

Pause 'n' Play: Formalizing Asynchronous C[#]

Gavin Bierman¹, Claudio Russo¹, Geoffrey Mainland¹,
Erik Meijer², and Mads Torgersen³

¹ Microsoft Research

² Microsoft Corp. and TU Delft

³ Microsoft Corporation

{gmb, crusso, gmainlan, emejjer, madst}@microsoft.com

Abstract. Writing applications that connect to external services and yet remain responsive and resource conscious is a difficult task. With the rise of web programming this has become a common problem. The solution lies in using asynchronous operations that separate issuing a request from waiting for its completion. However, doing so in common object-oriented languages is difficult and error prone. Asynchronous operations rely on callbacks, forcing the programmer to cede control. This inversion of control-flow impedes the use of structured control constructs, the staple of sequential code. In this paper, we describe the language support for asynchronous programming in the upcoming version of C[#]. The feature enables asynchronous programming using structured control constructs. Our main contribution is a precise mathematical description that is abstract (avoiding descriptions of compiler-generated state machines) and yet sufficiently concrete to allow important implementation properties to be identified and proved correct.

1 Introduction

Mainstream programmers are increasingly adopting asynchronous programming techniques once the preserve of hard-core systems programmers. This adoption is driven by a variety of reasons: hiding the latency of the network in distributed applications; maintaining the responsiveness of single-threaded applications or simply avoiding the resource cost of creating too many threads. To facilitate this programming style, operating systems and platforms have long provided non-blocking, asynchronous alternatives to possibly blocking, synchronous operations. While these have made asynchronous programming possible they have not made it easy.

The basic principle behind these asynchronous APIs is to decompose a synchronous operation that combines issuing the operation with a blocking wait for its completion, into a non-blocking initiation of the operation, that immediately returns control, and some mechanism for describing what to do with the operation's result once it has completed. The latter is typically described by a callback—a method or function. The callback is often supplied with the initiation as an additional argument. Alternatively, the initiation can return a handle which the client can use to selectively register an asynchronous callback or (synchronously) wait for the operation's result.

Whatever the mechanism, the difficulty with using these APIs is fundamentally this: to transform a particular synchronous call-site into an asynchronous call-site requires the programmer to represent the continuation of the original site as a callback. Moreover, for this callback to resume from where the synchronous call previously returned, it must preserve all of the state pertinent to the continuation of the call. Some aspects of the state

will be explicitly available (such as the values of local variables), but other aspects may not be. A prime example is the remainder of the current call stack. For languages that do not provide support for first-class continuations, like Java and C[#], accounting for this state often requires a manual transformation to continuation-passing-style of not only the enclosing method, but also all of its callers. Once reified as an explicit continuation, the state of a computation can be saved at the initiation of an asynchronous operation and restored on its completion by supplying it with a result.

The upcoming version of C[#] (and Visual Basic) features dedicated linguistic support for asynchronous programming that removes the need for explicit callbacks. C[#] 5.0 allows certain methods to pause and then later resume their computation, without blocking, at explicitly marked code points. The basic idea is to allow a method, designated as asynchronous, to *await* the completion of some other event, not by blocking its executing thread, but by *pausing* its own execution and releasing its thread to do further work. The caller of the paused method then receives a *task* representing the method's future result and is free to proceed. Subsequent completion of the awaited event causes the paused method to resume *playing* from where it left off. Since its original thread has carried on, the resumed method is played on some available thread. Depending on run-time context, this thread may be drawn from the .NET thread pool, or it may be the same, issuing thread but at a later opportunity (e.g., the resumed method might be re-enqueued in the user interface's message loop). The events that can be awaited are typically tasks returned by nested calls to asynchronous methods. They can also, more generally, belong to any (user-defined) *awaitable* type or primitive implementations provided by the framework. The aim of these new features is to make it easy to write asynchronous methods, without having to resort to continuation-passing-style and its debilitating inversion of control flow.

As we shall see, the new asynchronous features in C[#] 5.0 are quite subtle. Current, draft Microsoft specifications [16] describe the features using precise prose and by example, giving illustrative source-to-source translations from C[#] 5.0 to ordinary C[#] 4.0. Unfortunately, the translation is intricate—it compiles source to optimized, finite state machines—so its output is both verbose and difficult to comprehend. We believe a formal, mathematical approach can yield both a precise foundation for researchers, but also a better mental model for developers and compiler writers to justify the correctness of their translation. The primary contribution of this paper is to provide such a model: a direct, operational semantics of the feature in a representative fragment of C[#] 5.0. Our semantics both capture the intent of the feature and explain its performance-driven limitations, without appealing to low-level compiler output.

The paper is structured as follows. §2 gives an informal yet precise description of the feature. §2.1 presents a realistic example, re-coding a non-trivial synchronous method to an asynchronous one, first using the feature, then adding concurrency and finally contrasting with an equivalent, hand-crafted implementation. §3 formally describes our core fragment of C[#] 5.0 and presents both a type system and an operational semantics. §4 sketches some of the correctness properties that our formalization satisfies. §5 presents some extensions to our basic setting; in §5.1 we show how to develop our operational semantics to be less abstract, and much closer to the implementation but still without having to resort to a finite state machine translation. In §5.2 we show how to extend our formalization to capture the awaitable pattern. §6 surveys briefly related work on asynchronous programming; and §7 presents conclusions and some future work.

2 Background: Async C[#] Extensions in a Nutshell

Syntactically, the additions to C[#] are surprisingly few: one new modifier **async** to mark a member as asynchronous and one new *expression*, **await** *e*, for awaiting the result—control, a value or exception—of some *awaitable* expression. An **await** expression can also be used as a statement, **await** *e*;, discarding its value. The **async** modifier can be placed on methods (excluding iterators) and some other method-like constructs (anonymous, first-class methods, i.e., lambdas and delegates). An **await** expression can only appear in an **async** method; other occurrences are static errors.

Statically, an **async** method must have a *taskable* return type of `Task< σ >`, `Task` or `void`. The *return statements* of an **async** method with return type `Task< σ >` may only return values of type σ , never `Task< σ >`! The return statements in other **async** methods may only return control (but never a value).

The argument, *e*, of an **await** *e* expression must have an *awaitable* type. The concept of awaitable type is defined by a *pattern* (of available methods). A type is awaitable when it statically supports a `GetAwaiter()` instance method that returns some *awaiter* type (possibly the same type). In turn, the awaiter type must support:

- a boolean instance property `IsCompleted` testing if the awaiter has a result now.
- a `void`-returning instance method `OnCompleted(a)`, accepting a callback of delegate type `Action`.¹ `Action a` is a *one-shot* continuation; calling `a()` resumes the awaiting method. The action should be invoked *at most once*, on completion.
- a τ -returning instance method `GetResult()` for retrieving the result of a completed awaiter. `GetResult()` should either return control, some stored value, or throw some stored exception.

If the return type of `GetResult` is τ , then expression **await** *e* is an *expression* of type τ . All these operations should be (essentially) *non-blocking*.

Crucially, the types `Task< σ >` and `Task` are awaitable, allowing **async** methods to await the tasks of *nested* **async** method calls. A caller can use a returned task just like any other task (asynchronously awaiting its result, synchronously waiting for its completion or by registering an asynchronous callback). Asynchronous methods that return `void` cannot be awaited; such methods are intended for ‘fire-and-forget’ scenarios.

Dynamically, an **async** method executes like an ordinary method until it encounters an **await** on some value. If the value’s awaiter is complete, the method continues executing using the result of the awaiter as the result of the **await** expression. If the awaiter is incomplete, the method registers its continuation as a callback on the awaiter and then suspends its execution. Execution will resume with the result of the awaiter, on some thread,² when the callback is invoked. Invoking an **async** method allocates a fresh, incomplete task, representing this invocation, immediately enters the method (on the caller’s thread) and executes it until it encounters its first **await** on an incomplete awaiter. Exiting from an **async** method, either via `return` or throwing an exception, stores the result in its task, thus *completing* it. The first suspension of an **async** method call returns its incomplete task (or `void`) to its caller. If the call exits without ever suspending, it simply returns its completed task.

¹ Defined as `delegate void Action()`.

² We are deliberately vague here; the awaiter is free to choose how the method is resumed, catering for different behaviours in, for example, single- and multi-threaded applications.

The operational semantics is deliberately designed to minimize context switches. Continuing when an awaiter is already complete ensures methods only suspend when necessary. Dually, allowing an `async` method call to begin execution on its caller's thread gives the method the opportunity to enter-and-exit quickly when possible, without imposing the cost of a context switch just to get running in the first place. This design choice morally obliges the method not to block nor even spend too much time before ceding its caller's thread (by suspension).

2.1 Example

To both illustrate and motivate the feature we present an example adapted from the Async CTP.³ Consider the following synchronous copy function, which incrementally copies an input stream to an output stream in manageable chunks.

```
public static long CopyTo(Stream src, Stream dst) {
    var buffer = new byte[0x1000]; int bytesRead; long totalRead = 0;
    while ((bytesRead = src.Read(buffer, 0, buffer.Length)) > 0) {
        dst.Write(buffer, 0, bytesRead);
        totalRead += bytesRead; }
    return totalRead; }
```

Depending on their receivers, the calls to `Stream` methods `Read` and `Write` may well be blocking IO operations. Since this method could spend much of its time blocked, one might prefer an asynchronous variant. One way to achieve this is by replacing the synchronous calls to `Read` and `Write` with their asynchronous counterparts:

```
Task<int> ReadAsync(byte[] buffer, int offset, int count);
Task WriteAsync(byte[] buffer, int offset, int count);
```

The asynchronous variants initiate an asynchronous operation and immediately return a task representing its completion. Method `ReadAsync` begins a read and immediately returns a `Task<int>` that, once completed, will record the number of bytes actually read. The method `WriteAsync` begins a write and returns a `Task` that just tracks its completion. Tasks are completed at most once with some result. The result may be a value or some exception.

Asynchronous clients can register zero or more callbacks on a task, to be executed (on some thread) once the task has completed with a result, e.g.:

```
Task<int> rdTask = src.ReadAsync(buffer, 0, buffer.Length);
rdTask.ContinueWith((Task<int> completedRdTask) => {
    int bytesRead = completedRdTask.Result; /* won't block */
    dst.WriteAsync(buffer, 0, bytesRead); });
// do some work now
```

Since `ContinueWith` is also non-blocking, the client will quickly proceed with its work. The call to property `completedRdTask.Result` in the callback is guaranteed not to block because its task (aliasing `rdTask`) must, by causality, already be completed.

Synchronous clients can also access a task's `Result` property or just `Wait()` for its completion; these calls will block until or unless the task has completed:

³ <http://msdn.microsoft.com/vstudio/async>

```

Task<int> rdTask = src.ReadAsync(buffer, 0, buffer.Length);
// do some work now
int bytesRead = rdTask.Result; /* may block */
Task wrTask = dst.WriteAsync(buffer, 0, bytesRead);
// do some more work now
wrTask.Wait(); /* may block */

```

Now let us show how C# 5.0 enables the simple implementation of an asynchronous version of CopyTo. First we mark the method as **async**, and then simply await the results of `src.ReadAsync` and `dst.WriteAsync`. Otherwise *the code remains the same*.

```

public static async Task<long> CopyToAsync(Stream src, Stream dst) {
    var buffer = new byte[0x1000]; int bytesRead; long totalRead = 0;
    while ((bytesRead = await src.ReadAsync(buffer,0,buffer.Length)) > 0) {
        await dst.WriteAsync(buffer, 0, bytesRead);
        totalRead += bytesRead; }
    return totalRead; }

```

Awaiting the task returned by `ReadAsync` pauses the method unless the read has completed; when play is resumed, the `await` expression extracts the integer value of the task. `WriteAsync` returns a non-generic `Task`; the `await` statement pauses the method unless the write has completed; when it is played again, execution proceeds from the next statement. Note that the return type of our asynchronous method is `Task<long>` even though its body **returns** a `long`. Clearly, this is no ordinary **return** statement.

Though the code is almost identical to the original, the behaviour and resource consumption is quite different. A call to `CopyTo` will repeatedly block (in the kernel) on each call to `Read` and `Write`, tying up the resources dedicated to that thread. A call to `CopyToAsync`, on the other hand, will never block; instead, each continuation of an `await` will be executed on demand, on some available thread in the .NET thread pool.⁴

`CopyToAsync`'s reads and writes are still being executed sequentially, but from different threads rather than a single one, so we would not expect to gain any performance from the asynchronous implementation. Indeed, given the additional scheduling and compilation overheads, the synchronous `CopyTo` is likely to execute faster.

However, we are now in a good position to *overlap* the last write with the next read, leading to this potentially faster, concurrent implementation:

```

public static async Task<long> CopyToConcurrent(Stream src, Stream dst) {
    var buffer = new byte[0x1000]; var oldbuffer = new byte[0x1000];
    int bytesRead; long totalRead = 0; Task lastwrite = null;
    while ((bytesRead = await src.ReadAsync(buffer,0,buffer.Length)) > 0) {
        if (lastwrite != null) await lastwrite; // wait later
        lastwrite = dst.WriteAsync(buffer, 0, bytesRead); // issue now
        totalRead += bytesRead;
        { var tmp = buffer; buffer = oldbuffer; oldbuffer = tmp; }; }
    if (lastwrite != null) await lastwrite;
    return totalRead; }

```

In order to achieve this, we exploit the ability to separate the initiation of a task from the act of awaiting its completion, so that we can issue the next read *during* the last

⁴ If invoked from a user interface thread, each continuation will be scheduled on that thread's event queue.

write. The example illustrates the crucial advantage of allowing asynchronous methods to return incomplete, concurrently executing tasks, not just completed results. Though we emphasize concurrency, the reads and writes could also be executing in parallel, depending on the underlying streams.

In comparison, TaskJava [8] also provides constructs to avoid inversion of control while programming with asynchronous APIs. Like C[#] 5.0, TaskJava compiles straight-line code to a state machine, but its syntax is slightly more heavyweight and requires explicit calls to an event scheduler. TaskJava does not make a distinction between invoking and awaiting an asynchronous operation, so although it presents a pleasant programming model for those who wish to invoke asynchronous APIs sequentially, it is of no use to a programmer who must write code that executes multiple operations concurrently, like CopyToConcurrent. C[#] 5.0 does not sacrifice concurrency for convenience.

To appreciate the concision of CopyToAsync, let us contrast it with a representative hand-crafted version of CopyToAsync, CopyToManual, written in C[#] 4.0. Actually, this code is very close to the decompiled code emitted by the C[#] 5.0 compiler and described in the feature documentation. As mentioned earlier, the aim of our work is to eliminate the need to understand this compilation strategy.

```
public static Task<long> CopyToManual(Stream src, Stream dst) {
    var tcs = new TaskCompletionSource<long>(); // tcs.Task new & incomplete
    var state = 0; TaskAwaiter<int> readAwaiter; TaskAwaiter writeAwaiter;
    byte[] buffer = null; int bytesRead = 0; long totalRead = 0;
    Action act = null; act = () => {
        while (true) switch (state++) {
            case 0: buffer = new byte[0x1000]; totalRead = 0; continue;
            case 1: readAwaiter=src.ReadAsync(buffer,0,buffer.Length).GetAwaiter();
                    if (readAwaiter.IsCompleted) continue; // goto post-read
                    else { readAwaiter.OnCompleted(act); return;} // suspend at 2
            case 2: if ((bytesRead = readAwaiter.GetResult()) > 0) {
                    writeAwaiter=dst.WriteAsync(buffer,0,bytesRead).GetAwaiter();
                    if (writeAwaiter.IsCompleted) continue; // goto post-write
                    else { writeAwaiter.OnCompleted(act); return;} // suspend at 3
                } else { state = 4; continue;} // goto post-while
            case 3: writeAwaiter.GetResult();
                    totalRead += bytesRead;
                    state = 1; continue; // goto pre-while
            case 4: tcs.SetResult(totalRead); // complete tcs.Task & "return"
                    return; // exit machine
        }
    }; // end of act delegate
    act(); // start the machine on this thread
    return tcs.Task; } // on first suspend or exit from machine
```

Without going into too many details, notice how the control flow has been obscured by encoding the continuation of each await as states (here 2 & 3) of a finite state machine. The original locals, arguments and internal state of the method are (implicitly) allocated on the heap. (Note that C[#] lambdas such as act close over L-values, not R-values, automatically placing them on the heap; updates to those locations persist across lambda invocations.) State 0 is the initial state that sets up locals; state 1 is the **while** loop header, states 2 and 3 are the continuations of the **await** statements; state 4 is the final state and the continuation of the original **while** statement. State 4 exits the machine,

setting the result in the task held by shared variable `task` (completing the task). The finite state machine only suspends (by calling `return` without completing the task) in states 2 and 3 (just after an `await`); the other states encode internal control flow points.

3 Formalization: Featherweight C[#] 5.0

In the rest of the paper we study the essence of the new asynchronous features of C[#]. To do so we take a formal, mathematical approach and define an idealized fragment, Featherweight C[#] 5.0, or FC₅[#] for short. Whilst FC₅[#] programs remain syntactically valid C[#], it is a heavily restricted fragment—any language feature that is not needed to demonstrate the essence of the asynchronous features has been removed, and the resulting fragment has been further refactored to allow for a more succinct presentation.

As the new asynchronous features predominantly affect the control flow of C[#] programs, most of our attention is on the operational semantics. In contrast, other minimal fragments such as, for example, Featherweight Java [13], and Classic Java [10], are primarily concerned with typing issues. The asynchronous features in C[#] 5.0 have almost no impact on the type system. Consequently we have stripped the type system of FC₅[#] to the core: we have only simple non-generic classes, some value types and no subtyping at all! However, we emphasize that our formalization exposes enough of the inner workings of C[#] 5.0 to allow the reader to reason about how aspects of the language, like the split between invocation and awaiting, affect the concurrent execution of multiple threads of control. There is a danger in cutting out too much of a language during formalization. For example, Fischer et al. [8] give a semantics for CJT, a simplified version of TaskJava, that avoids a heap at the expense of any ability to model communication between tasks, including a spawned thread signaling completion to its parent. In contrast, our semantics models what we claim to be the essential features of task-based programming: concurrent execution of multiple, *communicating* tasks that are invoked by the same thread of control.

FC₅[#] programs and types:

$p ::= \overline{cd} mb$	Program
$cd ::= \text{public class } C \{ \overline{fd} \overline{md} \}$	Class declaration
$fd ::= \text{public } \sigma f;$	Field declaration
$md ::= \text{public } \phi m(\overline{\sigma} \overline{x}) mb \mid \text{async public } \psi m(\overline{\sigma} \overline{x}) mb$	Method declaration
$mb ::= \{ \overline{\sigma} \overline{x}; \overline{s} \}$	Method body
$\phi ::= \sigma \mid \text{void}$	Return type
$\sigma, \tau ::= \gamma \mid \rho$	Type
$\gamma ::= \text{bool} \mid \text{int}$	Value type
$\rho ::= C \mid \text{Task} \langle \sigma \rangle$	Reference type
$\psi ::= \text{Task} \langle \sigma \rangle$	Taskable return type

Our formalization makes heavy use of the Featherweight Java [13] overbar notation, i.e., we write \overline{x} for a possibly empty sequence x_1, \dots, x_n . We write the empty sequence as ϵ . We abbreviate operations on pairs of sequences, writing for example $\overline{\sigma} \overline{x}$ for the sequence $\sigma_1 x_1, \dots, \sigma_n x_n$, similarly $\overline{\sigma} \overline{x}$; for the sequence of variable declarations $\sigma_1 x_1; \dots, \sigma_n x_n$; and finally $f(\overline{\sigma})$ for the sequence $f(\sigma_1), \dots, f(\sigma_n)$.

A FC_5^\sharp program consists of a collection of class declarations and a single method body (C^\sharp 's main method). A FC_5^\sharp class declaration `public class C { \overline{fd} \overline{md} }` introduces a class C . We repeat that FC_5^\sharp does not support any form of subtyping so class declarations do not specify a superclass. This is a valid declaration in full C^\sharp as all classes inherit from `object` by default, but we do not even support the `object` class in FC_5^\sharp . Subtyping and inheritance are orthogonal to the new features in C^\sharp 5.0 and so we removed them from our fragment to concentrate solely on the support for asynchronous programming.⁵ The class C has fields \overline{f} with types $\overline{\sigma}$ and a collection of methods \overline{md} .

Method declarations can be either synchronous or asynchronous. A synchronous method `public ϕ m($\overline{\sigma}$ \overline{x}) mb` declares a public method m with return type ϕ , formal parameters \overline{x} of type $\overline{\sigma}$ and a method body mb . Methods may be `void`-returning, i.e., they return control not a value. Method bodies are constrained to be of a particular form: $\overline{\sigma}$ \overline{x} ; \overline{s} , i.e., they must declare all their local variables \overline{x} at the start of the method, and then contain a sequence of statements \overline{s} .

An asynchronous method is marked with the `async` keyword and is syntactically the same as a synchronous method, although it is type checked differently. The return type of an asynchronous method must be of a so-called taskable type. For FC_5^\sharp this means it must be of the form `Task< σ >`. C^\sharp 5.0 also classifies the non-generic class `Task` and `void` as taskable return types as discussed in §2.

FC_5^\sharp types are a simple subset of the C^\sharp types. Note that FC_5^\sharp does not support user-defined generics; again these are orthogonal to asynchrony and have been removed. For simplicity, we assume that `Task< σ >` is the only generic type.

FC_5^\sharp expressions and statements:

$e ::=$	Expressions
c	Constant (boolean b , integer i or <code>null</code>)
$x \oplus y$	Built-in operator
x	Variable
$x.f$	Field access
$x.m(\overline{y})$	Method invocation
<code>new C()</code>	Object creation
<code>await x</code>	Await expression
<code>Task.AsyncIO<γ>()</code>	Async primitive
$s, t ::=$	Statement
$x=e;$	Assignment statement
<code>if (x) {\overline{s}} else {\overline{t}}</code>	Conditional statement
<code>while (x) {\overline{s}}</code>	Iteration
$x.f = y;$	Field assignment statement
$x.m(\overline{y});$	Method invocation statement
<code>return;</code>	Return statement
<code>return x;</code>	Return value statement

FC_5^\sharp expressions are restricted to a form that we call *statement normal form (SNF)*. SNF forces all subexpressions to be named; i.e., all subexpressions are simply variables. SNF

⁵ The extensions to support single inheritance, overloading, constructor methods and many of the complications of the full C^\sharp type system have appeared elsewhere [2,3].

is the natural analogue to the A-normal-form popular in functional languages [9]. This regularity makes the presentation of the operational semantics (and the type system) much simpler at no cost to expressivity.

$FC_5^\#$ expressions include constants, c , which can be an integer, \underline{i} , a boolean, \underline{b} , or the literal **null**. We assume a number of built-in primitive operators, such as $=$, $<$, $>$ and so on. In the grammar we write $x \oplus y$, where \oplus denotes an instance of one of these operators. We do not specify operators further as their meaning is clear. We assume that x, y, z range over variable names, f ranges over field names and m ranges over method names. We assume that the set of variables includes the special variable **this**, which cannot be used as a formal parameter of a method declaration or declared as a local.

$FC_5^\#$ supports awaitable expressions, written **await** x . To get things off the ground we assume an in-built asynchronous method `Task.AsyncIO< γ >()` that spawns a thread and immediately returns a task. The thread may complete the task, depending on the scheduler, with some result of value type γ .

$FC_5^\#$ statements are standard. In what follows we assume that $FC_5^\#$ programs are well-formed, e.g., the identifier **this** does not appear as a formal parameter, all control paths in a method body contain a **return** statement, etc. These conditions can be easily formalized and are identical to restrictions on earlier fragments of C[#] but we suppress the details for lack of space. The only new well-formedness condition is that **await** expressions are only allowed to appear inside *asynchronous* method declarations.

We assume that a correct program induces a number of utility functions that we will use in the typing rules. First, we assume the partial function $f\text{type}$, which is a map from a type and a field name to a type. Thus $f\text{type}(\sigma, f)$ returns the type of field f in type σ . Second, we assume a partial function $m\text{type}$ that is a map from a type and a method name to a *type signature*. For example, we write $m\text{type}(\mathcal{C}, m) = (\bar{\tau}) \rightarrow \phi$ when class \mathcal{C} contains a method m with formal parameters of type $\bar{\tau}$ and return type ϕ .

The type system for full C[#] is actually a bidirectional type system [18] consisting of two typing relations: a type conversion relation and a type synthesis relation [3], along with a number of conversion (subtyping) judgements. However, the extreme parsimony of $FC_5^\#$ means that we have no subtyping judgements, and we need only a single judgement for type checking an expression. The judgement is written $\Gamma \vdash e : \sigma$ where Γ is a function from variables to types. We extend the overbar notation and write $\Gamma \vdash \bar{e} : \bar{\sigma}$ to mean the judgements $\Gamma \vdash e_1 : \sigma_1, \dots, \Gamma \vdash e_n : \sigma_n$.

$FC_5^\#$ expression type checking:

$[C\text{-Bool}] \frac{}{\Gamma \vdash \underline{b} : \text{bool}}$	$[C\text{-Int}] \frac{}{\Gamma \vdash \underline{i} : \text{int}}$	$[C\text{-Null}] \frac{}{\Gamma \vdash \text{null} : \rho}$
$[C\text{-Op}] \frac{\Gamma \vdash x : \sigma_0 \quad \Gamma \vdash y : \sigma_1 \quad \oplus : \sigma_0 \times \sigma_1 \rightarrow \tau}{\Gamma \vdash x \oplus y : \tau}$		
$[C\text{-New}] \frac{}{\Gamma \vdash \text{new } \mathcal{C}() : \mathcal{C}}$	$[C\text{-Var}] \frac{}{\Gamma, x : \tau \vdash x : \tau}$	$[C\text{-Field}] \frac{f\text{type}(\sigma, f) = \tau}{\Gamma, x : \sigma \vdash x.f : \tau}$
$[C\text{-MethInv}] \frac{m\text{type}(\sigma_0, m) = (\bar{\tau}) \rightarrow \sigma_1 \quad \Gamma, x : \sigma_0 \vdash \bar{y} : \bar{\tau}}{\Gamma, x : \sigma_0 \vdash x.m(\bar{y}) : \sigma_1}$		
$[C\text{-Await}] \frac{}{\Gamma, x : \text{Task}\langle\sigma\rangle \vdash \text{await } x : \sigma}$	$[C\text{-IO}] \frac{}{\Gamma \vdash \text{Task.AsyncIO}\langle\gamma\rangle() : \text{Task}\langle\gamma\rangle}$	

Most of the type checking rules are quite standard, but there are two new rules for dealing with asynchronous methods. Rule [C-Await] states that if x is of type $\text{Task}\langle\sigma\rangle$ then awaiting x results in a value of type σ . As discussed earlier, Rule [C-IO] states that $\text{Task.AsyncIO}\langle\gamma\rangle()$ returns a value of type $\text{Task}\langle\gamma\rangle$.

FC_5^\sharp statement type checking:

$$\begin{array}{c}
 \text{[C-Asn]} \frac{\Gamma, x: \sigma \vdash e: \sigma}{\Gamma, x: \sigma \vdash x = e;: \phi} \quad \text{[C-Cond]} \frac{\Gamma, x: \text{bool} \vdash \bar{s}: \phi \quad \Gamma, x: \text{bool} \vdash \bar{t}: \phi}{\Gamma, x: \text{bool} \vdash \text{if } (x) \{ \bar{s} \} \text{ else } \{ \bar{t} \}: \phi} \\
 \\
 \text{[C-While]} \frac{\Gamma, x: \text{bool} \vdash \bar{s}: \phi}{\Gamma, x: \text{bool} \vdash \text{while } (x) \{ \bar{s} \}: \phi} \\
 \\
 \text{[C-FAsn]} \frac{\text{ftype}(\sigma_0, f) = \sigma_1 \quad \Gamma, x: \sigma_0 \vdash y: \sigma_1}{\Gamma, x: \sigma_0 \vdash x.f=y;: \phi} \\
 \\
 \text{[C-MInv]} \frac{\text{mtype}(\sigma_0, m) = (\bar{\tau}) \rightarrow \text{void} \quad \Gamma, x: \sigma_0 \vdash \bar{y}: \bar{\tau}}{\Gamma, x: \sigma_0 \vdash x.m(\bar{y});: \phi} \\
 \\
 \text{[C-Return]} \frac{}{\Gamma \vdash \text{return};: \text{void}} \quad \text{[C-ReturnExp]} \frac{}{\Gamma, x: \sigma \vdash \text{return } x;: \sigma}
 \end{array}$$

As for full C^\sharp , we give type checking rules for statements; the judgement is written $\Gamma \vdash s: \phi$. The key rules are [C-ReturnExp] that asserts that the statement $\text{return } x$; is of return type σ if x is of type σ and [C-Return] that asserts that the statement $\text{return};$ is of return type void . In other words, the role of the type ϕ in the judgement $\Gamma \vdash s: \phi$ is to check any return statement. Again, we adopt an overbar notation and write $\Gamma \vdash \bar{s}: \phi$ to denote the judgements $\Gamma \vdash s_1: \phi, \dots, \Gamma \vdash s_n: \phi$.

FC_5^\sharp method and class typing (rule for programs omitted due to space):

$$\begin{array}{c}
 \text{[Class-OK]} \frac{C \vdash \overline{md} \text{ ok}}{\vdash \text{public class } C \{ \overline{fd} \overline{md} \} \text{ ok}} \\
 \\
 \text{[Meth-OK]} \frac{\bar{x}: \bar{\sigma}, \bar{y}: \bar{\tau}, \text{this}: C \vdash \bar{s}: \phi \quad \forall e \in \bar{s}, e \neq \text{await } _}{C \vdash \text{public } \phi \text{ m}(\bar{\sigma} \bar{x}) \{ \bar{\tau} \bar{y}; \bar{s} \} \text{ ok}} \\
 \\
 \text{[AsyncMeth-OK]} \frac{\bar{x}: \bar{\sigma}, \bar{y}: \bar{\tau}, \text{this}: C \vdash \bar{s}: \sigma_0}{C \vdash \text{async public Task}\langle\sigma_0\rangle \text{ m}(\bar{\sigma} \bar{x}) \{ \bar{\tau} \bar{y}; \bar{s} \} \text{ ok}}
 \end{array}$$

Rule [Class-OK] asserts that a class declaration is well-typed provided that all its method declarations are well-typed. Rule [Meth-OK] asserts that the (synchronous) method declaration $\text{public } \phi \text{ m}(\bar{\sigma} \bar{x}) \{ \bar{\tau} \bar{y}; \bar{s} \}$ is well-typed in class C provided that the statements \bar{s} can be typed at return type ϕ in the context $\bar{x}: \bar{\sigma}, \bar{y}: \bar{\tau}, \text{this}: C$. Moreover, \bar{s} cannot contain `await`. Rule [AsyncMeth-OK] asserts that the asynchronous method $\text{async public Task}\langle\sigma_0\rangle \text{ m}(\bar{\sigma} \bar{x}) \{ \bar{\tau} \bar{y}; \bar{s} \}$ is well-typed if the statements \bar{s} can be typed at return type σ_0 (not $\text{Task}\langle\sigma_0\rangle$) in the context $\bar{x}: \bar{\sigma}, \bar{y}: \bar{\tau}, \text{this}: C$.

We assume two methods on the $\text{Task}\langle\sigma\rangle$ type, `Result` and `GetResult`, which both take no argument and return a value of type σ , i.e., $\text{mtype}(\text{Task}\langle\sigma\rangle, _) = () \rightarrow \sigma$. Both these methods return the result of a complete task object, but will differ operationally on an incomplete task. In C^\sharp , `Result` is actually a method-like *property*.

3.1 Operational Semantics

The key contribution of this paper is a precise description of the operational behaviour of the new asynchronous features in C[#]. The syntactic restrictions of FC₅[#] mean that the operational semantics can be given as single-step transition rules between configurations.

A *heap*, H , is a partial map from an object identifier (ranged over by o) to a heap object. A *heap object* can be one of three forms: $\langle C, FM \rangle$ denotes a non-task object of class C with a *field map*, FM , which is a partial map from fields f to values. A task heap object is either of the form $\langle \text{Task}\langle\sigma\rangle, \text{running}(\overline{F}) \rangle$ or $\langle \text{Task}\langle\sigma\rangle, \text{done}(v) \rangle$. We explain these forms later. A *value*, v is either a constant, c , or an object identifier (the *address* of an object in the heap).

A *frame*, F , is written $\langle L, \overline{s} \rangle^\ell$ and consists of a locals stack, L , and a sequence of statements, \overline{s} , along with a frame label, ℓ . A *locals stack* is a partial map from local variables to values. A *frame label*, ℓ , is either s to denote a synchronous frame, or $a(o)$ for an asynchronous frame whose associated Task is stored at heap address o . A *frame stack*, FS , is essentially a list of frames. An empty frame stack is written ϵ , and we write $F \circ FS$ to denote a frame stack whose head is a frame F and tail is the frame stack FS . A process, P , is a collection of frame stacks, written $\{FS_1, \dots, FS_n\}$.

We factor the transition rules into three relations describing the small step evaluation of frames (method bodies), frame stacks (corresponding to individual threads) and collections of frame stacks (corresponding to a process, i.e., a pool of threads mutating a shared heap). Thus, a frame configuration is written $H \triangleright F$ and the transition relation between frame configurations is written $H_1 \triangleright F_1 \rightarrow H_2 \triangleright F_2$. A frame stack configuration is written $H \triangleright FS$ and the transition relation between frame stack configurations is written $H_1 \triangleright FS_1 \rightarrow H_2 \triangleright FS_2$. Finally, a process configuration is written $H \triangleright P$ and the transition relation between process configurations is written $H_1 \triangleright P_1 \rightsquigarrow H_2 \triangleright P_2$.

Simple frame transition rules:

$H \triangleright \langle L, x=c; \overline{s} \rangle^\ell \rightarrow H \triangleright \langle L[x \mapsto c], \overline{s} \rangle^\ell$	[E-Constant]
$H \triangleright \langle L, x=y; \overline{s} \rangle^\ell \rightarrow H \triangleright \langle L[x \mapsto L(y)], \overline{s} \rangle^\ell$	[E-Var]
$H \triangleright \langle L, x=y \oplus z; \overline{s} \rangle^\ell \rightarrow H \triangleright \langle L[x \mapsto L(y) \oplus L(z)], \overline{s} \rangle^\ell$	[E-Op]
$H \triangleright \langle L, x=y.f; \overline{s} \rangle^\ell \rightarrow H \triangleright \langle L[x \mapsto FM(f)], \overline{s} \rangle^\ell$ where $H(L(y)) = \langle \rho, FM \rangle$	[E-Field]
$H \triangleright \langle L, \text{if } (x) \{ \overline{s} \} \text{ else } \{ \overline{t} \} \overline{u} \rangle^\ell \rightarrow H \triangleright \langle L, \overline{s} \overline{u} \rangle^\ell$ where $L(x) = \text{true}$	[E-CondEq]
$H \triangleright \langle L, \text{if } (x) \{ \overline{s} \} \text{ else } \{ \overline{t} \} \overline{u} \rangle^\ell \rightarrow H \triangleright \langle L, \overline{t} \overline{u} \rangle^\ell$ where $L(x) = \text{false}$	
$H \triangleright \langle L, \text{while } (x) \{ \overline{s} \} \overline{t} \rangle^\ell \rightarrow H \triangleright \langle L, \overline{s} \text{while } (x) \{ \overline{s} \} \overline{t} \rangle^\ell$ where $L(x) = \text{true}$	[E-While]
$H \triangleright \langle L, \text{while } (x) \{ \overline{s} \} \overline{t} \rangle^\ell \rightarrow H \triangleright \langle L, \overline{t} \rangle^\ell$ where $L(x) = \text{false}$	
$H_0 \triangleright \langle L, x.f=y; \overline{s} \rangle^\ell \rightarrow H_1 \triangleright \langle L, \overline{s} \rangle^\ell$ where $L(x) = o$, $H_0(o) = \langle \sigma, FM \rangle$ and $H_1 = H_0[o \mapsto \langle \sigma, FM[f \mapsto L(y)] \rangle]$	[E-Asn]
$H_0 \triangleright \langle L, x=\text{new } C(); \overline{s} \rangle^\ell \rightarrow H_1 \triangleright \langle L[x \mapsto o], \overline{s} \rangle^\ell$ where $\text{fields}(C) = \overline{f}$, $o \notin \text{dom}(H_0)$ and $H_1 = H_0[o \mapsto \langle C, \overline{f} \mapsto \text{default}(\overline{f}) \rangle]$	[E-New]

In these transition rules the frames are labeled with meta-variable ℓ : they apply for both synchronous and asynchronous frames, factoring common semantics. Our transition rules $H_0 \triangleright \langle L_0, \bar{s} \rangle^\ell \rightarrow H_1 \triangleright \langle L_1, \bar{t} \rangle^\ell$ always preserve labels, i.e., a synchronous frame transitions to another synchronous frame, and an asynchronous frame transitions to an asynchronous frame with the *same* task. In rule [E-New] we use an auxiliary function, *default* that returns a default constant for a given type. This notion is taken from full C^\sharp [12, §5.2] but for FC_5^\sharp it simply maps type `int` to the value 0, type `bool` to the value `false` and all other types to the `null` literal. These simple transition rules are quite standard and for space reasons we do not elaborate on them further.

Next we consider the evaluation of a synchronous method call and returning from a synchronous method. For a `return` (though not a call) the label on the frame is important; as we shall see, the `return` rule is different for asynchronous frames.

Synchronous method call/return transition rules:

$$H_0 \triangleright F_0 \circ FS \rightarrow H_1 \triangleright F_1 \circ FS \quad \text{if } H_0 \triangleright F_0 \rightarrow H_1 \triangleright F_1 \quad [\text{E-Frame}]$$

$$H \triangleright \langle L_0, y_0=y_1 . m(\bar{z}) ; \bar{s} \rangle^\ell \circ FS \rightarrow H \triangleright \langle L_1, \bar{t} \rangle^s \circ \langle L_0, \bar{s} \rangle_{y_0}^\ell \circ FS \quad [\text{E-Method-Exp}]$$

where $H(L_0(y_1)) = \langle \rho, FM \rangle$, $mbody(\rho, m) = mb: (\bar{\sigma} \bar{x}) \rightarrow^s \sigma_1$, $mb = \bar{\tau} \bar{y}$; \bar{t} and $L_1 = [\bar{x} \mapsto L_0(\bar{z}), \bar{y} \mapsto \text{default}(\bar{\tau}), \text{this} \mapsto L_0(y_1)]$

$$H \triangleright \langle L_0, x . m(\bar{y}) ; \bar{s} \rangle^\ell \circ FS \rightarrow H \triangleright \langle L_1, \bar{t} \rangle^s \circ \langle L_0, \bar{s} \rangle^\ell \circ FS \quad [\text{E-Method-Stmt}]$$

where $H(L_0(x)) = \langle \rho, FM \rangle$, $mbody(\rho, m) = mb: (\bar{\sigma} \bar{x}) \rightarrow^s \text{void}$, $mb = \bar{\tau} \bar{z}$; \bar{t} and $L_1 = [\bar{x} \mapsto L_0(\bar{z}), \bar{z} \mapsto \text{default}(\bar{\tau}), \text{this} \mapsto L_0(x)]$

$$H \triangleright \langle L_0, \text{return } y ; \bar{s} \rangle^s \circ \langle L_1, \bar{t} \rangle_x^\ell \circ FS \rightarrow H \triangleright \langle L_1[x \mapsto L_0(y)], \bar{t} \rangle^\ell \circ FS \quad [\text{E-Return-Val}]$$

$$H \triangleright \langle L_0, \text{return} ; \bar{s} \rangle^s \circ \langle L_1, \bar{t} \rangle^\ell \circ FS \rightarrow H \triangleright \langle L_1, \bar{t} \rangle^\ell \circ FS \quad [\text{E-Return}]$$

These transition rules are also quite standard. Rule [E-Frame] transitions the top-most, *active* frame of a frame stack. Rule [E-Method-Exp] transitions a method invocation. It first looks up in the heap the runtime type of the receiver. We make use of another auxiliary function induced by correct program: *mbody* is a map from a type and a method name to a method body and an annotated type signature. For example, we write $mbody(C, m) = mb: (\bar{\sigma} \bar{x}) \rightarrow^s \phi$, when the method m in class C is a synchronous method, with formal parameters $\bar{\sigma} \bar{x}$, return type ϕ , and method body mb .

Rule [E-Method-Exp] applies when the receiver object supports method m and m is a synchronous method. In this case, we push a new synchronous frame (labeled s) on to the frame stack to execute the method body. Notice that we annotate the caller frame with the identifier that is waiting for the return value (this will be used in rule [E-Return-Val]). Rule [E-Method-Stmt] is similar except that m is a `void`-returning method, returning control. Note that the semantics of synchronous calls are the same whether issued from a synchronous or asynchronous frame (ℓ can be any label).

Rule [E-Return-Val] shows how a synchronous method returns a value to its caller. The caller frame, $\langle L_1, \bar{t} \rangle_x^\ell$, is waiting for a value for local identifier x . The active synchronous frame is popped and the caller frame becomes active and assigns the return value to x . Rule [E-Return] is similar except that no value is returned and the caller frame is not annotated with an identifier: the caller only expects control, not a value.

Asynchronous method call/return transition rules:

$H_0 \triangleright \langle L_0, y_0=y_1.m(\bar{z}); \bar{s} \rangle^\ell \circ FS$	[E-Async-Method]
$\rightarrow H_1 \triangleright \langle L_1, \bar{t} \rangle^{a(o)} \circ \langle L_0[y_0 \mapsto o], \bar{s} \rangle^\ell \circ FS$	
where $H_0(L_0(y_1)) = \langle \rho, FM \rangle$, $mbody(\rho, m) = mb: (\bar{\sigma} \bar{x}) \rightarrow^a \psi$, and $mb = \bar{\tau} \bar{y}; \bar{t}$ $o \notin dom(H_0), H_1 = H_0[o \mapsto \langle \psi, running(\epsilon) \rangle]$ $L_1 = [\bar{x} \mapsto L_0(\bar{z}), \bar{y} \mapsto default(\bar{\tau}), \mathbf{this} \mapsto L_0(y_1)]$	
$H_0 \triangleright \{ \langle L, \mathbf{return} y; \bar{s} \rangle^{a(o)} \circ FS \} \cup P$	[E-Async-Return]
$\sim H_1 \triangleright \{ FS \} \cup resume(\bar{F}) \cup P$	
where $H_0(o) = \langle \mathbf{Task}\langle \sigma \rangle, running(\bar{F}) \rangle$ and $H_1 = H_0[o \mapsto \langle \mathbf{Task}\langle \sigma \rangle, done(L(y)) \rangle]$	
$H \triangleright \langle L, x=\mathbf{await} y; \bar{s} \rangle^{a(o)} \circ FS \rightarrow H \triangleright \langle L[x \mapsto v], \bar{s} \rangle^{a(o)} \circ FS$	[E-Await-Continue]
where $H(L(y)) = \langle \mathbf{Task}\langle \sigma \rangle, done(v) \rangle$	
$H_0 \triangleright \langle L, x=\mathbf{await} y; \bar{s} \rangle^{a(o)} \circ FS \rightarrow H_1 \triangleright FS$	[E-Await]
where $L(y) = o_1, H_0(o_1) = \langle \mathbf{Task}\langle \sigma \rangle, running(\bar{F}) \rangle$ $H_1 = H_0[o_1 \mapsto \langle \mathbf{Task}\langle \sigma \rangle, running(\langle L, x=y.GetResult(); \bar{s} \rangle^{a(o)}, \bar{F}) \rangle]$	
$H \triangleright \langle L, x=y.Result(); \bar{s} \rangle^\ell \rightarrow H \triangleright \langle L[x \mapsto v], \bar{s} \rangle^\ell$	[E-Result]
where $H(L(y)) = \langle \mathbf{Task}\langle \sigma \rangle, done(v) \rangle$	
$H \triangleright \langle L, x=y.Result(); \bar{s} \rangle^\ell \rightarrow H \triangleright \langle L, x=y.Result(); \bar{s} \rangle^\ell$	[E-Result-Block]
where $H(L(y)) = \langle \mathbf{Task}\langle \sigma \rangle, running(\bar{F}) \rangle$	
$H \triangleright \langle L, x=y.GetResult(); \bar{s} \rangle^\ell \rightarrow H \triangleright \langle L[x \mapsto v], \bar{s} \rangle^\ell$	[E-GetResult]
where $H(L(y)) = \langle \mathbf{Task}\langle \sigma \rangle, done(v) \rangle$	
$H_0 \triangleright \{ \langle L, x=\mathbf{Task.AsyncIO}\langle \gamma \rangle(); \bar{s} \rangle^\ell \circ FS \} \cup P$	[E-Async-IO]
$\sim H_1 \triangleright \{ \langle L[x \mapsto o], \bar{s} \rangle^\ell \circ FS \} \cup P \cup \{ \{ \{ y \mapsto v \}, \mathbf{return} y; \}^{a(o)} \circ \epsilon \}$	
where $o \notin dom(H_0), H_1 = H_0[o \mapsto \langle \mathbf{Task}\langle \gamma \rangle, running(\epsilon) \rangle]$ and $v \in Values(\gamma)$	

These transition rules cover the new asynchronous features in C[#] 5.0. First we recall that task heap objects are of the form $\langle \mathbf{Task}\langle \sigma \rangle, done(v) \rangle$ for some value v , or $\langle \mathbf{Task}\langle \sigma \rangle, running(\bar{F}) \rangle$ where \bar{F} is a sequence of frames—we will refer to this sequence as the *running state* of the task heap object. (In reality these two forms are encoded as conventional objects using delegates for the frames.) Tasks are *stateful*: a task heap object is created in initial state $\langle \mathbf{Task}\langle \sigma \rangle, running(\epsilon) \rangle$, with no waiters; can transition from state $\langle \mathbf{Task}\langle \sigma \rangle, running(\bar{F}) \rangle$ to $\langle \mathbf{Task}\langle \sigma \rangle, running(F_o, \bar{F}) \rangle$, adding one waiter, and may terminate in a *completed* state $\langle \mathbf{Task}\langle \sigma \rangle, done(v) \rangle$ for some value v of type σ . Once completed, a task cannot change state again.

Rule [E-Async-Method] shows how to transition a call to an asynchronous method. We create a fresh Task object in the heap (at address o), and set its state to be running. Initially, there are no waiters for this task, so its running sequence is empty. We push a new frame containing the method body on the frame stack and label it as asynchronous, i.e., with the label $a(o)$. The caller frame is updated with the heap address of the task in its locals stack. Notice that the calling frame is *not* awaiting a value, just control.

[E-Async-Return] pops the active asynchronous frame, storing the return value in the task. It also resumes any waiters (there may be zero or more). The operation

$resume(\overline{F})$ is used to resume a sequence of suspended frames. It creates a bag of singleton frame stacks and is defined as $resume(\overline{F}) \stackrel{\text{def}}{=} \{\langle L, \overline{s} \rangle^\ell \circ \epsilon \mid \langle L, \overline{s} \rangle^\ell \in \overline{F}\}$.

Rule [E-Await-Continue] covers the case when a task being awaited is already completed. In this case we simply read out the value from the task and continue. Rule [E-Await] covers the case when the task being awaited is still running. In this case we need to pause the asynchronous method. Thus we pop the active asynchronous frame from the frame stack and add it to the sequence of awaiters of the incomplete task. Notice that we unfold the `await` y to $y.GetResult()$ —once resumed, the first thing the frame will do is read the value from y 's (completed) task object in the shared heap.

Rule [E-Result] and [E-Result-Block] implement the in-built method `Result` on tasks. If the task is completed then it returns the result; if it is running then it ‘blocks’ (which for simplicity we simulate by spinning, i.e. by transitioning to itself). In contrast, rule [E-GetResult] implements `GetResult`. It too returns the result if the task is completed. However, if the task is incomplete, no rule applies and the configuration is stuck (the implementation raises an exception). `GetResult` is *non-blocking* and *partial*.

Rule [E-Async-IO] models a prototypical asynchronous method. It immediately returns a fresh, running task to be completed, with some value v , by a separate thread.

Process transition rules:

$H \triangleright \{\epsilon\} \cup P \rightsquigarrow H \triangleright P$	[E-Exit]
$H_0 \triangleright \{FS_0\} \cup P \rightsquigarrow H_1 \triangleright \{FS_1\} \cup P$	if $H_0 \triangleright FS_0 \twoheadrightarrow H_1 \triangleright FS_1$ [E-Schedule]

Recall that a process is a collection of frame stacks, i.e., threads. Rule [E-Exit] deletes an empty frame stack from the process. Rather than formalizing a particular scheduler, rule [E-Schedule] simply transitions a process by non-deterministically selecting and transitioning a thread, possibly side-effecting the shared heap. Our semantics is an interleaved semantics, allowing preemption at every atomic statement.

4 Correctness Properties

Given our formalization of FC_5^\sharp we are able to prove some important correctness properties; specifically, type soundness. Interestingly, establishing these properties involves non-trivial extensions of the conventional techniques [4]. In this section we give some details of these extensions and the precise forms of the correctness properties; complete details are given in a technical report.

The typical approach to proving type soundness involves extending the notion of type checking to configurations, and then establishing preservation and progress properties. However, for FC_5^\sharp this is not strong enough—in particular to establish progress—we have to consider not only type correctness but also crucial non-interference properties of tasks; both those being executed on framestacks and also those that are waiters on others tasks (and so are suspended in the heap). We also need to establish that the stateful protocol of tasks described in §3.1—that tasks begin in an empty running state, acquire waiters and then terminate in a done state (and never transition once in a done state)—is preserved too.

Rather than combine the typing and non-interference properties into a single relation, we keep them separate (at the expense of more verbose theorem statements). The rules

for non-interference for processes, framestacks, frames, heaps and heap objects are as follows.

Non-interference properties:

$$\begin{array}{c}
 \text{[Proc-ok]} \frac{\vdash (H \triangleright FS_0) \text{ ok} \quad \dots \quad \vdash (H \triangleright FS_n) \text{ ok} \quad \forall i \neq j \in \{0..n\}. \text{taskIds}(FS_i) \# \text{taskIds}(FS_j)}{\vdash (H \triangleright \{FS_0, \dots, FS_n\}) \text{ ok}} \\
 \\
 \text{[EmpFS-ok]} \frac{}{H \vdash \epsilon \text{ ok}} \quad \text{[FS-ok]} \frac{H \vdash F \text{ ok} \quad H \vdash FS \text{ ok} \quad \text{taskIds}(F) \# \text{taskIds}(FS)}{H \vdash F \circ FS \text{ ok}} \\
 \\
 \text{[SF-ok]} \frac{}{H \vdash F^s \text{ ok}} \quad \text{[CSF-ok]} \frac{H \vdash F^s \text{ ok}}{H \vdash F_x^s \text{ ok}} \\
 \\
 \text{[AF-ok]} \frac{\text{Running}(H(o)) \quad \forall o_1 \in \text{dom}(H). o \notin \text{runningIds}(H(o_1))}{H \vdash F^{a(o)} \text{ ok}} \\
 \\
 \text{[H-ok]} \frac{\forall o \in \text{dom}(H). H \vdash H(o) \text{ ok} \quad \forall o_1 \neq o_2 \in \text{dom}(H). \text{runningIds}(H(o_1)) \# \text{runningIds}(H(o_2))}{\vdash H \text{ ok}} \\
 \\
 \text{[HO-ok]} \frac{}{H \vdash \langle \mathcal{C}, FM \rangle \text{ ok}} \quad \text{[DTHO-ok]} \frac{}{H \vdash \langle \text{Task}\langle \sigma \rangle, \text{done}(v) \rangle \text{ ok}} \\
 \\
 \text{[RTHO-ok]} \frac{\forall i \neq j \in \{0..n\}. \text{taskIds}(F_i) \# \text{taskIds}(F_j) \quad \forall i \in \{0..n\}. \forall o \in \text{taskIds}(F_i). \text{Running}(H(o))}{H \vdash \langle \text{Task}\langle \sigma \rangle, \text{running}(F_0, \dots, F_n) \rangle \text{ ok}}
 \end{array}$$

We use a function $\text{taskIds}(FS)$ which returns the task ids of a frame stack FS (i.e., the set of all object ids o found in asynchronous frame labels $a(o)$ in the frame stack). We also overload this function over frames. We use a function runningIds that returns the task ids of the running state of a given task heap object. The predicate Running tests whether the state of a task heap object is currently running.

Rule [Proc-ok] ensures that the task ids in the frame stacks in a process are pairwise disjoint (we use the symbol $\#$ to denote disjointness). Rule [FS-ok] ensures that in a frame stack the task ids are all distinct. Rule [AF-ok] ensures that any task id in an asynchronous frame label is not included in the task ids of any running state in the heap. Rule [H-ok] ensures that for all the task heap objects in the heap, the task ids of the running states are disjoint. Rule [RTHO-ok] ensures that a given running task heap object has no duplicate task ids in its running state, and that all task ids in its running state refer to running (non-completed) tasks.

We also define a relation between heaps that preserves the typing of the heap objects and also enforce non-interference of any new running state whilst bounding the task ids of any new running state.

Definition 1 (Heap evolution). *Heap H_0 evolves to H_1 wrt a set of task ids S , written $H_0 \leq_S H_1$ if (i) $\forall o \in \text{dom}(H_1)$. if $o \notin \text{dom}(H_0)$ and $H_1(o_1) = \langle \psi, \text{running}(\overline{F}) \rangle$ then $\overline{F} = \epsilon$, and (ii) $\forall o \in \text{dom}(H_0)$. if $H_0(o) = \langle \mathcal{C}, FM_0 \rangle$ then $H_1(o) = \langle \mathcal{C}, FM_1 \rangle$, if $H_0(o) = \langle \psi, \text{done}(v) \rangle$ then $H_1(o) = \langle \psi, \text{done}(v) \rangle$, and if $H_0(o) = \langle \psi, \text{running}(\overline{F}_0) \rangle$ then $H_1(o) = \langle \psi, \text{running}(\overline{F}_1, \overline{F}_0) \rangle$, $\text{taskIds}(\overline{F}_0) \# \text{taskIds}(\overline{F}_1)$ and $\text{taskIds}(F_1) \subseteq S$.*

We also have typing relations for processes, framestacks, frames and heaps, written $\vdash (H \triangleright P): \star$, $H \vdash FS: \phi_0 \rightarrow \phi_1$, $H \vdash F: \phi_0 \rightarrow \phi_1$ and $\vdash H: \star$, respectively. Space prevents us from giving definitions of these relations, but they are routine.

Theorem 1 (Preservation). *If $\vdash H_0: \star$ and $\vdash H_0 \text{ ok}$ then:*

1. *If $\Gamma; H_0 \vdash F_0: \phi_0 \rightarrow \phi_1$, $H_0 \vdash F_0 \text{ ok}$ and $H_0 \triangleright F_0 \rightarrow H_1 \triangleright F_1$ then $\vdash H_1: \star$, $\vdash H_1 \text{ ok}$, $\Gamma; H_1 \vdash F_1: \phi_0 \rightarrow \phi_1$, $H_1 \vdash F_1 \text{ ok}$ and $\forall S. H_0 \leq_S H_1$.*
2. *If $H_0 \vdash FS_0: \phi_0 \rightarrow \phi_1$, $H_0 \vdash FS_0 \text{ ok}$ and $H_0 \triangleright FS_0 \rightarrow H_1 \triangleright FS_1$ then $\vdash H_1: \star$, $\vdash H_1 \text{ ok}$, $H_1 \vdash FS_1: \phi_0 \rightarrow \phi_1$, $H_1 \vdash FS_1 \text{ ok}$ and $H_0 \leq_{taskIds(FS_0)} H_1$.*
3. *If $\vdash (H_0 \triangleright P_0): \star$, $\vdash (H_0 \triangleright P_0) \text{ ok}$ and $H_0 \triangleright P_0 \rightsquigarrow H_1 \triangleright P_1$ then $\vdash H_1: \star$, $\vdash H_1 \text{ ok}$, $\vdash (H_1 \triangleright P_1): \star$ and $\vdash (H_1 \triangleright P_1) \text{ ok}$.*

Proof. Part (1) is proved by case analysis on $H_0 \triangleright F_0 \rightarrow H_1 \triangleright F_1$. Part (2) is proved by induction on the derivation of $H_0 \triangleright FS_0 \rightarrow H_1 \triangleright FS_1$ and part (1), and part (3) by induction on the derivation of $H_0 \triangleright P_0 \rightsquigarrow H_1 \triangleright P_1$ and part (2).

Theorem 2 (Progress). *If $\vdash (H_0 \triangleright P_0): \star$ and $\vdash (H_0 \triangleright P_0) \text{ ok}$ then*

1. *$H_0 \triangleright P_0 \rightsquigarrow H_1 \triangleright P_1$, for some H_1, P_1 ; or*
2. *for all $FS \in P_0$, one of the following holds:*
 - (a) *$FS = \langle L, \text{return } x; \bar{t} \rangle^s \circ \epsilon$.*
 - (b) *$FS = \langle L, y=x.m(\bar{z}); \bar{t} \rangle^l \circ FS'$, or $FS = \langle L, x.m(\bar{z}); \bar{t} \rangle^l \circ FS'$, or $FS = \langle L, y=x.f; \bar{t} \rangle^l \circ FS'$, or $FS = \langle L, x.f=y; \bar{t} \rangle^l \circ FS'$, where $L(x) = \text{null}$.*
 - (c) *$FS = \langle L, y=\text{await } x; \bar{t} \rangle^{a(o)} \circ FS'$, where $L(x) = \text{null}$.*
 - (d) *$FS = \langle L, \epsilon \rangle^l \circ FS'$.*
 - (e) *$FS = \langle L, y=x.\text{GetResult}(); \bar{t} \rangle^l \circ FS'$, where $L(x) = o$ and $H(L(x)) = \langle \text{Task}\langle \tau \rangle, \text{running}(\bar{F}) \rangle$.*

The progress theorem states a well-formed process can either transition or must entirely consist of stacks in terminal or stuck states (the latter includes the case $P_0 = \{\}$). Case 2a, a terminal state, can only arise from finishing a call to a program's non-void main method. Cases 2b–2c are familiar and new expected stuck states due to null references. Case 2d is excluded by applying C^\sharp 's restriction that all control paths in a (non-void) method body contain a return statement [12, §8.1]. Note that an asynchronous return $x; \bar{t}$ is never stuck due to the enclosing frame's task being in an unexpected done($_$) state; this potential case is ruled out by $\vdash H_0 \triangleright P_0 \text{ ok}$. Interestingly, we could also rule out case 2e by simply excluding any occurrences of GetResult in the original program; although the formal details are beyond the scope of this paper. With this restriction, the only occurrences of GetResult arise from the [E-Await] transition. These frames are only resumed by the rule [E-Async-Return] which also transitions the state of the task object to done(v). We can also show a property that no transition rule changes the state of a task that is completed back to running. These two properties allow us to show that case 2e does not arise for GetResult-free source programs.

5 Extensions

5.1 Extension 1: Optimized, One-Shot Semantics

The semantics presented so far is idealized: when an asynchronous frame is suspended to await a task, rule [E-Await] appends a copy of the frame to the task's list of waiters.

At first glance, the act of copying the frame appears to require an expensive allocation of a *fresh* frame to store its contents. Notice, however, that frames are never duplicated: after copying the frame, [E-Await] pops the active frame, discarding it to proceed with its continuation, the calling stack *FS*. Since frames are used in a *linear* fashion, the expensive allocation on each suspend is entirely avoidable. The trick to avoiding repeated allocation is to allocate just one container for each asynchronous frame and destructively update its contents at each suspension of that frame.

The C[#] 5.0 implementation does just this, representing a suspended frame on the heap as a “stateful” *delegate* of type `Action`. Delegates [12, Chapter 15] are just closures, containing the address of some environment and the address of some static code taking the environment as a first argument. Both addresses are immutable. The state of the frame is therefore maintained, not directly in the closure, but in its environment. To achieve this, the environment itself has mutable fields that store the current values of the frame’s locals, its associated task, and the current state of the finite state machine. All read and writes of locals in the original code are compiled to indirected operations on fields of the environment. The delegate’s code pointer just contains the fixed code interpreting the frame’s state machine.

In this section, we formalize a high-level abstraction of this implementation. Our formalization makes the more efficient, destructive update explicit without descending all the way to the low-level representation of closures used in the concrete implementation. To do so, we require a new reference type, the delegate type `Action`. In our semantics, if not in the actual implementation, the heap representation of an action is just an object whose mutable state is a frame, containing some locals and statements. The locals map contains the *current* values of local variables. The statements represent the frame’s original body in *some state of unfolding*, i.e., the frame’s current “program counter”. This allows us to adequately represent a paused frame, without exposing the compilation details of its encoding as a C[#] 4.0 delegate with a fixed pointer to a mutable environment and static code. Making this change also paves the way for our formalization of the *awaitable pattern* in §5.2.

First, we must extend $FC_5^\#$ with the `Action` type and syntax for invoking an action:

$FC_5^\#$ additional types and statements:

$\rho ::= \dots \mid \text{Action}$	Delegate reference type
$s ::= \dots \mid a();$	Action invocation statement

We also need to extend and adjust our run-time representations. `Action` is a new reference type so action values are just addresses of objects in the heap. An `Action` object, $\langle \text{Action}, F \rangle$, contains a (mutable) frame F , storing locals, statements and label of a suspended frame. We also need to modify tasks to track, not waiting frames (running(\overline{F})), but waiting *actions*, represented as a sequence of *addresses* (running(\overline{o})). Thus a running task will have representation $\langle \text{Task} \langle \sigma \rangle, \text{running}(\overline{o}) \rangle$; completed tasks remains the same. The form of an asynchronous label, placed on frames, is now $a(o_1, o_2)$. The new label carries not one but *two* addresses: the address of the frame’s task, o_1 , as before, and a second address, o_2 , of an action. The action stores the previous state of the frame; recording its address in the frame label indicates where to save the next state of the frame prior to suspending.

Completing a task will need to resume a list of actions, not frames, so we adapt the definition of $resume(\bar{o})$ to set up appropriate synchronous stubs, one per action in \bar{o} :
 $resume(\bar{o}) \stackrel{\text{def}}{=} \{ \langle \{x \mapsto o_i\}, x(); \mathbf{return}; \rangle^s \circ \epsilon \mid o_i \in \bar{o} \}$.

Asynchronous method transition rules (One-shot semantics):

$$H \triangleright \langle L, x(); \bar{s} \rangle^\ell \circ FS \rightsquigarrow H \triangleright F \circ \langle L, \bar{s} \rangle^\ell \circ FS \quad [\text{E-Action-Invoke}]$$

where $H(L(x)) = \langle \text{Action}, F \rangle$

$$H_0 \triangleright \langle L_0, y_0=y_1.m(\bar{z}); \bar{s} \rangle^\ell \circ FS \quad [\text{E-Async-MethodOS}]$$

$$\rightarrow H_1 \triangleright \langle L_1, \bar{t} \rangle^{a(o_1, o_2)} \circ \langle L_0[y_0 \mapsto o_1], \bar{s} \rangle^\ell \circ FS$$

where $H_0(L_0(y_1)) = \langle \sigma_0, FM \rangle$ and $mbody(\sigma_0, m) = mb: (\bar{\sigma} \bar{x}) \rightarrow^a \psi$
 $mb = \bar{\tau} \bar{y}; \bar{t}$ and $o_1, o_2 \notin dom(H_0), o_1 \neq o_2$
 $H_1 = H_0[o_1 \mapsto \langle \psi, \text{running}(\epsilon) \rangle], o_2 \mapsto \langle \cdot, \text{Action} \langle L_1, \bar{t} \rangle^{a(o_1, o_2)} \rangle]$
 $L_1 = [\bar{x} \mapsto L_0(\bar{z}), \bar{y} \mapsto \text{default}(\bar{\tau}), \mathbf{this} \mapsto L_0(y_1)]$

$$H_0 \triangleright \{ \langle L, \mathbf{return} y; \bar{s} \rangle^{a(o_1, o_2)} \circ FS \} \cup P \quad [\text{E-Async-ReturnOS}]$$

$$\rightsquigarrow H_1 \triangleright \{ FS \} \cup resume(\bar{o}) \cup P$$

where $H_0(o_1) = \langle \text{Task} \langle \sigma \rangle, \text{running}(\bar{o}) \rangle$
 $H_1 = H_0[o_1 \mapsto \langle \text{Task} \langle \sigma \rangle, \text{done}(L(y)) \rangle, o_2 \mapsto \langle \text{Action}, \langle L, \bar{s} \rangle^{a(o_1, o_2)} \rangle]$

$$H \triangleright \langle L, x=\mathbf{await} y; \bar{s} \rangle^{a(o_1, o_2)} \circ FS \quad [\text{E-Await-ContinueOS}]$$

$$\rightarrow H \triangleright \langle L[x \mapsto v], \bar{s} \rangle^{a(o_1, o_2)} \circ FS \quad \text{where } H(L(y)) = \langle \text{Task} \langle \sigma \rangle, \text{done}(v) \rangle$$

$$H_0 \triangleright \langle L, x=\mathbf{await} y; \bar{s} \rangle^{a(o_1, o_2)} \circ FS \rightsquigarrow H_1 \triangleright FS \quad [\text{E-AwaitOS}]$$

where $L(y) = o_3, H_0(o_3) = \langle \text{Task} \langle \sigma \rangle, \text{running}(\bar{o}) \rangle$
 $H_1 = H_0[o_3 \mapsto \langle \text{Task} \langle \sigma \rangle, \text{running}(o_2, \bar{o}) \rangle,$
 $o_2 \mapsto \langle \text{Action}, \langle L, x= \text{GetResult}(y); \bar{s} \rangle^{a(o_1, o_2)} \rangle]$

Rule [E-Action-Invoke] formalizes the invocation of an action, similar to a method call. Notice that the entire frame, including label, is restored from the heap. In particular, an asynchronous frame will continue to signal completion through its task and have access to its action (for future suspension, if needed).

Rule [E-Async-MethodOS] is similar to [E-Async-Method] but it additionally allocates a new `Action`, storing the initial state of the asynchronous method. The address of the action, o_2 , is recorded in the extended label of the pushed frame.

Rule [E-Async-ReturnOS] is similar to [E-Async-Return], completing the asynchronous frame's task. Though it is not necessary, we save the current locals and continuation of the `return`, \bar{s} , in the frame's `Action`. For this simple semantics, it should be possible to show that this action can never be invoked again.⁶

Rule [E-Await-ContinueOS] is almost identical to [E-Await-Continue], continuing execution of the current frame with the argument's result. The only difference is the extended label. There is no need to update the value of o_2 at this point. Rule [E-AwaitOS] is similar to [E-Await], but the suspend mechanism is different. This rule writes the frame's current state, locals and continuation, to its associated action, stored at address o_2 , available from the frame's label. It then adds the address of that action to the

⁶ When we add support for the awaitable pattern, the potential for abuse of the awaitable protocol, will mean that this property no longer generally holds.

incomplete task's list of waiters. Notice how the state of the action in the heap is destructively modified - there is no way to “go back” to a previous state of this frame.

Consider rule [E-Async-MethodOS]. It directly pushes a new asynchronous frame and assigns its task, o_1 , to the caller's variable, y_0 . An alternative formulation would be to push a synchronous stub that invokes the new action, o_2 , and then returns the task, o_1 , to the waiting caller. This would be less direct, but equivalent, and somewhat more faithful to the actual implementation. For example, the implementation of `CopyToManual` from §2.1 is essentially a stub method that, when called, invokes its internal delegate, `act()`, before returning its task.

At this point, the change to using mutable state to represent suspended frames is just an optimization. The reason is that user-code is never provided with access to a suspended frame, so the change in semantics cannot be observed.

5.2 Extension 2: The Awaitable Pattern

As detailed in §2, in C[#] 5.0 it is possible to await not just tasks, but values of any *awaitable* type. Our formalization has assumed that the only awaitable type is `Task< σ >`. In this section, we embrace the full awaitable pattern, replacing rule [C-Await] with:

New typing rule for awaitable expressions ([C-Awaitable])

$$\frac{\begin{array}{l} mtype(\sigma_0, \text{GetAwaiter}) = () \rightarrow \sigma_1 \quad mtype(\sigma_1, \text{IsCompleted}) = () \rightarrow \text{bool} \\ mtype(\sigma_1, \text{OnCompleted}) = (\text{Action}) \rightarrow \text{void} \quad mtype(\sigma_1, \text{GetResult}) = () \rightarrow \sigma_2 \end{array}}{\Gamma, x: \sigma_0 \vdash \text{await } x: \sigma_2}$$

We simplify C[#] 5.0 and assume the property `IsCompleted` is an ordinary *method*; the distinction between methods and properties is entirely cosmetic so nothing is lost.

In the transition semantics `await` expressions can no longer transition atomically but must, instead, be evaluated in multiple steps. These steps commence with obtaining the argument's awaiter and proceed with calls to the awaiter's members, thus interleaving (potentially) user-defined code with the semantics of the `await` construct. Rule [C-Awaitable] statically ensures that these dynamic unfoldings are well-typed.

But first, we need to arrange that tasks are awaitable and implement the remaining requirements of the awaitable pattern. Our system already provides an appropriate `GetResult` for tasks; we are left with providing `GetAwaiter`, `IsCompleted` and `OnCompleted`, ascribed with the following types:

$$\begin{array}{l} mtype(\text{Task}<\sigma>, \text{GetAwaiter}) = () \rightarrow \text{Task}<\sigma> \\ mtype(\text{Task}<\sigma>, \text{IsCompleted}) = () \rightarrow \text{bool} \\ mtype(\text{Task}<\sigma>, \text{OnCompleted}) = (\text{Action}) \rightarrow \text{void} \\ mtype(\text{Task}<\sigma>, \text{GetResult}) = () \rightarrow \sigma \end{array}$$

To avoid hard-wiring C[#] 5.0's generic `TaskAwaiter< σ >` type, we simplify the C[#] 5.0 design and assume that `Task< σ >` is self-sufficient and serves as its *own* awaiter type. Correspondingly, `x.GetAwaiter()`'s type is just the type of task `x`; its implementation, by rule [E-Task-GetAwaiter] below, just returns the receiver.

Additional transition rules for Task's awaitable operations:

$$H \triangleright \langle L, x=y.\text{GetAwaiter}(); \bar{s} \rangle^\ell \quad [\text{E-Task-GetAwaiter}]$$

$$\rightarrow H \triangleright \langle L[x \mapsto L(y)], \bar{s} \rangle^\ell \text{ where } H(L(y)) = \langle \text{Task}\langle \sigma \rangle, FM \rangle$$

$$H \triangleright \langle L, x=y.\text{IsCompleted}(); \bar{s} \rangle^\ell \quad [\text{E-Task-IsCompleted}]$$

$$\rightarrow H \triangleright \langle L[x \mapsto \text{true}], \bar{s} \rangle^\ell \text{ where } H(L(y)) = \langle \text{Task}\langle \sigma \rangle, \text{done}(v) \rangle$$

$$\rightarrow H \triangleright \langle L[x \mapsto \text{false}], \bar{s} \rangle^\ell \text{ where } H(L(y)) = \langle \text{Task}\langle \sigma \rangle, \text{running}(\bar{o}) \rangle$$

$$H_0 \triangleright \langle L, x.\text{OnCompleted}(y); \bar{s} \rangle^\ell \rightarrow H_1 \triangleright \langle L, \bar{s} \rangle^\ell \quad [\text{E-Task-OnCompleted-Suspend}]$$

where $L(x) = o_1$ and $H_0(o_1) = \langle \text{Task}\langle \sigma \rangle, \text{running}(\bar{o}) \rangle$
 $L(y) = o_2$ and $H_1 = H_0[o_1 \mapsto \langle \text{Task}\langle \sigma \rangle, \text{running}(o_2, \bar{o}) \rangle]$

$$H \triangleright \{ \langle L, x.\text{OnCompleted}(y); \bar{s} \rangle^\ell \circ FS \} \cup P \quad [\text{E-Task-OnCompleted-Resume}]$$

$$\rightsquigarrow H \triangleright \{ \langle L, \bar{s} \rangle^\ell \circ FS \} \cup \text{resume}(o) \cup P$$

where $H(L(x)) = \langle \text{Task}\langle \sigma \rangle, \text{done}(v) \rangle$ and $L(y) = o$

Task's implementation of `IsCompleted()` tests the state field of the receiver, returning `true` if and only if it is `done(-)`. The implementation of `OnCompleted(y)` adds its callback y (an Action), to the receiver's list of waiters. If the task is already completed, the action cannot be stored and must, instead, be resumed in the process. The latter rule is required since there is a race between testing that a task `IsCompleted()`, finding it is `false`, and calling `OnCompleted(y)`—some other thread could intervene and complete the task *before* `OnCompleted(y)` executes.

We can now formalize the operational semantics of `await` on any awaitable. Because we need to interleave the execution of methods from the awaitable pattern—which take several transitions and could be user-defined—with the semantics of `await`, we need to introduce two additional, transient *control* statements that can only appear within asynchronous frames.

FC_s[‡] additional control statements:

$$s ::= \dots \mid \text{suspend}; \mid \text{getcc}(\text{Action } a) \{ \bar{s} \}; \mid \text{suspend} \ \& \ \text{get-current-continuation statements}$$

Though artificial, these statements have direct interpretations as intermediate steps of a compiler generated finite-state-machine.⁷

We define $\text{unfold}(x = \text{await } y; \bar{s})^{(z,b)}$ to be the syntactic unfolding of an `await` as a new sequence of statements using temporaries z and b (note a is bound):

$$\text{unfold}(x = \text{await } y; \bar{s})^{(z,b)} \stackrel{\text{def}}{=} \begin{cases} z = y.\text{GetAwaiter}(); \\ b = z.\text{IsCompleted}(); \\ \text{if } (b) \{ \} \text{ else} \\ \{ \text{getcc}(\text{Action } a) \{ z.\text{OnCompleted}(a); \text{suspend}; \}; \\ x = z.\text{GetResult}(); \\ \bar{s} \end{cases}$$

The operation unfolds an `await` of an awaitable object y by first retrieving its awaiter z and setting b to determine if the awaiter is complete. If complete, the code falls through

⁷ For example, in our hand-coded `CopyToManual` from §2.1, `suspend` corresponds to a `return`; from the `act` delegate that pauses execution (cases 1 and 2 of the switch); `getcc(Action a) {s}`; corresponds to advancing the (shared) `state` variable to the next logical state (following `getcc(Action a) {s}`;) and accessing the task's state machine (`act`).

the conditional. If incomplete, the code transfers the current continuation of the `getcc` statement to z (through a) and suspends. The continuation of both the `true` branch and the `getcc` statement is just $x = z.\text{GetResult}(); \bar{s}$. It assigns the result of the awaiter to x , and proceeds with the original continuation \bar{s} of the `await`.

Awaitable pattern, asynchronous method transition rules:

$$H_0 \triangleright \langle L_0, x = \text{await } y; \bar{s} \rangle^{a(o_1, o_2)} \quad \text{[E-Awaitable]}$$

$$\rightarrow H_0 \vdash \langle L_1, \text{unfold}(x = \text{await } y; \bar{s})^{(z, b)} \rangle^{a(o_1, o_2)}$$

where $L_1 = L_0[z \mapsto \text{null}, b \mapsto \text{false}]$ and $z, b \notin \text{dom}(L_0), z \neq b$

$$H_0 \triangleright \langle L, \text{getcc}(\text{Action } a) \{ \bar{s}; \bar{t} \} \rangle^{a(o_1, o_2)} \rightarrow H_1 \triangleright \langle L[a \mapsto o_2], \bar{s} \rangle^{a(o_1, o_2)} \quad \text{[E-GetCC]}$$

where $a \notin \text{dom}(L)$ and $H_1 = H_0[o_2 \mapsto \langle \text{Action}, \langle L, \bar{t} \rangle \rangle^{a(o_1, o_2)}]$

$$H \triangleright \langle L, \text{suspend}; \bar{s} \rangle^{a(o_1, o_2)} \circ FS \rightarrow H \triangleright FS \quad \text{[E-Suspend]}$$

Rule [E-Awaitable] unfolds its await expression using fresh temporaries, z and b . In rule [E-GetCC], `getcc(Action a) { \bar{s} ; \bar{t} }` unfolds by first saving its continuation \bar{t} in the frame's task, o_2 , discarding it from the active frame, and then entering the body \bar{s} . When \bar{s} is just $z.\text{OnCompleted}(a); \text{suspend};$, as per rule [E-Awaitable], \bar{s} will *transfer* the current continuation to the awaiter and suspend.

In rule [E-Suspend] the suspend control statement pauses the asynchronous frame. This is similar to a `return`, but the frame's task is not marked completed and remains in its current state. One might expect this state to be `running(⌊)` but it may not be, depending on the semantics of `OnCompleted`.

Although our semantics unfolds `awaits` dynamically, it is possible to statically expand well-typed `await` expressions by a source-to-source translation, sketched here:

$$\begin{aligned} [x = \text{await } y; \bar{s}_0]^{\Gamma} &\stackrel{\text{def}}{=} (\Gamma_1, \bar{s}_2) \text{ where } \text{mtype}(\Gamma(y), \text{GetAwaiter}) = () \rightarrow \sigma_1 \\ &(\Gamma_1, \bar{s}_1) = [\bar{s}_0]^{\Gamma, z: \sigma_1, b: \text{bool}} \\ &\bar{s}_2 = \text{unfold}(x = \text{await } y; \bar{s}_1)^{(z, b)} \\ &b, z \notin \text{dom}(\Gamma), b \neq z \\ [s; \bar{s}]^{\Gamma} &\stackrel{\text{def}}{=} \dots \end{aligned}$$

This translation must be type-directed (in order to determine awaiter types) and needs to produce a new context as well as the list of statements in order to properly account for generated variables. Notice, however, that it is finite and does not need to duplicate the input continuation \bar{s}_0 , making it suitable for compile-time expansion.

OnCompleted's one-shot Restriction, Explained. Once we add the awaitable pattern to the mix, the optimization described in §5.1 becomes a proper change to the semantics, with observable consequences. The culprit is the awaitable pattern's `OnCompleted(a)` method since it provides user-code with access to the *one-shot* continuation, a , of the frame, represented not as a pure value but as a stateful object. Recall that our informal description of the awaiter pattern stipulated that implementations of `OnCompleted` are required to invoke their action *at most once*. The reason why should now be clear. Invoking the action will resume the frame and potentially modify the action's state. In our semantics, the update would happen at the next suspension. In the real implementation,

the update would happen at the next write to some notionally local, but actually shared, variable of the frame. Two concurrent invocations of the same action have unpredictable behaviour: each would race to save its next, possibly different state in the same action. Frame execution depends on the shared heap, modified non-deterministically by rule [E-Schedule], so two invocations could very easily reach different states.

As it happens, when extended with the awaitable pattern, even the simpler, copying semantics cannot tolerate multiple invocations of a continuation, but the reason is more subtle. In the copying semantics, even a copy of a frame is inherently stateful because its label will contain a reference to the original frame's task. Allocated on the heap, this task is shared state: several invocations of the same continuation would race, with possibly different results, to complete the very same task on exit. Part of the semantics of tasks is that they should complete at most once. This invariant is violated by any abuse of one-shot continuations, as enabled by the awaitable pattern.

6 Related Work

The debate regarding how asynchronous software should be structured is both old and ongoing. Lauer and Needham [14] noted that the thread-based and event-based models are dual; a program written in one style can be transformed into a program written in the other style. Though this establishes that the two models are equivalent in expressive power, it does not resolve the question of which model is easier to use or reason about.

Ousterhout [17] famously stated that “threads are a bad idea (for most programs).” His argument revolves around the claim that threads are more difficult to program than events because the programmer must reason about shared state, locks, race conditions, etc., and that they are only necessary when true concurrency—in contrast to asynchrony—is desired. Though he conflates the threaded model of programming, in which there is no inversion of control, with concurrency, his observation that the programmer should be able to reason about the operation of code is well-taken.

SEDA [25] demonstrates that the event model can be highly scalable. Servers designed using SEDA are broken into separate stages with associated event queues. Each stage automatically tunes its resource usage and computation to meet application-wide performance goals. Within each stage multiple threads may process events, but these threads are utilized only for concurrency. The programmer still has to manually manage the state associated with each event. SEDA's goal is to provide a self-tuning architecture that adapts to overload conditions, not to make programming servers easier.

The dual argument in favor of threads over events is made by von Behren et al. [22]. They tease apart the different aspects of threads that may make them undesirable and argue that most of these deficiencies are merely implementation artifacts. Capriccio [23] demonstrates that this is the case by providing a very efficient cooperative threading mechanism that avoids inversion of control and provides an efficient runtime. Because it uses cooperative threading, Capriccio avoids the overhead of concurrency, and code transformations to insert stack checks allow threads' stacks to grow without requiring large amounts of pre-allocated stack space. Like Capriccio, asynchronous C[#] allows programs to be written in a natural way while providing an efficient implementation. However, instead of attempting to provide a general cooperative threading mechanism, it permits programmers to write asynchronous code in a straight-line fashion by

automating stack-ripping [1] via compilation to a state machine. Like the state machine translation C[#] employs, it is reminiscent of simpler coroutine implementations [21,7] built on Duff's device [6] that, however, make no attempt to maintain local state.

The observation that continuations provide a natural substrate on which to build a threading mechanism was made by Wand [24]. The observation that more restrictive, but more efficient, one-shot continuations suffice for continuation based threading dates back to [5]. Li and Zdancewic [15] use continuations to unify the event and thread-based models. They leverage the continuation monad in Haskell to allow programmers to write straight line code that is desugared into continuation passing style (CPS), thus allowing it to be used in an event-based IO framework they construct.

The *computation expressions* [20] of F[#] provide a generalized monadic syntax, which is essentially an extended form of Haskell's do-notation. When specialised to F[#]'s *asynchronous workflow* monad—itself a continuation monad—they allow programmers to write monadic code that is syntactically expanded to explicit continuation-passing-code. This offers much of the legibility of programming in direct-style while, at the same time, providing access to the implicit continuations (as F[#] functions) whenever required (e.g., when supplying callbacks to asynchronous calls). The generality of computation expressions has a cost: each continuation of a monadic let is a heap-allocated function; and every wait on an asynchronous value typically requires an expensive allocation of a fresh continuation. This is similar to our idealized semantics in §3. The upshot is that these continuations can, in principle, be invoked several times, allowing the encoding of a much wider range of control operators than the one-shot actions of C[#]'s feature. But there are more differences. In F[#], computation expressions produce *inert* values that are easily composed but must be explicitly run to produce a result. In C[#], on the other hand, each task returned by an `async` method call represents a *running* computation. This makes it easier to initiate asynchronous work but, perhaps, harder to define combinators that compose asynchronous methods. Though inspired by F[#], C[#] 5.0 support for asynchrony is quite different in performance, expressivity and usage.

Scala actors [11] provide an asynchronous, message passing model for programming with concurrent processes. Asynchronous programs may be written in terms of *receive*, which suspends the current thread until a message is received, or *react*, which enqueues a continuation that is called when a message is received; *receive* provides a thread-based interface and *react* provides an event-based interface to an underlying message passing framework. Although the two programming models are made similar through the use of various combinators, the programmer must still significantly modify code to move between styles. Rompf et al. [19] use a type and effect system to selectively CPS convert Scala programs, providing a less onerous path from threads to event-based asynchronous code, but C[#] 5.0's `async` keyword is even more lightweight.

Despite Ousterhout's early admonition that reasoning about threads is difficult and error-prone, none of the work mentioned makes an explicit attempt to provide programmers with a set of reasoning principles for asynchronous code. Although we believe C[#]'s support for asynchrony exists at a useful point in the design space, our focus is on providing these reasoning principles. Threads or events, manual stack ripping or CPS, a programmer must have clear ways to reason about code behavior in order to build correct systems of any kind.

7 Conclusions

Real-world software construction demands effective methods for dealing with asynchrony. For such a method to be termed “effective,” it must not require large-scale, manual code transformations such as stack ripping. Sequential computations should be expressible with sequential code, even if individual operations may execute asynchronously. Splitting sequential code up into a series of callbacks or explicitly rewriting it as a state machine is a steep price to pay, making code difficult to write, difficult to read, and difficult to reason about; if in doubt, contrast the synchronous, `async`-enabled, and hand-written state machine versions of the stream copying function from §2.1. While previous work has provided syntactic and library support for dealing with asynchrony, `C# 5.0` brings this support to a widely-deployed, mainstream language.

One deficiency of this previous work is a lack of reasoning principles for asynchronous code. Our primary contribution is an operational semantics for `C# 5.0` that allows programmers to answer questions about the code they write and make conclusions about the impact of adding asynchrony to their code. For example, using our semantics, the programmer can see that calling an `async` method does not spawn a new thread, but instead executes the method on the current stack. With the optimized semantics in §5.1, one can even begin to reason about space usage by, e.g., observing that the state of an `async` method is always stored in the same `Action`, allocated just once.

We plan to continue our formalization of `C# 5.0` by incorporating additional language features, such as cancellation tokens and synchronization context object—we have already formalized exceptions and their interaction with `async`, but due to space restrictions this formalization, as well as an asynchronous tail-call optimization, is only available in a separate tech report. Our semantics have been translated to Coq. We will use this as a foundation to validate a translation from Featherweight `C# 5.0`—including the `async` construct—to Featherweight CIL, an idealized version of the bytecode targeted by the `C# 5.0` compiler. Validation of this translation will prove that programmers can reason in terms of our high-level operational semantics even though the high-level program has been translated to bytecode and it is the bytecode that is actually executed.

While syntax is important for easing the pain of writing asynchronous code, a corresponding semantics is vital for writing correct software. With our semantics, `C# 5.0` both provides relief and the necessary tools for thinking carefully about the remedy.

Acknowledgements. We thank the `C#` and Visual Basic teams for their collaboration, especially Lucian Wischik who led much of the design and implementation effort.

References

1. Adya, A., Howell, J., Theimer, M., Bolosky, W.J., Douceur, J.R.: Cooperative task management without manual stack management. In: Proceedings of USENIX (2002)
2. Bierman, G., Meijer, E., Torgersen, M.: Lost in translation: Formalizing proposed extensions to `C#`. In: Proceedings of OOPSLA (2007)
3. Bierman, G., Meijer, E., Torgersen, M.: Adding Dynamic Types to `C#`. In: D’Hondt, T. (ed.) ECOOP 2010. LNCS, vol. 6183, pp. 76–100. Springer, Heidelberg (2010)
4. Bierman, G., Parkinson, M., Pitts, A.: MJ: An imperative core calculus for Java and Java with effects. Technical Report 563, University of Cambridge Computer Laboratory (2003)

5. Bruggeman, C., Waddell, O., Dybvig, R.K.: Representing control in the presence of one-shot continuations. In: Proceedings of PLDI (1996)
6. Duff, T.: Re: Explanation, please! (August 1988), USENET Article
7. Dunkels, A., Schmidt, O., Voigt, T., Ali, M.: Protothreads: Simplifying Event-Driven programming of Memory-Constrained embedded systems. In: Proceedings of SenSys (2006)
8. Fischer, J., Majumdar, R., Millstein, T.: Tasks: language support for event-driven programming. In: Proceedings of PEPM (2007)
9. Flanagan, C., Sabry, A., Duba, B., Felleisen, M.: The essence of compiling with continuations. In: Proceedings of PLDI (1993)
10. Flatt, M., Krishnamurthi, S., Felleisen, M.: A programmer's reduction semantics for classes and mixins. Technical Report TR-97-293, Rice University (1997)
11. Haller, P., Odersky, M.: Scala actors: Unifying thread-based and event-based programming. *Theoretical Computer Science* 410(2-3), 202–220 (2009)
12. Hejlsberg, A., Torgersen, M., Wiltamuth, S., Golde, P.: *The C# Programming Language*, 4th edn. Addison-Wesley (2011)
13. Igarashi, A., Pierce, B., Wadler, P.: Featherweight Java: A minimal core calculus for Java and GJ. *ACM TOPLAS* 23(3), 396–450 (2001)
14. Lauer, H.C., Needham, R.M.: On the duality of operating system structures. *Operating Systems Review* 13(2), 3–19 (1979)
15. Li, P., Zdancewic, S.: Combining events and threads for scalable network services. In: Proceedings of PLDI (2007)
16. Microsoft Corporation. C# language specification for asynchronous functions (2011), <http://msdn.microsoft.com/en-us/vstudio/async>
17. Ousterhout, J.K.: Why threads are a bad idea (for most purposes). In: USENIX Winter Technical Conference, Invited Talk (June 1996)
18. Pierce, B., Turner, D.: Local type inference. In: Proceedings of POPL (1998)
19. Rompf, T., Maier, I., Odersky, M.: Implementing first-class polymorphic delimited continuations by a type-directed selective CPS-transform. In: Proceedings of ICFP (2009)
20. Syme, D., Petricek, T., Lomov, D.: The F# Asynchronous Programming Model. In: Rocha, R., Launchbury, J. (eds.) *PADL 2011*. LNCS, vol. 6539, pp. 175–189. Springer, Heidelberg (2011)
21. Tatham, S.: Coroutines in C (2000), <http://www.chiark.greenend.org.uk/~sgtatham/coroutines.html>
22. von Behren, R., Condit, J., Brewer, E.: Why events are a bad idea (for high-concurrency servers). In: Proceedings of HotOS (2003)
23. von Behren, R., Condit, J., Zhou, F., Necula, G.C., Brewer, E.: Capriccio: scalable threads for internet services. In: Proceedings of SOSP (2003)
24. Wand, M.: Continuation-based multiprocessing. In: Proceedings of LISP and Functional Programming (1980)
25. Welsh, M., Culler, D., Brewer, E.: SEDA: an architecture for well-conditioned, scalable internet services. In: Proceedings of SOSP (2001)