

Will it Fit? Verifying Heap Space Bounds of Concurrent Programs under Garbage Collection with Separation Logic

ALEXANDRE MOINE, Inria, France

ARTHUR CHARGUÉRAUD, Inria & Université de Strasbourg, CNRS, ICube, France

FRANÇOIS POTTIER, Inria, France

We present IrisFit, a Separation Logic with space credits for reasoning about heap space in a concurrent call-by-value language equipped with tracing garbage collection and shared mutable state. Space credits, a purely logical device, are consumed when a heap block is allocated and recovered when a block becomes provably unreachable. For each allocated address, “pointed-by-heap” and “pointed-by-thread” assertions record which heap blocks point to this address and which threads hold this address as a root.

We point out a fundamental difficulty in the analysis of the worst-case heap space complexity of concurrent programs in the presence of tracing garbage collection: if garbage collection phases and program steps can be arbitrarily interleaved, then there exist undesirable scenarios where a root held by a sleeping thread prevents a possibly large amount of memory from being freed. This phenomenon leads to heap space complexity bounds that are worse than one might naively expect. For example, in an implementation of Treiber’s lock-free stack, one cannot prove that “pop frees up one list cell worth of heap space”.

To remedy this problem and eliminate such undesirable scenarios, we propose three language features, which are intended to be jointly exploited by the compiler and by the programmer, as follows. First, to guarantee that the size of the heap cannot exceed a certain limit, we make *memory allocation* a possibly *blocking* instruction. If a memory allocation instruction would cause the size of the heap to exceed the limit, then this instruction must wait for enough space to become available. Second, to ensure progress, we introduce *polling points*. If any thread is waiting for space, then a polling point is a blocking instruction. Polling points are intended to be automatically inserted by the compiler. By inserting sufficiently many of them, the compiler can ensure that a thread that is waiting for space cannot be forever starved. Third, to eliminate the aforementioned undesirable scenarios, where a thread reaches a polling point while holding certain roots, we let the programmer delimit *protected sections* where polling points must not be inserted.

Beyond this language design aspect, the heart of our contribution is IrisFit, a novel program logic that can establish worst-case heap space complexity bounds and whose reasoning rules can take advantage of the presence of protected sections. We construct IrisFit inside the Coq proof assistant on top of the Iris Separation Logic framework. We prove that IrisFit offers both a safety guarantee—programs cannot crash and cannot exceed a heap space limit—and a liveness guarantee—provided enough polling points have been inserted, every memory allocation request is eventually satisfied. We illustrate the use of IrisFit via a number of case studies, including a version of Treiber’s stack that is correctly decorated with protected sections.

CCS Concepts: • **Theory of computation** → **Separation logic; Program verification.**

Additional Key Words and Phrases: separation logic, tracing garbage collection, concurrency, program verification

ACM Reference Format:

Alexandre Moine, Arthur Charguéraud, and François Pottier. 2018. Will it Fit? Verifying Heap Space Bounds of Concurrent Programs under Garbage Collection with Separation Logic. 1, 1 (November 2018), 70 pages. <https://doi.org/XXXXXXX.XXXXXXX>

Authors’ addresses: [Alexandre Moine](mailto:alexandre.moine@inria.fr), Inria, Paris, France, alexandre.moine@inria.fr; [Arthur Charguéraud](mailto:arthur.chargueraud@inria.fr), Inria & Université de Strasbourg, CNRS, ICube, Strasbourg, France, arthur.chargueraud@inria.fr; [François Pottier](mailto:francois.pottier@inria.fr), Inria, Paris, France, francois.pottier@inria.fr.

2018. ACM XXXX-XXXX/2018/11-ART
<https://doi.org/XXXXXXX.XXXXXXX>

1 INTRODUCTION

Program Verification. The most common aim of program verification is to establish the *safety* and *functional correctness* of a program, that is, to prove that this program does not crash and computes a correct result. In the area of deductive program verification [Filiâtre 2011], a program is usually verified with the help of a *program logic*, that is, a set of deduction rules whose logical soundness has been demonstrated once and for all. Separation Logic [Reynolds 2002] and Concurrent Separation Logic [Brookes and O’Hearn 2016; O’Hearn 2019; Jung et al. 2018b] are examples of program logics that allow compositional reasoning (that is, reasoning about a program component in isolation) in the presence of challenging features such as dynamic memory allocation, mutable state, and shared-memory concurrency.

Verification of Resource Bounds. Beyond safety and functional correctness, it may be desirable to establish bounds on *resource consumption*, that is, to prove that the resource requirements of a program (or program component) do not exceed a certain bound. Indeed, a program that requires an unexpectedly large amount of *time* may be unresponsive. A program that requires an unexpectedly large amount of *stack space* may crash with a stack overflow. A program that requires an unexpectedly large amount of *heap space* may exhaust the available memory and make the system unstable.

Assuming that one is able to tell where in the code the resource of interest is consumed and produced, and how much of it is consumed or produced, reasoning about resource consumption can be reduced to reasoning about safety. To do so, one can construct a variant of the program that is instrumented with a *resource meter*, that is, a global variable whose value indicates what amount of the resource remains available. In this instrumented program, one places assertions that cause a runtime failure if the value of the meter becomes negative. If one can verify that the instrumented program is safe, then one has effectively established a bound on the resource consumption of the original program.

The principle of a resource meter has been exploited in many papers, using various frameworks for establishing safety. For instance, Crary and Weirich [2000] exploit a dependent type system; Aspinall et al. [2007] exploit a VDM-style program logic; Carbonneaux et al. [2015] exploit a Hoare logic; He et al. [2009] exploit Separation Logic. The manner in which one reasons about the value of the meter depends on the chosen framework. In the most straightforward approach, the value of the meter is explicitly described in the pre- and postcondition of every function. This is the case, for instance, in He et al.’s work [2009], where two distinct meters are used to measure stack space and heap space. In a more elaborate approach, which is made possible by Separation Logic, the meter is not regarded as an integer value but as a bag of *credits* that can be individually *owned*. The sum of all credits in circulation corresponds to the value of the meter. This removes the need to refer to the absolute value of the meter: instead, the specification of a function may indicate that this function requires a certain number of credits and produces a certain number of credits.

Verification of Heap Space Bounds, without Garbage Collection. A programming language that does not have garbage collection usually offers an explicit memory deallocation instruction. Thus, it is easy to tell where heap space is consumed and produced: an allocation instruction consumes the amount of space that it receives as an argument; a deallocation instruction recovers the space occupied by the heap block that is about to be deallocated.

In such a setting, traditional Separation Logic, extended with space credits, can be used to establish verified heap space bounds. Atkey [2011] presents a Separation Logic with an abstract resource that is consumed and produced by two distinct primitive operations, akin to space credits with explicit allocation and deallocation. Hofmann’s work on the typed programming language

LFPL [2000] can be viewed as a precursor of this idea: LFPL has explicit allocation and deallocation, which consume and produce values of a linear type, written \diamond , whose inhabitants behave very much like space credits.

Verification of Heap Space Bounds, with Garbage Collection. In the presence of garbage collection, how does one reason about heap space? In this setting, the programming language does not have a memory deallocation instruction. Thus, it is not evident at which program points space can be reclaimed. A tracing garbage collector (GC) can be invoked at arbitrary points in time, and may deallocate any subset of the *unreachable blocks*. An unreachable block is a block that is not *reachable* from any *root* via a *path* in the heap. Thus, reasoning about heap space in the presence of garbage collection requires reasoning about roots and unreachability.

Madiot and Pottier [2022] make a first step towards addressing this problem. They extend Separation Logic with several concepts. To keep track of free space, they use space credits. They view memory deallocation as a *logical operation*: it is up to the person who verifies the program to decide at which points this operation must be used and which memory blocks must be *logically deallocated*. This decision is subject to a proof obligation: a memory block can be logically deallocated only if it is unreachable. Unfortunately, the concept of unreachability is not local: that is, this concept cannot easily be expressed in terms of traditional Separation Logic assertions. Therefore, Madiot and Pottier rephrase this proof obligation as follows: a memory block can be logically deallocated if it has no predecessors and is not a root. To record the predecessors of every memory block, they use *pointed-by* assertions [Kassios and Kritikos 2013]. To record which blocks are roots, they focus their attention on a low-level language, where the stack is explicitly represented in the heap as a collection of “stack cells”. Then, a block is a root if and only if it is a stack cell.

In previous work [Moine et al. 2023], we scale Madiot and Pottier’s ideas up to a high-level language, where the stack is implicit. We introduce *Stackable* assertions to implicitly record which memory locations are “invisible roots”, that is, which memory locations are roots because they appear in some indirect caller’s stack frame.

Neither of these papers focuses on concurrency. Madiot and Pottier [2022] technically support concurrency, but only for a low-level language with stack variables explicitly allocated in the heap, and without any concurrent example covered. Moine et al. [2023] do not support it. By design, their *Stackable* assertion keeps track of a single stack. Extending it with support for multiple stacks is a priori not straightforward.

Verification of Heap Space Bounds, With Garbage Collection and Concurrency. In the present paper, we target a high-level programming language equipped with garbage collection and shared-memory concurrency. In such a setting, multiple threads run concurrently. They share a common heap; each thread has its own implicit stack.

Our initial aim in this project was to *propose a program logic* that allows its user to reason about heap space, and to verify heap space complexity bounds, in a concurrent setting. However, in the course of this work, we came to realize that, unless some care is taken, concurrent programs can have bad worst-case heap space complexity—that is, worse complexity than one might naively imagine. Thus, in addition to our initial aim, which was to let our program logic serve as a tool to *describe* the worst-case scenarios, we decided to also propose *programming language features* that let the programmer *eliminate* some of the worst-case scenarios.

In terms of programming language design, our proposed changes and additions are as follows. First, we impose a limit on the size of the heap. To ensure that this limit is never exceeded, we make *memory allocation* a possibly *blocking* instruction: if a memory allocation instruction would cause the size of the heap to exceed the limit, then this instruction must wait for enough space to become available. Second, we introduce *polling points*. If any thread is waiting for space, then

a polling point blocks the current thread. Polling points are intended to be automatically inserted by the compiler. By inserting sufficiently many of them, the compiler can ensure that a thread that is waiting for space cannot be forever starved. Third, to eliminate undesirable scenarios where a thread reaches a polling point while holding certain roots (which would prevent the garbage collector from reclaiming certain data structures), we let the programmer delimit *protected sections* where polling points must not be inserted. To ensure that a protected section terminates in bounded time, we forbid blocking instructions (that is, memory allocations and polling points) and function calls within protected sections. By decorating our implementation of Treiber’s lock-free stack with protected sections, we are able to prove that “pop frees up one list cell worth of heap space”, as desired (§3.3).

Regarding the heap size limit, we propose two variants of the semantics. Our *default* semantics uses a fixed limit, which must be set before the program is executed. Our *growing* semantics lets this limit grow at runtime under certain conditions. The default semantics is slightly easier to understand, whereas the growing semantics is more realistic, as it does not require advance knowledge of a suitable limit.

We propose a program logic that lets users establish worst-case heap space complexity bounds. This program logic is independent of which variant of the semantics is chosen. It is sound, in a certain sense, with respect to each variant. Indeed, under each variant, we are able to establish a safety property: the heap size of a verified program cannot exceed a certain bound. Furthermore, under both variants, and under the assumption that enough polling points have been inserted, we establish a liveness property: no thread remains forever blocked. Our program logic includes reasoning rules that exploit the presence of protected sections to establish *improved* heap space complexity bounds: this is illustrated by the example of Treiber’s stack (§3.3, §11.5).

Although our program logic does include reasoning rules for polling points, these rules are not needed by the end user, who reasons about a source program in which polling points have not yet been inserted. This source program can be viewed as the description of a family of instrumented programs, which are obtained from the source program by inserting polling points in arbitrary places, outside of protected sections. If the source program can be verified (that is, if it abides by the rules of the program logic) then every instrumented program in this family can also be verified. Thus, the correctness and worst-case heap space complexity properties that are established by the user about the source program still apply after polling points are inserted by the compiler.

Contributions. The main contributions of this paper are the following:

- We present LambdaFit (§2, §4), an imperative language with shared-memory concurrency and tracing garbage collection. The novel aspects of LambdaFit include a heap size limit, blocking allocations, polling points, and protected sections.
- We introduce IrisFit (§5, §6, §10), a Separation Logic that allows establishing *safety*, *functional correctness*, and *worst-case heap space complexity* properties of concurrent programs, in the presence of garbage collection, and allows *compositional reasoning*.
- We prove the soundness of IrisFit (§8). More specifically, we establish two results. A *safety* theorem guarantees that a verified program cannot crash and that its heap size cannot exceed a certain bound. A *liveness* theorem guarantees that after enough polling points have been (automatically) inserted, no thread can be forever blocked by a memory allocation request. In fact, we prove a slightly stronger liveness statement: at all times, *in a bounded number of steps*, the system must reach a configuration where no thread is blocked by a memory allocation request. We establish this result *without* formulating a fairness hypothesis. We prove safety and liveness theorems for both variants of the semantics of LambdaFit, namely the fixed-limit variant and the growing-limit variant.

- We prove that IrisFit is sound also with respect to an *oblivious* semantics of LambdaFit, where there is no heap size limit, a memory allocation instruction is never blocked, and polling points have no effect. In such a setting, our safety theorem (§8.4) guarantees that a verified program cannot crash and that (when every thread is outside of a protected section) its *live* heap size cannot exceed a certain bound. There is no need in this setting for a liveness theorem. This result emphasizes that IrisFit is largely independent of the non-standard features of LambdaFit.
- We encode *closures* in LambdaFit and show how to reason about them with IrisFit (§9). Compared with our previous paper [Moine et al. 2023], we propose an improved treatment of closures: the *Spec* predicate, which describes the behavior of a closure, is persistent.
- We verify several case studies (§11), namely: an implementation of “fetch-and-add” as a CAS loop; a concurrent counter that is encapsulated as a pair of closures; a library for async/finish parallelism; and Treiber’s lock-free stack [1986]. This gallery of challenging examples illustrates the expressive power of IrisFit.

All of our results, including the validity of our reasoning rules, our soundness theorems, and our case studies, are mechanized using the Coq proof assistant and the Iris framework [Jung et al. 2018b]. For details, we refer the reader to our mechanization [Moine 2024b] and to the first author’s dissertation [Moine 2024a].

Because we wish to make the present paper self-contained, we borrow some text from our previous paper [Moine et al. 2023]. The re-used material amounts to roughly 8 pages in total. The main re-used passages are the beginning of this introduction, the design and explanation of the pointed-by-heap assertion (§5.6), the discussion and definition of the closure macros (§2.7, §9.2), the concept and presentation of triples with souvenir (§10), and part of the discussion of the related work (§12.3, §12.4, §12.5).

2 OVERVIEW

LambdaFit is a call-by-value language with dynamic memory allocation, mutable state, shared-memory concurrency, and tracing garbage collection. Its syntax and semantics are standard, save for a few original aspects.

First, LambdaFit is restricted to closed functions, also known as *code pointers*. We encode closures as heap-allocated objects that store code and data (§2.7).

Second, LambdaFit exhibits several non-standard features related with memory management. First, a memory allocation instruction that would cause the size of the heap to exceed a certain limit is blocked. Second, LambdaFit’s syntax includes three non-standard instructions, namely an instruction that represents a *polling point* and two instructions that mark the beginning and end of a *protected section*.

2.1 One Language, Several Semantics

In this paper, we define and use three semantics, or three variants of the semantics, of LambdaFit. Let us explain why each variant exists and what role it plays in the paper.

- (1) The most important semantics is the *default semantics*. This semantics involves the concept of a *root* and has explicit garbage collection steps (§2.2). It imposes a *fixed limit* S on the size of the heap. This limit is a parameter of the semantics. A “large” memory allocation request, which would cause this limit to be exceeded, is blocked (§2.3): we say that it is “waiting for space”. Furthermore, if any thread is waiting for space, then a polling point is blocked (§2.4). Protected sections (§2.5) serve mainly to delimit areas where polling points must not be inserted by the compiler. In the semantics, they play a minor role: within a protected section,

polling points, memory allocation instructions, “fork” instructions, and function calls are forbidden—they cause a crash.

This semantics (as well as its variants, discussed next) is non-deterministic. This is due not only to concurrency but also to the fact that garbage collection can take place at an undetermined point in time and can deallocate an undetermined number of unreachable blocks. Our soundness theorems guarantee that, in *every possible execution*, the heap space usage of a verified program respects the bound S . Thus, the non-determinism that is inherent in the semantics makes our soundness theorems stronger. If a specific implementation of LambdaFit exhibits fewer behaviors than permitted by the semantics, then our soundness theorems still hold for this implementation.

A potential criticism that one might formulate about the default semantics is the fact that it requires setting the heap limit to a suitable value S before the program is executed. If the program has been fully verified using IrisFit, then a suitable value of S may be known. However, it is more realistic to expect that only part of the program has been verified. Furthermore, if the program receives data from the outside (via the file system or the network) then its space requirement may depend on the input that it receives. For these reasons, it seems worthwhile to design and discuss also a variant of the default semantics that is equipped with a *growing limit* on the size of the heap.

- (2) In the *growing semantics*, as in the default semantics, there is a limit on the size of the heap, which serves to identify and block “large” memory allocation requests. However, in the growing semantics, the limit is not fixed. Instead, at runtime, if the current limit is found to be too low, then the limit is increased. Such an increase can take place only while no thread is inside a protected section.

A similar strategy for growing the size of the heap is used in real-world runtime systems, such as the OCaml runtime system [Madhavapeddy and Minsky 2022, §25.4.1] and the Haskell runtime system [Marlow et al. 2008, §5].

Our discussion of the growing semantics occupies only two sections of the paper (§4.2.10, §8.3). However, we believe that this discussion plays an essential role in arguing that a *practical* runtime system that obeys the design of LambdaFit can be implemented.

Finally, we define an *oblivious semantics*, which ignores the non-standard features of LambdaFit: in this semantics, memory allocation is never blocked, and polling points have no effect. This semantics serves two purposes. First, it enables us to state the *core soundness* guarantee that is offered by IrisFit in the absence of the non-standard features of LambdaFit. Second, it plays a key technical role. The soundness properties of IrisFit with respect to the default and growing semantics are obtained as corollaries of its core soundness property with respect to the oblivious semantics.

- (3) In the oblivious semantics, no instruction is ever blocked. This semantics does not impose a limit on the size of the heap. Therefore, it does not need and does not have garbage collection: an unreachable block remains allocated forever. In this semantics, one cannot bound the total size of the heap (including garbage). Still, under the assumption that every thread is outside a protected section, one can bound the live heap size, that is, the size of the reachable fragment of the heap.

Our discussion of the oblivious semantics is limited to two sections of the paper (§4.2.7, §8.4).

We wish to reassure a reader who might be worried about an apparent proliferation of semantics. First, the three semantics share many auxiliary definitions, so there is limited redundancy. Second, throughout the paper, we focus on the default semantics. Where one of the other two semantics is discussed, this is explicitly indicated. Last but not least, we propose just one program logic, IrisFit, whose reasoning rules (§2.6) are independent of which semantics is considered. For each of the three semantics, we are able to prove that if a program has been verified using IrisFit then its heap space usage is bounded in a certain sense (§8).

2.2 Roots and Garbage Collection

Garbage collection [Jones and Lins 1996] deallocates some or all unreachable memory blocks, where a block is *reachable* if there exists a path from some *root*, through the heap, to this block. In the oblivious semantics of LambdaFit (§2.1), garbage collection is not explicitly modeled; nevertheless, one can define the worst-case heap space complexity of a program as the maximum value (over all possible executions) of its live heap size, where (at a given point in time) the *live heap size* is the sum of the sizes of all reachable heap blocks. In the default semantics of LambdaFit (§2.1), garbage collection is explicitly modeled: a garbage collection step deallocates an arbitrary number of unreachable heap blocks. In either approach, it is necessary to answer to the question: what is a root?

How can the intuitive concept of a root be formally defined in the setting of a small-step, substitution-based operational semantics? Before addressing this question, let us recall a few fundamental aspects of such a semantics. In an *operational* semantics, a program state, which represents the state of a running program, is a syntactic object. Here, because we are interested in concurrent programs with dynamic memory allocation, a program state includes a thread pool (a list of threads) and a heap (a finite map of memory locations to memory blocks). In a *small-step* semantics, the manner in which the program state evolves over time is described by a reduction relation, that is, a binary relation on program states. In a *substitution-based* semantics, within the thread pool, each running thread is represented as a closed term, that is, a term without free variables. The reduction rules ensure that, whenever the scope of a variable is entered, a closed value is substituted for this variable. Thus, a closed term that represents a running thread describes both the code that this thread is about to execute and the data to which this thread has access. In particular, a memory location ℓ is a closed value, and a closed term that represents a running thread can contain memory locations.

In such a setting, what is a root? A simple, commonly agreed-upon answer is: *a root is a memory location ℓ that appears in at least one running thread t* . By this, we mean that the closed term t , which represents one of the currently running threads, literally contains one or more occurrences of the memory location ℓ .

This convention is known as the *free variable rule* (FVR) [Felleisen and Hieb 1992; Morrisett et al. 1995]. Intuitively, the FVR makes sense because the (computable) set of memory blocks that are reachable from the locations that the program knows about is a conservative approximation of the (uncomputable) set of memory blocks that might be accessed in the future by the program. However, one must keep in mind that the FVR is not a static approximation of the dynamic semantics. Instead, the FVR is *part of* the definition of the dynamic semantics. It defines the concept of root, which in turn is used to define reachability and garbage collection.

The reader may wonder whether real-world programming languages respect the FVR. As far as we know, many real-world implementations of garbage-collected languages, such as OCaml, SML, Haskell, Scala, Java, and more, are meant to respect the FVR. Unfortunately, this intention is often undocumented. A prominent example of a compiler that explicitly respects the FVR is the CakeML verified compiler. Gómez-Londoño et al. [2020] and Gómez-Londoño and Myreen [2021] prove that the CakeML compiler respects a cost model that is defined at the level of the intermediate language DataLang and that includes a form of the FVR.

2.3 Why Block Large Memory Allocation Requests

In the presence of tracing garbage collection, two measures of the size of the heap must be distinguished, namely the allocated heap size and the live heap size. The *allocated heap size*, or simply *heap size*, is the sum of the sizes of all allocated heap blocks. The *live heap size* is the sum

of the sizes of all *reachable* heap blocks. In other words, it is the allocated heap size minus the sizes of the unreachable blocks.

Our main goal in this paper is to bound the heap size¹ in a setting where LambdaFit is equipped with its *default* semantics (§2.1).

The default semantics is parameterized by a limit S and is designed in such a way that the heap size always remains less than or equal to S . This property, which is stated by Lemma 4.2 (§4.2.9), is enforced as follows. Let us say that a memory allocation request is *large* if it would cause the heap size to exceed S , that is, if the sum of the current heap size and the number of requested words exceeds S . Otherwise, let us say that the allocation is *small*. Then, a large memory allocation instruction is not allowed to proceed: it is *blocked*. Once garbage collection takes place and is able to free enough space in the heap, this memory allocation instruction may become small, therefore unblocked.

This explains one aspect of the design of the default semantics: by blocking large memory allocation instructions, we ensure that one kind of undesirable behavior, namely *growing the heap too large*, is eliminated a priori. Two kinds of undesirable behavior remain permitted by the default semantics, namely *crashes* and *deadlocks*: a thread can crash or become forever blocked. Under certain assumptions about the placement of polling points, our program logic statically guarantees that these undesirable behaviors cannot arise: this is stated by our *safety* and *liveness* theorems (Theorems 8.1 and 8.2).

An alternative point of view is offered by the *oblivious* semantics of LambdaFit (§2.1). In that semantics, there is no limit on the size of the heap, and no instruction is ever blocked. Thus, a different kind of undesirable behavior, namely *deadlocks*, is eliminated a priori. The undesirable behaviors that remain possible are *crashes* and *growing the heap too large*. With respect to the oblivious semantics, our program logic provides two guarantees: first, no thread can crash; second, if the program has been verified under the hypothesis that S space credits are initially granted, then, in every possible execution, provided every thread is currently outside a protected section, the current live heap size is at most S . These guarantees are stated in our *core soundness* theorem (Theorem 8.4). We emphasize that S is not a parameter of the oblivious semantics; it appears only in the statement of the core soundness theorem.

The points of view offered by the default semantics and by the oblivious semantics are subtly different; we believe that it is enlightening to understand them both. Furthermore, there is a technical connection between them: we first establish the core soundness theorem (Theorem 8.4), then use this theorem as a stepping stone in the proof of Theorems 8.1 and 8.2.

2.4 Polling Points

The existence of blocking memory allocation instructions endangers *liveness*: for some programs, there exist adversarial schedules where a large memory allocation request is forever blocked. This can happen, for example, if some other threads are always scheduled, so the garbage collector is never allowed to run. More specifically, imagine that thread A is blocked by a large memory allocation request, while thread B is in an infinite loop. Then, the scheduler can choose to always execute thread B and never execute the garbage collector, so thread A remains forever blocked. We wish to forbid this scenario and to formally establish a liveness guarantee of the form: *always*,

¹By adopting a measure whose definition is a sum of the sizes of the heap blocks that the program creates, we restrict our attention to a *logical* notion of heap size, that is, a notion that depends only on the program and on the semantics of the programming language. One might instead wish to control the *physical* heap size, that is, how much memory the runtime system requests from the operating system. This would require knowing which garbage collection technique is used, whether the garbage collector is able to perform compaction (so as to eliminate fragmentation), and so on. Our analysis does not require this information and does not bound the physical heap size.

eventually, every thread can make progress (Theorem 8.2). Furthermore, we wish to do so without formulating a fairness hypothesis about the scheduler.

To this end, we equip LambdaFit with *polling points*. A polling point is a synchronization instruction, a form of barrier. A thread may proceed past a polling point only if no large memory allocation request is currently outstanding. In other words, if any thread is currently waiting for space, then no thread can move past a polling point.

By inserting sufficiently many polling points into a program, one can ensure that every memory allocation request is eventually satisfied. Indeed, by inserting enough polling points, one can enforce the following property: as soon as one thread is blocked on a large memory allocation request, every thread must eventually reach a polling point or a large memory allocation request, where it, too, becomes blocked. At this point, no thread can make progress, so garbage collection must take place. Once garbage collection has made enough space available—which our program logic statically guarantees is possible!—all outstanding memory allocation requests can be satisfied.

In the example scenario that was outlined above, a polling point must be inserted in the infinite loop of thread *B*. Then, this thread must eventually reach a polling point, where it becomes blocked. The only permitted step is then a garbage collection step, which is expected to free up enough memory to satisfy thread *A*'s large allocation request. Consequently, the two threads become unblocked.

In principle, polling points could be manually inserted by the programmer, but that would be tedious. In practice, we expect a compiler to automatically insert polling points where needed. In §8.2, we prove that a particular polling point insertion strategy, inspired by that of the OCaml 5 compiler, does indeed insert enough polling points to guarantee liveness.

2.5 Protected Sections

Polling points are meant to be inserted by the compiler in arbitrary places. However, inserting a polling point in a badly-chosen place can impact (that is, increase) the worst-case heap space complexity of a data structure: this is illustrated later on by the example of Treiber's stack (§3). To remedy this problem, we equip LambdaFit with *protected sections*, that is, sections of the code in which the compiler must not insert polling points. We require protected sections to be “short”, that is, to terminate in bounded time. Thus, when one thread is blocked by a large memory allocation request (§2.3), every other thread must eventually exit its protected section (if it is currently within such a section) and reach a polling point.

Protected sections provide the programmer with a means of controlling where polling points should *not* be placed. In comparison with manually placing polling points, this approach is not only more concise, but also more informative, as it lets the programmer explicitly document the places where a thread must not be blocked. The intuitive reason why a thread must not be blocked (in a specific place) is that it must instead be allowed to make progress and release some of its roots, thereby allowing the garbage collector to reclaim more space.

A protected section is explicitly delimited by two special instructions, *enter* and *exit*, which mark the beginning and end of the section. A single well-balanced construct “protected $\{t\}$ ” would be insufficiently flexible, because a protected section typically has one entry point and multiple exit points. This is illustrated by the example of Treiber's stack (Figure 3).

Protected sections are subject to two restrictions. First, they cannot be nested. Second, a protected section must not contain a memory allocation instruction, a “fork” instruction,² a polling point (§2.4),

²In our operational semantics, “fork” does not allocate any memory in the heap. We could technically allow “fork” inside a protected section without breaking any of our results. In the real world, though, “fork” is likely to allocate memory. Because we forbid memory allocation inside a protected section, it seems natural to disallow “fork” inside protected sections as well.

or a function call.³ These restrictions ensure that a protected section cannot contain a blocking instruction and can be exited in a bounded number of steps. The syntax of LambdaFit does not enforce these restrictions; however, violating them causes a runtime error, and is statically forbidden by our program logic.

The reader may wonder whether protected sections are ghost code, that is, whether (after polling points have been inserted) protected sections can be erased. In the default semantics (§2.1), this is essentially true: provided protected sections are correctly used (that is, enter and exit instructions are well-balanced, and protected sections do not contain memory allocations, polling points, “fork” instructions, or function calls), they have no semantic effect. If such correct usage is statically enforced, then at runtime protected sections can be erased. In the growing semantics (§2.1), however, protected sections do play a role and cannot be erased. Indeed, the heap size limit can be increased *only when all threads are outside protected sections*. Without this restriction, we would not be able to establish Theorem 8.3, which states that the heap size remains bounded.

Earlier (§1), we have pointed out that a source program describes a family of instrumented programs, which are obtained from the source program by inserting polling points in arbitrary places outside of protected sections. By applying IrisFit to just the source program, the user obtains a heap space bound that holds regardless of where polling points are inserted; in other words, the user obtains a guarantee about the worst-case heap space complexity of every program in this family. In terms of space complexity, decorating the source program with protected sections is always beneficial. Indeed, adding or extending protected sections reduces the family of programs that the source program represents: this can reduce the worst-case heap space complexity of this family, and cannot increase it. This phenomenon is illustrated by the example of Treiber’s stack (§3).

On the flip side, in practice, protected sections can have a cost. While they can improve the worst-case heap space complexity of a program family, they can also increase *latency*. Indeed, larger protected sections potentially imply larger gaps between polling points, therefore a longer time until a thread that is waiting for space can be unblocked. In practice, it seems wise to use as few protected sections as possible and to keep them as short as possible.

2.6 A Concurrent Separation Logic for Heap Space

This paper presents IrisFit, a concurrent Separation Logic for LambdaFit. IrisFit shares many features with pre-existing Separation Logics. The behavior of a program fragment is described by a *triple*, an assertion whose parameters include a precondition (an assertion that describes the initial state), the program fragment of interest, and a postcondition (an assertion that describes the final state). In IrisFit, a triple also includes a thread identifier, as the logic assigns a unique name to each thread. A rich vocabulary of logical connectives, including *points-to* assertions, *separating conjunction*, and many more, is used to construct assertions, which encode both *knowledge* of the current state and *permission to update* this state in certain ways.

What sets IrisFit apart from traditional Separation Logics? IrisFit borrows ideas from previous Separation Logics equipped with support for reasoning about heap space in the presence of garbage collection [Madiot and Pottier 2022; Moine et al. 2023] and scales them up to a concurrent setting. *Space credits* keep track of available space and serve as permissions to allocate memory. Furthermore, several kinds of assertions record which memory locations are reachable and in what way they can be reached. *Pointed-by-heap* assertions [Madiot and Pottier 2022] keep track of predecessors of each location in the heap. *Pointed-by-thread* assertions (new in this paper) keep track of the threads in which each location is a root. Like previous logics [Madiot and Pottier 2022; Moine et al.

³Because loops are encoded as recursive functions, forbidding function calls inside protected sections also forbids loops inside protected sections.

2023], IrisFit features a *ghost deallocation rule*. Because the programming language does not have an explicit memory deallocation instruction, it is up to the user of the logic to decide where to apply this rule. This rule requires proof that the memory block of interest is unreachable. This proof takes the form of pointed-by-heap and pointed-by-thread assertions, which are consumed; space credits are produced in their stead. A novelty of this paper is that logical deallocation *does not require or consume the points-to assertion*.

A crucial novel aspect of IrisFit is its ability to take advantage of protected sections while reasoning. Indeed, IrisFit offers a relaxed way of keeping track of roots inside protected sections. Ordinarily, pointed-by-thread assertions record which locations are roots, and as long as a location is a root, this location cannot be logically deallocated. Inside a protected section, however, an exception to this regime is made: the logic keeps track of a set of *temporary* roots. The user can turn an ordinary root into a temporary root (and vice-versa). The logic requires that, by the time the protected section ends, no temporary roots remain. Thus, by that time, every temporary root must no longer be a root (or must have been turned back into an ordinary root). Crucially, inside a protected section, the condition under which logical deallocation is permitted is: *if a location ℓ is not an ordinary root in any thread, and if ℓ has no live heap predecessors, then it can be logically deallocated*. In other words, logical deallocation is oblivious to the existence of temporary roots.⁴ Finally, perhaps surprisingly, because the points-to assertion survives logical deallocation and enables read and write access, *a temporary root that has already been logically deallocated can still be accessed* before the protected section ends. This pattern appears while verifying lock-free data structures (§11.5).

This overview of IrisFit may raise two questions about our approach. Why bother with garbage collection? Assuming that a program has been verified using IrisFit, could one replace the logical deallocation points identified by the proof with actual calls to free, thereby removing the need for garbage collection? We answer these two questions separately.

First, garbage collection is a widely used memory management technique. It is used, among other examples, in Scheme, Java, Scala, Haskell, OCaml, C#, JavaScript, and Go. There are arguably strong reasons why a programmer might choose a programming language equipped with garbage collection. These reasons include simplicity (garbage collection enables an elegant, high-level programming style), safety (garbage collection removes a class of runtime errors, including double-free and use-after-free errors), and performance (garbage collection allows bulk deallocation, which in certain circumstances can be more efficient than individual object deallocation). That said, garbage collection does make space usage analysis more difficult. IrisFit is the first program logic that allows such an analysis in the presence of challenging features such as mutable state and concurrency. Moreover, we believe that, in the future, IrisFit is can serve as a logical foundation for automated space complexity analyses (§13).

Second, although in simpler settings [Madiot and Pottier 2022; Moine et al. 2023] it could make sense to transform logical deallocation points into physical memory deallocation instructions, IrisFit is a more advanced logic, where such a transformation can be unsound or impossible. Indeed, as explained earlier in this section, IrisFit allows for logically deallocating a block *in advance* inside a protected section, even if this block is still in use within this protected section. In such a situation, transforming a logical deallocation operation into a physical free instruction would introduce a use-after-free error. Furthermore, IrisFit allows one thread to logically deallocate a block whose address is available (at runtime) only to some other thread. This idiom, which we call

⁴Should space become scarce, in the worst-case scenario, the garbage collector will be invoked at a time where every thread is blocked at a memory allocation instruction or at a polling point. In such a situation, every thread is outside a protected section, so there are no temporary roots. Thus, the existence of temporary roots inside protected sections has no impact on the worst-case scenario.

logical deallocation by proxy, appears in our analysis of Treiber’s stack (§11.5). In such a situation, transforming logical deallocation into a free instruction is impossible, because the address of the object that must be deallocated is not at hand! In conclusion, we remark that, in concurrent code, placing free instructions in a correct and optimal way is known to be an extremely hard problem, which IrisFit does not solve. This problem is the *raison d’être* of safe memory reclamation (SMR) schemes (§12.6) and one key argument in favor of general-purpose concurrent garbage collectors, which subsume SMR schemes.

2.7 Closures

To model the space complexity of programs that involve closures [Landin 1964; Appel 1992], we must somehow reflect the fact that a closure is a heap-allocated object. It has an address, a size, and may hold pointers to other objects. Thus, a closure has both direct and indirect impacts on space complexity: it occupies some space; and, by pointing to other objects, it keeps these objects live (reachable), preventing the GC from reclaiming the space that they occupy.

Thus, we cannot use the standard small-step, substitution-based semantics of the λ -calculus, where a λ -abstraction is a value that does not have an address or a size. Instead, two approaches come to mind. One approach is to view a λ -abstraction as a primitive expression (not a value) whose evaluation causes the allocation of a closure. Another approach is to adopt a restricted calculus that offers only closed functions (as opposed to λ -abstractions with free variables) and to *define* closure construction and closure invocation as *macros*, or canned sequences of instructions, on top of this calculus. As shown by Paraskevopoulou and Appel [2019], these two approaches yield the same space cost model. Furthermore, provided suitable syntax is chosen, the end user does not see the difference: it is just a matter of presentation in the metatheory.

We choose the second approach, because we find it simpler. In so doing, we follow Gómez-Londoño et al. [2020], who define the CakeML cost model at the level of DataLang, the language that serves as the target of closure conversion.

Thus, we equip LambdaFit with *closed functions*, which we also refer to as *code pointers*. We write $\mu_{\text{ptr}.f}. \lambda \vec{x}. t$ for a (recursive, multi-argument) closed function, and write $(v \vec{u})_{\text{ptr}}$ for the invocation of the code pointer v with arguments \vec{u} . LambdaFit does not have primitive closures. This allows us to present a program logic for LambdaFit and to establish the soundness of this logic without worrying about closures. Once this is done, we define *closure construction* $\mu_{\text{clo}.f}. \lambda \vec{x}. t$ and *closure invocation* $(\ell \vec{u})_{\text{clo}}$ as macros, and we extend our program logic with high-level reasoning rules for closures (§9). This allows end users to reason about these macros without expanding them and without even knowing how they are defined. In summary, LambdaFit can macro-express closures, and our logic allows reasoning about closures in the same way as if they were primitive constructs.

Our construction of closures as macros is the same as in our previous paper [Moine et al. 2023]. Our treatment of closures in the logic, however, has been generalized to multiple threads and simplified by describing closures via persistent predicates (§9).

3 WHY TREIBER’S STACK NEEDS PROTECTED SECTIONS

To illustrate how protected sections can lead to tighter space bounds, we use the example of Treiber’s stack, a lock-free, linearizable stack [Treiber 1986]. We first present a naive implementation of this data structure without protected sections (§3.1). We point out that this implementation has unsatisfactory worst-case heap space complexity: there are scenarios where a successful pop operation does not allow any memory cell to be freed (§3.2). All memory *can* eventually be recovered, but this may require waiting until all threads have completed their interaction with the stack. This situation is unpleasant: pop cannot be given a simple logical specification of the form “a successful pop frees up one list cell worth of heap space”. We show that, by annotating the

```

1  let create () = ref nil
2
3  let rec push s v =
4    let h = !s in
5    let h' = new_cell () in
6    set_data h' v;
7    set_tail h' h;
8    if compare_and_swap s h h'
9    then ()
10   else push s v
11
12 let rec pop s =
13   let h = !s in
14   if is_nil h
15   then pop s
16   else
17     let h' = tail h in
18     if compare_and_swap s h h'
19     then data h
20     else pop s

```

Fig. 1. An unsafe-for-space implementation of Treiber’s stack

code with protected sections, one can eliminate these undesirable scenarios and obtain the desired specification (§3.3). Near the end of this paper (§11.5), we present the details of how we formally establish this specification in IrisFit.

3.1 Naive Implementation of Treiber’s Stack

Treiber’s stack is implemented as a mutable reference to an immutable linked list, whose head corresponds to the top of the stack. Pseudo-code is presented in Figure 1.

The function call `create()` creates a new stack, represented as a fresh reference to an empty list `nil`. The `nil` value takes up no heap space: it is in fact an integer value.

The functions `push` and `pop` make crucial use of the atomic *compare-and-swap* (CAS) instruction. Each of them is implemented as a “CAS loop”: it prepares an operation and attempts to atomically commit this operation using a CAS instruction. If the CAS succeeds, the function returns; otherwise, the loop continues with another attempt. Here, each loop is encoded as a tail-recursive function.

The function `push s v` inserts a new element `v` in a stack `s`. First, `s` is dereferenced (line 4) so as to obtain the address `h` of the head of the linked list. Then, a new list cell `h'` is allocated (line 5). The “data” and “tail” fields are initialized with `v` (line 6) and `h` (line 7). Then, a CAS instruction attempts to update the content of `s` from `h` to `h'` (line 8). If this attempt is successful, `push` returns (line 9); otherwise, it means that a concurrent `push` or `pop` has succeeded. In this case, another attempt is made (line 10).

The function `pop s` extracts the top element of the stack `s`. First, the head `h` of the linked list is read (line 12). If the list is empty, `pop` makes another attempt (line 14), waiting for the stack to become nonempty. Otherwise, the “tail” field of the cell `h` is read so as to obtain the address `h'` of the next list cell (line 16). Then, a CAS instruction attempts to update the content of `s` from `h` to `h'` (line 17). If this attempt is successful, `pop` reads the “data” field of the cell `h` and returns its value (line 18); otherwise, it means that a concurrent `push` or `pop` has succeeded. In this case, another attempt is made (line 19).

Treiber’s stack is *linearizable* [Herlihy and Wing 1990], in the sense that `push` and `pop` atomically take effect at a certain point between the function call and return.

3.2 Space Consumption of Treiber’s Stack without Protected Sections

What is the space consumption of `push` and `pop`? Let us write W for the number of memory words occupied by one list cell. A successful `push` operation consumes W memory words, as it allocates one single list cell. Symmetrically, one might expect a successful `pop` operation to free up W memory words. Indeed, the list cell that is extracted from the list becomes unused, so one might expect the garbage collector to be able to reclaim it.

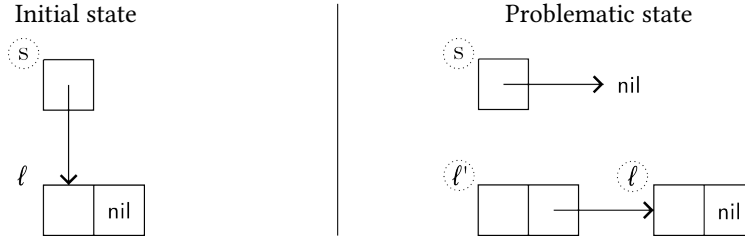


Fig. 2. Initial and problematic states of the example scenario for Treiber’s stack. A box represents a memory block, whose location appears at the top left. A circled location is a root.

However, this intuition is false: when pop returns, although the list cell that has just been extracted is indeed unused, it is not necessarily unreachable. Indeed, it might still be a root of other threads that are still in the process of executing a push or pop operation (which will fail) on this cell. In such a case, all descendants of this cell remain reachable as well. This issue can lead to bad worst-case heap space complexity, but this depends on the placement of polling points.

A Potentially Problematic Scenario and a Way Out. Let us examine the problematic scenario in greater detail. In this scenario, a cell that has been extracted by a successful pop operation remains reachable by other threads, preventing its immediate reclamation. Figure 2 depicts the initial state and the problematic state of this scenario. Suppose that the stack s consists of a single list cell whose address is ℓ . Suppose that thread A attempts to push a new value onto s , while thread B attempts to pop a value off s . Thread A starts making progress while thread B is asleep. Thread A begins to execute push. At line 4, its local variable h is bound to the address ℓ . At line 5, it allocates a new list cell at address ℓ' ; its local variable h' is bound to ℓ' . At line 7, the “tail” field of the new cell is set to ℓ . Then, suppose thread A falls asleep. Thread B wakes up and successfully pops one value off the stack. The reference s now stores the value nil . The cell ℓ has been extracted by pop and is no longer logically part of the stack. The cell ℓ' has not yet been inserted by push and is not logically part of the stack.

Because the cell ℓ has been extracted by a pop operation that has successfully completed, one might expect this cell to be now unreachable. However, this is not the case. Thread A has fallen asleep between lines 7 and 8. At this point, the local variables h and h' are still needed in the future: they occur on line 8. Therefore, the locations ℓ and ℓ' are roots in thread A . Besides, even if ℓ was not a root, it would still be reachable via the root ℓ' , since the “tail” field of the cell ℓ' contains the pointer ℓ . This is potentially problematic: a cell that has been extracted by pop is still reachable after pop has returned. So, *at this point, the garbage collector cannot reclaim this cell*. Therefore, the informal claim that “pop frees up W words of memory” seems compromised. Can it be formalized in such a way that it is actually true?

To answer this question in a positive way, a tempting first idea is to somehow forbid this undesirable scenario. For example, forbidding thread A from falling asleep at this particular point, between lines 7 and 8, might come to mind. However, this idea does not seem practical; we do not wish to impose restrictions on the scheduler.

Instead, we remark that *provided thread A is allowed to make progress until line 10*, there really is no problem at all. That is, *provided no polling point prevents thread A from reaching line 10*, the above claim is true, in a certain sense.

Indeed, should some other thread (say, thread Z) signal that it needs space, then the scheduler can choose to immediately invoke the garbage collector, which might or might not free up enough space; but, what is more interesting, the scheduler can also wake up thread A and let it make


```

1  let create () = ref nil
2
3  let rec push s v =
4    let h' = new_cell () in
5    set_data h' v;
6    enter; let h = !s in
7    set_tail h' h;
8    if compare_and_swap s h h'
9    then exit
10   else (exit; push s v)
11
12   let rec pop s =
13     enter; let h = !s in
14     if is_nil h
15     then (exit; pop s)
16     else
17       let h' = tail h in
18       if compare_and_swap s h h'
19       then (let v = data h in exit; v)
20       else (exit; pop s)

```

Fig. 3. A safe-for-space version of Treiber’s stack. Protected section entry and exit points are highlighted.

progress. Recall the scenario that we are considering: thread B has successfully executed `pop` after the location ℓ was read from s by thread A at line 4. Therefore, the CAS instruction in thread A must fail, and thread A must reach the second branch of the conditional construct, at line 10. At this point, the variables h and h' are no longer needed, so the locations ℓ and ℓ' are no longer roots in thread A . Moreover, ℓ' does not appear in the heap at all, and ℓ can be reached only via ℓ' : therefore, both ℓ and ℓ' are unreachable. Hence, at this point, these cells can be reclaimed. Thus, they can be reclaimed *before* the request for space by thread Z must be satisfied.

In summary, provided no polling point prevents thread A from reaching line 10, there is no problem: the informal claim that “pop frees up W words of memory”, is true, provided one understands that it means: “once pop has returned, *eventually*, W words of memory will be freed”.

On the other hand, if a polling point was inserted by the compiler in the code of `push` between lines 7 and 8, that would be problematic. Assuming that thread Z is waiting for space, there would be no way for thread A to make progress past this polling point. So, the location ℓ' would remain a root in thread A , and the garbage collector would be unable to reclaim the cell at address ℓ' . With some bad luck, though, perhaps reclaiming this cell is *necessary* in order to recover enough space to satisfy the memory allocation request by thread Z . In such an event, the system would be in a deadlock. It would be impossible to prove that “once pop has returned, eventually, W words of memory will be freed”.

3.3 Space Consumption of Treiber’s Stack with Protected Sections

In the previous section (§3.2), we have argued that the compiler must not insert a polling point between the point where the address of the head cell is read and a point where this address ceases to be a root. To forbid this, a protected section can and must be used. The modified pseudo-code in Figure 3 illustrates one way of doing so. With respect to the original code in Figure 1, two main changes are made. First, protected sections, delimited by `enter` and `exit` instructions, are inserted into the functions `push` and `pop`. Second, the allocation of a new list cell in `push` must be anticipated (moved higher up in the code), because memory allocation inside protected sections is forbidden (§2.5).

The protected sections in Figure 3 are placed in such a way that, outside of these protected sections, no list cell is a root. Thus, there is no danger of a polling point being inserted at a point where a list cell is a root.

With respect to this version of the code, we are able to formally prove that “pop frees up W words of memory”. Technically, the list cell addresses that are read inside protected sections are registered in our program logic as temporary roots. This allows these addresses to be logically deallocated by a successful `pop` operation. More details about this statement and about its proof are given later on (§11.5).

```

20 let rec push_aux s v h' =
21   set_data h' v;
22   enter; let h = !s in
23   set_tail h' h;
24   if compare_and_swap s h h'
25   then exit
26   else (exit; push_aux s v h')
27
28 let push' = push_aux s v (new_cell ())

```

Fig. 4. A tempting yet unsafe-for-space optimization

An Unsafe-For-Space Optimization. The code for push that we present in Figure 3 may seem naive. Indeed, a new list cell is allocated at each iteration of the CAS loop (that is, at each recursive call to push). A tempting optimization is to allocate a list cell outside of the loop and to reuse this cell at each recursive call. Such an optimization is standard [Herlihy and Shavit 2012, §11.2]. Figure 4 presents a version of push, named push', which incorporates it.

Unfortunately, this optimization re-introduces the space complexity problem that was discussed earlier (§3.2): it can prevent a concurrent pop operation from freeing up memory. Indeed, push' does *not* satisfy the property “outside protected sections, no list cell is a root”, which push satisfies, and which is sufficient (although not necessary) to achieve good space complexity. The culprit is the local variable h', which is mentioned on line 26, outside of the protected section. Let us write ℓ and ℓ' for the addresses to which the local variables h and h' are bound. Although the list cell at address ℓ' is not yet part of the data structure, it points to the list cell at address ℓ , which is or has been part of the data structure. Therefore, via the root ℓ' , the list cell ℓ and its descendants are reachable; this is undesirable.

To eliminate this problem, one solution would be to overwrite the pointer from ℓ' to ℓ in case the CAS instruction fails. One could do so by inserting a store operation just before the exit instruction on line 26. Then, the list cell at address ℓ' would still be a root outside of the protected section, but (because this cell has no successor) this would not cause any other cells to remain reachable. Thus, the desired heap space complexity would be achieved. However, we fear that, due to the extra store operation, this code may be more expensive in practice than the original unoptimized code, which allocates a new cell at each loop iteration. Furthermore, IrisFit, as presented in this paper, cannot verify this alternative version. Indeed, inside protected sections, IrisFit has a concept of temporary root, but does not have a concept of temporary heap-to-heap pointer.

Another tentative solution would be to place the entire loop inside a single protected section. Then, as in the unoptimized code, it would clearly be the case that, outside of a protected section, no list cell is a root. However, because we want to ensure that protected sections terminate in bounded time, we have forbidden loops inside protected sections. Even if we did allow provably terminating loops inside protected sections, that would not help: the loop in push_aux is potentially non-terminating.

4 SYNTAX AND SEMANTICS OF LAMBDAFIT

In this section, we formally present the syntax of LambdaFit (§4.1) and its small-step reduction relations (§4.2). These reduction relations define the three semantics of LambdaFit (§2.1).

4.1 Syntax

The syntax of LambdaFit appears in Figure 5. A *value* v is a piece of data that fits in one word of memory. A value can be the unit value $()$, a Boolean value b , an integer value z , a memory

Primitives	$\odot ::= \&\& \mid \ \mid + \mid - \mid \times \mid \div \mid =$			
Values	$v, w ::= () \mid b \in \{\text{false}, \text{true}\} \mid z \in \mathbb{Z} \mid \ell \in \mathcal{L} \mid \mu_{\text{ptr}} f. \lambda \vec{x}. t$ where $fv(t) \subseteq \{f\} \cup \vec{x}$			
Terms	$t, u ::= v$	value	$t[t]$	heap load
	x	variable	$t[t] \leftarrow t$	heap store
	let $x = t$ in t	sequencing	fork t	thread creation
	if t then t else t	conditional	CAS $t[t] t t$	compare-and-swap
	$(t \vec{u})_{\text{ptr}}$	code pointer invocation	enter	entering a protected section
	$t \odot t$	primitive operation	exit	exiting a protected section
	alloc t	heap allocation	poll	polling point
Contexts	$K ::= \text{let } x = \square \text{ in } t$	if \square then t else t	$\square \odot t$	$v \odot \square$
	alloc \square	$\square[t]$	$v[\square]$	$\square[t] \leftarrow t$
	$v[\square] \leftarrow t$	$v[v] \leftarrow \square$	$(\square \vec{u})_{\text{ptr}}$	$(v (\vec{v} \# \square \# \vec{u}))_{\text{ptr}}$
	CAS $\square[t] t t$	CAS $v[\square] t t$	CAS $v[v] \square t$	CAS $v[v] v \square$

Fig. 5. LambdaFit: syntax

location ℓ (drawn from an infinite set \mathcal{L}), or a code pointer $\mu_{\text{ptr}} f. \lambda \vec{x}. t$. Such a code pointer is a closed, recursive, multi-argument function. The side condition $fv(t) \subseteq \{f\} \cup \vec{x}$ ensures that the function is closed: that is, the only variables that may appear in the body of the function are f (a self-reference, allowing the function to invoke itself) and \vec{x} (the formal parameters).

The syntax of terms (also known as expressions) includes a number of standard sequential constructs, such as sequencing, conditionals, code pointer invocations, and primitive operations. The heap allocation expression $\text{alloc } n$ allocates a fresh memory block of size n and returns its address. The field at offset i in the memory block at address x is read by the “load” expression $x[i]$ and written by the “store” expression $x[i] \leftarrow y$.

Two standard concurrency-related constructs are “fork” and CAS. The expression $\text{fork } t$ spawns a new thread whose code is t . This is *unstructured concurrency*: there is no primitive operation to wait until a thread terminates. This approach contrasts with the less-expressive *structured concurrency*, such as fork/join or async/finish [Charles et al. 2005; Lee and Palsberg 2010], where there is a primitive operation that waits for a thread or a group of threads to terminate. As we demonstrate later in this paper (§11.4), structured concurrency (async/finish) can be encoded as a library on top of unstructured concurrency.

The compare-and-swap expression $\text{CAS } \ell[i] v v'$ atomically loads a value from block ℓ at offset i , compares this value with v , and, in case they are equal, overwrites this value with v' . Its Boolean result indicates whether the write took place.

The poll instruction is a polling point (§2.4). The instructions enter and exit mark the beginning and end of a protected section (§2.5).

4.2 Semantics

We now define the operational semantics of LambdaFit. We begin with our model of memory, that is, our view of the heap as a collection of memory blocks, and our notion of heap size (§4.2.1). We define thread pools and configurations (§4.2.2). Then, we introduce a series of reduction relations which, together, form the dynamic semantics of LambdaFit. The *head reduction* relation (§4.2.3) describes one elementary step of computation by one thread. The *step* relation (§4.2.4) allows head reduction to take place under an evaluation context. It represents one step of computation by one

$$\begin{array}{c}
\text{HEADLETVAL} \\
\text{let } x = v \text{ in } t / g / \sigma \xrightarrow{\text{head}} [v/x]t / g / \sigma / \varepsilon
\end{array}
\qquad
\begin{array}{c}
\text{HEADCALL} \\
\frac{v = \mu_{\text{ptr}} f. \lambda \vec{x}. t \quad |\vec{x}| = |\vec{w}|}{(v \vec{w})_{\text{ptr}} / \text{Out} / \sigma \xrightarrow{\text{head}} [v/f][\vec{w}/\vec{x}]t / \text{Out} / \sigma / \varepsilon}
\end{array}$$

$$\begin{array}{c}
\text{HEADIFTRUE} \\
\frac{\text{if true then } t_1 \text{ else } t_2 / g / \sigma}{\xrightarrow{\text{head}} t_1 / g / \sigma / \varepsilon}
\end{array}
\qquad
\begin{array}{c}
\text{HEADIFFALSE} \\
\frac{\text{if false then } t_1 \text{ else } t_2 / g / \sigma}{\xrightarrow{\text{head}} t_2 / g / \sigma / \varepsilon}
\end{array}
\qquad
\begin{array}{c}
\text{HEADENTER} \\
\frac{\text{enter} / \text{Out} / \sigma}{\xrightarrow{\text{head}} () / \text{In} / \sigma / \varepsilon}
\end{array}$$

$$\begin{array}{c}
\text{HEADEXIT} \\
\frac{\text{exit} / \text{In} / \sigma}{\xrightarrow{\text{head}} () / \text{Out} / \sigma / \varepsilon}
\end{array}
\qquad
\begin{array}{c}
\text{HEADPRIM} \\
\frac{v_1 \odot v_2 \xrightarrow{\text{pure}} v}{v_1 \odot v_2 / g / \sigma \xrightarrow{\text{head}} v / g / \sigma / \varepsilon}
\end{array}
\qquad
\begin{array}{c}
\text{HEADALLOC} \\
\frac{\ell \notin \text{dom}(\sigma) \quad 0 < n \quad \sigma' = [\ell := ()^n] \sigma}{\text{alloc } n / \text{Out} / \sigma \xrightarrow{\text{head}} \ell / \text{Out} / \sigma' / \varepsilon}
\end{array}$$

$$\begin{array}{c}
\text{HEADLOAD} \\
\frac{\sigma(\ell) = \vec{w} \quad 0 \leq i < |\vec{w}| \quad \vec{w}(i) = v}{\ell[i] / g / \sigma \xrightarrow{\text{head}} v / g / \sigma / \varepsilon}
\end{array}
\qquad
\begin{array}{c}
\text{HEADSTORE} \\
\frac{\sigma(\ell) = \vec{w} \quad 0 \leq i < |\vec{w}| \quad \sigma' = [\ell := [i := v] \vec{w}] \sigma}{\ell[i] \leftarrow v / g / \sigma \xrightarrow{\text{head}} () / g / \sigma' / \varepsilon}
\end{array}$$

$$\begin{array}{c}
\text{HEADCASFAILURE} \\
\frac{\sigma(\ell) = \vec{w} \quad 0 \leq i < |\vec{w}| \quad \vec{w}(i) \neq v}{\text{CAS } \ell[i] v v' / g / \sigma \xrightarrow{\text{head}} \text{false} / g / \sigma / \varepsilon}
\end{array}
\qquad
\begin{array}{c}
\text{HEADCASSUCCESS} \\
\frac{\sigma(\ell) = \vec{w} \quad 0 \leq i < |\vec{w}| \quad \vec{w}(i) = v \quad \sigma' = [\ell := [i := v'] \vec{w}] \sigma}{\text{CAS } \ell[i] v v' / g / \sigma \xrightarrow{\text{head}} \text{true} / g / \sigma' / \varepsilon}
\end{array}$$

$$\begin{array}{c}
\text{HEADPOLL} \\
\text{poll} / \text{Out} / \sigma \xrightarrow{\text{head}} () / \text{Out} / \sigma / \varepsilon
\end{array}
\qquad
\begin{array}{c}
\text{HEADFORK} \\
\text{fork } t / \text{Out} / \sigma \xrightarrow{\text{head}} () / \text{Out} / \sigma / t
\end{array}$$

Fig. 6. The head reduction relation

$$\begin{array}{c}
\text{STEPHEAD} \\
\frac{t / g / \sigma \xrightarrow{\text{head}} t' / g' / \sigma' / t^?}{t / g / \sigma \xrightarrow{\text{step}} t' / g' / \sigma' / t^?}
\end{array}
\qquad
\begin{array}{c}
\text{STEPCTX} \\
\frac{t / g / \sigma \xrightarrow{\text{step}} t' / g' / \sigma' / t^?}{K[t] / g / \sigma \xrightarrow{\text{step}} K[t'] / g' / \sigma' / t^?}
\end{array}$$

Fig. 7. The step relation

thread. The *garbage collection* relation (§4.2.5) describes the effect of the GC on the heap. The *action* relation (§4.2.6) combines computation steps and garbage collection steps. Allowing any computation steps to take place and preventing garbage collection yields the *oblivious semantics* (§4.2.7). Restricting the action relation to a subset of *enabled* actions (§4.2.8), which depends on a heap size limit S , yields the *default reduction* relation (§4.2.9), which is the main reduction relation of the default semantics. Allowing for the heap size limit S to grow over time yields the *growing reduction* relation (§4.2.10).

4.2.1 Memory Blocks, Stores, and Heap Size. A *memory block* is either a tuple of values, written \vec{v} , or a special deallocated block, written \clubsuit . A *store* σ (or *heap*) is a finite map of locations to memory blocks. We write \emptyset for the empty store.

Our semantics does not recycle memory locations. When a heap block at address ℓ is reclaimed by the GC, the store is updated with a mapping of ℓ to \clubsuit . The address ℓ continues to exist and is never re-used. Naturally, in an implementation, memory locations would be recycled. However, we

$$\begin{array}{c}
\text{EDGE} \\
\frac{\sigma(\ell) = \vec{w} \quad \vec{w}(i) = \ell'}{\ell \rightsquigarrow_{\sigma} \ell'} \\
\\
\text{GC} \\
\frac{\begin{array}{l} \forall \ell. \ell \in \text{dom}(\sigma) \implies \\ \left\{ \begin{array}{l} \sigma'(\ell) = \sigma(\ell) \\ \vee \sigma'(\ell) = \spadesuit \wedge \neg (\exists r \in R, r \rightsquigarrow_{\sigma}^* \ell) \end{array} \right. \\ \text{dom}(\sigma') = \text{dom}(\sigma) \end{array}}{R \vdash \sigma \xrightarrow{\text{GC}} \sigma'}
\end{array}$$

Fig. 8. The garbage collection relation

work at a higher level of abstraction. The reasoning rules of our program logic guarantee that a memory allocation always produces a fresh address.

We assume that the space usage (in words) of a block of n fields is $\text{size}(n)$, where size is a mathematical function of \mathbb{N} to \mathbb{N} . If, for instance, every memory block is preceded by a one-word header, then the function size would be defined by $\text{size}(n) = n + 1$. LambdaFit and IrisFit are independent of the definition of size . For our case studies (§11), we chose $\text{size}(n) = n$. We write $\text{size}(\vec{v})$ as a shorthand for $\text{size}(n)$, where n is the length of the list \vec{v} . By convention, we let $\text{size}(\spadesuit)$ be 0. This reflects the fact that a deallocated block occupies no space.

We define the size of a store σ as the sum of the sizes of its blocks. Thus, we do not measure the physical size of the heap, that is, how much memory has been borrowed from the operating system. Instead, we measure the total size of the memory blocks that are currently allocated. We ignore fragmentation.

4.2.2 Thread Pools and Configurations. A thread t is just a term. A thread's status g is either In or Out. The status records whether the thread is currently inside or outside a protected section. A thread pool θ is a list of pairs (t, g) of a thread t and its status g . A thread identifier π is an integer index into a thread pool.

A configuration c is a pair (θ, σ) of a thread pool θ and a store σ . The initial configuration for a program t consists of a thread pool that contains just the thread (t, Out) and the empty store \emptyset . We write $\text{init}(t)$ for this initial configuration. We define the heap size of a configuration as the size of its store: $\text{size}((\theta, \sigma)) = \text{size}(\sigma)$.

4.2.3 The Head Reduction Relation. The head reduction relation $t / g / \sigma \xrightarrow{\text{head}} t' / g' / \sigma' / t^?$ describes an evolution of the term t with status g and store σ to a term t' with status g' and store σ' , optionally forking off a new thread $t^?$. The metavariable $t^?$ denotes an optional term: it is either a term t or ε , which means that no thread was forked off.

The head reduction relation describes the reduction of a single isolated instruction. Reduction under an evaluation context is handled by the step relation introduced shortly afterwards (§4.2.4). Moreover, the head reduction relation describes how an instruction is executed under the assumption that this instruction is *enabled*, that is, not blocked. The definition of enabled instructions, which describes under what conditions an instruction is blocked, is given later on (§4.2.8).

The head reduction relation is defined by the rules in Figure 6.

HEADLETVAL, HEADIFTRUE, HEADIFFALSE, HEADPRIM are standard.

HEADLOAD, HEADSTORE, HEADCASSUCCESS and HEADCASFAILURE, which describe memory accesses, are also standard. These rules require that the memory location ℓ be valid: this is expressed by the premise $\sigma(\ell) = \vec{w}$. Furthermore, they require the integer value i to be a valid index into the memory block at address ℓ : this is expressed by the premise $0 \leq i < |\vec{w}|$. We write $\vec{w}(i)$ for the i -th value in the sequence \vec{w} , and $[i := v]\vec{w}$ for the sequence obtained by updating the sequence \vec{w} at index i with the value v . We write $[\ell := \vec{w}]\sigma$ for the store obtained by updating the store σ at

$$\begin{array}{c}
\text{ACTIONTHREAD} \\
\frac{\theta(\pi) = (t, g) \quad t / g / \sigma \xrightarrow{\text{step}} t' / g' / \sigma' / t^? \quad \theta' = [\pi := (t', g')] \theta ++ [(t^?, \text{Out})]}{(\theta, \sigma) \xrightarrow{\text{action}}_{\pi} (\theta', \sigma')} \\
\text{ACTIONGC} \\
\frac{\text{locs}(\theta) \vdash \sigma \xrightarrow{\text{gc}} \sigma' \quad \sigma \neq \sigma'}{(\theta, \sigma) \xrightarrow{\text{action}}_{\text{gc}} (\theta, \sigma')}
\end{array}$$

Fig. 9. The action relation

address ℓ with the block \vec{w} . Hence, $[\ell := [i := v] \vec{w}] \sigma$ describes an update of the i -th field of the block at location ℓ .

HEADENTER and **HEADEXIT** cause the thread to change its status from Out to In and vice-versa. By design, no reduction rule describes the effect of enter when the thread’s status is In or the effect of exit when the thread’s status is Out. Such a situation is considered a runtime error: the thread is *stuck*.

HEADCALL, **HEADALLOC**, **HEADFORK**, and **HEADPOLL** require the thread’s status to be Out. Thus, inside a protected section, a function call, a memory allocation request, a “fork” instruction, or a polling point causes a runtime error. Aside from this, **HEADCALL** and **HEADFORK** are standard. **HEADALLOC** allocates a block of n fields at a fresh memory location and initializes each field with a unit value. We write $()^n$ for a sequence of n unit values. **HEADPOLL** indicates that a polling point is a no-operation: poll acts as a form of barrier (§4.2.8), and is otherwise effectless.

4.2.4 The Step Relation. The *step* relation has the same shape as the head reduction relation (§4.2.3). It takes the form $t / g / \sigma \xrightarrow{\text{step}} t' / g' / \sigma' / t^?$. It is inductively defined by the rules **STEPHEAD** and **STEPCTX** in Figure 7. These rules allow one head reduction step under a stack of evaluation contexts. An *evaluation context* K is a term with a hole written \square at depth exactly 1. The syntax of evaluation contexts, presented in Figure 5, dictates a left-to-right, call-by-value evaluation strategy. We write $K[t]$ for the term obtained by filling the hole of the evaluation context K with the term t .

4.2.5 The Garbage Collection Relation. Several concepts related with garbage collection are defined in Figure 8. The *edge* relation $\ell \rightsquigarrow_{\sigma} \ell'$, defined by the rule **EDGE**, means that the block at location ℓ contains a pointer to location ℓ' .⁵ When this relation holds, we say that ℓ is a *predecessor* of ℓ' . The *reachability* relation $\ell \rightsquigarrow_{\sigma}^* \ell'$ is the reflexive-transitive closure of the edge relation.

The *garbage collection* relation $R \vdash \sigma \xrightarrow{\text{gc}} \sigma'$, defined by the rule **GC**, describes the effect of the GC. This relation means that a garbage collection phase can transform the store σ into a store σ' , while respecting the set of roots R , a set of memory locations. This relation is non-deterministic: the GC may reclaim any unreachable memory block, but need not reclaim every such block. According to the first premise of the rule **GC**, the stores σ and σ' have the same domain: garbage collection does not create or destroy any memory locations. According to the second premise, at each memory location ℓ , either nothing happens ($\sigma'(\ell) = \sigma(\ell)$) or a memory block becomes deallocated ($\sigma'(\ell) = \spadesuit$). The second case is permitted only if ℓ is not reachable from any of the roots in the set R .

4.2.6 The Action Relation. The relations defined so far describes how a thread makes a step (§4.2.4) and how the GC makes a step (§4.2.5). We now define a relation that interleaves these two kinds of steps. It is a labeled transition relation: each step is labeled with an *action* a , which is either a thread identifier π or the fixed token “gc”. The *action* relation $c \xrightarrow{\text{action}}_a c'$ relates two configurations c and c' and is labeled with an *action*. It is defined by the two rules in Figure 9. **ACTIONTHREAD** allows a step by one thread whose identifier is π . This thread evolves from (t, g) to (t', g') : the thread

⁵A value either is a location or contains no location at all. Thus, in **EDGE**, we write just $\vec{w}(i) = \ell'$ instead of the seemingly more general condition $\ell' \in \text{locs}(\vec{w}(i))$.

$$\frac{\text{OBLIVIOUS} \quad c \xrightarrow[\pi]{\text{action}} c'}{c \xrightarrow{\text{oblivious}} c'}$$

Fig. 10. The oblivious reduction relation

$$\begin{array}{c} \text{ISALLOCHEAD} \\ \text{IsAlloc } n \text{ (alloc } n) \end{array} \quad \frac{\text{ISALLOCTX} \quad \text{IsAlloc } n \ t}{\text{IsAlloc } n \ (K[t])} \quad \text{ISPOLLEHEAD} \quad \text{IsPoll } \text{poll} \quad \frac{\text{ISPOLLECTX} \quad \text{IsPoll } t}{\text{IsPoll } (K[t])}$$

$$\frac{\text{ALLOCFITS} \quad \forall n. \text{IsAlloc } n \ t \implies \text{size}(\sigma) + n \leq S}{\text{AllocFits}_S \ \sigma \ t} \quad \frac{\text{EVERYALLOCFITS} \quad \forall t \ g. (t, g) \in \theta \implies \text{AllocFits}_S \ \sigma \ t}{\text{EveryAllocFits}_S \ (\theta, \sigma)}$$

$$\frac{\text{ENABLEDTHREAD} \quad c = (\theta, \sigma) \quad \theta(\pi) = (t, g) \quad \text{AllocFits}_S \ \sigma \ t \quad \text{IsPoll } t \implies \text{EveryAllocFits}_S \ c}{\text{Enabled}_S \ c \ \pi} \quad \frac{\text{ENABLEDGC}}{\text{Enabled}_S \ c \ g}$$

Fig. 11. Enabled actions (and auxiliary predicates)

pool is updated accordingly. The heap, which is shared between all threads, evolves from σ to σ' . A new thread t^2 possibly appears: if so, the thread pool is extended with the new entry (t^2, Out) . **ACTIONGC** describes a garbage collection step. The roots provided to the GC are $\text{locs}(\theta)$, that is, the locations that occur in the thread pool: this is the FVR (§2.2). The side condition $\sigma \neq \sigma'$ forbids stuttering steps, where the GC is invoked but frees no memory. Without this side condition, we would be unable to establish our liveness theorem (Theorem 8.2), whose statement asserts that at all times, *in a bounded number of steps*, the system must reach a configuration where no thread is blocked by a memory allocation request.

4.2.7 The Oblivious Reduction Relation. The *oblivious reduction* relation $c \xrightarrow{\text{oblivious}} c'$ is defined by the rule **OBLIVIOUS** in Figure 10. This relation simply allows one action by an arbitrary thread π . Garbage collection steps are not permitted: in this semantics, there is no need for garbage collection. There is no limit on the size of the heap. This is the oblivious semantics of LambdaFit (§2.1).

4.2.8 Enabled Actions. In the default and growing semantics, two LambdaFit instructions can have a blocking behavior: a large memory allocation instruction is blocking (§2.3); if a large memory allocation request is outstanding, then a polling point is blocking (§2.4). To reflect this, we must define under what conditions an action is *enabled* (allowed to proceed) or *disabled* (blocked).

The distinction between *small* and *large* memory allocation requests depends on the heap size limit S (§2.3). Therefore, the notion of enabled action depends on the parameter S .

To define enabled actions, a few auxiliary predicates are needed. They appear in Figure 11.

The proposition $\text{IsAlloc } n \ t$ means that the next instruction of the thread t is “alloc n ”. In other words, this thread is now requesting a new memory block of n fields. Similarly, the proposition $\text{IsPoll } t$ means that the next instruction of the thread t is “poll”.

The proposition $\text{AllocFits}_S \ t \ \sigma$ means that, if the next instruction in thread t is an allocation request, then it is a small one: that is, there is currently enough free space in the store σ to satisfy it. When this is the case, we say that *thread t fits*. The proposition $\text{EveryAllocFits}_S \ c$ means that, in the configuration c , every thread fits. These propositions depend on the heap size limit S .

$$\begin{array}{c}
\text{ENABLEDACTION} \\
\frac{\text{Enabled}_S c a \quad c \xrightarrow{\text{action}}_a c'}{c \xrightarrow{\text{enabled actions}_S}_a c'} \\
\\
\text{DEFAULT} \\
\frac{c \xrightarrow{\text{enabled actions}_S}_a c'}{c \xrightarrow{\text{default}_S} c'}
\end{array}$$

Fig. 12. The default reduction relation

$$\begin{array}{c}
\text{ALLOUTSIDE} \\
\frac{\forall t g. (t, g) \in \theta \implies g = \text{Out}}{\text{AllOutside}(\theta, \sigma)} \\
\\
\text{GROWINGSTEP} \\
\frac{c \xrightarrow{\text{default}_S} c'}{(S, c) \xrightarrow{\text{growing}} (S, c')} \\
\\
\text{GROWINGINCREASELIMIT} \\
\frac{\text{AllOutside } c \quad \neg \text{EveryAllocFits}_S \text{ fullGC}(c)}{(S, c) \xrightarrow{\text{growing}} (\text{grow}(S), c)}
\end{array}$$

Fig. 13. The growing reduction relation

The proposition $\text{Enabled}_S c a$ means that, in the configuration c , action a is enabled. This proposition also depends on the heap size limit S . It is defined by the rules **ENABLEDTHREAD** and **ENABLEDGC** in Figure 11. For a thread π to be enabled, it must be the case that (1) thread π fits and (2) if thread π is at a polling point then every thread fits. **ENABLEDGC** states that garbage collection is enabled at all times.

The following simple lemma states that if every thread fits then every action is enabled. It is used in the proof of our liveness theorem (§8.2). In the following, we say that a thread identifier π is *valid with respect to the configuration* (θ, σ) if $0 \leq \pi < |\theta|$ holds.

LEMMA 4.1 (ALL ENABLED). *If $\text{EveryAllocFits}_S c$ holds, then, for every thread identifier π that is valid with respect to the configuration c , $\text{Enabled}_S c \pi$ holds.*

4.2.9 The Default Reduction Relation. The auxiliary relation $c \xrightarrow{\text{enabled actions}_S}_a c'$, which is defined in Figure 12, is the restriction of the action relation to enabled actions.

The *default reduction* relation $c \xrightarrow{\text{default}_S} c'$ is obtained from this auxiliary relation by abstracting away the action a . Thus, a step in the default reduction relation corresponds to an enabled action by some thread or by the GC. This reduction relation is the *default semantics* of LambdaFit.

By design of this semantics, the size of the heap can never exceed the limit S . This is an immediate consequence of the fact that large memory allocation requests are blocked.

LEMMA 4.2 (HEAP SIZE). *If $\text{size}(c) \leq S$ and $c \xrightarrow{\text{default}_S} c'$ then $\text{size}(c') \leq S$.*

This simple lemma is not used anywhere; it serves to document the design of the semantics.

4.2.10 The Growing Reduction Relation. The growing semantics (§2.1) is a variation of the default semantics where the limit on the size of the heap can be automatically adjusted at runtime. In this semantics, this limit is not fixed: instead, it is part of the current state of the machine. Thus, instead of relating two configurations c and c' (§4.2.2), the *growing reduction relation* relates two pairs (S, c) and (S', c') , where S and S' represent the value of the heap limit before and after the reduction step.

This relation is defined by the two rules in Figure 13. The rule **GROWINGSTEP** states that if *under the current heap limit S* a step is possible then this step can take place and the heap limit is unchanged. (This can be a computation step or a garbage collection step.) The rule **GROWINGINCREASELIMIT** states that if no thread is currently inside a protected section, and if in spite of the efforts of the

garbage collector the current heap limit blocks a memory allocation request, then the heap limit can be increased from S to $grow(S)$.⁶

The function $grow$ is fixed: it is a parameter of the semantics. We make the following assumption:

ASSUMPTION 4.2.1. *The function $grow : \mathbb{N} \rightarrow \mathbb{N}$ satisfies the following two properties:*

- $\forall x. x < grow(x)$
- $\forall x y. x \leq y \implies grow(x) \leq grow(y)$

A typical example would be $grow(S) = \max(2S, 1)$, which means that when the current limit is found to be too low, it is doubled (with special care for the case where the current limit is zero).

The default semantics and the growing semantics are close cousins: they are related by the following two lemmas. Lemma 4.3 states that if an execution under the growing semantics brings the heap limit up to the value S' , then the same execution is permitted by the default semantics with a fixed limit of S' . Conversely, Lemma 4.4 states that if an execution under the default semantics is possible, with a fixed limit of S , then the same execution is permitted by the growing semantics, without a need to increase the limit.

LEMMA 4.3 (GROWING TO DEFAULT). $(S, c) \xrightarrow{growing,*} (S', c')$ implies $c \xrightarrow{default_{S'}} c'$ and $S \leq S'$.

LEMMA 4.4 (DEFAULT TO GROWING). $c \xrightarrow{default_S,*} c'$ implies $(S, c) \xrightarrow{growing,*} (S, c')$.

Furthermore, in the growing semantics, by design, the size of the heap cannot exceed the current heap limit:

LEMMA 4.5 (HEAP SIZE). *If $size(c) \leq S$ and $(S, c) \xrightarrow{growing,*} (S', c')$ then $size(c') \leq S'$.*

Later on (§8.3), we prove that, under the growing semantics, a verified program can be executed in bounded space.

5 PROGRAM LOGIC: ASSERTIONS

This section offers an overview of the various kinds of assertions that play a role in IrisFit. We introduce the syntax of each assertion, its intuitive meaning, and the ghost reasoning rules that help understand this meaning, such as splitting and joining rules. We informally explain the life cycle of each assertion: where it typically appears, where it is exploited, and where it is consumed. A presentation of the reasoning rules for terms is deferred to the following section (§6).

We begin with a presentation of triples (§5.1) and ghost updates (§5.2). Then, we briefly present the standard points-to assertion (§5.3), the novel “*sizeof*” assertion (§5.4), and space credits (§5.5). We then devote our attention to the assertions that record reachability or unreachability information, namely the pointed-by-heap assertion (§5.6), the novel pointed-by-thread assertion (§5.7), the novel “*inside*” and “*outside*” assertions (§5.8), and deallocation witnesses (§5.9). Finally, we explain liveness-based cancellable invariants (§5.10), a useful idiom that expresses that a certain invariant holds as long as a certain location is live.

IrisFit is a variant of the Iris program logic [Jung et al. 2018b, §6–7] and is built on top of the Iris base logic [Jung et al. 2018b, §5]. We write Φ for assertions, $\ulcorner P \urcorner$ for a pure assertion, $\Phi * \Phi'$ for a separating conjunction, and $\Phi \multimap \Phi'$ for a separating implication. We express the logical equivalence of two assertions as $\Phi \equiv \Phi'$. A postcondition Ψ is a function of a value to an assertion: in other words, it is the form $\lambda v. \Phi$.

⁶The requirement that “no thread is inside is protected section” is expressed by the premise $\text{AllOutside } c$. We write $\text{fullGC}(c)$ for the configuration obtained from the configuration c by deallocating all unreachable heap blocks. The formal definition of the function fullGC can be found in our Coq mechanization [Moine 2024b].

$$\begin{array}{c}
\text{CONSEQUENCE} \\
\frac{\Phi \pi \Rightarrow^{locs(t)} \Phi' \quad \{\Phi'\} \pi: t \{\Psi'\} \quad \forall v. \Psi' v \pi \Rightarrow^{locs(v)} \Psi v}{\{\Phi\} \pi: t \{\Psi\}} \\
\text{FRAME} \\
\frac{\{\Phi\} \pi: t \{\Psi\}}{\{\Phi * \Phi'\} \pi: t \{\lambda v. \Psi v * \Phi'\}}
\end{array}$$

Fig. 14. Structural reasoning rules

5.1 Triples

A triple takes the form $\{\Phi\} \pi: t \{\Psi\}$. Its intuitive meaning is that if the store satisfies the assertion Φ then it is safe for thread π to execute the term t ; furthermore, if and when this computation terminates and produces a value v , then the store satisfies the assertion Ψv .

Even though the default reduction relation (§4.2.9) is parameterized with a heap size limit S , the meaning of triples is independent of S . Indeed, triples are internally defined in terms of the *oblivious* reduction relation (§8.4), which does not depend on S . Therefore, none of the reasoning rules mentions S . Our program logic is compositional: each program component can be verified in isolation and without knowledge of S .

Formally, a triple is also parameterized by a *mask* [Jung et al. 2018b, §2.2]. Masks prevent the user from opening an invariant twice. As our treatment of invariants and masks is standard, we omit masks everywhere. The interested reader is referred to our mechanization [Moine 2024b].

We write $\{\Phi\} \pi: t \{\lambda \ell. \Phi'\}$, where the metavariable ℓ denotes a memory location, as syntactic sugar for $\{\Phi\} \pi: t \{\lambda v. \exists \ell. \ulcorner v = \ell \urcorner * \Phi'\}$. We adopt the convention that multi-line assertions are implicitly joined by a separating conjunction.

5.2 Ghost Updates

Iris features *ghost state* and *ghost updates* [Jung et al. 2018b, §5.4]. A ghost update is written $\Phi \Rightarrow \Phi'$. It is an assertion, which means that (up to an update of the ghost state) the assertion Φ can be transformed into Φ' .

In IrisFit, it is sometimes necessary for a ghost update to refer to “the identifier of the current thread” or to “the roots of the current thread”. For this purpose, we introduce a *custom ghost update*, written $\Phi \pi \Rightarrow^V \Phi'$, whose extra parameters are a thread identifier π and a set of memory locations V . It is weaker than a standard ghost update: the law $(\Phi \Rightarrow \Phi') * (\Phi \pi \Rightarrow^V \Phi')$ is valid.

Custom ghost updates are exploited in the **CONSEQUENCE** rule, which appears in Figure 14. This rule allows strengthening the precondition and weakening the postcondition of a triple. Updating the precondition requires a custom ghost update where the parameter V is instantiated with $locs(t)$. Indeed, this set represents the roots at the point where this update takes place. Updating the postcondition requires a custom ghost update where V is instantiated with $locs(v)$, where v denotes the result value of the term t . Indeed, these are the roots at the point where that update takes place.

When a custom ghost update is independent of the parameters π and V , we omit them: we write $\Phi \Rightarrow \Phi'$ for $\forall \pi V. \Phi \pi \Rightarrow^V \Phi'$. Examples of custom ghost updates appear in Figures 18, 19, and 20 and are discussed in the following sections.

The **FRAME** rule, also shown in Figure 14, retains its standard form.

5.3 Points-to Assertions

IrisFit features standard points-to assertions of the form $\ell \mapsto_p \vec{w}$, where p is either a fraction in the semi-open interval $(0, 1]$ or the *discarded fraction* \square [Vindum and Birkedal 2021]. In the latter case, the points-to assertion is persistent.

$$\begin{array}{ccc}
\ell \mapsto_p \vec{w} * \ell \mapsto_p \vec{w} * \text{sizeof } \ell (\text{size}(\vec{w})) & \text{SIZEOFPOINTSTO} \\
\text{sizeof } \ell n * \text{sizeof } \ell m * \lceil n = m \rceil & \text{SIZEOFCONFRONT} \\
\text{sizeof } \ell n \text{ is persistent} & \text{SIZEOFPERSIST}
\end{array}$$

Fig. 15. Reasoning rules of the “sizeof” assertion

$$\begin{array}{ccc}
\lceil \text{True} \rceil \Rightarrow \diamond 0 & \text{ZEROSC} \\
\diamond(n_1 + n_2) \equiv \diamond n_1 * \diamond n_2 & \text{SPLITJOINSC}
\end{array}$$

Fig. 16. Reasoning rules for space credits

Rules. Points-to assertions can be split and joined in the usual way, and a points-to assertion that carries a fraction p can be permanently transformed into one that carries the discarded fraction \square . We do not show these standard rules.

Life cycle. A points-to assertion appears when a memory block is allocated. It is required (and possibly updated) when this block is accessed by a load, store, or CAS instruction (§6.2). It is *not* required or consumed when this block is logically deallocated (§6.1). This is an original feature of IrisFit.

5.4 Sizeof Assertions

The assertion $\text{sizeof } \ell n$ means that there is or there used to be a block of size n at address ℓ . It is persistent: indeed, once the size of a block has been fixed, it can never be changed.

Rules. Two reasoning rules allow introducing and exploiting “sizeof” assertions (Figure 15). **SIZEOFPOINTSTO** creates a “sizeof” assertion out of a points-to assertion. **SIZEOFCONFRONT** states that two “sizeof” assertions for the same address must agree on the size of the block at this address.

Life cycle. The “sizeof” assertion is produced by **SIZEOFPOINTSTO**. It is consulted by the logical deallocation rules (§6.1, §6.6) to determine the number of space credits that must be produced.

5.5 Space Credits

To reason about free space, we use *space credits* [Madiot and Pottier 2022; Moine et al. 2023]. The assertion $\diamond n$ denotes the unique ownership of n space credits. It can be understood as a permission to allocate n words of memory. At a lower level of understanding, this assertion means that n memory words *are currently free or can be freed* by the GC *once it is given a chance to run*. This interpretation of space credits is the same as in the earlier papers cited above.

Following Moine et al. [2023], space credits are measured using non-negative *rational* numbers. Of course, a physical word of memory cannot be split, so the total number of space credits in existence is a natural number; so are the numbers involved in the reasoning rules for memory allocation and deallocation. Still, rational numbers appear essential in certain amortized complexity analyses, as illustrated by the example of chunked stacks [Moine et al. 2023]. Rational credits also appear in amortized *time* complexity analyses [Charguéraud and Pottier 2019; Mével et al. 2019].

Rules. Figure 16 presents two basic reasoning rules about space credits. **ZEROSC** asserts that zero credits can be forged out of thin air. **SPLITJOINSC** asserts that space credits can be split and joined.

Life cycle. Space credits are consumed by memory allocation (§6.2) and produced by logical deallocation (§6.1). Because there is no way of creating space credits out of nothing, a program or program component is usually verified under the assumption that a number of space credits are

$$\begin{array}{ll}
(\ell \leftarrow_{q_1} L_1 * \ell \leftarrow_{q_2} L_2) * \ell \leftarrow_{q_1+q_2} (L_1 \uplus L_2) & \text{JOINPBHEAP} \\
\ell \leftarrow_{q_1+q_2} (L_1 \uplus L_2) * (\ell \leftarrow_{q_1} L_1 * \ell \leftarrow_{q_2} L_2) & \text{if } \begin{cases} q_1 = 0 \Rightarrow \text{NoPositive}(L_1) \\ q_2 = 0 \Rightarrow \text{NoPositive}(L_2) \end{cases} & \text{SPLITPBHEAP} \\
\ell \leftarrow_q L * \ell \leftarrow_q (L \uplus \{+\ell'\}) & \text{if } q > 0 & \text{COVPBHEAP}
\end{array}$$

Fig. 17. Reasoning rules for the pointed-by-heap assertion

provided. For example, our safety theorem for the default semantics (§8.1) states that if a (complete) program is verified under the precondition $\diamond S$, where S is the heap size limit, then this program can be safely executed. When a program component is considered in isolation, it is given a specification that does not mention S . For instance, our specification of Treiber’s stack (Figure 41) states that push consumes two space credits and that pop produces two space credits.

5.6 Pointed-By-Heap Assertions

Our *pointed-by-heap* assertions are the “pointed-by” assertions of our earlier paper [Moine et al. 2023]. The longer name “pointed-by-heap” avoids confusion with our novel “pointed-by-thread” assertions (§5.7). To make this paper self-contained, we recall what form these assertions take, what they mean, and what purpose they serve.

A *pointed-by-heap* assertion for the location ℓ' keeps track of a multiset L of predecessors of ℓ' (§4.2.5). It takes the form $\ell' \leftarrow_q L$, where L is a signed multiset of locations and q is a possibly-null fraction, that is, a rational number in the closed interval $[0; 1]$.

Signed multisets. Signed multisets [Hailperin 1986], also known as *generalized sets* [Whitney 1933; Blizard 1990] or *hybrid sets* [Loeb 1992], are a generalization of multisets: they allow an element to have *negative* multiplicity. A signed multiset is a total function of elements to \mathbb{Z} . The disjoint union operation \uplus is the pointwise addition of multiplicities. We write $+x$ for a positive occurrence of x and $-x$ for a negative occurrence of x . For example, $\{+x; +x\} \uplus \{-x\}$ is $\{+x\}$. We write $\text{NoNegative}(L)$ when no element has negative multiplicity in L . Symmetrically, we write $\text{NoPositive}(L)$ when no element has positive multiplicity in L .

Possibly-Null Fractions. In traditional Separation Logics with fractional permissions [Boyland 2003; Bornat et al. 2005], a fraction is a rational number in the semi-open interval $(0, 1]$. If there exists a share that carries the fraction 1, then no other shares can separately exist. With *possibly-null fractions*, the fraction 0 is allowed, so a full pointed-by-heap assertion $\ell' \leftarrow_1 L$ does *not* exclude the existence of a separate pointed-by-heap assertion with fraction zero, say $\ell' \leftarrow_0 L'$.

Nevertheless, we enforce the following *null-fraction invariant*: in a pointed-by-heap assertion $\ell' \leftarrow_q L$, if the fraction q is 0, then no location can have positive multiplicity in L ; or, in short, $q = 0$ implies $\text{NoPositive}(L)$.

Signed multisets and possibly-null fractions let us use the assertion $\ell' \leftarrow_0 \{-\ell\}$ as a *permission to remove one occurrence of ℓ from the predecessors of ℓ'* . This lets us formulate the reasoning rule for store instructions (§6.2) in a simpler way than would otherwise be possible.

Over-Approximation of Live Predecessors. We say that a location ℓ is *dead* if it has been allocated and logically deallocated already (§5.9, §6.1). We say that it is *live* if it has been allocated but not logically deallocated yet.

The true purpose of pointed-by-heap assertions is to keep track of *live* predecessors. A dead predecessor is irrelevant: increasing its multiplicity in a multiset of predecessors is sound; decreasing

$$\begin{array}{lll}
\ell \Leftarrow_{p_1+p_2} (\Pi_1 \cup \Pi_2) & \equiv & (\ell \Leftarrow_{p_1} \Pi_1 * \ell \Leftarrow_{p_2} \Pi_2) & \text{FRACPBTHREAD} \\
\ell \Leftarrow_p \Pi_1 & * & \ell \Leftarrow_p (\Pi_1 \cup \Pi_2) & \text{COVPBTHREAD} \\
\lceil \ell \notin V \rceil * \ell \Leftarrow_p \{\pi\} & \pi \Rightarrow^V & \ell \Leftarrow_p \emptyset & \text{TRIMPBTHREAD}
\end{array}$$

Fig. 18. Reasoning rules for the pointed-by-thread assertion

it is sound, too. As far as live predecessors are concerned, only over-approximation is permitted. Increasing the multiplicity of a live predecessor is sound; decreasing it is not.

In light of this, and in light of the null-fraction invariant, a *full* pointed-by-heap assertion $\ell' \Leftarrow_1 L$, where the fraction is 1, guarantees that the multiset L contains *all live predecessors* of the location ℓ' . In particular, the assertion $\ell' \Leftarrow_1 \emptyset$ guarantees that ℓ' has *no live predecessors*. Such full knowledge of the live predecessors is required by the logical deallocation rule (§6.1, §6.6).

Rules. Pointed-by-heap assertions obey the splitting, joining, and weakening rules in Figure 17. **JOINPBHEAP** joins two pointed-by-heap assertions by adding the fractions q_1 and q_2 and by adding the signed multisets L_1 and L_2 . In the reverse direction, **SPLITPBHEAP** splits a pointed-by-heap assertion. Its side condition ensures that the null-fraction invariant is preserved. **COVPBHEAP** asserts that a pointed-by-heap assertion (whose fraction is nonzero) is covariant in its multiset: that is, over-approximating the multiset of predecessors is sound. It is a direct consequence of **SPLITPBHEAP**, instantiated with $q_2 \triangleq 0$ and $L_2 \triangleq \{-\ell'\}$. In the reverse direction, the rule **CLEANPBHEAP**, which is discussed later on (§5.9), allows removing a dead predecessor from a multiset of predecessors.

Life cycle. A full pointed-by-heap assertion for the location ℓ appears when this location is allocated. Fractional pointed-by-heap assertions are required, updated, and produced by store instructions. For example, consider a store instruction that updates the field $\ell[i]$ and overwrites the value ℓ'_1 with the value ℓ'_2 . The reasoning rule for this instruction (§6.2) requires a pointed-by-heap assertion $\ell'_2 \Leftarrow_q \emptyset$, which it transforms into $\ell'_2 \Leftarrow_q \{+\ell\}$. Furthermore, the pointed-by-heap assertion $\ell'_1 \Leftarrow_0 \{-\ell\}$ is produced. A full pointed-by-heap assertion for the location ℓ is consumed when ℓ is logically deallocated.

Notation. We define a generalized pointed-by-heap assertion $v \Leftarrow_q L$ whose first argument is a value, as opposed to a memory location. If v is a location ℓ' , then this assertion is defined as $\ell' \Leftarrow_q L$. Otherwise, it is defined as $\lceil \text{True} \rceil$. Furthermore, we write $v \Leftarrow_q^0 L$ for the assertion $\lceil q > 0 \rceil * v \Leftarrow_q L$. This notation is used in the reasoning rule **STORE** (§6.2), among other places.

5.7 Pointed-By-Thread Assertions

The pointed-by-heap assertions presented in the previous section record *which heap blocks* contain pointers to a location ℓ . This information is useful but is not sufficient for our purposes. The logic must also record *which threads* have access to ℓ , that is, in which threads ℓ is a root. For this purpose, we introduce two distinct yet cooperating mechanisms. The first mechanism, presented here, is the pointed-by-thread assertion. The second mechanism, presented next (§5.8), is the “*inside*” assertion. When the fact that ℓ is a root in thread π is recorded by a pointed-by-thread assertion, we say that ℓ is an *ordinary root* in thread π ; when this fact is recorded by an “*inside*” assertion, we say that ℓ is a *temporary root* in thread π . The motivation for this distinction has been given earlier (§3, §2.5).

A *pointed-by-thread* assertion takes the form $\ell \Leftarrow_p \Pi$, where p is a fraction in the semi-open interval $(0; 1]$ and Π is a set of thread identifiers. These assertions intuitively generalize the *Stackable* assertions of our earlier paper [Moine et al. 2023] to a multi-threaded setting.

$inside\ \pi\ T\ *\ outside\ \pi$	\neg	$\lceil False \rceil$	INSIDENOTOUTSIDE
$inside\ \pi\ T\ *\ \ell \Leftarrow_p \{\pi\}$	\Rightarrow	$inside\ \pi\ (T \cup \{\ell\})\ *\ \ell \Leftarrow_p \emptyset$	ADDTEMPORARY
$inside\ \pi\ T\ *\ \ell \Leftarrow_p \emptyset$	\Rightarrow	$inside\ \pi\ (T \setminus \{\ell\})\ *\ \ell \Leftarrow_p \{\pi\}$	REMTEMPORARY
$inside\ \pi\ T$	$\pi \Rightarrow^V$	$inside\ \pi\ (T \cap V)$	TRIMINSIDE

Fig. 19. Reasoning rules for “inside” and “outside” assertions

A *full* pointed-by-thread assertion $\ell \Leftarrow_1 \Pi$, where the fraction is 1, guarantees that Π is the set of *all* threads in which ℓ is an ordinary root. Such full knowledge is required by the logical deallocation rule (§6.1, §6.6).

Rules. Figure 18 presents the splitting, joining, weakening, and trimming rules associated with the pointed-by-thread assertion. **FRACPBT** allows splitting and joining pointed-by-thread assertions. **COVPBT** asserts that a pointed-by-thread assertion is covariant in the set Π : that is, over-approximating Π is sound. **TRIMPBT** allows *trimming* a pointed-by-thread assertion, that is, removing the thread identifier π from a pointed-by-thread assertion for the location ℓ , provided it is evident that ℓ is no longer a root in thread π . This rule is expressed as a custom ghost update $\pi \Rightarrow^V$. It transforms $\ell \Leftarrow_p \{\pi\}$ into $\ell \Leftarrow_p \emptyset$, provided ℓ is not a member of the set V , which denotes the set of roots of the thread π (recall §5.2). The condition $\ell \notin V$ means indeed that ℓ is not a root in thread π . This condition explains why **TRIMPBT** must be expressed as a custom update $\pi \Rightarrow^V$ as opposed to a standard update \Rightarrow . Indeed, a standard update has no means of referring to “the identifier of the current thread” or “the roots of the current thread”.

A curious reader may wonder whether and why **TRIMPBT** remains sound in combination with the **BIND** rule. Indeed, **BIND** lets the user focus on a subterm, therefore implies that the set V is a strict *subset* of the set of all roots of the current thread. This aspect is explained later on (§6.4).

Life cycle. A full pointed-by-thread assertion $\ell \Leftarrow_1 \{\pi\}$ appears when a location ℓ is allocated by a thread π . A fractional pointed-by-thread assertion is ordinarily required and updated by load instructions: when a thread π obtains the location ℓ as the result of a load instruction, an assertion $\ell \Leftarrow_p \emptyset$ is updated to $\ell \Leftarrow_p \{\pi\}$. If the thread π is currently outside a protected section, such an update is mandatory. If the thread π is currently inside a protected section, then it can be avoided by recording ℓ as a temporary root (§6.3). Once ℓ is no longer a root in any thread, **TRIMPBT** can be used to obtain $\ell \Leftarrow_1 \emptyset$, which is consumed by the logical deallocation of ℓ . In fact, when reasoning outside of protected sections, **TRIMPBT** is the only way to trim an assertion $\ell \Leftarrow_1 \Pi$ into $\ell \Leftarrow_1 \emptyset$.

Notation. We define a generalized pointed-by-thread assertion $v \Leftarrow_p \Pi$, whose first argument is a value, as opposed to a memory location. If v is a location ℓ , then this assertion is defined as $\ell \Leftarrow_p \Pi$. Otherwise, it is defined as $\lceil True \rceil$. Besides, we write an iterated conjunction of pointed-by-thread assertions under the form $M \Leftarrow \Pi$, where M is a finite map of memory locations to fractions and Π is a set of thread identifiers. It is defined by $M \Leftarrow \Pi \triangleq \bigstar_{(\ell, p) \in M} (\ell \Leftarrow_p \Pi)$.

5.8 Inside and Outside Assertions

The assertion *outside* π means that the thread π is currently outside a protected section. The assertion *inside* $\pi\ T$ means that thread π is currently inside a protected section and that the set of its temporary roots (§2.6) is T . The set T is a set of memory locations.

Rules. Figure 19 presents a number of reasoning rules related to “inside” and “outside” assertions. **INSIDENOTOUTSIDE** states that a thread cannot be both inside and outside a protected section.

$\dagger \ell$	\Rightarrow	$\ell' \leftarrow_0 \{-\ell\}$	CLEANPBHEAP
$\dagger \ell * \ell \xrightarrow{q}_0 L$	\Rightarrow	$\lceil \text{False} \rceil$	DEADPBHEAP
$\dagger \ell * \ell \leftarrow_p \Pi$	\Rightarrow	$\lceil \text{False} \rceil$	DEADPBTHREAD
$\lceil \ell \in V \rceil * \dagger \ell * \text{outside } \pi$	$\pi \Rightarrow^V$	$\lceil \text{False} \rceil$	NODANGLINGROOTOUT
$\lceil \ell \in (V \setminus T) \rceil * \dagger \ell * \text{inside } \pi T$	$\pi \Rightarrow^V$	$\lceil \text{False} \rceil$	NODANGLINGROOTIN
$\dagger \ell$ is persistent			DEADPERSIST

Fig. 20. Reasoning rules for deallocation witnesses

ADDTEMPORARY converts an ordinary root to a temporary root. The pointed-by-thread assertion $\ell \leftarrow_p \{\pi\}$ is transformed to $\ell \leftarrow_p \emptyset$; meanwhile, ℓ is added to the set of temporary roots carried by the “inside” assertion. In the reverse direction, **REMPERMANENT** converts a temporary root to an ordinary root. **TRIMINSIDE** trims the set of temporary roots by removing any locations that are no longer roots in the current thread. It is analogous to **TRIMPBTHREAD**.

Life cycle. The assertion *outside* π appears when thread π is created and is consumed when this thread terminates. This will be visible in the statement of Theorem 8.1, which describes the creation and termination of the main thread, and in the reasoning rule for “fork” instructions (§6.2). The assertion *outside* π is required and preserved by the instructions that must not appear inside a protected section, namely memory allocations, function calls, “fork” instructions, and polling points. Entering a protected section transforms *outside* π into *inside* $\pi \emptyset$; exiting a protected section causes the reverse transformation.

5.9 Deallocation Witnesses

The persistent assertion $\dagger \ell$ is a *deallocation witness* for the location ℓ . This assertion guarantees that ℓ has been logically deallocated, that is, ℓ is dead.

The fact that ℓ is dead implies that ℓ cannot be reached from an ordinary root. However, this does not imply that ℓ is unreachable: indeed, it could still be reachable via a temporary root.

The assertion $\dagger \ell$ can be read as a permission to remove ℓ from the multiset of predecessors carried by a pointed-by-heap assertion. Indeed, the purpose of pointed-by-heap assertions is to keep track of live predecessors (§5.6).

A non-persistent deallocation witness $x \not\vdash$ appears in Incorrectness Separation Logic [Raad et al. 2020]. In RustBelt [Jung et al. 2018a], the fact that a lifetime κ has ended is expressed by a persistent assertion, known as a *dead token*, written $\lceil \dagger \kappa \rceil$. Persistent deallocation witnesses appear in Madiot and Pottier’s work [2022] and in our earlier paper [Moine et al. 2023]. These two papers do not have protected sections, therefore have no distinction between ordinary and temporary roots. There, a dead location is unreachable.

Rules. Figure 20 presents reasoning rules for deallocation witnesses. **CLEANPBHEAP** requires a deallocation witness for ℓ and produces the assertion $\ell' \leftarrow_0 \{-\ell\}$, a permission to remove ℓ from the predecessors of an arbitrary location ℓ' . **DEADPBHEAP** and **DEADPBTHREAD** reflect the fact that logical deallocation consumes full pointed-by-heap and pointed-by-thread assertions. Therefore, the assertions $\dagger \ell$ and $\ell \leftarrow_q L$ cannot coexist, except in the special case where q is zero, and the assertions $\dagger \ell$ and $\ell \leftarrow_p \Pi$ cannot coexist. However, in contrast with our earlier work [Madiot and Pottier 2022; Moine et al. 2023], our deallocation witness is compatible with the points-to assertion. Indeed, our logical deallocation rule does not consume the points-to assertion. **NODANGLINGROOTOUT** and **NODANGLINGROOTIN** both state that it is impossible for a dead location

to be an ordinary root. A dead location can, however, be a temporary root: indeed, our logical deallocation rule allows deallocating a temporary root (§6.1).

5.10 Liveness-Based Cancellable Invariants

An Iris *invariant* [Jung et al. 2018b, §2.2] is written in the form $\boxed{\Phi}$ ⁷. It is a persistent assertion, whose meaning is that the assertion Φ in the rectangular box holds at all times. The assertion Φ itself is usually not persistent. An invariant can be temporarily *accessed* so as to gain access to the assertion Φ .

A *cancellable invariant* [Jung et al. 2018b, §7.1.3] is an invariant that comes with a teardown mechanism, allowing the user to recover ownership of the assertion Φ once the invariant is canceled. This is a one-shot mechanism: once a cancellable invariant is torn down, it cannot be restored. Naturally, accessing a cancellable invariant requires proving that this invariant has not been torn down already.

In IrisFit, a form of *liveness-based cancellable invariants* (LCIs, for short) naturally arises. An LCI is tied to a memory location ℓ , and remains in force as long as this location is live. When the location ℓ is logically deallocated, all LCIs associated with ℓ are implicitly torn down. Therefore, to access an LCI associated with the location ℓ , one must prove that this location is still live: that is, one must prove that $\dagger \ell$ implies $\ulcorner \text{False} \urcorner$. This can be done using any of the rules `DEADPBHEAP`, `DEADPBTHREAD`, `NODANGLINGROOTOUT`, and `NODANGLINGROOTIN` in Figure 20. When the location ℓ is logically deallocated, the assertion Φ can be recovered at the same time. We have used LCIs to reason about closures (§9.5) and about Treiber’s stack (§11.5).

The implementation of LCIs is simple. A liveness-based cancellable invariant tied to the location ℓ , whose content is the assertion Φ , is just $\boxed{\dagger \ell \vee \Phi}$, that is, a plain Iris invariant whose content is the disjunction $\dagger \ell \vee \Phi$. By proving that $\dagger \ell$ is contradictory, the user excludes the left-hand disjunct, therefore obtains access to Φ . In particular, when one is about to logically deallocate ℓ , the assertion $\ell \Leftarrow \emptyset$ is at hand, so $\dagger \ell$ is excluded. One can therefore open the invariant, extract Φ , deallocate ℓ , and close the invariant by supplying $\dagger \ell$, keeping Φ . (This is a somewhat unusual variation on the “golden idol” technique [Kaiser et al. 2017], with the persistent assertion $\dagger \ell$ in the role of the “bag of sand”.)

6 PROGRAM LOGIC: REASONING RULES

In this section, we present the reasoning rules of IrisFit. Because most of our design is guided by the desire for a flexible logical deallocation rule, we begin with a presentation of this rule, in the simplified case where a single memory location is deallocated (§6.1). Then, we present the reasoning rules for terms (§6.2), devoting special attention to protected sections (§6.3) and to the `BIND` rule, whose form is non-standard (§6.4). The standard statement of the `BIND` rule can be recovered when the user enters a restricted mode where certain rules are disabled (§6.5). Finally, we present the general form of the logical deallocation rule, which can deallocate cycles (§6.6).

6.1 Logical Deallocation

As in the previous papers by Madiot and Pottier [2022] and Moine et al. [2023], a key aspect of IrisFit is to provide a *logical deallocation* rule. This rule produces space credits: by logically deallocating a memory block, the user recovers the space credits that were consumed when this block was allocated. It can be applied to a memory location ℓ as soon as one is able to prove that this memory location is eligible for collection *during the next garbage collection phase*.

⁷Formally, an invariant also carries a *namespace*, a technicality that prevents the user from accessing the invariant twice and obtaining two copies of Φ at the same time. For simplicity, we hide namespaces in this paper.

As in the previous work cited above, *if ℓ is unreachable* then it can be logically deallocated. Furthermore, what is new in this paper, *if ℓ is reachable only via temporary roots* (that is, via roots that will disappear by the time all protected sections are exited), then it can also be logically deallocated.

This reasoning rule may seem surprising, as it involves a form of anticipation: it exploits the fact that ℓ will be eligible for collection *once all protected sections have been exited*, yet it produces space credits *immediately*, at the point where the rule is applied. The intuitive reason why this is sound is related with the following property: “once a large allocation request is outstanding, eventually, either this allocation request becomes unblocked or every thread is outside of a protected section”.⁸ Hence, *only* in the ideal situation where every thread is outside of a protected section, must one prove that a large allocation request can be satisfied. In such a situation, there are no temporary roots; every location that has been logically deallocated is effectively unreachable.

In §6.6, we present the general form of the logical deallocation rule, which can deallocate multiple memory locations at once, even if they form a cycle. Here, we present **FREEONE**, a simplified rule that is also useful in practice and that deallocates a single location ℓ :

$$\text{sizeof } \ell n * \ell \leftarrow_1 \emptyset * \ell \leftarrow_1 \emptyset \Rightarrow \diamond \text{size}(n) * \dagger \ell \quad \text{FREEONE}$$

Because logical deallocation is a ghost operation, **FREEONE** is expressed as a ghost update. It consumes three assertions: the “sizeof” assertion $\text{sizeof } \ell n$, the pointed-by-heap assertion $\ell \leftarrow_1 \emptyset$, and the pointed-by-thread assertion $\ell \leftarrow_1 \emptyset$. The assertion $\text{sizeof } \ell n$ indicates that the memory block at address ℓ has size n . The assertion $\ell \leftarrow_1 \emptyset$ guarantees that ℓ has no predecessor in the heap, that is, no memory block contains the pointer ℓ . The assertion $\ell \leftarrow_1 \emptyset$ guarantees that ℓ is not an ordinary root of any thread: that is, if ℓ is a root at all in a thread π , then it must be a temporary root for this thread (§2.6, §5.8). Together, the last two assertions imply that ℓ will be eligible for collection in the next garbage collection phase.

On the right-hand side of the ghost update, **FREEONE** produces two assertions, namely the recovered space credits $\diamond n$ and the deallocation witness $\dagger \ell$. As noted earlier (§5.9), the latter assertion is a permission to remove ℓ from the predecessor multisets of other locations. Thus, by iterated application of **FREEONE**, acyclic chains of unreachable blocks can be logically deallocated.

FREEONE can be applied to a reachable location if this location is a temporary root inside a protected section. Our logic thereby allows such a location to be read or written *post mortem*, after it has been logically deallocated. This is made possible by the fact that the points-to assertion survives logical deallocation. This pattern appears, for example, in the verification of Treiber’s stack (§11.5).

Contrary to the logical deallocation rule presented by Moine et al. [2023], our rule does not consume or even mention a points-to assertion for the location ℓ . Indeed, the points-to assertion is not needed to guarantee that the location is unreachable, nor is it needed to prevent a location from being deallocated twice. The size of the deallocated block is obtained in this paper from the “sizeof” assertion, whereas in the previous paper this assertion did not exist, so the size was obtained from a points-to assertion.

6.2 Reasoning Rules for Terms

Figure 21 presents most of the reasoning rules that concern instructions, except for the rules that are specific to protected sections, which are presented later on (§6.3). The reasoning rule **BIND**, which

⁸This property itself holds, under the assumption that enough polling points have been inserted, because as soon as one thread is waiting for space, every thread must eventually reach a large memory allocation instruction or a polling point, where it becomes blocked; and, at this point, this thread must be outside of a protected section.

$$\begin{array}{c}
\text{IFTRUE} \\
\frac{\{\Phi\} \pi: t_1 \{\Psi\}}{\{\Phi\} \pi: \text{if true then } t_1 \text{ else } t_2 \{\Psi\}} \\
\\
\text{IFFALSE} \\
\frac{\{\Phi\} \pi: t_2 \{\Psi\}}{\{\Phi\} \pi: \text{if false then } t_1 \text{ else } t_2 \{\Psi\}} \\
\\
\text{LETVAL} \\
\frac{\{\Phi\} \pi: [v/x]t \{\Psi\}}{\{\Phi\} \pi: \text{let } x = v \text{ in } t \{\Psi\}} \\
\\
\text{PRIM} \\
\frac{v_1 \odot v_2 \xrightarrow{\text{pure}} w}{\{\ulcorner \text{True} \urcorner\} \pi: v_1 \odot v_2 \{\lambda v. \ulcorner v = w \urcorner\}} \\
\\
\text{VAL} \\
\{\ulcorner \text{True} \urcorner\} \pi: v \{\lambda v'. \ulcorner v' = v \urcorner\} \\
\\
\text{POLL} \\
\{\text{outside } \pi\} \pi: \text{poll } \{\lambda(). \text{outside } \pi\} \\
\\
\text{LOAD} \\
\frac{0 \leq i < |\vec{w}| \quad \vec{w}(i) = v}{\left\{ \begin{array}{l} \ell \mapsto_p \vec{w} \\ v \Leftarrow_{p'} \emptyset \end{array} \right\} \pi: \ell[i] \left\{ \begin{array}{l} \ulcorner v' = v \urcorner \\ \lambda v'. \ell \mapsto_p \vec{w} \\ v \Leftarrow_{p'} \{\pi\} \end{array} \right\}} \\
\\
\text{STORE} \\
\frac{0 \leq i < |\vec{w}| \quad \vec{w}(i) = v}{\left\{ \begin{array}{l} \ell \mapsto_1 \vec{w} \\ v' \Leftarrow_{q'}^0 \emptyset \end{array} \right\} \pi: \ell[i] \leftarrow v' \left\{ \begin{array}{l} \ell \mapsto_1 [i:=v']\vec{w} \\ \lambda(). v' \Leftarrow_{q'}^0 \{+\ell\} \\ v \Leftarrow_0 \{-\ell\} \end{array} \right\}} \\
\\
\text{CASSUCCESS} \\
\frac{0 \leq i < |\vec{w}| \quad \vec{w}(i) = v}{\left\{ \begin{array}{l} \ell \mapsto_1 \vec{w} \\ v' \Leftarrow_{q'}^0 \emptyset \end{array} \right\} \pi: \text{CAS } \ell[i] v v' \left\{ \begin{array}{l} \ulcorner b = \text{true} \urcorner \\ \lambda b. \ell \mapsto_1 [i:=v']\vec{w} \\ v' \Leftarrow_{q'}^0 \{+\ell\} \\ v \Leftarrow_0 \{-\ell\} \end{array} \right\}} \\
\\
\text{CASFAILURE} \\
\frac{0 \leq i < |\vec{w}| \quad \vec{w}(i) \neq v}{\left\{ \ell \mapsto_p \vec{w} \right\} \pi: \text{CAS } \ell[i] v v' \left\{ \begin{array}{l} \ulcorner b = \text{false} \urcorner \\ \lambda b. \ell \mapsto_p \vec{w} \end{array} \right\}} \\
\\
\text{FORK} \\
\frac{\text{dom}(M) = \text{locs}(t)}{(\forall \pi'. \{\text{outside } \pi' * M \Leftarrow \{\pi'\} * \Phi\} \pi': t \{\lambda(). \text{outside } \pi'\})} \\
\frac{}{\{\text{outside } \pi * M \Leftarrow \{\pi\} * \Phi\} \pi: \text{fork } t \{\lambda(). \text{outside } \pi\}}
\end{array}$$

Fig. 21. Syntax-directed reasoning rules, without protected-section-specific rules and without BIND

allows reasoning under an evaluation context, is presented after that (§6.4). In every rule, the thread identifier π represents the current thread, that is, the thread about which one is reasoning (§5.1).

IFTRUE, **IFFALSE**, **LETVAL**, **PRIM** and **VAL** are standard rules.

CALLPTR governs calls to (recursive, closed) functions, also known in this paper as code pointers. Its only unusual aspect is the presence of the assertion *outside* π , which ensures that the current thread is currently outside a protected section. The presence of this assertion forbids function calls inside protected sections.

Similarly, **POLL** forbids polling points inside a protected section. Outside of this aspect, a polling point is a no-operation.

ALLOC exhibits three differences with the allocation rule of Separation Logic. First, it requires and consumes *size*(n) space credits, so as to pay for the space occupied by the new block. Second, the presence of the assertion *outside* π forbids allocation inside a protected section. Third, in addition to a points-to assertion for the new block, allocation produces pointed-by-heap and pointed-by-thread

assertions. These assertions indicate that there is initially no pointer from the heap to the new block, and that this new block is a root for the current thread (and only for this thread).

As in standard Separation Logic, **LOAD** requires a (fractional) points-to assertion for the memory location ℓ that is accessed. Furthermore, it requires a pointed-by-thread assertion $v \Leftarrow_p \emptyset$ for the value v that is read from memory. This assertion is updated to $v \Leftarrow_p \{\pi\}$, reflecting the fact that the value v becomes a root for the current thread.

As in standard Separation Logic, **STORE** requires a full points-to assertion $\ell \mapsto_1 \vec{v}$ and produces an updated assertion $\ell \mapsto_1 [i:=v']\vec{v}$. Furthermore, it performs bookkeeping of predecessor multisets, so as to reflect the fact that the value v that was stored in the field $\ell[i]$ is overwritten with the value v' . First, to reflect the *creation* of an edge from ℓ to the value v' , an assertion of the form $v' \Leftarrow_q \emptyset$ is changed to $v' \Leftarrow_q \{+\ell\}$. Here, because ℓ has positive multiplicity in $\{+\ell\}$, the null-fraction invariant requires that q be positive; it cannot be 0. Second, to reflect the *deletion* of an edge from ℓ to the value v , the assertion $v \Leftarrow_0 \{-\ell\}$ appears in the postcondition. As explained earlier (§5.6), this assertion is a permission to remove one occurrence of ℓ from a multiset of predecessors of v .

CASSUCCESS is similar to **STORE**, but returns the Boolean value true rather than the unit value. Because a failed CAS does not modify the heap or create a new root, **CASFAILURE** is standard.

FORK reasons about the operation of spawning a new thread whose code is the term t . This operation must take place outside a protected section. Its impact on roots is as follows. Suppose, for a moment, that fork t is the last instruction in the parent thread. Then, the locations that occur in the term t cease to be roots of the parent thread π and become roots of the child thread π' . The reasoning rule reflects this intuition by updating a group of pointed-by-thread assertions. The iterated pointed-by-thread assertion $M \Leftarrow \{\pi\}$ is taken away from the parent thread, and the updated assertion $M \Leftarrow \{\pi'\}$ is transmitted to the child thread. M is a map of locations to fractions, whose domain is the set $\text{locs}(t)$. This is a form of *trimming*, similar in effect to the rules **TRIMPBTHREAD** and **TRIMINSIDE**.

If fork t is *not* the last instruction in the parent thread, then the user must use the reasoning rules **BIND** and **FORK** in combination. The interaction between the **BIND** rule and the “trimming” rules is discussed later on (§6.4, §6.5).

Still looking at **FORK**, an arbitrary assertion Φ is transmitted from the parent thread to the child thread. The assertion *outside* π' is made available in the child thread, reflecting the fact that a new thread initially runs outside a protected section. The child thread t must be verified with the nontrivial postcondition *outside* π' , thereby disallowing a thread to terminate while inside a protected section.

In our Coq mechanization, the postconditions of many reasoning rules contain a *later credit* [Spies et al. 2022]. Later credits play a role in eliminating the “later” modality. They are orthogonal to the main concern of this paper, namely the analysis of space complexity, so we hide them in the presentation of our reasoning rules. We do explain how later credits are used in our case study of the `async-finish` library (§11.4).

6.3 Reasoning about Protected Sections

Within a protected section, the reasoning rules presented in the previous section (§6.2) can still be used, except for **CALLPTR**, **ALLOC**, and **POLL**, which require the assertion *outside* π . In addition, a number of reasoning rules, shown in Figure 22, specifically concern protected sections.

ENTER allows entering a protected section. This rule transforms the assertion *outside* π into the assertion *inside* $\pi \emptyset$, thereby witnessing that the current thread is now inside a protected section and has no temporary roots.

Conversely, **EXIT** allows exiting a protected section. By consuming the assertion *inside* $\pi \emptyset$, this rule requires the user to prove that the current thread has no remaining temporary roots.

$$\begin{array}{c}
\text{ENTER} \\
\{ \textit{outside } \pi \} \pi : \textit{enter } \{ \lambda(). \textit{inside } \pi \emptyset \} \\
\\
\text{LOADINSIDE} \\
\frac{0 \leq i < |\vec{w}| \quad \vec{w}(i) = v}{\left\{ \begin{array}{l} \ell \mapsto_p \vec{w} \\ \textit{inside } \pi T \end{array} \right\} \pi : \ell[i] \left\{ \lambda v'. \ulcorner v' = v \urcorner * \ell \mapsto_p \vec{w} \right\} \textit{inside } \pi (\textit{locs}(v) \cup T)} \\
\\
\text{STOREDEAD} \\
\frac{0 \leq i < |\vec{w}|}{\left\{ \begin{array}{l} \ell \mapsto_1 \vec{w} \\ \dagger \ell \end{array} \right\} \pi : \ell[i] \leftarrow v' \{ \lambda(). \ell \mapsto_1 [i := v'] \vec{w} \}} \\
\\
\text{CASSUCCESSDEAD} \\
\frac{0 \leq i < |\vec{w}| \quad \vec{w}(i) = v}{\left\{ \begin{array}{l} \ell \mapsto_1 \vec{w} \\ \dagger \ell \end{array} \right\} \pi : \text{CAS } \ell[i] v v' \left\{ \begin{array}{l} \ulcorner b = \text{true} \urcorner \\ \lambda b. \ell \mapsto_1 [i := v'] \vec{w} \end{array} \right\}}
\end{array}$$

Fig. 22. Reasoning rules: protected-section-specific rules

$$\frac{\text{BIND} \quad \text{dom}(M) = \textit{locs}(K) \quad \{ \Phi \} \pi : t \{ \Psi' \} \quad \forall v. \{ M \Leftarrow \{ \pi \} * \Psi' v \} \pi : K[v] \{ \Psi \}}{\{ M \Leftarrow \{ \pi \} * \Phi \} \pi : K[t] \{ \Psi \}}$$

Fig. 23. Reasoning rules: the BIND rule of IrisFit

LOADINSIDE allows reading a value v from a location ℓ in the heap. The locations that appear in the value v become temporary roots of the current thread: the assertion $\textit{inside } \pi T$ is updated to $\textit{inside } \pi (T \cup \textit{locs}(v))$. In contrast with **LOAD**, no pointed-by-thread assertion is required or updated. In fact, the location ℓ or some locations in the set $\textit{locs}(v)$ might be logically deallocated already.

STOREDEAD allows writing a logically deallocated block. The rule requires and updates a points-to assertion. A deallocation witness $\dagger \ell$ is also required. Compared with **STORE**, no pointed-by-heap assertion is required or updated. Indeed, there is no need to do so. Pointed-by-heap assertions keep track of which blocks are reachable via ordinary roots; but, because the block at address ℓ is logically deallocated, it is not reachable via ordinary roots. This is reminiscent of **CLEANPBHEAP**.

Although **STOREDEAD** does not require an “*inside*” assertion, it can be used only inside a protected section. Indeed, the rule applies to a store instruction $\ell[i] \leftarrow v'$, where the address ℓ occurs. This means that ℓ is a root, yet ℓ is also logically deallocated. This is possible only if the current thread is currently inside a protected section. Indeed, outside a protected section, a logically deallocated location cannot be a root: the rule **NO DANGLING ROOT OUT** says so (§5.9).

CASSUCCESSDEAD is analogous to **STOREDEAD**. It concerns a successful CAS instruction on a logically deallocated location. Because a failed CAS does not write anything, the rule **CASFAILURE** can be applied to a logically deallocated location without change.

6.4 Reasoning under Evaluation Contexts

A proof in Separation Logic is traditionally carried out under an unknown context. That is, one reasons about a term t without knowing in what evaluation context K this term is placed. There are specific points in the proof where this unknown context grows and shrinks. As an archetypical example, consider the sequencing construct $\textit{let } x = t_1 \textit{ in } t_2$. To reason about this construct, one first focuses on the term t_1 , thereby temporarily forgetting the frame $\textit{let } x = \square \textit{ in } t_2$, which is pushed onto the unknown context. After the verification of t_1 is completed, this focusing step is reversed: the frame $\textit{let } x = \square \textit{ in } t_2$ is popped and one continues with the verification of t_2 . These focusing and defocusing steps are described by the “**BIND**” rule [Jung et al. 2018b, §6.2].

In our setting, however, a complication arises. An evaluation context contains memory locations. When one applies the **BIND** rule, so as to temporarily forget about this evaluation context, one must still somehow record that these locations are roots. We use pointed-by-thread assertions for this purpose.

Suppose we wish to decompose the sequence $\text{let } x = t_1 \text{ in } t_2$ into a subterm t_1 and an evaluation context $\text{let } x = \square \text{ in } t_2$. For simplicity, let us further assume that $\text{locs}(t_2)$ is a singleton set $\{\ell\}$. This implies that, while t_1 is being executed, the location ℓ is a root. In this specific case, our **BIND** rule takes the following form:

$$\frac{\text{PARTICULAR CASE OF BIND} \quad \text{locs}(t_2) = \{\ell\} \quad \{\Phi\} \pi : t_1 \{\Psi'\} \quad \forall v. \{\ell \Leftarrow_p \{\pi\}\} * \Psi' v \pi : [v/x]t_2 \{\Psi\}}{\{\ell \Leftarrow_p \{\pi\}\} * \{\Phi\} \pi : \text{let } x = t_1 \text{ in } t_2 \{\Psi\}}$$

What is unusual, compared with the standard **BIND** rule of Separation Logic, is that the fractional pointed-by-thread assertion $\ell \Leftarrow_p \{\pi\}$ is required in the beginning, taken away from the user while focusing on the term t_1 , and given back to the user once she is done reasoning about t_1 and ready to reason about t_2 . In other words, this assertion is *forcibly framed out* while reasoning about t_1 .

The assertion $\ell \Leftarrow_p \{\pi\}$ records that ℓ is a root in thread π . By taking it away from the user and by giving it back once she is done reasoning about t_1 , we ensure that the information that “ ℓ is a root in thread π ” is carried up to this point and cannot be prematurely destroyed.

What could go wrong if we did not do this? Then, the user would be allowed to keep the *full* pointed-by-thread assertion $\ell \Leftarrow_1 \{\pi\}$ while reasoning about t_1 . Technically, the user would do so by instantiating Φ with $\ell \Leftarrow_1 \{\pi\}$ in the **BIND** rule. Then, the user would focus on establishing the first premise, $\{\ell \Leftarrow_1 \{\pi\}\} \pi : t_1 \{\Psi'\}$. Now suppose $\ell \notin \text{locs}(t_1)$, that is, ℓ does not occur in t_1 . Then, the user could apply **TRIMPBTTHREAD** to transform the assertion $\ell \Leftarrow_1 \{\pi\}$ into $\ell \Leftarrow_1 \emptyset$. Oops! The assertion $\ell \Leftarrow_1 \emptyset$ means that ℓ is *not* a root. Yet ℓ really *is* still a root, as it occurs in the evaluation context that has been abstracted away, namely $\text{let } x = \square \text{ in } t_2$.

Besides **TRIMPBTTHREAD**, two reasoning rules, namely **FORK** and **TRIMINSIDE**, involve a form of “trimming” of sets of thread identifiers. The soundness of these rules relies on the fact that **BIND** forcibly frames out fractional pointed-by-thread assertions.

The general form of our **BIND** rule, shown in Figure 23, extends this idea to an arbitrary evaluation context K , in which an arbitrary number of locations may occur. Then, for every location in $\text{locs}(K)$, a fractional pointed-by-thread assertion is forcibly framed out.

6.5 Locally Trading Trimming for a Simpler and More Powerful Bind Rule

Forcing pointed-by-thread assertions to be framed out at each application of **BIND** is cumbersome, and can be restrictive, as there are situations where no pointed-by-thread assertion is at hand. (An example appears later on in this section.) Fortunately, such forced framing is unnecessary if the user promises not to exploit any of the trimming rules **TRIMPBTTHREAD**, **FORK** and **TRIMINSIDE**. Thus, we introduce a mode that the user may choose to enter at any time, in which the trimming rules are disabled and, in exchange, a simpler, more powerful **BIND** rule is made available.

We parameterize IrisFit triples with a *mode* m , which is either the normal mode \bowtie or the “no trim” mode \blacktimes . Thus, in general, our triples have the form $\{\Phi\} m/\pi : t \{\Psi\}$, and our custom ghost update has the form $\Phi \xrightarrow{\pi}^V_m \Phi'$. All of the reasoning rules presented so far are polymorphic in the mode, except for the trimming rules **TRIMPBTTHREAD**, **FORK**, and **TRIMINSIDE**, which are disabled in “no trim” mode. For example, **TRIMPBTTHREAD** is written $(\ulcorner \ell \notin V \urcorner * \ell \Leftarrow_p \{\pi\}) \xrightarrow{\pi}^V_{\bowtie} \ell \Leftarrow_p \emptyset$, which prevents its use when in the “no trim” mode \blacktimes . The public specification of a function is

$$\begin{array}{c}
\text{SWITCHMODE} \\
\frac{\{\Phi\} \star/\pi: t \{\Psi\}}{\{\Phi\} m/\pi: t \{\Psi\}}
\end{array}
\qquad
\begin{array}{c}
\text{BINDNO TRIM} \\
\frac{\{\Phi\} \star/\pi: t \{\Psi'\} \quad \forall v. \{\Psi' v\} m/\pi: K[v] \{\Psi\}}{\{\Phi\} m/\pi: K[t] \{\Psi\}}
\end{array}$$

Fig. 24. Reasoning rules: additional mode-specific rules

$$\begin{array}{c}
\lceil \text{True} \rceil \multimap \emptyset \bullet^0 \emptyset \qquad \text{CLOUDEMPTY} \\
\ell \Leftarrow_1 \emptyset \ * \ \ell \Leftarrow_1 L \ * \ \lceil \text{NoNegative}(L) \rceil \multimap (P \cup L) \bullet^{(n+m)} (D \cup \{\ell\}) \qquad \text{CLOUDADD} \\
\lceil P \subseteq D \rceil \ * \ P \bullet^n D \Rightarrow \diamond n \ * \ *_{\ell \in D} \dagger \ell \qquad \text{CLOUDFREE}
\end{array}$$

Fig. 25. Reasoning rules: logical deallocation

always stated in the normal mode. The “no trim” mode is intended for local use, inside the body of a function. It is an adaptation of Moine et al.’s “NOFREE” mode [2023].

Figure 24 presents two new reasoning rules, **SWITCHMODE** and **BINDNO TRIM**, which allow entering “no trim” mode and taking advantage of it.

When read from bottom to top, **SWITCHMODE** lets the user locally enter “no trim” mode, whenever she so wishes, in a subproof. When read from top to bottom, this rule asserts that if a triple holds in “no trim” mode then it also holds in normal mode. Indeed, every reasoning rule that is available in “no trim” mode is available in normal mode as well.

BINDNO TRIM is the standard **BIND** rule of Separation Logic, but imposes a switch to “no trim” mode \star in its left-hand premise. Thus, unlike our **BIND** rule, it does *not* force pointed-by-thread assertions to be framed out. Because of this, it must disable the trimming rules while the user reasons about the subterm t .

We remark that, inside a protected section, one can switch to “no trim” mode without loss of expressive power. Indeed, there, the trimming rules are never needed. **FORK** is forbidden inside protected sections; the effect of **TRIMPBTHREAD** can be simulated by **ADDTemporary**; and all uses of **TRIMINSIDE** can be postponed until the protected section is about to be exited.

At a high level, **BINDNO TRIM** is needed for reasoning about code that, within a protected section, reads or writes in a location after it has been logically deallocated. Indeed, in this case, **BIND** can be too restrictive. To illustrate this case, consider the following code, where we assume that the location r is not accessible via the heap and is not known to any thread other than the current thread:

```
enter ; (let  $x = t$  in  $x + r[0]$ ) ; exit
```

Just after entering the protected section, the user may wish to logically deallocate r , in order to recover the corresponding space credits without waiting for the end of the protected section. In this case, just after entering the protected section, she would use **ADDTemporary** to obtain a pointed-by-thread assertion $r \Leftarrow \emptyset$, then use **FREEONE** to logically deallocate r . **FREEONE** consumes this pointed-by-thread assertion but preserves the points-to assertion for r , which the load instruction $r[0]$ needs. Thereafter, the user may wish to decompose the let construct. Yet, the **BIND** rule cannot be used, as it would require a (fractional) pointed-by-thread assertion for r , which no longer exists, because the fraction 1 was consumed by **FREEONE**. Fortunately, **BINDNO TRIM** is applicable.

<pre> demo \triangleq $\mu_{\text{ptr}}.\lambda[]$. let $x = \text{alloc } 1$ in let $y = \text{alloc } 1$ in $x[0] \leftarrow y$; fork ($x[0] \leftarrow x$); (wait $[x]$)_{ptr} </pre>	<pre> wait \triangleq $\mu_{\text{ptr}}f.\lambda[x]$. if $x[0] = x$ then () else ($f [x]$)_{ptr} </pre>
---	---

Fig. 26. The demo function and its auxiliary function wait

6.6 Logical Deallocation of Cycles

Figure 25 presents our rules for deallocating an unreachable heap *fragment*, as opposed to a single location. This fragment may contain an arbitrary number of heap blocks, which may point to each other in arbitrary ways. In particular, these pointers may form one or more cycles.

These rules make use of the “cloud” assertion $P \blacklozenge^n D$, whose parameters P (for “predecessors”) and D (for “domain”) are sets of locations, and whose parameter n is a natural integer. This assertion means that the memory blocks at locations D have total size n , that the locations D are not roots in any thread, and that these locations can be reached only via the locations P . We refer to P also as the *entry points* of the cloud.

If $P \subseteq D$ holds, then the locations in the set D are reachable only via D itself. In other words, the set D is closed under predecessors. This means that the locations in the set D are in fact *unreachable*, and can safely be logically deallocated. This explains the side condition $P \subseteq D$ in the logical deallocation rule **CLOUDFREE**. We do not require $P \subseteq D$ to hold at all times: while constructing large “cloud” assertions out of smaller “cloud” assertions, one must allow the sets P and D to be unrelated.

Figure 25 presents two cloud construction rules as well as the logical deallocation rule, which consumes a cloud.

Out of nothing, **CLOUDEMPTY** creates an empty cloud $\emptyset \blacklozenge^0 \emptyset$.

CLOUDADD adds the memory block at location ℓ to an existing cloud $P \blacklozenge^n D$. This consumes the full pointed-by-thread assertion $\ell \Leftarrow_1 \emptyset$, which guarantees that ℓ is not a root in any thread, and the full pointed-by-heap assertion $\ell \Leftarrow_1 L$, which guarantees that L contains all of the predecessors of the location ℓ in the heap. A “*sizeof*” assertion determines the size m of the memory block at address ℓ . **CLOUDADD** produces an extended cloud, where L is added to the cloud’s entry points, m is added to the cloud’s size, and ℓ is added to the cloud’s domain.

CLOUDFREE logically deallocates a cloud that is closed under predecessors, that is, a cloud such that $P \subseteq D$ holds. The “cloud” assertion is consumed. In exchange for it, the rule produces n space credits, where n is the size of the cloud. Furthermore, it produces a deallocation witness for every location in the cloud.

The rule **FREEONE** that was presented earlier (§6.1) is easily derived from the rules in Figure 25.

7 INTERLUDE: VERIFYING A SMALL EXAMPLE

Before diving into the soundness statements of IrisFit, extensions of IrisFit, and case studies, let us showcase how one proves a small program, with concurrency, but without protected sections.

Demo Program. The program, named *demo*, appears in Figure 26. Following standard practice, we write $t_1 ; t_2$ as sugar for $\text{let } x = t_1 \text{ in } t_2$ where $x \notin \text{fv}(t_2)$. The demo program proceeds as follows. First, two blocks of size 1 are allocated. Their addresses, say ℓ_x and ℓ_y , are bound to the variables x and y . Then, the address ℓ_y is stored inside the block at address ℓ_x . Next, a new thread is forked.

This new thread executes the store instruction $\ell_x[0] \leftarrow \ell_x$: that is, it writes the address ℓ_x into the block at address ℓ_x , creating a cyclic pointer from this block to itself. Meanwhile, the main thread launches an active waiting loop, which runs until it observes that ℓ_x points to itself—that is, until it observes the effect of the store instruction executed by the child thread. This loop is implemented by means of the auxiliary function `wait`.

Specification. Our goal is to establish that executing `demo` requires two words of memory and that these two words are recovered once `demo` terminates. In `IrisFit`, this specification is expressed by this triple:

$$\{\textit{outside } \pi * \diamond 2\} \pi: (\textit{demo } [])_{\textit{ptr}} \{\lambda(). \textit{outside } \pi * \diamond 2\}.$$

The challenge is to prove that, by the time the main thread completes, the blocks ℓ_x and ℓ_y can be logically deallocated. To this end, we must argue that, upon completion of the main thread, the blocks ℓ_x and ℓ_y are not pointed by any other heap block and that they are no longer roots in any thread. We begin by presenting the high-level arguments.

Intuition for the Proof. Let us focus in turn on the addresses ℓ_x and ℓ_y . For each of these locations, let us examine in turn which heap blocks point to it and in which threads it is a root.

As soon as the child thread executes the instruction $\ell_x[0] \leftarrow \ell_x$, the location ℓ_x becomes stored in the heap. However, it is stored in the block ℓ_x itself, and nowhere else. The block ℓ_x is never pointed to by another heap block.

The location ℓ_x ceases to be a root for the main thread when the test $\ell_x[0] = \ell_x$ in `wait` succeeds. Indeed, the call to `wait` is the last instruction of the function `demo`, and the “then” branch in the function `wait` contains the trivial instruction `()`, which does not mention ℓ_x . Furthermore, by means of the invariant that is described further on, one can prove that the success of the test $\ell_x[0] = \ell_x$ guarantees that the child thread has executed its store instruction. Therefore, when the main thread completes, ℓ_x is no longer a root for the child thread. As a result, at that point, ℓ_x is no longer a root at all.

The location ℓ_y is initially stored in the block ℓ_x . However, as soon as the child thread executes its store instruction, which overwrites ℓ_y with ℓ_x , the location ℓ_y ceases to appear in the heap. Furthermore, as explained earlier, the main thread does not complete until the child thread has executed this store instruction. Therefore, when the main thread completes, the address ℓ_y is no longer stored in the heap.

The location ℓ_y is not mentioned in the child thread, and is not read from the heap by this thread, hence ℓ_y is never a root for the child thread. Furthermore, clearly, once the main thread completes, ℓ_y is no longer a root for the main thread.⁹

As mentioned above, we exploit an *Iris invariant* to transfer knowledge between the child thread and the main thread. An invariant can be thought of as a description of the states of a state machine, which governs how threads interact and what resources they exchange. The transitions of this state machine are not mentioned in the invariant: they are implicit. In a proof, an invariant is typically *opened*, *analyzed* so as determine which state (or states) may be current, then *closed* in the same state or in a new state.

⁹A look at the code suggests that the *variable* y ceases to be a root immediately after the instruction $x[0] \leftarrow y$ is executed. Indeed, beyond this point in the code, the *variable* y does not occur. However, this remark is irrelevant. The truly relevant question is not which *variables* are roots, but which *memory locations* are roots. Here, as long as the waiting loop runs, the *location* ℓ_y may be a root for the main thread, because the load instruction $\ell_x[0]$ in the function `wait` can read ℓ_y from the heap and return ℓ_y . Thus, ℓ_y definitively ceases to be a root for the main thread only once the load instruction $\ell_x[0]$ has returned ℓ_x .

The invariant involved in reasoning about demo is the following:

$$I \triangleq \begin{array}{l} \ell_x \mapsto_1 [\ell_y] \\ \vee \ell_x \mapsto_1 [\ell_x] * \ell_x \leftarrow_1 \{+\ell_x\} * \ell_x \leftarrow_{\frac{1}{2}} \emptyset * \ell_y \leftarrow_1 \emptyset \\ \vee \dagger \ell_y \end{array}$$

The first disjunct of the invariant I represents the initial state, in which the block at address ℓ_x contains a pointer to the block at address ℓ_y . The second disjunct represents the intermediate state in which the child thread has terminated and given up its resources, but the main thread has not yet observed that ℓ_x now points to itself. The third disjunct $\dagger \ell_y$ corresponds to the final state, which is reached at the point where the block ℓ_y can be logically deallocated. (The presence of a disjunct of the form $\dagger \ell_y$ makes I an example of a liveness-based cancellable invariant, in the sense of §5.10.)

Formal Arguments Involved in the Proof. Let us begin by reasoning on the instructions from demo. First, we apply the rule **CALLPTR** and enter the function body. Then, we face two allocations of blocks of size 1. We apply **SPLITJOINSC** to the assertion $\diamond 2$ to obtain $\diamond 1 * \diamond 1$. To be more precise, the first allocation that we face lies under a let binding. To enter the left-hand side of this let binding, we use **BINDNOTRIM**. Then, we apply **ALLOC**: we lose one space credit, we name the resulting location ℓ_x , and we obtain the assertions $\ell_x \mapsto_1 [()]$ and $\ell_x \leftarrow_1 \{\pi\}$ and $\ell_x \leftarrow_1 \emptyset$. We then use **LETVAL** to substitute ℓ_x for x in the remaining term. We repeat the exact same three steps to reason about the second allocation and name its result ℓ_y . We apply the rule **SIZEOFPOINTSTO** to obtain *sizeof* ℓ_y 1. We deduce from the two points-to assertions that ℓ_x and ℓ_y are distinct. Because ℓ_y is never read or written, we throw away its points-to assertion. At this point, we hold the permissions *sizeof* ℓ_y 1 and $\ell_y \leftarrow_1 \{\pi\}$ and $\ell_y \leftarrow_1 \emptyset$.

The term that remains to reason about is:

$$\ell_x[0] \leftarrow \ell_y; \text{fork}(\ell_x[0] \leftarrow \ell_x); (\text{wait}[\ell_x])_{\text{ptr}}.$$

We apply the rule **BINDNOTRIM** to focus on the store instruction $\ell_x[0] \leftarrow \ell_y$. For this instruction, we apply the rule **STORE**, thereby trading the assertion $\ell_x \mapsto_1 [()] * \ell_y \leftarrow_1 \emptyset$ for the assertion $\ell_x \mapsto_1 [\ell_y] * \ell_y \leftarrow_1 \{+\ell_x\}$.

We now reach the fork instruction. We initialize the aforementioned invariant I by entering the initial state, that is, by providing the assertion $\ell_x \mapsto_1 [\ell_y]$. At this stage, we need to split the assertion $\ell_x \leftarrow_1 \{\pi\}$ in two halves, because we will need one half to witness a root held by the context, and one half to transmit to the child thread. Concretely, we apply **FRACPBTHREAD** to this assertion and obtain $\ell_x \leftarrow_{\frac{1}{2}} \{\pi\} * \ell_x \leftarrow_{\frac{1}{2}} \{\pi\}$. To focus on the fork instruction, we apply the rule **BIND**. (We cannot exploit **BINDNOTRIM** because **FORK** involves trimming.) Applying **BIND** requires us to temporarily give up the fractional pointed-by-thread permission $\ell_x \leftarrow_{\frac{1}{2}} \{\pi\}$.

We next apply the rule **FORK**. Thereafter, we focus on the child thread, whose code is $\ell_x[0] \leftarrow \ell_x$. We name the new thread identifier π' . To the child thread, we transmit the invariant \boxed{I} as well as the assertions $\ell_x \leftarrow_1 \emptyset$ and $\ell_y \leftarrow_1 \{+\ell_x\}$. The application of **FORK** updates the assertion $\ell_x \leftarrow_{\frac{1}{2}} \{\pi\}$ into $\ell_x \leftarrow_{\frac{1}{2}} \{\pi'\}$, which is transmitted to the child thread.

To reason about the store instruction $\ell_x[0] \leftarrow \ell_x$, we first open the invariant \boxed{I} . We eliminate the third disjunct $\dagger \ell_y$ by using **DEADPBHEAP**. We also eliminate the second disjunct $(\dots * \ell_y \leftarrow_1 \emptyset)$ by using **JOINPBHEAP** together with fractional reasoning about the pointed-by-heap assertion for ℓ_y (a fraction cannot exceed 1). Only the first disjunct remains: so, we acquire the assertion $\ell_x \mapsto_1 [\ell_y]$. Then, we can apply the rule **STORE**, trading $\ell_x \mapsto_1 [\ell_y] * \ell_x \leftarrow_1 \emptyset$ for $\ell_x \mapsto_1 [\ell_x] * \ell_x \leftarrow_1 \{+\ell_x\} * \ell_y \leftarrow_1 \{-\ell_x\}$. We use **JOINPBHEAP** to transform $\ell_y \leftarrow_1 \{+\ell_x\} * \ell_y \leftarrow_1 \{-\ell_x\}$ into $\ell_y \leftarrow_1 \emptyset$. A store instruction returns the unit value $()$, so, after the store takes place, we use **CONSEQUENCE** and **TRIMPBTHREAD** to update the assertion $\ell_x \leftarrow_{\frac{1}{2}} \{\pi'\}$ into $\ell_x \leftarrow_{\frac{1}{2}} \emptyset$, witnessing that ℓ_x is not

a root of the child thread any more. We then close the invariant I by giving up all of the assertions at hand, forming the second disjunct of I .

We now turn our attention back to the main thread, whose sole remaining instruction is: $(\text{wait } [\ell_x])_{\text{ptr}}$. We exploit the fact that ℓ_y is not a root for this term, together with the rules **CONSEQUENCE** and **TRIMPBTHREAD**, to obtain $\ell_y \Leftarrow_1 \emptyset$. At this point, there remains to prove the following triple, which corresponds to a specification of the auxiliary function `wait`.

$$\{ \text{outside } \pi * \boxed{I} * \text{sizeof } \ell_y 1 * \ell_y \Leftarrow_1 \emptyset * \ell_x \Leftarrow_{\frac{1}{2}} \{\pi\} \} \pi : (\text{wait } [\ell_x])_{\text{ptr}} \{ \lambda(). \text{outside } \pi * \diamond 2 \}$$

We establish this triple using Löb induction [Jung et al. 2018b]. We use **CALLPTR** and enter the function body. We face the term: $\text{if } \ell_x[0] = \ell_x \text{ then } () \text{ else } (f \ [\ell_x])_{\text{ptr}}$. We first use **BINDNOTRIM** to focus on the condition of the if statement, $\ell_x[0] = \ell_x$, and use **BINDNOTRIM** again to focus on the load $\ell_x[0]$.

By making use of the information stored in the invariant, we will now prove that this load must return ℓ_x or ℓ_y . We open the invariant \boxed{I} and perform a case analysis on I , giving rise to three cases. In each of the first two cases, we explain how to close the invariant; the third case, we rule out.

- (1) In the first case, we have $\ell_x \mapsto_1 [\ell_y]$. We apply **LOAD**, which updates the assertion $\ell_y \Leftarrow_1 \emptyset$ into $\ell_y \Leftarrow_1 \{\pi\}$, and close the invariant in the same state as it was just opened.
- (2) In the second case, we have $\ell_x \mapsto_1 \ell_x * \ell_x \Leftarrow_1 \{+\ell_x\} * \ell_x \Leftarrow_{\frac{1}{2}} \emptyset * \ell_y \Leftarrow_1 \emptyset$. This case involves logical deallocation. First, we use **CONSEQUENCE** and **FREEONE** to logically deallocate ℓ_y . We obtain one space credit and a deallocation witness $\dagger \ell_y$. Then, we apply **LOAD**, which updates the assertion $\ell_x \Leftarrow_1 \emptyset$ into $\ell_x \Leftarrow_1 \{\pi\}$. We close the invariant using the third disjunct, that is, by providing $\dagger \ell_y$.
- (3) In the third case, we have $\dagger \ell_y$. This case is eliminated by using **DEADPBTHREAD**. Indeed, we hold the pointed-by-thread assertion $\ell_y \Leftarrow_1 \emptyset$ therefore ℓ_y cannot be deallocated.

Let ℓ_z be the result of the load $\ell_x[0]$. If we went through case (1) above, then we have $\ell_z = \ell_y$. If we went through case (2), then we have $\ell_z = \ell_x$. We apply the rule **PRIM** to reason about the test that compares ℓ_z with ℓ_x . Let us show, in each of the two cases, how to conclude the proof.

- (1) Case $\ell_z = \ell_y$. Because ℓ_y and ℓ_x are distinct addresses, the test must evaluate to false. We apply **IFFALSE** and enter the second branch of the conditional: we now face the recursive call $(f \ [\ell_x])_{\text{ptr}}$. Using **CONSEQUENCE** and **TRIMPBTHREAD**, we trim the pointed-by-thread assertion for ℓ_y , changing the assertion $\ell_y \Leftarrow_1 \{\pi\}$ back into $\ell_y \Leftarrow_1 \emptyset$. We conclude by applying the induction hypothesis.
- (2) Case $\ell_z = \ell_x$. This time, the test must evaluate to true. We apply **IFTRUE** and enter the first branch of the conditional. There remains to establish the following triple:

$$\left\{ \begin{array}{l} \text{outside } \pi * \diamond 1 \\ \ell_x \mapsto_1 [\ell_x] \\ \ell_x \Leftarrow_{\frac{1}{2}} \{\pi\} * \ell_x \Leftarrow_{\frac{1}{2}} \emptyset \\ \ell_x \Leftarrow_1 \{+\ell_x\} \end{array} \right\} \pi : () \{ \lambda(). \text{outside } \pi * \diamond 2 \}$$

Using **FRACPBTHREAD**, we obtain $\ell_x \Leftarrow_1 \{\pi\}$. Next, using **CONSEQUENCE** and **TRIMPBTHREAD**, we obtain $\ell_x \Leftarrow_1 \emptyset$. In order to logically deallocate the single-cell cycle ℓ_x , we use the cloud rules presented in §6.6 to construct the cloud assertion $\{\ell_x\} \clubsuit^1 \{\ell_x\}$. Then, using **CONSEQUENCE** and **CLOUDFREE**, we obtain one space credit. Finally, using **SPLITJOINSC**, we join the two space credits. We conclude using **VAL**.

$$\begin{array}{c}
\text{NOTSTUCKVAL} \\
\frac{\theta(\pi) = (v, \text{Out})}{\text{NotStuck}_S(\theta, \sigma) \pi}
\end{array}
\qquad
\begin{array}{c}
\text{NOTSTUCKSTEP} \\
\frac{c \xrightarrow{\text{enabled actions}}_{\pi} c'}{\text{NotStuck}_S c \pi}
\end{array}
\qquad
\begin{array}{c}
\text{SAFE} \\
\frac{\forall \pi. \text{Enabled}_S c \pi}{\text{Safe}_S c} \\
\implies \text{NotStuck}_S c \pi
\end{array}$$

Fig. 27. Predicates used in the statement of the safety theorem

8 SAFETY AND LIVENESS

In this section, we state several theorems about programs that have been verified using IrisFit. In short, we wish to establish three properties, namely *safety* (no thread can crash), *liveness* (no thread can be blocked forever), and *bounded space consumption* (the size of the heap cannot exceed a certain bound). We first state these properties about the default semantics (§4.2.9), then discuss the growing semantics (§4.2.10) and the oblivious semantics (§4.2.7).

The *safety* theorem (§8.1) guarantees that no thread crashes. More precisely, it states that if a thread is enabled (§4.2.8) then this thread is not stuck: either it has reached a value or it can make a step.

The *liveness* theorem (§8.2) guarantees that no thread can be blocked forever. More precisely, under the assumption that there is a polling point in front of every function call, we prove that every thread is eventually enabled. Furthermore, we prove that inserting a polling point in front of every function call preserves safety. Thus, after a source program without polling points has been verified with IrisFit, one can let a compiler automatically insert polling points, and obtain both safety and liveness for this instrumented program.

In the default semantics of LambdaFit, by design, the size of the heap cannot exceed the limit S (Lemma 4.2). Therefore, the *bounded space consumption* property comes for free.

With respect to the growing semantics, we are able to establish similar results. In this semantics, the heap size limit can grow at runtime, so the statement of the *bounded space consumption* property must be slightly relaxed (§8.3).

All of our results about the default and growing semantics follow from a single *core soundness* theorem stated with respect to the oblivious semantics (§8.4). This theorem spells out the guarantee that is offered by IrisFit when a LambdaFit program is executed with blocking instructions ignored and garbage collection disabled.

All of our results are mechanized using the Coq proof assistant. For more details about our proofs, the reader is referred to our mechanization [Moine 2024b] and to the first author’s dissertation [Moine 2024a].

8.1 Safety

A concurrent Separation Logic typically comes with a safety guarantee, formulated in the form: “no thread can crash”. A more precise statement is: “always, every thread is not stuck”. In other words, in every reachable configuration of the system, every thread either has terminated or is able to make a reduction step. A thread that has not reached a value and is unable to make a step is *stuck*: by convention, this is considered an undesirable situation, akin to a crash.

In our setting, however, this statement must be amended, because LambdaFit has blocking instructions. A blocking instruction is sometimes *disabled* (§4.2.8), therefore unable to make a step; yet, this situation is not considered a crash.

Our amended safety guarantee is qualified as follows: “always, every *enabled* thread is not stuck”. A thread that is not enabled is considered blocked: this is a normal situation.

$$\begin{array}{c}
\text{HOLDSNOW} \\
\frac{P c}{\text{AfterAtMost } (\longrightarrow) n P c}
\end{array}
\qquad
\begin{array}{c}
\text{HOLDSAFTER} \\
\frac{\exists c'. c \longrightarrow c' \quad \forall c'. c \longrightarrow c' \implies \text{AfterAtMost } (\longrightarrow) n P c'}{\text{AfterAtMost } (\longrightarrow) (n+1) P c}
\end{array}$$

$$\begin{array}{c}
\text{ALWAYS} \\
\frac{\forall c'. c \longrightarrow^* c' \implies P c'}{\text{Always } (\longrightarrow) P c}
\end{array}
\qquad
\begin{array}{c}
\text{EVENTUALLY} \\
\frac{\text{AfterAtMost } (\longrightarrow) n P c}{\text{Eventually } (\longrightarrow) P c}
\end{array}$$

Fig. 28. Temporal logic predicates

Figure 27 defines a few auxiliary predicates that appear in the statement of the safety theorem. The proposition $\text{NotStuck}_S c \pi$ means that, in the configuration c , the thread identified by π is not stuck. It is defined by two rules. **NOTSTUCKVAL** states that if a thread has reached a value and is outside a protected section, then it is not stuck. (Terminating inside a protected section is forbidden.) **NOTSTUCKSTEP** states that if a thread can take a step, then it is not stuck. The proposition $\text{Safe}_S c$, defined by the rule **SAFE**, means that no enabled thread in the configuration c is stuck.

The proposition $\text{Always } (\longrightarrow) P c$, which is defined by the rule **ALWAYS** in Figure 28, means that every configuration that is reachable from the configuration c via the reduction relation \longrightarrow satisfies the predicate P .

The safety theorem (Theorem 8.1) can be read as follows. Suppose that the program t has been verified using IrisFit, with an arbitrary identifier π for the main thread, under the precondition $\diamond S * \text{outside } \pi$ and the postcondition $\text{outside } \pi$. The precondition provides S space credits and guarantees that the main thread initially runs outside a protected section. The postcondition forbids termination inside a protected section. Then, the initial configuration $\text{init}(t)$ (§4.2.2) is *always safe*: that is, beginning in this configuration, running the program under the default semantics with heap limit S cannot reach a configuration where a thread is stuck.

THEOREM 8.1 (SAFETY). *Assume that the following triple holds:*

$$\forall \pi. \quad \{ \diamond S * \text{outside } \pi \} \pi : t \{ \lambda _ . \text{outside } \pi \}$$

Then $\text{Always } (\xrightarrow{\text{defaults}}) \text{Safe}_S (\text{init}(t))$ holds.

Although this theorem mentions S , the meaning of a triple is independent of S . Therefore, the reasoning rules are independent of S as well. One can verify a program component without mentioning S and without knowing its value. A concrete value of S must be chosen and fixed only when Theorem 8.1 is applied to a complete (closed) program.

8.2 Liveness

The safety theorem guarantees that no thread can crash, but allows a thread to become blocked. Therefore, a liveness guarantee is also desirable: one would like to be assured that *always, every thread is eventually enabled*. In other words, there is no execution scenario where, past a certain point, a thread remains forever blocked (i.e., is never enabled).

In fact, we are able to offer a stronger guarantee: we prove that *always, eventually, every allocation fits*. In other words, in every execution scenario, infinitely often, the system reaches a point where

no allocation request is blocked due to a lack of memory. This property is indeed stronger, because it guarantees that, at that point, *all* threads are simultaneously enabled.¹⁰

However, our liveness guarantee is subject to a condition: the program must contain *enough polling points*. To see why this is necessary, imagine a program where thread *A* is blocked on a large allocation request and thread *B* is running in an infinite loop, without allocating memory or encountering a polling point. Then, there exists a scenario where thread *B* runs forever, the garbage collector is never invoked, and thread *A* never becomes enabled. Thus, the desired liveness property does not hold. However, suppose that a polling point is inserted in the loop: thread *B* is not allowed to proceed past this polling point. Then, in every scenario, a garbage collection step eventually takes place, at which time both thread *A* and thread *B* become unblocked.

How can one tell whether a program has enough polling points? Or how can one tell where polling points must be inserted so that the program has enough polling points? We propose a simple approach, which is to *insert a polling point in front of every function call*.¹¹ This ensures that every thread must reach a polling point in a bounded number of steps. Up to an administrative side condition,¹² we prove that this polling point insertion strategy preserves safety and ensures liveness. We refer to this polling point insertion strategy as *addpp*. Thus, if *t* is a term, then *addpp(t)* is the term obtained by inserting a polling point in front of every function call in the term *t*.

Figure 28 introduces several auxiliary predicates that appear in the statement of the liveness theorem. The proposition *AfterAtMost* (\longrightarrow) *n P c* means that, out of the configuration *c*, every execution path via the reduction relation \longrightarrow reaches, *in at most n steps*, a configuration that satisfies *P*. This proposition is inductively defined by the rules **HOLDSNOW** and **HOLDSAFTER**. **HOLDSAFTER** guarantees not only that the predicate continues to hold after every possible step, but also that there exists such a step. The proposition *Eventually* (\longrightarrow) *P c*, defined by the rule **EVENTUALLY**, means that *in a bounded number of steps*, out of the configuration *c*, every execution path reaches a configuration that satisfies *P*. It is defined via an existential quantification over *n*.¹³

The following theorem combines a safety guarantee and a liveness guarantee. It states that if the program *t* has been verified using *IrisFit*, under the exact same conditions as in Theorem 8.1, then the program *addpp(t)*, in which enough polling points have been inserted, is safe and live.

THEOREM 8.2 (COMBINED SAFETY AND LIVENESS AFTER POLLING POINT INSERTION). *Let t be a term in administrative normal form. Assume that the following triple holds:*

$$\forall \pi. \{ \diamond S * \text{outside } \pi \} \pi : t \{ \lambda _ . \text{outside } \pi \}$$

¹⁰By Lemma 4.1, the property *always, eventually, every allocation fits* implies that *always, eventually, all threads are enabled* at the same time; which, in turn, implies that *always, every thread is eventually enabled*.

¹¹*LambdaFit* does not have loops: instead, loops must be simulated via tail-recursive functions. Thus, inserting a polling point in front of every function call effectively implies inserting a polling point inside every loop as well. Incidentally, because function calls are forbidden inside protected sections, a polling point is never inserted into a protected section, satisfying our restriction that polling points in protected sections are forbidden. Our polling point insertion strategy is loosely inspired by the (undocumented) polling point insertion strategy of the OCaml compiler. The OCaml compiler inserts a polling point at the beginning of every function (except possibly small leaf functions), inside every loop, and views memory allocation instructions as polling points.

¹²Prior to inserting polling points, we require the program to be in administrative normal form (ANF). That is, in every function call, we require the function itself and the actual arguments to be variables or values, as opposed to arbitrary expressions. This guarantees that the polling point that is inserted in front of the function call is executed *after* the actual arguments have been computed and *just before* the function is invoked.

¹³We propose a strong definition of *Eventually*, whose quantifier prefix is of the form $\exists \forall$: “there exists *n* such that every execution path reaches in at most *n* steps a point where *P* is satisfied.” A weaker definition would involve a quantifier prefix of the form $\forall \exists$: “every execution path eventually reaches a point where *P* is satisfied.” This alternative definition is strictly weaker, because an infinitely branching tree where each branch is finite does not necessarily have finite depth [Bertot and Castéran 2004]. Our reduction relations have infinite non-determinism because memory allocation picks an arbitrary fresh address.

$$\frac{\text{EveryAllocFits}_S c}{\text{EveryAllocFitsPair}(S, c)} \qquad \frac{\text{Safe}_S c}{\text{SafePair}(S, c)} \qquad \frac{S \leq S'}{\text{LimitsAtMost } S'(S, c)}$$

Fig. 29. Predicates used in the statements of soundness for the growing semantics

Let t' stand for the term $\text{addpp}(t)$. Then, the following propositions hold:

- (1) Always $(\xrightarrow{\text{defaults}})$ $\text{Safe}_S(\text{init}(t'))$
- (2) Always $(\xrightarrow{\text{defaults}})$ (Eventually $(\xrightarrow{\text{defaults}})$ $\text{EveryAllocFits}_S(\text{init}(t'))$).

This statement reflects how we envision the practical use of IrisFit. We expect the user to verify a program t in which polling points have not yet been inserted. Thus, the user need not know where polling points will be placed. The uninstrumented verified program t enjoys safety but not necessarily liveness. Nevertheless, the theorem guarantees that, once enough polling points have been inserted, the program enjoys both safety and liveness.

Although Theorem 8.2 makes use of polling points and fixes a specific polling point insertion strategy to ensure liveness, namely addpp , we do support other approaches. Our mechanization [Moine 2024b] includes a more general liveness theorem that splits the burden into (1) proving that always, at least one thread or the GC can take a step, and (2) proving that always, eventually, the program crashes or every allocation fits. Regarding addpp , under the assumption that the original program t has been verified, we prove that both properties hold for the transformed program $\text{addpp}(t)$. The more general liveness theorem and the proof that it implies Theorem 8.2 are presented in the first author's dissertation [Moine 2024a, §7.5.2].

8.3 Safety and Liveness for the Growing Semantics

For the growing semantics, we establish the following theorem, whose general structure is the same as that of Theorem 8.2. It uses several auxiliary predicates whose definitions appear in Figure 29.

THEOREM 8.3 (COMBINED SAFETY AND LIVENESS AFTER POLLING POINT INSERTION). *Let t be a term in administrative normal form. Assume that the following triple holds:*

$$\forall \pi. \{ \diamond S * \text{outside } \pi \} \pi : t \{ \lambda _ . \text{outside } \pi \}$$

Let t' stand for the term $\text{addpp}(t)$. Let ρ stand for the initial state $(0, \text{init}(t'))$, where the heap limit is 0, the heap is empty, and the program t' is ready to run. Then, the following propositions hold:

- (1) Always $(\xrightarrow{\text{growing}})$ $\text{SafePair } \rho$
- (2) Always $(\xrightarrow{\text{growing}})$ (Eventually $(\xrightarrow{\text{growing}})$ $\text{EveryAllocFitsPair}(\rho)$)
- (3) Always $(\xrightarrow{\text{growing}})$ ($\text{LimitsAtMost}(\text{grow}(S))$) ρ .

This theorem states that if the program t has been verified under the precondition $\diamond S$ then the program $\text{addpp}(t)$, in which polling points have been inserted, is safe (no thread can crash—item 1), is live (always, eventually, all threads are enabled—item 2), and never needs more than $\text{grow}(S)$ words of memory (item 3). Indeed, item 3 is a bounded space consumption property: it states that, always, the current heap limit is at most $\text{grow}(S)$. By Lemma 4.5, this implies that the size of the heap, too, is at most $\text{grow}(S)$.

In summary, even though the current heap limit is automatically increased when the runtime system finds that the current limit is too low, if one has (statically) verified (using IrisFit) that the program needs at most S words of memory, then (at runtime) the heap size is bounded by $\text{grow}(S)$. If one sets for example $\text{grow}(S) = \max(2S, 1)$, as suggested earlier (§4.2.10), then one finds that the heap size is $O(S)$.

$$\begin{array}{c}
\text{NOTSTUCKOBLIVIOUSVAL} \\
\theta(\pi) = (v, \text{Out}) \\
\hline
\text{NotStuckOblivious } (\theta, \sigma) \pi
\end{array}
\qquad
\begin{array}{c}
\text{NOTSTUCKOBLIVIOUSSTEP} \\
c \xrightarrow{\text{action}}_{\pi} c' \\
\hline
\text{NotStuckOblivious } c \pi
\end{array}$$

Fig. 30. Predicates used in the statement of the Core Soundness theorem

Our motivation for proposing the growing semantics is that it does not require a suitable value of the limit to be known before execution begins. Indeed, imagine that the program t has *not* been fully verified. Then, one does not know what value S is large enough to guarantee that the triple $\forall \pi. \{\diamond S * \text{outside } \pi\} \pi : t \{\lambda _ . \text{outside } \pi\}$ holds. Nevertheless, under the optimistic assumption that such a value S exists, one can be assured that running the program under the growing semantics will not require more than $\text{grow}(S)$ words of memory.

8.4 Core Soundness

A provocative yet fundamental remark is that IrisFit has nothing to do with garbage collection. Indeed, its deallocation rule is purely logical. More generally, its reasoning rules are independent of *when* garbage collection takes place, or *whether* it takes place at all. In reality, IrisFit is concerned with the *live heap size* of a program, that is, the sum of the sizes of the reachable blocks.

Our earlier results, namely Theorems 8.1, 8.2, and 8.3, follow from a *core soundness* result, which is expressed with respect to the *oblivious semantics*, a semantics in which no garbage collection takes place and no instructions are blocking (§2.1, §2.3, §4.2.7).

In this setting, we must redefine what it means for a thread to be *not stuck*. The proposition $\text{NotStuckOblivious } c \pi$, defined in Figure 30, serves this purpose. A thread is not stuck if either it has reached a value outside a protected section or it can make a step.

Let us write $\text{livesize}(R, \sigma)$ for the total size of the fragment of the store σ that is reachable from the set of roots R . Let us write $\text{livesize}(c)$ for the live heap size of the configuration c . It is defined by $\text{livesize}((\theta, \sigma)) = \text{livesize}(\text{locs}(\theta), \sigma)$.

Our core soundness theorem states that always (with respect to the oblivious semantics), the following two properties hold. First, no thread is stuck. Furthermore, if every thread is currently outside a protected section, then the live heap size is at most S , where S is the number of space credits that was granted when the program was statically verified.

THEOREM 8.4 (CORE SOUNDNESS). *Assume that the following triple holds:*

$$\forall \pi. \{\diamond S * \text{outside } \pi\} \pi : t \{\lambda _ . \text{outside } \pi\}$$

Then, for every configuration c such that $\text{init}(t) \xrightarrow{\text{oblivious}}^ c$,*

- (1) *for every identifier π of a thread in c , the property $\text{NotStuckOblivious } c \pi$ holds;*
- (2) *AllOutside c implies $\text{livesize}(c) \leq S$.*

This statement may seem surprisingly weak, as it offers no guarantee about $\text{livesize}(c)$ at a time where $\text{AllOutside } c$ does not hold, that is, at a time where at least one thread is inside a protected section. Moreover, this statement offers just a safety guarantee; it offers no liveness guarantee. Nevertheless, this core soundness theorem is sufficiently strong to derive Theorems 8.1, 8.2, and 8.3, which express the purpose of our logic in a different manner.

Our internal definition of IrisFit triples [Moine 2024b] is relative to the oblivious semantics. The proof of Theorem 8.4, as well as the proofs of our reasoning rules, involve the oblivious semantics only. Thus, in many of our proofs, there is no need for us to reason about garbage collection or about the distinction between enabled and disabled reduction steps.

9 CLOSURES

As explained earlier (§2.7), LambdaFit does not have primitive closures. Instead, we define *closure construction* $\mu_{\text{clo}}f. \lambda \vec{x}. t$ and *closure invocation* $(\ell \vec{u})_{\text{clo}}$ as macros, which expand to sequences of primitive LambdaFit instructions. These macros implement *flat closures* [Appel 1992, Chapter 10]. That is, a closure is represented as a record whose fields store a code pointer (at offset 0) and a series of values (at offset 1 and beyond). The implementation of these macros (§9.2) is the same as in our earlier paper [Moine et al. 2023]. Our reasoning rules for closure construction, invocation, and deallocation are improved versions of the rules presented in our earlier paper [Moine et al. 2023]. In particular, they involve *pending substitutions* (§9.3). The main improvement is that the assertions that describe closures are now *persistent*. From an end user’s point of view, this makes closures much easier to work with. Internally, this is made possible by using *liveness-based cancellable invariants* (§5.9).

Our reasoning rules for closures are abstract and do not reveal *how* closures are implemented. They reveal only how much space a closure occupies and which pointers it keeps live. A user can apply these rules without knowing how closures are internally represented.

Our construction of the reasoning rules for closures is in two layers. First, we introduce a low-level assertion $\text{Closure } E f \vec{x} t \ell$, which asserts that, at location ℓ in the heap, one finds a closure that behaves like the function $\mu f. \lambda \vec{x}. t$ under the environment E . Crucially, in this assertion, the term $\mu f. \lambda \vec{x}. t$ can have free variables, whose values are given by E . This assertion does not reveal how a closure is represented in memory, but does reveal its code. We give an overview of this low-level API (§9.4), then describe some details of its implementation (§9.5). Second, we define a high-level assertion $\text{Spec } n E P \ell$, which describes the behavior of a closure in a more abstract way. It asserts that, at location ℓ , one finds a closure that corresponds to a n -ary function, whose behavior is described by the predicate P , and whose environment is E . The exact type and meaning of P are explained later on; roughly speaking, it is a Hoare triple. Although the environment E does not participate in the description of the behavior of the closure, it remains needed in order to reason about the pointers that it contains and about the size of the closure block. We give an overview of this high-level API (§9.6), then describe its implementation (§9.7). Only the high-level layer is exposed to the end user; the low-level layer remains internal.

9.1 Environments

We write $\text{fvclo}(f, \vec{x}, t)$ for a list of the free variables of the function $\mu f. \lambda \vec{x}. t$, that is, for a list of the variables in the set $\text{fv}(t) \setminus \{f, \vec{x}\}$. The order in which the variables occur in this list does not matter, but is fixed: this is reflected in the fact that fvclo is a function of f, \vec{x} , and t .

An environment E is a list of pairs (v, q) of a value v and a nonzero fraction q . This fraction is used in a pointed-by-heap assertion, as follows: we write $E \leftarrow L$ for the conjunction $\bigstar_{(v, q) \in E} v \leftarrow_q L$. The assertion $E \leftarrow L$ can be understood as a collective fractional pointed-by-heap assertion that covers every memory location that occurs in the environment E .

The length and order of the list E are intended to match the length and order of the list $\text{fvclo}(f, \vec{x}, t)$. An environment E is not a runtime object: it is a mathematical object that we use as a parameter of the predicates Closure and Spec .

9.2 Closure Implementation

The definitions of the closure macros $\mu_{\text{clo}}f. \lambda \vec{x}. t$ and of $(\ell \vec{v})_{\text{clo}}$ appear in Figure 31. Both macros generate LambdaFit syntax: that is, the result of their expansion is a LambdaFit expression.

The code produced by the macro $\mu_{\text{clo}}f. \lambda \vec{x}. t$ allocates a block of size $n + 1$, stores a code pointer in the first field, stores the values currently bound to the variables y_0, \dots, y_{n-1} in the remaining

<p><i>Closure construction:</i></p> $\mu_{\text{clo}}f. \lambda \vec{x}. t \triangleq$ $\text{let } f = \text{alloc } (n + 1) \text{ in}$ $f[0] \leftarrow \text{codeclo}(f, \vec{x}, t);$ $f[i + 1] \leftarrow y_i; \quad \# \text{ for each } i \text{ in } [0, n)$ f <p><i>Closure invocation:</i></p> $(v \vec{w})_{\text{clo}} \triangleq$ $(v[0] (v :: \vec{w}))_{\text{ptr}}$	<p><i>Closure code pointer:</i></p> $\text{codeclo}(f, \vec{x}, t) \triangleq$ $\mu_{\text{ptr}}. \lambda (f :: \vec{x}).$ $\text{let } y_i = f[i + 1] \text{ in } \quad \# \text{ for each } i \text{ in } [0, n)$ t <p><i>Side condition:</i></p> $\text{fvclo}(f, \vec{x}, t) = [y_0; \dots; y_{n-1}]$
--	---

Fig. 31. Closures: macros for closure construction and invocation

fields, and returns the address of this block. The variables y_0, \dots, y_{n-1} are the free variables of the function $\mu f. \lambda \vec{x}. t$, that is, $\text{fvclo}(f, \vec{x}, t)$.

The code pointer is produced by the auxiliary macro $\text{codeclo}(f, \vec{x}, t)$. It is a closed function whose parameters are f (the closure itself) followed with \vec{x} . This function loads the values stored in the closure and binds them to the variables y_0, \dots, y_{n-1} before executing the body t .

The code produced by the closure invocation macro $(v \vec{w})_{\text{clo}}$ first fetches the code pointer that is stored in the first field of the closure, then invokes this code pointer, passing it the closure v itself as well as the actual arguments \vec{w} .

9.3 Pending Substitutions

Our specifications of closure construction (§9.4, §9.6) involve *pending substitutions*. A pending substitution, written $[\vec{v}/\vec{y}](\mu_{\text{clo}}f. \lambda \vec{x}. t)$, is the application of the substitution $[\vec{v}/\vec{y}]$ to the closure construction macro $\mu_{\text{clo}}f. \lambda \vec{x}. t$. In this substitution, we may assume that the variables \vec{y} are the free variables of the function $\mu f. \lambda \vec{x}. t$. Any other variables can be removed from the domain of the substitution, as they do not impact the closure. The reason why we must be prepared to reason about a pending substitution is that the premise of **LETVAL** gives rise to substitutions which (after being propagated down) become blocked in front of the *opaque* macro $\mu_{\text{clo}}f. \lambda \vec{x}. t$. The values \vec{v} that appear in this substitution are the values “captured” by the closure, that is, the values that are stored in the closure when it is constructed.

To illustrate these rules, let us take the example of a generator of integers:

$$\text{let } r = \text{alloc } 1 \text{ in } r[0] \leftarrow 0; (\mu_{\text{clo}}. \lambda_. \text{let } x = r[0] \text{ in } r[0] \leftarrow (x + 1); x)$$

This code allocates a reference r , initializes it to 0, then allocates a closure that loads the content of r , names it x , increments r , and returns x . This closure captures the free variable r .

Now, let us briefly describe which reasoning rules must be applied in order to verify this code, and how the code that appears in the goal evolves as the proof progresses. After applying **BIND** and **ALLOC** to reason about the memory allocation instruction, the term looks as follows, where we have named ℓ the memory location produced by the instruction $\text{alloc } 1$:

$$\text{let } r = \ell \text{ in } r[0] \leftarrow 0; (\mu_{\text{clo}}. \lambda_. \text{let } x = r[0] \text{ in } r[0] \leftarrow (x + 1); x).$$

Applying **LETVAL** gives rise to the substitution $[\ell/r]$, which is applied to the right-hand side of the let binding. Thus, after applying **BIND** and **STORE** to reason about the store instruction, the term that appears in the goal is:

$$[\ell/r](\mu_{\text{clo}}. \lambda_. \text{let } x = r[0] \text{ in } r[0] \leftarrow (x + 1); x)$$

$$\begin{array}{c}
\text{MkCLO} \\
\frac{\vec{y} = \text{fvclo}(f, \vec{x}, t) \quad E = \text{zip } \vec{v} \vec{q} \quad |\vec{v}| = |\vec{y}| \quad f \notin \vec{x}}{\left\{ \begin{array}{l} \diamond(\text{size}(1 + |E|)) * \text{outside } \pi \\ E \leftarrow \emptyset \end{array} \right\} \pi: [\vec{v}/\vec{y}] (\mu_{\text{clo}}f. \lambda \vec{x}. t) \left\{ \begin{array}{l} \lambda \ell. \text{outside } \pi * \text{Closure } E f \vec{x} t \ell \\ \ell \Leftarrow \{\pi\} * \ell \Leftarrow \emptyset \end{array} \right\}} \\
\\
\text{CALLCLO} \\
\frac{\vec{y} = \text{fvclo}(f, \vec{x}, t) \quad E = \text{zip } \vec{v} \vec{q} \quad |\vec{x}| = |\vec{w}| \quad \text{locs}(\vec{v}) = \text{dom}(M) \quad \{\text{outside } \pi * M \Leftarrow \{\pi\} * \Phi\} \pi: [\vec{v}/\vec{y}][\ell/f][\vec{w}/\vec{x}]t \{\Psi\}}{\{\text{Closure } E f \vec{x} t \ell * \text{outside } \pi * M \Leftarrow \{\pi\} * \Phi\} \pi: (\ell \vec{w})_{\text{clo}} \{\Psi\}} \\
\\
\text{Closure } E f \vec{x} t \ell * \ell \Leftarrow \emptyset * \ell \Leftarrow \emptyset \Rightarrow \diamond(\text{size}(1 + |E|)) * \dagger \ell * E \Leftarrow \emptyset \quad \text{CLOFREE} \\
\text{Closure } E f \vec{x} t \ell \text{ is persistent} \quad \text{CLOPERSIST}
\end{array}$$

Fig. 32. Closures: low-level API

As explained above, the macro $\mu_{\text{clo}}f. \lambda \vec{x}. t$ is opaque, so a substitution cannot be pushed into it. This explains why our reasoning rules for closure construction must allow reasoning about a term of the form $[\vec{v}/\vec{y}](\mu_{\text{clo}}f. \lambda \vec{x}. t)$.

9.4 Low-Level Closure API

Our low-level reasoning rules for closures, shown in Figure 32, involve the predicate *Closure*, which describes the layout of a closure in memory. Its definition is given in the next section (§9.5).

The rule **MkCLO** specifies a closure construction operation with a pending substitution. In the second premise of **MkCLO**, an environment E is built by pairing up the values \vec{v} with nonzero fractions \vec{q} . Then, according to the precondition in **MkCLO**, closure construction consumes $E \Leftarrow \emptyset$. In other words, for each memory location that occurs in E , it consumes a fractional pointed-by-heap assertion. This records the fact that there exists a pointer from the closure to each such memory location.

According to the precondition in **MkCLO**, closure construction consumes $\text{size}(1 + |E|)$ space credits, reflecting the space needed to store a code pointer and the values \vec{v} in a flat closure.

Because closure construction involves an allocation, **MkCLO** requires the thread π to be outside a protected section.

According to the postcondition in **MkCLO**, closure construction produces a memory location ℓ . Pointed-by-heap and pointed-by-thread assertions for this memory location are produced, indicating that this memory location is fresh. Furthermore, the assertion $\text{Closure } E f \vec{x} t \ell$, which guarantees that there is a well-formed closure at address ℓ , is also produced. In this paper, in contrast with our earlier work [Moine et al. 2023], this assertion is *persistent* [Jung et al. 2018b, §2.3]. This means that the knowledge that there is a closure at address ℓ can be shared without any restriction. The pointed-by-heap and pointed-by-thread assertions $\ell \Leftarrow \{\pi\} * \ell \Leftarrow \emptyset$ are *not* persistent. Indeed, these assertions allow deallocating the closure, and our program logic ensures that every object is deallocated at most once.

The rule **CALLCLO** closely resembles the rule **CALLPTR** for primitive function calls (Figure 21). One difference is that **CALLCLO** requires the assertion $\text{Closure } E f \vec{x} t \ell$, which describes the closure. Another difference is that, whereas a primitive function $\mu_{\text{ptr}}f. \lambda \vec{x}. t$ must be closed, a general function can have a nonempty list of free variables \vec{y} , an alias for $\text{fvclo}(f, \vec{x}, t)$. In the last premise of **CALLCLO**, which requires reasoning about the function's body, the variables \vec{y} are replaced with the values \vec{v} captured at closure construction time, which are recorded in the environment E .

$$\begin{aligned}
\text{Closure } E f \vec{x} t \ell &\triangleq \lceil f \notin \vec{x} \wedge |E| = |\text{fv}_{\text{clo}}(f, \vec{x}, t)| \rceil * \\
&\ell \mapsto_{\square} (\text{code}_{\text{clo}}(f, \vec{x}, t) :: \text{map } \text{fst } E) * \\
&\boxed{\dagger \ell \vee E \leftarrow \{+\ell\}}
\end{aligned}$$

Fig. 33. Definition of the predicate *Closure*

The precondition of **CALLCLO** requires a pointed-by-thread assertion $M \Leftarrow \{\pi\}$, where the domain of the map M includes all of the locations that appear in \vec{v} , that is, all of the locations that appear in the closure’s environment. This assertion is not consumed: it appears again in the precondition of the triple that forms the last premise of **CALLCLO**. In other words, it is transmitted from the caller to the callee. The presence of this assertion is imposed to us by the fact that, when the closure is invoked, these values are read from memory: the load instructions that appear in the definition of $\text{code}_{\text{clo}}(f, \vec{x}, t)$ in Figure 31 require pointed-by-thread assertions for the values that are read. If desired, the pointed-by-thread assertion $M \Leftarrow \{\pi\}$ can be transmitted back from the callee to the caller via a suitable instantiation of the postcondition Ψ . Alternatively, it may be consumed by the callee to justify a logical deallocation operation.

Together, the rules **MkCLO** and **CALLCLO** express the correctness of our closure construction and invocation macros. They guarantee that a closure at address ℓ constructed by $[\vec{v}/\vec{y}] \mu_{\text{clo}} f. \lambda \vec{x}. t$, when invoked with actual arguments \vec{w} , behaves like the term $[\vec{v}/\vec{y}][\ell/f][\vec{w}/\vec{x}]t$. This is the operational behavior that is expected of a closure.

CLOFREE logically deallocates a closure. It resembles **FREEONE**, but, instead of a “*sizeof*” assertion, requires the abstract assertion $\text{Closure } E f \vec{x} t \ell$. Like **FREEONE**, it produces space credits and a deallocation witness for the closure. Furthermore, **CLOFREE** lets the user recover the pointed-by-heap assertion $E \leftarrow \emptyset$, thereby undoing the effect of **MkCLO**.

9.5 Low-Level Closure API: Implementation Details

Figure 33 presents the internal definition of the assertion $\text{Closure } E f \vec{x} t \ell$. It records two pure facts: the name f is disjoint from the parameters \vec{x} and the length of the environment E matches the number of free variables of the closure.

Then, a points-to assertion states that the location ℓ points to a block of size $1 + |E|$, whose first field contains the code of the closure, $\text{code}_{\text{clo}}(f, \vec{x}, t)$, and whose remaining fields contain the values recorded in the environment E . Because this points-to assertion carries a *discarded fraction* \square [Vindum and Birkedal 2021], it is a *persistent points-to* assertion. This reflects the fact that the closure is immutable.

The last component in this definition is a liveness-based cancellable invariant (§5.10): a persistent assertion that we can tear down and regain full ownership when we deallocate ℓ .

Because every assertion involved in its definition is persistent, the assertion $\text{Closure } E f \vec{x} t \ell$ is itself persistent.

The liveness-based cancellable invariant contains the pointed-by-heap assertion $E \leftarrow \{+\ell\}$, which means that every memory location in E is pointed to by the closure. In the proof of the reasoning rule **CLOFREE**, we tear down the liveness-based cancellable invariant, and gain back the assertion $E \leftarrow \{+\ell\}$. Because ℓ is now dead, we use the **CLEANPBHEAP** rule to change $E \leftarrow \{+\ell\}$ into $E \leftarrow \emptyset$. This explains how, in the proof of **CLOFREE**, we are able to produce the assertion $E \leftarrow \emptyset$.

$$\begin{array}{c}
\text{MkSPEC} \\
\frac{\vec{y} = \text{fvclo}(f, \vec{x}, t) \quad E = \text{zip } \vec{v} \vec{q} \quad |\vec{v}| = |\vec{y}| \quad f \notin \vec{x} \quad n = |\vec{x}|}{\forall \vec{w}. \square (\text{Spec } n E P \ell \multimap P \ell \vec{w} ([\vec{v}/\vec{y}][\ell/f][\vec{w}/\vec{x}]t))} \\
\frac{\left\{ \begin{array}{l} \diamond(\text{size}(1 + |E|)) * \text{outside } \pi \\ E \leftarrow \emptyset \end{array} \right\} \pi: [\vec{v}/\vec{y}] (\mu_{\text{clo}f}. \lambda \vec{x}. t) \left\{ \begin{array}{l} \text{outside } \pi * \text{Spec } n E P \ell \\ \ell \Leftarrow \{\pi\} * \ell \Leftarrow \emptyset \end{array} \right\}}{\text{CALLSPEC}} \\
\frac{E = \text{zip } \vec{v} \vec{q} \quad \text{dom}(M) = \text{locs}(\vec{v}) \quad |\vec{w}| = n}{(\forall u. P \ell \vec{w} u \multimap \{ \text{outside } \pi * M \Leftarrow \{\pi\} * \Phi \} \pi: u \{ \Psi \})} \\
\{ \text{Spec } n E P \ell * \text{outside } \pi * M \Leftarrow \{\pi\} * \Phi \} \pi: (\ell \vec{w})_{\text{clo}} \{ \Psi \} \\
\\
\square (\forall \vec{w} t. P_1 \ell \vec{w} t \multimap P_2 \ell \vec{w} t) * \text{Spec } n E P_1 \ell \multimap \text{Spec } n E P_2 \ell \quad \text{SPECWEAK} \\
\text{Spec } n E P \ell * \ell \Leftarrow \emptyset * \ell \Leftarrow \emptyset \Rightarrow \diamond(\text{size}(1 + |E|)) * \dagger \ell * E \Leftarrow \emptyset \quad \text{SPECFREE} \\
\text{Spec } n E P \ell \text{ is persistent} \quad \text{SPECPERSIST}
\end{array}$$

Fig. 34. Closures: high-level API

9.6 High-Level Closure API

The user of a program logic is ultimately interested in the specification of a function, not in the details of its implementation. Yet, the predicate $\text{Closure } E f \vec{x} t \ell$ reveals the code of the closure. As a result, a user naturally wishes to hide this information via an existential quantification over this code. This pattern is common enough and technical enough that we offer a higher-level API where this existential quantification is built in. To this end, we introduce the assertion $\text{Spec } n E P \ell$ (defined further on in §9.7), where n is the arity of the function, E is the environment of the closure, P describes the behavior of the closure, and ℓ is the location of the closure in memory.

Like the Closure predicate (§9.4, §9.5), and unlike the Spec predicate presented in our previous paper [Moine et al. 2023], the predicate Spec is persistent. This enables a better separation of concerns between the persistent assertion $\text{Spec } n E P \ell$, which views the closure as an eternal service provider, and the affine assertion $\ell \Leftarrow \{\pi\} * \ell \Leftarrow \emptyset$, which views it as an object in memory, allowing it to participate in the object graph and (eventually) to be logically deallocated.

Figure 34 presents the reasoning rules associated with the Spec predicate. Let us first examine the rule **CALLSPEC**. In many ways, this rule is the same as the low-level rule **CALLCLO**. The main difference is that, to prove that the call $(\ell \vec{w})_{\text{clo}}$ admits the postcondition Ψ , the user must check that the entailment $\forall u. P \ell \vec{w} u \multimap \{ \text{outside } \pi * M \Leftarrow \{\pi\} * \Phi \} \pi: u \{ \Psi \}$ holds. Intuitively, u denotes the instantiated function body that was visible in **CALLCLO**; however, this function body is now abstracted away by the universal quantification over u . The predicate P represents the specification of the function, and is typically instantiated with a triple. For example, in the specification of a closure of arity 1 whose effect is to increment a reference r that it receives as an argument, the predicate P takes the form: $\lambda \ell \vec{w} u. \forall r n. \ulcorner \vec{w} = [r] \urcorner \multimap \{ r \mapsto [n] \} \pi: u \{ \lambda(). r \mapsto [n+1] \}$. In short, the user must prove an entailment stating that the specification needed by the caller follows from the specification P .

Let us now consider the rule **MkSPEC**. It is again quite similar to the low-level rule **MkCLO**. The premise on the second line ensures that P is a valid description of the behavior of the function body, whose concrete form $[\vec{v}/\vec{y}][\vec{w}/\vec{x}]t$ is visible. In comparison with the low-level API (§9.4), the work of reasoning about the function body is shifted from the closure invocation site to the closure construction site. Moreover, while establishing $P \ell \vec{w} ([\vec{v}/\vec{y}][\ell/f][\vec{w}/\vec{x}]t)$, the user is allowed to assume $\text{Spec } n E P \ell$: this allows verifying recursive calls.

$$\begin{aligned}
\text{Spec } nEP\ell &\triangleq \\
&\exists f \vec{x} t P'. \\
&\quad \ulcorner |\vec{x}| = n \urcorner * \text{Closure } E f \vec{x} t * \\
&\quad \text{let } \vec{v} = \text{mapfst } E \text{ in} \\
&\quad \text{let } \vec{y} = \text{fvclo}(f, \vec{x}, t) \text{ in} \\
&\quad \text{let } \text{body } \vec{w} = [\vec{v}/\vec{y}][\ell/f][\vec{w}/\vec{x}]t \text{ in} \\
&\quad \triangleright \square(\forall \vec{w}. \text{Spec } nEP'\ell \multimap P'\ell \vec{w} (\text{body } \vec{w})) * \\
&\quad \triangleright \square(\forall \vec{w} u. P'\ell \vec{w} u \multimap P\ell \vec{w} u)
\end{aligned}$$
Fig. 35. Definition of the predicate *Spec*

The rule **SPECWEAK** is a consequence rule: it allows weakening the assertion $\text{Spec } nEP_1\ell$ into $\text{Spec } nEP_2\ell$, under the hypothesis that P_1 is stronger than P_2 . This hypothesis is expressed as $\square(\forall \vec{w} t. P_1\ell \vec{w} t \multimap P_2\ell \vec{w} t)$, where \square is the *persistence modality* [Jung et al. 2018b, §5.3]. This modality requires the proof of the implication $(\forall \vec{w} t. P_1\ell \vec{w} t \multimap P_2\ell \vec{w} t)$ to depend only on persistent resources.

The rule **SPECFREE** is similar to the rule **CLOFREE**.

9.7 High-Level Closure API: Implementation Details

Figure 35 presents the definition of the assertion $\text{Spec } nEP\ell$. This is a guarded recursive definition: *Spec* appears (under a “later” modality) in its own definition. The definition is existentially quantified over the code of the closure, represented by f , \vec{x} , and t . It is also existentially quantified over a predicate P' that is required to be stronger than P . This lets us establish **SPECWEAK**.

10 TRIPLES WITH SOUVENIR

In this section, we propose *triples with souvenir* [Moine et al. 2023], a syntactic sugar that allows for simpler reasoning rules—in particular, a simpler **BIND** rule—while reasoning about code that lies outside a protected section. We first present the reasoning rules of triples with souvenir (§10.1), then cover how they are defined (§10.2).

10.1 Those Who Cannot Remember the Past Are Condemned to Repeat It

IrisFit, as presented until this point, can be cumbersome to use, for two unrelated reasons.

One reason is that the user must give up pointed-by-thread assertions at each application of **BIND**, even in the common case where such a fraction has been framed already at a previous application of **BIND**, which encloses the current application. This obligation to split off and give up pointed-by-thread assertions becomes especially heavy when a variable x denotes a location and has a long *live range*, that is, when this location remains a root throughout a long sequence of instructions. In such a situation, at each point in the sequence, the user is required to split off and give up a fractional pointed-by-thread assertion for x .¹⁴

A second reason is that, typically, the large majority of instructions are placed outside protected sections. Yet, the user must provide the assertion *outside* π at each application of the *outside rules* **ALLOC**, **CALLPTR**, **FORK**, **POLL**, **MRSPEC**, and **CALLSPEC**. This is not difficult, but the presence of this assertion creates visual clutter in pre- and postconditions.

To alleviate both problems at once, we follow Moine et al. [2023] and introduce *triples with souvenir*. A triple with souvenir takes the form $[R] \{\Phi\} \pi : t \{\Psi\}$, where R is a set of locations

¹⁴The problem is partly mitigated by the “no trim” mode \blackstar (§6.5). However, this mode is designed for very local use, and cannot be exploited if trimming is needed.

$$\begin{array}{c}
\text{BINDWITHSOUVENIR} \\
\frac{\text{dom}(M) = \text{locs}(K) \setminus R \quad [R \cup \text{locs}(K)] \{\Phi\} \pi: t \{\Psi'\} \quad \forall v. [R] \{M \Leftarrow \{\pi\} * \Psi' v\} \pi: K[v] \{\Psi\}}{[R] \{M \Leftarrow \{\pi\} * \Phi\} \pi: K[t] \{\Psi\}} \\
\\
\text{ADDSOUVENIR} \qquad \qquad \qquad \text{FORGETSOUVENIR} \\
\frac{[\{\ell\} \cup R] \{\Phi\} \pi: t \{\Psi\}}{[R] \{\ell \Leftarrow_p \{\pi\} * \Phi\} \pi: t \{\lambda v. \ell \Leftarrow_p \{\pi\} * \Psi v\}} \qquad \frac{R' \subseteq R \quad [R'] \{\Phi\} \pi: t \{\Psi\}}{[R] \{\Phi\} \pi: t \{\Psi\}} \\
\\
\text{EMPTYSOUVENIR} \\
[\emptyset] \{\Phi\} \pi: t \{\Psi\} \equiv \{\Phi * \text{outside } \pi\} \pi: t \{\lambda v. \Psi v * \text{outside } \pi\}
\end{array}$$

Fig. 36. Key reasoning rules for triples with souvenir

$$\begin{array}{l}
[R] \{\Phi\} \pi: t \{\Psi\} \triangleq \\
\forall M. R = \text{dom}(M) \implies \\
\{\Phi * \text{outside } \pi * M \Leftarrow \{\pi\}\} \pi: t \{\lambda v. \Psi v * \text{outside } \pi * M \Leftarrow \{\pi\}\}
\end{array}$$

Fig. 37. Definition of triples with souvenir

for which the user has already given up a pointed-by-thread assertion. Recording this *souvenir* (or remembrance) relieves the user from the obligation of giving up another pointed-by-thread assertion at future applications of the `BIND` rule. Furthermore, a triple with souvenir implicitly carries an *outside* π assertion: this allows for more concise statements of the “outside rules”.

For each reasoning rule in Figure 21, we provide a new rule (not shown) that operates on triples with souvenir and that is polymorphic in R . This is done simply by inserting $[R]$ in front every triple that appears in the rule. We do not provide new reasoning rules for protected sections, as triples with souvenir are applicable only outside protected sections.

The new reasoning rules that make use of souvenirs appear in Figure 36. `BINDWITHSOUVENIR` is what we aimed for: it is our motivation for introducing triples with souvenir. It closely resembles `BIND`, but does not require the user to give up pointed-by-thread assertions for the locations that are already part of the souvenir R . The first premise requires the domain of M (a map of locations to nonzero fractions) to cover all roots of the evaluation context K , except those that are already in the souvenir R . In other words, *if a location already appears in R then there is no need to again split off and give up a pointed-by-thread assertion for this location*. Furthermore, `BINDWITHSOUVENIR` augments the current souvenir by changing R to $R \cup \text{locs}(K)$ in its second premise. Thus, nested applications of this rule do not require repeatedly and redundantly giving up pointed-by-thread assertions. The rule `ADDSOUVENIR` extends the current souvenir with a location ℓ . This requires framing out (temporarily giving up) a pointed-by-thread assertion for ℓ . The rule `FORGETSOUVENIR` shrinks the current souvenir. The rule `EMPTYSOUVENIR` shows that a triple with an empty souvenir is equivalent to a triple without souvenir and with an “*outside*” assertion in its pre- and postcondition.

By exploiting triples with souvenir, each of the “outside rules” (`ALLOC`, `CALLPTR`, `FORK`, `POLL`, `MkSPEC` and `CALLSPEC`), can be given a more concise statement. For example, the reasoning rule `POLL` can be more concisely formulated as `POLLWITHSOUVENIR`:

$$\begin{array}{c}
\text{POLL} \\
\{\text{outside } \pi\} \pi: \text{poll } \{\lambda(). \text{outside } \pi\} \\
\\
\text{POLLWITHSOUVENIR} \\
[R] \{\ulcorner \text{True} \urcorner\} \pi: \text{poll } \{\lambda(). \ulcorner \text{True} \urcorner\}
\end{array}$$

10.2 Internals of Souvenirs

The definition of triples with souvenir appears in Figure 37. A triple with souvenir $[R] \{\Phi\} \pi: t \{\Psi\}$ is expressed as an ordinary triple where the assertions *outside* π and $M \Leftarrow \{\pi\}$ are framed out. That is, these assertions appear in the pre- and postcondition, so they are required and preserved, but they are not made available to a user who views a triple with souvenir as an abstract assertion. The domain of the map M is the set R : this ensures that, for every location in this set, a fractional pointed-by-thread assertion is indeed framed out.

A triple with souvenir describes a piece of code whose execution begins and ends outside a protected section: it cannot be used to describe a code fragment that lies inside a protected section. To establish a triple with souvenir about a whole protected section, the user must unfold the definition of triples with souvenir and drop down to the level of standard triples. Then, all of the reasoning rules for standard triples are applicable.

In our mechanization [Moine 2024b], we use a more general triple that allows both “no trim” mode (§6.5) without a souvenir and normal mode with a souvenir. This general triple always frames out an “*outside*” assertion. In our case studies, this is the triple that we use most of the time.

11 CASE STUDIES

We now showcase the expressiveness of IrisFit via a series of representative case studies. We first present *logically atomic triples* [da Rocha Pinto et al. 2014; Jung et al. 2015], a standard way of specifying operations on concurrent data structures. We begin our case studies with an encoding of the fetch-and-add operation in LambdaFit, which makes use of protected sections (§11.2). Then, we present an implementation of a concurrent counter object, implemented as a pair of closures that share an internal reference (§11.3). We continue with a library for async/finish parallelism, which exploits our implementation of fetch-and-add (§11.4). We conclude this section by presenting our version of Treiber’s stack (§11.5), which exploits protected sections, along the lines sketched earlier (§3). For each case study, we present the code, the specification, and some insights into the proof. For establishing concrete heap bounds, we pose in this section that a block of n fields is represented by n memory words, that is, we pose $size(n) = n$. Another practical choice such as $size(n) = n + 1$ would only affect the constant values that appear behind diamond symbols in specifications.

Our mechanization [Moine 2024b] contains additional case studies that we do not cover here. They include sequential examples (a sequential algorithm written in continuation-passing style; a sequential singly-linked circular list; three distinct implementations of sequential stacks) and concurrent examples (a spin lock; Michael and Scott’s lock-free queue, with protected sections).

11.1 Atomic triples

Our specifications for fetch-and-add (§11.2) and for Treiber’s stack (§11.5) involve *logically atomic triples*, also known simply as *atomic triples* [da Rocha Pinto et al. 2014; Jung et al. 2015]. In our work, an atomic triple takes the form:

$$[R] \left\langle \frac{\Phi_{private}}{\forall \vec{x}. \Phi_{public}} \right\rangle \pi: t \left\langle \frac{\lambda v. \Phi'_{private}}{\Phi'_{public}} \right\rangle$$

The parameter R between square brackets is a souvenir (§10). We construct our atomic triples on top of our triples with souvenir (§10) in the same way that atomic triples are usually constructed on top of ordinary triples. Intuitively, an atomic triple that carries a souvenir $[R]$ is an atomic triple whose private pre- and postconditions are extended with a pointed-by-thread assertion

```

faa  $\triangleq$   $\mu_{\text{ptr}} f. \lambda[l, i, n].$ 
  let  $m = l[i]$  in
  enter; if CAS  $l[i] m (m + n)$ 
  then (exit;  $m$ )
  else (exit; ( $f [l, i, n]$ )ptr)

```

$$\text{FAA} \left\langle \frac{\ell \Leftarrow_p \{\pi\}}{\forall \vec{v} m. \ulcorner \vec{v}(i) = m^\urcorner * \ell \mapsto \vec{v}} \right\rangle \pi : (\text{faa } [\ell, i, n])_{\text{ptr}} \left\langle \frac{\lambda m'. \ulcorner m' = m^\urcorner}{\ell \mapsto ([i := (m + n)]\vec{v}) * \ell \Leftarrow_p \emptyset} \right\rangle$$

Fig. 38. Code and specification of fetch-and-add

that covers R (that is, a pointed-by-thread assertion $M \Leftarrow \{\pi\}$ where $R = \text{dom}(M)$) and with the assertion *outside* π .

The private precondition Φ_{private} and the private postcondition $\lambda v. \Phi'_{\text{private}}$ play the same role as the precondition and postcondition of a standard triple. The private precondition is given up by thread π when the execution of the term t begins; the private postcondition is gained by thread π when the execution of the term t ends. They are *private* in the sense that they are invisible to other threads.

The characteristic feature of atomic triples is the presence of a public precondition Φ_{public} and of a public postcondition Φ'_{public} . An atomic triple guarantees that the public precondition Φ_{public} continuously holds until a certain point in time, the *linearization point* [Herlihy and Wing 1990], where it is atomically transformed into the public postcondition Φ'_{public} [Birkedal et al. 2021]. Technically, an atomic triple involves a quantification over a list of variables \vec{x} , whose scope is Φ_{public} , Φ'_{private} , and Φ'_{public} . The existentially quantified public precondition $\exists \vec{x}. \Phi_{\text{public}}$ continuously holds until the linearization point is reached. There, a specific instantiation of the variables \vec{x} becomes fixed. For this specific choice of \vec{x} , the public precondition is transformed into the public postcondition Φ'_{public} , and the value v that is eventually returned satisfies Φ'_{private} .

11.2 Fetch-and-Add

The “fetch-and-add” (FAA) operation atomically increments the content of an integer reference, and returns the previous content of the reference. Although this operation is commonly provided in hardware, implementing it in LambdaFit is a fairly instructive exercise. Indeed, this code and its proof offer a typical example of the use of protected sections.

Code. In our setting, FAA takes three parameters: an address l , an offset i , and the desired increment n , an integer value. We encode FAA as a tail-recursive function whose body contains a CAS instruction enclosed in a protected section. The code is shown in Figure 38. The recursive function is named f ; its parameters are l , i and n . Initially, the content of the memory at address l and offset i is loaded into the variable m . Then, a protected section is entered, and a CAS instruction attempts to update the content of the memory from m to $m + n$. In case of success, the protected section is exited and the value m is returned. In case of failure, the protected section is also exited, and a recursive call is performed, so as to try again.

Thanks to the protected section, as soon as the CAS instruction succeeds, the memory location l can be considered as a temporary root, as opposed to an ordinary root. Indeed, as soon as CAS succeeds, it is known that the first branch of the conditional construct will be taken, so the protected section will be exited via the first exit instruction, where l is no longer a root.

Without a protected section, at the program point that follows CAS and precedes the separation of the two branches, l would still be considered a root (that is to say, an ordinary root), because it occurs inside the “else” branch, and according to the FVR (§2.2), every location that occurs in the code that lies ahead is a root.

Specification. Our specification of FAA appears in Figure 38. The private precondition consumes a pointed-by-thread assertion for the location ℓ , carrying some fraction p and the current thread identifier π . The public precondition requires that ℓ point to a block \vec{v} and that the value stored at offset i in this block be m . The public postcondition asserts that FAA atomically updates m into $m + n$. Crucially, it also produces an updated pointed-by-thread assertion for ℓ , carrying the same fraction p and an *empty* set of thread identifiers. This means that as soon as the linearization point is reached, one can consider that ℓ is not a root in thread π . This turns out to be crucial while reasoning about our async/finish library (§11.4). The private postcondition asserts that the result of FAA is m .

Proof insights. Here is how we use the reasoning rules of protected sections (Figure 22) while verifying that FAA obeys its specification. Upon entering the protected section, we use **ENTER** and transform the assertion *outside* π into the assertion *inside* $\pi \emptyset$. Then, we face the CAS instruction, a possible linearization point. We open the public precondition, and gain the points-to assertion for ℓ . By case analysis on the value that is currently stored at address l and offset i , we consider the case where CAS succeeds and the case where it fails. Let us focus on the case where it succeeds. We use **CASSUCCESS**, which updates the points-to assertion, and effectively execute the linearization point. At this point, the atomic triple requires us to prove that the public postcondition holds. Using **ADDTEMPORARY**, we make ℓ a temporary root: this changes the assertions $\ell \Leftarrow_p \{\pi\}$ and *inside* $\pi \emptyset$ into $\ell \Leftarrow_p \emptyset$ and *inside* $\pi \{\ell\}$. By giving up the points-to and pointed-by-thread assertions, we fulfill the public postcondition. Then, we use **IFTRUE** and enter the first branch of the “if” statement. There, **TRIMINSIDE** lets us change the assertion *inside* $\pi \{\ell\}$ to *inside* $\pi \emptyset$. This allows us to exit the protected section using **EXIT**. We finish the proof with **VAL**.

11.3 A Concurrent Counter Object

Our next example is a concurrent monotonic “counter” object, whose internal state is stored in a mutable reference, and whose access is mediated by a pair of closures: a closure i *increments* the counter; a closure g *gets* its current value. This is an example of a procedural abstraction [Reynolds 1975], also known as an *object*: indeed, “an object is a value exporting a procedural interface to data or behavior” [Cook 2009]. Crucially, a counter can be used concurrently by several threads.

Code. The top of Figure 39 presents the code that we verify. The function call $(\text{ref } [x])_{\text{ptr}}$ allocates a mutable reference, that is, a block of size 1. The function call $(\text{pair } [x, y])_{\text{ptr}}$ allocates a mutable pair, that is, a block of size 2. The function call $(\text{ignore } [x])_{\text{ptr}}$ ignores its argument and returns the unit value. The function call $(\text{create } [])_{\text{ptr}}$ returns a fresh “counter”, that is, a pair of two closures i and g . Both closures point to an internal reference r , which is initialized to the value 0. The closure i uses our fetch-and-add function (§11.2) and ignores its result.

Specifications. Figure 39 presents the specification of our concurrent counter. It is inspired by a specification that appears in lecture notes [Birkedal and Bizjak 2023]. It relies on an abstract assertion *counter* $i g p n$ where i is the location of the “increment” closure, g is the location of the “get” closure, $p \in (0; 1]$ is a fraction that represents a *share* of the ownership of the counter, and n , a natural number, represents a *past contribution* to the current value of the counter. If p is 1 then the contribution n is in fact the current value of the counter.

$$\begin{array}{l}
\text{ref} \triangleq \mu_{\text{ptr}\dots} \lambda[x]. \\
\quad \text{let } r = \text{alloc } 1 \text{ in} \\
\quad r[0] \leftarrow x; r \\
\text{pair} \triangleq \mu_{\text{ptr}\dots} \lambda[x, y]. \\
\quad \text{let } r = \text{alloc } 2 \text{ in} \\
\quad r[0] \leftarrow x; r[1] \leftarrow y; r
\end{array}
\qquad
\begin{array}{l}
\text{ignore} \triangleq \mu_{\text{ptr}\dots} \lambda[x]. () \\
\text{create} \triangleq \mu_{\text{ptr}\dots} \lambda[[]]. \\
\quad \text{let } r = (\text{ref } [0])_{\text{ptr}} \text{ in} \\
\quad \text{let } i = \mu_{\text{clo}\dots} \lambda_. (\text{ignore } [(\text{faa } [r, 0, 1])_{\text{ptr}}])_{\text{ptr}} \text{ in} \\
\quad \text{let } g = \mu_{\text{clo}\dots} \lambda_. r[0] \text{ in} \\
\quad (\text{pair } [i \ g])_{\text{ptr}}
\end{array}$$

$$\begin{array}{l}
(\text{counter } i \ g \ (p_1 + p_2) \ (n_1 + n_2)) \quad \equiv \quad (\text{counter } i \ g \ p_1 \ n_1 \ * \ \text{counter } i \ g \ p_2 \ n_2) \\
[\emptyset] \{\diamond 7\} \ \pi: (\text{create } [])_{\text{ptr}} \left\{ \begin{array}{l} \ell \mapsto [i; g] \ * \ \text{counter } i \ g \ 1 \ 0 \\ \ell \Leftarrow \{\pi\} \ * \ \ell \Leftarrow \emptyset \\ i \Leftarrow \emptyset \ * \ i \Leftarrow \{+\ell\} \\ g \Leftarrow \emptyset \ * \ g \Leftarrow \{+\ell\} \end{array} \right\} \\
[\emptyset] \{\text{counter } i \ g \ p \ n\} \quad \pi: (i \ [])_{\text{clo}} \quad \{\lambda(). \text{counter } i \ g \ p \ (n + 1)\} \\
[\emptyset] \{\text{counter } i \ g \ p \ n\} \quad \pi: (g \ [])_{\text{clo}} \quad \left\{ \lambda m. \ulcorner n \leq m \wedge (p = 1 \implies n = m) \urcorner \right. \\
\left. \begin{array}{l} \text{counter } i \ g \ p \ n \\ \text{counter } i \ g \ 1 \ n \\ i \Leftarrow \emptyset \ * \ i \Leftarrow \emptyset \\ g \Leftarrow \emptyset \ * \ g \Leftarrow \emptyset \end{array} \right\} \quad \Rightarrow \quad (\diamond 5)
\end{array}$$

Fig. 39. Code and specification of a concurrent monotonic counter

The equivalence axiom in Figure 39 shows that “*counter*” assertions can be split and joined; both the fraction and the contribution are then split or joined by addition. This allows a counter to be used in a concurrent setting: the user can split the “*counter*” predicate into several parts and give a part to each participating thread. In the end, the user can gather all parts, draw conclusions about the final value of the counter, and logically deallocate the counter.

The specification of $(\text{create } [])_{\text{ptr}}$ states that this call consumes 7 space credits (1 credit for the shared reference, 2 credits for each closure, and 2 credits for the pair). It returns a pair ℓ of two locations i and g such that *counter* $i \ g \ 1 \ 0$ holds. This assertion captures the full ownership of the counter, and specifies that its current value is 0.

Figure 39 also shows the specifications of calls to i and g . Both calls require an assertion of the form *counter* $i \ g \ p \ n$. The postcondition of a call to the “increment” closure contains an updated assertion *counter* $i \ g \ p \ (n + 1)$. The postcondition of a call to the “get” closure contains an unmodified “*counter*” assertion. Furthermore, it guarantees that the natural number m that is returned by this call is no less than the past contribution n and, in the case where p is 1, is equal to the past contribution.

Last, Figure 39 shows the reasoning rule for deallocating a counter. This rule requires full ownership of the counter as well as pointed-by-heap and pointed-by-thread assertions for the closures i and g , with fraction 1 and empty sets—this witnesses that both closures are unreachable. In exchange, the rule produces 5 space credits. The 2 credits corresponding to the pair produced by *create* can be recovered independently.

Proof insights. The proof that the counter obeys its specification uses ghost state in a standard way [Birkedal and Bizjak 2023, §8.7]. The internal definition of the abstract predicate “*counter*” involves an existential quantification over the shared location r : indeed, this location does not appear in the specification. The assertion *counter* $i \ g \ p \ n$ contains an empty pointed-by-thread

$$\begin{array}{l}
\text{create} \triangleq \mu_{\text{ptr}} \lambda []. \\
\quad (\text{ref } [0])_{\text{ptr}} \\
\text{async} \triangleq \mu_{\text{ptr}} \lambda [l, f]. \\
\quad (\text{faa } [l, 0, 1])_{\text{ptr}} ; \\
\quad \text{fork} ((f [])_{\text{clo}} ; (\text{ignore } [(f \text{aa } [l, 0, -1])_{\text{ptr}}])_{\text{ptr}})
\end{array}
\qquad
\begin{array}{l}
\text{finish} \triangleq \mu_{\text{ptr}} f. \lambda [l]. \\
\quad \text{if } l[0] = 0 \\
\quad \text{then } () \\
\quad \text{else } (f [l])_{\text{ptr}}
\end{array}$$

$$\begin{array}{l}
\text{AFCREATE} \\
[\emptyset] \{ \diamond 1 \} \pi : (\text{create } [])_{\text{ptr}} \{ \lambda \ell. \text{AF } \ell * \ell \leftarrow_{\frac{1}{2}} \{ \pi \} * \ell \leftarrow_1 \emptyset \}
\end{array}$$

$$\frac{\text{AFASYNC} \quad \forall \pi'. \quad [\{ \ell \}] \{ f \Leftarrow_p \{ \pi' \} * \Phi \} \pi' : (f [])_{\text{clo}} \{ \lambda (). \Psi \}}{[\{ \ell \}] \{ \text{AF } \ell * f \Leftarrow_p \{ \pi \} * \Phi \} \pi : (\text{async } [l, f])_{\text{ptr}} \{ \lambda (). \text{spawned } \ell \Psi \}}$$

$$\begin{array}{l}
\text{AFFINISH} \\
[\emptyset] \{ \text{AF } \ell * \ell \leftarrow_{\frac{1}{2}} \{ \pi \} \} \pi : (\text{finish } [l])_{\text{ptr}} \{ \lambda (). \text{finished } \ell \}
\end{array}$$

$$\begin{array}{l}
\text{FINISHEDSPAWNED} \\
\text{finished } \ell * \text{spawned } \ell \Psi \quad \Rightarrow \quad \Psi
\end{array}
\qquad
\begin{array}{l}
\text{FINISHEDFREE} \\
\text{finished } \ell * \ell \leftarrow_1 \emptyset \quad \Rightarrow \quad \diamond 1
\end{array}$$

$$\begin{array}{l}
\text{AFPERSISTENT} \\
\text{AF } \ell \text{ is persistent}
\end{array}
\qquad
\begin{array}{l}
\text{FINISHEDPERSISTENT} \\
\text{finished } \ell \text{ is persistent}
\end{array}$$

Fig. 40. Code and specification of an async/finish library

assertion for the location r with fraction p . Moreover, the assertion *counter i g p n* contains *Spec* assertions (§9.6) for the closures i and g . The environments that appear in these *Spec* assertions map r to the fraction $\frac{1}{2}$, which means that each closure owns one half of the pointed-by-heap assertion for the location r .

The proof of the logical deallocation rule for a counter (that is, the last rule in Figure 39) is straightforward. We first deallocate the closures i and g , and recover 2×2 space credits as well as an empty pointed-by-heap assertion for the shared location r . Then, by exploiting the empty pointed-by-thread assertion for r , which is contained inside the *counter* assertion, we logically deallocate the location r , thereby recovering one more space credit. In total, 5 space credits are recovered, as expected.

11.4 An Async/Finish Library

The async/finish paradigm was introduced in X10 [Charles et al. 2005; Lee and Palsberg 2010] as a generalization of the spawn/sync mechanism of Cilk [Blumofe et al. 1996], which itself was a generalization of the fork/join paradigm. where exactly two child threads are spawned and awaited. The async/finish paradigm allows spawning an arbitrary number of tasks before waiting at a common join point. More precisely, “async” allows spawning new tasks, whereas “finish” performs synchronization: it blocks until all previously spawned tasks terminate. In this section, we encode these constructs in LambdaFit using a shared mutable reference that is updated via a fetch-and-add operation (§11.2). We then provide specifications in IrisFit, and show that the space

credits associated to the shared reference can be recovered as soon as “finish” returns.¹⁵ A strength of our specification is that it allows for *nested* spawns: a spawned task can itself spawn tasks.

Code. The code of our `async/finish` library is presented in the top part of Figure 40. The library uses a reference that we call the *session*. A session is a channel through which tasks communicate. It stores the number of currently running tasks.

The function `(create [])ptr` returns a fresh session, with zero running tasks.

The function `(async [l, f])ptr` expects a session l and a closure f as arguments. It first atomically increments the session, hence recording the existence of a new running task, then forks off a thread that invokes the closure f with no arguments. When this invocation terminates, it atomically decrements the session, thereby recording that this task is finished.

The function `(finish [l])ptr` consists of an active waiting loop. This loop ends when it observes that the session contains the value 0, which guarantees that all previously spawned tasks have terminated.

Specifications. The bottom part of Figure 40 presents the specification of our `async/finish` library.

According to `AFCREATE`, `(create [])ptr` consumes one space credit, which corresponds to the space occupied by the session, and returns a location ℓ such that `AF ℓ` holds. This persistent assertion guarantees that ℓ is a session. The postcondition also provides pointed-by-thread and pointed-by-heap assertions for the location ℓ . The pointed-by-heap assertion carries the fraction $\frac{1}{2}$; the other half is hidden from the user.

The specification of `(async [l, f])ptr` is stated as a triple featuring a souvenir on ℓ . This means that, for the duration of this call, ℓ is a root. The precondition requires ℓ to be a session. A fractional pointed-by-thread assertion for the closure f , as well as an arbitrary assertion Φ , are consumed and transmitted to the new task, which invokes the closure f . The premise of the rule `AFASYNC` requires the user to prove that, under an arbitrary thread identifier π' , this invocation is safe and satisfies some postcondition Ψ . The postcondition of `(async [l, f])ptr` provides a witness that this task was spawned, in the form of the assertion `spawned ℓ Ψ` . This assertion is not persistent: it can be understood as a unique permission to collect Ψ once the task is finished.

The specification of f in the premise of `AFASYNC` is again a triple with a souvenir of ℓ . This formulation allows f to itself use `async`. Using an ordinary triple there would place a stronger requirement on f and would forbid the use of `async` inside f .

According to `AFFINISH`, `(finish [l])ptr` consumes the pointed-by-thread assertion that was produced by `create`. This forbids any further use of the session ℓ : indeed, both `AFASYNC` and `AFFINISH` require a pointed-by-thread assertion for ℓ .¹⁶ The postcondition contains the persistent assertion `finished ℓ` , which witnesses that this session has been ended.

The ghost update `FINISHEDSPAWNED` states that if the witness `finished ℓ` is at hand then the assertion `spawned ℓ Ψ` can be converted to Ψ . This reflects the idea that if the session has been ended, then all tasks must have terminated: so, a permission to collect Ψ can indeed be converted to Ψ . The ghost update `FINISHEDFREE` states that if the session has ended then abandoning the pointed-by-heap assertion for ℓ allows recovering the space credit associated with the session ℓ .

Proof insights. The assertion `AF ℓ` is internally defined as an Iris invariant, with a part consisting of a liveness-based cancellable invariant (§5.10). Among other things, this invariant imposes a protocol on the pointed-by-thread assertion for the session ℓ . Initially, the invariant contains a pointed-by-thread assertion carrying the fraction $\frac{1}{2}$ and an empty set; the other half is given to the user by

¹⁵That is to say, as soon as every task reaches the linearization point of the fetch-and-add operation to signal that it is done. A task can still execute some code past the linearization point before actually terminating.

¹⁶In the case of `AFASYNC`, this is implicit in the fact that the conclusion of the rule is a triple with a souvenir on ℓ .

$$\begin{array}{l}
\text{STACKCREATE} \\
[\emptyset]\{\diamond 1\} \pi : (\text{create } [])_{\text{ptr}} \{ \lambda \ell. \text{stack } \ell [] * \ell \Leftarrow \{\pi\} * \ell \Leftarrow \emptyset \} \\
\\
\text{STACKPUSH} \\
[\{\ell\}] \left\langle \frac{\diamond 2 * v \Leftarrow_p \{\pi\} * v \Leftarrow_q^0 \emptyset}{\forall vpqs. \text{stack } \ell \text{ } vpqs} \right\rangle \pi : (\text{push } [\ell; v])_{\text{ptr}} \left\langle \frac{\lambda(). \lceil \text{True} \rceil}{\text{stack } \ell ((v, p, q) :: vpqs)} \right\rangle \\
\\
\text{STACKPOP} \\
[\{\ell\}] \left\langle \frac{\lceil \text{True} \rceil}{\forall v p q vpqs. \text{stack } \ell ((v, p, q) :: vpqs)} \right\rangle \pi : (\text{pop } [\ell])_{\text{ptr}} \left\langle \frac{\lambda w. \lceil w = v \rceil * v \Leftarrow_p \{\pi\}}{\text{stack } \ell \text{ } vpqs * \diamond 2 * v \Leftarrow_q^0 \emptyset} \right\rangle \\
\\
\text{STACKFREE} \\
\text{stack } \ell \text{ } vpqs * \ell \Leftarrow \emptyset * \ell \Leftarrow \emptyset \Rightarrow \diamond(1 + 2 \times |vpqs|) * \bigstar_{(v,p,q) \in vpqs} (v \Leftarrow_p \emptyset * v \Leftarrow_q^0 \emptyset)
\end{array}$$

Fig. 41. Specification of Treiber's Stack

AFCREATE. Each spawned task gets a fraction of this assertion: indeed, spawning a task involves “fork”, and our **FORK** rule requires updating a pointed-by-thread assertion so as to reflect the fact that ℓ is a root of the new thread. When a task signals that it is finished, it surrenders its fractional pointed-by-thread assertion, carrying an *empty* set of thread identifiers. Hence, once every task has terminated, the invariant again contains $\ell \Leftarrow_{\frac{1}{2}} \emptyset$.

How and when exactly does a task signal that it is finished? This is done via a fetch-and-add (FAA) operation, which decrements the count of active tasks, and takes effect precisely at the linearization point of this FAA operation. Hence, as soon as this linearization point is reached, the invariant requires this task to surrender its fractional pointed-by-thread assertion. Fortunately, our specification of FAA (§11.2) allows this: the pointed-by-thread assertion $\ell \Leftarrow_p \emptyset$ appears in the public postcondition in **FAA**.

The absence of a “later” modality in front of Ψ in **FINISHEDSPAWNED** may seem surprising. As the assertion Ψ has transited through an invariant, an Iris expert might expect it to be guarded by such a modality. The usual way to eliminate a “later” modality is through a physical step, yet this rule is a ghost update. Fortunately, IrisFit supports and takes advantage of *later credits* (§6.2). A later credit is a piece of ghost state that is produced by a physical step and that can later be used to eliminate a “later” modality. With each spawned task, we are able to internally associate one later credit, which we obtain from the function call $(\text{async } [\ell, f])_{\text{ptr}}$. By exploiting this later credit, we can eliminate the “later” modality in front of Ψ before giving this assertion back to the user.

11.5 Treiber's Stack

Code. The code that we verify is the code of Figure 3, translated to LambdaFit syntax. A reference is a block of size 1 and a list cell is a block of size 2.

Specifications. Figure 41 presents our specification of Treiber's stack. The stack is described in terms of the abstract predicate $\text{stack } \ell \text{ } vpqs$, where ℓ is the location of the stack and $vpqs$ is its mathematical model. This model is a list of triples (v, p, q) of a value v and two positive fractions p and q . The list of the values v describes the content of the stack. For each value v , the fractions p and q describe what quantity of the pointed-by-thread and pointed-by-heap assertions for the value v have been acquired by the stack. Having the stack acquire a fractional pointed-by-heap assertion for the value v lets us record that this value is pointed to by a list cell without revealing or even mentioning the address of this cell. Having the stack acquire a fractional pointed-by-thread

assertion for the value v lets us express a plausible specification for “pop”. Indeed, “pop” needs to read the value v from the heap: then, the **LOAD** rule requires (and updates) a fractional pointed-by-thread assertion for v . Expecting the caller to supply this assertion seems impractical, so it must be found in the stack itself.

The assertion $stack\ \ell\ v p q s$ is not fractional: it represents the full ownership of the stack. To allow the stack to be accessed by several concurrent threads, the user must share this assertion. This is typically achieved via an Iris invariant [Birkedal and Bizjak 2023].

According to **STACKCREATE**, creating a new stack consumes one space credit. This is the size of the reference that holds the address of the top list cell. The result is a fresh location ℓ that represents an empty stack.

The specification of $(push\ [\ell; v])_{ptr}$, expressed by **STACKPUSH**, is an atomic triple with a souvenir on ℓ . The private precondition requires two space credits, which is the size of a new list cell, as well as fractional pointed-by-heap and pointed-by-thread assertions for the value v that is pushed onto the stack. Together, the public precondition and postcondition indicate that the model of the stack is atomically updated from $v p q s$ updated to $(v, p, q) :: v p q s$ at the linearization point.

The specification of $(pop\ [\ell])_{ptr}$, expressed by **STACKPOP**, is also an atomic triple with a souvenir on ℓ . The public precondition and postcondition indicate that the model of the stack is atomically updated from $(v, p, q) :: v p q s$ to $v p q s$. Furthermore, according to the public postcondition, at the linearization point, two space credits are produced, and a pointed-by-heap assertion for v , carrying an empty multiset of predecessors, is produced as well, as a pointer from the stack to v has been destroyed.

Our specification of “pop” exhibits a certain asymmetry: whereas the space credits and the pointed-by-heap assertion appear in the *public* postcondition, which means that they are produced at the linearization point, the pointed-by-thread assertion appears in the *private* postcondition, which means that it is produced when the function returns. The space credits and the pointed-by-heap assertion can be produced at the linearization point because there we are already able to logically deallocate the list cell and to argue that a pointer from the stack to v has been destroyed. However, the pointed-by-thread assertion cannot be surrendered as part of the public postcondition, because the value v is read from the heap *after* the linearization point has been passed.

The last rule in Figure 41, **STACKFREE**, logically deallocates a (possibly nonempty) stack. The assertion $stack\ \ell\ v p q s$, as well as empty pointed-by-thread and pointed-by-heap assertions for ℓ , are consumed. A number of space credits are produced, which reflect the overall size occupied by the stack data structure in the heap: one credit for the toplevel reference, plus two credits per list cell. The pointed-by-thread and pointed-by-heap assertions associated with every triple (v, p, q) in the stack are also produced. Of course, in the common case where $v p q s$ is an empty list, this rule can be significantly simplified.

Proof insights. As argued earlier (§3), the main difficulty of the proof is to produce space credits when a “pop” operation succeeds. This requires logically deallocating the list cell that is being extracted. This in turn requires exhibiting both an empty pointed-by-thread assertion and an empty pointed-by-heap assertion for this cell. Yet, neither of these assertions is easy to obtain.

Let us discuss the pointed-by-thread assertion first. The difficulty is that “push” and “pop” are *invisible readers* [Alistarh et al. 2018]: these operations read the top of the stack (that is, the address of a list cell) without synchronization. Such a read normally requires updating a pointed-by-thread assertion for the cell whose address is thus obtained. However, here, we do not wish to record that this cell is pointed to by the current thread. Fortunately, these reads occur inside protected sections. Hence, we use **LOADINSIDE**, which updates an “inside” assertion instead of a pointed-by-thread assertion. This allows the stack’s invariant to keep an *empty* pointed-by-thread assertion, at all

times, for every list cell. This in turn allows a successful “pop” operation to extract this empty pointed-by-thread assertion out of the invariant. Maintaining empty pointed-by-thread assertions for locations that are acquired only inside protected sections is a typical idiom.

Next, let us discuss the pointed-by-heap assertion. Here, the difficulty is that a list cell ℓ may be pointed to by a new cell ℓ' that has just been allocated by an ongoing “push” operation. This scenario was discussed earlier (§3.2). Hence, each ongoing “push” holds an assertion $\ell \leftarrow_p \{+\ell'\}$, where ℓ is the list cell that “pop” is attempting to extract and ℓ' is the new list cell that “push” has allocated. Now, how can “pop” obtain the assertion $\ell \leftarrow_1 \emptyset$ that is required to allow logical deallocation? We answer this question via an original technique that we dub *logical deallocation by proxy*: the thread that successfully pops the list cell ℓ also takes care of logically deallocating the predecessor cells ℓ' that have been allocated by ongoing “push” operations.¹⁷ The logical deallocation of these locations is made possible by the protected section in “push”. This approach has a somewhat strange consequence: in the proof of “push”, it may be the case that the cell ℓ' has been logically deallocated by another thread, yet “push” still needs to access this cell. Fortunately, IrisFit allows this: for example, the proof of “push” makes use of the rule `STOREDEAD`.

12 RELATED WORK

12.1 Polling Points

A stop-the-world event may be viewed as an asynchronous interruption: a thread that emits such an event stops the execution of all other threads. Such an interruption can be implemented using hardware interrupts, but this scheme can be expensive and non-portable [Feeley 1993]. Another approach is to let the compiler insert explicit tests for interruptions into the code. These tests appear in the literature under various names, including *polling points* [Feeley 1993], *GC points* [Agesen 1998], *yield points* [Lin et al. 2015], and *safe points* [Sivaramakrishnan et al. 2020]. Let us refer to them collectively as *safe points*. Safe points are typically inserted by the compiler in such a way that no computation can run forever without encountering a safe point. When a thread encounters a safe point, it tests whether some other thread has requested garbage collection. If so, it pauses and passes control to the runtime system. Once all threads have paused in this way, the runtime system performs a global garbage collection phase.

Safe points are used in the Jalapeño/Jikes RVM [Alpern et al. 1999, 2005] and in OCaml 5 [Sivaramakrishnan et al. 2020]. The existence of safe points is not revealed to the programmer, who is not expected to know about their existence and is given no means of controlling their placement. As an experimental feature, the OCaml 5 compiler does offer a `[@poll error]` attribute [Jaffer 2021]. This attribute is placed on a function definition. An attempt by the compiler to insert a safe point into a function that carries this attribute causes a compile-time error. This lets the programmer check that a function body does not contain any safe point, therefore is (de facto) a protected section. At this time, there is not a clear consensus whether this feature is useful and corresponds to the needs of expert programmers.

Safe points, as described above, and polling points, as proposed in this paper, are two related yet distinct concepts. Indeed, in our view, safe points play two distinct roles. On the one hand, they are points where a thread *must* stop and allow garbage collection to take place. On the other hand, they dictate where garbage *can* take place: indeed, the GC cannot run unless every thread has reached a safe point. Because of this dual role, safe points do *not* enjoy a monotonicity property. When a safe point is inserted in a program, the set of possible behaviors of this program is neither shrunk nor enlarged; it is transformed into an incomparable set. We believe that our design, where polling points and protected sections are separate concepts, is better behaved. In particular, it enjoys several

¹⁷Note that these ongoing “push” operations will fail, because the top list cell previously observed has been replaced.

monotonicity properties. Inserting a new polling point shrinks the set of possible behaviors of the program.¹⁸ Creating or enlarging a protected section shrinks the family of programs that is implicitly described by a source program.

In our approach, the user *explicitly* inserts enough protected sections to (verifiably) obtain the desired worst-case heap space complexity, then lets the compiler *implicitly* insert enough polling points to guarantee liveness, without endangering the program’s space complexity. This is expressed by Theorem 8.2.

12.2 Protected Sections

In the production systems that we are aware of, the concept that seems closest to our protected sections appears in the .NET runtime system, where it was introduced in 2015, with performance in mind [Lander 2015]. The API of the GC module [Microsoft 2024] provides a method `TryStartNoGCRegion(Int64)` and a method `EndNoGCRegion()`. A “NoGC region” is not quite a protected section in our sense, though, as allocation is permitted inside a “NoGC region”. The integer parameter of the method `TryStartNoGCRegion` is a request for a certain amount of free heap space: garbage collection takes place at this point so as to guarantee that this much free space exists. Allocation requests within the “NoGC region” are then served out of this pre-allocated free space. However, if the runtime system runs out of free space while some thread is inside a “NoGC region”, then garbage collection will take place.

Feeley [1993, §1.2.1] discusses why “critical sections”—sections in which the GC must not run—may be needed for safety reasons. He takes the example of a store instruction that stores a 64-bit pointer into memory and that is decomposed into two 32-bit stores. In between the two stores, the memory is in an inconsistent state and must not be read by the GC.

To the best of our knowledge, our paper is the first where a notion of protected section is introduced for complexity reasons, that is, with the aim of guaranteeing tighter worst-case heap space complexity bounds.

12.3 Reasoning about Space without a GC

Hofmann [1999, 2003] introduces space credits in the setting of an affine type system for the λ -calculus. Hofmann [2000] and Aspinall and Hofmann [2002] adapt the idea to LFPL, a first-order functional programming language without GC and with explicit destructive pattern matching. There, a value of type \diamond exists at runtime and can be understood as a pointer to a free block in the heap. Subsequent work aims at automating space complexity analyses. In particular, Hofmann and Jost [2003] propose an affine type system where types carry space credits. Hofmann and Jost [2006]; Hofmann and Rodriguez [2009, 2013] analyze a variant of Java where garbage collection has been replaced with explicit deallocation. RaML [Hoffmann et al. 2012a,b, 2017] analyzes a fragment of OCaml, also without GC and with explicit destructive pattern matching. Niu and Hoffmann [2018] present a type-based amortized space analysis for a pure, first-order programming language where destructive pattern matching can be applied to shared objects, an unusual feature. Their system performs significant over-approximations: when a data structure becomes shared, the logic charges the cost of creating a copy of this data structure. As far as we understand, this analysis can be used to reason in a sound yet very conservative way about a programming language with GC. Kahn and Hoffmann [2021] present a system that is equipped with more flexible typing rules than its predecessors and therefore can derive tighter resource consumption bounds. Hoffmann and Jost [2022] offer a survey of two decades of work on automated amortized resource analysis (AARA).

¹⁸In our setting, inserting a new polling point does not create a new opportunity for the GC to run, because the GC is always enabled anyway.

Following the ideas of LFPL, Lorenzen et al. [2023] introduce a calculus with “reuse” credits. Explicit destructive pattern matching produces reuse credits, which can be used to satisfy a new allocation. Because the system allows fragmentation, reuse credits cannot be joined. The goal of Lorenzen et al. [2023] is to statically detect *fully in-place* functions—that is, functions that do not need to allocate new memory. This includes, for example, functions that reuse the heap space occupied by their arguments.

Chin et al. [2005, 2008] present a type system that automatically keeps track of data structure sizes. The type system incorporates an alias analysis, which distinguishes between shared and unique objects and allows unique objects to be explicitly deallocated. Shared objects can never be logically deallocated. Specifications indicate how much memory a method may need (a high-water mark) and how much memory it releases, in terms of the sizes of the arguments and results.

Compared with type systems, program logics offer weaker automation but greater expressiveness. Aspinall et al. [2007] propose a VDM-style program logic, where postconditions depend not only on the pre-state, post-state, and return value, but also on a cost. Atkey [2011] proposes an extension of Separation Logic with an abstract notion of resource, such as time or space, and introduces an assertion that denotes the ownership of a certain amount of resources.

All of the work cited above concerns languages with explicit memory deallocation, where there is no need to reason about unreachability. Reasoning about unreachability in the setting of a static analysis or program logic is a central challenge.

12.4 Reasoning about Space with a GC

Hur et al. [2011] propose a Separation Logic for the combination of a low-level language with explicit deallocation and a high-level language with a GC. They are interested in verifying just safety, not space complexity.

Madiot and Pottier [2022] and Moine et al. [2023] propose Separation Logics that allow reasoning about space in the presence of a GC.

The logic presented by Madiot and Pottier [2022] concerns a low-level language where roots are explicitly marked as such inside the heap via so-called “stack cells”. Programming in this artificial language is impractical because local variables must be manually registered and unregistered as stack cells. Madiot and Pottier’s language and program logic technically support concurrency, but the paper does not provide any case study.

The logic presented in our previous paper [Moine et al. 2023] concerns a high-level language, where the roots and the call stack are implicit, but is restricted to a sequential setting. This paper also introduces support for closures. The logic relies on a distinction between *visible roots*—the roots of the term under focus—and *invisible roots*—the roots of the evaluation context. The logic keeps track of invisible roots using a *Stackable* assertion, and introduces the idea that *Stackable* assertions must be “forcibly framed out” at applications of the `BIND` rule. We re-use this idea in our own `BIND` rule (§6.4), but replace *Stackable* assertions with pointed-by-thread assertions, which are better suited to a concurrent setting. In so doing, we remove the distinction between visible roots and invisible roots, which does not seem to make sense in a concurrent setting; our pointed-by-thread assertions keep track of all (ordinary) roots. In contrast, Moine et al. [2023] do not keep track of visible roots via an a dedicated assertion: indeed, in their setting, it suffices to inspect the term under focus to determine the set of visible roots. This allows them to offer a standard `LOAD` rule, whereas our `LOAD` rule updates a pointed-by-thread assertion for the value that is loaded (§6.2).

Our mechanization [Moine 2024b] includes an encoding inside IrisFit of our previous logic for sequential programs [Moine et al. 2023]. This encoding demonstrates that our concurrent program logic can be used to reason about sequential programs with no overhead.

12.5 Space-Related Results for Compilers

Paraskevopoulou and Appel [2019] prove that, in the presence of a GC, closure conversion is safe for space: that is, it does not change the space consumption of a program. They view closure conversion as a transformation from a CPS-style λ -calculus into itself. This calculus is equipped with two different environment-based big-step operational semantics. The “source” semantics implicitly constructs a closure for each function definition by capturing the relevant part of the environment and storing it in the heap. The “target” semantics performs no such construction: it requires every function to be closed. In either semantics, the roots are defined as the locations that occur in the environment. Up to the stylistic difference between a substitution-based semantics and an environment-based semantics, this definition is equivalent to the “free variable rule” (FVR) [Morrisett et al. 1995].

Besson et al. [2019] prove that (an enhanced version of) CompCert [Leroy 2024] preserves memory consumption when compiling C programs.

In a sequential setting, Gómez-Londoño et al. [2020] prove that the CakeML compiler respects a cost model that is defined at the level of the intermediate language DataLang, which serves as the target of closure conversion. Our cost model is analogous to theirs. Our work and theirs are complementary: whereas they prove that the CakeML compiler respects the DataLang cost model, we show how to establish space complexity bounds about source programs, based on a similar cost model. One could in principle adapt IrisFit to DataLang. Then, one would be able to use IrisFit to establish a space complexity bound about a source CakeML program, to compile this program down to machine code using the CakeML compiler, and to obtain a machine-checked space complexity guarantee about the compiled code.

12.6 Safe Memory Reclamation Schemes

Manual memory management can be so difficult in a concurrent setting that programmers often rely on semi-automatic *safe memory reclamation* (SMR) schemes. Two main families exist, namely hazard pointers [Michael 2004; Michael et al. 2023] and read-copy-update (RCU) [McKenney 2004; McKenney et al. 2023]. The two families offer roughly similar APIs. First, the user declares *hazardous* locations for a delimited scope. While it is marked hazardous, a location is not deallocated. Second, the user can *retire* a location to indicate that this location is no longer needed. The SMR implementation deallocates a retired location once it is not marked hazardous by any thread.

Hazard pointers operates selectively: the user manually marks and unmarks pointers as hazardous. On the contrary, RCU declares *every* non-retired pointer hazardous inside a certain section of the code. RCU seems close to our concept of a protected section. Indeed, RCU delays the deallocation of a retired pointer until every RCU section that mentions this pointer terminates. Yet, there is not a perfect analogy between RCU sections and protected sections. Indeed, garbage collection provides a strong guarantee: *no dangling pointer can exist*. SMR schemes, on the contrary, can create dangling pointers, which client data structures often tolerate. For example, with RCU, a location that is mentioned in the code, but is not read or written, does not need to be protected. For example, in the “push” operation of Treiber’s stack (Figure 1), the address of the first cell of the internal list (named `h` in the code) is loaded from the heap and is a root throughout a certain section of the code, but is never dereferenced. Hence, “push” *can* tolerate the deallocation of this cell by a successful concurrent pop operation, and “push” does *not* need an RCU section [Jung et al. 2023, mechanization]. On the contrary, the “pop” operation does need a RCU section. Indeed, the address of the first cell of the internal list (again named `h` in the code) is dereferenced twice: once during the call to the function `is_nil` and once by the CAS instruction.

Equipping SMR schemes with abstract Separation Logic specifications and verifying them has long been a challenge. Treiber’s stack has been the first data structure based on hazard pointers to be verified. This task was tackled several times using different variants of Concurrent Separation Logic [Parkinson et al. 2007; Fu et al. 2010]. Tofan et al. [2011] verify Treiber’s stack both with hazard pointers and with garbage collection (though without a heap space complexity analysis). They show that a large part of the main invariant can be shared between the two proofs. Gotsman et al. [2013] provide the first general framework for verifying programs using SMR schemes in Separation Logic, making use of temporal logic reasoning. Jung et al. [2023] provide a more abstract framework, where temporal reasoning is replaced with ownership arguments. Their work unveils a close relationship between RCU and garbage collection. Indeed, RCU allows accessing any location that was *not* retired when the current RCU section was entered. (There is a loose analogy with our liveness-based cancellable invariants: to access such an invariant, one must eliminate the case where ℓ has been logically deallocated.) To prove that a location is *not* retired at a certain point in time, Jung et al. [2023] express the topology of data structures using pointed-by-heap assertions, which they borrow from our prior paper [Moine et al. 2023]. Like us, when retiring a location, they require the predecessors of this location to have been previously retired.

Outside the Separation Logic world, Meyer and Wolff [2019] propose an API for SMR schemes, in the form of an observer automaton, inspired by the temporal reasoning of Gotsman et al. [2013]. Meyer and Wolff [2019] make use of the observer automaton to de-correlate the verification of lock-free data structures from the SMR implementation, allowing them to develop an automatic linearizability checker.

13 CONCLUSION

We have presented LambdaFit, a language with shared-memory concurrency and tracing garbage collection. In particular, LambdaFit is equipped with protected sections, a new, realistic construct that programmers can and sometimes must exploit to ensure that fine-grained concurrent data structures have the desired worst-case heap space complexity. We believe that protected sections could be a useful part of a concurrent programmer’s toolbox, and that they should be considered for inclusion in high-level languages.

Furthermore, we have presented IrisFit, a Concurrent Separation Logic with space credits, which allows expressing and verifying worst-case heap space bounds about LambdaFit programs. IrisFit features pointed-by-heap and pointed-by-thread assertions, which offer a compositional means of keeping track of the various ways through which a memory block is reachable. These assertions can be used to prove that a block is unreachable, or more accurately, that by the time the garbage collector is allowed to run, this block will be unreachable. IrisFit provides special treatment of temporary roots within protected sections and is thereby able to take advantage of protected sections to establish stronger worst-case heap space bounds.

All of our results are mechanized in the Coq proof assistant using the Iris library [Jung et al. 2018b] and its dedicated Proof Mode [Krebbbers et al. 2018]. Our definitions and proofs are available in electronic form [Moine 2024b]. Discounting blank lines and comments, the definition of LambdaFit and of its oblivious semantics occupy roughly 2800LOC; the construction of IrisFit, including the reasoning rules and the core soundness theorem, represent 9200LOC; the definition of the default semantics of LambdaFit and the proof of the safety and liveness theorems take up 4500LOC; and the verification of the case studies represents 6400LOC. In addition to these numbers, we re-use about 3700LOC of proofs from Madiot and Pottier [2022] and from our own previous work [Moine et al. 2023]. We provide tactics that facilitate reasoning with IrisFit and achieve a basic level of automation thanks to the Diaframe library [Mulder et al. 2022].

14 FUTURE WORK

We now propose avenues for future work on LambdaFit and IrisFit.

More Liberal Protected Sections. LambdaFit forbids function calls inside protected sections. This is a simple way of ensuring that a protected section is exited in a bounded number of steps. One might wish to relax this restriction and tolerate calls to “small” functions (that is, functions whose execution requires a bounded number of steps) inside protected sections. The obligation of proving that every protected section terminates in bounded time would then have to be delegated to the user of IrisFit, perhaps via a variation on time credits [Charguéraud and Pottier 2019]. One could relax this restriction even further by requiring only that every protected section terminates (without a statically known bound). The proof obligation would again be delegated to the user of IrisFit. This would allow traversing a data structure of unbounded size inside a protected section. Harris’s list [Harris 2001; Michael 2002] is an example where this is required. Indeed, the *search* function searches for an element in a linked list. As long as this function is running, the location of at least one internal list cell is a root. Hence, in order to respect the principle that “an internal cell may be a root only inside a protected section”, the whole *search* loop should be wrapped inside a protected section.

If the restrictions that bear on protected sections are relaxed, then the polling point insertion strategy must be adapted accordingly. Our current polling point insertion strategy, *addpp*, inserts a polling point in front of each function call. This is a simple way of ensuring that every execution path must, in a bounded number of steps, reach a polling point. However, polling points must not be inserted inside protected sections. So, if a function call appears inside a protected section, then it must not be preceded by a polling point.

Additional Case Studies. We would like to apply IrisFit to more ambitious case studies. This includes larger examples as well as subtler concurrent examples, including multi-CAS algorithms such as RDCSS [Harris et al. 2002], Harris’s list [2001], or its variant due to Michael [2002]. As explained in the previous paragraph, Harris’s list would require unbounded protected sections. Additionally, Harris’s list features *lazy deletion*, in which a node is first marked as deleted, before an attempt is made to physically unlink it from the structure. If this attempt fails, another function call may perform this unlinking operation. It is presently unclear to us what the specification of Harris’s list delete function would be and whether protected sections would allow this function to be verified.

Immutable Data Structures. We would like to determine whether immutable data structures could be specified and verified in a more pleasant and lightweight manner. At present, IrisFit offers no special support for immutable data structures: every memory block is considered mutable by default, and it is up to the user to exploit the logical tools offered by Iris, such as invariants, to indicate that a memory block is immutable. In this paper, we have done so in the special case of closures: we have been able to describe the behavior of a closure via a *persistent* predicate, while still allowing for its deallocation. We would like to investigate whether this approach can be extended to all immutable data structures.

IrisFit as a Foundation for Type Systems or Static Analyses. We would like to draw upon our experience with IrisFit to investigate automated static analyses of the worst-case heap space complexity of a program in the presence of garbage collection. As far as we know, relatively few such analyses have been presented in the literature [Braberman et al. 2008; Unnikrishnan and Stoller 2009; Albert et al. 2013; Niu and Hoffmann 2018] and none of them is justified by a machine-checked

argument. It would be interesting to justify existing analyses by reduction to the reasoning rules of IrisFit or to draw inspiration from these rules to design new analyses.

ACKNOWLEDGMENTS

We wish to thank the anonymous reviewers, whose insights and perseverance have significantly helped improve the quality and presentation of this paper.

REFERENCES

- Ole Agesen. 1998. *GC Points in a Threaded Environment*. Technical Report SMLI TR-98-70. Sun Microsystems, Inc.
- Elvira Albert, Samir Genaim, and Miguel Gómez-Zamalloa. 2013. [Heap space analysis for garbage collected languages](#). *Science of Computer Programming* 78, 9 (2013), 1427–1448.
- Dan Alistarh, William Leiserson, Alexander Matveev, and Nir Shavit. 2018. [ThreadScan: Automatic and Scalable Memory Reclamation](#). *ACM Trans. Parallel Comput.* 4, 4, Article 18 (May 2018).
- Bowen Alpern, C. Richard Attanasio, John J. Barton, Anthony Cocchi, Susan Flynn Hummel, Derek Lieber, Ton Ngo, Mark F. Mergen, Janice C. Shepherd, and Stephen E. Smith. 1999. [Implementing Jalapeño in Java](#). In *Object-Oriented Programming Systems, Languages & Applications (OOPSLA)*. 314–324.
- Bowen Alpern, Steve Augart, Stephen M. Blackburn, Maria A. Butrico, Anthony Cocchi, Perry Cheng, Julian Dolby, Stephen J. Fink, David Grove, Michael Hind, Kathryn S. McKinley, Mark F. Mergen, J. Eliot B. Moss, Ton Anh Ngo, Vivek Sarkar, and Martin Trapp. 2005. [The Jikes Research Virtual Machine project: Building an open-source research community](#). *IBM Syst. J.* 44, 2 (2005), 399–418.
- Andrew W. Appel. 1992. *Compiling with Continuations*. Cambridge University Press.
- David Aspinall, Lennart Beringer, Martin Hofmann, Hans-Wolfgang Loidl, and Alberto Momigliano. 2007. [A program logic for resources](#). *Theoretical Computer Science* 389, 3 (2007), 411–445.
- David Aspinall and Martin Hofmann. 2002. [Another Type System for In-Place Update](#). In *European Symposium on Programming (ESOP) (Lecture Notes in Computer Science, Vol. 2305)*. Springer, 36–52.
- Robert Atkey. 2011. [Amortised Resource Analysis with Separation Logic](#). *Logical Methods in Computer Science* 7, 2:17 (2011), 1–33.
- Yves Bertot and Pierre Castéran. 2004. *Interactive Theorem Proving and Program Development – Coq’Art: The Calculus of Inductive Constructions*. Springer.
- Frédéric Besson, Sandrine Blazy, and Pierre Wilke. 2019. [CompCertS: a Memory-Aware Verified C Compiler Using a Pointer as Integer Semantics](#). *Journal of Automated Reasoning* 63, 2 (2019), 369–392.
- Lars Birkedal and Aleš Bizjak. 2023. [Lecture notes on Iris: Higher-order concurrent separation logic](#). (2023). Unpublished.
- Lars Birkedal, Thomas Dinsdale-Young, Armaël Guéneau, Guilhem Jaber, Kasper Svendsen, and Nikos Tzevelekos. 2021. [Theorems for free from separation logic specifications](#). *Proceedings of the ACM on Programming Languages* 5, ICFP (2021), 1–29.
- Wayne D. Blizard. 1990. [Negative membership](#). *Notre Dame Journal of Formal Logic* 31, 3 (1990), 346–368.
- Robert D. Blumofe, Christopher F. Joerg, Bradley C. Kuszmaul, Charles E. Leiserson, Keith H. Randall, and Yuli Zhou. 1996. [Cilk: An Efficient Multithreaded Runtime System](#). *J. Parallel Distributed Comput.* 37, 1 (1996), 55–69.
- Richard Bornat, Cristiano Calcagno, Peter O’Hearn, and Matthew Parkinson. 2005. [Permission accounting in separation logic](#). In *Principles of Programming Languages (POPL)*. 259–270.
- John Boyland. 2003. [Checking Interference with Fractional Permissions](#). In *Static Analysis Symposium (SAS) (Lecture Notes in Computer Science, Vol. 2694)*. Springer, 55–72.
- Victor A. Braberman, Federico Javier Fernández, Diego Garbervetsky, and Sergio Yovine. 2008. [Parametric prediction of heap memory requirements](#). In *International Symposium on Memory Management*. 141–150.
- Stephen Brookes and Peter W. O’Hearn. 2016. [Concurrent separation logic](#). *SIGLOG News* 3, 3 (2016), 47–65.
- Quentin Carbonneaux, Jan Hoffmann, and Zhong Shao. 2015. [Compositional certified resource bounds](#). In *Programming Language Design and Implementation (PLDI)*. 467–478.
- Arthur Charguéraud and François Pottier. 2019. [Verifying the Correctness and Amortized Complexity of a Union-Find Implementation in Separation Logic with Time Credits](#). *Journal of Automated Reasoning* 62, 3 (March 2019), 331–365.
- Philippe Charles, Christian Grothoff, Vijay Saraswat, Christopher Donawa, Allan Kielstra, Kemal Ebcioglu, Christoph von Praun, and Vivek Sarkar. 2005. [X10: an object-oriented approach to non-uniform cluster computing](#). In *Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. 519–538.
- Wei-Ngan Chin, Huu Hai Nguyen, Corneliu Popeea, and Shengchao Qin. 2008. [Analysing memory resource bounds for low-level programs](#). In *International Symposium on Memory Management*. 151–160.
- Wei-Ngan Chin, Huu Hai Nguyen, Shengchao Qin, and Martin C. Rinard. 2005. [Memory Usage Verification for OO Programs](#). In *Static Analysis Symposium (SAS) (Lecture Notes in Computer Science, Vol. 3672)*. Springer, 70–86.

- William R. Cook. 2009. [On understanding data abstraction, revisited](#). In *Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*. 557–572.
- Karl Cray and Stephanie Weirich. 2000. [Resource bound certification](#). In *Principles of Programming Languages (POPL)*. 184–198.
- Pedro da Rocha Pinto, Thomas Dinsdale-Young, and Philippa Gardner. 2014. [TaDA: A Logic for Time and Data Abstraction](#). In *European Conference on Object-Oriented Programming (ECOOP) (Lecture Notes in Computer Science, Vol. 8586)*, Richard E. Jones (Ed.). Springer, 207–231.
- Marc Feeley. 1993. [Polling Efficiently on Stock Hardware](#). In *Functional programming languages and computer architecture (FPCA)*. 179–190.
- Matthias Felleisen and Robert Hieb. 1992. [The Revised Report on the Syntactic Theories of Sequential Control and State](#). *Theoretical Computer Science* 103, 2 (1992), 235–271.
- Jean-Christophe Filiâtre. 2011. [Deductive software verification](#). *Software Tools for Technology Transfer* 13, 5 (2011), 397–403.
- Ming Fu, Yong Li, Xinyu Feng, Zhong Shao, and Yu Zhang. 2010. [Reasoning about Optimistic Concurrency Using a Program Logic for History](#). In *International Conference on Concurrency Theory (CONCUR) (Lecture Notes in Computer Science, Vol. 6269)*. Springer, 388–402.
- Alejandro Gómez-Londoño and Magnus O. Myreen. 2021. [A flat reachability-based measure for CakeML’s cost semantics](#). In *Implementation of Functional Languages (IFL)*. 1–9.
- Alejandro Gómez-Londoño, Johannes Áman Pohjola, Hira Taqdees Syeda, Magnus O. Myreen, and Yong Kiam Tan. 2020. [Do you have space for dessert? A verified space cost semantics for CakeML programs](#). *Proceedings of the ACM on Programming Languages* 4, OOPSLA (2020), 204:1–204:29.
- Alexey Gotsman, Noam Rinetzkky, and Hongseok Yang. 2013. [Verifying Concurrent Memory Reclamation Algorithms with Grace](#). In *European Symposium on Programming (ESOP) (Lecture Notes in Computer Science, Vol. 7792)*. Springer, 249–269.
- Theodore Hailperin. 1986. [Formalization of Boole’s Logic](#). In *Boole’s Logic and Probability*. Studies in Logic and the Foundations of Mathematics, Vol. 85. Elsevier, 135–172.
- Timothy L. Harris. 2001. [A Pragmatic Implementation of Non-blocking Linked-Lists](#). In *Proceedings of the 15th International Conference on Distributed Computing (DISC ’01)*. Springer-Verlag, Berlin, Heidelberg, 300–314.
- Timothy L. Harris, Keir Fraser, and Ian A. Pratt. 2002. [A Practical Multi-word Compare-and-Swap Operation](#). In *Distributed Computing*, Dahlia Malkhi (Ed.). Springer Berlin Heidelberg, Berlin, Heidelberg, 265–279.
- Guanhua He, Shengchao Qin, Chenguang Luo, and Wei-Ngan Chin. 2009. [Memory Usage Verification Using Hip/Sleek](#). In *Automated Technology for Verification and Analysis (ATVA) (Lecture Notes in Computer Science, Vol. 5799)*. Springer, 166–181.
- Maurice Herlihy and Nir Shavit. 2012. *The Art of Multiprocessor Programming, Revised Reprint* (1st ed.). Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- Maurice P. Herlihy and Jeannette M. Wing. 1990. [Linearizability: a correctness condition for concurrent objects](#). *ACM Transactions on Programming Languages and Systems* 12, 3 (July 1990), 463–492.
- Jan Hoffmann, Klaus Aehlig, and Martin Hofmann. 2012a. [Multivariate amortized resource analysis](#). *ACM Transactions on Programming Languages and Systems* 34, 3 (2012), 14:1–14:62.
- Jan Hoffmann, Klaus Aehlig, and Martin Hofmann. 2012b. [Resource Aware ML](#). In *Computer Aided Verification (CAV) (Lecture Notes in Computer Science, Vol. 7358)*. Springer, 781–786.
- Jan Hoffmann, Ankush Das, and Shu-Chun Weng. 2017. [Towards automatic resource bound analysis for OCaml](#). In *Principles of Programming Languages (POPL)*. 359–373.
- Jan Hoffmann and Steffen Jost. 2022. [Two decades of automatic amortized resource analysis](#). *Mathematical Structures in Computer Science* 32, 6 (2022), 729–759.
- Martin Hofmann. 1999. [Linear Types and Non-Size-Increasing Polynomial Time Computation](#). In *Logic in Computer Science (LICS)*. 464–473.
- Martin Hofmann. 2000. [A type system for bounded space and functional in-place update](#). *Nordic Journal of Computing* 7, 4 (2000), 258–289.
- Martin Hofmann. 2003. [Linear types and non-size-increasing polynomial time computation](#). *Information and Computation* 183, 1 (2003), 57–85.
- Martin Hofmann and Steffen Jost. 2003. [Static prediction of heap space usage for first-order functional programs](#). In *Principles of Programming Languages (POPL)*. 185–197.
- Martin Hofmann and Steffen Jost. 2006. [Type-Based Amortised Heap-Space Analysis](#). In *European Symposium on Programming (ESOP) (Lecture Notes in Computer Science, Vol. 3924)*. Springer, 22–37.
- Martin Hofmann and Dulma Rodriguez. 2009. [Efficient Type-Checking for Amortised Heap-Space Analysis](#). In *Computer Science Logic (Lecture Notes in Computer Science, Vol. 5771)*. Springer, 317–331.
- Martin Hofmann and Dulma Rodriguez. 2013. [Automatic Type Inference for Amortised Heap-Space Analysis](#). In *European Symposium on Programming (ESOP) (Lecture Notes in Computer Science, Vol. 7792)*. Springer, 593–613.

- Chung-Kil Hur, Derek Dreyer, and Viktor Vafeiadis. 2011. *Separation Logic in the Presence of Garbage Collection*. In *Logic in Computer Science (LICS)*. 247–256.
- Sadiq Jaffer. 2021. OCaml Compiler Pull Request 10462: Add [@poll error] attribute. <https://github.com/ocaml/ocaml/pull/10462>.
- Richard E. Jones and Rafael Dueire Lins. 1996. *Garbage collection – algorithms for automatic dynamic memory management*. Wiley.
- Jaehwang Jung, Janggun Lee, Jaemin Choi, Jaewoo Kim, Sunho Park, and Jeehoon Kang. 2023. *Modular Verification of Safe Memory Reclamation in Concurrent Separation Logic*. *Proceedings of the ACM on Programming Languages* 7, OOPSLA2 (2023), 828–856.
- Ralf Jung, Jacques-Henri Jourdan, Robbert Krebbers, and Derek Dreyer. 2018a. *RustBelt: Securing the Foundations of the Rust Programming Language*. *Proceedings of the ACM on Programming Languages* 2, POPL (2018), 66:1–66:34.
- Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Aleš Bizjak, Lars Birkedal, and Derek Dreyer. 2018b. *Iris from the ground up: A modular foundation for higher-order concurrent separation logic*. *Journal of Functional Programming* 28 (2018), e20.
- Ralf Jung, David Swasey, Filip Sieczkowski, Kasper Svendsen, Aaron Turon, Lars Birkedal, and Derek Dreyer. 2015. *Iris: monoids and invariants as an orthogonal basis for concurrent reasoning*. In *Principles of Programming Languages (POPL)*. 637–650.
- David M. Kahn and Jan Hoffmann. 2021. *Automatic amortized resource analysis with the quantum physicist’s method*. *Proceedings of the ACM on Programming Languages* 5, ICFP (2021), 1–29.
- Jan-Oliver Kaiser, Hoang-Hai Dang, Derek Dreyer, Ori Lahav, and Viktor Vafeiadis. 2017. *Strong Logic for Weak Memory: Reasoning About Release-Acquire Consistency in Iris*. In *European Conference on Object-Oriented Programming (ECOOP)*. 17:1–17:29.
- Ioannis T. Kassios and Eleftherios Kritikos. 2013. *A Discipline for Program Verification Based on Backpointers and Its Use in Observational Disjointness*. In *European Symposium on Programming (ESOP) (Lecture Notes in Computer Science, Vol. 7792)*. Springer, 149–168.
- Robbert Krebbers, Jacques-Henri Jourdan, Ralf Jung, Joseph Tassarotti, Jan-Oliver Kaiser, Amin Timany, Arthur Charguéraud, and Derek Dreyer. 2018. *MoSeL: a general, extensible modal framework for interactive proofs in separation logic*. *Proceedings of the ACM on Programming Languages* 2, ICFP (2018), 77:1–77:30.
- Rich Lander. 2015. *Announcing .NET Framework 4.6*. <https://devblogs.microsoft.com/dotnet/announcing-net-framework-4-6/>.
- Peter J. Landin. 1964. *The Mechanical Evaluation of Expressions*. *Computer Journal* 6, 4 (Jan. 1964), 308–320.
- Jonathan K. Lee and Jens Palsberg. 2010. *Featherweight X10: a core calculus for async-finish parallelism*. In *Principles and Practice of Parallel Programming (PPoPP)*. 25–36.
- Xavier Leroy. 2024. *The CompCert C compiler*. <http://compcert.org/>.
- Yi Lin, Kunshan Wang, Stephen M. Blackburn, Antony L. Hosking, and Michael Norrish. 2015. *Stop and go: understanding yieldpoint behavior*. In *Symposium on Memory Management (ISMM)*. 70–80.
- Daniel Loeb. 1992. *Sets with a negative number of elements*. *Advances in Mathematics* 91, 1 (1992), 64–74.
- Anton Lorenzen, Daan Leijen, and Wouter Swierstra. 2023. *FP²: Fully in-Place Functional Programming*. *Proceedings of the ACM on Programming Languages* 7, ICFP (Aug. 2023), 275–304.
- Anil Madhavapeddy and Yaron Minsky. 2022. *Real World OCaml: Functional programming for the masses* (2 ed.). Cambridge University Press.
- Jean-Marie Madiot and François Pottier. 2022. *A Separation Logic for Heap Space under Garbage Collection*. *Proceedings of the ACM on Programming Languages* 6, POPL (Jan. 2022), 718–747.
- Simon Marlow, Tim Harris, Roshan P. James, and Simon Peyton Jones. 2008. *Parallel generational-copying garbage collection with a block-structured heap*. In *Proceedings of the 7th International Symposium on Memory Management (Tucson, AZ, USA) (ISMM)*. Association for Computing Machinery, New York, NY, USA, 11–20. <https://doi.org/10.1145/1375634.1375637>
- Paul McKenney, Michael Wong, Maged M. Michael, Andrew Hunter, Daisy Hollman, JF Bastien, Hans Boehm, David Goldblatt, Frank Birkbacher, Erik Rigtorp, Tomasz Kamiński, Olivier Giroux, David Vernet, and Timur Doumler. 2023. *Read-Copy Update (RCU)*. P2545R4 <https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2023/p2545r4.pdf>.
- Paul E. McKenney. 2004. *Exploiting deferred destruction: an analysis of read-copy-update techniques in operating system kernels*. Ph.D. Dissertation. Oregon Health & Science University.
- Roland Meyer and Sebastian Wolff. 2019. *Decoupling lock-free data structures from memory reclamation for static analysis*. *Proc. ACM Program. Lang.* 3, POPL, Article 58 (jan 2019), 31 pages.
- Maged M. Michael. 2002. *High performance dynamic lock-free hash tables and list-based sets*. In *Proceedings of the Fourteenth Annual ACM Symposium on Parallel Algorithms and Architectures (Winnipeg, Manitoba, Canada) (SPAA)*. Association for Computing Machinery, New York, NY, USA, 73–82. <https://doi.org/10.1145/564870.564881>

- Maged M. Michael. 2004. *Hazard Pointers: Safe Memory Reclamation for Lock-Free Objects*. *IEEE Transactions on Parallel and Distributed Systems* 15, 6 (2004), 491–504.
- Maged M. Michael, Michael Wong, Paul McKenney, Andrew Hunter, Daisy Hollman, JF Bastien, Hans Boehm, David Goldblatt, Frank Birkbacher, and Mathias Stearn. 2023. Hazard Pointers for C++26. P2530R3 <https://www.open-std.org/jtc1/sc22/wg21/docs/papers/2023/p2530r3.pdf>.
- Microsoft. 2024. *Documentation of the GC class of the .NET 8.0 framework*.
- Alexandre Moine. 2024a. *Formal Verification of Heap Space Bounds under Garbage Collection*. Ph. D. Dissertation. Université Paris Cité.
- Alexandre Moine. 2024b. *Will it Fit? Verifying Heap Space Bounds for Concurrent Programs under Garbage Collection with Separation Logic (Artifact)*.
- Alexandre Moine, Arthur Charguéraud, and François Pottier. 2023. *A High-Level Separation Logic for Heap Space under Garbage Collection*. *Proceedings of the ACM on Programming Languages* 7, POPL (Jan. 2023), 718–747.
- J. Gregory Morrisett, Matthias Felleisen, and Robert Harper. 1995. *Abstract Models of Memory Management*. In *Functional Programming Languages and Computer Architecture (FPCA)*. 66–77.
- Ike Mulder, Robbert Krebbers, and Herman Geuvers. 2022. *Diaframe: automated verification of fine-grained concurrent programs in Iris*. In *Programming Language Design and Implementation (PLDI)*. 809–824.
- Glen Mével, Jacques-Henri Jourdan, and François Pottier. 2019. *Time credits and time receipts in Iris*. In *European Symposium on Programming (ESOP) (Lecture Notes in Computer Science, Vol. 11423)*. Springer, 1–27.
- Yue Niu and Jan Hoffmann. 2018. *Automatic Space Bound Analysis for Functional Programs with Garbage Collection*. In *Logic for Programming Artificial Intelligence and Reasoning (LPAR) (EPIc Series in Computing, Vol. 57)*. 543–563.
- Peter W. O’Hearn. 2019. *Separation logic*. *Commun. ACM* 62, 2 (2019), 86–95.
- Zoe Paraskevopoulou and Andrew W. Appel. 2019. *Closure conversion is safe for space*. *Proceedings of the ACM on Programming Languages* 3, ICFP (2019), 83:1–83:29.
- Matthew J. Parkinson, Richard Bornat, and Peter W. O’Hearn. 2007. *Modular verification of a non-blocking stack*. In *Principles of Programming Languages (POPL)*. 297–302.
- Azalea Raad, Josh Berdine, Hoang-Hai Dang, Derek Dreyer, Peter W. O’Hearn, and Jules Villard. 2020. *Local Reasoning About the Presence of Bugs: Incorrectness Separation Logic*. In *Computer Aided Verification (CAV) (Lecture Notes in Computer Science, Vol. 12225)*. Springer, 225–252.
- John C. Reynolds. 1975. *User-defined Types and Procedural Data Structures as Complementary Approaches to Data Abstraction*. Technical Report 1278. Carnegie Mellon University.
- John C. Reynolds. 2002. *Separation Logic: A Logic for Shared Mutable Data Structures*. In *Logic in Computer Science (LICS)*. 55–74.
- K. C. Sivaramakrishnan, Stephen Dolan, Leo White, Sadiq Jaffer, Tom Kelly, Anmol Sahoo, Sudha Parimala, Atul Dhiman, and Anil Madhavapeddy. 2020. *Retrofitting Parallelism onto OCaml*. *Proceedings of the ACM on Programming Languages* 4, ICFP (Aug. 2020), 113:1–113:30.
- Simon Spies, Lennard Gäher, Joseph Tassarotti, Ralf Jung, Robbert Krebbers, Lars Birkedal, and Derek Dreyer. 2022. *Later credits: resourceful reasoning for the later modality*. *Proceedings of the ACM on Programming Languages* 6, ICFP (2022), 283–311.
- Bogdan Tofan, Gerhard Schellhorn, and Wolfgang Reif. 2011. *Formal Verification of a Lock-Free Stack with Hazard Pointers*. In *Theoretical Aspects of Computing (ICTAC) (Lecture Notes in Computer Science, Vol. 6916)*. Springer, 239–255.
- R. Kent Treiber. 1986. *Systems programming: Coping with parallelism*.
- Leena Unnikrishnan and Scott D. Stoller. 2009. *Parametric heap usage analysis for functional programs*. In *International Symposium on Memory Management*. 139–148.
- Simon Friis Vindum and Lars Birkedal. 2021. *Contextual refinement of the Michael-Scott queue*. In *Certified Programs and Proofs (CPP)*. 76–90.
- Hassler Whitney. 1933. *Characteristic Functions and the Algebra of Logic*. *Annals of Mathematics* 34, 3 (1933), 405–414.