

Unifying Tables, Objects and Documents

Erik Meijer and Wolfram Schulte
Microsoft Corporation

Abstract

This paper proposes a number of type-system and language extensions to natively support relational and hierarchical data within a statically typed object-oriented setting. In our approach SQL tables and XML documents become first class citizens that benefit from the full range of features available in a modern programming language like *C#* or Java. This allows objects, tables and documents to be constructed, loaded, passed, transformed, updated, and queried in a unified and type-safe manner.

1 Introduction

The most important current open problem in programming language research is to increase programmers productivity, that is to make it easier and faster to write correct programs [36]. The integration of data access in mainstream programming languages is of particular importance — millions of programmers struggle with this every day. Data sources and sinks are typically XML documents and SQL tables, but they don't merge nicely into a statically typed object-oriented setting in which most programs are written.

This paper addresses how to integrate tables and documents into modern object-oriented languages by providing a novel type-system and corresponding language extensions.

1.1 The Need for a Unification

Distributed web-based applications are predominantly structured using a three-tier model that consists of a *middle tier* containing the business logic that extracts relational data from a *data services tier* and munches it into hierarchical data that is displayed in the *user interface tier*. The middle tier is often programmed in an object-oriented language such as Java or *C#*.

As a consequence, middle tier programs have to deal with relational data (SQL tables), object graphs, and hierarchical data (HTML, XML). Unfortunately these three different worlds are not very well integrated. As the following ADO.Net based example shows, access to a database in this

style involves sending a string representation of a SQL query over an explicit connection via a stateful API and then iterating over a weakly typed representation of the result set:

```
SqlConnection Conn = new SqlConnection(...);
SqlCommand Cmd = new SqlCommand
    ("SELECT Name, HP FROM Pokedex", Conn);
Conn.Open();
SqlDataReader Rdr = Cmd.ExecuteReader();
```

Creating HTML or XML documents is then done by emitting document fragments in string form, without separating the model and presentation:

```
while (Rdr.Read()) {
    Response.Write("<tr><td>");
    Response.Write(Rdr.GetInt32(0));
    Response.Write("</td><td>");
    Response.Write(Rdr.GetString(1));
    Response.Write("</td></tr>");
}
```

Communication between the different tiers using untyped strings is obviously very brittle with lots of opportunities for errors and zero probability for static checking. The cynical thing is that due to the poor integration, the performance suffers badly as well.

The next code fragment rewrites the same functionality using a hypothetical language that unifies objects, tables and documents.

```
tr* pokemon =
    select <tr>
        <td>{Name}</td><td>{HP}</td>
    </tr>
    from Pokedex;
```

```
Table t =
    <table>
        <tr><th>Name</th><th>HP</th></tr>
        {pokemon}
    </table>;
```

```
Response.Write(t);
```

In this case, strongly typed XML values are first-class citizens (i.e. the XML literal `<table>...</table>` has type `static Table`) and SQL-style `select` queries are build-in. There is ample opportunity for static checking, and because the SQL and XML type-systems are integrated into the language, the compiler can do a better job in generating efficient code.

1.2 Growing a Language

It is easy to criticize the current lack of integration between tables, objects and documents, but it is much harder to come up with a design that gracefully unifies these separate worlds. No main-stream programming language has yet emerged that realizes this vision [7].

Often language integration only deals with SQL *or* with XML but not with both [11, 24, 14, 17, 2]. Alternatively they start from a completely new language such as XQuery or XDuce [6, 22, 10]. Approaches based on language binding using some kind of pre-compiler such as XSD.exe, Castor, or JAXB [28, 1] do not achieve a real semantic integration. The impedance mismatch between the different type-systems then leads to strange anomalies or unnatural mappings. Another popular route to integrate XML and SQL is by means of domain specific embedded languages [23] using functional language such as Scheme or Haskell [32, 33, 30, 31, 27, 25, 18, 40, 42, 11] as the host. In our experience however, the embedded DSL approach does not scale very well, and it is particularly difficult to encode the domain specific type-systems [38] and syntax into the host language.

In his invited talk at OOPSLA98 [20], Guy Steele remarked that

... from now on, a main goal in designing a language should be to plan for growth. The language should start small, and the language must grow as the set of users grows.

This paper shows how to grow a modern object-oriented language (we take C^\sharp as the host language, but the same approach will work with Java, Visual Basic, C++, etc.) to encompass the worlds of tables and documents by adding new types and expressions. In the remainder of this paper we will discuss:

Streams (Section 2) Streams are homogenous sequences of values of variable size. A database table consists of zero or more tuples; in the document world nodes can have zero or more sub-documents of the same kind, and in the object world we often work with (lazy) streams of values.

Tuples (Section 3) Tuples are heterogeneous sequences of values of fixed size. As we have just noticed, a database table is a stream of tuples; in the document world the **sequence** construct is used to model groups of sub-documents that must be present in a particular order, and finally several proposals have been made to extend Java and other object-oriented languages with tuples [41, 26].

Unions (Section 4) Unions represent a choice between values of different type. They play a very important role in semi-structured documents [8] and many schemas use the **choice** construct to model alternatives. Union types also occur naturally in the result-types of queries.

Content Classes (Section 5) Content classes are ordinary classes whose members can be anonymous (unnamed). We use content classes to model top-level elements and complex types in document schemas.

Queries (Section 6) Finally we will extend our repertoire of accessors of our new types to match the expressive power of XPath and SQL queries. These accessors include implicit (homomorphic extension) and explicit (apply-to-all) mapping over streams, filtering, transitive member access, and relational select and join.

The growth of our experimental language is controlled by applying the following design principles:

Denotables values should be (easily) expressible If programmers can declare a variable of a certain type, it must be possible to write an expression of that type in a convenient way.

Expressible values should be denotable If programmers can write an expression of a certain type, it must be possible to declare a variable whose static type precisely matches that of the expression.

No forced identity Programmers should never be forced to introduce either nominal identity of types, or object identity of values (aliasing).

Orthogonality There should be no special cases that discriminate between tables, documents and objects. Operations should work uniformly across the three worlds.

Flexibility The new types should have rich subtyping relationships that ease in writing type correct and evolvable software [9].

2 Streams

Streams are *typed refinements of C^\sharp 's iterators*. Iterators encapsulate the logic for enumerating elements of collections.

The **foreach** loop of C^\sharp makes it very convenient to *consume* values of type **IEnumerable** (in Java there is no interface that corresponds to **IEnumerable** and collections directly implement the **Iterator** interface instead). For instance, since type **string** implements the **IEnumerable** interface, we can iterate over all the characters in a string using a simple **foreach** loop:

```
foreach(char c in s) Console.WriteLine(c);
```

The **foreach** loop in C^\sharp is syntactic sugar for the following (simplified) **while** loop that calls into the **IEnumerable** and **IEnumerator** interfaces:

```
IEnumerator e = ((IEnumerable)s).GetEnumerator();
while (e.MoveNext()) {
    char c = (char)e.Current;
    Console.WriteLine(c);
}
```

While consuming an iterator is easy, it is much more difficult to write a generator that *implements* the **IEnumerable** (or the underlying **IEnumerator**) interface. In order to implement the **IEnumerable** interface on type **string** for instance, we have to manually create a state-machine that iterates over the individual characters in the string via **MoveNext** and exposes the current character via the **Current** property:

```

class String : IEnumerable {
    IEnumerator GetEnumerator() {
        return new Chars(this);
    }

    private class Chars : IEnumerator {
        private String s; private int i = 0;
        private char c;

        Chars(String s) { this.s = s; }

        public bool MoveNext() {
            if (i < s.Length) {
                c = s[i++]; return true;
            } else {
                return false;
            }
        }

        public char Current { get { return c; } }
    }
}

```

Note that this implementation does not correctly handle the border cases of calling `Current` before the first call to `GetEnumerator` and calling it after `GetEnumerator` has returned `false`.

In C^\sharp and Java iterators are denotable, but not easily expressible. Moreover, the type `IEnumerable` is not very accurate since it does not convey the element type of the iterator. In other words, iterators of a particular type are expressible, but not precisely denotable.

We remedy both problems by introducing a new type of streams and a new statement to generate streams.

- The type T^* denotes homogenous streams of arbitrary length with elements of type T . Type T^* is a subtype of both `IEnumerable` and `IEnumerator`.
- Stream generators are like ordinary methods except that they may yield multiple values instead of returning a single time. The `yield e` statement returns the value of expression e into the `Current` property of its corresponding stream and suspends execution until `MoveNext` is called at which time execution resumes. Upon termination of the iterator `MoveNext` returns `false`.

Using streams and generators it becomes much simpler to enumerate all the characters in a string. The helper method `char* explode(String s)` generates the stream of the individual characters of string s . The `GetEnumerator` method of class `string` then simply explodes itself:

```

class String : IEnumerable {
    private char* explode() {
        int e = this.Length;
        for(int i = 0; i < e; i++) yield s[i];
    };

    public IEnumerator GetEnumerator() {

```

```

        return this.explode();
    }
}

```

In this case maintaining the state is implicit in the control-flow of the `explode` function and in particular the borderline cases are handled correctly by definition.

Streams and generators are not new concepts. They are supported by a wide range of languages in various forms [19, 5, 37, 26, 29, 34]. Our approach is a little different in that:

- We classify streams into a hierarchy of streams of different length (see below).
- We automatically flatten streams of streams (see Section 2.1).
- We identify the value `null` with the empty stream (see Section 2.2).

To keep type-checking tractable, we restrict ourselves to the following four stream types: T^* denotes possibly empty and unbounded streams, T^+ denotes non-empty possibly unbounded streams, $T^?$ denotes streams of at most one element, and $T!$ denotes streams with exactly one element. We will use $T^?$ to represent optional values, where the non-existence is represented by the value `null` and analogously we use $T!$ to represent non-null values.

The different stream types form a natural subtype hierarchy, where subtyping corresponds to stream inclusion:

$$\begin{aligned}
 T! &<: T+ \\
 T+ &<: T^* \\
 T^? &<: T^*
 \end{aligned}$$

For instance $T! <: T^+$ reflects the fact that a stream of exactly one element is also a stream of at least one element.

We embed non-stream types T into the hierarchy by placing them between non-null values $T!$ and possibly null values $T^?$:

$$\begin{aligned}
 T! &<: T \\
 T &<: T^?
 \end{aligned}$$

The next two rules reflect the facts that `null` (we use $\emptyset^?$ for the null-type) is a possible value of any reference type, but that value types are never `null`:

$$\begin{aligned}
 \emptyset^? &<: T, \quad T \text{ is a reference type} \\
 T &<: T!, \quad T \text{ is a value type}
 \end{aligned}$$

Like arrays, streams are covariant. This means that subtyping on the element types is lifted to subtyping on streams. The special case for the `null` type says that possibly-null values can be `null`:

$$\begin{aligned}
 S &<: T \\
 \hline
 S^* &<: T^* \\
 \emptyset^? &<: T^?
 \end{aligned}$$

Let `Button` be a subtype of `Control`, then the first rule says that `Button*` is a subtype of a stream of controls `Control*`. The second rule says for instance that `null` can be assigned to a variable of type `int?`.

2.1 Coherence and Flattening

Note that the shape of our type-hierarchy allows types to be converted to T^* in different ways, which introduces the danger of incoherence [35].

The combination of covariance and embedding non-stream types provides a particularly interesting scenario. Let us write values of type T^* as $[\dots, t, \dots]$. If we start with a stream $[1, 2, 3, 4]$ of type `int*` and upcast it to type $(\text{int}^*)^*$ by using covariance, we obtain a stream of singleton streams of integers $[[1], [2], [3], [4]]$. When on the other hand we upcast the stream as a whole, we obtain a singleton stream of a stream of integers $[[1, 2, 3, 4]]$. The two different ways of upcasting a stream of type `int*` to a stream of type `int**` resulted in two different streams of streams and our type-system is incoherent!

To restore coherence we do not allow streams to be elements of other streams, in other words, both differently nested streams $[[1], [2], [3], [4]]$ and $[[1, 2, 3, 4]]$ are flattened into the non-nested stream $[1, 2, 3, 4]$. Table 2.1 gives the general flattening rules for nested streams T^{ij} of all possible combinations of stream constructors: The rule

T^{ij}	$j = !$	$?$	$+$	$*$
$i = !$!	?	+	*
$?$!	?	+	*
$+$	+	*	+	*
$*$	+	*	+	*

Figure 1: Flattening rules for streams

$T^{*+} = T^+$, for instance, reflects the fact that a non-empty stream of possibly empty stream flattens into a non-empty stream, while $T^{*+} = T^*$ reflects that a possibly empty stream of non-empty streams flattens to a possibly empty stream.

Besides flattening, the fact that `object` is at the root of the type-hierarchy also introduces a potential incoherence problem. For example since $T <: T?$ and both $T? <: \text{object}$ and $T <: \text{object}$ the coercion from $T?$ to `object` must “undo” the injection of T into $T?$ such that the two conversions end up with the same `object` value.

2.2 Non-nullness

The type $T!$ denotes streams with exactly one element, and since we identify `null` with the empty stream, this implies that values of type $T!$ can never be `null`.

Being able to express that a value cannot be `null` via the type system allows *static* checking for `null` pointers (see [15, 16] for more examples). This turns many (potentially unhandled) dynamic errors into compile-time errors.

One of the several methods in the .NET base class library that throws an `ArgumentNullException` when its argument is `null` is the function `IPAddress.Parse`. Consequently, the

implementation of `IPAddress.Parse` needs an explicit `null` check:

```
public static IPAddress Parse(string ipString) {
    if (ipString == null)
        throw new ArgumentNullException("ipString");
    ...
}
```

Dually, clients of `IPAddress.Parse` must be prepared to catch an `ArgumentNullException`. Nothing of this is apparent in the type of the `Parse` method in C^\sharp . In Java the signature of `Parse` would at least show that it possibly throws an exception.

It would be much cleaner if the type of `IPAddress.Parse` indicated that it expects its `string` argument to be non-`null`:

```
public static IPAddress Parse(string! a);
```

Now, the type-checker statically rejects any attempt to pass a string that might be `null` to `IPAddress.Parse`.

The proof obligation for returning a non-`null` stream $T!$ or T^+ is similar to proving the definite assignment rule in C^\sharp or Java. For statement blocks that return or yield non-empty streams, each non-exceptional execution path should return or yield at least one non-`null` value. The type-checker will therefore accept the first definition of `FromTo` but will reject the second:

```
int+ FromTo(int s, int d, int e) {
    yield s; while(s <= e) yield s += d;
}
```

```
// Type error
int+ FromTo(int s, int d, int e) {
    while(s <= e){ yield s; s += d; }
}
```

Non-empty streams T^+ are implicitly convertible to possibly empty streams T^* ; we can forget the fact that a stream has at least one element. It is in general *not* safe to downcast from a possibly empty stream T^* to a non-empty stream T^+ . At first sight we might think that testing if the stream contains at least one non-`null` value would suffice. Alas this is not true. By cunningly using side-effects, the generator function `OnlyOnce()` only yields 4711 the first time it is evaluated and every subsequent evaluation produces an empty stream:

```
bool Done = false;
int* OnlyOnce() {
    if(!Done){ Done = true; yield 4711; }
};
int+ xs = (int+)OnlyOnce(); // 1. cast succeeds
int+ xs = (int+)OnlyOnce(); // 2. cast fails
```

To prevent such loopholes, downcasting from T^* to T^+ will only succeed if the dynamic type of the underlying stream is T^+ .

3 Tuples

Tuples are *heterogeneous* sequences of *optionally labelled* values of *fixed* length. Another way of viewing tuples is as anonymous structs whose members are ordered, in particular tuples have no object identity.

The function `DivMod` returns the quotient and remainder of its arguments as a tuple that contains two named integer fields `sequence{int Div, Mod;}`:

```
sequence{int Div, Mod;} DivMod(int x, int y) {
    return new(Div = x/y, Mod = x%y);
}
```

The members of a tuple do not need to be labelled, for example, we can create a tuple consisting of a labelled `Button` and an unlabelled `TextBox` as follows:

```
sequence{Button enter; TextBox;} x =
    new(enter=new Button(), new TextBox());
```

An unlabelled member of a *nominal* type is a shorthand for the same member implicitly labelled with its type.

Tuples can be picked apart by position, or by member access, provided of course that that member is labelled:

```
int m = DivMod(47,11)[0];
Button b = x.enter;
```

Like streams, tuples are subject to a rich subtype hierarchy. The first subtype relation for tuples formalizes the fact that labels are optional and that we can forget them by upcasting:

$$\text{sequence}\{\dots; T\ m; \dots\} <: \text{sequence}\{\dots; T; \dots\}$$

Using this rule we see that we can assign `DivMod(47, 11)` to an unlabelled pair of integers of type `sequence{int; int;}`.

The next subtype relation for tuples allows us to forget nesting of inner tuples:

$$\begin{aligned} &\text{sequence}\{\dots; \text{sequence}\{\dots; F; \dots\}; \dots\} \\ &<: \\ &\text{sequence}\{\dots; \dots; F; \dots; \dots\} \end{aligned}$$

Finally, we can forget both labels *and* nesting by upcasting a homogeneous tuple all whose elements are of type T to a stream. The special cases give tighter types for the empty tuple (which gets converted to the empty stream `null`) and singleton tuple (which gets converted to its underlying value):

$$\begin{aligned} \text{sequence}\{T\} &<: \emptyset? \\ \text{sequence}\{T\} &<: T \\ \text{sequence}\{\dots; T; \dots\} &<: T* \end{aligned}$$

Together with the union introduction rules, which we will introduce in Section 4, this means that we can enumerate the values of *any* tuple as a stream, i.e. the tuple `new(4711,true,'z',3.14)` can be converted into the stream `[4711,true,'z',3.14]`.

3.1 Coherence and Non-Nullness

Even though tuples have no object identity, the fact that they are convertible to streams makes them subtly different from nominal value types.

Suppose that we would add the rule that tuples are not null, i.e., `sequence{...} <: sequence{...}!`. Then by applying this rule in combination with the singleton rule `sequence{T} <: T` we can assign the value `null` to a variable of non-null type `Button!`:

```
// Type error
sequence{Button;} a = new(null);
sequence{Button;}! b = a;
Button! c = b; // c = null
```

To maintain coherence we have a weaker rule that states that a tuple is non-null if it has at least one member that is non-null. This guarantees that when the tuple is converted to a stream the resulting stream has the right cardinality. For singleton sequences the conversion also holds in the reverse direction:

$$\begin{aligned} \text{sequence}\{\dots; T!\ m; \dots\} &<: \text{sequence}\{\dots; T\ m; \dots\}! \\ \text{sequence}\{T\ m; \dots\}! &<: \text{sequence}\{T!\ m; \dots\} \end{aligned}$$

By applying this rule in combination with the fact that `int <: int!`, we can show that the sequence of integers `new(1)` is convertible into a non-empty stream of type `sequence{int;} <: sequence{int!;} <: sequence{int;}! <: int*! <: int+`.

3.2 Streams+Tuples = Tables

Relational data is stored in tables, which are sets of tuples. Sets can be represented by streams, thus streams and tuples together can be used to model relational data.

The table below contains some basic facts about Pokemon characters such as their name, their strength, their kind, and the Pokemon from which they evolved (see <http://www.pokemon.com/pokedex/> for more details about these interesting creatures).

Name	HP	Kind	Evolved
Meowth	50	Normal	
Rapidash	70	Fire	Ponyta
Charmelon	80	Fire	Charmander
Zubat	40	Plant	
Poliwag	40	Water	
Weepinbell	70	Plant	Bellsprout
Ponyta	40	Fire	

Each row in this table is a value of type `Pokemon` and the table itself is modelled as a variable `Pokedex` of type `Pokemon*`. The keyword `type` identifies the name on the left with the type expression on the right. It is just an abbreviation mechanism.

```
enum Kind {Water, Fire, Plant, Normal, Rock}
```

```
type Pokemon = sequence{
    string Name; int HP; Kind Kind; string? Evolved;
```

```
}

```

```
Pokemon* Pokedex;
```

The fact that basic Pokemon are not evolutions of other Pokemon shows up in that the `Evolved` column has type `string?`.

Representing tables is necessary for the integration of relational data, but it is not sufficient: we also have to provide operations that work on tables. We will introduce such query expressions in Section 6.3.

4 Unions

Union types often appear in content classes (see section 5 below). The type `Address` uses a union type to allow either a member `Street` of type `string` or a member `POBox` of type `int`:

```
class Address {
  sequence{
    choice{ string Street; int POBox; };
    string City; string? State; int Zip;
    string Country;
  };
}
```

The second situation in which union types are used is in the result types of generalized member access (see Section 6). For example, when `p` has type `Pokemon`, the expression `p.*` returns a stream containing all the members of a `Pokemon` and has type `choice{string; int; Kind; string?}*`. Using the subtype rules for `choice` and streams this is equivalent to `choice{string; int; Kind;}*`.

Finally, union types enable any tuple to be converted into streams, i.e. the pair `sequence{Button b; TextBox}` can be upcast to a stream of type `choice{Button; TextBox;}*` that contains both `Button` and `TextBox` elements by using the fact that both `Button` and `TextBox` are subtypes of `choice{Button; TextBox;}`. However, as we will discuss below, in general the conversion from a type `S` to the union type `choice{S; T;}` is not unique and hence leads to possible loss of coherence.

Choice types are idempotent (duplicates are removed), associative and commutative (nesting and order are ignored):

```
choice{...; F; F; ...} = choice{...; F; ...}
choice{...; choice{...}; ...} = choice{...; ...; ...}
choice{...; F; G; ...} = choice{...; G; F; ...}
```

We can safely eliminate a union by converting it to a common supertype of all its members:

$$\frac{\forall i :: S_i <: T}{\text{choice}\{ \dots; S_i; \dots \} <: T}$$

This ensures that `choice{...} <: object`.

Streams distribute over unions. Non-nullness and possibly nullness distribute in both ways, and any inner streams gets

absorbed by an outer `+` or `*`:

```
choice{...; T; ...}! = choice{...; T!; ...}!
choice{...; T; ...}? = choice{...; T?; ...}?
choice{...; Ti; ...}+ = choice{...; T; ...}+
choice{...; Ti; ...}* = choice{...; T; ...}*
```

where i is any stream functor.

The flattening and distribution rules allow us to normalize streams of choices: inner stream functors can either be eliminated completely or can be moved out of the choice.

4.1 Unions, overloading and legalized incoherence

The reason that we have not yet defined an introduction rule for unions has everything to do with coherence. Suppose that we have a class `C` that implements two independent interfaces `IA` and `IB`, then we can convert `C` into `choice{IA; IB;}` by converting to either `IA` or `IB`.

This kind of ambiguity is similar to what happens for overloaded methods in Java and C^\sharp . In the example below, the type-checker cannot decide between converting `C` to `IA` or to `IB` and therefor marks the call to `f(new C())` as ambiguous:

```
static void f(IA a){ Console.WriteLine("IA"); }
static void f(IB a){ Console.WriteLine("IB"); }
```

If `IB` extends `IA`, then the type-checking rules select the conversion from `C` to `IB` as a better match and the overloading is resolved to call `void f(IB)`.

Overloading in Java and C^\sharp is a form of legalized incoherence, and we will use the same overloading resolution rules for upcasting to a union type. Hence if `IB <: IA` then converting type `C` to `choice{IA; IB;}` will choose the best matching alternative `IB`, but without providing any coherence guarantees.

$$S <: \text{choice}\{ \dots; S; \dots \}$$

It follows that `choice` is also covariant, however, this can also introduce incoherence:

$$\frac{S <: T}{\text{choice}\{ \dots; S; \dots \} <: \text{choice}\{ \dots; T; \dots \}}$$

4.2 Restoring coherence

As the previous section shows, conversions into union types may cause incoherence. The problem is that a value `new C()` can be coerced into two different values of type `choice{IA; IB;}`, i.e. it *nondeterministic*. However, the situation is not completely hopeless, since we can still prove that conversions into *deterministic* types are coherent.

Each conversion $S <: T$ has a corresponding witness function $Tf(S)$ that coerces values of type `S` into values of type `T`. Hence every value t of type `T` is of the form $f(s)$ for some explicitly constructed (i.e., not obtained by means of a coercion) *seed* value s of type `S` such that $S <: T$ via witness function f . The function *unwrap* takes a value $f(s)$ where f is a (composition) of conversions and retrieves the

underlying seed s . The *unwrap* function correspond to the intuition that coercions may change the representation of a value, but that they never change their actual value. In the actual implementation this corresponds to the fact that conversions are either identity preserving or implemented using wrappers.

A type T is called *deterministic* if $unwrap(t_1) = unwrap(t_2) \implies t_1 = t_2$ for all values t_1 and t_2 in T . Our coherence theorem says that *all conversions $S <: T$ where T is deterministic are coherent*. The precise formulation of the coherence theorem and its formal proof are subject of a forthcoming paper.

The theorem implies that even if a chain of conversions involves (non-disjoint) union types, the total conversion is coherent once we reach a deterministic type. For example it is the case that the two conversions $C <: IA|IB <: IA$ and $C <: IA$ are coherent.

5 Content Classes, XSDs and XML

Now that we have introduced streams, tuples, and unions, our type system is rich enough to model a large part of the XSD schema language; our aim is to cover as much of the essence of XSD [39] as possible while avoiding most of its complexity.

The correspondence between XSD particles such as `<sequence>` and `<choice>` with local element declarations and the type constructors `sequence` and `choice` with (labelled) fields should be intuitively clear. Likewise, the relationship of XSD particles with occurrence constraints to streams is unmistakable. For T^* the attribute pair (`minOccurs`, `maxOccurs`) is (0, `unbounded`), for T^+ it is (1, `unbounded`), for $T^?$ it is (0, 1), and for $T!$ it is (1,1).

The content class `Address` that we defined in Section 4 corresponds to the XSD schema `Address` below:

```
<element name="Address"><complexType><sequence>
  <choice>
    <element name="Street" type="string">
    <element name="POBox" type="integer">
  </choice>
  <element name="City" type="string">
  <element name="State" type="string"
    minOccurs="0"/>
  <element name="Zip" type="integer"/>
  <element name="Country" type="string"/>
</sequence></complexType></element>
```

The only difference between a content class and a normal $C^\#$ class is the fact that the content class does not have a field label. As a consequence, the content can only ever be accessed via its individually named children, which allows the compiler to choose the most efficient data layout.

The next example schema defines two top level elements `Author` and `Book` where `Book` elements can have zero or more `Author` members:

```
<element name="Author"><complexType><sequence>
  <element name="Name" type="string"/>
</sequence></complexType></element>

<element name="Book"><complexType><sequence>
```

```
<element name="Title" type="string"/>
<element ref="Author" minOccurs="0"
  maxOccurs="unbounded"/>
</sequence></complexType></element>
```

In this case, the local element reference is modelled by an unlabelled field and the schema is mapped onto the following two content type declarations:

```
class Author { string Name; }
class Book { sequence{ string Title; Author*; } }
```

All groups such as the one used in the following schema for the complex type `Name`

```
<complexType name="Name"><all>
  <element name="First" type="string"/>
  <element name="Last" type="string"/>
</complexType></all>
```

are mapped to ordinary fields of the containing type, i.e. without a `sequence`:

```
class Name { string First; string Last; }
```

As these examples show, both top-level element declarations and named complex type declarations are mapped to top-level types. This allows us to unify derivation of complex types and substitution groups of elements using standard inheritance.

5.1 XML Literals

XML literals are an intuitive way to construct instances of content classes. We can define an `Address` instance simply by writing an XML document that conforms to the schema for `Address`:

```
Address Microsoft =
  <Address>
    <Street>One Microsoft Way</Street>
    <City>Redmond</City><Zip>98052</Zip>
    <Country>USA</Country>
  </Address>;
```

XML literals can also have placeholders to describe *dynamic* content (anti-quoting). We use the XQuery convention whereby an arbitrary expression or statement block can be embedded inside an element by escaping it with curly braces.

```
Author NewAuthor(string name) {
  return <Author>{name.ToUpper()}</Author>;
}
```

Embedded expressions must return or yield values of the required type (in this case `string`). Validation of XML literals with placeholders is non-trivial and is the subject of a forthcoming paper.

Note that XML literals are just object constructors, there is nothing special about content classes. In fact, we can write XML literals to construct values of any type, for example, the assignment

```

Button b = <Button>
    <Text>Click Me</Text>
</Button>;

```

creates an instance of the standard `Button` class and sets its `Text` field to the string "Click Me".

6 Generalized Member Access

In the previous sections we have concentrated on the type-system extensions to our theoretical language. In this section we will extend our repertoire of *expressions* to transform and query values of these new types.

6.1 Apply-to-All, Lifting and Filter

To make the creation of streams as concise as possible, we allow statement blocks (anonymous method bodies) as expressions. In the example below we assign the (conceptually) infinite stream of positive integers to the variable `nats`:

```

// 0, 1, 2, ...
int* nats = { int i=0; while(true) yield i++; };

```

Our stream constructors (`*`, `+`, `?`, `!`) are functors, and hence we *implicitly lift operations on the element type of a stream* (such as member or property access and method calls) *over the stream itself*. For instance, to convert each individual string in a stream `Ss` of strings to uppercase, we can simply write `ss.ToUpper()`:

```

string* Ss = { yield "Hello"; yield "World!"; };
string* SS = Ss.ToUpper();

```

This only works since streams do not have an `ToUpper` method. If both streams and elements have the same member, no lifting takes place and member access on the whole stream is the best match.

If we nevertheless want to lift member access over a stream, we can use an *apply-to-all* block. For example, both streams and integers have a `GetHashCode` method. To convert the stream of integers `nats` to a stream of their hash-codes, we write `nats.{ return it.GetHashCode(); }`:

```

int hc = nats.GetHashCode();
int* hcs = nats.{ return it.GetHashCode(); };

```

The implicit argument `it` refers successively to each element of the stream `nats`.

As the next example shows, the *apply-to-all* block itself can yield a stream, in which case the resulting nested stream is flattened according to the rules of table 2.1

```

// 1, 2,2, 3,3,3, 4,4,4,4, ...
int* rs = nats.{ for(i=0; i<it; i++) yield it; };

```

If an *apply-to-all* block returns `void`, no new stream is constructed and the block is eagerly applied to all elements of the stream. For example to print all the elements of a stream we can just write:

```

nats.{ Console.WriteLine(it); };

```

Apply-to-all blocks can be stateful, so we can use them to do reductions (in the functional community called *folds*). For example, we can sum all integers in an integer stream `xs` as follows:

```

int sum(int* xs){
    int s = 0; xs.{ s += it; }; return s;
}

```

We need to be careful when lifting over non-null types, since the fact that the receiver object is not `null` does not imply that its members are not `null` either:

```

Button! b = <Button/>;
Control p = b.Parent; // Parent might be null

```

Hence the return type of lifting over a non-null type is not guaranteed to return a non-null type.

Optional types provide a standard implementation of the null design pattern [21]; when a receiver of type `T?` is `null`, accessing any of its members returns `null`:

```

string? t = null;
int? n = t.Length; // n = null

```

In Objective-C [3] this is the standard behavior for any object that can be `null`.

The table 6.1 show how lifting of member-access interacts with streams types. Let T^j be a stream type, and m of type S^i be a member of the element type T that we want to lift over the stream. The result type of lifting m is then given by $s^{i \oplus j}$ (here $_$ denotes a naked type):

\oplus	$j=_$!	?	+	*
$i=_$?	?	*	*
!	!	!	?	*	*
?	?	?	?	*	*
+	+	+	*	+	*
*	*	*	*	*	*

Figure 2: Lifting over streams

Member access is not only lifted over streams, but over all structural types. For example the expression `xs.x` will return the stream `true, 1, 2` of type `choice{bool;int;}+` when `xs` is defined as:

```

sequence{ bool x; sequence{ int x; }*; } xs =
    new( x=true
        , { yield new(x=1); yield new(x=2); }
        );

```

Lifting over union types introduces a possibility of nullness for members that are not in all of the alternatives.

Suppose `x` has type `choice{ int; string; }`. Since only `string` has a `Length` member, the type of `x.Length` is `int?` which reflects the fact that in case the dynamic type of `x` is `int`, the result of `x.Length` will be `null`. Since `int` and `string` both have a member `GetType()`, the return type of `x.GetType()` is `Type`:

```

choice{ int; string; } x = 4711;
int? n = x.Length; // null
Type t = x.GetType(); // System.Int32

```

In case the alternatives of a union have a member of different type in common, the result type is the union of the types of the respective members.

Often we want to filter a stream according to some predicate on the elements of the stream. For example, to construct a stream with only odd numbers, we filter out all even numbers from the stream `nats` of natural numbers using the filter expression `nats[it %2 ==1]`. For each element in the stream to be filtered, the predicate is evaluated with that element as `it`. Only if the predicate is true the element becomes part of the new stream.

```
int* odds1 = nats[it%2 == 1];
```

On closer inspection, we realize that filters are just abbreviations of an apply-to-all-block:

```
int* odds2 = nats.{if (it%2 == 1) yield it;};
```

Hence `odds1` and `odds2` denote streams that both have the same elements in the same order.

6.2 Wildcard, Transitive and Type-based Member-access

The only query form available in object-oriented languages is member access. But that is rather restrictive. To allow for more flexible forms of member access, we provide *wildcard*, *transitive* and *type-based* access. These forms are similar to the concepts of *nametest*, abbreviated relative location paths and name filters in XPath [13]. However we adapted them to work uniformly on object graphs.

Wildcards allow to access all accessible members of a type without having to know their names. Suppose that we want to have all fields of an `Address`, then we can write:

```
choice{string; int;}* addressfields = Microsoft.*;
```

The wild-card expression returns the content of all accessible fields and properties of the variable `Microsoft` in their declaration order. In this case "One Microsoft Way", "Redmond", 98052, "USA".

Transitive member-access, written as `e...m`, returns all accessible members `m` that are transitively reachable from `e` in depth-first order. The following declaration of `authors` (lazily) returns a stream containing all `Author` of all `Books` in the source stream `books`:

```
Book F = <Book>
  <Title>Faust</Title>
  <Author>Goethe</Author>
</Book>;
Book K = <Book>
  <Title>De Klompeniers</Title>
  <Author>Jac. Broersen</Author>
</Book>;
```

```
Book* books = { yield F; yield K; };
string* authors = books...Author;
```

Transitive member access allows to abstract from the concrete representation of a document; as long as the mentioned member is reachable and accessible, its values are returned.

Looking for just a field name might not be enough, especially for transitive queries where there might be several reachable members with the same name, but of different type. In that case we can add an additional type-test to restrict the matching members. A type-test on `T` selects only those members whose static type is a subtype of `T`.

For instance, if we are only interested in Microsoft's POBox number, and Zip code, we can write the transitive query `Microsoft...int::*`.

Note that type based access is also useful for unnamed members, since even if they have no name, they do have a static type.

6.3 Select and Join

The previous sections presented our solutions to querying documents. However for accessing relational data, which we have modelled as streams of tuples, simpler SQL queries are sufficient. Here we only show the integration of the SQL `select-from-where` clause, and defer the discussion of more advanced features such as data manipulation and transactions to a future paper.

The fundamental operations of relational algebra are *selection*, *projection*, *union*, *difference* and *join*. Selection is similar to filter and transforms one stream of tuples into another stream of tuples. Here are two variations of selection:

```
Pokemon* normalPokemons1 =
  select * from Pokedex where Kind == Normal;
Pokemon* normalPokemons2 =
  select it from (Pokemon it in Pokedex)
  where it.Kind == Normal;
```

The first example uses the familiar SQL syntax. Its meaning is provided by the second form, which uses the explicit iterator variable `it` as we have seen before.

We use similar sugar to introduce names for projection. Projection produces a stream of tuples by selecting only certain columns in its input stream:

```
sequence{string Name; Kind Kind;}*
pokemonAbstract1 =
  select Name, Kind from Pokedex;
sequence{string Name; Kind Kind;}*
pokemonAbstract2 =
  select new(Name= it.Name, Kind=it.Kind)
  from (Pokemon it in Pokedex);
```

Again, the first declaration shows the traditional SQL syntax, where the second shows the unsugared representation, which explicitly builds the resulting tuple by projecting the required members.

In practice, the result types of SQL queries can be quite involved and hence it becomes painful for programmers to explicitly specify types. Since the compiler already knows the types of sub-expressions, the result types of queries can be inferred automatically. Providing type declarations for method local variables is not necessary, and we can simply write:

```
pokemonAbstract3 =
  select Name, Kind from Pokedex;
```

without having to declare the type of `pokemonAbstract3`.

Union and difference present no difficulty in our framework. They can easily be handled with existing operations on streams. Union concatenates two streams into a single stream. Difference takes two streams, and returns a new stream that contains all values that appear in the first but not in the second stream.

The real power of `select-from-where` comes from join. Join takes two input streams and creates a third stream whose values are composed by combining members from the two input streams. For example, here is an expression that selects pairs of Pokemons which have evolved from each other:

```
select p.Name, q.Name
from p in Pokedex, q in Pokedex
where p.Evolved == q
```

Again, we would like to stress the fact that everything fits together. The select expression works on arbitrary streams, whether in memory or on the harddisk; streams simply virtualize data access. Strong typing makes data access secure. But there is no burden for the programmer since the result types of queries are inferred.

7 Conclusion

The language extensions proposed in this paper support both the SQL [4] and the XML schema type system [39] to a large degree, but we have not dealt with all of the SQL features such as (unique) keys, and the more esoteric XSD features such as redefine. Similarly, we already covered much of the expressive power of XPath [13], XQuery [6] and XSLT[12], but we do not support the full set of XPath axis. We are able to deal smoothly with namespaces, attributes, blocking, and facets however. Currently we are investigating whether and which additional features need to be added to our language.

Summarizing, we have shown that it is possible to have both SQL tables and XML documents as first order citizen in an object-oriented language. Only a bridge between the type worlds is needed. Building the bridge is mainly an engineering task. But once it is available, it offers the best of three worlds

8 Acknowledgments

We would like to acknowledge the support, encouragement, and feedback from Mike Barnett, Nick Benton, Don Box, Luca Cardelli, Bill Gates, Steve Lucco, Chris Lucas, Todd Proebstring, Dave Reed, and Clemens Szyperski and the hard work of the WebData languages team consisting of William Adams, Joyce Chen, Kirill Gavrylyuk, David Hicks, Steve Lindeman, Chris Lovett, Frank Mantek, Wolfgang Manousek, Neetu Rajpal, Herman Venter, and Matt Warren. Special thanks to Rostislav Yavorskiy for helping to formulate and prove the coherence theorem.

References

- [1] The castor project. <http://castor.exolab.org/>.
- [2] The navision x++ programming language.
- [3] The Objective-C Programming Language.
- [4] G. Bierman and A. Trigoni. Towards a Formal Type System for ODMG OQL. Technical Report 497, University of Cambridge, Computer Laboratory, 2000.
- [5] B.Liskov, R.Atkinson, T. Bloom, E. Moss, J.C.Schaffert, R. Scheifler, and A. Snyder. *CLU Reference Manual*. LNCS 114. Springer=Verlag, 1981.
- [6] S. Boag, D. Chamberlin, M. F. Fernandez, D. Florescu, J. Robie, and J. Siméon. XQuery 1.0: An XML Query Language. Technical report, W3C, November 2002.
- [7] T. Bray. XML Is Too Hard For Programmers. <http://www.tbray.org/ongoing>.
- [8] P. Buneman and B. Pierce. Union Types for Semistructured Data. In *International Database Programming Languages Workshop*, LNCS 1949, 2000.
- [9] L. Cardelli. Types for data-oriented languages. In *EDBT*, 1988.
- [10] L. Cardelli, P. Gardner, and G. Ghelli. Manipulating trees with hidden labels. In *FOSSACS*, 2003.
- [11] A. S. Christensen, A. Muller, and M. I. Schwartzbach. Static Analysis for Dynamic XML. In *PlanX*, 2002.
- [12] J. Clark. XSL Transformations (XSLT). Technical report, W3C, November 1999.
- [13] J. Clark and S. DeRose. XML Path Language (XPath) Version 1.0. Technical report, W3C, November 1999.
- [14] R. Connor, D. Lievens, and F. Simeoni. Projector: a partially typed language for querying XML. In *PlanX*, 2002.
- [15] M. Fahndrich and R. M. Leino. Non-Null Types in an Object-Oriented Language.
- [16] C. Flanagan, K. Leino, M. Lillibridge, C. Nelson, J. Saxe, and R. Stata. Extended Static Checking for Java. In *PLDI*, 2002.
- [17] V. Gapeyev and B. Pierce. Regular object types. In *ECOOP*, 2003.
- [18] P. Graunke, S. Krishnamurthi, S. V. D. Hoeven, and M. Felleisen. Programming the Web with High-Level Programming Languages. In *Automated Software Engineering*, LNCS 2028, 2001.
- [19] R. Griswold and M. Griswold. *The Icon Programming Language (2nd edition)*. Prentice Hall, 1990.
- [20] J. Guy L. Steele. Growing a Language. *Journal of Higher-Order and Symbolic Computation*, 12(3), 1999.
- [21] K. Henney. Null Object, Something for Nothing. In *Seventh European Conference on Pattern Languages of Programs*, 2002.

- [22] H. Hosoya and B. C. Pierce. XDuce: A Typed XML Processing Language. In *International Workshop on the Web and Databases (WebDB)*, 2000.
- [23] P. Hudak. Building Domain Specific Embedded Languages. *ACM Computing Surveys*, 28(4), 1996.
- [24] M. Kempa and V. Linnemann. On XML Objects. In *PlanX*, 2002.
- [25] O. Kiselyov and S. Krishnamurthi. SXSLT: Manipulation Language for XML. In *PADL*, LNCS 2562, 2003.
- [26] A. Krall and J. Vitek. On Extending Java. In *JMLC*, 1997.
- [27] D. Leijen and E. Meijer. Domain Specific Embedded Compilers. In *2nd USENIX Conference on Domain-Specific Languages*, 1999.
- [28] T.-W. Lin. Java architecture for xml binding (jaxb): A primer, 2002.
- [29] B. Liskov, M. Day, M. Herlihy, P. Johnson, and G. Leavens. ARGUS Reference Manual. Technical Report MIT/LCS/TR-400, 1987.
- [30] E. Meijer. Server Side Web Scripting in Haskell. *Journal of Functional Programming*, 10(1), January 2000.
- [31] E. Meijer, D. Leijen, and J. Hook. Client Side Web Scripting with HaskellScript. In *PADL*, 2002.
- [32] E. Meijer and M. Shields. XMLambda: a Functional Language for Constructing and Manipulating XML Documents, 1999. <http://www.cse.ogi.edu/~mbs>.
- [33] E. Meijer and D. van Velzen. Haskell Server Pages. In *Haskell Workshop 2000*, 2000.
- [34] S. Murer, S. Omohundro, D. Stoutamire, and C. Szyperski. Iteration abstraction in Sather. *ACM Transactions on Programming Languages and Systems*, 18(1):1–15, January 1996.
- [35] B. C. Pierce. *Types and Programming Languages*. The MIT Press, 2002.
- [36] T. Proebstring. Disruptive Language Technologies.
- [37] N. Schemenauer, T. Peters, and M. L. Hetland. Simple Generators. PEP-0255.
- [38] M. Shields and E. Meijer. Type-indexed Rows. In *POPL*, 2001.
- [39] J. Simeon and P. Wadler. The Essence of XML. In *POPL*, 2003.
- [40] P. Thiemann. WASH/CGI: Server Side Web Scripting with Sessions and Typed, Compositional Forms. In *Practical Aspects of Declarative Languages*, 2002.
- [41] C. van Reeuwijk and H. Sips. Adding tuples to java: a study in lightweight data structures. In *ACM Java Grande/ISCOPE*, 2002.
- [42] N. Welsh, F. Solsona, and I. Glover. SchemeUnit and SchemeQL: Two Little Languages. In *Third Workshop on Scheme and Functional Programming*, 2002.