

On Binary Methods

Kim Bruce, Luca Cardelli, Giuseppe Castagna,
The Hopkins Objects Group,
Gary T. Leavens, and Benjamin Pierce

TR #95-08
May 1995

Keywords: object-oriented programming, type checking, binary method, matching, subtyping, precise type, multi-method, behavioral subtyping, covariance, contravariance.

1994 CR Categories: D.1.5 [*Programming Techniques*] Object-oriented Programming; D.2.2 [*Software Engineering*] Tools and Techniques — modules and interfaces; D.3.1 [*Programming Languages*] Formal Definitions and Theory — semantics; D.3.2 [*Programming Languages*] Language Classifications — object-oriented languages; D.3.3 [*Programming Languages*] Language Constructs — Abstract data types, modules, packages; F.3.1 [*Logics and Meanings of Programs*] Specifying and Verifying and Reasoning about Programs — logics of programs; F.3.2 [*Logics and Meanings of Programs*] Studies of Program Constructs — type structure.

Submitted for publication.

© Kim Bruce, Luca Cardelli, Giuseppe Castagna, The Hopkins Objects Group, Gary T. Leavens, and Benjamin Pierce, 1995. The copyright will eventually be given to the publisher of the final paper.

Department of Computer Science
226 Atanasoff Hall
Iowa State University
Ames, Iowa 50011-1040, USA

On Binary Methods

Kim Bruce* Luca Cardelli† Giuseppe Castagna‡ The Hopkins Objects Group§
Gary T. Leavens¶ Benjamin Piercel||

May 16, 1995

Abstract

Giving types to binary methods causes significant problems for object-oriented language designers and programmers. This paper offers a comprehensive description of the problems arising from typing binary methods, and collects and contrasts diverse views and solutions. It is intended to expose a wide audience of readers to the current debate on this question.

1 Introduction

Binary methods have caused great difficulty for designers of strongly typed object-oriented languages and for programmers using those languages. In this paper we study the sources of these problems and compare and contrast a variety of solutions.

The authors of this paper have differing views on what the most appropriate solutions are. We have attempted here to collect together the solutions that individuals among us advocate and to present a consensus on what can be fairly stated as strengths and weaknesses of each approach. This paper grew from presentations and discussions at the 2nd Workshop on Foundations of Object-Oriented Languages, which was sponsored by NSF and ESPRIT and held in Paris in June, 1994 [18].

Informally, a simple binary method of some object is a method that expects to be passed another object of the same class as an argument. Such a method is *binary* in the sense that it acts on two objects: the object passed as argument and the receiving object itself. In general, a binary method could also include other arguments; by a standard abuse of terminology we still refer to these as binary methods. Simple examples of binary methods include arithmetic operations on number objects, as well as binary relations such as = and <, and set operations like subset and union.

We can characterize binary methods more precisely in terms of types. A method m occurring in some object of type τ is a binary method if m has a type where one or more arguments are of type τ , the type of the object itself.

*Department of Computer Science, Williams College, Williamstown, Massachusetts 01267, USA. Internet: kim@cs.williams.edu. Bruce's research was partially supported by NSF grant CCR-9121778.

†Digital Equipment Corporation, Systems Research Center, 130 Lytton Ave, Palo Alto, CA 94301, USA. Internet: luca@src.dec.com.

‡(CNRS) LIENS-DMI, 45 rue d'Ulm, 75005 Paris, FRANCE. Internet: castagna@dmi.ens.fr

§Jonathan Eifrig, Scott Smith, Valery Trifonov. Contact Scott Smith, Department of Computer Science, The Johns Hopkins University, Baltimore, Maryland 21218. Internet: scott@cs.jhu.edu. Research partially supported by NSF grant CCR-9301340 and AFOSR grant F49620-93-1-0169.

¶229 Atanasoff Hall, Department of Computer Science, Iowa State University, Ames Iowa, 50011, USA. Internet: leavens@cs.iastate.edu. Leavens's research was partially supported by NSF grant CCR-9108654.

||Computer Laboratory, New Museums Site, Pembroke Street, Cambridge CB2 3QG, United Kingdom. Internet: benjamin.pierce@cl.cam.ac.uk

```

class PointClass
  instance variables
    xValue: real
    yValue: real
  methods
    x: real is return(xValue)
    y: real is return(yValue)
    equal(p: Point): bool is return( (xValue==p.x) && (yValue==p.y) )
end class

```

Figure 1: The class *PointClass*.

The most significant problem with binary methods lies in their typing in the presence of inheritance. There is a conflict between inheritance in the presence of binary methods and subtyping [23]; binary methods seem to require that one of the two be sacrificed. A second problem is the asymmetry of a binary method: the method may have privileged access to only one of the two objects the method is invoked on. These two problems are described in more detail in Section 2.

Sections 3 and 4 concentrate on solutions to the problem of typing binary methods in the presence of inheritance. We consider the question from two sides: in Section 3, we reflect on whether it actually need be solved at all (i.e., whether binary operations might best *not* be treated as methods); in Section 4, we meet the problem head-on and review some solutions that have been proposed.

Turning to the problem of privileged access, Section 5 sketches a technique by which object-style data encapsulation can be blended with conventional ADT-style encapsulation to allow implementation of binary operations with privileged access to object representations. Section 6 addresses the problem of binary methods from the point of view of behavioral specification. Section 7 offers concluding remarks.

2 The Problem of Binary Methods

This section describes the problems caused by binary methods. The first section describes typing problems in the presence of inheritance, and the second describes problems with privileged access.

2.1 Typing Binary Methods in the Presence of Inheritance

In procedural or functional languages, the type of a binary function that takes two arguments of type τ and returns a value of type σ is written $\tau \times \tau \rightarrow \sigma$. In an object-oriented language, functions or procedures are typically replaced by *methods* belonging to a class corresponding to one of the arguments. Figure 1 shows a standard example of a class with a binary method.¹ In *PointClass*, the method *equal*, which tests for equality with another instance of *PointClass*, is written with a single parameter of type *Point*. As may be seen in this example, binary operations—when regarded as

¹A few notes on our notation follow. Methods are functions or procedures whose body occurs after the keyword *is*. We write parameterless functions and procedures by omitting the parentheses. We write the type of parameterless functions as if they were variables of their return type. That is, we omit an implicit *unit* \rightarrow before the result type. Commented text is preceded by --.

```

class ColorPointClass subclass of PointClass
  instance variables -- xValue and yValue are inherited
    cValue : string
  methods -- x and y are inherited
    c: string is return(cValue)
    equal(p: ColorPoint): bool is
      return( (cValue==p.c) && (xValue==p.x) && (yValue==p.y) )
end class

```

Figure 2: The class *ColorPointClass*.

methods—are asymmetric: the receiver plays a role somewhat different than the parameter. This distinction is highlighted when we define a subclass of a class with a binary method.

Figure 2 defines a subclass *ColorPointClass* of *PointClass*. In *ColorPointClass*, the type of the parameter of *equal* is changed to *ColorPoint* to match the type of the receiver, allowing two *ColorPoint* objects to be compared by the *equal* method. If we model objects as recursive records with instance variables hidden, then instances of *PointClass* and *ColorPointClass* have the following record types:

$$\begin{aligned}
 \textit{Point} &\equiv \langle\langle x: \textit{real}; y: \textit{real}; \textit{equal}: \textit{Point} \rightarrow \textit{bool} \rangle\rangle \\
 \textit{ColorPoint} &\equiv \langle\langle x: \textit{real}; y: \textit{real}; c: \textit{string}; \textit{equal}: \textit{ColorPoint} \rightarrow \textit{bool} \rangle\rangle
 \end{aligned}$$

Note that both of these definitions are recursive: the type being defined appears on the right-hand side of the \equiv . Note also that the instance variables (such as *xValue* and *yValue*) are not visible in the type of objects generated from a class.

Informally, a type σ is a *subtype* of τ , written $\sigma <: \tau$, if an expression of type σ can be used in any context that expects an expression of type τ (cf. [12, 14, 49]). Associated with subtyping is the principle of *subsumption* (subtype polymorphism): if $\sigma <: \tau$ and a program fragment has type σ , it also has type τ . A simple example of subtyping in object-oriented programming is that an object type is a subtype of the type with some methods removed, as any context that expects the object with fewer methods will not directly use the extra methods and thus no type errors will occur. Similar reasoning can be applied to record types: a record r can be used where another record r' is expected only if r has at least the same fields as r' and the types of these fields are subtypes of the types of the corresponding fields of r' . In other words $\langle\langle \ell_1: S_1, \dots, \ell_n: S_n, \dots, \ell_{n+k}: S_{n+k} \rangle\rangle <: \langle\langle \ell_1: T_1, \dots, \ell_n: T_n \rangle\rangle$ (with $k \geq 0$) if and only if, for each $i \in \{1..n\}$, $S_i <: T_i$.

The rule for subtyping functions states that $\sigma \rightarrow \tau <: \sigma' \rightarrow \tau'$ if and only if $\sigma' <: \sigma$ and $\tau <: \tau'$ [12]. (This is sometimes called the “contravariant rule” because it is contravariant in the left argument of \rightarrow .) This rule is informally justified by the following. If f is expected to have type $\sigma' \rightarrow \tau'$, but actually has type $\sigma \rightarrow \tau$, then f can be passed an argument of type σ' when (by subsumption) $\sigma' <: \sigma$; furthermore, the result of such a call will have type τ , which (by subsumption) can be considered to be of type τ' . Hence all functions of type $\sigma \rightarrow \tau$ can be used as if they had type $\sigma' \rightarrow \tau'$ without type error.

Subtype polymorphism is a useful feature of object-oriented programming: if subclasses corresponded to subtypes, a subclass object could always be passed to a function or method expecting a superclass object, allowing re-use of code. Unfortunately, subclasses do not always generate subtypes, if the types of methods need to change in subclasses (and as the example illustrates,

```

procedure breakit(p: Point)
  var
    nuPt: Point
  begin
    nuPt := new PointClass(3.2, 4.5)
    if p.equal(nuPt) then
      ...
  end

```

Figure 3: The procedure **breakit**.

the types of binary methods evidently do need to change). Because of the contravariance of the subtyping relation on the domain of *equal*, *ColorPoint* is not a subtype of *Point*. For the subtype relation to hold, the type of *equal* in *ColorPoint* would have to be a subtype of the type in *Point*. Thus $ColorPoint \rightarrow bool$ must be a subtype of $Point \rightarrow bool$. But, by the subtyping principle for functions, this requires *Point* to be a subtype of *ColorPoint*, exactly the *opposite* of what we are after and clearly untrue.

This loss of subtyping in this case is not an artifact of the definition of subtyping for functions; the procedure **breakit** of Figure 3 illustrates how allowing this subtyping would be unsound.² When **breakit** is applied to an actual parameter of type *Point*, there is no problem. However if the actual parameter is a *ColorPoint*, a run-time error will occur when *p.equal(nuPt)* is evaluated. Since the value of *p* will be a *ColorPoint*, the code for *equal* in *ColorPointClass* will be executed. When *nuPt* is sent the message *c*, it will fail because it has no corresponding method. Thus, in a sound type system, a call of **breakit** with an actual parameter of type *ColorPoint* must not type check.

Most statically-typed object-oriented languages require subclasses to generate subtypes, even in the presence of binary methods. One type requirement that has been used to this end is that the types of methods may not be changed upon inheritance; this is done, for example, in C++ [51] Object Pascal [52], and Modula-3 [45]. In such languages, the one cannot write *ColorPointClass* as in Figure 2, with the typing discussed above.

Eiffel does allow these so-called “covariant” changes to the types of parameters redefined in subclasses, and compensates for the resulting insecurity in the type system by performing a link-time data-flow analysis of the program (called a system validity check) in order to catch possible type errors [41]. Under this view, the “subtype” relation has no clear meaning: Eiffel would claim *ColorPoint* to be a subtype of *Point*, but would not allow anything but a *Point* to be passed to **breakit**. So even though Eiffel judges *ColorPoint* to be a “subtype” of *Point*, *ColorPoint* objects cannot be used in all contexts where *Point* objects can be used.

Another problem related to the failure of $ColorPoint <: Point$ in the presence of the binary *equal* method is a methodological one. Suppose that version 0.9 of a commercial library includes *PointClass* and *ColorPointClass*, but that neither has an *equal* method yet. In this version of the library we have $ColorPoint <: Point$. Suppose also that the library includes a procedure *usePoint(p:Point)*, whose code of course does not involve equality.

Programmers using version 0.9 may take advantage of subsumption by treating a *ColorPoint*

²More on our notation: The expression **new** *OtypeClass*(*args*) results in the creation of a new object of the corresponding type whose instance variables are initialized to *args*. As a notational convention, we write *OtypeClass* for the class generating objects of type *Otype*.

object as a *Point*. For example, they might write:

```
var nuCpt: ColorPoint;
nuCpt := new ColorPointClass(8.7, 9.1, "red");
library.usePoint(nuCpt);
```

A few years later, Version 1.0 is released. Both *PointClass* and *ColorPointClass* have been extended to include the binary *equal* method. Now the inclusion *ColorPoint* <: *Point* fails, and client code that tries to call *usePoint* with a *ColorPoint* no longer typechecks.

Extending classes with new functionality is extremely common in large and evolving software systems. One often relies on the intuitive property that a simple extension will not cause existing client code to fail to typecheck. In the presence of binary methods this property fails. (Of course, other kinds of changes can also break client code. For example, one might happen to add a new attribute to a library class whose name already appears in a client subclass. But this kind of clash is much easier to repair.)

The *Point/ColorPoint* example illustrates some but not all of the problems that arise in typing binary methods in the presence of inheritance. Further examples that illustrate additional problems will be presented in the sections below.

2.2 Privileged Access to Object Representations

A completely different kind of problem with binary operations on objects—whether they are methods or free-standing procedures—is that they must often be given privileged access to the instance variables of *both* of their arguments.

The equality methods of points and colored points are examples of the simpler case where this need does not arise—the necessary attributes of the argument are already publicly available through existing methods. In order to write the *equal* method for the point class, we only needed to compare the receiver’s instance variables *xValue* and *yValue* to the values returned by the *x* and *y* methods of the argument *p*. There is no need to access *p*’s instance variables directly. Indeed, *p* might not even have instance variables named *xValue* and *yValue*; there is no need to know anything at all about its internal representation. The situation is similar for the *equal* of the colored point class.

On the other hand, suppose we want to write a class definition for simple integer set objects with the following interface type:

$$\begin{aligned} \text{IntSet} \equiv \langle\langle \text{add}: \text{int} \rightarrow \text{unit}; \text{member}: \text{int} \rightarrow \text{bool}; \\ \text{union}: \text{IntSet} \rightarrow \text{IntSet}; \text{superSetOf}: \text{IntSet} \rightarrow \text{bool} \rangle\rangle \end{aligned}$$

We can easily choose a representation for integer sets, (say, as lists of integers) and implement the *add* and *member* methods as in Figure 4. But when we come to implementing the *union* and *superSetOf* methods, we get stuck: given the interface type we have chosen for sets, there is no way to find out what elements a given set contains.

The obvious thing to do is to extend the public interface of sets with an *enumerate* method that (for example) returns a list of the elements of the set. But suppose we want to use a more efficient internal representation for sets, storing the elements in a balanced tree. We would certainly expect not only the *add* and *member* methods to be efficient, but *superSetOf* and *union* as well. But, to achieve good performance, *union* needs to work directly with the balanced tree representations of the two sets, so the *enumerate* method has to be replaced by an *asBalancedTree* method that returns the underlying representation. Doing so is unsatisfactory, because it makes representation details visible to users.

```

class IntSetClass
  instance variables
    elts: IntList
  methods
    add(i: int): unit is elts := elts.cons(i)
    member(i: int): bool is return(elts.memq(i))
    union(s: IntSet): IntSet is ???
    superSetOf(s: IntSet): bool is ???
end class

```

Figure 4: The class *IntSetClass*, for which writing *superSetOf* and *union* is problematic.

In order to handle binary operations like *equal* as methods, we needed a way of constraining the *type* of a parameter of a method to be the same as the type of the receiving object; for methods like *superSetOf* and *union*, we need an additional mechanism for constraining the *implementation* of a parameter to be the same as the receiver's. Indeed, such a mechanism is required whether or not we want to consider *union* a proper method of set objects: an external procedure for computing the union of two set objects will also need to gain privileged access to the internal representations of both of its arguments.

3 Avoiding Binary Methods

Sometimes the simplest solution to a problem is to ignore it. Here, one might take the position that binary operations like *equal*, *union*, and $+$ should not be regarded as *methods* of either of their argument objects, thus sidestepping the thorny typing issues raised so far.

There are some theoretical benefits to taking this step. For example, aside from binary methods, the types of methods are always *positive*, in the sense that the object type itself appears only in result positions. In this case, the classic encoding of object types as recursive records (as sketched above) may be replaced by an encoding where objects are modeled by existential types [48, 32, 21].

It may also be argued that keeping binary operations separate from their arguments avoids conceptual confusion. Turning a symmetric operation like $+$ into a method gives one of its arguments an artificially special status, requiring programmers to think in terms of contorted locutions like “Ask the number a to add itself to b and send back the result,” instead of the more straightforward “Compute the sum of a and b .” (However, having said this, it is only fair to give the methodological counterargument: An important property of objects is their appearance as active entities that encapsulate both data and the code acting on that data. Removing binary methods from objects disrupts this property, requiring an additional layer of module structure to encapsulate the binary methods with their class. Section 5 suggests that when binary methods require privileged access to both object states, such additional encapsulation may be needed anyway.)

With these arguments in mind, we consider in this section some alternatives to binary methods.

3.1 Using Functions Instead of Binary Methods

In languages that provide both objects and conventional procedural abstraction, an alternative to using binary methods is simply to make binary operations into functions. These binary functions

can be defined outside of classes, and can be applied to pairs of arguments as usual.

```
function eqPoint(p1,p2: Point): bool is return( (p1.x == p2.x) && (p1.y == p2.y));  
function eqColorPoint(cp1,cp2: ColorPoint):bool is  
    return( (cp1.x == cp2.x) && (cp1.y == cp2.y) && (cp1.c == cp2.c))
```

Ordinarily, one advantage of using methods instead of functions is dynamic dispatch: each class can choose its own code to execute in response to a given message. Therefore, moving from binary methods to binary functions may seem a step backwards. The programmer must now know when to apply `eqPoint` and `eqColorPoint`, instead of relying on the objects themselves “knowing” which equality is appropriate. But, when binary methods are used in conventional class-based object-oriented languages, the objects themselves often have no choice. Consider an arbitrary proper subclass of *PointClass*, such as *ColorPointClass*. Since *ColorPoint* has *equal* as a binary method, it is not a subtype of *Point*, as we have seen. To be a subtype of *Point* it would have to have an *equal* method whose argument type is *Point*. It is thus an operation on a *Point* and a *ColorPoint*, so two *ColorPoint* objects could not be properly compared. Thus it would seem that we can statically determine which equality test is performed. Consider a program variable `pt` of static type *Point*. What values could `pt` acquire during program execution? It could be a *Point* of course. It could also be a *ColorPoint*, but only in the case that *ColorPoint*’s *equal* method had argument type *Point*. Therefore, just from the static type of `pt`, we infer that the only binary equality that is soundly applicable to it is `eqPoint`. So in this example, there is no potential for dynamic dispatch of binary equality.

However, in general there are cases where dynamic dispatch on a binary method will occur; later, in Section 4.3, we present one such example.

The somewhat weaker “static dispatch” given by static overload resolution (as in C++ or Ada) may also be desirable: the programmer would like to be able to simply type *equal*, not *eqPoint* or *eqColorPoint*, even if the code that is invoked is statically determined.

A more serious problem with this approach is as follows: any method that in its body uses a binary method that is overridden cannot be properly inherited. This can give rise to unnecessary code duplication. Figure 5 is an example that illustrates this problem. *LinkClass* is a simple class of linked list objects, and *DoubleLinkClass* is a subclass that uses double links (a more complete implementation would include methods such as *reverse*, *map*, and *length*). The type *MyType* given to variables `next` and `link` in the example represents the type of objects of the current class. That is, it means *Link* in the class *LinkClass*, but means *DoubleLink* in the class *DoubleLinkClass* and in the instance variables and methods it inherits from *LinkClass*. *MyType* will be discussed in more detail in Section 4.1 below; also c.f. [41, 10, 27]. The objects now have only one interesting method, *append*, which is inherited by *DoubleLinkClass*. This method uses *setNext*, a binary method, to set the pointer `next`, and *setNext* is overridden in *DoubleLinkClass* to also properly maintain the `prev` link to the previous object. In a hypothetical function encoding, the *setNext* method would be replaced by functions `setNextLink` and `setNextDoubleLink` that lie outside the class definition (ignoring for now questions of privileged access). However, since *append* invokes *setNext*, it must be re-written as two almost identical functions, one invoking `setNextLink` and the other invoking `setNextDoubleLink`, causing unnecessary code duplication. An in-place *reverse* method of no arguments is another method for which inheritance would suffer under this encoding. Thus dispatch can be statically resolved, but only at the cost of code duplication if this scheme is used.

```

class LinkClass
  instance variables
    value: integer
    next: MyType
  methods
    getValue: integer is return(value)
    getNext: MyType is return(next)
    setValue(n: integer): unit is value := n
    setNext(link: MyType): unit is next := link
    append(link: MyType): unit is
      if next == nil then self.setNext(link) else next.append(link)
end class

class DoubleLinkClass subclass of LinkClass
  instance variables -- value and next are inherited
    prev: MyType
  methods -- getValue, getNext, setValue, and append are inherited; setNext is overridden
    getPrev: MyType is return(prev)
    setNext(link: MyType): unit is next := link; link.setPrev(self)
    setPrev(link: MyType): unit is prev := link
end class

```

Figure 5: The classes *LinkClass* and *DoubleLinkClass*.

3.2 Making Both Arguments into One Object

Even in a “purist” object-oriented language where every operation is treated as a message sent to some object, we may place binary operations outside of the objects on which they operate by turning the two argument objects into a single pair object and invoking the method on the pair. To see how this would work, imagine that the types *Point* and *ColorPoint* do not have any binary methods. For example, they could be:

$$\begin{aligned}
 \textit{Point} &\equiv \langle\langle x: \textit{real}; y: \textit{real} \rangle\rangle \\
 \textit{ColorPoint} &\equiv \langle\langle x: \textit{real}; y: \textit{real}; c: \textit{string} \rangle\rangle
 \end{aligned}$$

With this definition, *ColorPoint* would be a subtype of *Point*.

Now define two new classes, *PointPairClass* and *ColorPointPairClass*, with two methods named *equal*, as shown in Figure 6.

So the former binary methods are now unary methods of these new classes. What would originally have been written as:

```
aCPoint.equal(anotherCPoint)
```

to compare two colored points, will now be written with these new classes as:

```
(new PointPairClass(aCPoint, anotherCPoint)).equal
```

If the types of *aCPoint* and *anotherCPoint* are both *ColorPoint*, then one might wish instead to compare them as colored points, in which case one would write:

```

class PointPairClass
  instance variables
    p1: Point
    p2: Point
  methods
    equal: bool is return( (p1.x==p2.x) && (p1.y==p2.y) )

class ColorPointPairClass
  instance variables
    p1: ColorPoint
    p2: ColorPoint
  methods
    equal: bool is return( (p1.c==p2.c) && (p1.x==p2.x) && (p1.y==p2.y) )

```

Figure 6: The classes *PointPairClass* and *ColorPointPairClass*.

```
(new ColorPointPairClass(aCPoint, anotherCPoint)).equal
```

There is a benefit in making these pair objects: it clarifies the perspective desired for the equality comparison. When one creates a *PointPair* object, it is clear what behavior is expected from its *equal* method; this expectation is borne out even when the two points that make up the *PointPair* object are actually *ColorPoint* objects.

The types, *PointPair* and *ColorPointPair*, of objects of these pair classes, are as follows.

```

PointPair ≡ ⟨⟨equal: bool⟩⟩
ColorPointPair ≡ ⟨⟨equal: bool⟩⟩

```

Note that *ColorPointPair* is a subtype of the type *PointPair*. Therefore, one can have a variable *myPointPair* of type *PointPair* that denotes an object of type *ColorPointPair*.

```
var myPointPair: PointPair := new ColorPointPairClass(aCPoint, anotherCPoint)
```

In this case a message send such as “*myPointPair.equal*” results in the invocation of the *equal* method defined in class *ColorPointPairClass*. Thus, sending the *equal* message to a pair object gets the view with which the pair was created, regardless of the dynamic type of the points in the pair object. This should be contrasted with the function call “*eqPoint*(*aCPoint*, *anotherCPoint*),” which always compares its two arguments as points. It can also be contrasted with a message-send of the form “*aCPoint.equal*(*anotherCPoint*),” which always uses the *equal* code of *ColorPointClass*.

This approach has problems similar to the function approach that was discussed previously—the *LinkClass* example of Figure 5 would require code duplication for inherited methods such as *append*.

4 Embracing Binary Methods

Having presented some arguments that binary methods can be avoided, we now consider the typing mechanisms that must come into play if we choose *not* to avoid them.

Two important solutions have been proposed to the typing problems posed by binary methods. One solution, first proposed by the Abel project at HP labs [23], develops a method that partially solves the *Point/ColorPoint* problem by relaxing the requirement that subclasses generate subtypes. As they put it, “Inheritance is not subtyping.” They did not, however, propose a concrete mechanism for realizing their ideas in an object-oriented language. In Section 4.1, we show one way this may be done using the concept of *matching* [9].

The other important solution was presented in two papers at the 1991 OOPSLA conference [2, 28]. These papers deal with the static type-checking of languages with *multi-methods* (also called *generic functions* or *overloaded functions*). Multi-methods as in CLOS allow, as we show in Section 4.2.1, the *Point/ColorPoint* example to be typed preserving the subtyping of the two classes. But this is obtained at the expense of the encapsulation of the methods, since the generic functions, like the functions in Section 3.1, are separated from objects (objects encapsulate only data). In Section 4.2.2 we show how to reconcile multi-methods with objects encapsulating data and code [15, 44].

Closely related to the solutions of Section 4.2 is Ingalls’ solution to the multiple dispatch problem [33]. He presented his solution in an untyped framework, but it can be adapted to a typed language, as Section 4.3 shows.

Lastly, we show in Section 4.4 how a general principle of giving more “precise” types to binary methods produces more flexible typings across a range of approaches, even in the case where binary operations are not treated as methods.

4.1 Matching

This section describes how a relation called “matching,” which is weaker than subtyping, can replace subtyping in many situations [9]. In particular, we will see below that this generalization of subtyping provides us with the ability to handle binary methods smoothly.

4.1.1 Generalizing subtyping to matching

As seen in Section 2.1, languages that insist that subclasses generate subtypes often compensate for the resulting type problems by restricting the programmer’s ability to change the types of parameters of inherited methods. This effectively eliminates the use of binary methods in these cases. If one feels that binary methods are important, then an obvious solution is to give up the identification of subclasses with subtypes. An important advantage of this decision, not discussed further here, is to separate the notion of interface (type) from that of implementation (class). In the remainder of this section we assume such a separation has been made, and thus that the notions of subtyping and matching (defined below) depend only on the interfaces of objects, not the classes generating them.

Most object-oriented languages provide a name for the receiver of a message (e.g., **self** or **this**), which can be used inside method bodies. Similarly, we use *MyType* as a keyword that denotes the *type* of the receiver [50]. It may be used in the definition of methods whose parameters or return types should be the same as that of the receiver. By convention, the recursive record type in the following is simply an abbreviation for the type *Point* given above, and it could also be the type of objects of a polar implementation of points.

$$Point \equiv \langle\langle x: real; y: real; equal: MyType \rightarrow bool \rangle\rangle \equiv PolarPoint$$

One advantage of *MyType* is that it makes it easier for human readers to compare types like *Point* and *PolarPoint*. A more important advantage is that it works well with inheritance of

methods, because its meaning changes in the subclass. For example, when *MyType* is used in the definition of *ColorPointClass*, all occurrences of *MyType* in the methods automatically represent *ColorPoint* rather than *Point*.

$$\textit{ColorPoint} \equiv \langle\langle x:\textit{real}; y:\textit{real}; c:\textit{string}; \textit{equal}:\textit{MyType} \rightarrow \textit{bool} \rangle\rangle$$

Since this is just an abbreviation for the type of *ColorPoint* given earlier, the type *ColorPoint* is still not a subtype of *Point*. However, there is a relationship between the types *ColorPoint* and *Point*, which is clearly apparent when looking at their types written using *MyType*. One can see that the only difference is the addition of a new method, *c*, to *ColorPoint*.

We say one object type *matches* another if the first has at least the methods of the second and the corresponding method types are the same, considering *MyType* in one to be “the same” as *MyType* in the other. We use $\langle\#$ to denote this relationship. In symbols,

$$\langle\langle m_1:\tau_1; \dots; m_n:\tau_n \rangle\rangle \langle\# \langle\langle m_1:\tau_1; \dots; m_k:\tau_k \rangle\rangle$$

holds iff $k \leq n$. (In fact, a more general definition is possible in which the types of corresponding methods of the first are all subtypes of the corresponding types of the second [9]. This means that the corresponding result types are subtypes—vary in a covariant way—while the corresponding parameter types are supertypes—vary in a contravariant way. However, this more general relation will not be needed here.) Section 5 gives an interpretation of matching as subtyping on higher-order operators.

Because the meaning of *MyType* changes in subclasses, the meanings of the types of methods in subclasses need *not* be the same as those of the corresponding methods in the superclass. However, type-safe rules for defining subclasses can ensure that the types of the objects from the subclass always match the types of the objects generated from the superclass. In order to obtain type safety, it is necessary to type check the methods of a class under the assumption that *MyType* only matches the type of objects being defined by the class. This ensures that these methods will continue to be type-safe when inherited in subclasses [9]. While some routines will not type check with this assumption, even though they would have passed under the stronger assumption that *MyType* is exactly the type of objects generated by this class, in actual practice very few routines fail.

What about the earlier remark that *ColorPoint* is not a subtype of *Point*? Nothing has changed that fact. *ColorPoint* and *Point* provide an example of two types which match, but are not subtypes. Basically, if a class has a binary method, that is, a method with a parameter of type *MyType*, subclasses of that class that add new methods will not generate subtypes. On the other hand if a method’s return type is *MyType*, this will not stand in the way of subtyping. Both of these follow easily from the subtyping rule for recursive types mentioned in Section 2.1, and the fact that *MyType* is used as an abbreviation for a recursive definition of types.

What if we want to use a *ColorPoint* as an actual parameter in a procedure or function that originally expected a *Point* parameter? Since the example of **breakit** in the introduction showed this could not always be done, another, more restrictive, construct is needed.

We can introduce a language feature to support a form of bounded polymorphism using matching. With this feature, functions can be specified to take type parameters whose values are restricted to “match” another type. Of course, unrestricted type parameters can also be provided, but in a large number of situations some sort of restriction is necessary.

As an example, suppose we wish to write a routine to sort an array whose elements are drawn from some ordered set. In an object-oriented language, the requirement that the elements be ordered can be modeled by demanding that they support (at least) *less_than* and *equal* methods. Define:

$$\textit{Comparable} \equiv \langle\langle \textit{less_than}:\textit{MyType} \rightarrow \textit{bool}, \textit{equal}:\textit{MyType} \rightarrow \textit{bool} \rangle\rangle$$

With this definition, the header of our polymorphic sort routine is as follows, where the notation “ $T \triangleleft\# Comparable$ ” means that the type parameter T must match the type $Comparable$:

```
procedure sort( $T \triangleleft\# Comparable$ ; A: Array of  $T$ );
```

And the function then has type ³

```
sort :  $All(T \triangleleft\# Comparable)$  (Array of  $T$ )  $\rightarrow$  unit.
```

If `PhoneEntry` is an object type supporting at least methods `less_than` and `equal` of type

$$PhoneEntry \rightarrow bool,$$

and if `PArray` is an array of elements of type `PhoneEntry`, then `sort(PhoneEntry, PArray)` is a legal call of `sort`.

It is worth noting here that the type `Comparable` has no useful proper subtypes because of the appearance of `MyType` as the type of a parameter in its methods. Thus, if the bounds on type parameters were only expressed in terms of subtyping, it would be impossible to apply the `sort` routine to any interesting arguments.

The use of bounded matching is equivalent to the use of F-bounded polymorphism suggested in [11]. It is also very similar in effect to the restrictions on type parameters expressible in CLU and Ada (as well as the type classes of Haskell). For example, in Ada one would write the `sort` routine as:

```
generic
  type t is private;
  with function "<"(x,y: t) return BOOLEAN is <>;
  with function "="(x,y: t) return BOOLEAN is <>;
procedure sort (A: in out array (<>) of t) is ...
```

This is similar to the `sort` procedure written with bounded matching. Object-oriented languages containing similar constructs are Emerald [8] and Theta [38].

Returning to our example with `Point`, if `f(p: Point)` is a function accepting an argument of type `Point` then it can often be rewritten in the form `f(T $\triangleleft\# Point$; p:T)` so that it accepts a type parameter matching `Point` and an object of that type. If this type checks, then it will be possible to apply it to the type `ColorPoint` as well as an object of type `ColorPoint`. Of course this rewriting will not succeed in all cases—`breakit` being a prime example. The reason this transformation will not succeed in `breakit` is that the formal parameter `nuPt` will be of some type $T \triangleleft\# Point$, while `p` will always be of type `Point`. Thus we cannot guarantee that the type of the argument to `equal` in the body of `breakit` will be the same as the type of `p`, and the type check must fail.

What can actually be done with the information that one type matches another? The matching relation guarantees that certain messages may be sent to an object. If $T \triangleleft\# Comparable$ then objects of type T can be sent messages `less_than` and `equal` (and their parameters must also be of type T). It turns out that for most situations this is all that is needed in order to ensure that the object is usable. The stronger information that a type is actually a subtype of another is generally not needed.

In particular, bounded matching can be viewed as an explicit, weakened (and hence more generally applicable) form of subtyping. If subsumption were necessary to type a function call, the

³The notation $All(S \triangleleft\# T)E(S)$ is the universally polymorphic type that can be instantiated to $E(S)$, for all S such that $S \triangleleft\# T$.

code could be re-written so the function constrains the type parameter, like *Comparable* above, and function invocations explicitly pass the “smaller” type as argument. Simple subtyping is handled by the case where the type constraint contains no occurrences of *MyType*. The disadvantage of this encoding of subtyping is that all subtypings must be explicitly given in the program.

In general the use of bounded matching requires one to “plan ahead,” by identifying the type parameter to be matched against. This was illustrated in the `sort` example: the type *Comparable* needs to be discovered by the programmer, and every use of `sort` requires that an explicit type parameter be passed. This is in contrast to subtyping, which is implicit and, as mentioned above, does not require any explicit type instantiation to be given in the program. (It is an interesting open problem to show how to automatically infer declarations that use matching.) If we were to decide to eliminate subtyping altogether in favor of matching, then all object subtypings would have to be recast as bounded matchings. Moreover, since matching only applies to object types, we would not be able to capture the use of subtyping on other types without extending the definition of matching.

Another difficulty with relying only on matching is that it is not type-safe to perform an assignment to a variable of an object whose type only matches that of the variable. For example, imagine a framework for graphical user interfaces, in which one creates a main window as a subclass of some framework class, and has to store the window in some variable. In this case the type of the subclass objects has to be a subtype of the declared type of the variable in the framework. Subtyping seems to be required for this sort of cross-type assignment. While this can be worked around by using type parameters to designate the types of instance variables, it does limit flexibility in handling heterogeneous data structures, all of whose elements are subtypes of a given type.

The use of *MyType* is sufficient to write examples such as linked lists or trees, where methods for attaching a node to another or returning an adjoining node must be binary methods. The types of the instance variables of these nodes also can be written in terms of *MyType*. If the definition of singly-linked node is written using *MyType* in this way, it is easy to define a doubly-linked node as a subclass of singly-linked node. Figure 5 presents such an example. As expected, the type of a doubly-linked node is not a subtype of singly-linked node, but it does match. It is then relatively easy to write an implementation for lists which takes a type parameter which matches singly-linked node. By applying this to either the type for singly-linked node or doubly-linked node, the corresponding kind of list can be generated without code duplication. (See [10] for the details of this parameterized example.)

The use of *MyType* in class definitions makes it easier to write useful subclasses in statically typed object-oriented languages, especially when the superclasses contain binary methods. As illustrated in the sorting example above, the matching relation is very useful in defining bounded polymorphic functions. In fact, the use of these two features should provide a type-safe replacement for the (unsafe) uses of covariant typing in languages like Eiffel, while providing comparable expressiveness. The object-oriented language LOOP [27] has no matching relation *per se*, but has similar expressivity, achieved by circular subtype assertions $\tau <: \sigma$ where τ and σ may share free type variable X . This can be viewed as a form of operator subtyping $\tau(X) <: \sigma(X)$. Addition of a matching relation is thus one, but not the only, solution to this problem.

In the next two subsections we explore the mathematical aspects of the matching relation.

4.1.2 A Higher-Order Interpretation of Matching

We have seen that there exists a useful relation, matching, that holds between certain object types even when the subtyping relation does not hold between them. Unlike subtyping, which can be defined for all type constructions, matching is defined only for object types. Is matching, then,

an ad hoc relation, or can it be understood in some larger context?

We show in this section how matching can be seen as operator subtyping [13, 46, 48]. It also may be seen as a form of F-bounded subtyping [23, 27]. With respect to these interpretations, matching appears more elementary, and therefore is somewhat more appealing as a language construct. However, the type rules for matching are not trivially determined, and one can acquire confidence by deriving them from better-known constructions. This argument is explained in more detail in [1].

We assert that the relationship between *Point* and *ColorPoint* can best be seen by transforming them from types into type operators, as in [43, 48, 32]. The following definitions are obtained by uniformly abstracting the occurrences of *MyType* (or recursion variables) as type parameters.

$$\begin{aligned} \textit{PointOperator} &\equiv \lambda P : \text{Type}. \langle\langle x : \textit{real}; y : \textit{real}; \textit{equal} : P \rightarrow \textit{bool} \rangle\rangle \\ \textit{ColorPointOperator} &\equiv \lambda CP : \text{Type}. \langle\langle x : \textit{real}; y : \textit{real}; c : \textit{string}; \textit{equal} : CP \rightarrow \textit{bool} \rangle\rangle \end{aligned}$$

A type operator is a function from types to types—a function that, when applied to a type, yields another type. For example, we have the following:

$$\begin{aligned} \textit{PointOperator}(\textit{Point}) &\equiv \langle\langle x : \textit{real}; y : \textit{real}; \textit{equal} : \textit{Point} \rightarrow \textit{bool} \rangle\rangle \\ \textit{PointOperator}(\textit{ColorPoint}) &\equiv \langle\langle x : \textit{real}; y : \textit{real}; \textit{equal} : \textit{ColorPoint} \rightarrow \textit{bool} \rangle\rangle \end{aligned}$$

The original *Point* and *ColorPoint* types can be recovered by taking the fixed points of the respective operators⁴:

$$\begin{aligned} \textit{Point} &\equiv \text{Fix}(\textit{PointOperator}) \\ \textit{ColorPoint} &\equiv \text{Fix}(\textit{ColorPointOperator}) \end{aligned}$$

To define the operator translation we first need to extend subtyping from types to operators (we use the same symbol, $<$, for all these relations). Given two operators P and Q we define:

$$P <: Q \stackrel{\text{def}}{=} \text{for all types } X, P(X) <: Q(X).$$

We can verify that

$$\textit{ColorPointOperator} <: \textit{PointOperator}$$

because, for any type X :

$$\langle\langle x : \textit{real}; y : \textit{real}; c : \textit{string}; \textit{equal} : X \rightarrow \textit{bool} \rangle\rangle <: \langle\langle x : \textit{real}; y : \textit{real}; \textit{equal} : X \rightarrow \textit{bool} \rangle\rangle$$

The unknown type X here has exactly the same purpose as *MyType*—it is an undetermined type that is only a subtype of itself.

Thus, *ColorPoint* and *Point* are related by higher-order subtyping on their respective operators. The matching relation can be defined, in general, as follows⁵:

$$A <\# B \stackrel{\text{def}}{=} A\textit{Operator} <: B\textit{Operator}$$

Correspondingly, match-bounded quantification can be defined as:

$$\textit{All}(X <\# B)C(X_{ty}, X_{op}) \stackrel{\text{def}}{=} \textit{All}(F <: B\textit{Operator})C(\text{Fix}(F), F)$$

where $C(X_{ty}, X_{op})$ is a type where X_{ty} denotes occurrences of X that are used as types (as in $X \rightarrow \textit{Int}$), and X_{op} denotes ones used as operators (as in $\textit{All}(Y <\# X)D$); the former are replaced by fixed points of operators. Details appear in [1].

We can derive rules for matching from this higher-order interpretation, obtaining rules closely approximating those of PolyTOIL [10].

⁴The notation $\text{Fix}(F)$ means the least fixpoint of F .

⁵The notation $\textit{All}(S <: T)E(S)$ is the universally polymorphic type that can be instantiated to $E(S)$, for all S such that $S <: T$ [14].

4.1.3 Some Mathematical Weaknesses of Matching

We previously discussed a few practical weaknesses of matching, including having to plan ahead and difficulties with assignment. The matching relation also has some undesirable mathematical properties.

It is reasonable to allow recursive types to “unfold” by replacing *MyType* in an object type with the type itself. However, matching does not respect this operation: one can construct types T , T' and T'' such that $T \prec\# T'$ and T' unfolds to T'' , but it would be inconsistent to allow $T \prec\# T''$. The reason for this is matching is “really” a relation on type operators, not on types themselves, as suggested by the translation of the previous section. Inspection of the PolyTOIL rules in [10] shows that every time matching is used, the upper bound must explicitly be an object type, and this type can also be viewed as a type function applied to *MyType*.

Here is an example of the problems resulting from unrestricted use of unfolding of object types. Define:

$$\begin{aligned} F &= \lambda X : \text{Type}. \langle\langle p : X \rightarrow \text{Int}, q : \text{Int} \rangle\rangle \\ G &= \lambda X : \text{Type}. \langle\langle p : X \rightarrow \text{Int} \rangle\rangle \\ H &= \lambda X : \text{Type}. \langle\langle p : \text{Fix}(G) \rightarrow \text{Int} \rangle\rangle \end{aligned}$$

As before, the corresponding object types are found as fixed points: $Fobj = \text{Fix}(F)$, $Gobj = \text{Fix}(G)$, and $Hobj = \text{Fix}(H)$. Then $Fobj \prec\# Gobj$, and $Gobj \equiv Hobj$ by unfolding, but it is not the case that $Fobj \prec\# Hobj$. As a result of this problem with unrestricted unfolding, PolyTOIL does not have an “unfold” rule. Instead, objects are automatically unfolded when a message is sent to them.

4.2 Multi-methods

A different solution whereby binary methods can be embraced is to use multi-methods. Contrary to matching, this solution does not introduce a new relation on types, since with multi-methods one can have both type safety and subtyping relations such as $ColorPoint <: Point$.

A *multi-method* is a collection of method bodies associated with one message name. The selection of which method body to execute depends on the classes of one or more of the parameters of the method (rather than just on the class of the receiver as in ordinary object-oriented languages).

In this survey we distinguish two different kinds of multi-methods: the ones used by the language CLOS [24], and the *encapsulated multi-methods* of [15, 44]. A unified analysis of both kinds of multi-methods is given in [15]. We now describe each kind in turn.

4.2.1 Multi-methods à la CLOS

Intuitively the idea is to consider (multi-)methods (in CLOS jargon, *generic functions*) as global functions that are dynamically bound to different method bodies according to the classes of the actual arguments. An object does not encapsulate its methods, just the data (the instance variables). There no longer exists the notion of a privileged receiver for a method (the one that encapsulates it, usually denoted by **self** or **this**) since a multi-method is applied to several arguments that equally participate in the selection of the body. In this case we talk of “multiple dispatching” languages, in antithesis to “single dispatching” ones where a privileged receiver is used. A class of objects is then characterized just by the internal variables of its instances. For example, in a typed multi-method-based language, the classes given in Figures 1 and 2 would be defined as in Figure 7.

```

class PointClass
  includes
    xValue: real
    yValue: real
end class

class ColorPointClass subclass of PointClass
  includes -- xValue and yValue are inherited
    cValue : string
end class

method x(p: PointClass):real is return(p.xValue)

method y(p: PointClass):real is return(p.yValue)

method c(p: ColorPointClass):real is return(p.cValue)

method equal(p:PointClass,q:PointClass):bool is
  return( (x(p)==x(q)) && (y(p)==y(q)) )

method equal(p:ColorPointClass,q:ColorPointClass): bool is
  return( (x(p)==x(q)) && (y(p)==y(q)) && (c(p)==c(q)) )

```

Figure 7: *PointClass* and *ColorPointClass* written using multi-methods à la CLOS.

In order to simplify the exposition in this section, we identify classes and types, in the sense that the name of a class is used as the type of its instances; therefore in this section (and in this section only) $p : PointClass$ will *also* mean “ p is an instance of *PointClass*.” Thus, when discussing multi-methods à la CLOS, we write class names where types would otherwise appear.⁶ This allows one to consider multi-methods as overloaded functions, whose actual code is dynamically selected according to the type (i.e., the class) of the arguments they are applied to.

The definitions of the methods in Figure 7 are completely disconnected from those of classes. There are two distinct definitions for *equal*, one for arguments of types *PointClass* \times *PointClass* and the other for arguments of type *ColorPointClass* \times *ColorPointClass*. We say that the message *equal* denotes a multi-method (or a generic function, or an overloaded function) formed by two *branches* (or method bodies). The type of a multi-method is the set of the types of its branches; thus *equal* has type:

$$(PointClass \times PointClass \rightarrow bool, ColorPointClass \times ColorPointClass \rightarrow bool)$$

When *equal* is applied to a pair of arguments, the system executes the branch defined for those parameters whose type “best matches” the type of the arguments. For example if *equal* is applied to two arguments in which at least one of them is of type *PointClass* and the other is a subtype of it, then the first definition of *equal* is executed; if both arguments have as type a subtype of *ColorPointClass* then the second definition is selected. More generally, when a multi-method of

⁶If we were to distinguish between types and classes (i.e. between interfaces and implementations: cf 4.1.1), then a new notation would be needed to specify both a class and a type parameter for multi-methods. One possibility is to use the notation of Cecil [19, 20], which does separate these concepts.

type

$$(S_1 \rightarrow T_1, S_2 \rightarrow T_2, \dots, S_n \rightarrow T_n)$$

is applied to an argument of type S , the system executes the body defined for the parameter of type $S_j = \min_{i=1..n} \{S_i \mid S <: S_i\}$. This selection is performed at run-time. In this way one obtains dynamic dispatch. Note that in this paradigm binary methods are really binary, since the implicit argument given by the receiver of the message is, in this case, explicit.

In [17] it is proved that to have type safety it suffices that every multi-method of type $(S_1 \rightarrow T_1, S_2 \rightarrow T_2, \dots, S_n \rightarrow T_n)$ satisfies the following condition.⁷

$$\forall i, j \in [1..n] \quad \text{if } S_i <: S_j \text{ then } T_i <: T_j \tag{1}$$

(This is similar to the monotonicity condition of [49] and, independently, [40].) Note that all the multi-methods defined in Figure 7 (and in particular *equal*) satisfy this condition. Therefore *ColorPointClass* <: *PointClass* does not cause type insecurities.

Intuitively, the problem with binary methods is that, in general, it is not possible to choose the branch to execute according to the type of just one argument. To determine which method body must be executed one needs to know the types of all the arguments of the method. In single dispatching the branch selection is based only on one argument—the receiver; therefore binary methods and subtyping cannot be type safely combined. On the contrary by a multi-method we can refine the selection by considering all the arguments. Thus it need never happen that the argument of a method has a supertype of the type of the corresponding parameter (as in the case of *breakit*).

Note also that multi-methods allow one to specialize *equal* in a different way for each possible combination of arguments. It suffices to add the branches for the remaining cases:

```
method equal(p: PointClass, q: ColorPointClass):bool is ...
method equal(p: ColorPointClass, q: PointClass): bool is ...
```

As we have seen, CLOS’s multi-methods induce an object-oriented style of programming that is rather different from the one of traditional single dispatching object-oriented languages. Most of the languages that use multi-methods are untyped (e.g. CLOS [24], Dylan [5], which use classes instead of types to drive the selection of multi-methods). The only strongly-typed languages in our ken that use multi-methods are Cecil [20], and Polyglot [2].

The lack of encapsulation in multi-methods is both an advantage and a drawback. The drawback is methodological: an object (or a class of objects) is no longer associated with a fixed set of methods that have privileged access to its internal representation. The usual rule is that any method with a formal parameter of a given class can access the instance variables of the actual parameter object. The advantage is that this solves the privileged access problem described in Section 2.2, because a binary method can gain privileged access to both its arguments. However, because such methods can be defined anywhere in the program, one cannot restrict direct access to instance variables to a small area of the program text. One way to fix such problems may be to add a separate module system to control instance variable access [20]. Instead of pursuing that idea, in the next subsection, we show how to apply the ideas of multi-methods in more traditional object-oriented languages with single dispatching and classes.

Moreover, multi-method dispatch as in CLOS is more expensive than single-dispatch. This is because it is more expensive to compute the branch of the multi-method that matches the arguments

⁷Some further conditions are required to assure that a best matching branch always exists for the selection (see [2], [20], and [17]).

```

class ColorPointClass subclass of PointClass
  instance variables -- xValue and yValue are inherited
    cValue : string
  methods
    c:string is return(cValue)
    equal(p: Point):bool is return( (xValue==p.x) && (yValue==p.y) )
    equal(p: ColorPoint):bool is
      return( (cValue==p.c) && (xValue==p.x) && (yValue==p.y) )
end class

```

Figure 8: The class *ColorPointClass* written using encapsulated multi-methods.

best, whereas with single-dispatching, a single look-up suffices. A final drawback of multi-methods à la CLOS is the difficulty of combining independently developed systems of multi-methods [22]. While other ways to solve this problem have been studied [20], the problem nearly disappears when multi-methods are combined with single dispatching, as described next.

4.2.2 Encapsulated multi-methods

To solve the encapsulation problems of multi-methods à la CLOS, we seek to emulate the Smalltalk model, where every method is the method of one object. Thus each method has a privileged receiver argument (**self**), which is the only argument whose internal state can be accessed by the method. Instead of defining multi-methods as global functions, the idea is to use them to define the bodies of some methods in a class definition [15]. In this way a multi-method is always associated to a message m of a class C . When m is sent to an object of class C , it is dispatched to the corresponding method. If this method happens to be a multi-method, then the branch is selected according to the types of the further arguments of m . Thus, the selection of the method is still based on the receiver, but the actual code is selected among several bodies that are encapsulated inside the object. Inside these bodies, the receiver is still denoted by the keyword **self** (or **this**). Encapsulated multi-methods are to be distinguished from static overloading (as found in Ada, Haskell, C++, and other languages), because the selection of code must be made dynamically.

As an example of this technique, take the class *PointClass* as in Figure 1 and rewrite the class *ColorPointClass* as in Figure 8. In that Figure there are two definitions for *equal*: the first is executed when the argument of *equal* is of type *Point*, the other when it is of type *ColorPoint*. The selection of the appropriate definition is done at run-time when the argument of *equal* has been fully evaluated. The selection is based on the type of the additional argument. In other words, we have transformed the method associated to *equal* into a multi-method, where arguments of different types are associated to different codes.

There are two differences from multi-methods à la CLOS. The first is that multi-methods are defined in particular classes, whereas in CLOS they are a property of a global (generic function) names. This solves the encapsulation problems of CLOS multi-methods, because access to instance variables is restricted to the methods of a class, as only the receiver's instance variables can be accessed. The second difference is that dispatch is not based on actual argument classes, but rather on actual argument types. This is possible because no privileged access is obtained to the additional arguments. However, unless types are equated with classes (as in the previous section), the technique cannot solve the problem of privileged access to other arguments (of Section 2.2), be-

cause several different classes might implement the same type. (If a method is defined to work with objects of a given interface (type) and it is selected according to the type, it cannot access the internal implementation of these objects, since the same type may have many different implementations (classes) active in a program.)

The type of a multi-method is the set of the types of its different codes. Thus the type of an instance of *ColorPointClass* now becomes

$$\textit{ColorPoint} \equiv \langle\langle x: \textit{real}; y: \textit{real}; c: \textit{string}; \textit{equal}: (\textit{Point} \rightarrow \textit{bool}, \textit{ColorPoint} \rightarrow \textit{bool}) \rangle\rangle$$

and $\textit{ColorPoint} <: \textit{Point}$ holds, since in the type system for multi-methods (see [15]) one can deduce: $(\textit{Point} \rightarrow \textit{bool}, \textit{ColorPoint} \rightarrow \textit{bool}) <: (\textit{Point} \rightarrow \textit{bool})$.

More precisely, the subtyping relation between sets of types states that one set of types is smaller than another if and only if for every type contained in the latter there exists a type in the former smaller than it. This fits the intuition that one multi-method can be replaced by another multi-method of different type when for every branch that can be selected in the former there is one branch in the latter that can replace it (for the subtyping, ordinary methods are considered as multi-methods with just one branch: their type is a singleton set).

Thus, if in writing a subclass one wants the type of the instances to be a subtype of the type of the instances of the superclass, then some care in overriding binary methods is required. Indeed, the rule of the thumb for this approach is that to override a binary method one must use an (encapsulated) multi-method with (at least) two branches: one with a parameter whose type is the type of the instances of the class being defined, the other with a parameter whose type is the type of the instances of the original superclass in which the message associated with the binary method has been first defined. Thus, when a binary method is overridden in a new class, it is not enough to specify what the new method has to do with the objects of the new class. It is also necessary to specify what it has to do when the argument is an object of a superclass. Fortunately, this does not require a large amount of extra programming. The number of branches that suffice to override a binary (or n -ary) method is independent of both the size and the depth of the inheritance hierarchy; indeed, it is always equal to two. For example, suppose that we further specialize our *Point* hierarchy by adding further dimensions:

```
class 3DPointClass subclass of PointClass
  instance variables x3Value: real
  methods ...
```

```
class 4DPointClass subclass of 3DPointClass
  instance variables x4Value: real
  methods ...
```

...and so on, up to a dimension n . The new classes form a chain in the inheritance hierarchy. If we want to override *equal*, what do we have to do in order for this to be a chain of the subtyping hierarchy too (i.e., $nDPoint <: \dots <: 4DPoint <: 3DPoint <: Point$)? If we want to override *equal* in *nDPointClass* (thus we want that in the description of *nDPointClass* a definition of the form $\textit{equal}(\mathbf{p} : nDPoint) = \dots$ appears), then the first idea is to write for the class *nDPointClass* a multi-method with $n - 1$ branches, one for each class in the chain⁸. This is possible, but for type

⁸Of course, if in the definition of *nDPointClass* we do not give any definition for *equal* then *nDPointClass* inherits the last (multi-)method defined for *equal* in the upper hierarchy. It is important to be clear that, in the formalization we use, a new definition of a (multi-)method completely overrides the old one (i.e. it is not possible to inherit some branches and override others: this could be obtained by adding some extra syntax.)

safety a multi-method with two branches is enough: one for arguments of type $nDPoint$, which is the one we want to define, and the other for arguments of type $Point$, which will handle all the arguments of a supertype of $nDPoint$. Type safety stems from the observation that, according to the subtyping rule given above, if for all $i \in 1..n$, $T_i <: T$ then $(T \rightarrow S) <: (T_1 \rightarrow S, \dots, T_n \rightarrow S)$ (and of course $(T_n \rightarrow S, T \rightarrow S) <: (T \rightarrow S)$); take $Point$ for T , $bool$ for S and $iDPoint$ for T_i and the result follows.

A final remark is in order. The different branches that compose a single multi-method are not required to return the same type. For type safety it suffices to have the condition (1) as for multi-methods à la CLOS: for each pair of multi-method branches c_1, c_2 with the same name and number of arguments,⁹ if the parameter types of c_1 are smaller than the corresponding parameter types of c_2 , then the result type of c_1 must be smaller than the result type of c_2 [49, 17]. Some further consistency conditions are required in case of multiple inheritance [28, 44, 17, 20].

One of the main advantages of this approach is that the extra branch required to assure type safety of subtyping can be generated in an automatic way. Therefore this technique can be embedded directly in the technology of the compiler, and used to “patch” the already existing code of languages that use covariant specialization, like Eiffel and O₂ [6]. Thus, like the solution given in the next section, this solution can be directly applied to languages with covariant specialization without requiring any modification of the code: a recompilation of existing code will suffice (see [16]).

On the other hand this approach has some disadvantages. One disadvantage compared to multi-methods à la CLOS is that it does not solve the problem of obtaining privileged access to other arguments in a binary method. Another disadvantage of this approach is that in case of multiple inheritance additional type checking constraints are needed. The problem is that when multiple inheritance is used, the notion of a “best matching branch” to select or to inherit may be lost. Consequently, unconstrained use of multi-methods can break the modularity of programming [22], since the addition of a new class to the system might require the addition of some new code in a different class to assure the existence of the best branch (see, for example, [20]). However the problem with modularity is less critical than in the case of multi-method à la CLOS. An additional disadvantage is again the performance penalty imposed by multi-methods. One extra test and branch is required to decide which code is to be executed. The overhead to resolve uses of encapsulated multi-methods is however smaller than in the case of CLOS multi-methods since there is no special lookup needed for the privileged receiver.

There are also some less important disadvantages. The first one is that, as it depends on an *avant garde* type theory, the interactions of this theory with fairly standard features like polymorphism (both implicit and explicit) are not yet clear. (Models based on records have been more deeply studied than those based on overloading.) Also, with multiple inheritance, the automatic generation of the code to “patch” existing programs is not as satisfactory as in case of single inheritance. As a consequence, in some pathological cases the automatically generated code may not have the expected behavior. And, even though there is not a blowup of the number of extra method bodies that must be written, there is at least a doubling of the number of method bodies that must be written each time a binary method is overridden. Some further negative remarks are to be found at the end of the next section.

⁹Indeed multi-methods may have more than one parameter (this allows us to deal with n -ary methods), and the multi-method branches are not all required to have the same number of parameters.

```

class PointClass
  ...
  methods
    ...
    equal(p: Point): bool is return( p.equalPoint(self) )
    equalPoint(p: Point): bool is return( xValue==p.x && (yValue==p.y) )
    equalColorPoint(p: ColorPoint): bool is return( self.equalPoint(p) )

class ColorPointClass subclass of PointClass
  ...
  methods
    equal(p: Point): bool is return( p.equalColorPoint(self) )
    -- equalPoint is inherited
    equalColorPoint(p: ColorPoint): bool is
      return( (cValue==p.c) && (xValue==p.x) && (yValue==p.y) )

```

Figure 9: Ingalls' simulation of multi-methods.

4.3 Simulating Multi-methods in a Single-Dispatching Language

Ingalls offered a solution to what he called the problem of “multiple polymorphism” in the first OOPSLA conference [33]. His solution to the binary method problem, offered in the context of single-dispatching languages such as Smalltalk-80 [29], was to use two message dispatches, one to resolve the polymorphism of each argument.

In the example of points, colored points, and equality, the *equal* method would be coded as in Figure 9. As usual, the class *ColorPointClass* inherits the method *equalPoint* from the class *PointClass*. Now the (recursive) types of the instances of *PointClass* and *ColorPointClass* are:

```

Point ≡ ⟨⟨x: real; y: real;
  equal: Point → bool;
  equalPoint: Point → bool;
  equalColorPoint: ColorPoint → bool⟩⟩
ColorPoint ≡ ⟨⟨x: real; y: real; c: string;
  equal: Point → bool;
  equalPoint: Point → bool;
  equalColorPoint: ColorPoint → bool⟩⟩

```

Notice that, with this typing, *ColorPoint* is a subtype of *Point*. Also, *equal* in *ColorPoint* is a binary method, since by subsumption it can have argument type *ColorPoint* as well. This typing can be said to be more *precise* than the typing of *ColorPoint* given in the introduction; the general issue of the use of more precise typings is taken up in Section 4.4.

The solution offered by Ingalls is probably the best-known way to simulate multiple dispatch in a language with only single dispatch. With respect to true multiple-dispatch, the Ingalls simulation is a more exact simulation than the function simulation offered in Section 3.1, since it can arrange for *equal* with two arguments whose dynamic type is *ColorPoint* to always take color into account, regardless of the static types of the argument expressions. This is because of the second dynamic dispatch in the *equal* method. Such a result is not possible with the function simulation

of Section 3.1: one will always be able to apply the `eqPoint` function to two *ColorPoint* objects and lose exact type information. This example is thus one case for which dynamic dispatch on binary methods occurs. In this respect, the simulation of multiple-dispatch with external functions is less faithful and flexible than the Ingalls simulation.

This translation can also be contrasted with encapsulated multi-methods as described in Section 4.2.2. Ingalls’ translation lacks modularity in that it requires *equalColorPoint* to be added to the *PointClass* class when *ColorPointClass* is defined. With multi-methods, modularity can be preserved since the redefinition of the *equal* method inside *ColorPointClass* does not require any modification of the code for *PointClass*; however, this introduces an unnatural asymmetry, since the redefinition of *equal* requires one to write code for how a *ColorPoint* behaves when its *equal* method is passed a *Point*, but not vice-versa. The natural symmetry cannot be restored except by breaking the modularity of the multi-method solution.

It should be pointed out that the above argument only holds if we require (multi-)methods to be written in classes, as in Section 4.2.2. For multi-methods à la CLOS there is no problem of asymmetry, although there is still a modularity problem. However, the multi-method approach still requires one to go back and add code for types that appeared to have been completed earlier.

Ingalls’ solution is surprisingly general—by overriding *equalPoint* in *ColorPointClass*, a different method can be executed for all four combinations of *Point* and *ColorPoint*. Ingalls’ solution could in fact be used as one technique for implementing encapsulated multi-methods in a compiler, provided the compiler had access to all of the code at compilation time.

Finally, for large inheritance hierarchies the number of cases required by Ingalls’ solution can, in principle, become quite cumbersome.

4.4 Precise Typings

It is sometimes advantageous to use more precise typings for methods. A binary method only needs its argument to have the methods that are explicitly used. Generally this is a weaker requirement than having the argument be an object of the current class, and it may allow for a “larger” (with respect to the $<$: relation) type of the argument of this method; by the contravariant subtyping rule for functions this produces a smaller type for the method. The informal idea is thus to give methods smaller types [7, 8]. By subsumption, these types can always be lifted to “true binary” form, allowing objects of the same class to be passed as arguments to the method. Thus, specifying a smaller type of a method can only increase its usability.

Ingalls’ solution in Section 4.3 in fact depends on the use of precise types, for the key to its typability is the use of *Point* for the type of the argument of *equal* in *ColorPointClass*. This gives the method a smaller type than if the argument were of type *ColorPoint*. In this section we elaborate on this technique.

By way of illustration, consider the original *Point/ColorPoint* example of Figures 1 and 2. Since neither *equal* method calls *equal* recursively, the types

$$\begin{aligned} Point_{min} &\equiv \langle\langle x: real; y: real; equal: \langle\langle x: real; y: real \rangle\rangle \rightarrow bool \rangle\rangle \\ ColorPoint_{min} &\equiv \langle\langle x: real; y: real; c: string; equal: \langle\langle x: real; y: real; c: string \rangle\rangle \rightarrow bool \rangle\rangle \end{aligned}$$

may also be given. These types are subtypes of the types given originally. Note that the objects passed to *equal* themselves require no *equal* method be present. Since *Point_{min}* is a subtype of $\langle\langle x: real; y: real \rangle\rangle$ and similarly for *ColorPoint_{min}*, it is easy to see that

$$\begin{aligned} Point_{min} &<: \langle\langle x: real; y: real; equal: Point_{min} \rightarrow bool \rangle\rangle \\ ColorPoint_{min} &<: \langle\langle x: real; y: real; c: string; equal: ColorPoint_{min} \rightarrow bool \rangle\rangle \end{aligned}$$

```

class MPointClass
  ...
  methods
    ...
    max(p: MPoint): MPoint is return( p.maxMPoint(self) )
    maxMPoint(p: MPoint): MPoint is
      if xValue **2 + yValue **2 < p.x**2+ p.y**2
      then return(p) else return(self)
    maxColorMPoint(p: ColorMPoint): MPoint is return( self.maxMPoint(p) )
class ColorMPointClass subclass of MPointClass
  ...
  methods
    max(p: MPoint): MPoint is return( p.maxColorMPoint(self) )
    -- maxMPoint is inherited
    maxColorMPoint(p: ColorMPoint): ColorMPoint is
      if (xValue **2 + yValue **2) * brightness(cValue)
      < (p.x**2+ p.y**2) * brightness(p.c)
      then return(p) else return(self)

```

Figure 10: Ingalls simulation of points with a *max* method: first attempt.

so *equal* is indeed a binary method, and no typings are lost in this approach. In fact, something is gained over the matching interpretation described in Section 4.1: it is possible to invoke the *equal* method of a $Point_{(min)}$ with a $ColorPoint_{(min)}$ as argument. Typing this “heterogeneous” invocation is crucial for a class defining binary methods intended to be inherited without redefinition and able to take as arguments objects of any subclass [25]. In a type system based on matching, a method declared to take arguments of type *MyType* cannot, in general, accept an object of a subclass as argument; it is necessary to use bounded matching to realize this (see the discussion at the end of Section 4.1.1). Precise types here provide a simpler solution based on subtyping. Note that $ColorPoint_{min}$ does not match, nor is it a subtype of, $Point_{min}$.

Use of precise types also overcomes some problems with modularity. Continuing an example given in Section 2.1, the library procedure `usePoint`, which in version 0.9 did not use the *equal* method for Points, can be given a precise type in which its parameter is explicitly required to provide certain methods that Points have, but not *equal*. Thus the addition of *equal* to the methods of Point and ColorPoint in version 1.0 will have no effect on the applications of `usePoint` to objects of these classes — they will still type-check.

As shown in Section 4.3, typing Ingalls’ solution when *MyType* appears only in the types of method parameters is possible simply by using subsumption, *e.g.*, to lift a $ColorPoint$ to a $Point$ in `cp1.equal(cp2)`, where both `cp1` and `cp2` are objects of type $ColorPoint$. However this technique cannot be directly applied to binary methods with result of type *MyType* (or involving *MyType*), because subsumption on the type of the argument may cause loss of interesting type information. Consider the example in Figure 10, defining classes of points and colored points with a method *max* which among its argument and the receiver (**self**) returns the one further from the origin. Objects of *MPointClass* and *ColorMPointClass* could be given the following types, which are simple modifications of the types of *Points* in Section 4.3.

```

class MPointClass
  ...
  max[X:Type](p:⟨⟨maxMPoint:MPoint → X; maxColorMPoint:ColorMPoint → X⟩⟩):X is
    return( p.maxMPoint(self) )
  ...
class ColorMPointClass subclass of MPointClass
  ...
  max[X: Type](p: ⟨⟨maxColorMPoint:ColorMPoint → X⟩⟩): X is
    return( p.maxColorMPoint(self) )
  ...

```

Figure 11: Ingalls simulation of points with a *max* method: precise typing.

```

MPoint ≡ ⟨⟨x:real; y:real;
  max:MPoint → MPoint;
  maxMPoint:MPoint → MPoint;
  maxColorMPoint:ColorMPoint → MPoint⟩⟩

```

```

ColorMPoint ≡ ⟨⟨x:real; y:real; c:string;
  max:MPoint → MPoint;
  maxMPoint:MPoint → MPoint;
  maxColorMPoint:ColorMPoint → ColorMPoint⟩⟩

```

The subtyping $ColorMPoint <: MPoint$ still holds, but note that the result of method *max* of *ColorMPointClass* is of type *MPoint*; type checking would fail if we assigned this method the type $MPoint \rightarrow ColorMPoint$. Thus the static type of taking the *max* of two *ColorMPoints* will have to be merely *MPoint* (unless the method *maxColorMPoint* was used explicitly). True multi-methods do not suffer from this shortcoming.

We can overcome this problem in a more expressive type system that provides for polymorphism in addition to recursive types. The idea is to make the type of *max* more precise, and in this case, polymorphic. The code for the *max* methods with their new type annotations is given in Figure 11. This modification yields the following types for the objects of *MPointClass* and *ColorMPointClass*.

```

MPoint ≡ ⟨⟨x:real; y:real;
  max:All(X)⟨⟨maxMPoint:MPoint → X; maxColorMPoint:ColorMPoint → X⟩⟩ → X;
  maxMPoint:MPoint → MPoint;
  maxColorMPoint:ColorMPoint → MPoint⟩⟩

```

```

ColorMPoint ≡ ⟨⟨x:real; y:real; c:string;
  max:All(X)⟨⟨maxColorMPoint:ColorMPoint → X⟩⟩ → X;
  maxMPoint:MPoint → MPoint;
  maxColorMPoint:ColorMPoint → ColorMPoint⟩⟩

```

If *p* is a *MPoint*, its *max* method can still be specialized to a binary method: *p.max*[*MPoint*] is of type $MPoint \rightarrow MPoint$, and similarly for the *max* method of a *ColorMPoint*. The relation with the types of the “true binary” methods is more direct in an implicitly typed language, where the precise types are smaller [26, 25].

With this typing, taking the *max* of two elements of *ColorMPoint* returns a *ColorMPoint*; any other combination returns a *MPoint*, the best static type possible. Note that *ColorMPoint* is still

a subtype of *MPoint* in a system with implicit unfolding of recursive types. So, this typing has all the desired properties of the typing via pure multi-methods of Section 4.2, giving more situations in which Ingalls’ method may be usefully applied.

SOOP and PolyTOIL are two languages in which all of the precise typings of this section may be expressed. Precise types are complex, however, and it is difficult to imagine programmers writing them routinely. A solution to this problem is to automatically infer *minimal* types. See [25] for a type inference algorithm for the I-LOOP object-oriented language. The algorithm infers a form of F-bounded polymorphic type for classes and objects. It infers minimal types for the original *Point/ColorPoint* example that are very similar to the “small” types presented above. The types inferred for objects of *MPointClass* and *ColorMPointClass* are slightly more general than the form above.

To summarize some of the advantages and disadvantages of precise typing:

- + Precise types allow more flexibility in typing than matching alone. They may be expressed using bounded matching, but bounded matching requires explicit quantification and instantiation where subtyping alone may suffice.
- + Precise types are a critical component of a typed version of Ingalls’ solution.
- + More precise types in module interfaces can be used to overcome some of the limitations of matching.
- The generally more complicated form of the precise types suggests that a type inference algorithm may be the only practical alternative.
- In defining a subclass, one may have to go back and modify the type annotations of the superclass (and, in general, the superclass of the superclass, etc.) in order to generate subtypes. This may be seen as another argument in favor of type inference, since no modifications will be required in an implicitly typed language.

5 Privileged Access to Object Representations

In Section 2.2, we saw that the problems of typing binary methods are often accompanied by difficulties in implementing binary operations without exposing object internals to public view. This section sketches a technique whereby such “overexposed objects” can be wrapped in an additional layer of abstraction, creating a limited scope in which their internal structure is visible. The technique was developed by Pierce and Turner [47] and by Katiyar, Luckham, and Mitchell [34]; we refer the reader to these papers for further details. In particular, [47] demonstrates that the mechanism shown here is compatible with inheritance (though it requires some additional machinery). These ideas give a semantic basis for some aspects of the encapsulation via friends found in C++ and the encapsulation in Cecil [19]. Returning to the example of integer set objects (and dropping the *union* method, for brevity), it is clear that the typing

$$IntSet \equiv \ll add: int \rightarrow unit; member: int \rightarrow bool; superSetOf: IntSet \rightarrow bool \gg$$

does not provide a sufficiently rich protocol to allow the *superSetOf* method to be implemented: there is no way to find out what are the elements of the other set (the one provided as argument to *superSetOf*). We have no choice but to extend the interface of set objects with a method that

```

class IntSetExposedClass
  instance variables
    elts: IntList
  methods
    add(i: int): unit is elts := elts.cons(i)
    member(i: int): bool is return(elts.memq(i))
    superSetOf(s: IntSet): bool is return(elts.superListOf(s.rep))
    rep: IntList is return(elts)
end class

```

Figure 12: The class *IntSetExposedClass*, for which writing *superSetOf* is straightforward

provides access to this information; let us call it *rep*, as a reminder that, in general, it may need to provide access to the whole internal representation of the object.

$$\text{IntSetExposed} \equiv \langle\langle \text{add}: \text{int} \rightarrow \text{unit}; \text{member}: \text{int} \rightarrow \text{bool}; \text{superSetOf}: \text{IntSetExposed} \rightarrow \text{bool}; \text{rep}: \text{IntList} \rangle\rangle$$

Now we can easily implement all the methods of *IntSetExposedClass*, as shown in Figure 12.

It remains to show how to package the class *IntSetExposedClass* so that the *rep* method can only be called by other instances of the same class. For this, we generalize Mitchell and Plotkin’s motto that “abstract types have existential type” [42], combining it with the idea of object interfaces as type operators from Section 4.1.2 and Cardelli and Wegner’s *partially abstract types* [14].

Using the notation from Section 4.1.2, the interface of exposed integer set objects can be written:

$$\text{IntSetExposedOperator}(S) \equiv \langle\langle \text{add}: \text{int} \rightarrow \text{unit}; \text{member}: \text{int} \rightarrow \text{bool}; \text{superSetOf}: S \rightarrow \text{bool}; \text{rep}: \text{IntList} \rangle\rangle$$

Similarly, the interface of ordinary integer set objects (without *rep*) can be written:

$$\text{IntSetOperator}(S) \equiv \langle\langle \text{add}: \text{int} \rightarrow \text{unit}; \text{member}: \text{int} \rightarrow \text{bool}; \text{superSetOf}: S \rightarrow \text{bool} \rangle\rangle$$

Now comes the key point. Instead of defining $\text{IntSet} = \text{Fix}(\text{IntSetOperator})$ as we did before, we build an abstract data type (ADT) and then open it to obtain *IntSet*. The implementation of the ADT uses *IntSetExposedOperator*, so that *superSetOf* makes sense, but the *rep* method is hidden from public view.

The integer set package (or module) is defined in Figure 13. To verify that its type is

$$\text{intSetPackage} : \text{Some}(\text{ISOp} <: \text{IntSetOperator}) \langle\langle \text{newIntSet} : \text{Fix}(\text{ISOp}) \rangle\rangle$$

we need only check that when the hidden “witness type” *IntSetExposedOperator* is replaced by the abstract placeholder *ISOp* in the type of the body of the package

$$\langle\langle \text{newIntSet} : \text{Fix}(\text{IntSetExposedOperator}) \rangle\rangle$$

we obtain the body of the abstract type:

$$\langle\langle \text{newIntSet} : \text{Fix}(\text{ISOp}) \rangle\rangle.$$

```

intSetPackage =
  pack
    procedure newIntSet() is
      var
        nuSet: Fix(IntSetExposedOperator)
      begin
        nuSet := new IntSetExposedClass();
        return(nuSet)
      end
    as
      Some(ISOp <: IntSetOperator) «newIntSet : Fix(ISOp)»
    hiding
      IntSetExposedOperator
    end

```

Figure 13: The package `intSetPackage`.

Having built `intSetPackage`, we can open it to obtain the creation procedure `newIntSet` and the abstract interface `ISOp`, from which we define the type `IntSet`:

```

open intSetPackage
as ISOp with «newIntSet»

type IntSet = Fix(ISOp)

```

In the remainder of the program, objects created using `newIntSet` have type `IntSet`. In particular, they can be sent the `superSetOf` message.

In effect, what we have accomplished is to blend object- and ADT-style abstraction mechanisms. The primary mechanism is objects: both ordinary (unary) operations like `add` and binary operations like `superSetOf` are methods of objects rather than free-standing procedures. The extra layer of packaging guarantees that elements of `IntSet` can only be created by calling `newIntSet`—i.e., that every element of `IntSet` is actually an instance of `IntSetExposedClass`, and hence supports the `rep` message.

6 Binary Methods and Behavioral Specification

Although our example of `ColorPoint`, `Point`, and their `equal` methods is standard, it can be criticized on methodological grounds. To understand the criticism, consider one of two ways one might possibly specify the behavior of the type `Point`, which is shown in Figure 14. In the figure, the abstract values [31] of the type `Point` are defined in the Larch Shared Language (LSL) [30] trait “PointTrait” of Figure 15. The postconditions of the methods follow the keyword **ensures**. They are stated using the mathematical vocabulary defined in PointTrait, and the keyword **result**, which stands for the value returned by the method. With this specification, one can verify that objects of class `PointClass` in Figure 1 correctly implement the specification of the type `Point`.

```

type specification Point
  uses PointTrait
  methods
    x: real
      ensures result = xCoord(self)
    y: real
      ensures result = yCoord(self)
    equal(p: Point): bool
      ensures result  $\equiv$  (xCoord(self) = xCoord(p)
                              $\wedge$  yCoord(self) = yCoord(p))
  end type specification

```

Figure 14: A first specification of the behavior of the type *Point*, which has a strong specification of the *equal* method.

```

PointTrait : trait
  includes Rational(real for Q)
  introduces
    makePoint: real, real  $\rightarrow$  Point
    xCoord, yCoord: Point  $\rightarrow$  real
  asserts Point generated by makePoint
   $\forall x, y$ :real
    xCoord(makePoint(x,y)) == x
    yCoord(makePoint(x,y)) == y

```

Figure 15: The trait *PointTrait*, which defines a mathematical model of points.

```

function testPointEqual(p1: Point, p2: Point): bool is
  return( ((p1.x == p2.x) && (p1.y ==p2.y))
           == p1.equal(p2) )

```

Figure 16: The function `testPointEqual`.

The mechanisms of Sections 4.2, 4.3, and 4.4, above allow *ColorPoint* to be a subtype of *Point*. Suppose that the code in Figure 8 correctly implements a specification of the type *ColorPoint*. Then by virtue of the subtyping, one can use instances of *ColorPointClass* where objects of type *Point* are required. For example, one could call the function `testPointEqual` given in Figure 16, and pass it two instances of *ColorPointClass*. Looking just at the specification of the type *Point* and the code in the function `testPointEqual`, one would conclude that it should always return *true*. However, if one passes it two instances of *ColorPointClass* with the same *x* and *y* coordinates, but with different colors, then it will return *false*. This happens even if the type theory, as in Section 4.2, allows the type *ColorPoint* to be a subtype of *Point*. After briefly introducing some terminology, we return to the problem of what to do about such examples below.

It is useful to distinguish among three relationships: subclassing (or inheritance), subtyping, and behavioral subtyping. We have consistently distinguished subclasses from subtypes in this paper. For example, we noted that, without multi-methods, *ColorPoint* is not a subtype of *Point*, even when *ColorPointClass* is a subclass of *PointClass*.

A subtype relationship means that objects of the subtype can be used in any situation where objects of the supertype can be used, without the possibility of encountering type errors. Behavioral subtyping is a stronger relationship than subtyping, and, in addition to guarantees about lack of type errors, makes behavioral guarantees.

In order to decide when one ADT¹⁰ is a behavioral subtype of another, one needs behavioral specifications of the types in question. A type *S* is a *behavioral subtype* of *T* [3, 4, 36, 35, 39, 37] when objects of type *S* can be manipulated according to the specification of type *T*, and when this is done, they behave as dictated by the specification of type *T*. For example, one can compare the types *Point* as specified in Figure 14 and *ColorPoint* as specified in Figure 17 (and the trait in Figure 18). For these two types, *ColorPoint* is not a behavioral subtype of *Point*. The reason is that *ColorPoint* objects can be observed to differ from the specification in Figure 14; one such observation is given by the function `testPointEqual` of Figure 16. (If one passes two *ColorPoint* objects with the same *x* and *y* coordinates, but with different colors, to this function, it returns *false*, which cannot happen according to the specification of the type *Point*.) Hence *ColorPoint* is not a behavioral subtype of *Point* as specified in Figure 14.

If an ADT *S* has an interface that is not a subtype of *T*'s, then objects of type *S* cannot be manipulated as objects of type *T*, and so *S* cannot be a behavioral subtype of *T*. However, in Sections 4.2, 4.3, and 4.4, (the interface type) of *ColorPoint* is a subtype of *Point*. Hence this is an example of two types that are in a subtype relationship but not a behavioral subtype relationship.

There are (at least) two conclusions one can draw from such examples. The first is that struc-

¹⁰In this section we use the term “ADT” with a different meaning than earlier. Before we discussed ADT-style encapsulation, by which we meant encapsulation enforced by a module system or a type system as in Ada or CLU. In this section ADT means what it does in software engineering: a set of objects whose behavior is characterized by a specification. The behavioral specifications of types such as record and function types are fixed by the semantics of the programming language, and hence we use the word “types” in this section to mean both ADTs and composite or higher-order types built from them.

```

type specification ColorPoint
uses ColorPointTrait
methods
  x: real
    ensures result = xCoord(self)
  y: real
    ensures result = yCoord(self)
  c: string
    ensures result = colorOf(self)
  equal(p: ColorPoint): bool
    ensures result  $\equiv$  (xCoord(self) = xCoord(p)
                           $\wedge$  yCoord(self) = yCoord(p)
                           $\wedge$  colorOf(self) = colorOf(p))
end type specification

```

Figure 17: A specification of the behavior of the type *ColorPoint*.

```

ColorPointTrait : trait
includes Rational(real for Q), String(string for C)
introduces
  makeColorPoint: real, real, string  $\rightarrow$  ColorPoint
  xCoord, yCoord: ColorPoint  $\rightarrow$  real
  colorOf: ColorPoint  $\rightarrow$  string
asserts ColorPoint generated by makeColorPoint
   $\forall x, y:\text{real}, s:\text{string}$ 
  xCoord(makeColorPoint(x,y,s)) == x
  yCoord(makeColorPoint(x,y,s)) == y
  colorOf(makeColorPoint(x,y,s)) == s

```

Figure 18: The trait *ColorPointTrait*, which defines a mathematical model of colored points.

```

type specification Point
  uses PointTrait
  methods
    x: real
      ensures result = xCoord(self)
    y: real
      ensures result = yCoord(self)
    equal(p: Point): bool
      ensures result  $\Rightarrow$  (xCoord(self) = xCoord(p)
         $\wedge$  yCoord(self) = yCoord(p))
end type specification

```

Figure 19: A second specification of the behavior of the type *Point*, which has a weaker specification of the *equal* method.

tural type information alone is not enough to decide behavioral subtype relationships, especially for ADTs. For example, if one intends *ColorPoint* and *Point* to be ADTs, and not simply record types, then different type rules might be desired. Since ADTs are defined by type specifications, one might not want their types to be identified with the types of their interfaces (a record type for their operations). Then one would base type checking for ADT objects on the name of the ADTs, and not on the structure of the interface record type. Similarly, subtyping for ADTs might be based on declared subtype relationships. (This is done, for example, in the subtyping rule for the witness types of existential types [14, 12].) The advantage of declaring subtype relationships for ADTs is that the type system can be prevented from using a subtype relationship that the programmer knows is not a behavioral subtype relationship. That is, if the programmer can declare subtype relationships among ADTs, the programmer can use the type system to enforce behavioral subtyping, and thus the type system can aid in reasoning about programs that use behavioral subtyping [36, 35, 37]. The disadvantage, however, is that the programmer must then declare all subtype relationships among ADTs: the language would not be allowed to infer any.

The second conclusion is that one should think about potential behavioral subtypes when specifying types, especially those with binary operations. For example, if the type *Point* is specified as in Figure 19, then the *ColorPoint* specified in Figure 17 can be a behavioral subtype (when the type theory allows it to be a subtype). This is because of the weaker specification of the *equal*, which only requires that when the result is *true*, that the *x* and *y* coordinates be equal, not vice versa. Hence a behavioral subtype, such as *ColorPoint*, can have an *equal* method that takes extra information (such as color) into account.

The disadvantage of such a weak specification of *Point* and its method *equal* is that it may be too weak to prove some desired properties of programs. For example, assuming extra methods to mutate points, one might write a loop that compares a given point to the origin; if the two points are really *ColorPoint* objects, then the comparison may fail even though the *x* and *y* coordinates are the same.

One can also examine the question of whether to avoid binary methods from the perspective of behavioral subtyping. Depending on how one writes type specifications, having binary methods may limit behavioral subtyping, whereas not having binary methods, such as *equal*, allows more behavioral subtype relationships. Thus, it might sometimes be wise to avoid binary methods as a

way of promoting behavioral subtyping among ADTs.

7 Summary and Conclusions

Binary methods pose real problems in object-oriented programming languages. There is a typing problem because types with binary methods have few interesting subtypes, and there is a problem obtaining privileged access to additional arguments in binary methods. Because of the latter problem, it is not always possible to avoid binary methods.

We discussed the following solutions to the typing problem for binary methods.

- Using a notion of matching, which is weaker than subtyping. This allows more polymorphism in the presence of types with binary methods. However, it seems to require programmers to plan ahead more than they would using subtyping, and its flexibility is not as great as with multi-methods.
- Using multi-methods, either as a basis for object-oriented programming, or as a solution within the framework of single-dispatched languages. This gives the programmer more flexibility in programming binary methods, and consequently allows more subtyping. However, there are modularity and efficiency problems with these approaches.
- Using more precise typings for methods (including the Ingalls simulation of multi-methods). This allows more flexibility than matching. However, it seems to require type inference to be practical [25], and the resulting types may be more complicated than programmers want to see.

To solve the problem of privileged access to additional arguments, we discussed adding additional layers of abstraction. However, this by itself does not seem to solve the typing problems of binary methods, and so would have to be combined with one of the previous solutions.

So, which solution is the best? If the authors agreed on the answer, this paper would not have been written to begin with. None of the solutions is perfect, and for practical programming languages the question may boil down to what sort of inconvenience the programmer is most likely to tolerate. Some work also remains in determining if some of the solutions will scale up to full-featured languages.

Acknowledgements

Thanks to the US National Science Foundation and ESPRIT for their support of the workshop that resulted in this paper.

References

- [1] Martín Abadi and Luca Cardelli. On subtyping and matching. In *Proceedings ECOOP '95*, 1995. To appear.
- [2] Rakesh Agrawal, Lindga G. DeMichiel, and Bruce G. Lindsay. Static type checking of multi-methods. *ACM SIGPLAN Notices*, 26(11):113–128, November 1991. OOPSLA '91 Conference Proceedings, Andreas Paepcke (editor), October 1991, Phoenix, Arizona.

- [3] Pierre America. Inheritance and subtyping in a parallel object-oriented language. In Jean Bezivin et al., editors, *ECOOP '87, European Conference on Object-Oriented Programming, Paris, France*, pages 234–242, New York, NY, June 1987. Springer-Verlag. Lecture Notes in Computer Science, Volume 276.
- [4] Pierre America. Designing an object-oriented programming language with behavioural subtyping. In J. W. de Bakker, W. P. de Roever, and G. Rozenberg, editors, *Foundations of Object-Oriented Languages, REX School/Workshop, Noordwijkerhout, The Netherlands, May/June 1990*, volume 489 of *Lecture Notes in Computer Science*, pages 60–90. Springer-Verlag, New York, NY, 1991.
- [5] Apple Computer Inc., Eastern Research and Technology. *Dylan: an object-oriented dynamic language*, April 1992.
- [6] François Bancilhon, Claude Delobel, and Paris Kanellakis (eds.). *Implementing an Object-Oriented database system: The story of O₂*. Morgan Kaufmann, 1992.
- [7] Andrew Black, Norman Hutchinson, Eric Jul, and Henry Levy. Object structure in the Emerald system. *ACM SIGPLAN Notices*, 21(11):78–86, November 1986. OOPSLA '86 Conference Proceedings, Norman Meyrowitz (editor), September 1986, Portland, Oregon.
- [8] Andrew P. Black and Norman Hutchinson. Typechecking polymorphism in Emerald. Technical Report CRL 91/1 (Revised), Digital Equipment Corporation, Cambridge Research Lab, Cambridge, Mass., July 1991.
- [9] Kim Bruce. A paradigmatic object-oriented programming language: design, static typing and semantics. *Journal of Functional Programming*, 4(2):127–206, 1994.
- [10] Kim B. Bruce, Angela Schuett, and Robert van Gent. PolyTOIL: A type-safe polymorphic object-oriented language. In *Proceedings ECOOP '95*, 1995. To appear.
- [11] Peter Canning, William Cook, Walter Hill, Walter Olthoff, and John Mitchell. F-bounded quantification for object-oriented programming. In *Fourth International Conference on Functional Programming Languages and Computer Architecture*, pages 273–280, September 1989.
- [12] Luca Cardelli. A semantics of multiple inheritance. In G. Kahn, D. MacQueen, and G. Plotkin, editors, *Semantics of Data Types*, volume 173 of *Lecture Notes in Computer Science*, pages 51–67. Springer-Verlag, 1984. Full version in *Information and Computation* 76(2/3):138–164, 1988.
- [13] Luca Cardelli. Notes about F_{\leq}^{ω} . Unpublished manuscript, October 1990.
- [14] Luca Cardelli and Peter Wegner. On understanding types, data abstraction, and polymorphism. *Computing Surveys*, 17(4), December 1985.
- [15] Giuseppe Castagna. Covariance and contravariance: conflict without a cause. *ACM Transactions on Programming Languages and Systems*, 17(3), 1995.
- [16] Giuseppe Castagna. A proposal for making O₂ more type safe. Technical Report LIENS-95-4, Laboratoire d'Informatique de l'École Normale Supérieure, February 1995. To appear in *BDA '95*. Currently available by anonymous ftp from `ftp.ens.fr` in file `/pub/dmi/users/castagna/o2.dvi.Z`.

- [17] Giuseppe Castagna, Giorgio Ghelli, and Giuseppe Longo. A calculus for overloaded functions with subtyping. *Information and Computation*, 117(1):115–135, February 1995. A preliminary version appeared in *ACM Conference on LISP and Functional Programming*, June 1992 (pp. 182–192).
- [18] Giuseppe Castagna and Gary T. Leavens. Foundation of object-oriented languages: 2nd workshop report. *SIGPLAN Notices*, 30(2):5–11, February 1995.
- [19] Craig Chambers. Object-oriented multi-methods in Cecil. In Ole Lehrmann Madsen, editor, *ECOOP '92, European Conference on Object-Oriented Programming, Utrecht, The Netherlands*, volume 615 of *Lecture Notes in Computer Science*, pages 33–56. Springer-Verlag, New York, NY, 1992.
- [20] Craig Chambers and Gary T. Leavens. Typechecking and modules for multi-methods. *ACM SIGPLAN Notices*, 29(10):1–15, October 1994. OOPSLA '94 Conference Proceedings, October 1994, Portland, Oregon.
- [21] Adriana B. Compagnoni and Benjamin C. Pierce. Multiple inheritance via intersection types. *Mathematical Structures in Computer Science*, 1995. To appear. Preliminary version available as University of Edinburgh technical report ECS-LFCS-93-275 and Catholic University Nijmegen computer science technical report 93-18, Aug. 1993.
- [22] William R. Cook. Object-oriented programming versus abstract data types. In J. W. de Bakker, W. P. de Roever, and G. Rozenberg, editors, *Foundations of Object-Oriented Languages, REX School/Workshop, Noordwijkerhout, The Netherlands, May/June 1990*, volume 489 of *Lecture Notes in Computer Science*, pages 151–178. Springer-Verlag, New York, NY, 1991.
- [23] William R. Cook, Walter L. Hill, and Peter S. Canning. Inheritance is not subtyping. In *Proc. 17th ACM Symp. on Principles of Programming Languages*, pages 125–135, January 1990.
- [24] L.G. DeMichiel and R.P. Gabriel. Common Lisp Object System overview. In Bézivin, Hullot, Cointe, and Lieberman, editors, *Proc. of ECOOP '87 European Conference on Object-Oriented Programming*, number 276 in LNCS, pages 151–170, Paris, France, June 1987. Springer-Verlag.
- [25] Jonathan Eifrig, Scott Smith, and Valery Trifonov. Sound polymorphic type inference for objects. In *OOPSLA*, 1995. To appear. Currently available as `ftp://ftp.cs.jhu.edu/pub/scott/sptio.ps.Z`.
- [26] Jonathan Eifrig, Scott Smith, and Valery Trifonov. Type inference for recursively constrained types and its application to OOP. In *Mathematical Foundations of Programming Semantics, New Orleans*, volume 1 of *Electronic Notes in Theoretical Computer Science*. Elsevier, 1995. To Appear. Currently available as `ftp://ftp.cs.jhu.edu/pub/scott/ooinfer.ps.Z`.
- [27] Jonathan Eifrig, Scott Smith, Valery Trifonov, and Amy Zwarico. Application of OOP type theory: State, decidability, integration. In *OOPSLA*, 1994.
- [28] Giorgio Ghelli. A static type system for message passing. In *Proc. OOPSLA*, pages 129–145, 1991.
- [29] Adele Goldberg and David Robson. *Smalltalk-80: The Language and Its Implementation*. Addison-Wesley, Reading, MA, 1983.

- [30] John V. Guttag, James J. Horning, S.J. Garland, K.D. Jones, A. Modet, and J.M. Wing. *Larch: Languages and Tools for Formal Specification*. Springer-Verlag, New York, NY, 1993.
- [31] C. A. R. Hoare. Proof of correctness of data representations. *Acta Informatica*, 1(4):271–281, 1972.
- [32] Martin Hofmann and Benjamin Pierce. A unifying type-theoretic framework for objects. *Journal of Functional Programming*, 1995. To appear. Previous versions appeared in the Symposium on Theoretical Aspects of Computer Science, 1994, and, under the title “An Abstract View of Objects and Subtyping (Preliminary Report),” as University of Edinburgh, LFCS technical report ECS-LFCS-92-226, 1992.
- [33] Daniel H. H. Ingalls. A simple technique for handling multiple polymorphism. In Norman Meyrowitz, editor, *OOPSLA '86 Conference Proceedings, Portland, Oregon, September 1986*, volume 21(11) of *ACM SIGPLAN Notices*, pages 347–349. ACM, November 1986.
- [34] Dinesh Katiyar, David Luckham, and John Mitchell. A type system for prototyping languages. In *Conference Record of POPL '94: 21st ACM SIGPLAN–SIGACT Symposium of Principles of Programming Languages, Portland, Oregon*, pages 138–150. ACM, January 1994.
- [35] Gary T. Leavens. Modular specification and verification of object-oriented programs. *IEEE Software*, 8(4):72–80, July 1991.
- [36] Gary T. Leavens and William E. Weihl. Reasoning about object-oriented programs that use subtypes (extended abstract). In N. Meyrowitz, editor, *OOPSLA ECOOP '90 Proceedings*, volume 25(10) of *ACM SIGPLAN Notices*, pages 212–223. ACM, October 1990.
- [37] Gary T. Leavens and William E. Weihl. Specification and verification of object-oriented programs using supertype abstraction. *Acta Informatica*, 1994. To appear. An expanded version is Department of Computer Science, Iowa State University, Technical Report 92-28d, August 1994.
- [38] Barbara Liskov, Dorothy Curtis, Mark Day, Sanjay Ghemawat, Robert Gruber, Paul Johnson, and Andrew C. Myers. Theta reference manual. Technical Report Programming Methodology Group Memo 88, MIT, February 1995.
- [39] Barbara Liskov and Jeannette Wing. A behavioral notion of subtyping. *ACM Transactions on Programming Languages and Systems*, 16(6):1811–1841, November 1994.
- [40] Narciso Martí-Oliet and José Meseguer. Inclusions and subtypes. Technical Report SRI-CSL-90-16, Computer Science Laboratory, SRI International, December 1990.
- [41] Bertrand Meyer. *Eiffel: the language*. Prentice-Hall, 1992.
- [42] John Mitchell and Gordon Plotkin. Abstract types have existential type. *ACM Transactions on Programming Languages and Systems*, 10(3), July 1988.
- [43] John C. Mitchell. Toward a typed foundation for method specialization and inheritance. In *Proceedings of the 17th ACM Symposium on Principles of Programming Languages*, pages 109–124, January 1990. Also in Carl A. Gunter and John C. Mitchell, editors, *Theoretical Aspects of Object-Oriented Programming: Types, Semantics, and Language Design* (MIT Press, 1994).

- [44] W. B. Mugridge, J. G. Hosking, and J. Hamer. Multi-methods in a statically-typed programming language. In Pierre America, editor, *ECOOP '91 Conference Proceedings, Geneva, Switzerland*, volume 512 of *Lecture Notes in Computer Science*. Springer-Verlag, 1991.
- [45] Greg Nelson, editor. *Systems Programming with Modula-3*. Prentice Hall, 1991.
- [46] Benjamin Pierce and Martin Steffen. Higher-order subtyping. Submitted for publication. A preliminary version appeared in IFIP Working Conference on Programming Concepts, Methods and Calculi (PROCOMET), June 1994, and as University of Edinburgh technical report ECS-LFCS-94-280 and Universität Erlangen-Nürnberg Interner Bericht IMMD7-01/94, January 1994.
- [47] Benjamin C. Pierce and David N. Turner. Statically typed friendly functions via partially abstract types. Technical Report ECS-LFCS-93-256, University of Edinburgh, LFCS, April 1993. Also available as INRIA-Rocquencourt Rapport de Recherche No. 1899.
- [48] Benjamin C. Pierce and David N. Turner. Simple type-theoretic foundations for object-oriented programming. *Journal of Functional Programming*, 4(2):207–247, April 1994. A preliminary version appeared in *Principles of Programming Languages*, 1993, and as University of Edinburgh technical report ECS-LFCS-92-225, under the title “Object-Oriented Programming Without Recursive Types”.
- [49] John Reynolds. Three approaches to type structure. In *Mathematical Foundations of Software Development*. Springer-Verlag, 1985. Lecture Notes in Computer Science No. 185.
- [50] Craig Schaffert, Topher Cooper, Bruce Bullis, Mike Kilian, and Carrie Wilpolt. An introduction to Trellis/Owl. In Norman Meyrowitz, editor, *OOPSLA '86 Conference Proceedings, Portland, Oregon, September 1986*, volume 21(11) of *ACM SIGPLAN Notices*, pages 9–16. ACM, November 1986.
- [51] Bjarne Stroustrup. *The C++ Programming Language*. Addison-Wesley, Reading, Mass, 1986.
- [52] Larry Tesler. Object Pascal report. Technical Report 1, Apple Computer, 1985.