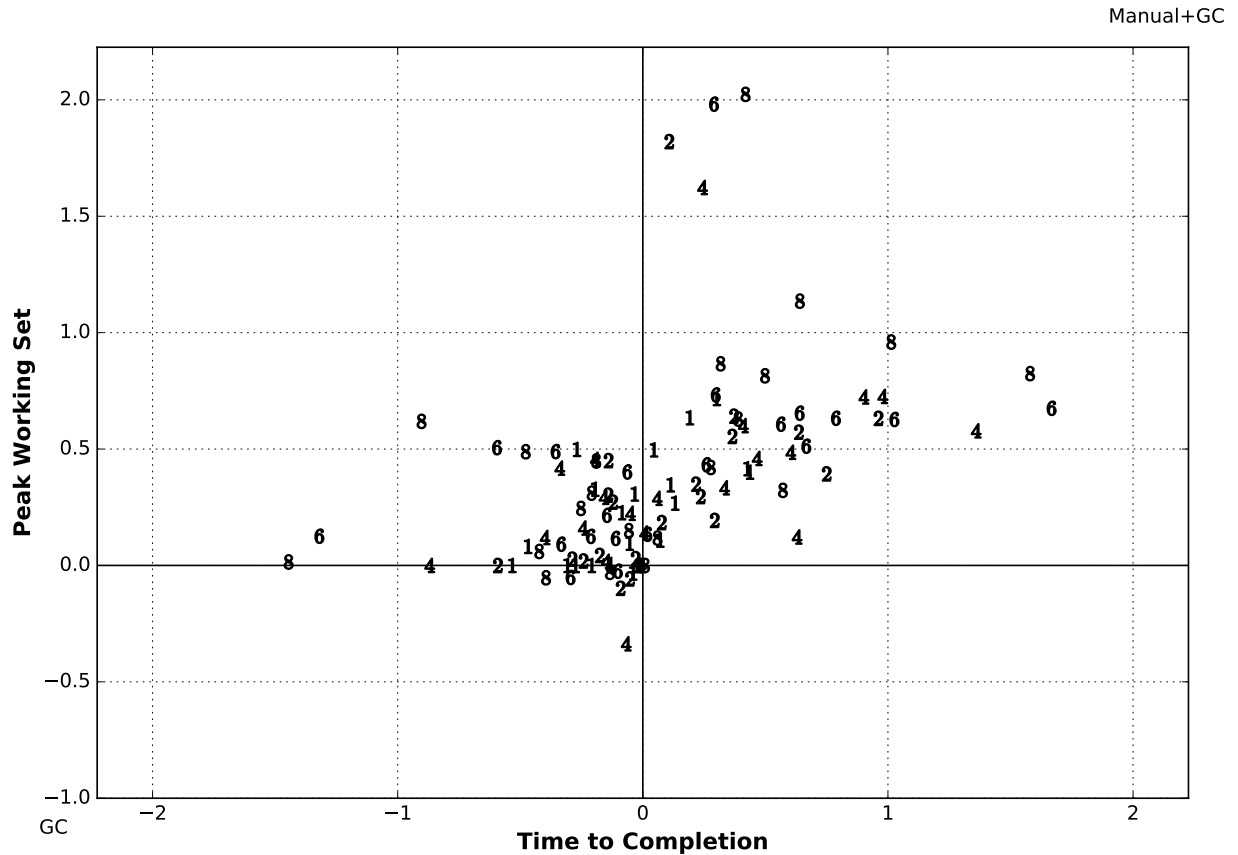


Figure 13: Comparison of peak working set and time to completion on the TPCH large dataset using `System.Linq.Manual`.

For the `BinTree`, `BinTreeLive`, `BinTreeGrow` and `RedBlackTree` benchmarks, we show results using both Workstation and Server GC. However, for the `DirectedGraph`, `DirectedGraphReplace` and `HavlakPool` benchmarks, we only show results using Workstation GC, since the integration with Server GC is not yet stable enough to run these benchmarks. We present the following figures:

DirectedGraph Figure 26 compares the time to completion in msec (log scale) on the left hand side, and the peak working set in GB (log scale) on the right hand side.

DirectedGraphReplace Figure 27 compares the time to completion in msec (log scale) on the left hand side, and the peak working set in GB (log scale) on the right hand side.

BinTree Figure 28 compares the time to completion in msec (log scale). Figure 29 compares the peak working set in GB.

BinTreeLive Figure 30 compares the time to completion in msec. Figure 31 compares the peak working set in GB. Both figures are in log scale.

BinTreeGrow Figure 32 compares the time to completion in msec (log scale). Figure 33 compares the peak working set in GB.

RedBlackTree Figure 34 compares the time to completion in msec. Figure 35 compares the peak working set in GB. Both figures are in log scale.

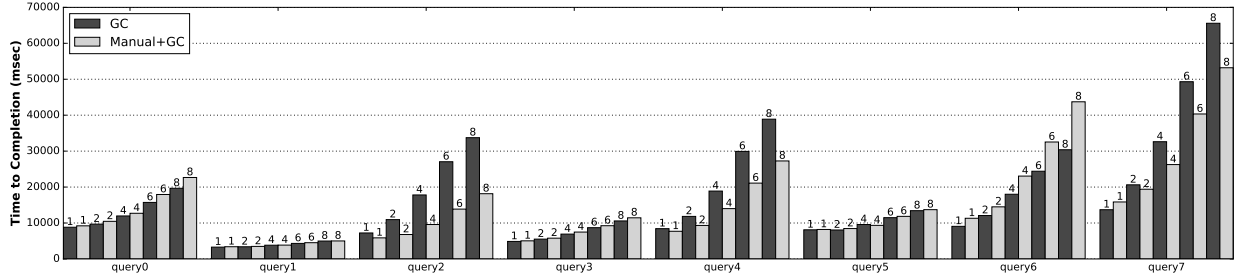


Figure 14: Comparison of time to completion on the TPCCH small dataset (queries 0 to 7) using System.Linq.Manual.

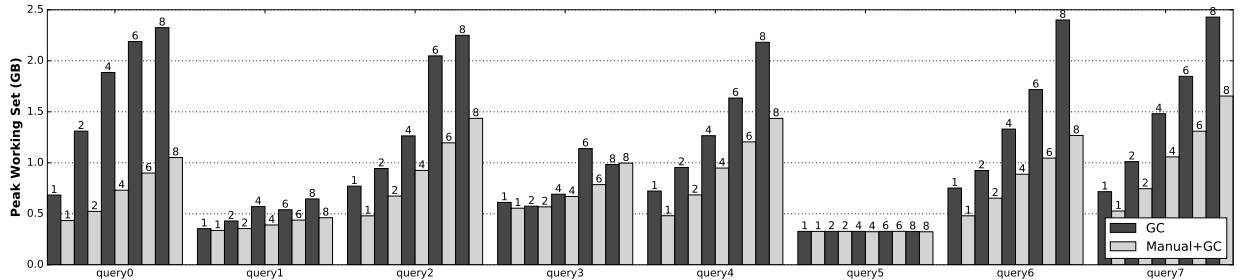


Figure 15: Comparison of peak working set on the TPCCH small dataset (queries 0 to 7) using System.Linq.Manual.

HavlakPool Figure 36 compares the end-to-end time (in msec), which includes the construction of the input CFG and the time taken to run the Havlak loop finding algorithm on the specified CFG. Figure 37 compares only the time (in msec) taken to complete the Havlak loop finding algorithm. Figure 38 compares the peak working set in GB.

Table 4 presents the full GC collections and pauses data from running the HavlakPool benchmark using 1 thread and the small, medium and large configurations. HavlakPool is an interesting example and deserves some attention. It has a bi-modal behaviour: in the first allocation-intensive phase of the algorithm a huge control flow graph is constructed. This is an allocation-heavy phase with the vast majority of objects surviving to the oldest generations. After that phase, the request phase is trying to identify loops in the control flow graph by using a pool of objects (and a pool of owner objects in our case) to improve throughput – not all allocations are eliminated though. We see in Figure 38 that PWS are not substantially different, we in fact have a slightly higher PWS. In terms of runtime, we are also slightly slower than the purely GC-based version; however the breakdown to the total time (sum of two phases, Figure 36), versus just the loop finding time (Figure 37) reveals that our loop finding is much faster once the CFG has been initialized in memory. The reason that the first phase (CFG construction) is so slow is a combination of (a) lower jemalloc throughput, (b) inefficiencies in scanning the cards corresponding to the manual heap, (c) lack of support for asynchronous background collection from the manual heap roots. In fact Table 4 shows substantially higher collection pauses, although the number of collections is lower. We have discussed this phenomenon in the main paper and suggested directions for improvement.

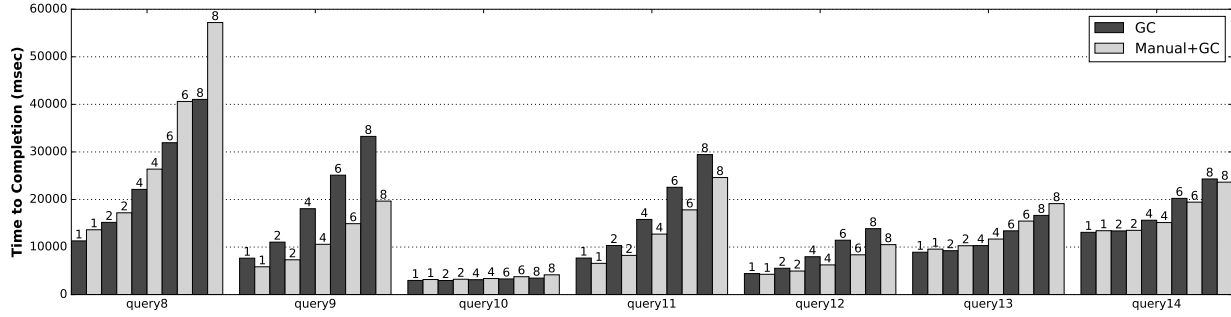


Figure 16: Comparison of time to completion on the TPCCH small dataset (queries 8 to 14) using System.Linq.Manual.

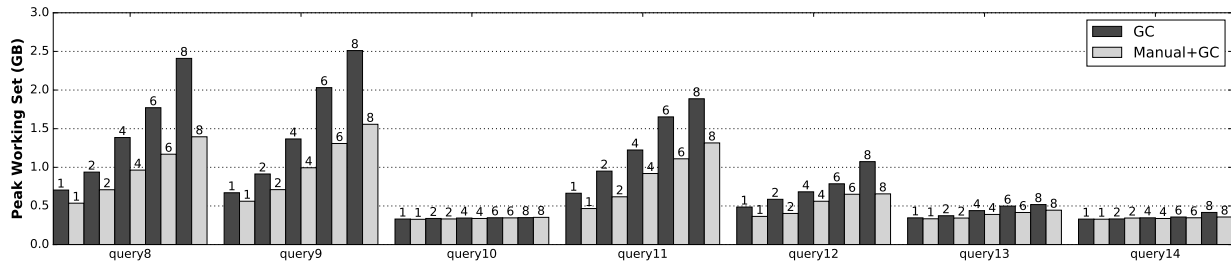


Figure 17: Comparison of peak working set on the TPCCH small dataset (queries 8 to 14) using System.Linq.Manual.

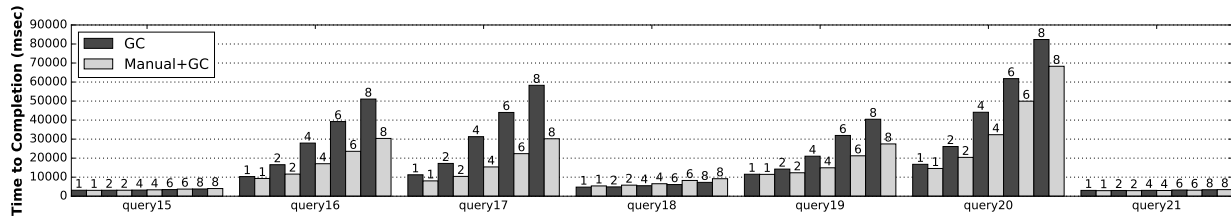


Figure 18: Comparison of time to completion on the TPCCH small dataset (queries 15 to 21) using System.Linq.Manual.

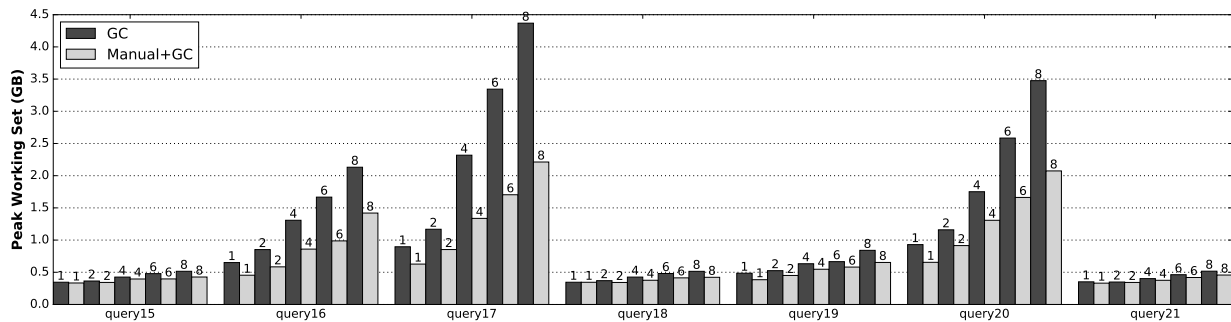


Figure 19: Comparison of peak working set on the TPCCH small dataset (queries 15 to 21) using System.Linq.Manual.

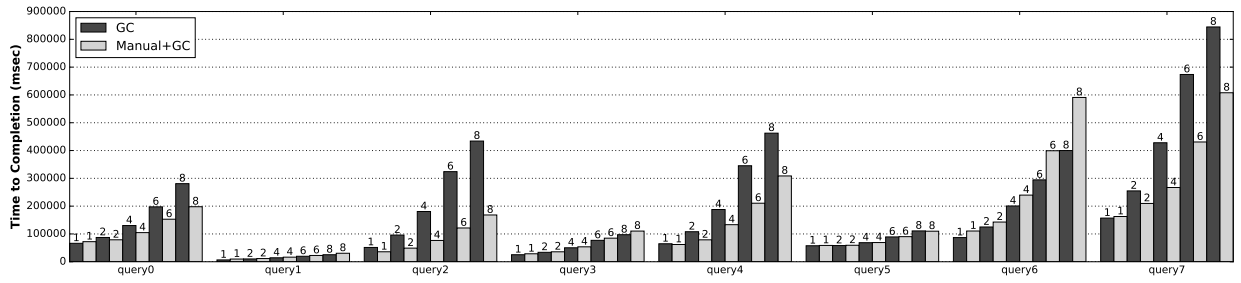


Figure 20: Comparison of time to completion on the TPCCH large dataset (queries 0 to 7) using System.Linq.Manual.

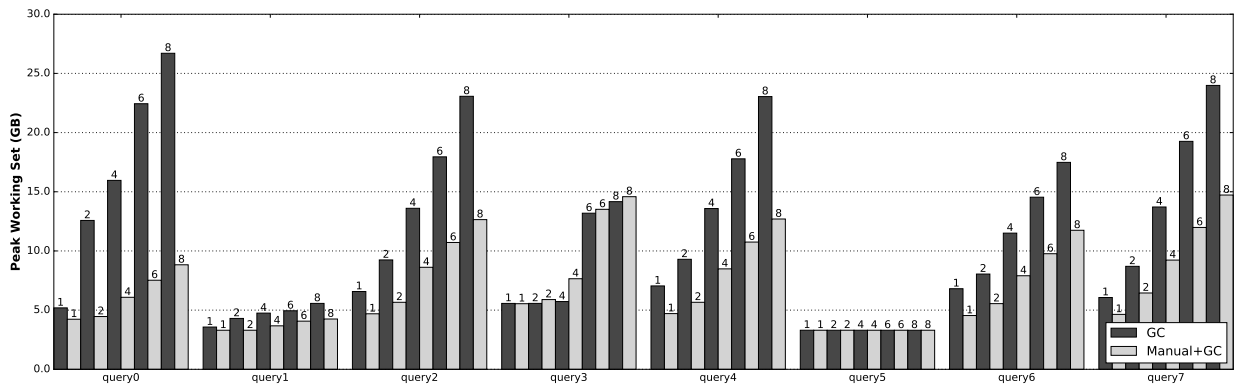


Figure 21: Comparison of peak working set on the TPCCH large dataset (queries 0 to 7) using System.Linq.Manual.

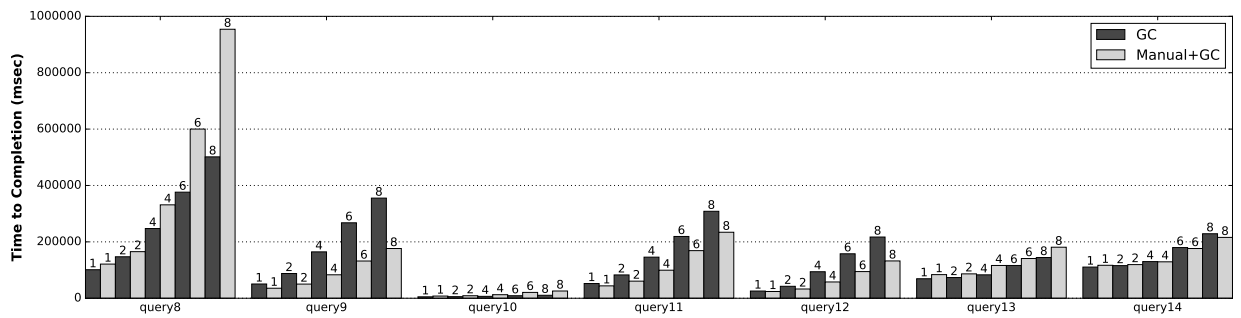


Figure 22: Comparison of time to completion on the TPCCH large dataset (queries 8 to 14) using System.Linq.Manual.

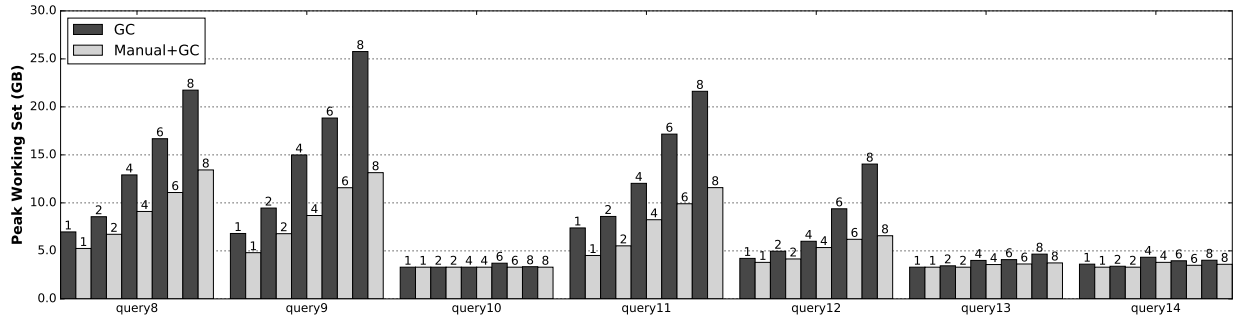


Figure 23: Comparison of peak working set on the TPC-H large dataset (queries 8 to 14) using System.Linq.Manual.

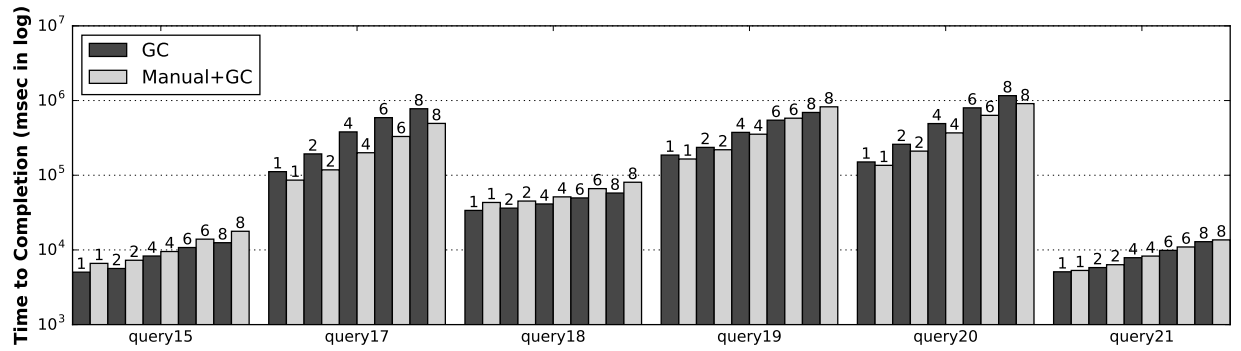


Figure 24: Comparison of time to completion (log scale) on the TPC-H large dataset (queries 15 to 21) using System.Linq.Manual.

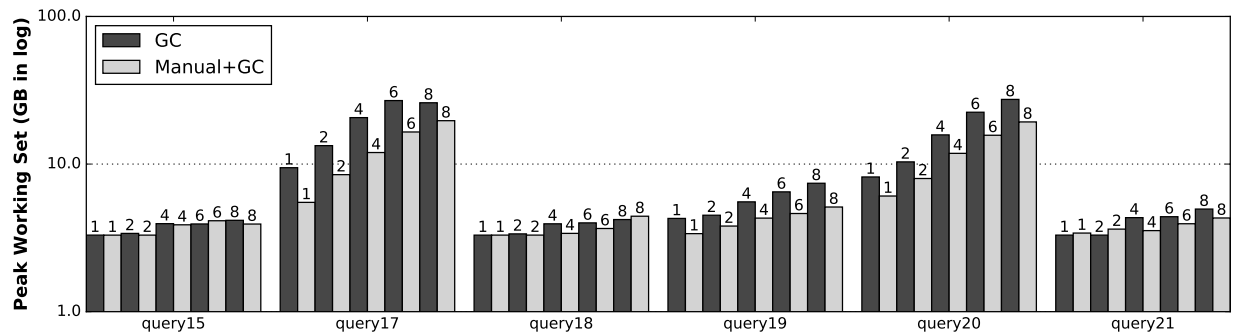


Figure 25: Comparison of peak working set (log scale) on the TPC-H large dataset (queries 15 to 21) using System.Linq.Manual.

Table 2: Collections and GC pauses in the 22 TPCB queries using `System.Linq.Manual` (1 thread).

Query	Heap	%GC	Gen0			Gen1			Gen2		
			Count	Pause Time		Count	Pause Time		Count	Pause Time	
				Mean	Max		Mean	Max		Mean	Max
0	GC	13.4%	225	1.2	8.2	72	3.3	7.7	15	1.0	2.2
0	Man+GC	4.9%	224	1.2	4.5	71	3.3	7.0	7	0.4	0.8
1	GC	15.3%	233	1.3	4.0	82	3.8	10.5	8	0.6	1.7
1	Man+GC	13.2%	233	1.3	3.8	73	3.4	7.5	8	0.7	2.3
2	GC	34.6%	411	1.5	15.1	165	9.8	25.3	16	1.3	3.8
2	Man+GC	11.0%	252	1.5	6.3	89	4.3	24.3	8	0.8	3.5
3	GC	15.5%	213	1.5	6.8	81	4.6	19.4	10	0.9	1.8
3	Man+GC	14.1%	212	1.5	6.8	77	3.6	8.5	10	1.4	4.5
4	GC	33.1%	542	1.1	22.6	165	9.8	25.6	16	1.1	3.0
4	Man+GC	9.8%	367	1.3	6.7	103	3.8	29.2	8	0.7	2.9
5	GC	5.5%	199	1.3	4.8	70	3.3	7.0	7	0.5	0.8
5	Man+GC	5.5%	199	1.3	3.8	70	3.3	7.0	7	0.5	0.7
6	GC	20.8%	687	0.7	19.7	160	5.7	38.3	15	0.8	2.0
6	Man+GC	9.0%	561	1.0	7.6	124	3.2	25.1	8	0.7	2.4
7	GC	21.7%	778	0.7	7.4	292	6.3	34.6	17	0.9	1.5
7	Man+GC	6.8%	732	1.2	7.6	83	4.7	37.2	8	1.0	4.1
8	GC	22.3%	664	1.2	19.4	173	6.9	39.6	15	0.7	1.1
8	Man+GC	12.1%	585	2.1	7.8	84	5.4	34.0	8	1.6	9.5
9	GC	36.0%	363	1.7	6.6	159	11.3	37.8	16	1.3	3.5
9	Man+GC	11.8%	218	1.7	8.7	83	5.2	30.9	8	1.0	4.2
10	GC	12.3%	235	1.2	4.3	72	3.3	7.3	7	0.5	0.8
10	Man+GC	11.8%	240	1.2	4.0	71	3.4	6.6	7	0.5	0.9
11	GC	31.6%	500	1.1	3.7	153	9.5	25.2	15	1.0	1.9
11	Man+GC	11.0%	345	1.1	6.5	100	3.8	26.3	8	0.7	1.9
12	GC	20.0%	217	2.2	16.6	86	6.6	46.5	8	0.9	4.4
12	Man+GC	11.5%	207	1.6	10.7	79	4.4	19.2	8	1.0	4.1
13	GC	7.3%	360	0.9	4.6	86	3.3	6.9	8	0.6	1.3
13	Man+GC	5.9%	370	1.0	4.8	74	3.4	8.7	7	0.5	0.8
14	GC	4.4%	207	1.3	3.7	69	3.6	10.0	8	0.6	1.3
14	Man+GC	3.7%	201	1.3	4.4	72	3.3	9.7	7	0.5	0.8
15	GC	14.3%	227	1.3	3.4	82	3.2	6.4	8	0.6	1.7
15	Man+GC	12.1%	233	1.2	6.3	74	3.3	8.8	7	0.7	1.6
16	GC	28.2%	541	1.1	8.4	188	7.5	28.9	19	1.4	8.3
16	Man+GC	8.5%	412	1.3	7.9	76	4.0	26.7	8	0.8	2.2
17	GC	41.5%	518	1.5	16.6	261	11.3	26.7	20	1.2	3.2
17	Man+GC	10.0%	247	1.8	8.5	89	5.1	21.9	8	1.3	6.3
18	GC	11.3%	367	0.9	4.1	86	3.1	6.6	8	0.5	1.3
18	Man+GC	9.0%	376	0.9	4.9	74	3.3	8.7	7	0.5	0.8
19	GC	14.5%	350	2.5	8.8	131	6.9	31.9	9	0.9	4.1
19	Man+GC	7.9%	381	1.5	5.4	87	3.8	8.7	8	0.6	1.7
20	GC	43.1%	895	1.4	51.2	398	10.2	51.1	37	2.4	11.2
20	Man+GC	15.1%	716	2.2	10.9	100	9.1	40.7	10	1.3	5.4
21	GC	12.1%	207	1.3	3.6	73	3.5	7.9	7	0.5	0.7
21	Man+GC	11.8%	204	1.3	4.3	73	3.4	9.0	7	0.4	0.7

Table 3: Collections and GC pauses in the Infer.NET machine learning benchmark (1 thread).

Heap	%GC	Gen0			Gen1			Gen2		
		Count	Pause Time		Count	Pause Time		Count	Pause Time	
			Mean	Max		Mean	Max		Mean	Max
GC	41.2%	655	5.4	9.7	375	20.5	44.5	174	12.2	107.5
Manual+GC	13.0%	229	8.4	15.3	139	19.5	51.2	73	22.1	61.9

Table 4: Collections and GC pauses in the HavlakPool benchmark (1 thread).

Config	Heap	%GC	Gen0			Gen1			Gen2		
			#	Pause Time		#	Pause Time		#	Pause Time	
				Mean	Max		Mean	Max		Mean	Max
small	GC	51.9%	27	10.6	25.8	24	15.4	33.8	6	3.4	7.1
small	Man+GC	12.3%	5	18.7	27.5	4	24.9	34.6	2	50.2	69.5
medium	GC	57.5%	53	9.7	39.2	49	16.9	54.3	7	3.7	6.1
medium	Man+GC	17.4%	10	24.2	33.9	8	29.8	38.2	3	79.3	138.6
large	GC	59.0%	77	9.5	33.0	76	15.9	35.3	8	4.6	6.8
large	Man+GC	24.2%	15	37.1	122.6	12	35.3	45.6	3	82.9	149.4

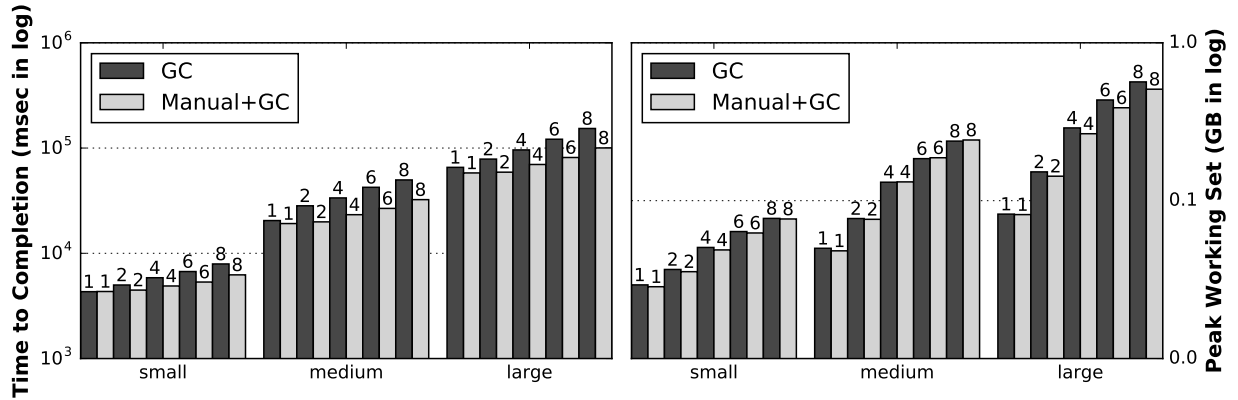


Figure 26: Comparison of time to completion (left hand side) and peak working set (right hand side) in the DirectedGraph micro-benchmark using Workstation GC. All results are in log scale.

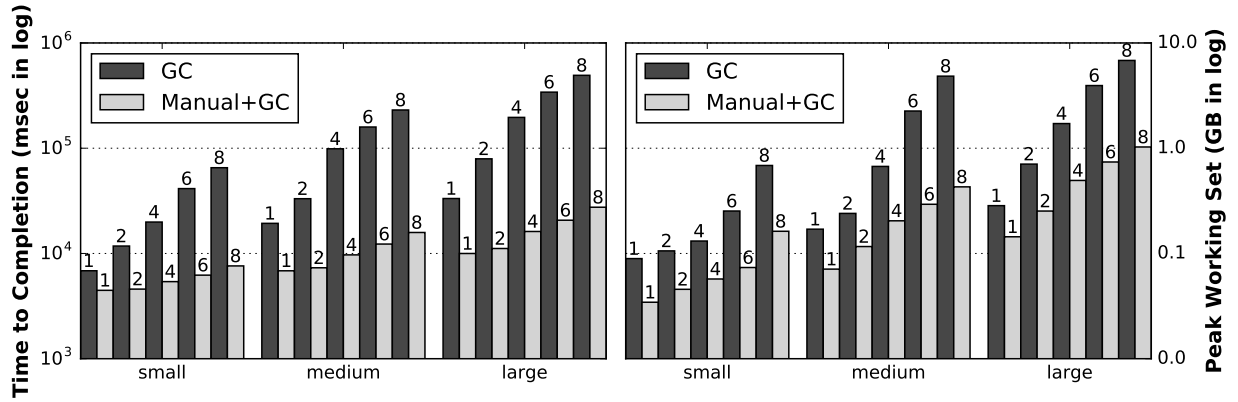


Figure 27: Comparison of time to completion (left hand side) and peak working set (right hand side) in the DirectedGraphReplace micro-benchmark using Workstation GC. All results are in log scale.

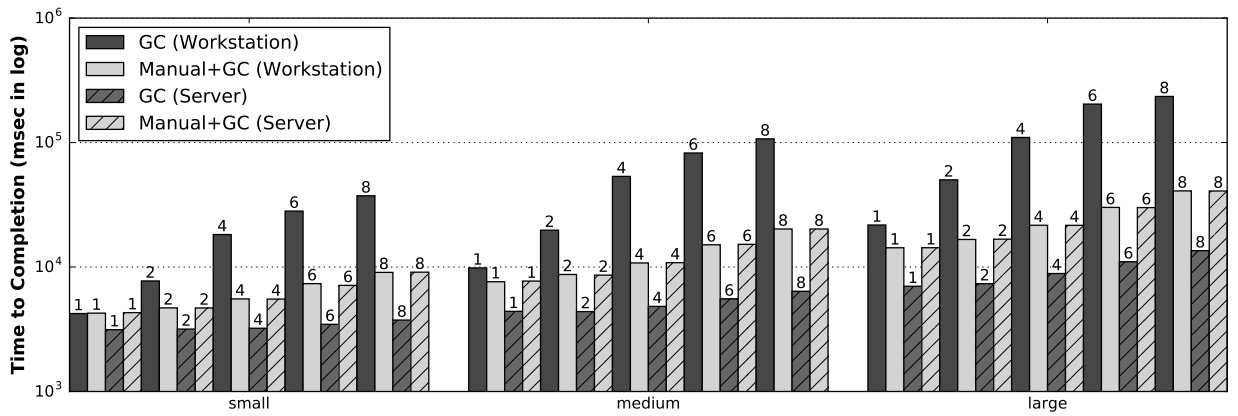


Figure 28: Comparison of time to completion (log scale) in the BinTree micro-benchmark using both Workstation and Server GC.

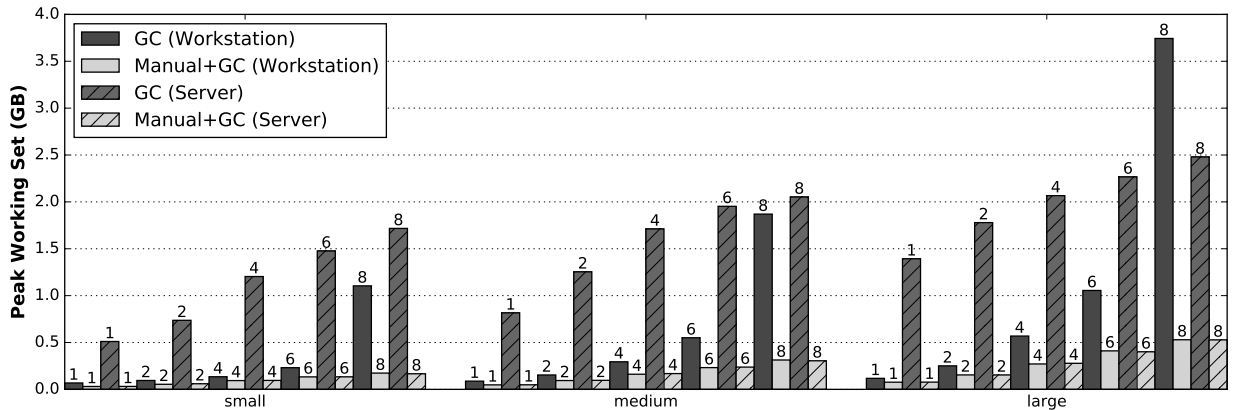


Figure 29: Comparison of peak working set in the BinTree micro-benchmark using both Workstation and Server GC.

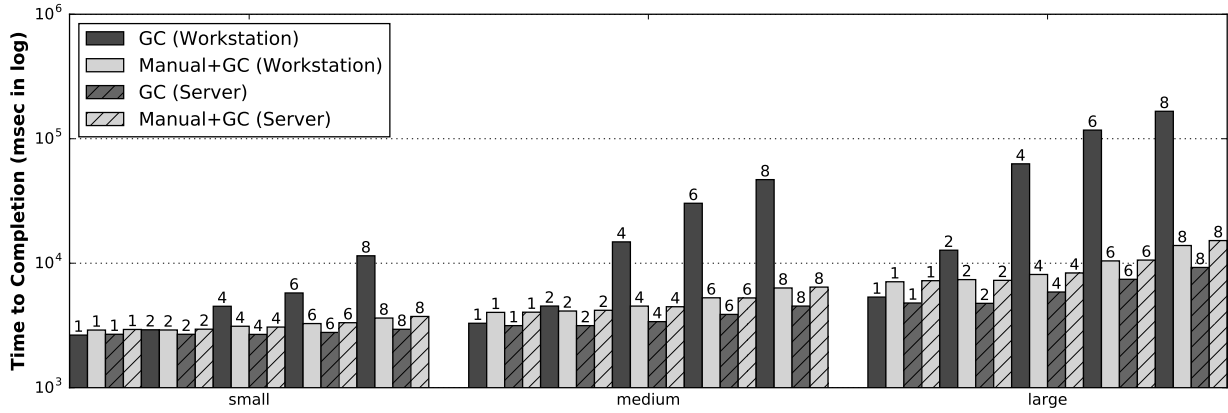


Figure 30: Comparison of time to completion (log scale) in the BinTreeLive micro-benchmark using both Workstation and Server GC.

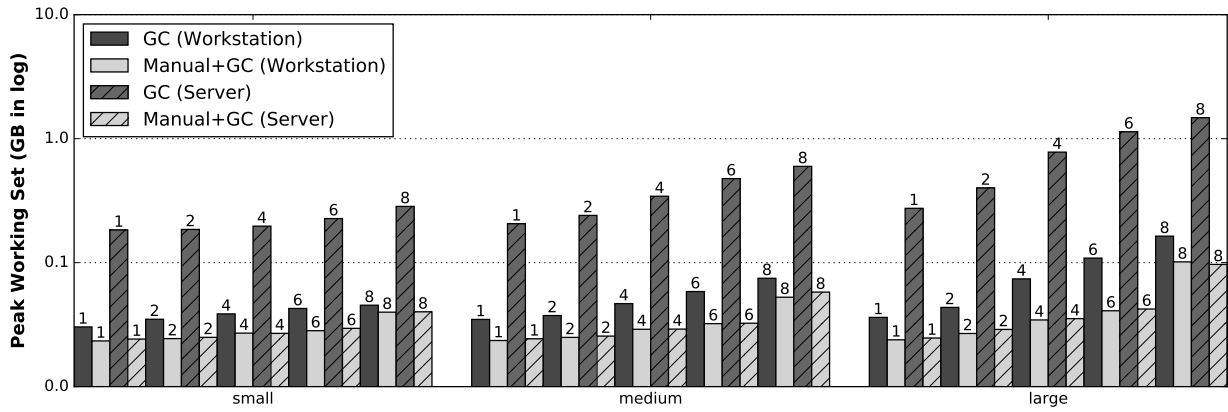


Figure 31: Comparison of peak working set (log scale) in the BinTreeLive micro-benchmark using both Workstation and Server GC.

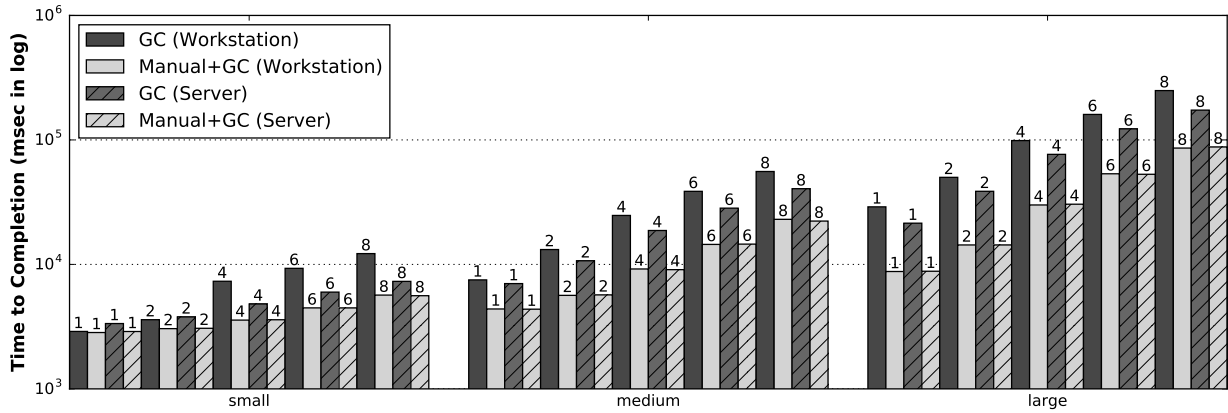


Figure 32: Comparison of time to completion (log scale) in the BinTreeGrow micro-benchmark using both Workstation and Server GC.

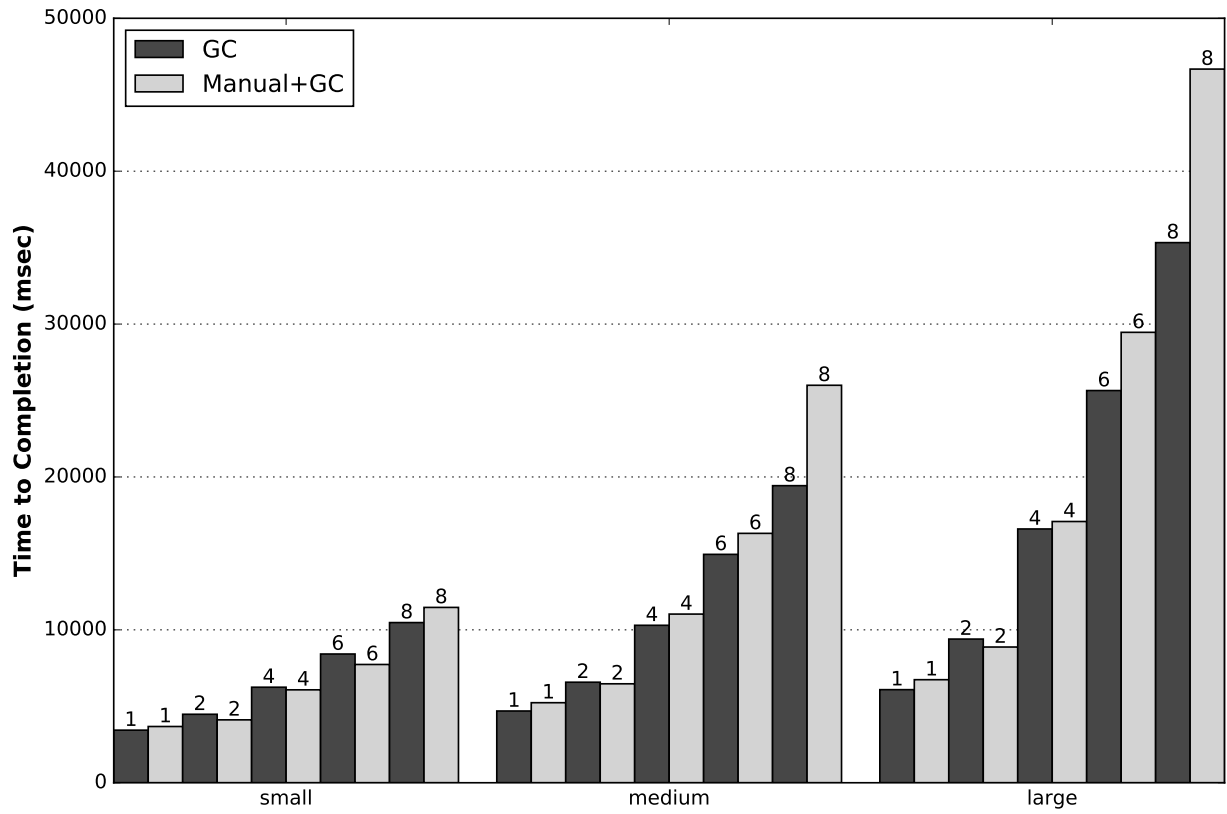


Figure 36: Comparison of end-to-end time (includes CFG construction and the Havlak loop finding algorithm) in the HavlakPool micro-benchmark using Workstation GC.

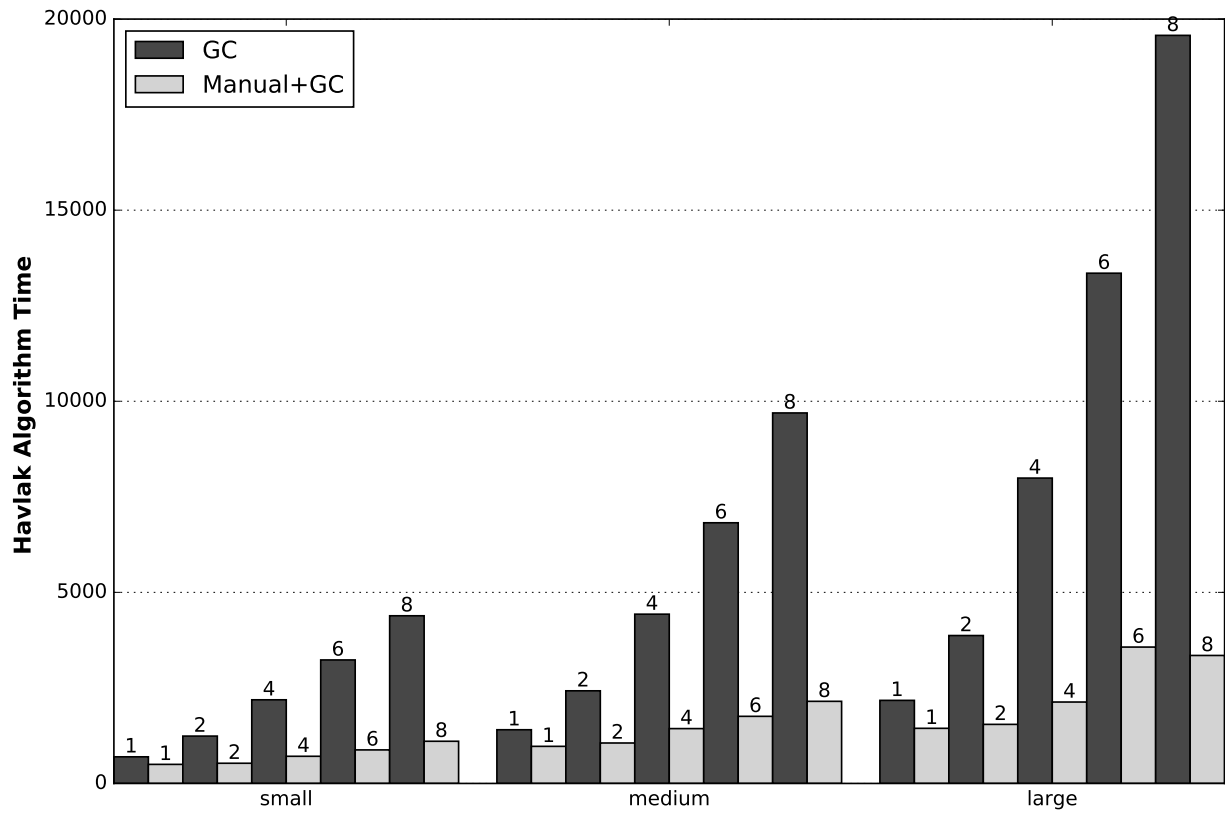


Figure 37: Comparison of time to complete the Havlak loop finding algorithm in the HavlakPool micro-benchmark using Workstation GC.

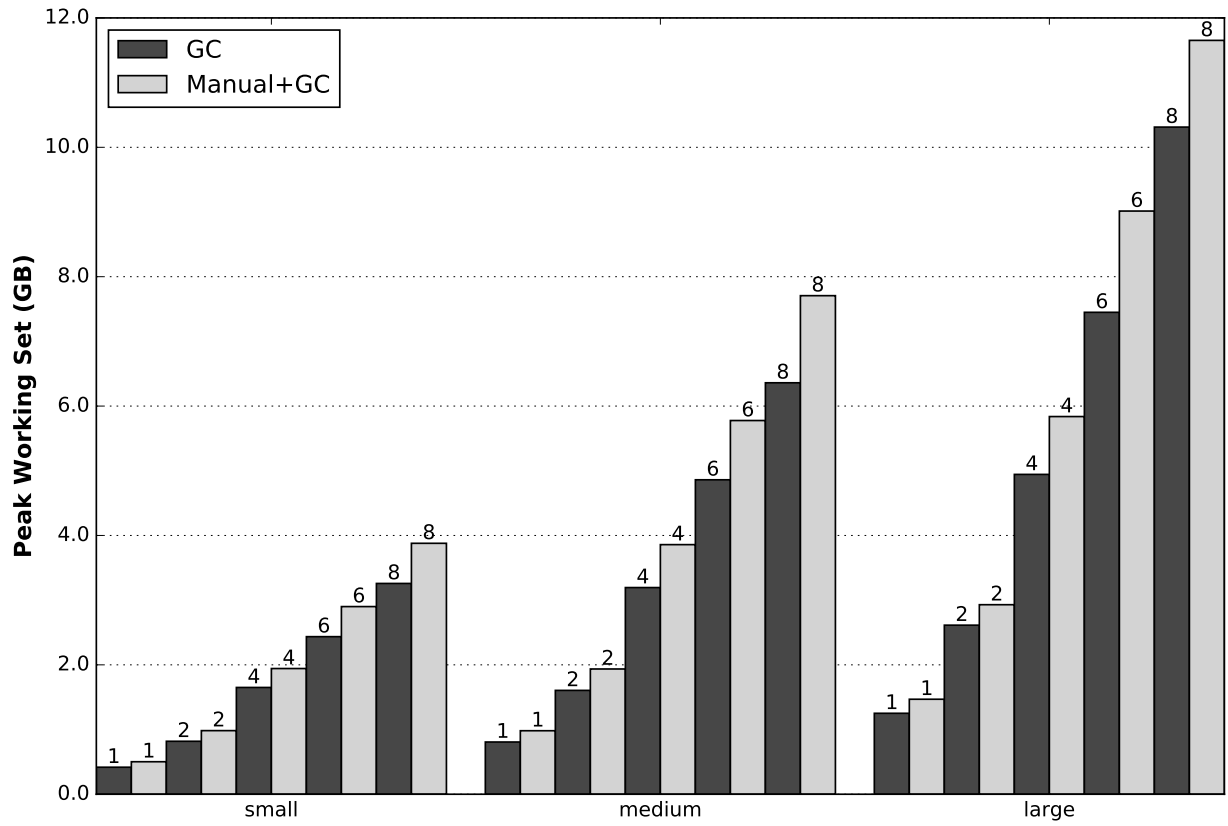


Figure 38: Comparison of peak working set in the Havlkak loop finding algorithm using Workstation GC.