
Right and wrong: ten choices in language design

Bertrand Meyer

Constructor University and Eiffel Software

Bertrand.Meyer@inf.ethz.ch

ABSTRACT

A description of language design choices that have a profound effect on software quality, criticism of how ordinary OO languages addressed them, and explanation of the thinking behind Eiffel's corresponding mechanisms.

A INTRODUCTION

Object-oriented programming rules the world. Or so it seems.

In reality, while languages presented as object-oriented (OO) have come to enjoy widespread usage, their usage frequently shows OO methodology being misunderstood and misapplied. Some of the blame lies with the design of the languages themselves.

Eiffel as a programming language is simply an attempt to help programmers apply the methodology of OO analysis, design and programming, including Design by Contract as expounded in [4] and [10]. Eiffel's design is a combination of object-oriented concepts with ideas from formal methods (particularly axiomatic semantics), "correctness by construction", modularity, and more generally modern software engineering. At the core of this approach are a set of clearly expressed principles, from Design by Contract to the Open-Closed and Single Choice principles and a good dozen others, which the supporting methodological publications explain in detail.

Other OO languages are not so firmly rooted in methodology. C++, in particular, had compatibility with C as one of its key design goals and presents itself as a "multi-paradigm language" (its creator also described it as a "Swiss knife") rather than strictly an OO language. Its successors retain a C-like style and a design guided by other criteria.

Half-embracing the object-oriented paradigm has negative consequences on software quality: OO design is a consistent discipline which needs to be applied consistently to yield advantages of extendibility, reusability, scalability and reliability. The programming language can support this goal, or not. This article identifies Eiffel's resolution of ten individual language issues, and compares some of these choices with those of other widely used OO languages, showing how their misunderstanding of methodological principles can lead programmers to produce bad designs and programs of subpar quality.

A.1 Background

The need for this article arose in part from the new difficulty of teaching programming.

It used to be that few people in society knew anything about programming. Among those who did have programming experience, even fewer knew anything about OO programming when it started attracting attention. At that time, as one now realizes, teaching the concepts was fairly easy: they are consistent, powerful and convincing. No preconception on the learners' side stood in the way of accepting them with an open mind.

Today, having some programming experience is no longer a rare skill. Increasingly, pre-university curricula include some programming. Many people have even encountered an OO language; they include most students starting a university computer science program. Language exposure is not by itself the problem; the trouble is that together with the language they were taught bogus principles, such as "*multiple inheritance is bad*" (often just "*inheritance is bad*"), or "*you should make fields private and write getter functions to provide access to them*", or "*you can return from anywhere in a function*", or "*writing interfaces helps OO design*" or "*you should use defensive programming*" and other atrocities apparently entered into young minds at the same time and with the same authority as the information that the earth revolves around the sun.

Teaching is hard enough, but the first task when dealing with such preconceptions is to un-teach, which is even harder. This article is an attempt to set basic choices straight.

A.2 Design guidelines

One of the goals of the following discussion is to highlight the mode of thinking that has guided the design of Eiffel.

There seems to be a general trend in modern language design of proposing clever constructs that offer impressive modes of expression for certain cases. While this quest for programmer convenience is legitimate, it should only be one of the guiding forces in language design. The criterion that plays the principal role in the design and evolution of Eiffel is the constant effort to help programmers produce software that meets fundamental software quality criteria [4] [10]: correctness, robustness, extendibility, reusability.

Not *preventing* programmers from writing *bad* programs (that is not possible — at best the programming language may be able to make bad programs a trifle harder to produce) but *helping* them write *good* programs if they are determined to do so.

This observation also serves to qualify the focus on quality: Eiffel is in no way intended as a straitjacket, a holier-than-thou injunction to programmers to go through endless hurdles in pursuit of correctness. To the contrary, the language should be supportive of the programmer, devoid of verbosity, and pleasant to write, which is where the expressiveness criterion comes back. Sections 1.5 and 7 will explore this aspect, showing how Eiffel reconciles advanced OO concepts with concise traditional notations.

A.3 Discussion style

The presentation follows the tradition of articles of a generation ago, discussing concepts of programming methodology and programming languages.

It does not resemble articles that find their way to publication today, and would be desk-rejected by all the significant venues because of its absence of an empirical evaluation. While quantitative assessments are of great value, there remains room for discussions of a purely conceptual nature, based on logical rather than numerical arguments.

A.4 Language criticism

The strong terms in which the discussion assesses some design decisions of dominant languages do not imply disrespect for either the designs or the designers. The languages' success is evidence enough that they filled a need. But it does not forestall criticism.

Discussions of such matters are healthy. They are also important: at issue are not individual opinions and susceptibility, but the quality of software driving today's societal processes. A harmful language design decision leads to buggy programs and to malfunctions that ultimately affect people in ways superficial or existential.

A.5 Scope

After introducing a few mathematical conventions in section B, the discussion covers the following ten design choices:

- Uniform access, or how to use object properties abstractly (section 1).
- Control structures of structured programming (section 2).
- Overloading and constructors (section 3).
- Inheritance, single and multiple, and the noxious notion of interface (section 4).
- Design by Contract mechanisms (section 5).
- How to handle static and shared information (section 6).
- How to provide classic syntax for novel concepts (section 7).
- The proper role of exception handling (section 8).
- Concurrent programming in an OO context: SCOOP (section 9).
- Removing null pointer dereferencing (section 10).

In addition, Appendix A has some observations about language definition and Appendix B about language definition.

Another general comment is in order. The essence of Eiffel, one might almost say its soul, is its type system (complemented by Design by Contract concepts). The very first publication about Eiffel [2] described the subtle combination of genericity and (multiple) inheritance that forms the backbone of that system. All major subsequent evolutions, including some of the original language traits described in this article, from export controls to creation, conversions, concurrency (SCOOP) and void safety, are extensions and tightenings of the type system, devised to achieve a proper balance of expressive power and safety. Programming in Eiffel largely consists of putting that type system to work.

A.6 Precise object-oriented terminology

A *class* (static mechanism) defines a set of potential run-time *objects* (dynamic concept), its *instances*. A class has *features* (or “members”) specifying operations applicable to its instances. A feature may be a *command*, changing objects, or a *query*, providing information about objects. A command is also known as a *procedure*. A query may be implemented as a *function*, computing the required information, or an *attribute*, which defines a *field* present in every object and containing that information. (In the Java/C# context the term “field” also covers attributes, but it is preferable to distinguish the static and dynamic views: a class has attributes, defining fields in its instances.) Procedures and functions are both defined by algorithms and are known as *routines* (or “methods”).

A program element can be correct at three levels (corresponding to syntax, static semantics and semantics), each only defined if the preceding one is: it is *syntactically correct* if it is built according to the structure of the language (for example, matching parentheses); *valid* if it satisfies additional static rules not expressible in syntax formalisms (for example, type rules); and *correct* if its execution produces the expected results.

B MATHEMATICAL NOTATIONS

While largely informal, the discussion will at times reflect the underlying mathematical concepts, particularly for the discussion of inheritance in section 4.4 (see [21] for the outline of a general mathematical basis for programming), using the following conventions.

B.1 Functions

A function is a set of pairs, for example $\{[Elizabeth, Philip], [Charles, Diana], [Ann, Mark]\}$ (since it is a set, not a sequence, the order in which we list these pairs is immaterial), such that for any given element at most one pair has it as its first element. Removing this restriction, for example by including a pair $[Charles, Camilla]$, still yields a *relation*, but it is no longer a function.

The sets from which a function’s pairs take their first and second values (respectively) are its *source set* and *target set*. The actual first and second elements of pairs in the function make up its *domain* and *range* (respective subsets of the source and target sets); if we call the preceding function *spouse*, its domain is $\{Elizabeth, Charles, Ann\}$ and its range $\{Philip, Diana, Mark\}$. These concepts also apply to relations that are not functions.

The defining property of functions makes *functional notation* possible: if x is an element of the domain of *spouse*, we may use *spouse* (x) to denote the single element y such that there is a pair $[x, y]$. For example *spouse* (*Elizabeth*) is *Philip*.

We are particularly interested in functions with a *finite* domain, and hence finite range; the set of such functions with source and target sets A and B will be written $A \mapsto B$. “Mapping” is a synonym for “function” but will be restricted in this presentation to finite functions only.

If f and g are functions, their composition $f ; g$ is the function h such that $h(x) = g(f(x))$ wherever defined. (It is also commonly written $g \circ f$, but the “ \circ ” notation is convenient since it lists operands in order of their application.)

B.2 Overriding union

An important mathematical operator for the modeling of programming concepts is “overriding union” Ψ , which yields a function from two functions.

The union of two functions is a relation but not necessarily a function; for example if *spouse_new* is $\{[Edward, Sophie], [Charles, Camilla]\}$, then *spouse* \cup *spouse_new* has two pairs starting with *Charles*, making it violate the defining property of functions.

To guarantee that the result is a function, we may instead of union use the Ψ operator. It acts like \cup , but in case of a pair clash it only retains the pair from the second operand, so that it always yields a function. (The notation reflects this property by including a “ Ψ ” which “leans” towards the second operand.) Here *spouse* Ψ *spouse_new* is the function $\{[Elizabeth, Philip], [Charles, Camilla], [Ann, Mark], [Edward, Sophie]\}$. Overriding union is useful to model situations, such as inheritance, in which a new mapping both complements and overrides an existing one.

B.3 Finite permutations

We will need the notion of finite permutation on a set A . It is the usual notion of permutation (a total function σ from A to A that is one-to-one, meaning that both σ and its inverse σ^{-1} are injective functions of domain A), but with the condition that it is the identity except on a finite subset of σ . An example is $\{[a, b], [b, a], [c, d]\} \cup \{[x, x] / x \notin \{a, b, c\}\}$; in other words, a function that maps a to b , b to a , c to d , and any value other than a, b, c (including d) to itself. Another way of stating that convention is that we only consider permutations that substitute a finite subset of elements, extending them with the identity on the rest of A to so as to get total functions.

1 UNIFORM ACCESS

We start with a matter on which the original OO language, Simula 67, had things right and the design of C++, faithfully followed by today’s dominant languages, introduced a conceptual mistake, breaking the consistency and simplicity of OO principles. Even UML, which as a modeling language should rely on abstract concepts rather than imple-

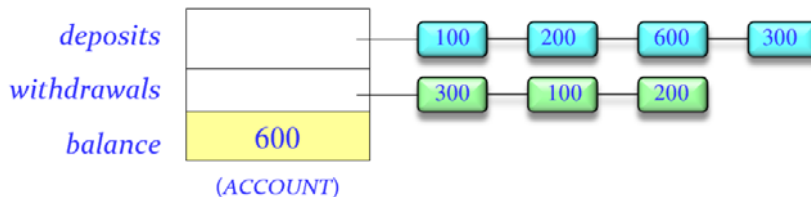
mentation nuances, perpetuates the flaw. In some cases, such as C#, the designers realized the awkwardness of the result but instead of adopting the obvious simple solution they devised a contorted workaround.

1.1 Computing, or looking up

Uniform access, the elegant principle violated by these designs, is the idea [4] [10] that “*it does not matter whether you look up or compute*”. Client modules needing to obtain some information should be able to access it in the same way regardless of the implementation choice: *storage* (the element is recorded in the corresponding data structure, obtaining it is a simple lookup); or *computation* (the element can be computed from other elements through an algorithm, obtaining it requires executing the corresponding operation).

Uniform access is a natural consequence of the general principle of data abstraction, which lies at the basis of object technology. The principle (also known as encapsulation, and closely related to another fundamental guideline, information hiding) states that clients of a module should be able to access its available elements of functionality through the specification of their abstract properties, without having to use properties of their implementation. The abstract properties take the form of features (operations), known to client programmers through their specifications (not the way in which they are implemented). Features are of two kinds: commands, which can change the state of objects; and queries, which return information about objects. Uniform access governs queries.

A typical query, for a bank account class, is the operation *balance* returning the balance of an account. Consider an implementation in which a bank account object looks like this (“Representation A”):



For every bank account object, the representation includes (in addition to other properties not shown such as account number and owner), a field denoting the current balance and two lists of banking operations, one for deposits and the other for withdrawals. Each element of these lists includes the amount involved (plus any other information not shown, such as the operation’s date). Representation A has the following properties:

- The fields are not independent. They are constrained by a property which we must include in the class invariant: $balance = deposits.total - withdrawals.total$ (with a function *total* yielding the accumulated value of such lists of banking operations).
- Commands *deposit* and *withdraw* must maintain this property: they have to update both the *balance* field and the appropriate list. For example, *withdraw* (100) must not

only decrease the *balance* field by 100 but also add an element with value 100 to the list *withdrawals*.

- To obtain the balance, it suffices to look up the *balance* field, which is always up to date thanks to the invariant property.

In another approach, which we may call “Representation B”, the representation of an account does not include the *balance* field. Representation B can still provide all the functionality of the original but internally does things in a different way:

- Commands *deposit* and *withdraw* only need to add an element to the respective list. There is no field to update.
- Finding out the balance of an account, on the other hand, now requires a computation, returning the value of *deposits.total* – *withdrawals.total*. Computing *total* involves (in a simple implementation) traversing the respective list to sum its element values.

Uniform access means protecting client modules from such internal implementation decisions as the choice between Representation A and Representation B. Languages such as C++, Java and C# force a different syntax to access the balance of an account *a*:

- *a.balance* for a field representation (A).
- *a.balance ()* for an representation by a function (B).

While the difference is small, it forces knowledge of suppliers’ implementation onto client programmers and negates the benefits of information hiding. Changes of representation between variants such as A and B occur routinely during the development of a system; preserving clients against having to update each time is part of the general discipline of data abstraction and key to the stability of the development.

It gets worse.

1.2 What do we do with fields?

Even the most uneducated users of languages such as Java, C++ and C# have heard about information hiding and realize that giving clients unfettered access to fields of objects (attributes of the corresponding classes) is not good for their karma. Hence the advice universally taught in the manuals for such languages (which the recipients mistake for good OO methodology, rather than for what it is — punishing programmers with the consequences of flawed language design decisions): never export an attribute. Instead, write a “getter function” of the form

```
balance_getter: INTEGER                                     -- [P1]
    -- Balance of this account.
do
    Result := balance
end
```

and export only *balance_getter*, keeping *balance* secret. (**Result** denotes the result to be returned by the function; see below.)

Why this self-imposed rule? The reason is a property that is more crazy than anything seen so far (although more follows): if you export an attribute in C++ and its successors, you export it with *full rights for the client*. So if *balance* were exported, any client code could use, for an account *a*:

```
a.balance := 5000 -- [P2]
```

which causes a gross violation of information hiding since it allows clients to play at will with the internals of a supplier object (in a way similar to what happens if instead of using an electronic device through its buttons you open it up and play with the wires, which typically voids the warranty).

The example is typical of why such unfettered export is wrong: the designer of class *ACCOUNT* has clearly decided, regardless of the internal choice of implementation, that clients should not directly modify *balance* but go through *deposit* and *withdraw* which may perform various checks and updates. If nothing else, both of these routines should only accept non-negative arguments (enforced in Eiffel through a precondition *require amount >= 0*). With remote assignments to *a.balance* allowed, all methodological protections remain.

No wonder then that exporting attributes is widely presented as a cardinal sin. This systematic advice is, in itself, a bad sign: while it is normal to associate methodological guidelines with language constructs, a language design that includes a mechanism together with an express injunction not to use it can only be described as sado-masochistic.

1.3 Exporting queries as queries

The proper approach is of course to allow the exporting of a query, whether an attribute or a function — but not *as an attribute* or *as a function*, which is none of the clients' business: as a query. In the interface specification of an Eiffel class (produced automatically from the class text, as will be seen in 4.8 below), you will see

```
balance: INTEGER  
-- Current balance of this account.
```

whether the class is implemented as Representation A or Representation B (or another). This specification is generated automatically from the class text, through environment tools that produce the “contract view” (also known as “short form”) of a class. It is the same for an attribute or a function and does not betray which of these cases hold, but in both cases will include abstract properties from preconditions, postconditions and invariants, for example the invariant clause

```
balance ≥ authorized_overdraft -- [P3]
```


In either case, attribute or function, the notation $a.balance$ is permitted in a client — if *ACCOUNT* exports *balance* to it — and denotes an expression. One cannot assign to an expression ($a + b := 1$ makes no sense in usual programming languages), so the remote assignment $a.balance := 5000$ is trivially incorrect at the syntactic level, not even requiring a validity rule. (See below for a legitimate re-interpretation of this notation.)

1.4 Getters and setters

Once one has understood that it is OK to export a query — attribute or function — when exporting it *only* as a query, the notion of a getter function becomes useless. Yet decently written OO code in C++ and its successors is peppered with getter functions such as *balance_getter* above. Such functions are useless, constituting pure code bloat.

It is remarkable that recent programming languages have tended to terseness in their syntax, favoring braces and other special symbols, in the C tradition, over keywords — but then lead to considerable amounts of such boilerplate code, hence doubly decreasing readability. Code bloat is always detrimental, leading not only to waste of time in the initial coding effort but to an increased maintenance burden (“technical debt”). Concision is not served by cryptic syntax but by getting rid of useless code.

1.5 Assigner commands

Getter functions are always useless and harmful. Setters are another matter: for an exported attribute, the class author needs to control what kind of change to permit clients to perform. In the balance case for a bank account, it is unlikely that the class author would permit clients to set the corresponding field of an *ACCOUNT* object to an arbitrary value; the appropriate policy is probably, as assumed above, to define two setter commands *deposit* and *withdraw*, which do modify balance but each in its own controlled way — add or subtract a non-negative value, within certain bounds. Note the importance of equipping such setter with contracts, for example the precondition in

```
-- [P4]
```

```
withdraw (amount: INTEGER)
```

```
require
```

```
     $amount \leq 0$ 
```

```
     $amount \geq balance - authorized\_overdraft$ 
```

```
    ...
```

In this case the only way to modify the balance of an account is to call *withdraw* or *deposit*, which perform modification of a very specific type.

In other cases the author of a class may want to allow clients to set the value directly; for example a class *HEATER* may have an attribute *temperature* and a setter command

```
set_temperature (t: REAL)                                -- [P5]
    -- Change the temperature to t.
require
    t ≥ 15
    t ≤ 25
do
    temperature := t
ensure
    temperature = t
end
```

Even when conferring field-setting privileges to their clients, classes often restrict the values as in this example. Contracts (particularly preconditions) are essential.

In such cases a notational facility is available, preserving an intuitive notation as a shortcut for a methodologically sound operation. If you add an “assign” specification to the declaration of *temperature*, making it

```
temperature: REAL assign set_temperature
```

you are associating *set_temperature* as an “assigner command” for *temperature*, with the effect that for the setting call

```
my_heater.set_temperature (21.5)                          -- [P6]
```

another synonymous syntax is permitted:

```
my_heater.temperature := 21.5                            -- [P7]
```

[P7] looks like an assignment (a remote assignment to a field of an object) but is not an assignment. Known as an “assigner call”, it is actually a procedure call, like [P6], but using a different syntax for exactly the same semantics. In particular, it is bound by the same type rules (for the type of the argument in [P6] and the right-side expression in [P7]) and, just as importantly, the contract: the precondition defined in the declaration of *set_temperature* continues to apply. (C# has a facility in the same spirit, “properties”, but it does not remove the need for getter functions, and of course does not support any notion of contract.)

The use of classical syntax to represent constructs that fit in OO methodology is part of a systematic language design pattern, discussed further in section 7.

2 DIJKSTRA WAS RIGHT

Before continuing with object-oriented abstraction mechanisms we briefly stop with a plague that should have been eradicated by our grandparents but still affects us. It is a

common myth that “*Python (like almost every programming language today) supports structured programming which controls flow using if/then/else, loop and subroutines*” (from a comment in a Stack Overflow discussion at bit.ly/3ErsDPZ). “Myth” is a polite way to say “lie”: all major imperative programming languages today other than Eiffel have an essentially full-fledged goto, and programmers use it generously. It is not expressed as a “goto” but takes the form of:

- Instructions to break out of a loop: break and continue.
- Worse yet, the “return” instruction in routines.

(Plus exceptions, discussed in section 8.) It is hypocritical and downright misleading to pretend that these languages forsake the goto and apply “structured programming” since these three constructs correspond to all practical uses of the “goto” (other than totally crazy ones). They provide a cynical violation of the principles enunciated in Edsger Dijkstra’s famous 1968 goto diatribe and the structured programming principle of organizing the control flow into one-entry, one-exit blocks (OEOE).

OEOE is particularly important in the case of routines. An instruction (in Python, Java, C++, C#...) of the form

return *exp*

may appear anywhere in the body of a routine, breaking out of all remaining computations to return (for a function) the value of *exp*.

Eiffel applies a strict OEOE policy; there is no goto of any kind or disguise. The matter of defining return values of functions is handled in a simple fashion: functions may use the predefined local variable **Result**, initialized to the default value of the function’s return type (0 for integers, false for booleans and so on). Whatever value Result has on the (single) point of routine exit is the value returned by the function. This convention avoids the need to check that a function’s result is always defined (as required with goto-style “return” instructions).

Dijkstra’s arguments were informal, largely an appeal to gut feeling (“*for a number of years I have been familiar with the observation that the quality of programmers is a decreasing function of the density of go to statements they produce*” — one can hear the referees screaming “*where is the empirical study?*”). Even the more respectable argument that structured-programming proponents made at the time — that **goto** complicated formal axiomatic reasoning — went away when someone pointed out that it suffices to perform, at each node of the control flow, a logical **or** of the verification conditions at all predecessor nodes. The present discussion also makes no representation of a scientific demonstration but simply on the author’s (and the Eiffel community’s) many decades of programming without goto instructions. Not as a discipline (which anyone can violate secretly) but as a strict impossibility since no goto, actual or ersatz, is physically possible.

The lesson, even if not scientifically established, is clear: there are only benefits in sticking to OEOE. Program flow becomes simpler and clearer. When you have something more complicated, and your Java programmer friends would tell you “oh, just break out of the loop”, you know better. You either:

- Introduce a boolean local variable to express the exit/continuation condition, which makes the program reader clearly understand what is going on logically.
- Realize that the structure is getting too complicated and modularize.

These observations apply both to loops and to routines. In the case of routines, the guarantee that the code is executed to the end, and that there is no need to look for possible short-circuits or early exits, considerably simplifies building, reasoning and debugging.

3 OVERLOADING AND CONSTRUCTORS

Name overloading, which we may more precisely call *syntactic* overloading because it is a notational convenience not adding any expressive power to the language, is the presence of a language-defined disambiguating mechanism enabling the programmer to give the same name to two routines appearing in the same syntactical context, with the proviso that the argument signature (the number and types of arguments) is different, so that in any call the actual arguments will define which version is intended.

3.1 Syntactic overloading: an example

If we have two routines

open_by_name (*employee_name*: *STRING*)

open_by_number (*employee_id*: *INTEGER*)

(both meant to open the record associated with an employee), a language with overloading allows us to give them a single name such as *open*. Then a call *open* ("*Jane*") will use the first version and *open* (*2654*) will use the second one.

The Ada language popularized name overloading, which brought a minor advantage of flexibility in a non-OO context (Ada did not introduce OO mechanisms until later versions). Combined with object-oriented programming, syntactic overloading is not only unneeded but harmful, as it breaks the essential simplicity of the OO scheme.

3.2 Semantic overloading through polymorphism and dynamic binding

One of the key OO mechanisms is the combination of polymorphism and dynamic binding, which provides a form of *semantic* overloading, far more interesting than the syntactic version and rendering it largely useless. A variable *x* can be polymorphic, meaning that it can at run time denote objects of different types defined by their respective classes (all, however, descendants of a common ancestor in the inheritance hierarchy). Then if some of these classes have different versions of a routine with the same name, the version triggered by any call *x.f* is the one corresponding to the type of the object that *x*

denotes at that particular time during execution. The standard example, excellent in spite of being standard, is a graphical operation

fig.display

which applies the appropriate version of `display`, different for each kind of figure, where `fig` is polymorphic — dynamically representing a rectangle, or a circle etc., with a specific `display` algorithm in each case. (Polymorphism and dynamic binding are distinct mechanisms, although they go well with each other; many people confuse them. The functional programming community adds to the confusion by the gratuitous introduction of a pompous term, “parametric polymorphism”, for the simple concept of a type defined with generic parameters, which has nothing to do with polymorphism.)

Unlike syntactic overloading, the combination of polymorphism and dynamic binding truly provides a new level of expressive power, since a run-time mechanism is involved: the selection of the right version of the operation. From a software architecture perspective, it is a natural and necessary extension of information hiding: a client programmer can call an operation such as `display` not only without knowing about its implementation, but also without knowing *which version* of the operation, among several competing ones, will be triggered in each case.

3.3 Syntactic overloading is useless and harmful in an OO context

In contrast with the OO combination of polymorphism and dynamic binding, syntactical overloading is a mere notational convenience. The difference between a purely syntactic mechanism and a semantic one (increasing expressive power) is clear if we consider the consequences of removing the mechanism:

- If name overloading is not available, the author of a class simply needs to devise a different name — such as the original names above, *open_by_name* and *open_by_number* — for each variant of an operation; then client programmers must use the appropriate name in each call. The burden is light; in fact this solution enforces clarity and helps avoid errors by forcing programmers to specify which version they want. (In the overloaded version, *open (2654)* mistakenly used instead of *open ("2654")* would be accepted but cause wrong behavior, referring to a file identifier rather than a file name.)
- With dynamic binding removed from OO languages, each use of a routine with multiple versions needs to discriminate explicitly between the various cases: if we have a circle, apply the circle algorithm; else if we have a rectangle, apply the rectangle algorithm; and so on. Every client module using these facilities must know the full list of cases and undergo an update in the case of changes of additions to this list. In contrast, the OO idiom *fig.display* with polymorphism and dynamic binding frees clients from having to know which kind of figure `fig` will denote in any particular run-time case, and rely instead on an automatic mechanism to select the right version of the operation in each case.

These observations suggest that syntactic overloading is useless in an OO context. It is, in fact, harmful. In the pursuit of the central issue of programming methodology, maintaining a correspondence between what the program says (the program text) and what the program does (its myriad possible executions), overloading destroys the fundamental simplicity of the notion of class. A class maps a set of operation names to a set of operations. A mathematical model will associate such a mapping with every class:

features: Name \mapsto *Feature*

Overloading invalidates this view: *features* is no longer a function. To obtain the correspondence between names and features, we need an intermediate mapping of names into some internal disambiguating names invisible to the programmer. This approach considerably complicates the conceptual framework, with — as seen in the next section — particularly damaging consequences for the inheritance mechanism.

Even as a notational simplification, overloading often fails in its objective because the disambiguation criterion — the argument signature — is not necessarily the right one, as can be seen in one of the simplest examples that come to mind. To initialize a complex number, 2D-vector or 2D-point, two operations are equally available:

set_cartesian (*x, y: REAL*) -- Initialize by cartesian coordinates

set_polar (*ro, theta: REAL*) -- Initialize by polar coordinates

Would it not be nice to use the same name? Unfortunately, the signatures are the same, so it does not work.

Although the need to use different names is — as noted above — at worst a minor nuisance (and at best actually a benefit), this example causes more than notational trouble because of a crucial language mechanism (in C++-like languages) relies on it.

3.4 Overloaded constructors

While it would seem appropriate, when using C++, Java or C#, to stay away from overloading, these languages do not leave programmers the choice because they *require* overloading in the case of “constructors”.

The goal of a constructor in a class is to initialize its instances. More precisely (but unknown to most OO programmers, who unless they learn the approach through Eiffel are not taught any formal concepts) is to make sure that newly created objects satisfy the class invariant. To see how to achieve this goal in a reasonable way, it is useful to look at how Eiffel handles creation:

- A basic creation instruction, for *x* of type *C*, reads create *x*. It initializes all fields of the new object to default values (0 for integers and so on). This form assumes that *C*'s invariant is trivial since default fields will satisfy it.
- To override the default initialization, the instruction takes the form create *x.make* (...) where *make* is a *creation procedure* of the class, which performs the needed initializations. “Creation procedure” means the same as “constructor”, but the term

emphasizes that such an operation is just a procedure, used here for the particular purpose of initialization.

- For the preceding creation instruction to be valid, `make` must appear in a `create` clause listing the allowed creation procedures. In the point example the clause would list `set_cartesian` and `set_polar`, allowing such creations as `create x.set_cartesian (1, 0)`.
- One of the creation procedures can be `default_create`, a feature of all classes (inherited from the top-level class of the inheritance structure, *ANY*), allowing the simplified form `create x`. (A class can redefine `default_create` to override the default initializations.) The absence of a `create` clause is considered a shorthand for a clause of the form `create x.default_create`.
- A creation procedure is simply a procedure; in addition to being used for creation, as in `create x.set_cartesian (1, 0)`, it can also, if it has the appropriate export status, be available for calls, as in `x.set_cartesian (1, 0)`, which resets an existing object.
- That possibility is open but not always desired. The class author has full flexibility: exporting a procedure such as `set_cartesian` for creation only (it is not available for normal call), for call only, or for both. (Note that the Eiffel export mechanism is fine-grained: you can export features to all classes, or to specific classes and their descendants, or to none at all.)
- The proper way to choose between these possibilities is the underlying theoretical criterion, unfortunately absent from usual presentations of object-oriented programming: the class invariant. For a procedure (such as `set_cartesian`) to be usable as a creation procedure, its body *b* must (in a simplified analysis, ignoring arguments) satisfy the Hoare-triple property $\{DEF\} b \{INV \wedge POST\}$ where *INV* is the class invariant, *POST* the procedure's postcondition and *DEF* the assertion expressing that all fields have their default values. For it to be usable in a normal call, it must satisfy $\{INV \wedge PRE\} b \{INV \wedge POST\}$ where *PRE* is the precondition. (The more detailed rule appears as part of a full discussion of the semantics of class invariants in [22]). Depending on which of these properties the routine satisfies, it may be usable for creation, for call, or both.

While programming language details can vary, an approach such as this one follows directly from the principles of data abstraction and program correctness in object-oriented programming. In contrast, C++ and its successors have adopted a bizarre set of conventions for object creation. The most bizarre is the rule that instead of creation procedures, treated as features of the class, creation relies on a set of “constructors” which all have the same name — the name of the class, for example *POINT*. There is no creation instruction, but creation expressions of the form `POINT (arg1, arg2, ...)`, which return a newly created object. In the tradition of overloading, the constructors' signatures determine which one is meant in each case.

The absence of a creation instruction forces creation always to repeat the type, even in a typed language environment in which the type of the target (such as the type of *x* in

create `x` or create `x.make (...)` has been declared anyway; it leads to the common and strange Java stuttering idiom `Point p = new point (...)`. (Eiffel does have creation expressions too, useful to avoid declaring a variable, for example to create an object and pass it as an argument to a routine in `rou (create {C}.make (...)`). But to create an object and associate it with a variable there is no need to repeated the variable's declared type.)

The slightly masochistic practice of giving the same name to different constructors makes it impossible to have different constructors of the same structure, as in the frequent Eiffel case of permitting both

```
create point1.make_cartesian (0, 1)
create point2.make_polar (1, Pi/2)    -- Creating identical points
```

The effect is impossible to obtain in C++, Java or C# since all the constructors would have the same name POINT and the same signature (two arguments of REAL type or equivalent). The only workaround would be to have a single constructors that handles all cases thanks to a third argument, possibly boolean in this case, indicating which version is desired. Such a contorted solution goes against all principles of good programming.

Programming language sites confirm the impossibility of same-signature constructors, dismiss it in blame-the-victim approach. For example, the answer on a Stack Overflow page (bit.ly/3yIdkjY) states that “*there are very few use cases for such an idea even if there was a technically valid way to do it (which in Java there isn't)*”. It is a typical cop-out in response to comments about a language's inability to support a programming need: just pronounce that (in Agile development parlance) YAGNI, “You Ain't Gonna Need It”! As the example of cartesian and polar initializations for points illustrates, the need is normal and reasonable. What is unreasonable and unjustifiable is the C++/Java reliance on a single name to cause ambiguity, and an artificial criterion, signatures, to remove the ambiguity (and give up if impossible).

Just as unreasonable is the resulting damage to program clarity: to understand what a particular creation expression `new C (x, y, z)` means, the client programmer has to go to the text of the class and peruse the list of constructors to see which one has a signature matching the types of `x`, `y`, `z`. Also note that inheritance-based subtyping can make the disambiguation rules complex.

What is most puzzling here is the absence of any need, or justification, for using overloading. Even the most unexperienced parents know to give their children different names. In programming too, different operations appearing in the same class should have different names, so that calls (such as those in creations using `make_cartesian` and `make_polar` above) are clearly differentiated. (Of course, different classes can use the same names — in the same way that different families can all have a daughter called Jill —, opening the way to the truly useful mechanism of dynamic binding.)

By breaking the fundamental simplicity of a class defining an injective mapping in *Name* \mapsto *Feature*, overloading also jeopardizes, as we will now see, the fundamental OO mechanism of inheritance.

4 INHERITANCE AND INTERFACES

Inheritance is another area where a simple approach, directly following data abstraction principles and seeking simplicity (including simplicity of the mathematical model) makes it possible to get rid of accumulated sediments of language obfuscation.

Misconceptions often arise because of an unduly restrictive view of inheritance. An earlier article [9] explored the wealth of applications of this technique.

4.1 Redefinition

In the OO mechanism of inheritance, a class *D* may be defined as inheriting from a class *C*, to express that *D* has all the features and invariant clauses of *C*, in addition to any that it may define for itself.

Fundamentally for the expressive power of OO design and its approach to reusability, it is also possible for *D* to redefine, or “override”, a feature inherited from *C*, by giving it a new implementation. For example a class *RECTANGLE* retains the function *perimeter* from its parent *POLYGON* but can provide a better implementation, specific to polygons that are rectangles.

If the single-mapping property holds (and hence no overloading is permitted), the situation is simple and the mechanism easy to explain: the presence of a new feature with the same name as an inherited one means that it overrides it, in the sense of the overriding union operator. It must have a conforming signature. To avoid any uncertainty or confusion, Eiffel requires the inheritance clause (the place at the start of the class where *D* states it inherits from *C*) to specify **redefine** *f, g, ...* for redefined features. (Otherwise the class is invalid since it includes different features with the same name.)

4.2 Renaming

Another useful facility is renaming, which enables a descendant to inherit a feature under a different name.

The mathematical operation is also very simple: if *m* is the names-to-features mapping (an element of *Name* \mapsto *Feature*) defined by class *C*, and *D* performs a set of renamings, mathematically a finite permutation σ (as defined in A.6), then the name-feature mapping associated with *D* is simply $\sigma^{-1} ; m$.

By the definition will only be valid if it does not introduce overloading, meaning that σ^{-1} is injective.

Comparing renaming and redefinition is illuminating, and combining them is often useful. Renaming changes the name, without changing the feature. Redefining changes

the feature, without changing its name. You may need one of these operation, or the other, or both $((\sigma^{-1} ; m) \cup m')$.

It is a source of constant wonder that programmers and students do not understand the fundamental beauty and simplicity of these ideas, and imagine all kinds of complicated scenarios, by failing to understand the basic distinction between a feature and its name. They are of course not to blame, but have fallen victim to approaches that have turned simple matters into complicated ones.

4.3 Multiple inheritance

Inheritance supports a refinement process whereby we build new abstractions on the basis of existing ones. It is an essential mechanism of object-oriented programming, applying to program construction the basic scientific concept of taxonomy.

Often, we need to use not only one existing abstraction but several. In a hierarchical windowing system, a window is both a graphical object and a tree. It has all the prerogatives (in the sense of features and semantic properties including invariants) of both of these notions. In such examples, which constantly arrive in the practice of object-oriented programming, multiple inheritance is the simple and natural solution, offering an advanced form of reuse.

Unfortunately, a botched design for multiple inheritance in C++ gave this technique a bad name. The first culprit is once again overloading. In the simple overloading-free world of OO programming, two classes can only be combined through multiple inheritance if their feature names do not clash. But a name clash is easy to resolve: just use renaming. So if both B and C have a feature called *f*, we may define

```
class D inherit
  A rename f as A_f end -- Obviously it would be enough to
  B rename f as B_f end -- rename one of the two.
feature
  ... New features of D ...
end
```

With overloading, both inherited versions can coexist if their signatures are different. No wonder that with such complicated semantics multiple inheritance gets a bad rap.

This area is one where wrong preconceptions are so entrenched that discussions become impossible, so convinced are programmers (and even beginning students) that this simple matter has to be utterly complex. Those trained in C++ immediately go into discussions of data structures — “Vtables” and other tables of pointers, techniques that were in fact pioneered by Eiffel implementations but should be the concern of the small community of compiler writers, not general programmers (who in Eiffel are thankfully spared from having to know them). A little as if we had to explain conditional instructions by talking about Intel’s *jne* and *jnz* CPU instructions. As to Java/C# programmers,

they have already been instructed that multiple inheritance is bad (because C++ could not get it right) and “interfaces” solve the problem. They do not, by a long shot.

4.4 Interfaces

Early attempts to introduce principles of data abstraction and information hiding into programming languages led (even though they were posterior to the invention of OO programming in Simula 67, which they ignored) to modular mechanisms that required a module to be defined in two separate parts: the interface (or “specification”) part and the implementation (or “body”) part. The interface part only includes the list of features (taken here in the general, non-necessarily-OO sense of functionality elements) with their signatures; the implementation contains the actual algorithms and data representation choices. The typical “*modular languages*” applying this idea are Modula-2 and Ada. (It may be noted that the C language, while not providing an actual module construct above the concept of routine — called “function” — already promoted a similar discipline by directing programmers to define every program element *prog* through two files, *prog.h* for its functions’ interfaces and *prog.c* for their implementation.)

4.5 Limitations of interfaces in pre-OO modular languages

While historically a healthy step towards proper program structuring, the interface/implementation requirement of modular languages suffers from important limitations and defeats modern design methodology.

The most obvious deficiency is that these interfaces are not real specifications since they are purely structural, specifying only the names and signatures of the features and saying nothing of their semantics. An operative notion of interface would have to include Design by Contract mechanisms, specifying semantic properties — preconditions and postconditions of features, invariants — that would be binding on the implementations. Without these properties, the specification given by an interface in C, Ada, Modula-2 (as well as Java and C#) does not provide the client programmer with the information that would be necessary to apply the supposed principle of information hiding: being able to use functionality without knowing its implementation.

There is an even worse problem with the interface-implementation duopoly. A dogmatic split into just two categories misses the need for proper models to include a full spectrum between the most abstract descriptions and full implementations. Here inheritance leaves the interface/implementation concept in the dust.

4.6 Programs with holes

A brilliant invention of OO programming is the use of inheritance to provide a full spectrum of refinements, through the technique called “programs with holes” in [4] and [10]. It reflects a basic difference between definitions in mathematics and programming.

Consider the notion of a total order relation:

- In mathematics, we can define it by specifying that a certain operation, “<”, satisfies certain properties (transitive, irreflexive, total). We can then deduce other properties, for example that it is also asymmetric, and that it has relations “≤”, “>”, “≥” with their own properties (for example is transitive, symmetric, reflexive and total). We could in fact define the total order concept from any of them, but for the definition we only use properties that are strictly necessary; all other properties are theorems, not part of the definition (but consequences of it).
- To define the corresponding structure in programming, the perspective is different. The four operations are equally important as part of the interface; there is no reason to divide them into definitions and consequences.

In a corresponding class *COMPARABLE*, we may still choose one as unimplemented, or “deferred”, and the others “effective”, that is to say defined from it. For example:

```
lesser alias "<" (other: COMPARABLE): BOOLEAN
```

```
-- Is current object strictly less than other?
```

```
deferred end
```

```
greater alias ">" (other: COMPARABLE): BOOLEAN
```

```
-- Is current object strictly greater than other?
```

```
do
```

```
  Result := other < Current
```

```
ensure
```

```
  definition: Result = (other < Current)
```

```
end
```

```
... And similarly for “≤” and “≥” ...
```

(The alias clauses make it possible to call the features in infix form, for example $a < b$ as a synonym for $a.\textit{lesser}(b)$.) Feature *lesser* is deferred, but all the others are effective, defined in terms of it (“do” followed by an implementation, instead of “deferred”).

Class *COMPARABLE* is deferred, as is the case with any class that has at least one deferred feature. But it is not fully deferred: some of its features, such as *greater* and the others listed above, are effective, defined in terms of it. *COMPARABLE* is a “class with holes”: the implementation of *lesser* is left open, a “hole” for descendants to fill, but the others are already filled. A descendant class will effect (implement) *lesser*, as in

```
class STRING inherit COMPARABLE feature
```

```
  lesser alias "<" (other: COMPARABLE): BOOLEAN
```

```
    -- Is current string strictly less than other?
```

```
    do ... Implementation of lexicographic order on strings ... end
```

```
  ...
```

```
end
```

Key here is the property that such classes do not need to, and should not, do anything about the other ordering features, “ \leq ”, “ $>$ ” and “ \geq ”. They were defined in an effective form in the general class, *COMPARABLE*, and remain unchanged, although of course their implementation depended on “ $<$ ” and follows its own implementation in *STRING*.

4.7 Combining abstractions

Such a *COMPARABLE* class exists in the Eiffel Kernel Library, which in a similar way includes a class *NUMERIC* describing number-like objects; more formally, members of a commutative ring, with operations “+”, “-”, “*” and “/” and elements zero and one. Like *COMPARABLE*, *NUMERIC* may define effective features based on the deferred ones, for example “squared” in terms of “*”.

To define numerical classes such as those defining integers and reals, it suffices to use multiple inheritance from *COMPARABLE* and *NUMERIC*. As in other cases of inheriting from these classes, the effective classes provide implementations of the features inherited in deferred form, such as “ $<$ ” and “+” — but those only, since the effective ones automatically follow suit.

In Java and C#, classes are limited to single inheritance. It is a widely held *belief* (the word seems the most appropriate) that this limitation does not matter because a class can inherit from multiple interfaces. But as examples such as the above show, this possibility is a lame imitation of OO techniques. If a class conceptually needs to inherit from two others, it will have to use a kludge replacing one of them by an interface. If *COMPARABLE*, for example, becomes an interface:

- Features that should normally be declared as effective (“ \leq ”, “ $>$ ” and “ \geq ”) must be declared with their signature only, as if they were deferred.
- They will have to be implemented in every descendant class, like “ $<$ ”. But unlike with “ $<$ ” those implementations are always the same (as given for “ $>$ ” above).
- Somehow the interface must then come with a kind of user’s manual: if you implement the interface, always implement these features in a specified way.
- There is no guarantee that programmers of descendant classes will follow that injunction. Errors can creep in. In particular, Java interfaces, like Java itself, cannot include implementations. So the requirement just mentioned that implementations of “ $>$ ” etc. “are always be the same” is only wishful thinking, or documentation at best; authors of descendant classes can miss the injunction and introduce implementations that are incompatible with “ \leq ”.
- Contracts could alleviate the problem: in Eiffel, you can give “ $>$ ” (for example) a postcondition stating **Result** = *other* < *Current* (I am greater than you only if you are less than me). The rules on contract inheritance make such constraints binding on descendants. (In Eiffel they open the possibility of actually changing the default implementation for a routine like “ $>$ ”, overriding the default defined above from “ $<$ ” in *COMPARABLE*, if desirable for performance or any other reasons, while preserv-

ing the correct semantics.) But most languages do not have contract mechanisms — an observation which brings us to our next topic.

4.8 Interfaces as automatically derived artifacts

While interfaces in the Ada or Java sense are not a useful *design* tool, the need remains for *documenting* existing classes (deferred or effective) in an abstract form, without implementation details. One of the flaws of interfaces, not yet mentioned, is that they cause code duplication (always bad in software): the abstract information — names of features, their type signatures, header comments, contracts if supported by the language — appears in the concrete module as well. A proper approach is to *extract* that information from classes, for purposes of documentation.

The EiffelStudio environment provides that functionality through the notion of *views*, in particular the **contract view**, which extracts the elements just cited (name, signature etc.) from a class text, leaving out all routine bodies and all non-exported (private) features. In accordance with the Uniform Access principle (section 1, particularly 1.3), functions and attributes appear in the same form in that view.

The result is an effective form of documentation, abstract but guaranteed to be in sync with the actual code since it is extracted from it. (One of the worst problems with documentation arises when the software evolves and its documentation is no longer accurate.) It can be produced in various formats, including HTML and PDF.

The contract view is one of several such forms produced by the environment. The **flat view** is a reconstructed form of the class text including all inherited features, and applying all renaming and redefinition clauses. The **interface view** is the contract view of the flat view; as the name suggests, it provides the true interface documentation of the class, as seen by client programmers (who do not need to distinguish between features introduced in the class itself and those inherited from ancestors). As an example, here is the extract of the interface view for the library class *LINKED_LIST*, providing the specification of a specific routine:

```
copy (other: like Current)
    -- Update current object using fields of object attached
    -- to `other`, so as to yield equal objects.
require -- from ANY
    other_not_void: other /= Void
    type_identity: same_type (other)
ensure -- from ANY
    is_equal: Current ~ other
```

Note the inherited assertions, with a comment (added by the tool) indicating the class where they originated.

5 CONTRACT MECHANISMS

Design by Contract is the language trait that most people who have heard of Eiffel can cite. It is also among the most misunderstood, for example by being associated with defensive programming, of which it is the exact opposite.

This discussion will not present the basic concepts, under the assumption that the reader has heard of preconditions, postconditions, class invariants, loop invariants and loop variants. It limits itself to dispelling a few common misconceptions.

5.1 A close integration with the language

The first mistake is the often heard comment that “*one can use Design by Contract in any languages*” (often intended, in the author’s personal experience of many years, as a compliment, as in “*I like your ideas, but I apply them in ...*”, to which the only reasonable reaction is to smile politely without mentioning the absurdity of the whole idea). The evidence massively disproves this view: while dozens — possibly hundreds — of DbC extensions have been proposed for just about every significant programming language, none has gained wide acceptance. The only exceptions are Ada Spark and JML, used successfully but in the narrow field of formally verified mission-critical applications (and largely, in the second case, for research).

Why? We may immediately dismiss a folk explanation: that the concepts are too hard. This old canard has not basis in fact. Contracts are boolean conditions. Anyone who can write a conditional instruction — that is to say, any programmer — understands this notion and can use it for specification just as well as for implementation. The long experience of many educational institutions in teaching the concepts, including 14 consecutive years of the introductory programming course at ETH Zurich by the author and the resulting textbook [15], are evidence enough.

The actual reason is simple: Design by Contract can only work if tightly integrated with the fabric of the language and the development environment (IDE). For example, the rules about inheritance specify that the precondition of a redeclared routine must be kept or weakened, and its postcondition kept or strengthened. (They are commonly and bizarrely attributed to the “Liskov Substitution Principle” although they date back to Eiffel publications in 1986-1988 [1][3][4].) The Eiffel mechanism enforces this methodological rule by prohibiting a redeclared version of a routine from using the normal **require** or **ensure**, permitting only no clause — in which case the routine retains the original pre- or post-condition — or clauses of the forms **require else new** and **ensure then new**, which respectively “or” and “and” the *new* part with the original, automatically yielding a weaker resp. stronger assertion. Without this kind of built-in language rule, the use of pre- and postconditions is largely pointless since you have to duplicate the original assertion or risk inconsistencies.

5.2 Class invariants

Contracts without built-in language support become even more unrealistic in the case of class invariants. Dating back to an article by Hoare in 1972, class invariants (see [22] for a recent extensive review of the concept) can and should play a major role in software design by expressing the consistency conditions that underlie any significant class. But a class invariant can only be expressed at the level of a class.

Without a dedicated language construct, the invariant has to be repeated (causing massive code duplication, which is widely known as one of the major obstacles to software quality) as an addition to both the precondition and the postcondition of every exported routine, and its every redeclaration in a descendant! We are bordering on craziness here, and of course no one, including people who claim to “*apply the Design by Contract methodology but use an ordinary language*”, uses these facilities, crucial to obtaining any benefits from DbC.

5.3 A plethora of run-time checks for free

Here is a concrete example of these advantages. In section 7.2 below we will see the full version of a class invariant [P14], from a real program by the author, of the form

$$\forall x: A \wedge \forall y: H(x) \wedge \forall z: L(y) \wedge \text{some_property}(z)$$

where A is an array of hash tables $H(x)$ (for every element x of the array) of lists $L(y)$ (for every element y of such a table). One of the applications of contracts is the ability, provided by the compiler and IDE, to evaluate assertions as part of the testing and debugging process. The class invariant clause expresses a subtle consistency condition affecting elements of a data structures that in practical application may total tens of thousands of elements or more. Such contract elements are in fact easy for the programmer to write, since they just express the assumptions that preside over the writing of the code. In normal development — even without the prospect of program proofs — every execution of an exported routine of the class will cause the conditions to be evaluated, resulting in an enormous amount of automatic consistency checks. Once runs of the code proceed without violating any of these conditions, the programmer’s confidence in the code, while not absolute (see the discussion of proofs below), is far higher than any traditional approach to testing can provide.

The effect of these techniques on rooting out errors and more generally on transforming the process of writing software and getting it to work correctly is hard to imagine for people who have not programmed in this environment

5.4 At the very core of inheritance concepts

Coming back to the methodological side, we note, in addition to the invariant accumulation rule : the methodology includes the rule that the invariant of any class includes (in the sense of logical “and”) the invariants of its parents under inheritance. Here too no

manual or tool-based addition to a non-contract-equipped language can realistically apply this rule.

The inheritance rules of DbC — weakening/strengthening or pre/post-conditions, accumulation of invariants — are not just interesting ideas but part of the very definition of inheritance. Without the “and”-ing of invariants, inheritance is just a way to reuse some facilities, little more than the reference to an “include file” in C. The invariant rule shows what inheritance truly is: the “is-a” (subtyping) relation, not just providing the descendants with the ancestor’s features but also constraining them by the same semantic rules as the ancestors. The weakening/strengthening rules are similarly critical to understanding the fundamental OO tools of polymorphism and dynamic binding. If you redefine r in a descendant class, meaning that $x.r(\dots)$ will trigger a different operation depending on the type of x (ancestor or descendant) in any particular execution, you have a very powerful mechanism, but also a truly dangerous one: what prevents a descendant from redefining r in a way that violates the original intent? The only way to answer that question is to define that original intent and enforce it in descendants; in other words, use a precondition and a postcondition, and rely on the weakening/strengthening rules, which mean that the new version can do better — be more tolerant of the input, or produce better output — but can never do worse.

This approach makes it possible to teach inheritance properly. One can only regret that the vast majority of OO programmers have never heard of these concepts and hence lack any methodological guidance to use the powerful but risky tools of object technology.

(Just as they do with inheritance, the concepts of Design by Contract provide the proper perspective to discuss another important and delicate concept: exception handling. Section 8 will cover this topic.)

5.5 Understanding the Design by Contract methodology

In Eiffel, the use of contracts is neither abstract advice nor a constraint, but a powerful tool for building software that will be correct by construction. People who heard about the ideas but did not study them often confuse them with “defensive programming”, their opposite. Defensive programming is the practice of protecting operations, over and beyond, with checks of their applicability. The more checks, the safer the software (supposedly).

Contracts are entirely different. The idea is, as a human contract, to *distribute* the responsibility for consistency conditions between the various parties in the myriad interactions that take place between software elements — clients and suppliers. Distribute, but never replicate. If negative and non-negative values are to be handled differently for an operation, then either:

- The operation is responsible: it processes all values, producing different results in the two cases. The postcondition specifies the results in these various cases.

- The clients are responsible: for example the routine accepts non-negative values only, and puts that requirement in its precondition.

Defensive programming would encourage checking on both sides. This approach is haphazard and downright dangerous. If everyone is responsible, no one is responsible; errors will creep in anyway. There is no specification (in the form of a contract), so it is easy for each of the parties to believe that the other is in charge of a particular condition.

(On this particular topic, ample empirical evidence exists. For many years we conducted the Distributed Software Laboratory course, also known as DOSE, involving a project conducted by cooperating student teams from universities in many countries, often far from each other. We closely monitored student activity and systematically collected project information. In one of the first iterations we “encouraged” students to use contracts. Many did not, and the result was a plethora of conditions left unchecked out of the assumption that the other team was responsible. The second year we *mandated* the systematic use of Eiffel contracts and these problems disappeared entirely. Several publications present the empirical results; see [18] and the URL given in that reference.)

Another problem with defensive programming is that it adds code, often redundant. Code bloat is not the solution, but only compounds the problem. It may introduce further bugs. In fact, redundant code is more likely to contain bugs since it typically does not get exercised during testing. Duplication bugs can be particularly subtle, as the conditions checked redundantly may not be identical, particularly after the software undergoes change. Contracts involve no duplication or redundancy: they *complement* the code, which expresses the “how”, by stating the intent, the “what” (sometimes even the “why”). That is the reason why run-time contract monitoring (5.3) is effective as a testing technique: it compares reality to intent, expressed at a different level of abstraction.

A good criterion for finding out if a team is pretending to use Design by Contract but in reality just adding layers of (potentially noxious) defensive code is to ask if routines check their preconditions (as in a square root routine requiring a non-negative argument x in its precondition, and also including **if $x < 0$ then ...** in its code). With Design by Contract this sloppy defensive practice is never acceptable. (A Microsoft Research team working on the future “Code Contracts” library once asked to meet the author to present its work proudly. A quick examination showed that preconditions — for example, that a queue is not empty when you want to get its first element — gave rise to checks in the code that the conditions held, and raising an exception otherwise. End of story.)

Any redundancy useful for testing and debugging will be provided not by duplicating code but through the automatic tools of run-time contract monitoring (5.3).

Professional, correctness-by-construction software development does not consist of blindly throwing extra “just in case” consistency checks all over the code (which one can compare to adding props of all kinds to a structurally deficient building). The proper approach is to identify the consistency conditions that the computation requires, assign

each of them (on the basis of architectural considerations) which of the parts should be responsible for ensuring it, and stick to the decision.

5.6 An integral part of the approach

What makes Design by Contract practical is its applicability at many levels, supported by all elements of the Eiffel methodology, language and IDE:

- In the methodology, contracts provide a strong guidance to the programmer, a form of correctness by construction as just discussed.
- The language closely integrates the concept, woven into the rest of its mechanisms, particularly inheritance as seen above.
- Contracts provide the basic documentation tool, in the form of “contract views” and their variants, produced by the environment.
- The environment also supports run-time assertion monitoring, a key tool for testing and debugging as discussed above. Modern testing frameworks of the “xUnit” style also rely on assertions as test oracles, but the assertions are written only for testing. Eiffel’s assertions are an integral part of the program, written with it, not a posteriori, express fundamental semantic properties (preconditions, postconditions, invariants) and not just testing oracles, and serve all the other goals discussed here. Note that an extensive effort around the AutoTest framework [16] has taken advantage of built-in contracts to extend the xUnit ideas and, in particular, generate tests automatically.

For verification, run-time contract monitoring is a powerful technique, but one would ideally expect static proofs of correctness (defined as the conformance of the implementation to the contracts). The AutoProof program proving environment for Eiffel [20], pursues its goal. While still a research tool, it has reached significant capabilities — which anyone can try out in a browser at autoproof.sit.org — and has been used (in Nadia Polikarpova’s ETH PhD thesis) to prove formally and mechanically the correctness of a full library of fundamental data structures and algorithms, EiffelBase 2.

6 STATIC CLASSES AND SHARED INFORMATION

Object-oriented programming is the quintessentially modular method: it organizes systems into moderately sized units (classes), each built around a well-defined abstraction. The decentralization of the resulting software architecture is key to obtaining the method’s benefits of extensibility (ease of adapting software to changing needs or circumstances) and reusability (avoidance of repetitive or duplicated programming work).

Not everything can be decentralized: even the most modular architecture retains the need for some shared information. The most obvious example is console output: all parts of a program may need to write to a common medium.

6.1 Class methods and static classes

Smalltalk addressed this need with “class methods”, a technique that assumes that classes can at run time be treated as objects; this approach leads to interesting conceptual developments (“Metaobject protocols”) but seems like overkill for solving a simple need. C++ and its successors introduced the concept of a static class, which cannot be instantiated.

6.2 Once routines

The principal issue with shared information is its initialization. A class — a normal one, not “static” — describes a set of potential run-time objects, and specifies mechanisms to initialize them: creation procedures or constructors (discussed in section 3 above). But with a static class there is no object.

Eiffel introduced an original mechanism for this purpose: once routines. If a routine is declared with the keyword “**once**” rather than “**do**”, its body is executed not for every call to the routine, but only for the first call. In the case of a procedure, this mechanism is useful for initializing a shared resource on-demand, without having to know in advance which caller — if any — will need it. For example various routines of a graphical library may need the guarantee that the display will have been set up. The library may have a procedure

```
setup_display
once
    ... Operations to set up the display ...
ensure
    is_initialized
end
```

Then every routine of the library that needs to find the display initialized will start by calling *setup_display*. Only the first to do so will actually cause the setup operation; every subsequent call will have no effect.

In the case of a function, the result will be computed the first time around, and returned again by every subsequent call. This facility works well for shared objects, as in

```
Console: IO_MEDIUM
once
    ... Preparatory operations ...
    create Result.make (...)
end
```

Note the use of a creation instruction of target **Result**, typical of this design pattern. The first call creates the console, all subsequent ones use the same object. It does not matter where that first call comes from. If no element of the software calls console, the object

will not be created; this property is important as a system may need potentially require many such shared objects, but it would be a waste of resources to create all of them in all cases. In addition, since they are created individually and only when first needed, they avoid delaying system startup through a whole set of object creations.

6.3 Instance-free features

In an object-oriented architecture where every element of information appears in a class, where should shared information, such as once functions, appear?

One solution is simply to use inheritance. A class gathering shared information can serve as parent for any other class that needs access to the corresponding features. It is common to frown on such uses of inheritance, but they are perfectly legitimate since the corresponding classes embody valid abstractions; the real reason for criticism is the absence of usable multiple inheritance in most OO languages (section 4). An alternative to inheritance in this case is simply to qualify the feature with the class name, as in

```
{IO}.Console.write("message")
```

where *IO* is the class containing the above declaration of console. *{IO}.Console* is known as a “non-object call” since it does not involve an object, only a class.

It is valid to use a feature (here *Console*) only if it is instance-free, meaning that its implementation only depends on the properties of the class as a whole and not of any specific instance. A class feature may not refer to any attributes, except for constant attributes (declared in the form *Pi: REAL = 3.1415926535*). It may of course use local variables, including *Result* in the case of a function such as *Console*.

These constraints are implementation properties, not visible to client programmers through the interface specification. The definition is that a feature is a class feature — hence usable in a non-object call — if and only if it has a postcondition (part of its **ensure** clause) of the form

```
instance_free: class
```

which appears as part of the routine’s specification (which, as the reader will remember from 1.3, is generated automatically as part of the “contract view” of the class).

6.4 Varieties of once routines

The meaning of “once” is subject to parameterization using a “once key”. By default once is a synonym for *once ("PROCESS")*, but in a multithreaded environment you can use *once ("THREAD")*.

An occasionally useful mechanism is “once per object”: a routine specified *once ("OBJECT")* is executed on its first call, if any, on any given instance of the class. As a typical example, an object representing a currency exchange value may have an associated history — the list of previous values over the past five years — recorded in the database, accessible through a query *history*. That query should not be an attribute: then all histories of all currencies would be loaded in memory at all times. But it should

also not be a normal do function, since it would be executed (loading all relevant history values) each time it is needed, even though the result is always the same. The proper implementation is as a once-per-object function, which achieves the proper tradeoff between attributes and functions (reinforcing once again the concept of query and the Uniform Access principle). If a particular object never needs its history, the history will not be loaded. If it requires it once or many times, it will be automatically loaded the first time around, then kept in the object structure.

7 OLD SYNTAX FOR NEW SEMANTICS

The example of assigner commands and assigner calls illustrates a guiding principle for the syntax of Eiffel: maintain the integrity of the methodology, including design principles (the syntax), while accommodating well-established notations (the syntax). OO methodology is ambitious: it covers

This approach stands in contrast of the solutions taken by other OO languages, which compromise OO principles to retain compatibility with other approaches. A typical example is the case of basic types, which instead of being part of the OO type system remain in C++ special types with special properties. Java offers both OO versions and (supposedly for performance) traditional versions, with a tricky process of going back-and-forth between them. The Eiffel approach is different. All types, basic or not, are defined as classes. For example integer classes are part of the basic library and can be displayed in the EiffelStudio environment, where they start like this:

```
frozen expanded class INTEGER_16 inherit
```

```

INTEGER_16_REF
  redefine
    is_less,
    plus,
    minus,
    product,
    quotient,
    power,
    integer_quotient,
    integer_remainder,
    opposite,
    identity,
    as_natural_8,
    as_natural_16,
    -- ...

```

(different sizes are provided but you can use just *INTEGER*), with features such as

```
feature -- Basic operations
```

```

plus alias "+" (other: INTEGER_16): INTEGER_16
  -- Sum with `other`
  external
    "built_in"
  end

minus alias "-" alias "-" (other: INTEGER_16): INTEGER_16
  -- Result of subtracting `other`
  external
    "built_in"
  end

```

and class invariants indicating traditional properties. These classes are only “special” in that you cannot modify them (see the “built-in” mentions) and compilers handles them

specially, allowing the same efficiency as in a language like C where they directly correspond to machine operations.

All operations conform to the OO model of computation: an addition is, conceptually, an operation of the form *a.plus (b)*, and can be written that way. Such syntax is not, of course, what most people (at least most people having learned arithmetic before learning object-oriented programming) want to write; but then the **alias "+"** specification in the declaration of plus makes it possible to write *a + b* as a synonym for *a.plus (b)*. Not a step out of OO programming and back into C-style programming, but a different syntax for the semantics of a function call (which the compiler is knowledgeable enough to compile into the same code as its C counterpart).

Many Eiffel mechanisms follow this spirit of providing convenient syntax, often reflecting traditional modes of expression, for mechanisms that fully respect the object-oriented paradigm. Here are a few examples.

7.1 Aliases of various kinds

The alias mechanism makes it possible to define prefix or infix equivalents for functions with (respectively) 0 or 1 argument. There may be more than one alias, accommodating in particular Unicode operators, as in

```
quotient alias "/" alias "÷" (other: INTEGER_16): REAL_64
  -- Division by `other'
  external
    "built_in"
  end
```

(To insert Unicode characters with an ordinary keyboard it suffices to use a pull-down menu in the IDE.)

Specific aliases include brackets, enabling standard array notation. For example the basic element access feature for arrays is declared as

```
item alias "[" (i: INTEGER): G assign put -- [P8]
```

making it possible for an array *a* and an integer *i* to use *a [i]* as a synonym for *a.item (i)*, which yields the element of index *i* in *a*. Thanks to the assign part, the notation *a [i] := x* is permitted, as a synonym for *a.put (x, i)*, which updates the entry. The instruction

```
a [i] := a [i] + 1
```

is easier to read than its classic-OO-style equivalent *a.put (a.item (i) + 1, i)* (let alone the version using *plus* instead of “+”).

The use of *a [i]* on the left side requires the assigner command mechanism (here the specification of put as an assigner command for item in [P8]).

In these and other facilities a general principle of the language’s design is that whenever a simplification mechanism is available to the language designer it is also offered (if potentially useful) to the programmer — the class designer. Allowing *a [i]* as a shortcut for is not a special facility for placating programmers used to traditional array syntax.

It is a full-fledged “bracket alias” mechanism that any class can use. For example a hash table class (such as *HASH_TABLE* in the EiffelBase library) can define a bracket alias for the operation that accesses an element through its key, and associate it with the appropriate assigner command, allowing the writing of an insertion as

```
age ["Jill Smith"] := 29
```

7.2 Iterators

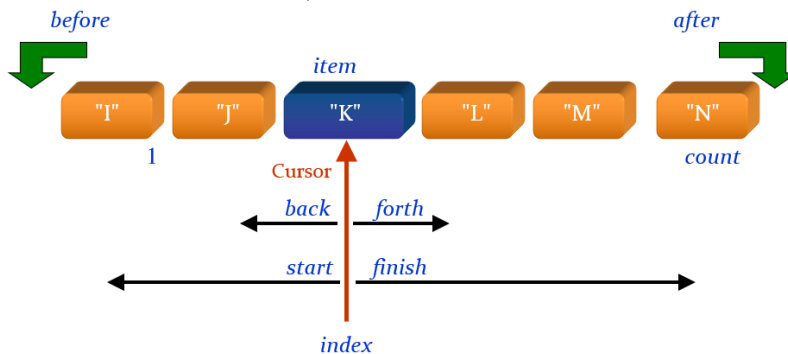
Many languages have special constructs for iterating over specific data structures. Eiffel conceptually has a single kind of loop

```
from Init invariant Inv until Cond loop Body variant var end           -- [P9]
```

(where the **invariant** and **variant** clause, and they only, are optional). It can serve to iterate over a data structure *struct* to which a cursor *curs* has been associated:

```
from curs.start until curs.after loop                                     -- [P10]
    some_operation (curs.item) ; curs.forth
end
```

with operations on cursors as illustrated below. (In early Eiffel-related literature these were operations on *struct* itself, using its own built-in internal cursor, but it is better, in particular in a concurrent programming context, to use an external cursor, so that various clients using a data structure can each maintain their own cursors, keeping different views of the same structure.)



This form assumes that the structure under iteration inherits from the Kernel Library class *LINEAR*, with the associated class for the cursor. Many usual classes such as *ARRAY* and list classes are descendants of *LINEAR*, and of course any programmer-defined class can inherit from it, defining its own implementations of the needed features *start*, *before*, *item* etc.

Form [P10] is common enough to justify an abbreviation hiding the details of cursor manipulation:

```
across struct as x loop some_operation (x) end                       -- [P11]
```


The name x serves to denote the current element, like the bound variable i in an expression of the form $\sum_{i: 1..n} p(i)$, or

$\forall i: E \text{ / } \textit{property}(i)$ -- [P12]

Quantified expressions such as [P12] are commonly needed, particularly in assertions (contracts). In fact [P12] can be written exactly as given, as well as its variant using \exists . In line with previous ideas, [P12] is an abbreviation for the keyword-based form

across *struct* **as** x **all** *property* (x) **end** -- [P13]

It is interesting to note the power of expression that these mechanisms provide. A class invariant clause in a program recently written by the author, handling a complex data structure — an array of hash tables, with integer keys, of linked lists of triples of integers — reads (given here with the original identifiers and comments):

-- [P14]

$\forall \textit{htpl}: \textit{triple_successors} \text{ /}$ -- *htpl.item*: hash table of list of triples
 $\quad \forall \textit{tpl}: \textit{htpl} \text{ /}$ -- *tpl*: list of triples (on the other hand,
-- $\textit{@tpl.key}$ is a key (a tag) for the hash table)
 $\quad \forall \textit{tp}: \textit{tpl} \text{ /}$ -- *tp*: triple
 $\quad \quad \textit{tp.tag} = \textit{@tpl.key} \wedge \textit{tp.source} = \textit{@htpl.cursor_index}$

(Note the use of “ \wedge ”, defined as an alias for **and** in the library class *BOOLEAN*. For an iteration variable c denoting a cursor, $\textit{@c}$ denotes its index in the iteration — 1 when the cursor is on the first element, 2 when it is on the second one and so on.) What is remarkable here, and not possible in any ordinary programming languages (including functional ones) is the ability to use the full power of logic (first-order), with usual mathematical notation, provided not in the form of special constructs tailor-made for particular needs but as notations for normal OO mechanisms, in this case loops iterating over any sequential data structure.

7.3 Conversions

Another application of the principle of providing general mechanisms rather than ad hoc solutions is the conversion facility. In typed languages, conversions sit somewhat outside of the type system mentioned in section A.5. Most languages allow the assignment $r := i$ for a real target r and an integer source i , causing a conversion; such cases, however, do not fit in the general type system and are treated as special.

Conversions can actually be useful for many types, not just predefined ones such as *REAL* and *INTEGER*. The principles behind the Eiffel conversion mechanism are the following:

- A conversion is an abbreviation for an object creation. You create, for example, a real from an integer, using specific rules.

- The creation, like any other, uses a creation procedure (3.4). You can mark some creation procedures as being also “conversion procedures”.
- There is also a need for “creation functions”. The reason is that instead of converting “from” we may also want to convert “to” an existing type. For example, we might want to convert a file path name, an instance of one of our classes *PATH_NAME*, to a string. *STRING*, however, is a predefined library class; we cannot add a creation (and conversion) procedure to it. Instead we define in *PATH_NAME* a function that returns a string from an instance of this class, and mark it as a conversion function.
- We say that a type converts to another if there is a conversion procedure or a conversion function available between them.
- If an assignment *target := source*, or the corresponding argument passing, would not be valid according to the other rules of the language, but the type of source converts to the type of target, then — in the spirit of the previous mechanisms, providing simpler syntax for important special cases — we interpret this assignment-like operation as a creation: an abbreviation for create *target.conv_proc(source)*, where *conv_proc* is the conversion procedure between the two types (or *target := source.conf_func* in the case of conversion through a function).
- Strict rules guarantee that there is no ambiguity: a type can convert to another in at most one way (either only one conversion procedure or only one conversion function); and, most importantly, the conversion principle: a type may not both conform and convert to another. “Conformance” is the usual type relation based on inheritance. The principle states for example that if *RECTANGLE* inherits from *POLYGON*, and hence conforms to it, then there cannot be a conversion procedure or function from *RECTANGLE* to *POLYGON*. As a result, no assignment will ever be ambiguous for either the compiler or a programmer.

This solution accommodates an important classical notion within the general design rules of the language and methodology.

7.4 Unfolded forms

Such uses of convenient syntactic variants for mechanisms that actually fit in the methodological framework of OO design with contracts is directly reflected in the language definition [13] through a device that keeps the semantic specification: unfolded forms.

The idea is that basic forms, for example the standard feature call $x.f(args)$, get a full-fledged semantic definition, and the language definition gives for any other syntactical variants an “unfolded form” expressed in terms of form basic form. For example an expression of any kind, such as an operator expression $a + b$, has an “Equivalent Dot Form” which unfolds it into a feature call of the basic kind, in this case $a.plus(b)$. Similarly, a \forall or \exists expression unfolds into an **across ... all ...** loop, and further into the basic **from** form [P10].

This language description technique makes it possible to keep the semantic definition simple, since it needs only specific a core set of constructs. Other variants are defined by translation (unfolding) to the basic forms.

7.5 Agents and functional-style programming

One area in which the general unfolding idea provides a particularly clear and concise mode of expression is the provision of functional-programming-like facilities, embedded in the OO framework, through the concept of agent, introduced in Eiffel in 1999 [11] and a precursor to others such as delegates in C# and lambdas in Java.

As an example of its use, consider a routine which (conceptually) takes another routine as an argument. For example, in an interactive system using the “command-pattern” (introduced in [4]) to provide an undo-redo mechanism, in a modern version avoiding an explosion of small classes, there can be a routine

```
do_undo (doer, undoer: ...) -- [P15]
```

whose arguments represent the operations to execute and cancel (respectively) a certain user command. (Their types are left blank above and will be specified below.) If we knew the command and hence the associated *doer* and *undoer*, the body of *do_undo* could just call them directly, as in

```
insert_line (editor) -- [P16]
cancel_line_insertion (editor)
```

but we do not know them: *doer* and *undoer* are arguments representing the corresponding do and undo operations for an arbitrary command (that is the whole idea of the pattern). Even so it would be desirable to use the same syntax

```
doer (editor) -- [P17]
undoer (editor)
```

where *doer* and *undoer* are, as noted, formal arguments of the enclosing routine [P15].

We note in passing that C comes close to allowing such a notation. In C the arguments would simply represent addresses for the corresponding code; we would call *do_undo* with actual arguments such as

```
&insert_line, &cancel_line_insertion -- [P18]
```

then *do_undo* itself can call the associated routine as

```
(*doer)(editor) -- [P19]
```

The arguments in this low-level mechanism are just addresses and it is easy to circumvent type restrictions (at the address given by the argument the memory could contain some unwanted code, or no code at all). Higher-level mechanisms such as Eiffel’s agents and C#’s delegates provide type-safe versions of this facility.

How do we allow a routine (such as `do_undo`) to receive routines as arguments? In an object-oriented context the only values programs ever manipulates are objects, or references to objects. A routine — a piece of computation — is not an object; in fact the object-computation duality lies at the basis of the object-oriented method (which defines the architecture of software systems by attaching computations to object types). The solution then is to define objects that represent computation, and hence can be used in data structures. Such objects are called agents.

The common type of all agents is

ROUTINE [*ARGS* → *TUPLE*]

where *ROUTINE* is a Kernel Library class devised for this purpose. *ARGS*, the generic parameter, represents the types of arguments of an associated routine. For example instances of *ROUTINE* [*TUPLE* [*INTEGER*, *INTEGER*]] are agents associated with routines of two integer arguments. To understand this class specification note that:

- A generic class can have a constrained generic parameter: if its declaration starts with class *C* [*G* → *CT*], where *CT* is a type, then generic derivation *C* [*T*] are only valid if the actual generic parameter *T*, also a type, conforms to *CT* (its base class is a descendant of the base class of *CT*). For *ROUTINE* the constraining type is *TUPLE*.
- *TUPLE* is a language-defined type describing sequences (tuples) of values. There are an infinity of tuple types: *TUPLE* covers all tuples; *TUPLE* [*T*], for a type *T*, covers tuples of at least of element, the first of which is of a type conforming to *T*; *TUPLE* [*T*, *U*], tuples of at least two elements starting with values conforming to *T* and *U*; and so on. The type conformance relation is hence that all tuple types including *TUPLE* [*T*] conform to *TUPLE*, then *TUPLE* [*T*, *U*] conforms to *TUPLE* [*T*] and so on.

ROUTINE is generally not used directly but through its descendant classes

PROCEDURE [*ARGS* → *TUPLE*]

FUNCTION [*ARGS* → *TUPLE*, *RES*]

representing procedures and functions respectively. Values of those types are obtained through the **agent** notation; for example with *p* a procedure with two integer arguments and *f* a function with an integer argument and a boolean result we may use

agent *p* -- of type *PROCEDURE* [*INTEGER*, *INTEGER*] -- [P20]
agent *f* -- of type *FUNCTION* [*INTEGER*, *BOOLEAN*]

The basic activity of object-oriented design and programming— look for data abstractions — again pays off with agents: an agent such as **agent** *p* is not a spiffed-up costume for a *&p* function address but truly an object in the OO sense, representing a routine equipped with applicable operations; one of these operations, as seen next, is to call the associated routine, but it is not the only one. Classes *ROUTINE* and its descendants con-

tain other features describing properties the associated routines, and classes covering sequential data structures contain iterator features such as

```
do_if (action: PROCEDURE [G]; test: FUNCTION [G, BOOLEAN])
  -- Apply action on all elements of the structure satisfying test.
```

Details of the agent and tuple mechanisms, and in particular the associated type system properties, are not essential for the rest of this presentation, but let us focus again on notational simplifications. For an agent *a*, the basic features include *call*, which calls the associated routine, and, for a function, *last_result* (which returns the last result produced by a call) and *item* (which calls call and returns the value of *last_result*). So if we have a routine with arguments of the signature given in the example of [P20]

```
r (a: PROCEDURE [INTEGER, INTEGER];
   b: FUNCTION [INTEGER, BOOLEAN])
```

its body can call *a.call* ([0, 1]) and *b.item* ([0]) (the latter returning a boolean). Similarly, for the *do_undo* routine [P15], which we can now specify with an appropriate signature

```
do_undo (doer, undoer: PROCEDURE [TUPLE [EDITOR]]) -- [P21]
```

a typical call will be

```
do_undo (agent insert_line, agent cancel_insert_line) -- [P22]
```

The body of *do_undo* may itself contain calls such as

```
doer.call ([editor]) -- [P23]
```

which when executed as part of the call [P20] will mean

```
(agent insert_line).call ([editor]) -- [P24]
```

and in this case has the same effect as the direct call

```
insert_line (editor) -- [P25]
```

although of course we do not know that in the body of *do_undo*, where *doer* could represent any other routine with the appropriate signature.

The notations for using the mechanism [P21] and [P20], are consistent and the underlying structures clear, but the notation is heavy when used extensively, with significant noise (*TUPLE*, brackets). It forces the programmer to remember all the time the underlying intellectual construction, even when the goal is simply to pass a routine *p* to a routine *r* and let *r* call *p*. Once again syntactical simplifications intervene. The first is the permission to omit the brackets and the *TUPLE* specification in [P20], giving simply

```
do_undo (doer, undoer: PROCEDURE [EDITOR]) -- [P26]
```

Since *PROCEDURE* expects a tuple type as generic parameter, the syntactical simplification does not introduce any ambiguity; [P26] is simply an abbreviation for [P21], which remains usable for someone preferring to use a purist form. Form carries the right

impression: that *do_undo* expects two arguments, either of which represents a procedure with an argument of type *EDITOR*. The role of tuples is essential to make the mechanism type-safe, and consistent with the rest of the OO type system, but there is no need to force the programmer to remember every time that arguments will be conceptually grouped into a tuple.

In the same spirit, since *call* always expects a tuple argument, we can drop the brackets on the call side, in [P23], giving

doer.call (editor) -- [P27]

The same observation applies: for the programmer, there are no tuples, only arguments, in this case just one, *editor*, with no need to remember that it is internally treated as a one-element tuple. (Precautions are taken to ensure that there is no ambiguity in the rare case of a routine with an argument that is actually of a tuple type. Here we reap the advantages of a strong type system which enables the analysis to avoid any confusion.)

The presence of “. *call*” is still in some sense noise since it forces the programmer to remember that the argument *doer* is of an agent type with its specific properties. In simple cases the programmer simply wants to think of *doer* as a routine (rather than the correct and more general but also more elaborate notion of an object encapsulating a routine). Once again the type system helps: we just allow, with the guarantee of non-ambiguity enforced by type analysis, dropping the “. *call*” part to allow, as a synonym for [P27] and hence another synonym for the fully explicit version [P23], the form

doer (editor) -- [P28]

which is what we wanted in the first place [P17], with the same simple form as a direct call to a known routine [P25], and similar to a call to a routine in C [P18][P19] but with full type safety, a higher level of reasoning, and no need for messy pointer manipulations through the *&* and *** symbols.

This example again illustrates the general idea of devising mechanisms that take full advantage of OO concepts, particularly the constant reliance on data abstraction, while providing convenient and concise syntactical forms.

7.6 About the role of agents: functional mechanisms in an OO shell

Beyond notation, some comments on agents are appropriate.

Functional languages have an enthusiastic following, due to their closeness to mathematical concepts and their consequent rejection of imperative constructs. It is important to note, however, that functional programming as practiced typically involves imperative aspects, if only to model stateful operations of the program’s domain (printing is an elementary but typical example). To unabashed imperative programmers, the benefit of rejecting assignment — a notion that after all has clear axiomatic models — in the name of simplicity and elegance, only later to bring in such notions as monads, possibly not the most intuitive possible for non-computer-science-PhDs, sounds debatable. These

matters are subject to opinion, however, and the elegance of a functional solution to certain algorithmic problems is undeniable.

In understanding how OO and functional ideas can coexist and support each other, it is important to understand where their strengths lie. Object-oriented concepts are primarily a system structuring mechanism, based on the identification and classification of data abstractions as the backbone of a software architecture. They are unrivaled when it comes to building large systems which will be maintainable and extendible, and their parts reusable. Functional programming has no particular contribution here. The proper relationship — as discussed in the detailed analysis of an earlier article [14] — is to use functional elements in an OO architecture. Agents provide these elements.

Since agents define abstractions of operations, including functions, much of what one can do in a functional programming language — define and manipulate functions, including higher-level functionals — is readily doable in Eiffel. So you get the best of all worlds. Not *all* of the second world; there is no direct equivalent, in particular, of the elegant mechanism of modern functional languages, from Miranda on, for defining functions by inductive pattern matching rather than programming-style recursion. (Here is the result if the argument is an empty list, and there it is in terms of a and L if it is the concatenation of an element a and a list L .)

8 EXCEPTION HANDLING

Modern imperative programming languages all have an exception mechanism, typically based on the `try A catch B` scheme of C++, which means: execute A , but if an exception occurs handle it in B .

8.1 Modes of exception use

Try-catch mechanisms provide a reasonable way to cover the role of exceptions, if we define that role as handling abnormal run-time events. Such a definition, however, is so vague as to be almost useless. What is an “abnormal event”? What is normal to you can be highly abnormal to you. The observation of both programming textbooks and programming practice in try-catch languages shows a wildly diverse range of uses of exceptions:

- At one end, exceptions serve to recover from truly unpredictable run-time situations, such as arithmetic overflow (hard to predict in most programming languages), running out of memory, a file suddenly becoming unavailable while being read.
- At the other extreme, exceptions are just a control structure, another way of providing a `goto` while pretending all along (section 2) that the language supports “structured programming”. An example — not made up! — is a programming style in which the routine searching for an element in a hash table produces an exception when it does not find it.

While the last example is clearly an aberration, it is not always so clear to a practitioner when exceptions are a legitimate solution in preference to more sober control structures. Programmers need clear methodological guidelines defining the role of exceptions in the construction of programs.

Java makes a step in this direction by enabling and in fact requiring the author of a routine (method) to specify which types of exception it may trigger and not handle by itself; then any calling routine must handle the corresponding exception types (by naming them in a `catch` clause). Unfortunately these obligations only apply to programmer-raised (“check”) exceptions — those actually used as substitutes for `goto` instructions or multiple returns — and do not help for the crucial cases, those involving exceptions occurring unpredictably at run time.

8.2 A clear conceptual framework for exception handling

Eiffel introduced a methodology for using exceptions properly. (The first presentation in 1988 could not find a publication outlet and remained a technical report [5]; the ideas were explained in chapter of [10].) The idea is simple. Any operation is characterized by a contract: it assumes certain initial properties and commits, under the assumption that they initially hold, to delivering certain final properties. The contract is the combination of these initial and final properties (precondition and postcondition). That characterization applies to operations at any level, from the topmost — say, launching a rocket, a highly complex operation — down to a single routine with its explicit precondition and postcondition in the code and to a basic operation such as adding two integers (with the implicit precondition that the sum fits in the computer representation, and the postcondition that the result is the mathematical sum). The need for exception handling arises when the operation is — for any reason — not able to meet this commitment. The first step towards a methodology is to distinguish between two concepts:

- An **exception** is a run-time event that prevents an operation from executing the scenario — specifically in software, the algorithm — in which it is engaged to fulfill its contract (the just mentioned commitment).
- A **failure** is an operation’s inability to fulfill its contract altogether.

The reason the two concepts are different is that there may be more than one way to fulfill the contract. The operation may have, metaphorically, lost a battle but not the war. The reason exception mechanisms exist is that the designer of the operation may have provided alternative scenarios to achieve the contract in case one of them fails.

The next methodological step is to examine what courses of action make sense in the system element — in software, a routine — whose execution encounters an exception. In light of the previous observations, there are only two allowable responses:

- **Resumption**: start an another scenario (another way to attempt to ensure the contract).
- **Failure**: Give up and concede failure (accepting the loss of not only the specific battle but the routine’s “war”).

(A third strategy, “**oblivion**”, is in principle possible: ignore the exception and return as if nothing had happened. It can only be justified in rare cases; think for example of a spurious exception that the the window system generates if the interactive user resized the window during a computation, if that computation does not care about the user interface. The policy should be reserved for exception types that have been explicitly marked as “ignorable”. Eiffel does support it as explained below.)

The third and final methodological rule is that the failure of a routine, the second case above, causes an exception in the routine’s caller. (If there is no caller, meaning that the failing routine is the top-level element of the execution, the program execution as a whole fails — terminates abnormally.) Any exception is of a certain type, such as arithmetic overflow; the secondary exception created in the routine failure case is generally of a different type from the original exception, the type “routine failure”.

8.3 Eiffel exception mechanism

The language mechanism directly follows from the analysis. It consists of two parts:

- An effective routine can include, after its body (**do** clause), a **rescue** clause, whose instructions take over when an execution of its body triggers an exception. Termination of the rescue clause causes either failure or resumption.
- To decide between those two possibilities, a integer local variable **Retried** (pre-defined, like **Result** for functions from section 2) is available. If its value is 0, as it is by default, the **rescue** clause ends in the routine failing (and, as explained above, triggering an exception in its caller). Otherwise, the routine executes its body again, *without* reinitializing the local variables (including **Result** and of course **Retried** itself).

Here is a typical example of use:

```

transmit_if_possible (f: FILE; max: NATURAL)                                -- [P29]
  -- Attempt to transmit f over an unreliable line, in at most max attempts;
  -- produce a message if impossible.
do
  if Retried < max then
    attempt_transmission (f)
    -- We assume that attempt_transmission, a low-level transmission
    -- routine, may cause an exception.
  else
    print_message
  end
rescue
  Retried := Retried + 1
end

```

This routine never fails since it can always satisfy its contract, which states that it should either transmit the file or print a message. In contrast, the following shorter version (same header) can fail:

```

do
    attempt_transmission (f)
    -- We assume that attempt_transmission, a low-level transmission
    -- routine, may cause an exception.
rescue
    if Retried < max then
        Retried := Retried + 1
    else
        Retried := 0
    end
end
end
-- [P30]
```

If **Retried** reaches the maximum the rescue clause resets it zero and hence terminates with a zero value for it, causing the routine to fail.

Fine control of the mechanism is available to the programmer:

- Each exception has a type. There are a number of system-defined exception types, such as “arithmetic overflow”, can programmers can define their own.
- An operation is available (through a library routine *raise*) to trigger an exception of a programmer-defined type.
- The program can discriminate between exceptions of different types, to handle them differently.
- If an exception of a certain type, such as “arithmetic overflow”, causes a routine to fail (in the absence of resumption), the result, as seen above, is to pass an exception to the calling routine. That exception is of the type “routine failure”. The program has access both to the original exception type and to the derived one.
- Exceptions of certain types are ignorable. The programmer can specify, again through a library mechanism, that execution will actually ignore (or not) exceptions of an ignorable type, thereby allowing the “oblivion” policy (8.2) when harmless.

It is also interesting to see how the methodological reasoning extends to a concurrent context, raising among others the question of what should happen to an exception raised by a routine whose caller was asynchronous (in another thread or process) and may even have ceased its execution. See [19] for such an extension to the concurrency mechanism described in the next section.

Clearly any try-catch-style exception handling can be emulated in the Eiffel rescue-retry style. The other way around too, although only by writing fairly complicated control structures to represent resumption. The spirit is, however, different.

Like other mechanisms reviewed in this article, the exception handling constructs address an important design and implementation need in conformance with principles of modern software development, and provide a clear methodological framework for applying the ideas.

9 CONCURRENCY

Concurrent programming used to be a specialist endeavor for people writing operating systems or networking applications, but is practiced today much more widely, thanks in particular to the spread of multi-threading technology. Here too it is important to maintain a consistent level of abstraction, by providing mechanisms that are compatible with modern principles of methodology, rather than taking programmers through a roller-coaster ride between wildly different level of abstractions.

9.1 Semantic mismatches

An example of the roller coaster is the typical multithreading mechanism of many OO languages. While some modern languages (such as Go) were designed specifically for concurrency, much of the concurrent programming performed in object-oriented languages uses “thread libraries”, which provide low-level mechanisms; semaphores are the principal means of synchronization.

The semantic gap is enormous. Object-oriented languages emphasize abstraction; their control structures, even with remnants of the goto, limit the potential for “spaghetti code”; if they are statically typed, the type system provides powerful means of expression, and avoids many errors. In contrast, the semaphore is a concept from the 1960s, a goto-like low-level construct leaving opening the way for many concurrency errors. It makes it hard in particular to avoid **race conditions**: cases of data corruption, leading to and incorrect results, due to different threads writing and accessing the same data, each in an order that is correct for the thread taken independently but incorrect in combination with the actions of the other threads (as in a carelessly designed reservation system in which two customers find the same airplane seat initially available and both reserve it).

9.2 An example: what happens to preconditions?

Initial enquiries that led to the design of Eiffel’s concurrency mechanism, SCOOP, examined [8] the following question, a specific case of the race condition issue: what

does a precondition mean in a concurrent context? Consider a typical precondition-equipped routine, inserting an element into some data structure *s*:

```

put (s: STRUCTURE; e: ELEMENT)                                -- [P31]
    require
        not s.is_full
    do
        ... Implementation of the insertion (see below)...
    ensure
        s.has (e)
        not s.is_empty
    end

```

The precondition in [P31] entitles the code for “Implementation of the insertion” to include operations such as

```

s.insert (x)                                                    -- [P32]

```

where *insert* itself also has the precondition **not** *s.is_full*, since this operation is executed as part of *put* whose precondition in [P31] guarantees that same property.

Under the basic ideas of Design by Contract, the precondition of a routine such as *put* is not just a requirement on the caller but also a guarantee to the caller that it is the *only* requirement: the supplier (here the routine *r*) *commits* to handling correctly any call that satisfies the precondition. Examples of such calls include

```

if not s.is_full then s.put (x) end                            -- [P33]

```

and

```

s.remove ; s.put (x)                                           -- [P34]

```

assuming that *remove* (dually with *put*) has the postcondition **not** *is_full*. Specifically, the code corresponding to “where insert has the precondition , since that precondition is guaranteed by the caller

Such modes of reasoning are fundamental to our ability to write meaningful programs: if a supplier states a condition for doing its work, and you can satisfy that condition, then you will get the promised result.

In a concurrent settings, this pillar of civilization crumbles!

If the client (the caller executing [P33] or [P34]) and the supplier (the code for the routine, executing operations on *s*, such as *insert* in [P32]) are running in different threads, the guarantee evaporates. Between the time you ascertain **not** *s.is_full* in [P33] (or call *remove* in [P34]) and the time you call *s.put* (*x*), some other thread may access *s* and make it full again.

The solution to this issue in SCOOP appears in section 9.8 below. More generally, this discussion illustrates the spirit that presided over the design of a concurrency mechanism for Eiffel, in line with the general goals stated in the introduction

9.3 Sequential versus concurrent

Some approaches to concurrency start from the position that everything should be concurrent by default; sequential computing is just a special case. While this argument is by itself correct, in practice concurrent reasoning is difficult. The development of programming methodology over the past decades has reached a stage at which we master the processes of sequential programming to a good extent; bringing unrestrained concurrency into the picture makes things far more difficult.

The precondition example is a typical case. We are used to thinking that if it has been ascertained that operation *B* will perform correctly under condition *C*, and that *A* ensures *C*, we may comfortably execute *A* then *B*. With the interleaving of operations permitted by concurrency, in particular by multithreading, this reasoning no longer holds. No other form of reasoning presents itself as an obvious replacement.

As a consequence of these observations SCOOP lets the programmer use, as much as possible, the modes of reasoning that work in sequential computing, circumscribing the concurrent aspects to specific combinations of sequential mechanisms. In addition, SCOOP provides expressive power that is seldom available in other frameworks.

9.4 Simultaneous resource reservation

An important aspect of the approach is the ability to reserve several resources at once. The implementation of the classic dining-philosophers problem reads, in the scheme for each philosopher:

```

separate left, right as l, r do                                     -- [P35]
    l.pick ; r.pick ; eat (l, r) ; l.put_down ; r.put_down
end

```

The **separate** construct gets exclusive hold of the objects listed, giving them local names (*l* and *r*). The body of the construct can then manipulate these objects safely — as in a sequential context — since they are under its exclusive control. (For completeness of the example we may assume that *eat* has preconditions *l.picked* and *r.picked*, *pick* has a postcondition *picked*, and *put_down* a postcondition **not** *picked*.)

[P35] (embedded in a loop for each philosopher) is essentially the entire solution to the dining-philosopher problem, shorter and simpler than any of those usually proposed. The ability to obtain exclusive hold of an arbitrary set of objects frees the programmer from complicated, error-prone (in particular, deadlock-prone) algorithms; the mechanism is handled by the SCOOP runtime, providing programmers with a significant gain in expressive power.

9.5 Processors and regions

“Resource” as used above is not a special concurrent concept. Neither does SCOOP need any special notion of “concurrent object”. SCOOP simply uses Eiffel object; the only fundamental new concept for concurrency is that of a “processor”.

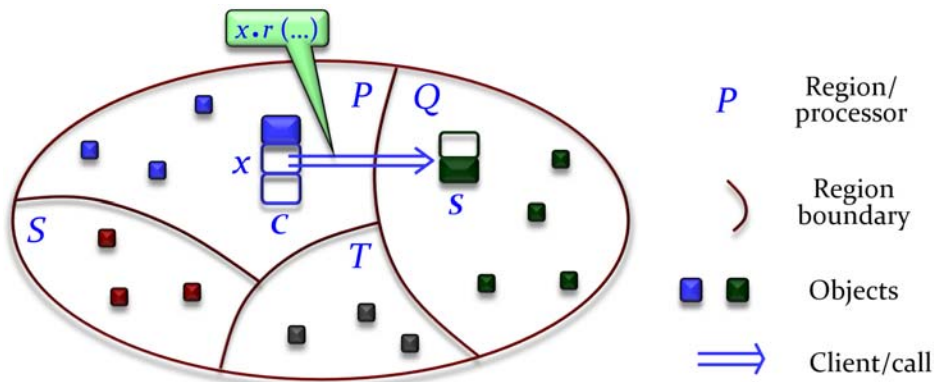
Introducing active objects, as in some OO concurrency frameworks, violates the basic scheme of object-oriented programming, in which an object is a kind of server, providing a number of services (its features) to other objects, its clients, through feature calls. Such an object cannot have its own scenario (as an active object would have), since it would immediately clash with the scenarios of its clients when they call its features; these clashes lead to complicated synchronization mechanisms. Concurrent OO programming does not need any of that. It suffices to extend the sequential framework by adding that each object has an associated **processor**, a sequential mechanism in charge of executing all the features called on the object. The model is general enough to accommodate various realizations of the notion of processor; it could be a physical CPU or a software process. In current implementation, processors are threads. A processor is sequential, like a traditional computer; concurrency comes from having several. The first major difference between sequential and concurrent computation is this move from one to several processors.

9.6 Synchrony, asynchrony and wait by necessity

Since calls on any particular object are handled by its processor, we may see the set of objects as partitioned into a number of “regions”, each associated with a processor. The second major difference between sequential and concurrent comes with the semantics of the basic OO operation, feature call:

$x.r(\dots)$ -- [P36]

This call will be executed by a processor P , handling a certain client object C ; it requests the execution of the routine r , on the supplier object S denoted by x (where x may be a field of the client object, as in the figure). Let Q be the processor handling that object. We say that the call is non-separate if Q is P (the two objects are in the same region), and separate if they are different processors, as shown.



Considering the effect of the call:

- For a non-separate call, only one processor is available, the only possibility is the same as in sequential programming (which has only one processor available for the entire computation). The call execution has to be *synchronous*: the client does not proceed until the body of r has been executed.
- For a separate call, *asynchrony* is possible: the client computation can continue, without waiting for the execution of r to complete.

The asynchronous case leads to a distinction (due to Benjamin Morandi) between two parts of the semantics of a call such as [P36]: *call* (proper) and *application*. The caller executes the call, which means that it registers its request to have r executed; the supplier executes, or “applies”, the body of r . For synchronous calls, including sequential programming, the difference does not matter, since call causes immediate application; with asynchrony, the client can proceed with further instructions immediately after having executed the call. The relationship is subject to the following consistency rules:

- **Correctness**: at the time an application of a routine r ends, the postcondition of r holds.
- **Causality**: the application of a call follows it. (The application may occur at any time after the call, but not before.)
- **Order preservation**: if two calls and involve the same client processor and the same supplier processor, their applications happen in the same order as the calls.

Order preservation is limited to this case of interaction between the same two processors: no order guarantee exists for the application of calls from different clients and suppliers. In other words, there is no global clock.

Separate calls permit asynchrony, but do not require it: a synchronous semantics would still be possible (although it might lead to deadlocks). When are calls actually asynchronous? The answer makes it possible to answer another question: how does a client catch up with an asynchronous supplier? If the client started a call, such as *New_York.sell* (“*Microsoft*”, 300), it is for a purpose (in this case, getting rid of some shares); the client may need at some point to make sure that the purpose has been achieved, in other words that the application — and not just the call — has executed to completion. It does so by executing a *query* (A.6, 1.1), for example *New_York.portfolio* (“*Microsoft*”). The rule, known as **wait by necessity** (it could also be called “lazy wait” and is originally due to Denis Caromel) is that waiting for resynchronization occurs in this case, and in this case only. The call will cause the client to wait until all previous separate calls with the same target (here *New_York*) have completed their application, and the application of the current query (*portfolio*) is also complete.

For a separate call:

- If the routine is a command (or procedure, see A.6), the call is asynchronous.
- If it is a query, the call is synchronous.

A typical pattern consists of firing a succession of command calls to tell suppliers to perform operations, then executing query calls to take advantage of the results.

9.7 Language support

Since separate and non-separate calls have different semantics (possibly asynchronous versus guaranteed synchronous), the source text must express which version applies. The type system takes care of this matter (as of many others). A declaration of the form

x: **separate** *T* -- [P37]

(as opposed to the default *x*: *T*) states that the object denoted by *x* may be separate (reside in a different region, handled by a different processor). Procedure calls of target *x* may then be asynchronous (their application is decoupled from the call), and query calls of target *x* will force the client to wait until the query completes.

A **separate** declaration such as [P37] does not specify the processor of *x* and does not even require it to be a different processor: separate means “potentially separate”. The processor could actually be the same in some executions. What the declaration specifies is that asynchronous execution is possible for calls of target *x*.

The extension of Eiffel from a sequential programming language to a concurrent one uses only one keyword, **separate** (also used in the resource reservation construct [P35], although, as seen next, only as an abbreviation in that case).

In line with the general observation that the properties of the language are largely embodied in the type system (A.5), type rules are associated with the preceding concepts: in particular, you can assign from non-separate to separate (remember that separate means possibly in a different region) but not the other way around. Associated type rules apply to calls and were elegantly defined in Piotr Nienaltowski’s ETH PhD thesis (se.ethz.ch/people/nienaltowski/papers/thesis.pdf).

9.8 Synchronization

Consider a routine whose formal arguments include some of separate types, as in an adaptation of [P31] using a separate argument:

```

put (s: separate STRUCTURE; e: ELEMENT) -- [P38]
require
  not s.is_full
... The rest as in [P31] ...

```

Then a call such as *put* (*struct*, *x*) will proceed when both:

- S1 The processors corresponding to all separate arguments, in this case *s*, are available.
- S2 The corresponding separate preconditions are satisfied.

A “separate precondition” is a precondition clause containing an expression *x.some_property* where *x* is a separate formal argument, as is the case with the precondition here.

This rule provides the basic synchronization mechanism. It is in line with the discussion of section 9.2 and answers the question raised there, “what happens to preconditions?”. Non-separate preconditions retain their sequential semantics of correctness conditions, but separate preconditions such as the one for *put* in [P38] become wait conditions. It is important to point out that this decision is not arbitrary; it appears in fact to be the only reasonable choice, since the standard correctness-condition semantics is just impossible to apply, as shown in the discussion of 9.2. It also provides a powerful synchronization mechanism; note in particular that the scheme for *put* in [P38], together with the corresponding one for a remove operation

```

put (s: separate STRUCTURE)
  require
    not s.is_empty
  do
    ... Implementation of the removal...
  ensure
    not s.is_empty
end

```

provides the basic implementation of a producer-consumer system with a shared buffer, expressed more simply than solutions in the classical literature.

An important component of this semantic specification is that the reservation of resources corresponding to rule S1 above (wait until all processors for separate arguments are available) applies to *any number* of separate arguments. A philosopher object may for example call *step* (*left*, *right*) with the routine definition

```

step (l, r: FORK) -- [P39]
  -- Execute the basic cycle step of a philosopher with l and r.
  do l.pick ; r.pick ; eat (l, r) ; l.put_down ; r.put_down end

```

providing a variant of the implementation using the **separate** construct in [P35]. After a call to *step*, application of the routine’s body will only proceed when it gets exclusive hold on the processors corresponding to the two arguments. As was noted in the presentation of the **separate** construct, this facility achieves high expressive power by providing a mechanism for simultaneous reservation of multiple resources. The soundness of the mechanism is guaranteed by the rule which actually *requires* a separate target *x*, in a call *x.r(...)*, to be a separate formal argument of the enclosing routine, or a local name defined in an enclosing **separate** construct.

The form with the feature call, [P39], is actually the basic construct. The **separate** construct (in the [P35] example, **separate left, right as l, r do eat (l, r) end**) is actually just an abbreviation for the call *step (left, right)* and the routine definition of *step* in [P39], in the spirit of the many “unfolded forms” of Eiffel (discussed in section 7.4). It makes it possible, when an operation needs to work on separate expressions, to write a wrapper routine, such as *step*, for the operation. But the semantics is the same. The routine form remains necessary in the case of wait conditions, which the routine will express as separate preconditions.

SCOOP has proved to be a practical and safe solution to concurrent programming using threads, letting programmers apply the same principles as in the rest of OO design and programming with Eiffel.

10 VOID SAFETY

The description of the last mechanism in this review will be only a sketch, even though it is one of the most important characteristics of today’s Eiffel.

Eiffel is (but was not always) void-safe: a null-pointer dereference cannot occur.

The dereferencing of a null pointer (C terminology), or void reference (OO terminology for a more abstract notion of pointer) is an attempt by the execution of a program to execute or evaluate $x.a$ for an attribute a or $x.r(\dots)$ for a routine r in state in which x , a pointer or reference is null, or void, meaning that it does not denote any object. Null-pointer dereferencing is also known as a “void call”.

In C and other languages that do not have a full-fledged exception mechanism, the consequence is to crash the program. In exception-equipped languages a void call causes an exception, which often has the same result anyway in the absence of a handler specifically written for that case. (The reason programs typically do *not* catch such exceptions is that if the programmer realizes x can be void it is much easier to fix that problem, handling the case of a void x through an ordinary test in a conditional instruction, than to let the exception happen and try to recover from it. As a result, the void calls that do remain are those the programmers did not foresee, directly leading to a crash.)

In the Common Vulnerabilities and Exposure (CVE) database, hundreds of the most egregious cases of security violations, some in widely used products from such companies as Apple, Cisco, Google, IBM and Microsoft and technologies such as JPEG involve null-pointer dereferencing, as analyzed in Alexander Kogtenkov’s ETH PhD thesis (se.ethz.ch/people/kogtenkov/thesis.pdf). CVE mostly refers to C programs, but the problem persists in all the dominant OO languages, including the most strictly statically-typed ones. Some recent designs, such as Kotlin, provide a partial solution, but we know of no other widely available OO language than Eiffel that provides a full solution, mechanically proven correct as part of the aforementioned thesis.

One should note in passing that functional programming does not solve the problem. Functional languages do address it in part by removing the notion of pointer or reference and providing instead built-in classes for common data structures such as lists, hiding pointer manipulations from the programmer. But any modeling of complex data structures not covered by the library (general graphs) would require the introduction of a pointer-like mechanism.

The Eiffel solution (which was initially influenced by early attempts in research languages, notably C-omega from Microsoft Research) was initially described in in [12], with the ironed-out and fully-implemented version explained a few years later in [17] (and Kogtenkov’s thesis). Like many other critical mechanisms, it essentially consists of enriching and tightening the type system (A.5).

The starting point was the realization that for reference types void (or null) is an anomaly, not a basic case. The reason programming languages typically accept that references can be void is computer scientists’ familiarity with linked data structures, for which it is customary (and often appropriate) to terminate a structure by linking the last element to void. But in application domains, there is usually no real need or room for void values. If we have a class *EMPLOYEE* in a payroll system, what is the point of an void *EMPLOYEE* value? In fact, none at all.

As a result of this observation we consider that reference values should by default be “attached”, meaning guaranteed always to denote an object at run time; the case of a “detachable” expression, which can be void, remains the exception. So if we declare

e: *EMPLOYEE* -- [P40]

we are making a *commitment* that the corresponding object or objects will always exist. If we want to allow an expression to become void, for example for a list cell, we say it explicitly by declaring it as

next: **detachable** *CELL* -- [P41]

(In we may use the explicit form **attached** *EMPLOYEE*, but this is not necessary since “attached” is the default for the reasons noted.)

In the attached case, [P39], it is our responsibility to meet the commitment, and of course the compiler’s responsibility (backed by the language definition and the aforementioned formal proof of soundness) to enforce it, to guarantee that *e* will never be void when execution attempts to apply a feature to it. The precise rules will not be described here (see the cited references and the language standard), but it is important to mention the two principal properties:

- Once again, the type system is key. The basic rule (compare with the SCOOP rule in 9.7) is that you can assign from attached to detachable but not the other way around.
- The most delicate part has to do with initialization. Every object should be initialized upon creation; for fields not explicitly initialized by the creation procedure, the lan-

guage relies on default values for each type, but for an attached type other than basic types there is no universal default; the program must provide it.

The crux of the mechanism’s viability is its convenience. An extreme solution, protecting every call $x.r$ by a test for the possibly vacuity of x , might be theoretically suitable (although it only pushes the problem, as you have to write the “else” part of the condition instruction), but is of the “straitjacket” kind and not palatable in practice. On the other hand, we cannot just rely on the assumption that the compiler is “smart” and finds out by itself which cases are harmless (not void-safe-prone): the language is defined by a precise specification, not the behavior of one compiler or another. The rules have to be simple, and equally clear to the programmer and the compiler writer.

This need for a delicate balance between expressiveness and safety, constrained by the requirement to have a crystal-clear language specification and a provably sound mechanism, explains the delay of about a decade between the time the essential components of the solution were first published [12] and when the mechanism became a fully-integrated part of the Eiffel approach. The challenge was to guarantee void safety without imposing an undue burden on the programmer.

The discussion was biased initially by the difficulty of *converting* existing non-void-safe code to void safety. Once the key elements of existing software, in particular many thousands of lines of library code, were made void-safe, the focus could turn to the more important and, as it turns out, more easily attainable goal: writing code that is void-safe in the first place. Once one has got the knack of void-safe programming, this task becomes natural. While it requires a specific mode of thinking, with a new tack on “plain old” OO programming, it is not hard to practice and leads to programs that are more readable and convincing — in addition to providing a major guarantee against a risk that existed in pre-void-safe Eiffel and still does in other approaches: that a program that seems to have been carefully designed, seriously reviewed and extensively tested will some day, unpredictably, make a catastrophic attempt to dereference a null pointer.

APPENDIX A: LANGUAGE DESCRIPTION

A distinctive feature of Eiffel concerns, rather than a language construct, the way the language is defined. The Ecma and ISO standard [13] follows a strict discipline of distinguishing between:

- The lexical level: individual tokens.
- The syntactic level: program structure (for a program made of correct tokens). The syntax description uses BNF-E, a variant of BNF restricting, for simplicity and readability, the definition of each construct through a single production which is one of: choice, concatenation or iteration. (Usual BNF can mix these variants, for example by defining a construct as a choice between concatenations, as in $Z = (U V) | (X Y)$,

which in BNF-E requires three productions, one defining Z as $A \mid B$, one A as $U V$ and one B as $X Y$.)

- The static semantic level: validity constraints (for a program that is syntactically correct). A validity constraint is a static requirement on program texts that is not expressible (or not conveniently expressible) in BNF. Type rules are a typical example.
- The semantic level: effect of constructs (if they are valid).

The figure below, extracted from the standard, shows a typical use of the last three levels. Two further kinds of elements appear in the standard:

- Informative text: explanations, often including programming examples. Although not formally part of the language definition, they help make the standard readable.
- Definitions (not appearing in the example), introducing auxiliary notions — at any of the above four levels — in a precise way.

8.9.7	Syntax: “Old” postcondition expressions Old \triangleq old Expression	
8.9.8	Validity: Old Expression rule An Old expression oe of the form old e is valid if and only if it satisfies the following conditions:	Validity code: VAOX
	<ol style="list-style-type: none"> 1 It appears in a Postcondition part $post$ of a feature. 2 It does not involve Result. 3 Replacing oe by e in $post$ yields a valid Postcondition. 	
	<i>Informative text</i>	
	<p>Result is otherwise permitted in postconditions, but condition 2 rules it out since its value is meaningless on entry to the routine. Condition 3 simply states that old e is valid in a postcondition if e itself is. The expression e may not, for example, involve any local variables (although it might include Result were it not for condition 2), but may refer to features of the class and formal arguments of the routine.</p>	
	<i>End</i>	
8.9.9	Semantics: Old Expression Semantics, associated variable, associated exception marker The effect of including an Old expression oe in a Postcondition of an <u>effective feature</u> f is equivalent to replacing the semantics of its Feature_body by the effect of a call to a fictitious routine possessing a local variable av , called the associated variable of oe , and semantics defined by the following succession of steps:	
	<ol style="list-style-type: none"> 1 Evaluate oe. 2 If this evaluation <u>triggers</u> an <u>exception</u>, record this event in an associated exception marker for oe. 3 Otherwise, assign the value of oe to av. 4 Proceed with the original semantics. 	

This specification discipline particularly applies to the validity constraints (static semantics), where it prescribes an “if and only if” style as illustrated in the extract shown for the “Old expression” construct. All programming languages (even those without a particularly strong type system) have validity rules, stating that certain syntactically correct permutations are permitted and others not. It is easy and tempting for the language

designer to focus on the “only if” part, by stating certain obligations that the program must meet. For example, the type of e in an assignment $v := e$ must conform to the type of v . It is also natural to give examples that are valid and others that are not. But the demands on a language specification are higher: a rule must specify precisely — for the benefit of both programmers and compiler when exactly a tentative program element is valid, and when not. Hence the “if and only if” style, illustrated by the VAOX rule above which states under what condition an “Old expression” is valid.

While it is easy to tell programmers pieces of what they may do and especially may not do (use an Old expression outside of a postcondition, use Result in it), the “if” part is much harder since it represents a commitment by the language definition, which the compiler writer must honor: *if* you meet all my conditions, *then* I promise, as part of my (the compiler writer’s) contract with you (the programmer) that I will be able to process your program properly so that its execution will conform to the specification appearing in the “Semantics” part.

While tough to apply, this discipline is beneficial not only for programmers but also for the quality language definition, which it encourages to leave no stone unturned.

Examination of language definitions for today’s principal languages shows that this style is not common and that the more usual approach is to list some conditions, in an “only if” style, without a guarantee that they are all the conditions. In general, language standards tend to be less rigorous. The Eiffel specification is not formal in the sense of a mathematical specification, but through the traits discussed here it is strongly influenced by practice with formal specifications and retain some of their benefits.

APPENDIX B: LANGUAGE EVOLUTION

The final observation on Eiffel’s conceptual underpinning addresses change. Discussions of language design typically present a static view: how to build the ideal language? In practice, of course, no one reaches perfection, or even an earthly approximation thereof, the first time around. Unless the language is a one-off shot intended to produce papers rather than programs, it continues to evolve, as programmers complain about infelicitous features, designers realize that some of their decisions were not optimal, other languages come up with mechanisms worthy of attention, the evolution of the IT industry brings in new possibilities and needs, and more generally as new ideas come up. Languages such as Eiffel, C++ and Java have all been around for several decades, inevitably undergoing significant evolution in the process.

The key question, if a language has an installed base, is how to balance innovation with compatibility. The dilemma is not specific to programming languages but plagues the IT industry as a whole; it could be phrased cynically as how not to punish your users too harshly for their crime of trusting you too early.

The larger the installed base, and the more innovative the new solutions, the harder the decisions. (This matter was discussed at some length in the preface of the original language book [6].) Often, the dilemma is resolved in the direction of compatibility: do not harm the installed base. While the reasons are understandable and even commendable, the result after years of evolution is that modern languages continue to carry the accumulated sediments of many earlier eras.

It will come as no surprise that a number of the mechanisms in today's Eiffel, including some discussed in this article, were not present in the initial design. The most significant update, discussed in section 10, was to make Eiffel a void-safe language. But others, if less momentous, were also major innovations at the time of their introduction. While maintaining compatibility has been a constant concern in the evolution of the language, another strong principle is that it is important, in the spirit of the scientific method, to admit mistakes — meaning here design decisions that in retrospect turn out not to be ideal — and correct them, rather to cling to old ways of doing things. This approach has turned out to work provided changes apply the following principles:

- **Certainty.** Only introduce an incompatible change if you are absolutely certain that the new way is better than the old; not only better, but significantly better. (If not, “leave good enough alone” can be a reasonable guideline.) Otherwise, language evolution would go back and forth between contradictory decisions. Another way to state this principle is that compatibility remains the norm; incompatible change is the rare exception which must be fully justified.
- **Discussion.** Conduct a dialog with the user community and integrate its feedback. In our experience that feedback has sometimes led to abandoning proposed changes, but for the retained ones has often led to crucial improvements of the proposals.
- **Education.** Explain and justify again and again.
- **Transition.** Whenever possible, provide the old mechanism along with the new, at least for a generously determined transition periods (up to two years). In general it is preferable to cut that transition eventually, for fear of causing the sedimentation problem mentioned above, but give users time.
- **Conversion.** Offer tools to update existing software automatically, particularly for changes that do risk breaking existing code.

The application of these principles has allowed a bolder than usual approach to language evolution, helping programmers to move with the times while preserving the essential simplicity and consistency of the language, embodied by the design principles of which this article has presented some of the most significant applications.

BIBLIOGRAPHY

The list below only includes publications by the author and collaborators, since they are the only ones cited in the text. Those publications themselves rely on many contributions from

many authors over many decades. The bibliography of [10], for example, contains 300 references to contributions of other authors.

The references throughout the text to properties of languages such as Simula, Smalltalk, C++, Java and C# come from the standard documents and textbooks for these languages. References to Eiffel properties come from the ISO standard and the language site eiffel.org. Information about the EiffelStudio IDE is available at eiffel.com.

- [1] Bertrand Meyer: *Design by Contract*, Technical Report TR-EI-12/CO, Interactive Software Engineering Inc., 1986. Republished as [7].
- [2] Bertrand Meyer, *Genericity versus inheritance*, in ACM Conference on Object-Oriented Programming Systems, Languages and Applications (OOPSLA 1), Portland, Sept. 29 - Oct. 2, 1986, pages 391-405, se.ethz.ch/~meyer/publications/acm/geninh.pdf.
- [3] Bertrand Meyer, Jean-Marc Nerson and Masanobu Matsuo: *Safe and Reusable Programming using Eiffel*, in Proc. First European Software Engineering Conference (ESEC 87), Strasbourg (France), September 8-11, 1987, LNCS, Springer, 1987, se.ethz.ch/~meyer/publications/eiffel/eiffel_esec.pdf.
- [4] Bertrand Meyer: *Object-Oriented Software Construction*, Prentice Hall, 1988. (See also [10].)
- [5] Bertrand Meyer: *Disciplined Exceptions*, Technical Report TR-EI-13/DE, Interactive Software Engineering Inc., May 1988, se.ethz.ch/~meyer/publications/methodology/exceptions.pdf.
- [6] Bertrand Meyer: *Eiffel: The Language*, Prentice Hall, 1990.
- [7] Bertrand Meyer: *Design by Contract*, in *Advances in Object-Oriented Software Engineering*, eds. D. Mandrioli and B. Meyer, Prentice Hall, 1991, pp. 1-50.
- [8] Bertrand Meyer, *Systematic Concurrent Object-Oriented Programming*, in *Communications of the ACM*, vol. 36, no. 9, September 1993, pp. 56-80.
- [9] Bertrand Meyer: *The Many Faces of Inheritance: A Taxonomy of Taxonomy*, in *Computer* (IEEE), vol. 29, no. 5, May 1996, pp. 105-108, se.ethz.ch/~meyer/publications/computer/taxonomy.pdf. (Similar material in section 24.5 of [10].)
- [10] Bertrand Meyer: *Object-Oriented Software Construction*, second edition, Prentice Hall, 1997. (Considerably extended revision of [4].) A high-quality PDF version of the full text was recently made available at bertrandmeyer.com/OOSC2.
- [11] Paul Dubois, Mark Howard, Bertrand Meyer, Michael Schweitzer and Emmanuel Stapf: *From Calls to Agents*, in *Journal of Object-Oriented Programming* (JOOP), vol. 12, no 6, June 1999, se.ethz.ch/~meyer/publications/joop/agent.pdf.
- [12] Bertrand Meyer: *Attached Types and their Application to Three Open Problems of Object-Oriented Programming*, in ECOOP 2005 (Proceedings of European Conference on Object-Oriented Programming, Edinburgh, 25-29 July 2005), ed. Andrew Black, LNCS 3586, Springer, 2005, pages 1-32, se.ethz.ch/~meyer/publications/lncs/attached.pdf.

- [13] Ecma International: *ECMA standard: Eiffel Analysis, Design and Programming Language* (ed. Bertrand Meyer), approved as International Standard 367 by ECMA International, 21 June 2005; revised edition, December 2006, approved by the International Standards Organization as the ISO standard ISO/IEC 25436:2006. Text available at www.ecma-international.org/publications/standards/Ecma-367.htm.
- [14] Bertrand Meyer: *Software Architecture: Functional vs. Object-Oriented Design*, in *Beautiful Architecture*, eds. Diomidis Spinellis and Georgios Gousios, O'Reilly, 2009, pages 315-348, se.ethz.ch/~meyer/publications/functional/meyer_functional_oo.pdf.
- [15] Bertrand Meyer: *Touch of Class: Learning to Program Well Using Object Technology and Design by Contract*, Springer, 2009.
- [16] Bertrand Meyer, Ilinca Ciupa, Andreas Leitner, Arno Fiva, Yi Wei and Emmanuel Stempf: *Programs that Test Themselves*, in *IEEE Computer*, vol. 42, no. 9, pages 46-55, September 2009, se.ethz.ch/~meyer/publications/computer/test_themselves.pdf. See also AutoTest page at se.ethz.ch/research/autotest/.
- [17] Bertrand Meyer, Alexander Kogtenkov and Emmanuel Stempf: *Avoid a Void: The Eradication of Null Dereferencing*, in *Reflections on the Work of C.A.R. Hoare*, eds. C. B. Jones, A.W. Roscoe and K.R. Wood, Springer, 2010, pages 189-211, available at se.ethz.ch/~meyer/publications/lncs/attached.pdf.
- [18] Martin Nordio, Roman Mitin and Bertrand Meyer: *Advanced Hands-on Training for Distributed and Outsourced Software Engineering*, in ICSE 2010, Cape Town, IEEE Computer Society Press, 2010, se.ethz.ch/~meyer/publications/teaching/dose_icse.pdf. More publications and results on the DOSE distributed course and project at se.ethz.ch/research/dose.
- [19] Benjamin Morandi, Sebastian Nanz and Bertrand Meyer: *Can Asynchronous Exceptions Expire?*, in Proc. of 5th International Workshop on Exception Handling (WEH 2012), ICSE, Zurich, June 2012, IEEE Computer Press, 2012, available at se.ethz.ch/~meyer/publications/concurrency/exceptions_expire.pdf.
- [20] Julian Tschannen, Carlo A. Furia, Martin Nordio and Bertrand Meyer: *Automatic Verification of Advanced Object-Oriented Features: The AutoProof Approach*, in *Tools for Practical Software Verification; International Summer School, LASER 2011*, eds. Bertrand Meyer and Martin Nordio [350], LNCS 7682, Springer, December 2011, se.ethz.ch/~meyer/publications/proofs/autoproof.pdf. Web-based version of AutoProof and extensive further publications at autoproof.sit.org.
- [21] Bertrand Meyer: *Theory of Programs*, in Proc. LASER summer school on Software 2013-2014 2007/2008, eds. B. Meyer and M. Nordio, Springer LNCS 8987, 2015, se.ethz.ch/~meyer/publications/proofs/top.pdf.
- [22] Bertrand Meyer, Alisa Arkadova and Alexander Kogtenkov: *The Concept of Class Invariant in Object-Oriented Programming*, submitted for publication, preprint available at arxiv.org/abs/2109.06557, August 2022.

