

The essence of data access in $C\omega$

The power is in the dot!

Gavin Bierman
Microsoft Research
gmb@microsoft.com

Erik Meijer
Microsoft Corporation
emeijer@microsoft.com

Wolfram Schulte
Microsoft Research
schulte@microsoft.com

Abstract

In this paper we describe the data access features of $C\omega$, an experimental programming language under development at Microsoft Research. $C\omega$ targets distributed, data-intensive applications and accordingly extends C^\sharp 's support of both data and control. In the data dimension it provides a type-theoretic integration of the three prevalent data models, namely the object, relational, and semi-structured models of data. In the control dimension $C\omega$ provides elegant primitives for asynchronous communication.

In this paper we concentrate on the data dimension. Our aim is to describe the *essence* of these extensions; by which we mean we identify, exemplify and formalize their essential features. Our tool is a small core language, $FC\omega$, which is a valid subset of the full $C\omega$ language. Using this core language we are able to formalize both the type system and the operational semantics of the data access fragment of $C\omega$.

1 Introduction

Programming languages, like living organisms, need to continuously evolve in response to their changing environment. These evolutionary steps are typically quite modest: most commonly the provision of better or reorganized APIs. Occasionally a more radical evolutionary step is taken. One such example is the addition of generic classes to both Java [5] and C^\sharp [19].

We should like to argue that the time has come for another large evolutionary step to be taken. Much software is now intended for distributed, web-based scenarios. It is typically structured using a three-tier model consisting of a *middle tier* containing the business logic that extracts relational data

from a *data services tier* (a database) and processes it to produce semi-structured data (typically XML) to be displayed in the *user interface tier*.

It is the writing of these middle tier applications that we should like to address. These applications are most commonly written in an object-oriented language such as Java or C^\sharp and have to deal with relational data (essentially SQL tables), object graphs, and semi-structured data (XML, HTML).

In addition, these applications are fundamentally concurrent. Because of the inherent latency in network communication, the more natural model of concurrency is asynchronous. Accordingly, $C\omega$ provides a simple model of asynchronous (one-way) concurrency based on the join calculus [10]. For the rest of this paper, we shall focus exclusively on the data access aspects of $C\omega$; the concurrency primitives have been discussed elsewhere [1]. Thus when we write $C\omega$, we mean the language excluding the concurrency primitives.

Unfortunately programming language support for data access has barely evolved at all. All that exists is naïve access via simple APIs. Consider the following fragment of Java that uses JDBC to query a SQL database (a user-supplied country is stored in variable `input`).

```
Connection con = DriverManager.getConnection(...);
Statement stmt = con.createStatement();
String query = "SELECT * FROM COFFEES"+
               "WHERE Country='"+input+"'";
ResultSet rs = stmt.executeQuery(query);
while (rs.next()) {
    String s = rs.getString("Cof_Name");
    float n = rs.getFloat("Price");
    System.out.println(s+" - "+n);
}
```

Using strings to represent SQL queries is not only clumsy but also removes any possibility for static checking. The impedance mismatch between the language and the relational data is quite striking; e.g. a value is projected out of a row by passing a string denoting the column name and using the appropriate conversion function. Perhaps most seri-

ously, the passing of queries as strings is often a security risk (the “script code injection” problem—e.g. consider the case when the variable input is the string " ' OR 1=1; --" [15]).

Unfortunately API support in both Java and C[#] for XML and XPath/XQuery is depressingly similar (even those APIs that map XML values tightly to an object representation, still offer querying facilities by string passing). In summary, our contention is that common object-oriented languages need to evolve to support data access satisfactorily. Our observation is that what is missing is the support of the rich structure of both relational and semi-structured data. Our solution is then to enrich object-oriented languages with the structure inherent in relational and semi-structured data, and to enhance them with familiar query-like capabilities.

Design objectives of C ω We have seen above the need for first-class language support for the manipulation of relational and semi-structured data. The question remains how to provide it. One possibility is to design a special purpose language [13, 2], but for the applications we have in mind this is impractical. Instead, we choose to evolve an existing language, C[#], into a new language that we call C ω . Although we have started with C[#], our extensions apply equally well to other object-oriented languages, including Java.

Addressing the title of our paper, the essence of C ω is twofold: its extensions to the C[#] type system and, perhaps more importantly, the elegant provision of query-like capabilities (the sub-title of our paper). C ω has been carefully designed around a set of core design principals.

1. C ω is a coherent extension of C[#], i.e. C[#] programs should be valid C ω programs with the same behaviour.
2. The type system of C ω is intended to be both as simple as possible and closely aligned to the type system in the XPath/XQuery standard. Our intended users are C[#] programmers who are familiar with XPath/XQuery.
3. From a programming perspective, the real power of C ω comes from its elegant query-like capabilities. These have been achieved by generalizing member access to allow simple XPath-like queries.

Paper organization The rest of the paper is organized as follows. In §2 we give a comprehensive overview to the C ω programming language. In §3.1 we identify and formalize FC ω , a core calculus of C ω . In §3.2 we detail a simpler fragment, IC ω , and in §3.4 show how FC ω can be compiled to IC ω . Using this compilation, we are able to show a number of properties of FC ω in §3.5, including a type soundness theorem. We briefly discuss some related work in §4 and conclude in §5.

2 An introduction to C ω

Our design goal was to evolve C[#] to provide an integration of the object, relational and semi-structured data models. One possibility would be to add these data models to our programming language in an orthogonal way, e.g. by including new types XML<S> and TABLE<R>, where S and R are XML and relational schema respectively. Rather we have sought to integrate these models by *generalization*, rather than by ad-hoc specializations. In the rest of this section we shall present the key ideas behind C ω , and give a number of small programs to illustrate these ideas. This section should serve as a programmer’s introduction to C ω . We assume that the reader is familiar with C[#]/Java-like languages.

2.1 New types

C ω is an extension of C[#], so the familiar primitive types such as integers, booleans, floats are present, as well as classes and interfaces. In this section we shall consider in turn the extensions to the type system—streams, anonymous structs, discriminated unions, and content classes—and for each consider the new query capabilities.

Streams The first structural type we add is a stream type; streams represent ordered homogeneous collections of zero or more values. For example, int* is the type for homogeneous sequences of integers. Streams in C ω are aligned with iterators, which will appear in C[#] 2.0. C ω streams are typically generated using iterators, which are blocks that contain yield statements. For example, the FromTo method:

```
virtual int* FromTo(int b, int e){
    for (i = b; i <= e; i++) yield return i;
}
```

generates a finite, increasing stream of integers. Importantly, it should be noted that, just as for C[#], invoking such a method body does *not* immediately execute the iterator code, but rather immediately returns a closure. (Thus C ω streams are essentially lazy lists, in the Haskell sense.) This closure is consumed by the foreach statement, e.g. the following code fragment builds a finite stream and then iterates over the elements, printing each one to the screen.

```
int* OneToHundred = FromTo(1,100);
foreach (int i in OneToHundred) Console.WriteLine(i);
```

A vital aspect of C ω streams is that they are always *flattened*; there are no nested streams of streams. C ω streams thus coincide with XPath/XQuery sequences which are also flattened. This alignment is a key design decision for C ω : it enables the semantics of our generalized member access to

match the path selection of XQuery. We give further details later.

In addition, flattening of stream types also allows us to efficiently deal with recursively defined streams. Consider the following recursive variation of the function `FromTo` that we defined previously:

```
virtual int* FromTo2(int b, int e){
    if (b>e) yield break;
    yield return b;
    yield return FromTo2(b+1,e);
}
```

The statement `yield break;` returns the empty stream. The non-recursive call `yield return b` yields a single integer. The recursive call `yield return FromTo2(b+1,n);` yields a stream of integers. As the type system treats the types `int*` and `int**` as equivalent this is type correct.

Without flattening we would be forced to copy the stream produced by the recursive invocation, leading to a quadratic instead of a linear number of yields:

```
virtual int* FromTo3(int b, int e){
    if (b>e) yield break;
    yield return b;
    foreach (int i in FromTo3(b+1,e)) yield return i;
}
```

Note that $C\omega$'s flattening of stream types does *not* imply that the underlying stream is flattened via some coercion; every element in a stream is `yield`-ed at most once. As we will see in the operational semantics (§3.3), iterating over a stream will effectively perform a depth-first traversal over the n -ary tree produced by the iterator.

$C\omega$ offers a limited but extremely useful form of *covariance* for streams. Covariance is allowed provided that the conversion on the element type is the identity; for example `Button*` is a subtype of `object*` whereas `int*` is *not* (as the conversion from `int` to `object` involves boxing). This notion is a simple extension of the notion of covariance for arrays in C^\sharp , although it is safe (unlike array covariance) as we can not overwrite elements of streams.

The rationale for this is that implicit conversions should be limited to constant-time operations. Coercing a stream of type `int*` to a stream of type `object*` would be linear in the length of the stream, as the boxing conversion from `int` to `object` is not the identity.

A key programming feature of $C\omega$ is generalized member access; as the subtitle suggests the familiar 'dot' operator is now much more powerful. Thus if the receiver is a stream the member access is mapped over the elements, e.g. `OneToHundred.ToString()` implicitly maps the method call over the elements of the stream

`OneToHundred` and returns a value of type `string*`. This feature significantly reduces the burden on the programmer. Moreover, member access has been generalized so it behaves like a *path expression*. For example, `OneToHundred.ToString().PadLeft(10)` converts all the elements of the stream `OneToHundred` to a string, and then pads each string, returning a stream of these padded strings.

Sometimes one wishes to map more than a simple member access over the elements of a stream. $C\omega$ offers a convenient shorthand called an *apply-to-all expression*, written $e.\{\bar{s}\}$, which applies the block $\{\bar{s}\}$ to each element in the stream e .¹ The block may contain the variable `it` which plays a similar role as the implicit receiver argument `this` in method bodies and is bound to each successive element of the iterated stream. For example, the following code first creates the stream of natural numbers from 1 to 256, converts each of the elements to a hex string, converts each of these to upper case, and then applies an apply-to-all expression to print the elements to the screen:

```
FromTo(1,256).ToString("x").ToUpper().
    { Console.WriteLine(it); };
```

Anonymous structs The second structural type we add are anonymous structs, which encapsulate heterogeneous ordered collections of values. An anonymous struct is like a tuple in ML or Haskell and is written as `struct{int i; Button;}` for example. A value of this type contains a member `i` of type `int` and an unlabelled member of type `Button`. We can construct a value of this type with the following expression:

```
new{i=42,new Button()}
```

To access components of anonymous structs we (again) generalize the notion of member access. Thus assuming a value `x` of the previous type, we write `x.i` to access the integer value. Unlabelled members are accessed by their position; for example `x[1]` returns the `Button` member. As for streams, member access is lifted over unlabelled members of anonymous structs. To access the `BackColor` property of the `Button` component in variable `x` we can just write `x.BackColor`, which is equivalent to `x[1].BackColor`.

At this point we can reveal even more of the power of $C\omega$'s generalized member access. Given a stream `friends` of type `struct{string name;int age;}*`, the expression `friends.age` returns a stream of integers. The member access has been lifted over *both* structural types. The following query-like statement prints the names of one's friends:

```
friends.name.{ Console.WriteLine(it);};
```

¹We shall adopt the FJ shorthand and write \bar{x} to mean a sequence of x .

Interestingly, $C\omega$ also allows repeated occurrences of the same member name within an anonymous struct type, even at different types. For example, assume the following declaration: `struct{int i; Button; float i;} z;` Then `z.i` projects the two `i` members of `z` into a new anonymous struct that is equivalent to `new{z[0], z[2]}` and of type `struct{int; float;}`.

$C\omega$ provides a limited form of covariance for anonymous structs, just as for streams. For example, the anonymous struct `struct{int; Button;}` is a subtype of `struct{int; Control;}`. However it is *not* a subtype of `struct{object; Control;}` since the conversion from `int` to `object` is not an identity conversion. $C\omega$ does not support width subtyping for anonymous structs.

Choice types The third structural type we add is a particular form of discriminated union type, which we call a choice type. This is written, for example, `choice{int; bool;}`. As the name suggests, a value of this type is either an integer or a boolean, and may hold either at any one time. Unlike unions in C/C++ and variant records in Pascal where users have to keep track of which type is present, values of an discriminated unions in $C\omega$ are implicitly tagged with the static type of the chosen alternative, much like unions in Algol68. In other words, discriminated union values are essentially a pair of a value and its static type.

There is no syntax for creating choice values; the injection is implicit (i.e. it is generated by the compiler).

```
choice{int;Button;} x = 3;
choice{int;Button;} y = new Button();
```

$C\omega$ provides a test, `e was τ` , on choice values to test the value's *static* type. Thus `x was int` would return `true`, whereas `y was int` would return `false`.

Assuming that an expression `e` is of type `choice{ $\bar{\tau}$ }`, the expression `e was τ` is true for *exactly one* τ in $\bar{\tau}$. This invariant is maintained by the type system. The only slight complication arises from subtyping, e.g.

```
choice{Control; object;} z = new Button();
```

As `Button` is a subtype of both `Control` and `object`, which type tag is generated by the compiler? The answer should be obvious to the experienced Java/C# programmer: a choice type can be thought of as providing a family of overloaded constructor methods, one for each component type. Just as for standard object creation, the *best* constructor method is chosen. In the example above, clearly `Control` is better than `object`. Thus `z was Control` returns `true`. The notion of “best” for $C\omega$ is the routine extension of that for C#.

As the reader may have guessed, member access has also been generalized over discriminated unions. Here the behaviour of member access is less obvious, and has been

designed to coincide with XPath. Consider a value `w` of type `choice{char; Button;}`. The member access `w.GetHashCode()` succeeds irrespective of whether the value is an character or a `Button` object. In this case the type of the expression `w.GetHashCode()` is `int`.

However the member may not be supported by all the possible component types, e.g. `w.BackgroundColor`. Classic treatments of union types would probably consider this to be type incorrect [18, p.207]. However, $C\omega$'s choice types follow the semantics of XPath where, for example, the query `foo/bar` returns the `bar` nodes under the `foo` node if any exist, and *the empty sequence* if none exist. Thus in $C\omega$, the expression `w.BackgroundColor` is well-typed, and will return a value of type `Color?`. This is another new type in $C\omega$ and is a variant of the nullable type to appear in C# 2.0. A value of type `Color?` can be thought of as a singleton stream, thus it is either empty (and equal to `null`), or contains a single `Color` value (when `w` contains a `Button`). Again, we emphasize that this behaviour precisely matches that of XPath.

$C\omega$ follows the design of C# in allowing all values to be boxed and hence all value types are a subtype of the super-type `object`. Thus both anonymous structs and choice types are considered to be subtypes of the class `object`.

Content classes To allow close integration with XSD and other XML schema languages, we have included the notion of a *content class* in $C\omega$. A content class is a normal class that has a single *unlabelled* type that describes the content of that class, as opposed to the more familiar (named) fields. The following is a simple example of a content class.

```
class friend{
    struct{ string name; int age; };
    void incAge(){...}
}
```

Again we have generalized member access over content classes. Thus the expression `Bill.age` returns an integer, where `Bill` is a value of type `friend`.

From an XSD perspective, classes correspond to global element declarations, while the content type of classes correspond to complex types. Further comparisons with the XML data model are immediately below, but a more comprehensive study can be found elsewhere [17].

2.2 XML programming

It should be clear that the new type structures of $C\omega$ are sufficient to model simple XML schema. For example, the following XSD schema

```
<element name="Address"><complexType><sequence>
    <choice>
```

```

    <element name="Street" type="string"/>
    <element name="POBox" type="int"/>
  </choice>
  <element name="City" type="string"/>
</sequence></complexType></element>

```

can be represented (somewhat more succinctly!) as the $C\omega$ content class declaration:

```

class Address {
  struct{
    choice{ string Street; int POBox; };
    string City;
  };
}

```

The full $C\omega$ language supports XML literals as syntactic sugar for serialized object graphs. For example, we can create an instance of the Address type from the introduction using the following literal:

```

Address a = <Address>
  <Street>13 Elm St</Street>
  <City>Hollywood</City>
</Address>;

```

The $C\omega$ compiler contains a validating XML parser that deserializes the above literal into normal constructor calls. XML literals can also contain typed holes, much as in XQuery, that allow us to embed expressions to compute part of the literal. This is especially convenient for generating streams.

The inclusion of XML literals and the semantics of the generalized member access mean that XQuery code can be almost directly written in $C\omega$. For example, consider one of XQuery Use Cases [7], that processes a bibliography file (assume that this is stored in variable bs) and for each book in the bibliography, lists the title and authors, grouped inside a result element. The suggested solution written in XQuery is as follows.

```

for $b in $bs/book
return <result>{$b/title}{$b/author}</result>

```

The $C\omega$ solution is almost identical:

```

foreach (b in bs.book)
yield return <result>{b.title}{b.author}</result>;

```

The full $C\omega$ language adds several more powerful query expressions to those discussed in this paper. For instance, filter expressions $e[e']$ return the elements in the stream e that satisfy the boolean expression e' . As labels can be duplicated in anonymous structs and discriminated unions, the full language also allows type-based selection. For example, given a value x of type $\text{struct}\{ \text{int } a; \text{struct}\{\text{string}$

$a; \}; \}$ we can select only the string member a by writing $x.\text{string}::a$.

Transitive queries are also supported in the full $C\omega$ language: the expression $e \dots \tau :: m$ selects all members m of type τ that are transitively reachable from e . Transitive queries are inspired by the XPath descendant axis.

2.3 Database programming

Relational tables are merely streams of anonymous structs. For example, the relational table created with the SQL declaration:

```
CREATE TABLE Customer (name string, custid int);
```

can be represented as the $C\omega$ declaration:

```
struct{string name; int custid}* Customer;
```

In addition to path-like queries, the full $C\omega$ language also supports familiar SQL expressions, notably `select-from-where`. For example, one of the XQuery use-cases [7] asks to list the title prices for each book that is sold by both book-sellers A and BN. Using a select statement and XML-literals, this query can be written in $C\omega$ as the following expression:

```

select
  <book-with-prices>
  <title>{a.title}</title>
  <price-A>{a.price}</price-A>
  <price-BN>{bn.price}</price-BN>
</book-with-prices>
from book a in A.book, book bn in BN.book
where a.title == bn.title

```

Note the use of XML placeholders $\{a.\text{title}\}$ and $\{bn.\text{price}\}$: when this code is evaluated new titles and new prices are computed from the bindings of the `select-from-where` clause.

3 The essence of $C\omega$

In the rest of this paper we study formally the essence of $C\omega$, by which we mean we identify its essential features. We adopt a formal, mathematical approach and define a core calculus, Featherweight $C\omega$, or $FC\omega$ for short, similar to core subsets of Java such as FJ [16], MJ [4] and ClassicJava [9]. This core calculus, whilst lightweight, offers a similar computational “feel” to the full $C\omega$ language: it supports the new type constructors and generalized member access. $FC\omega$ is a completely valid subset of $C\omega$ in that every $FC\omega$ program is literally an executable $C\omega$ program.

The rest of this section is organized as follows. In §3.1 we define the syntax and type system for $FC\omega$. Rather than

give an operational semantics directly for $FC\omega$ we prefer to first “compile out” some of its features, in particular generalized member access. This both greatly simplifies the resulting operational semantics and demonstrates that $C\omega$'s features do not require extensive new machinery. Thus in §3.2 we define a target language, Inner $C\omega$, or $IC\omega$, for this “compilation”. $IC\omega$ is essentially the same language, but for a handful of new language constructs and a dramatically simpler type system. In §3.3 we give an operational semantics directly on $IC\omega$ programs. In §3.4 we specify the compilation of $FC\omega$ programs into $IC\omega$ programs. This translation is, on the whole, quite straightforward. We conclude the section in §3.5 by stating some properties of our calculi. Most important is the type soundness property for $FC\omega$. Space prevents us from providing any details of the proofs, but they are proved using standard techniques and are similar to analogous theorems for Java [16, 4].

3.1 A core calculus: $FC\omega$

Syntax An $FC\omega$ program consists of one or more class declarations. Each class declaration defines zero or more methods and contains exactly one unlabelled type that we call the *content type*. (We can code up a class declaration with a number of field declarations using an anonymous struct.) $FC\omega$ follows C^\sharp and requires methods to be explicitly marked as `virtual` or `override`. Given a program we assume that there is a unique designated method within the class declarations that serves as the entry point.

Program

$$p ::= \overline{cd}$$

Class Definition

$$cd ::= \text{class } c : c \{ \tau ; \overline{md} \}$$

Method Definition

$$md ::= \text{virtual } \tau m(\overline{\tau} \overline{x}) \{ \overline{s} \} \\ | \text{override } \tau m(\overline{\tau} \overline{x}) \{ \overline{s} \}$$

$FC\omega$ supports two main kinds of types: *value types* and *reference types*. As usual, the distinguished type `void` is used for methods that do not return anything; `null` is only used to type null references, as with C^\sharp . Value types include the base types `bool` and `int` and the structural types: anonymous structs and discriminated unions. Reference types are either class types or streams. As usual only reference types have object identity and are represented at runtime by references into the heap. We assume a designated special class object.

Types

$\tau ::=$	γ	Value types
	ρ	Reference types
	<code>void</code> <code>null</code>	Void and null types

Value Types

$\gamma ::=$	b	Base types
	<code>struct</code> { \overline{fd} }	Anonymous structs
	<code>choice</code> { $\overline{\tau}_c$ }	Disjoint union types

Base Types

$$b ::= \text{bool} \mid \text{int}$$

Reference Types

$\rho ::=$	c	Classes
	τ_s^*	Stream types
	$\tau_n^?$	Singleton stream type

Field Definition

$fd ::=$	$\tau f;$	Named member
	$\tau;$	Unnamed member

(We employ the shorthand τ_c, τ_s, τ_n to denote any type *except* a choice type, stream type, singleton stream type, respectively.) As $C\omega$ flattens stream types, we have made the simplification to $FC\omega$ of removing nested stream types altogether from the type grammar. We have also simplified $FC\omega$ choice types so that the members are unlabelled and we also exclude (for simplification) nested choice types. These can be coded up in $FC\omega$ using unlabelled anonymous structs.

$FC\omega$ expressions, as for C^\sharp , are split into ordinary expressions and promotable expressions. Promotable expressions are expressions that can be used as statements. We assume a number of built-in primitive operators, such as `==`, `||` and `&&`. In the grammar we write $e \oplus e$, where \oplus denotes an instance of one of these operators. We do not formalize these operators further as their meaning is clear.

Expression

$e ::=$	$b \mid i$	Literals
	$e \oplus e$	Built-in operators
	x	Variable
	<code>null</code>	Null
	$(\tau) e$	Cast
	$e \text{ is } \tau$	Dynamic typecheck
	$e \text{ was } \tau$	Static typecheck for choice values
	<code>new</code> $\tau(e)$	Object creation
	<code>new</code> { \overline{be} }	Anonymous struct creation
	$e.f$	Field access
	$e[i]$	Field access by position
	pe	Promotable expression

Promotable expression

$pe ::=$	$x = e$	Variable assignment
	$e.m(\overline{e})$	Method invocation
	$e.\{e\}$	Apply-to-all

Binding expression

$be ::=$	$f = e$	Named binding
	e	Unnamed binding

We have made a simplification in the interests of space to restrict apply-to-all expressions to contain an expression, rather than a sequence of statements. This simplifies the typing rules, but as apply-to-all expressions can be coded using `foreach` loops, it is not a serious restriction.

Statements in $FC\omega$ are standard. As mentioned earlier we

have adopted the `yield` statement that will appear in $C^\#$ 2.0 to generate streams.

Statement

$s ::=$	<code>;</code>	Skip
	<code>pe;</code>	Promoted expression
	<code>if (e) s else s</code>	Conditional
	<code>τ x = e;</code>	Variable declaration
	<code>return e;</code>	Return statement
	<code>return;</code>	Empty return
	<code>yield return e;</code>	Yield statement
	<code>yield break;</code>	End of stream
	<code>foreach (τ x in e) s</code>	Foreach loop
	<code>while (e) s</code>	While loop
	<code>{\bar{s}}</code>	Block

In what follows we assume that $FC\omega$ programs are well-formed, e.g. no cyclic class hierarchies, correct method body construction, etc. These conditions can be easily formalized but we surpress such detail for lack of space.

Subtyping Before we define the typing judgements for $FC\omega$ programs we first need to define a number of auxiliary relations. First we define the subtyping relation. We write $\tau <: \tau'$ to mean that type τ is a subtype of type τ' . The rules defining this relation are given in Figure 1.

Most of these rules are straightforward. The rules **[Stream]** and **[Struct]** make use of a predicate *covarOK*, which captures the restricted form of covariance described earlier. In this short paper we shall not give its straightforward definition.

Generalized member access As we have seen a key programming feature of $C\omega$ is generalized member access. Capturing this behaviour in the type system can be tricky, but we have adopted a rather elegant solution, whereby we define two auxiliary relations. The first, written $\tau.f : \tau'$, tells us that given a value of type τ accessing member f will return a value of type τ' . We define a similar relation for function member access, written $\tau.m(\bar{\tau}') : \tau''$. Having generalized member access captured by a separate typing relation greatly simplifies the typing judgements for expressions. As it is a core feature of $C\omega$, we shall give it in detail.

The definition of this relation over stream types is as follows.

$$\frac{\tau.f : \tau'_s}{\tau*.f : \tau'_s*} \quad \frac{\tau.f : \tau' *}{\tau*.f : \tau' *}$$

$$\frac{\tau.m(\bar{\tau}') : \tau''_s}{\tau*.m(\bar{\tau}') : \tau''_s*} \quad \frac{\tau.m(\bar{\tau}') : \tau'' *}{\tau*.m(\bar{\tau}') : \tau'' *} \quad \frac{\tau.m(\bar{\tau}') : \text{void}}{\tau*.m(\bar{\tau}') : \text{void}}$$

The first two rules map the member access over the stream elements, making sure that we do not create a nested stream type. The next two rules for function member access are

similar. The last rule captures the intuition that mapping a void-valued method over a stream, forces the evaluation of the stream and does not return a value.

Before defining the rules for member access over anonymous structs, we need to define rules for member access over named field definitions. This is pretty straightforward and as follows.

$$\frac{}{\tau.f.f : \tau} \quad \frac{\tau.m(\bar{\tau}') : \tau''}{\tau.f.m(\bar{\tau}') : \tau''}$$

Now we consider the rules for generalized member access over anonymous structs. First we give the degenerate cases where only one component supports the member access.

$$\frac{\exists!k \in \{1 \dots n\}. fd_k.f : \tau_k}{\text{struct}\{fd_1; \dots fd_n; \}.f : \tau_k} \quad \frac{\exists!k \in \{1 \dots n\}. fd_k.m(\bar{\tau}') : \tau''}{\text{struct}\{fd_1; \dots fd_n; \}.m(\bar{\tau}') : \tau''}$$

The non-degenerate cases are then as follows.

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. fd_{S_k}.f : \tau_k}{\text{struct}\{fd_1; \dots fd_n; \}.f : \text{struct}\{\tau_1; \dots \tau_p; \}}$$

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. fd_{S_k}.m(\bar{\tau}') : \tau'_k}{\text{struct}\{fd_1; \dots fd_n; \}.m(\bar{\tau}') : \text{struct}\{\tau'_1; \dots \tau'_p; \}}$$

Thus a subset, S , of the components support the member, and we map the member access over these components in order. The overall return type is an anonymous struct of the component return types.

We now consider the rules for generalized member access over choice types. Again we consider these rules depending on how many components support the member access. First we give the simple case when *all* possible components support the member access.

$$\frac{\forall k \in \{1 \dots n\}. \tau_k.f : \tau}{\text{choice}\{\tau_1; \dots \tau_n; \}.f : \tau} \quad \frac{\forall k \in \{1 \dots n\}. \tau_k.m(\bar{\tau}') : \tau}{\text{choice}\{\tau_1; \dots \tau_n; \}.m(\bar{\tau}') : \tau}$$

We also have the case when only one of the possible components supports the member access. These rules are as follows.

$$\frac{\exists!k \in \{1 \dots n\}. \tau_k.f : \tau \quad n > 1}{\text{choice}\{\tau_1; \dots \tau_n; \}.f : \tau?}$$

$$\frac{\exists!k \in \{1 \dots n\}. \tau_k.m(\bar{\tau}') : \tau \quad n > 1}{\text{choice}\{\tau_1; \dots \tau_n; \}.m(\bar{\tau}') : \tau?}$$

The reader will recall that the return of this generalized member access involves a singleton stream type. Finally we gives the cases where more than one of the possible components supports the member access.

$$\begin{array}{c}
\frac{}{\tau <: \tau} \text{[Ref]} \quad \frac{\tau <: \tau' \quad \tau' <: \tau''}{\tau <: \tau''} \text{[Trans]} \quad \frac{}{\gamma <: \text{object}} \text{[Box]} \quad \frac{\text{class } c : c'}{c <: c'} \text{[Sub]} \quad \frac{}{\text{null} <: \rho} \text{[Null]} \quad \frac{\tau <: \tau' \quad f = f'}{\tau f <: \tau' f'} \text{[FD]} \\
\\
\frac{\tau <: \tau' \quad \text{covarOK}(\tau, \tau')}{\tau * <: \tau' *} \text{[Stream]} \quad \frac{}{\tau * <: \text{object}} \text{[SBox]} \quad \frac{}{\tau ? <: \tau *} \text{[SSub]} \quad \frac{}{\tau_n <: \tau_n ?} \text{[Sing]} \quad \frac{}{\text{null} <: \tau_n ?} \text{[NullSing]} \\
\\
\frac{\overline{fd} <: \overline{fd'} \quad \text{covarOK}(\overline{fd}, \overline{fd'})}{\text{struct}\{\overline{fd}\} <: \text{struct}\{\overline{fd'}\}} \text{[Struct]} \quad \frac{\tau \neq \text{choice}\{\overline{\tau''}\}}{\tau <: \text{choice}\{\tau; \overline{\tau''}\}} \text{[SubChoice]} \quad \frac{}{\text{choice}\{\overline{\tau}\} <: \text{choice}\{\overline{\tau} \overline{\tau'}\}} \text{[Choice]}
\end{array}$$

Figure 1: FC ω subtyping

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. \tau_{S_k}.f : \tau'_k}{\text{choice}\{\tau_1; \dots \tau_n\}.f : \text{choice}\{\tau'_1; \dots \tau'_p\} ?}$$

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. \tau_{S_k}.m(\overline{\tau'}) : \tau''_k}{\text{choice}\{\tau_1; \dots \tau_n\}.m(\overline{\tau'}) : \text{choice}\{\tau''_0; \dots \tau''_p\} ?}$$

Generalized member access over singleton streams is relatively straightforward; the only complication being again to ensure that no nested streams are generated.

$$\frac{\tau.f : \tau'_n \quad \tau.f : \tau' ?}{\tau?.f : \tau'_n ?} \quad \frac{\tau.m(\overline{\tau'}) : \tau''_n \quad \tau.m(\overline{\tau'}) : \tau'' ?}{\tau?.m(\overline{\tau'}) : \tau''_n ?}$$

Finally we need to define rules for generalized member access over classes. Clearly these need to reflect the standard C \sharp semantics: function member access on classes searches the class hierarchy until a matching method is found. If we find a matching method $\tau'.m(\overline{\tau''})$ in class c , we need to check the actual types of the arguments to the types expected by m . This behaviour is given by the following two rules.

$$\frac{\text{class } c : c'\{\tau; \overline{md}\} \quad \tau'.m(\overline{\tau''}) \in \overline{md} \quad \overline{\tau} <: \overline{\tau''}}{c.m(\overline{\tau}) : \tau'} \\
\frac{\text{class } c : c'\{\tau; \overline{md}\} \quad \tau'.m(\overline{\tau''}) \notin \overline{md} \quad c'.m(\overline{\tau}) : \tau'}{c.m(\overline{\tau}) : \tau'}$$

Next we consider the rules for generalized field access. There is a small subtlety here concerning recursive class definitions; consider the following recursive class List of lists of integers.

```
class List { struct{ int head; List; } }
```

Given an instance xs of type List, we do not want $xs.head$ to recursively select all head fields in xs . However simply unfolding the content type and using the rules given earlier for generalized access over anonymous structs that is precisely what would happen!

There are a number of solutions, but in our determination to make the C ω type system as simple as possible, we follow e.g. Haskell and SML and break recursive cycles at nominal types. In our setting that means that we simply do not perform member lookup on nominal members of the content of nominal types. Using these refined rules, the result type of $xs.head$ is int .

Formalizing this is trivial but time-consuming. We define another family of generalized member access judgements, written $\tau \bullet f : \tau'$, which is identical to the previous rules except they are not defined for nominal types. We eschew the definitions here.

To define field access on nominal types, we first define formally the content type of a class, written $\text{content}(c)$ for some class c , as follows.

$$\frac{\text{class } c : \text{object}\{\tau; \overline{md}\}}{\text{content}(c) = \tau} \quad \frac{\text{class } c : c'\{\tau; \overline{md}\}}{\text{content}(c) = \text{struct}\{\tau'; \tau; \}}$$

The rule for generalized member access over classes simply searches for the member f on the content type of class c , and is given by the following rule.

$$\frac{\text{content}(c) = \tau \quad \tau \bullet f : \tau'}{c.f : \tau'}$$

Generalized index access As we mentioned earlier, elements of anonymous structs can be accessed by position. This is captured by the following rule.

$$\frac{\text{type}(fd_i) = \tau_i}{\text{struct}\{fd_1; \dots fd_n\}[i] : \tau_i}$$

As the reader might have expected, this index access is generalized over the other types. In the interests of space, we elide the rather routine details.

Typing judgements We are now able to define typing judgements for $FC\omega$. We define three relations corresponding to the three syntactic categories of expressions, promotable expressions and statements. For all three judgements we write Γ to mean a partial function from program identifiers to types. The judgements for expressions and promotable expressions are written $\Gamma \vdash e : \tau$ and $\Gamma \vdash pe : \tau$, respectively. These are given in Figure 2.

Most of these rules are routine; we shall discuss a few of the more interesting details here. In the rule [TStruct], we have made use of a typing judgement for a binding expression. This is defined as follows:

$$\frac{\Gamma \vdash e : \tau}{\Gamma \vdash f = e : \tau f}$$

The compactness of the rules [TField], [TIndex] and [TMeth] shows the elegance of having captured generalized member access with auxiliary relations.

The rules [TAAExp1] and [TAAExp2] ensure that the return type of apply-to-all expressions are not nested. The rule [TAAExp3] captures the intuition that applying a void-typed expression to a stream forces the evaluation of that stream and hence the overall type is also void.

The typing judgement for $FC\omega$ statements is written $\Gamma; \tau \vdash s$ and is intended to mean that a statement s is well-typed in the typing environment Γ . If it returns a value (either via a normal return or a `yield return`) then that value has to be of type τ .

The rules [TForEach1] and [TForEach2] reflect the fact that the type of the stream elements can be cast to the type of the bound variable. This can be either via an upcast ([TForEach1]) or a downcast ([TForEach2]).

3.2 An inner calculus: $IC\omega$

Rather than consider further our featherweight calculus $FC\omega$, we shall in fact define another core calculus for $C\omega$. This inner calculus, called $IC\omega$, is intended to be similar but lower-level than $FC\omega$; it can be thought of as the internal language of a compiler.

The chief simplification in $IC\omega$ is that its type system does *not* support generalized member access. The intention is that we compile out generalized member access when translating $FC\omega$ programs into $IC\omega$ programs. We give some details of this compilation in §3.4. Apart from a simplified type system, we can define quite simply an operational semantics for $IC\omega$; this is given in §3.3.

The grammar of $IC\omega$ is then a modest extension of the grammar for $FC\omega$. Some extra expression and statement forms are added (which reflects the lower-level nature of $IC\omega$) and

likewise a couple are removed from the grammar as they are redundant. We do not expect these new syntactic forms to be made available to the $C\omega$ programmer (although they could be). The extensions are as follows:

Expression		
$e ::=$...	
	<code>new $\tau(\bar{s})$</code>	Closure creation
	<code>new (τ, e)</code>	Choice creation
	<code>e.Content</code>	Class content
	<code>e at τ</code>	Choice content
	<code>new $\tau?(e)$</code>	Nullable creation
	<code>e.Value</code>	Nullable content
	<code>e.HasValue</code>	Nullable tag
Promotable expression		
$pe ::=$...	
	<code>$\tau(\{\bar{s}\})$</code>	Block expression
Statement		
$s ::=$...	
	<code>yield return (τ, e)</code>	Typed yield

Thus $IC\omega$ includes expressions to create closure, choice and nullable elements. We include an operator `e .Content` to extract the content element from an object e . Given an element e of a choice type, we add an operation `e at τ` to extract its τ -valued content. (If it is of another type, this will raise an exception.) We also use two members `Value` and `HasValue` for nullable values; these will appear in $C\omega$ 2.0. The expression `e .HasValue` returns a boolean depending on whether the expression e is non-null. The expression `e .Value` returns the value of expression e if it is non-null and raises an exception it is null. We add (typed) block expressions to $IC\omega$, and in addition we provide a typed `yield` statement.

The two syntactic forms that we removed from the grammar of $FC\omega$ are: (1) We remove field accesses `$e.f$` completely; they are replaced by positional access, i.e. `$e[i]$` ; and (2) We remove the untyped `yield` statement; all `yields` in $IC\omega$ are explicitly typed.

We can define typing judgements for $IC\omega$ expressions and statements, which are written $\Gamma \triangleright e : \tau$ and $\Gamma; \tau \triangleright s$, respectively. Most of these rules are identical to those for $FC\omega$; we shall just give the rules for the new syntactic forms. The rules for creating closure, choice and nullable elements are as follows:

$$\frac{\Gamma; \tau * \triangleright \bar{s}}{\Gamma \triangleright \text{new } \tau * (\bar{s}) : \tau *}$$

$$\frac{\Gamma \triangleright e : \tau'_c \quad \tau'_c <: \tau}{\Gamma \triangleright \text{new } (\tau, e) : \text{choice}\{\tau; \}}$$

$$\frac{\Gamma \triangleright e : \text{null}}{\Gamma \triangleright \text{new } \tau?(e) : \tau?}$$

$$\frac{\Gamma \triangleright e : \tau_n}{\Gamma \triangleright \text{new } \tau_n?(e) : \tau_n?}$$

The typing rules for extracting the content of content class and choice elements are as follows:

$$\frac{\Gamma \triangleright e : c}{\Gamma \triangleright e.\text{Content} : \text{content}(c)}$$

$$\frac{\Gamma \triangleright e : \text{choice}\{\tau; \bar{\tau}'\}}{\Gamma \triangleright e \text{ at } \tau : \tau}$$

The typing rule for block expressions is as follows:

$\Gamma \vdash e : \tau$

$$\begin{array}{c}
\frac{}{\Gamma \vdash i : \text{int}} [\text{TInt}] \quad \frac{}{\Gamma \vdash b : \text{bool}} [\text{TBool}] \quad \frac{}{\Gamma, x : \tau \vdash x : \tau} [\text{TId}] \quad \frac{}{\Gamma \vdash \text{null} : \text{null}} [\text{TNull}] \quad \frac{\Gamma \vdash e : \tau' \quad (\tau' <: \tau) \vee (\tau <: \tau')}{\Gamma \vdash (\tau) e : \tau} [\text{TSub}] \\
\frac{\Gamma \vdash e : \tau' \quad (\tau' <: \tau) \vee (\tau <: \tau')}{\Gamma \vdash e \text{ is } \tau : \text{bool}} [\text{TIIs}] \quad \frac{\Gamma \vdash e : \text{choice}\{\overline{\tau'} \tau; \overline{\tau''}\}}{\Gamma \vdash e \text{ was } \tau : \text{bool}} [\text{TWas}] \quad \frac{\Gamma \vdash \overline{be} : \overline{fd}}{\Gamma \vdash \text{new}\{\overline{be}\} : \text{struct}\{\overline{fd}\}} [\text{TStruct}] \\
\frac{\Gamma \vdash e : \tau \quad \tau <: \text{content}(c)}{\Gamma \vdash \text{new } c(e) : c} [\text{TNew}] \quad \frac{\Gamma \vdash e : \tau \quad \tau.f : \tau'}{\Gamma \vdash e.f : \tau'} [\text{TField}] \quad \frac{\Gamma \vdash e : \tau \quad \tau[i] : \tau'}{\Gamma \vdash e[i] : \tau'} [\text{TIndex}]
\end{array}$$

$\Gamma \vdash pe : \tau$

$$\begin{array}{c}
\frac{\Gamma \vdash x : \tau \quad \Gamma \vdash e : \tau' \quad \tau' <: \tau}{\Gamma \vdash x = e : \tau} [\text{TAss}] \quad \frac{\Gamma \vdash e : \tau \quad \Gamma \vdash \overline{e'} : \overline{\tau'} \quad \tau.m(\overline{\tau'}) : \tau''}{\Gamma \vdash e.m(\overline{e'}) : \tau''} [\text{TMeth}] \quad \frac{\Gamma \vdash e : \tau^* \quad \Gamma, \text{it} : \tau \vdash e' : \tau'_s}{\Gamma \vdash e.\{e'\} : \tau'_s^*} [\text{TAAExp1}] \\
\frac{\Gamma \vdash e : \tau^* \quad \Gamma, \text{it} : \tau \vdash e' : \tau'_s^*}{\Gamma \vdash e.\{e'\} : \tau'_s^*} [\text{TAAExp2}] \quad \frac{\Gamma \vdash e : \tau^* \quad \Gamma, \text{it} : \tau \vdash e' : \text{void}}{\Gamma \vdash e.\{e'\} : \text{void}} [\text{TAAExp3}]
\end{array}$$

$\Gamma; \tau \vdash s$

$$\begin{array}{c}
\frac{}{\Gamma; \tau \vdash ;} [\text{TSkip}] \quad \frac{\Gamma; \tau \vdash \overline{s}}{\Gamma; \tau \vdash \{\overline{s}\}} [\text{TNest}] \quad \frac{\Gamma \vdash pe : \tau}{\Gamma; \tau' \vdash pe;} [\text{TProm}] \quad \frac{\Gamma \vdash e : \text{bool} \quad \Gamma; \tau \vdash s_1 \quad \Gamma; \tau \vdash s_2}{\Gamma; \tau \vdash \text{if } (e) s_1 \text{ else } s_2} [\text{TIf}] \\
\frac{\Gamma \vdash e : \text{bool} \quad \Gamma; \tau \vdash s}{\Gamma; \tau \vdash \text{while } (e) s} [\text{TWhile}] \quad \frac{}{\Gamma; \text{void} \vdash \text{return};} [\text{TRetV}] \quad \frac{\Gamma \vdash e : \tau \quad \tau <: \tau}{\Gamma; \tau \vdash \text{return } e;} [\text{TRet}] \\
\frac{}{\Gamma; \tau^* \vdash \text{yield break};} [\text{TYieldB}] \quad \frac{\Gamma \vdash e : \tau'_s \quad \tau' <: \tau}{\Gamma; \tau^* \vdash \text{yield return } e;} [\text{TYield1}] \quad \frac{\Gamma \vdash e : \tau'^* \quad \tau' <: \tau}{\Gamma; \tau^* \vdash \text{yield return } e;} [\text{TYield2}] \\
\frac{\Gamma \vdash e : \tau'^* \quad \tau' <: \tau'' \quad \Gamma, x : \tau''; \tau \vdash s}{\Gamma; \tau \vdash \text{foreach } (\tau'' x \text{ in } e) s} [\text{TForEach1}] \quad \frac{\Gamma \vdash e : \tau'^* \quad \tau'' <: \tau' \quad \Gamma, x : \tau''; \tau \vdash s}{\Gamma; \tau \vdash \text{foreach } (\tau'' x \text{ in } e) s} [\text{TForEach2}]
\end{array}$$

Figure 2: Typing judgements for $\text{FC}\omega$ expressions, statements, class declarations, method definitions

$$\frac{\Gamma; \tau_s \triangleright \overline{s} \quad \tau \neq \text{void}}{\Gamma \triangleright \tau_s(\{\overline{s}\}) : \tau_s}$$

The typing rules for the typed yield statement are as follows:

$$\frac{\Gamma \triangleright e : \tau'_s \quad \tau' <: \tau}{\Gamma; \tau^* \triangleright \text{yield return } (\tau'_s, e);} \quad \frac{\Gamma \triangleright e : \tau'^* \quad \tau' <: \tau}{\Gamma; \tau^* \triangleright \text{yield return } (\tau'^*, e);}$$

3.3 Operational semantics for $\text{IC}\omega$

In this section we formalize the dynamics of $\text{IC}\omega$ by defining an operational semantics. We follow FJ [16] and MJ [4] and give this in the form of a small-step reduction relation, although a big-step evaluation relation can easily be defined.

First we define the value forms of $\text{IC}\omega$ expressions and statements (where bv is the value form of a binding expression):

Expression values

$v ::=$	$b \mid i \mid \text{null} \mid \text{void}$	Basic values
	r	Reference
	$\text{new } \{\overline{bv}\}$	Struct value
	$\text{new } (\tau, v)$	Choice value
	$\text{new } \tau?(v)$	Nullable value

Statement values

$sv ::=$	$;$	Skip
	$\text{return } v;$	Return value
	$\text{return};$	
	$\text{yield return } (\tau, v);$	Typed yield value
	$\text{yield break};$	End of stream value

Evaluation of $\text{IC}\omega$ expressions and statements takes place in the context of a state, which is a pair (H, R) , where H is a heap and R is a stack frame. A heap is represented as a

finite partial map from references r to runtime objects, and a stack frame is a finite partial map from variable identifiers to values. A runtime object, as for \mathbb{C}^\sharp , is a pair (τ, cn) where τ is a type and cn is a canonical, which is either a value or a closure. A closure is the runtime representation of a stream and is written as a pair (R, \bar{s}) where R is a stack frame and \bar{s} is a statement sequence. In what follows we assume that expressions and statements are well-typed.

We give the complete set of reduction rules in Figure 3. We use evaluation contexts to encode the evaluation strategy in the familiar way [9]—the definition of evaluation contexts is routine and omitted for space.

Figure 3 defines evaluation relations for the three syntactic classes of expressions, promotable expressions and statements. We consider each of these in turn.

The evaluation relation for $\text{IC}\omega$ expressions is written $S, e \rightarrow S', e'$ which means that given a state S , expression e reduces by one or possibly more steps to e' and a (possibly updated) state S' . (We use an auxiliary function *value* defined as follows: $\text{value}(f = v) \stackrel{\text{def}}{=} v$, $\text{value}(v) \stackrel{\text{def}}{=} v$.) These rules are routine.

The evaluation relation for $\text{IC}\omega$ promotable expressions is written $S, pe \rightarrow S', pe'$. The rules for method invocation deserve some explanation: they are differentiated according to whether the method is `void`-returning. If it is not then the method body is unfolded, and executed until it is of the form `return v`; where v is a value. This value is then the result of the method invocation. If the method is `void`-valued, then we unfold the method body and execute it until it is of the form `return ;`. The result is the special value `null`.

The evaluation relation for statements is written $S, s \rightarrow S', s'$. Most of these rules are standard. The last two rules are most interesting as they embody the flattening of streams. To evaluate a `foreach` loop we first evaluate the stream until it yields a value. If that value is itself a stream, then we should first execute the `foreach` loop on this stream.

As is usual we have a number of cases that lead to a predictable error state, e.g. following a dereference of a `null` object. These errors in $\text{IC}\omega$ are *CastX*, *ChoiceX*, *NullX* and *NbleX*. We say that a pair S, e is *terminal* if e is one of these errors, or it is a value.

3.4 Compiling $\text{FC}\omega$ to $\text{IC}\omega$

In this section we give some details of the compilation of $\text{FC}\omega$ into $\text{IC}\omega$. Much of this compilation is routine, so in the interests of space we shall concentrate only on the most interesting aspect: generalized member access.

We employ a “coercion” technique, in that we translate the *implicit* generalized member access of $\text{FC}\omega$ into an *ex-*

PLICIT $\text{IC}\omega$ code fragment. In Figure 4 we give the complete compilation of generalized member access (GMA). This is expressed as an inductively defined relation, written $|\tau.f:\tau'| \rightsquigarrow g$ and $|\tau.m(\bar{\tau}'):\tau''| \rightsquigarrow g$ for member and function member access respectively. A judgement $|\tau.f:\tau'| \rightsquigarrow g$ is intended to mean that if invoking a member f on an element of type τ returns an element of type τ' , then g is the $\text{IC}\omega$ coercion that will explicitly access the appropriate member. In the definition we have employed a function-like syntax for coercions, although they are really contexts.

We can compile an instance of member access in $\text{FC}\omega$, $e.f$, as follows: we first compile the expression e into $\text{IC}\omega$, yielding e' , and also generate a coercion, g , corresponding to the member access. The result of the compilation of $e.f$ is then simply $g(e')$. We write the compilation of, e.g. an expression, e , as $|\Gamma \vdash e:\tau| \rightsquigarrow e'$.

3.4.1 Incoherence by design

Java and \mathbb{C}^\sharp are by design incoherent. Both languages use a notion of “best” conversion when there is more than one conversion between two types. If there does not exist a best conversion, a compile-time error is generated. In compiling $\text{FC}\omega$ to $\text{IC}\omega$ we use this notion of a best conversion when dealing with rules that use subtyping. We do not formalize this notion of “best” here; both the Java and \mathbb{C}^\sharp language specifications give precise details. The new types in $\mathbb{C}\omega$ do not complicate this notion greatly: For example, there are two conversions between `int` and `object`: one using the rule [Box], the other using the rules [SubChoice] and [Box] along with [Trans] (i.e. `int <: choice{int;string;} <: object`). It is clear that the first conversion is better. The other critical pairs are similarly easy to resolve.

3.5 Properties of $\text{FC}\omega$ and $\text{IC}\omega$

In this section we briefly mention some properties of $\text{FC}\omega$ and $\text{IC}\omega$, including type soundness. We do not give any details of the proofs, as they are standard and follow analogous theorems for Java [16, 4].

First we formalize that $\text{IC}\omega$ is type-sound, which is captured by the following properties. (We use generalized judgements, e.g. $\Gamma \triangleright (S, e):\tau$ to mean that the expression e is well-typed and also that the state S is well-formed with respect to Γ , in the familiar way. As is usual [16] we need to add “stupid” rules to the type system for the formal proof.)

Theorem 1 (Type soundness for $\text{IC}\omega$)

1. If $\Gamma \triangleright (S, e/pe):\tau$ and $(S, e/pe) \rightarrow (S', e'/pe')$ then $\exists \tau'$ such that $\Gamma \triangleright (S', e'):\tau'$.

Expressions:

$$\begin{array}{c}
\frac{}{(H, R), x \rightarrow (H, R), R(x)} \quad \frac{H(r) = (\tau', cn) \quad \tau <: \tau'}{(H, R), (\tau)r \rightarrow (H, R), r} \quad \frac{H(r) = (\tau', cn) \quad \neg(\tau <: \tau')}{(H, R), (\tau)r \rightarrow (H, R), CastX} \\
\frac{H(r) = (\tau', cn) \quad \tau' <: \tau}{(H, R), r \text{ is } \tau \rightarrow (H, R), \text{true}} \quad \frac{H(r) = (\tau', cn) \quad \tau' \not<: \tau}{(H, R), r \text{ is } \tau \rightarrow (H, R), \text{false}} \quad \frac{}{S, \text{new } (\tau, v) \text{ was } \tau \rightarrow S, \text{true}} \\
\frac{\tau \neq \tau''}{S, \text{new } (\tau, v) \text{ was } \tau'' \rightarrow S, \text{false}} \quad \frac{r \notin \text{dom}(H)}{(H, R), \text{new } \tau(v) \rightarrow (H[r \mapsto (\tau, v)], R), r} \quad \frac{r \notin \text{dom}(H)}{(H, R), \text{new } \tau(\bar{s}) \rightarrow (H[r \mapsto (\tau, (R, \bar{s})]), R), r} \\
\frac{H(r) = (c, cn)}{(H, R), r.\text{content} \rightarrow (H, R), cn} \quad \frac{0 \leq i \leq n}{S, \text{new } \{bv_0, \dots, bv_n\}[i] \rightarrow S, \text{value}(bv_i)} \quad \frac{}{S, \text{new } (\tau, v) \text{ at } \tau \rightarrow S, v} \\
\frac{\tau \neq \tau'}{S, \text{new } (\tau, v) \text{ at } \tau' \rightarrow S, \text{ChoiceX}} \quad \frac{}{S, \text{null}.\text{content} \rightarrow S, \text{NullX}} \quad \frac{}{S, \text{new } \gamma?(\text{null}).\text{HasValue} \rightarrow S, \text{false}} \\
\frac{v \neq \text{null}}{S, \text{new } \gamma?(v).\text{HasValue} \rightarrow S, \text{true}} \quad \frac{v \neq \text{null}}{S, \text{new } \gamma?(v).\text{Value} \rightarrow S, v} \quad \frac{}{S, \text{new } \gamma?(\text{null}).\text{Value} \rightarrow S, \text{NbleX}}
\end{array}$$

Promotable expressions:

$$\begin{array}{c}
\frac{}{(H, R), x = v \rightarrow (H, R[x \mapsto v]), v} \quad \frac{}{S, \text{null}.\text{m}(\bar{v}) \rightarrow S, \text{NullX}} \quad \frac{(H, R), \bar{s} \rightarrow^* (H', R'), \text{return } v; \bar{s}'}{(H, R), \tau(\{\bar{s}\}) \rightarrow (H', R'), v} \\
\frac{H(r) = (c, _) \quad \text{method}(m, c) = \tau'(\bar{\tau} \bar{x})\{\bar{s}\} \quad \tau' \neq \text{void}}{(H, []), \{c \text{ this } = r; \bar{\tau} \bar{x} = \bar{v}; \bar{s}\} \rightarrow^* (H', R'), \text{return } v'; \bar{s}'} \quad \frac{H(r) = (c, _) \quad \text{method}(m, c) = \text{void}(\bar{\tau} \bar{x})\{\bar{s}\}}{(H, []), \{c \text{ this } = r; \bar{\tau} \bar{x} = \bar{v}; \bar{s}\} \rightarrow^* (H', R'), \text{return } ; \bar{s}'} \\
\frac{}{(H, R), r.\text{m}(\bar{v}) \rightarrow (H', R), v'} \quad \frac{}{(H, R), r.\text{m}(\bar{v}) \rightarrow (H', R), \text{void}}
\end{array}$$

Statements:

$$\begin{array}{c}
\frac{}{S, \text{if } (\text{true}) s_1 \text{ else } s_2 \rightarrow S, s_1} \quad \frac{}{S, \text{if } (\text{false}) s_1 \text{ else } s_2 \rightarrow S, s_2} \quad \frac{}{S, \text{while } (e)\{\bar{s}\} \rightarrow S, \text{if } (e)\{\bar{s} \text{ while } (e)\{\bar{s}\}\} \text{ else } \{; \}} \\
\frac{}{(H, R), \tau x = v; \rightarrow (H, R[x \mapsto v]), ;} \quad \frac{}{S, \{\bar{s}\} \rightarrow S, \bar{s}} \quad \frac{}{S, v; \rightarrow S, ;} \quad \frac{}{S, \text{foreach } (\tau x \text{ in null}) s \rightarrow S, ;} \\
\frac{H(r) = (\tau', (R', \bar{s}'))}{(H, R'), \bar{s}' \rightarrow^* (H', R''), \text{yield break } ; \bar{s}''} \quad \frac{H(r) = (\tau', (R', \bar{s}'))}{(H, R'), \bar{s}' \rightarrow^* (H', R''), \text{yield return } (\tau, \text{null}); \bar{s}''} \\
\frac{}{(H, R), \text{foreach } (\tau x \text{ in } r) s \rightarrow (H'[r \mapsto (\tau', (R'', ;))], R), ;} \quad \frac{}{(H, R), \text{foreach } (\tau x \text{ in } r) s \rightarrow (H'[r \mapsto (\tau', (R'', \bar{s}''))], R), \text{foreach } (\tau x \text{ in } r) s} \\
\frac{H(r) = (\tau', (R', \bar{s}'))}{(H, R'), \bar{s}' \rightarrow^* (H', R''), \text{yield return } (\tau_s'', v); \bar{s}''} \quad \frac{}{v \neq \text{null}} \\
\frac{}{(H, R), \text{foreach } (\tau x \text{ in } r) s \rightarrow (H'[r \mapsto (\tau', (R'', \bar{s}''))], R), \{\{\tau x = v; s\} \text{foreach } (\tau x \text{ in } r) s\}} \\
\frac{H(r) = (\tau', (R', \bar{s}'))}{(H, R'), \bar{s}' \rightarrow^* (H', R''), \text{yield return } (\tau''*, v); \bar{s}''} \quad \frac{}{v \neq \text{null}} \\
\frac{}{(H, R), \text{foreach } (\tau x \text{ in } r) s \rightarrow (H'[r \mapsto (\tau', (R'', \bar{s}''))], R), \{\text{foreach } (\tau x \text{ in } v) s \text{foreach } (\tau x \text{ in } r) s\}}
\end{array}$$

Figure 3: Evaluation rules for IC ω expressions, promotable expressions and statements

Compiling GMA over streams

$$\frac{|\tau.f:\tau'| \rightsquigarrow g}{|\tau*.f:\tau'*\rightsquigarrow z \mapsto z.\{g(\text{it})\}} \quad \frac{|\tau.m(\bar{\tau}'):\tau''| \rightsquigarrow g}{|\tau*.m(\bar{\tau}'):\tau''*\rightsquigarrow (z, \bar{a}) \mapsto z.\{g(\text{it}, \bar{a})\}} \quad \frac{|\tau.m(\bar{\tau}'):\text{void}| \rightsquigarrow g}{|\tau*.m(\bar{\tau}'):\text{void}| \rightsquigarrow (z, \bar{a}) \mapsto \text{foreach}(\tau \text{ it in } z) g(\text{it}, \bar{a});}$$

Compiling GMA over anonymous structs

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. |fd_{S_k}.f : \tau_k| \rightsquigarrow g_k}{|\text{struct}\{fd_1; \dots fd_n\}.f : \text{struct}\{\tau_1; \dots \tau_p\}| \rightsquigarrow z \mapsto \text{new}(g_1(z[1]), \dots, g_p(z[p]))} \quad \frac{\exists! k \in \{1 \dots n\}. |fd_k.f : \tau_k| \rightsquigarrow g}{|\text{struct}\{fd_1; \dots fd_n\}.f : \tau_k| \rightsquigarrow z \mapsto g(z[k])}$$

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. |fd_{S_k}.m(\bar{\tau}') : \tau'_k| \rightsquigarrow g_k}{|\text{struct}\{fd_1; \dots fd_n\}.m(\bar{\tau}') : \text{struct}\{\tau'_1; \dots \tau'_p\}| \rightsquigarrow (z, \bar{a}) \mapsto \text{new}(g_1(z[1], \bar{a}), \dots, g_p(z[p], \bar{a}))} \quad \frac{\exists! k \in \{1 \dots n\}. |fd_k.m(\bar{\tau}'):\tau''| \rightsquigarrow g}{|\text{struct}\{fd_1; \dots fd_n\}.m(\bar{\tau}'):\tau''| \rightsquigarrow (z, \bar{a}) \mapsto g(z[k], \bar{a})}$$

Compiling GMA over choice types

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. |\tau_{S_k}.f : \tau'_k| \rightsquigarrow g_k}{|\text{choice}\{\tau_1; \dots \tau_n\}.f : \text{choice}\{\tau'_1; \dots \tau'_p\};?| \rightsquigarrow z \mapsto (\{\text{if}(z \text{ was } \tau_{S_1}) \text{ return new choice}\{\tau'_1; \dots \tau'_p\};?(new(\tau_{S_1}, g_1(z \text{ at } \tau_{S_1})))\}; \dots \text{if}(z \text{ was } \tau_{S_p}) \text{ return new choice}\{\tau'_1; \dots \tau'_p\};?(new(\tau_{S_p}, g_p(z \text{ at } \tau_{S_p}))) \text{ else return new choice}\{\tau'_1; \dots \tau'_p\};?(null)\})}$$

$$\frac{|\tau_i.f : \tau'| \rightsquigarrow g_i \quad 1 \leq i \leq k}{|\text{choice}\{\tau_1; \dots \tau_k\}.f : \tau'| \rightsquigarrow z \mapsto (\{\text{if}(z \text{ was } \tau_1) \text{ return } g_1(z \text{ at } \tau_1); \dots \text{if}(z \text{ was } \tau_k) \text{ return } g_k(z \text{ at } \tau_k)\})} \quad \frac{\exists! k \in \{1 \dots n\}. |\tau_k.f : \tau| \rightsquigarrow g \quad n > 1}{|\text{choice}\{\tau_1; \dots \tau_n\}.f : \tau?| \rightsquigarrow z \mapsto (\{\text{if}(z \text{ was } \tau_k) \text{ return new } \tau?(g(z \text{ at } \tau_k)) \text{ else return new } \tau?(null)\})}$$

$$\frac{\exists S \subseteq \{1 \dots n\}. |S| \geq 2 \wedge p = |S| \wedge \forall k \in [1..p]. |\tau_{S_k}.m(\bar{\tau}') : \tau''_k| \rightsquigarrow g_k}{|\text{choice}\{\tau_1; \dots \tau_n\}.m(\bar{\tau}') : \text{choice}\{\tau''_1; \dots \tau''_p\};?| \rightsquigarrow (z, \bar{a}) \mapsto (\{\text{if}(z \text{ was } \tau_{S_1}) \text{ return new choice}\{\tau''_1; \dots \tau''_p\};?(new(\tau_{S_1}, g_1(z \text{ at } \tau_{S_1}, \bar{a})))\}; \dots \text{if}(z \text{ was } \tau_{S_p}) \text{ return new choice}\{\tau''_1; \dots \tau''_p\};?(new(\tau_{S_p}, g_p(z \text{ at } \tau_{S_p}, \bar{a}))) \text{ else return new choice}\{\tau''_1; \dots \tau''_p\};?(null)\})}$$

$$\frac{|\tau_i.m(\bar{\tau}'') : \tau'| \rightsquigarrow g_i \quad 1 \leq i \leq k}{|\text{choice}\{\tau_1; \dots \tau_k\}.m(\bar{\tau}'') : \tau'| \rightsquigarrow (z, \bar{a}) \mapsto (\{\text{if}(z \text{ was } \tau_1) \text{ return } g_1(z \text{ at } \tau_1, \bar{a}); \dots \text{if}(z \text{ was } \tau_k) \text{ return } g_k(z \text{ at } \tau_k, \bar{a})\})} \quad \frac{\exists! k \in \{1 \dots n\}. |\tau_k.m(\bar{\tau}'') : \tau| \rightsquigarrow g \quad n > 1}{|\text{choice}\{\tau_1; \dots \tau_n\}.m(\bar{\tau}'') : \tau?| \rightsquigarrow (z, \bar{a}) \mapsto (\{\text{if}(z \text{ was } \tau_k) \text{ return new } \tau?(g(z \text{ at } \tau_k, \bar{a})) \text{ else return new } \tau?(null)\})}$$

Compiling GMA over nullable types

$$\frac{|\tau.f:\tau'| \rightsquigarrow g}{|\tau?.f:\tau'?| \rightsquigarrow z \mapsto (\{\text{if} (z.\text{HasValue}) \text{ return new } \tau'(g(z.\text{Value})) \text{ else return new } \tau'?(null)\})}$$

$$\frac{|\tau.f:\tau'?| \rightsquigarrow g}{|\tau?.f:\tau'?| \rightsquigarrow z \mapsto (\{\text{if} (z.\text{HasValue}) \text{ return } g(z.\text{Value}) \text{ else return new } \tau'?(null)\})}$$

$$\frac{|\tau.m(\bar{\tau}'):\tau''| \rightsquigarrow g}{|\tau?.m(\bar{\tau}'):\tau''?| \rightsquigarrow (z, \bar{a}) \mapsto (\{\text{if} (z.\text{HasValue}) \text{ return new } \tau''?(g(z.\text{Value}, \bar{a})) \text{ else return new } \tau''?(null)\})}$$

$$\frac{|\tau.m(\bar{\tau}'):\tau''?| \rightsquigarrow g}{|\tau?.m(\bar{\tau}'):\tau''?| \rightsquigarrow (z, \bar{a}) \mapsto (\{\text{if} (z.\text{HasValue}) \text{ return } g(z.\text{Value}, \bar{a}) \text{ else return new } \tau''?(null)\})}$$

Figure 4: Compilation of Generalized Member Access

2. If $\Gamma; \tau \triangleright s$ and $(S, s) \rightarrow (S', s')$ then $\exists \tau'$ such that $\Gamma; \tau' \triangleright (S', s')$.
3. If $\Gamma \triangleright (S, e/pe): \tau$ then either $(S, e/pe)$ is a terminal state or $\exists S', e'/pe'$ such that $(S, e/pe) \rightarrow (S', e'/pe')$.
4. If $\Gamma; \tau \triangleright (S, s)$ then either (S, s) is a terminal state or $\exists S', s'$ such that $(S, s) \rightarrow (S', s')$.

We can also prove that our compilation preserves the typing, for instance if an $FC\omega$ expression e in environment Γ has type τ , then there is a compilation of e resulting in an $IC\omega$ expression e' , such that e' in Γ also has type τ .

Theorem 2

1. If $\Gamma \vdash e/pe: \tau$ then $\exists e'/pe'$ such that $|\Gamma \vdash e/pe: \tau| \rightsquigarrow e'/pe'$ and $\Gamma \triangleright e'/pe': \tau$.
2. If $\Gamma; \tau \vdash s$ then $\exists s'$ such that $|\Gamma; \tau \vdash s| \rightsquigarrow s'$ and $\Gamma; \tau \triangleright s'$.

Together with additional properties about the compilation and reduction, the later for example has to preserve the typing of the phrase to be reduced, it is possible to show the type soundness of $FC\omega$. The formalization is only complicated by the need to refer to a “best” conversion (as discussed in §3.4.1). Otherwise, it is analogous to Theorem 1.

4 Related work

Numerous languages have been proposed for manipulating relational and semi-structured data. For reason of space we focus here only on designs that can handle XML.

A number of special-purpose functional languages [13, 3, 8] have been proposed for processing XML values. This stands in contrast to our approach, which aimed at extending an existing widely-used imperative programming language.

The languages most similar to $C\omega$ are XJ [12] and Xtatic [11]. XJ adds XML and XPath as a first-class construct to Java, and uses logical XML classes to represent XSDs. In this way XJ allows compile time checking of XML fragments; however since the impedance mismatch between XML and objects is quite large, it does not deal with a mix of data from the the object and the XML world. One consequence is, for example, that XPath queries are restricted to work on XML data only. A unique feature of XJ is its support for updates on XML data. For $FC\omega$ we have formalized, but not yet implemented the semantics of updates for our type-system extensions.

Xtatic extends C^\sharp with an ingenious integration of regular expression types [14]. Subtyping is structural. While this

gives a lot of flexibility this neither conforms with XML Schema, where subtyping is defined by name through restrictions and extensions, nor does it allow a free mix of objects and XML. Further, Xtatic uses pattern matching for XML projections, which fits well with the chosen type system but lacks first-class queries. Updates on XML data are not supported.

In contrast to XJ and Xtatic, $C\omega$ has a much simpler type system. Its ingenuity lies in the uniform integration of the new stream, choice and struct types into the existing types and the generalization of member access—“the power is in the dot”. In fact, generalized member access in $C\omega$ achieves many of the benefits that these and other type systems try to solve. For example, a long standing problem is how to write a query over data that comes from two sources that are similar, modulo some distribution rules, but not the same [6]. The type algebra of regular expression types often allows a factorization which makes this scenario possible. Generalized member access, on the other hand, handles this problem itself, without the need for distribution rules at the type level.

Another popular approach to deal with XML in an object-oriented language is by using so called data-bindings. A data-binding generates some strongly typed object representation from a given XML schema (XSD). JAXB for Java and xsd.exe in the .NET framework generate classes from a given XSD. However, it is often impossible to generate reasonable bindings, since the rich type system of XSDs cannot adequately be mapped onto classes and interfaces. As a consequence the resulting mappings are often weakly typed.

$C\omega$ takes a simpler view: XML is considered to be a serialization syntax for the rich type system of $C\omega$. We are not tied to a particular XML data model. While $C\omega$ by design doesn’t support the entirety of the full XML stack, in our experience $C\omega$ ’s type system and language extensions are rich enough to support realistic applications. We have written a large number of applications, including the complete set of XQuery Use Cases, several XSL stylesheets, and a substantial application (50KLOC) to manage TV listings.

5 Conclusions and future work

In this paper we have considered the problem of manipulating relational and semi-structured data within common object-oriented languages. We propose that existing methods using APIs provide poor support for these common application scenarios. We have proposed a series of elegant extensions to C^\sharp that provides type-safe, first-class access to these forms of data. We have built a full compiler that implements our design. In this paper we have studied these extensions formally.

This work represents a industrial application of formal meth-

ods; on the whole, we found the process of formalizing our intuitions extremely useful, and indeed we managed to trap a number of subtle design flaws in the process. That said, we also found it useful to be simultaneously developing a compiler. On a small number of occasions we found that our formalization was too high-level, in that it failed to capture some lower-level issues. Also whilst $FC\omega$ is small enough to prove theorems about by hand, we should have liked to formalize a larger fragment of the language. At the moment, this seems unrealistic without more highly developed machine assistance.

One aspect of this project that we should like to consider further is the compilation. The Common Type System (CTS) for the Common Language Runtime (CLR) whilst general, lacks support for structural types. This means that the choice and anonymous structs have to be “simulated”. In future work, we plan to study extending the CLR with structural types. This would also enable more effective compilation of other languages that offer structural types, such as functional languages. It would also be interesting to study whether the lightweight covariance of $C\omega$ could be added to the CTS and other languages.

Implementation status A prototype $C\omega$ compiler is freely available.² The compiler covers the entirety of C^\sharp excluding unsafe code. It includes all the data access features described in this paper (and more) and also the “polyphonic” asynchronous concurrency primitives [1].

Acknowledgements We should like to thank members of the WebData team in Redmond, and the PPT group in Cambridge. We are also grateful to Sophia Drossopoulou for suggesting a number of improvements to an earlier draft.

References

- [1] N. Benton, L. Cardelli, and C. Fournet. Modern concurrency abstractions for C^\sharp . In *Proceedings of ECOOP*, 2002.
- [2] V. Benzaken, G. Castagna, and A. Frisch. CDuce: An XML-centric general-purpose language. In *Proceedings of ICFP*, 2003.
- [3] V. Benzaken, G. Gastagna, and A. Frisch. Cduce: An xml-centric general-purpose language. In *Proceedings of ICFP*, 2003.
- [4] G.M. Bierman, M.J. Parkinson, and A.M. Pitts. MJ: An imperative core calculus for Java and Java with effects. Technical Report 563, University of Cambridge Computer Lab, 2003.
- [5] G. Bracha, M. Odersky, D. Stoutamire, and P. Wadler. Making the future safe for the past: Adding genericity to Java. In *Proceedings of OOPSLA*, pages 183–200, 1998.
- [6] P. Buneman and B.C. Pierce. Union types for semistructured data. In *Proceedings of IDPL*, 1998.
- [7] D. Chamberlin et al. XQuery use cases. <http://www.w3.org/TR/xquery-use-cases/>.
- [8] S. Boag et al. XQuery. <http://www.w3.org/TR/xquery>.
- [9] M. Flatt, S. Krishnamurthi, and M. Felleisen. Classes and mixins. In *Proceedings of POPL*, 1998.
- [10] C. Fournet and G. Gonthier. The reflexive chemical abstract machine and the join-calculus. In *Proceedings of POPL*, 1996.
- [11] V. Gapeyev and B.C. Pierce. Regular object types. In *Proceedings of ECOOP*, 2003.
- [12] M. Harren, M. Raghavachari, O. Shmueli, M. Burke, V. Sarkar, and R. Bordawekar. XJ: Integration of XML processing into Java. Technical report, IBM Research, 2003.
- [13] H. Hosoya and B.C. Pierce. XDuce: A typed XML processing language. In *Proceedings of WebDB*, 2000.
- [14] H. Hosoya, J. Vouillon, and B.C. Pierce. Regular expression types for XML. In *Proceedings of ICFP*, 2000.
- [15] M. Howard and D. LeBlanc. *Writing Secure Code*. Microsoft Press, 2003.
- [16] A. Igarashi, B.C. Pierce, and P. Wadler. Featherweight Java: A minimal core calculus for Java and GJ. *ACM TOPLAS*, 23(3):396–450, 2001.
- [17] E. Meijer, W. Schulte, and G.M. Bierman. Programming with circles, triangles and rectangles. In *Proceedings of XML*, 2003.
- [18] B.C. Pierce. *Types and programming languages*. MIT Press, 2002.
- [19] D. Yu, A. Kennedy, and D. Syme. Formalization of generics for the .NET common language runtime. In *Proceedings of POPL*, 2004.

²<http://research.microsoft.com/comega>