

*Bertrand Meyer*

# SYSTEMATIC CONCURRENT OBJECT-ORIENTED PROGRAMMING

**J**udging by the looks of the two parties, the marriage between concurrent computation and object-oriented programming—a union much desired by practitioners in such fields as telecommunications, high-performance computing, banking and operating systems—appears easy enough to arrange. This appearance, however, is deceptive: the problem is a hard one.

This article points the way toward a possible solution. The precise problem examined here is restricted to:

*What is the simplest, smallest and most convincing extension to the method of systematic object-oriented software construction that can address the needs of concurrent and distributed computing as well as those of sequential computation?*

The article does not claim to discuss concurrency and distribution in a general and unbiased way. Rather, it takes the object-oriented paradigm as a given (on the basis of its contributions to the production of quality software) and investigates how best to adapt it so it covers both concurrent and sequential applications.

The word “systematic” as used in the title of the article will provide strong guidance to our search for this minimal extension. We are interested in an approach that makes it possible to reason about software systems in a precise way, by extending to the concurrent case the systematic techniques, known as ‘design by contract’, which can be applied to the systematic (although not necessarily fully formal) development of sequential object-oriented software.

A word of warning as to the actual ambition of the model discussed here. No claim is made that the work as reported is final, and a number of possible criticisms are discussed in the final section. I do think, however, that in its discussion

of how the mutual attraction between object orientation and concurrency can be turned into a reasonably happy marriage, this article raises a number of questions. Although crucial for both theoretical understanding and practical implementation of concurrent object-oriented computation, these questions have not been addressed (or in some cases even mentioned) by previous work on the subject, and will have to be resolved before a solution can be widely accepted and applied.

### **Similarities and Contradictions**

The property which initially suggests an easy match between the ideas of concurrency and object orientation is the remarkable similarity between the basic constructs of both object orientation and concurrency. It is hard to miss the analogies between objects and processes, or more accurately between the underlying abstractions: classes and process types. Both categories of constructs support:

- Local variables (attributes of a class, variables of a process or process type)
- Persistent data, keeping their value between successive activations.
- Encapsulated behavior (a single cycle for a process; any number of routines for a class)
- Heavy restrictions on how modules can exchange information
- A communication mechanism usually based on some form of message passing.

It is not surprising, then, that researchers have tried to unite the two areas. But although existing designs (for surveys see [1, 19] as well as a recent thesis [17]) have introduced many productive ideas, it is fair to state that so far none has succeeded in providing a widely accepted mechanism for concurrent object-oriented programming. The primary reason is probably the undue complexity of most of the proposed solutions, which tend to add the

full power of an independent concurrency facility to an object-oriented language, or vice versa. Another reason is that little of the existing literature devotes much attention to correctness issues, and more generally to the possibility of systematic reasoning about concurrent object-oriented programs.

The mechanism described here attempts to remedy these limitations at least in part. The syntactical extension, which it brings to an object-oriented language (Eiffel) is the smallest feasible: one new keyword. A new library class with procedures for setting two options is also provided to adjust the behavior in special cases. The mechanism makes the greatest possible use of existing object-oriented facilities—classes, inheritance, assertions, argument passing, deferred classes—to cover such concepts of concurrent computation as exclusive access, processes and synchronization.

### Criteria

The design of a proper concurrent object-oriented mechanism must satisfy a number of criteria. The combination of these criteria, as defined in the following subsections, places rather strict constraints on the possible concurrency mechanism—to the extent that one might wonder at first whether any solution is possible.

### Minimality of Mechanism

Object-oriented software construction is a rich and powerful paradigm, which, as noted previously, would intuitively seem to be ready for supporting concurrency. It is essential, then, to aim for the smallest possible extension. Minimalism here is not just a question of good language design. If the concurrent extension is not minimal, some concurrency constructs will be redundant with the object-oriented constructs, or conflict with them, making the programmer's task difficult or impossible. To avoid such a situation, we must find the smallest syntactic and semantic epsilon that will give concurrent execution capabilities to our object-oriented programs.

### Full Use of Inheritance and Other Object-Oriented Techniques

It would be unacceptable to have a concurrent object-oriented mechanism that does not take advantage of all object-oriented techniques, in particular inheritance.

One of the most interesting contributions of object-oriented technology is its ability to support many different patterns of computation. Once a useful type of behavior is identified, it can be encapsulated in a deferred class (see "Deferred Classes and Features" sidebar) from which any class that uses that behavior will inherit. Such a class is similar to a process type in concurrent programming and its general form may be expressed as:

*deferred class* *PROCESS* *feature*

*live is*

— General structure with variants.

**do**

*from setup until over loop*

*step*

**end;**

*finalize*

**end;**

*feature* {*NONE*}

*setup is deferred end;*

*over: BOOLEAN is deferred end;*

*step is deferred end;*

*finalize is deferred end*

**end**

(The clause *feature* {*NONE*} introduces features that are not available to clients, being meant for internal use only. Features declared in a clause beginning with just *feature* without qualification are available for calls by any client and are said to be exported.)

Since routines *initialize*, *over*, *step* and *finalize* are deferred, descendants of *PROCESS* may provide effective implementations of these routines, corresponding to individual variants. There may be as many such classes as variants are needed. The overall behavior, however, is the same for all variants; it is determined by the effective routine *live*. In the example the structure of *live* involves a loop, but the same ideas are applicable to any other structure. Also, there will often be more than one routine such as *live*, covering several patterns of behavior that are known at the level of the deferred class.

This technique and many others (such as the use of polymorphism, static typing and dynamic binding to obtain flexible and safe software architectures) are essential to the object-oriented approach, and are potentially beneficial to concurrent programming as well. This has been recognized for example by recent changes to the parallel object-oriented language (POOL) [3], early versions of which did not support inheritance. The model presented here makes full use of all object-oriented techniques.

### Compatibility with Design by Contract

Another central idea of what may be called the Eiffel methodology of object-oriented software construction is the notion of Design by Contract [11, 15]. According to this view, the design of a reliable software system should use a number of components (classes) that communicate with one another on the basis of precise definitions of obligations and benefits: contracts. The obligations and benefits are documented in the text of the software itself, where they appear in the form of **assertions**.

A contract governs the relation between client objects requesting services (by calling a feature) and supplier objects providing services. The services are expressed by the routines of the supplier's class. For each routine, an assertion known as the precondition expresses the input requirements; it binds the clients and protects the supplier. Another assertion, the postcondition, expresses the properties ensured by any call to the routine; it binds the supplier and guarantees a certain result to the clients. Beyond the precondition and postcondition of its individual routines, a class is also characterized by its invariant, which applies to all exported routines, being in effect added to both their preconditions and postconditions.

A convenient example that will aid in the search for the proper concurrent extension is a class describing queues (first-in-first-out container structures) of bounded size. The corresponding concurrent notion, developed in the following program



text, will describe bounded buffers. Following is a sketch of the class, with the implementation details omitted. Preconditions are introduced by *require*, postconditions by *ensure*.

```

class BOUNDED_QUEUE[G] feature
  empty: BOOLEAN is
    --Is there no accessible element?
  do Result: = ... end;
  full: BOOLEAN is
    --Is there no room left
    --for more elements?
  do Result: = ... end;
  put (x:G) is
    -- Add x as newest element.
  require
    not full
  do
    ...
  ensure
    not empty
  end;
  remove is
    -- Remove oldest element.
  require
    not empty
  do
    ...
  ensure
    not full
  end;
  item is
    --Oldest element
    --not yet consumed
  require
    not empty
  do Result: = ... end
feature {NONE}
  ... Secret features (used for
  the implementation) ...
invariant
  ... See below ...
end

```

Through its assertions, class *BOUNDED\_QUEUE* describes a number of contracts governing the use of its features by clients. For example, the contract for *put (x)* is expressed in Table 1.

An important property of preconditions, which we will have to examine again in the concurrent context, is that they express the sole obligations the client has to meet. This may be called the **no hidden clauses** rule of Design by Contract. For example, a call to *put* is guaranteed to succeed if the client has ensured the precondition, perhaps by writing it (with *q* of type *BOUNDED\_QUEUE [X]* and *a* of type *X* for some type *X*) as:

```

-- /1/
if not q.full then
  q.put (a)
end
or by prefacing it with a call to
remove, whose postcondition implies
the precondition of put (assuming
that the call to remove is itself correct,
that is to say, the precondition not
q.empty initially holds):
-- /2/
q.remove; q.put (x)

```

Another important use of assertions is the class invariant, which characterizes the consistency of a class. Assume for example that an implementation of *BOUNDED\_QUEUE* relies on the well-known technique of using an array of *array\_count* elements, managed in a circular way as illustrated in Figure 1. The details of the implementation are left to the reader; the following invariant will express its fundamental properties, in particular the need to keep one position unused in order to be able to distinguish between the empty and full cases:

```

array_count = capacity - 1;
abs (next - oldest) < capacity;
0 <= oldest; oldest <= capacity;
0 <= next; next <= capacity;

```

Any call to an exported routine may assume that the invariant is initially satisfied (i.e., it may expect to find the object in a consistent state); but it must also restore the invariant, in addition to ensuring its postcondition, on exit (i.e., it must leave the object in a consistent state). We may thus picture the life of the object as a sequence of transitions between consistent states. The shaded squares in Figure 1 represent states satisfying the invariant; the transitions represent calls to exported features, executed by clients. These states are the only ones in which the object is accessible to a new client.

The notion of invariant is essential to the design of a proper exception mechanism. In the Design by Contract approach, an exception occurs when a routine is unable to fulfill its contract through the initially planned strategy. The routine may try again through a Retry instruc-

tion, usually after attempting to correct the source of trouble. Otherwise, the routine will execute to the end its (implicit or explicit) Rescue clause, fail, and cause an exception in the client, which will then be faced with the same choice—Retry or failure. In case of failure the Rescue clause is not required to fulfill the contract as expressed by the postcondition (this would be success, not failure!), but it must restore the invariant, leaving the object in a consistent state for any later attempt at Retry.

Another reason why the notion of invariant is important for concurrency is that it enables us to put the notion of “express message” [19] into a proper perspective. Express messages, as proposed, allow interrupting the execution of a routine on a certain supplier object, on behalf of some client, when a call comes in from another client which is deemed more important. An incoming express message will then get served right away; only then will the original client’s execution resume. Defined in this way, however, such a mechanism conflicts with correctness requirements: if we allow executions to be interrupted, we cannot guarantee they will preserve the invariant. Producing an object which does not satisfy the invariant of its own class is probably the worst disaster that may occur during the execution of an object-oriented program. (Another well-known source of such a situation is static binding.)

As a consequence of this analysis, the mechanism described does not support express messages in the sense of [19]. It will, however, allow a VIP client to interrupt an earlier client, causing an exception in that client. Such a facility is compatible with Design by Contract, since the invariant will be restored as a result of the exception handling.

The relevance of invariants to concurrent programming was expressed by Hailpern [8], who uses a concept of “monitor invariant, [which] must be true when no process owns the monitor.” Starting from an object-oriented basis, we do not need a special notion; class invariants will provide a way to characterize the invariant properties associated with

monitors, processes and other constructs of nonobject-oriented concurrent programming.

**Provability**

Design by Contract provides the starting point for a potential formal approach to object-oriented computation. Assuming proof rules were available for the inner details of an object-oriented language, we could use a general proof rule for calls to prove entire object-oriented systems. The proof rule for calls is fundamental because at the heart of object-oriented computation lie operations of the form

$t.f(\dots, a, \dots)$

which call a feature  $f$ , possibly with arguments such as  $a$ , on a target  $t$  representing an object. The basic proof rule may be informally stated as follows:

If we can prove that the body of  $f$ , started in a state satisfying the precondition of  $f$ , terminates in a state satisfying the postcondition, then we can deduce the same property for the preceding call, with actual arguments such as  $a$  substituted for the corresponding formal arguments, and every clause of the form *some\_boolean\_property* in the assertions replaced by the corresponding property on  $t$ , of the form *t.some\_boolean\_property*.

For example, if we are able to prove that the actual implementation of *put* in class *BOUNDED\_QUEUE*, assuming *not full* initially, produces a state satisfying *not empty*, then for any queue  $q$  and element  $a$  the rule allows us to deduce

$\{\text{not } q.\text{full}\} q.\text{put}(x) \{\text{not } q.\text{empty}\}$

where the assertions in braces describe the input and output assumptions respectively.

The proof rule, an adaptation to the object-oriented form of computation of Hoare's inference rule for procedures [9, 13] may be expressed as shown in Figure 3, where *INV* is the class invariant, *Pre*( $f$ ) is the set of precondition clauses of  $f$  and *Post*( $f$ ) the set of its postcondition clauses. Recall that an assertion is the conjunction of a set of clauses, of the form

## Basic Terminology

**O**bject-oriented development is based on the notion of class. A class describes a set of potential run-time objects, defining them entirely by the applicable operations, or features. Any object created at run time from the pattern defined by the class is called an instance of the class.

The features of a class are of two kinds: routines, which describe computations applicable to instances of the class; and attributes, which describe data fields associated with instances of the class. For example, a class *CAR* may have attributes *weight* and *speed*, and routines *start*, *stop*, *accelerate*, *average\_speed*. A routine is either a procedure, which may change the object to which it is applied but does not directly return a result, or a function, which computes some information about the object and returns that information as a result. In both cases, the routine may have one or more arguments. In the example *start*, *stop*, *accelerate* will be procedures, and *average\_speed* will be a function, computing the average speed since a certain starting time; the function will have one argument, representing that time.

The basic computational mechanism is feature call, which applies a feature to a certain object, known through a name in the software text, or entity. All entities are declared. For example, with the declarations

$c: \text{CAR}; x: \text{REAL}$

a call to feature *average\_speed*, used here in an assignment, could have the form

$x := c.\text{average\_speed}(0)$

which assigns to  $x$  the value of the average speed of the object attached to  $c$  (an instance of *CAR*) since time 0. The notation used here for calls is dot notation, which includes the following components: an entity representing the target object of the call; a dot (.); a feature name; a list of actual arguments, if any, in parentheses.

Other examples of calls are

$c.\text{accelerate}(20);$

$x := c.\text{speed}$

A call such as any of the preceding ones is executed as part of the text of some routine—itself executed on a certain object *C\_OBJ*, the client object of the call. The object to which a call applies (the object attached to  $c$  in the examples) is called the supplier object.

A feature such as *speed* may be implemented as either an attribute or a function. Principles of information hiding and uniform access imply that this choice of implementation should make no difference to clients. The notation for feature calls, and the standard form for class documentation, known as the short form of a class, are indeed the same in both cases.

**Table 1.** The contract for *put(x)*

	Obligations	Benefits
Client	Queue not full	Queue not empty; $x$ inserted
Supplier	Must insert $x$	Queue not full (some space left)

$clause_1; \dots clause_n$

The large  $\wedge$  signs in Figure 3 indicate conjunction of all the given clauses. The actual arguments of  $f$  have not been explicitly included in the call, but the primed expressions such as  $t.q'$  indicate substitution of the actual arguments of the call for the formal arguments of  $f$ . The rule is stated in Figure 3 in the form which does not support proofs of recursive routines. Adding such support, however, does not affect the present discussion. For details of handling recursion, see [9] or [13].

The reason for considering the assertion clauses separately and then “anding” them is that this form prepares the rule’s adaptation to “separate” calls in the concurrent extension. Another property of the nonconcurrent rule which is of interest as preparation for that extension is the presence of the invariant  $INV$  in the proof of the routine body (above the line), with no directly visible benefit for the proof of the call (below the line). More assertions with that property will appear in the concurrent rule.

The provability requirement and the criterion of compatibility with Design by Contract have a major immediate consequence for concurrent computation. The bad news about proving the correctness of a class is that for every exported routine  $f$  you must prove a property of the form

$$\{INV \wedge allpre(f)\} \text{ body } (f) \{INV \wedge allpost(f)\}$$

The good news, however, is that you only have as many properties of this kind to prove as the class has exported routines. (Attributes do not require any specific proof, but they may be involved in the invariant.) This provides a very strong guideline for choosing the granularity of exclusive object access in concurrent computation. If we allowed a supplier object to accept a new call from a client while a call from another client is in progress, the “good news” would not hold any more; the property depicted in Figure 2, where the shaded states (before and after execution of exported features) are the

only ones in which an object is accessible from the outside, would also fail to be satisfied. Consequently, a proof of correctness would have to take into account all possible interleavings, producing a combinatorial explosion of properties to prove. This would in fact remove any hope of producing realistic proofs, even informal ones, or just of being able to reason about software texts in a systematic fashion.

As a result, any client accessing an object through a feature must be guaranteed exclusive access to the object throughout the duration of the call. The smallest permissible level of granularity for exclusive access to an object is the execution of a call to an exported feature.

This observation generalizes the comments made about express messages. Nothing prevents us from interrupting an ongoing call as long as this is done properly: the previous client must receive an exception, so it will be forced either to fail or to take corrective action. A facility supporting such a scheme, which we may call a *duel*, will be part of the mechanism.

### Support for Command-Query Distinction

An important part of the object-oriented method which we should try to preserve in concurrent programming is the necessity to maintain, whenever possible, a strict distinction between two kinds of features: commands and queries. A query, implemented as an attribute or function, is a feature that returns some information about an object. A command, implemented as a procedure, is a feature which may modify the state of an object.

The command-query distinction directs programmers to refrain from using a programming style that has become popular in recent years, especially in connection with the spread of C: functions that produce visible side effects. This practice endangers referential transparency (substitutivity of equal for equals).

The attempt to maintain a strict command-query distinction explains why class *BOUNDED\_QUEUE* introduced previously has two separate features, *remove* and *item*. Function

*item*, a query, accesses the oldest element not yet removed, but does not remove it; successive calls to this feature will return the same value. Procedure *remove*, a command, removes the oldest element. A designer who does not apply the command-query distinction might be tempted to replace these routines by a side-effect-producing function *get*, which both removes an element and returns its value. The presence of several concurrent clients accessing the same supplier will initially make it difficult to ensure the command-query distinction; but this rule will in fact suggest some of the important properties of the concurrent mechanism.

### Applicability to Many Forms of Concurrency

A general criterion for the design of a concurrent mechanism is that it should support many different forms of concurrency: shared memory, multitasking, network programming, distributed processing, real-time applications.

With such a broad set of application areas, a language mechanism cannot be expected to provide all the answers. But it should lend itself to adaptation to all the intended forms of concurrency. Mechanisms external to the language itself will make it possible to describe a mapping between the physically available processing units and the abstract threads of control needed by the software text.

### Support for Coroutine Programming

An interesting form of computation is the use of coroutines. Present in the first object-oriented language, Simula 67 [5, 18], coroutines (see Figure 4) are sequential program units that communicate on an equal basis (rather than the master-subordinate relationship of a calling unit and a routine). Between its successive reactivations, a coroutine retains the value of its data, and the location of its active program counter—as opposed to a routine, which is always restarted at the beginning. Coroutines restart each other explicitly through operations called *resume* in Simula.

Coroutine computation may be viewed as an extreme form of concurrency in which only one processing unit is available. Any general-purpose concurrency mechanism should reduce to a coroutine mechanism in this case. An important consequence of this property will be to ease the transition from a simulation to the actual system. To write a computer program simulating a real-world system that involves concurrent activities (as most do), it is often convenient to use a coroutine-based scheme, such as Simula's discrete event simulation facilities. If the same language mechanism is used for coroutines and for actual concurrency, distinguished only by the number of available processing units, writing the simulation helps write the actual system; conversely, if the system has already been written, it is easier to write a simulation.

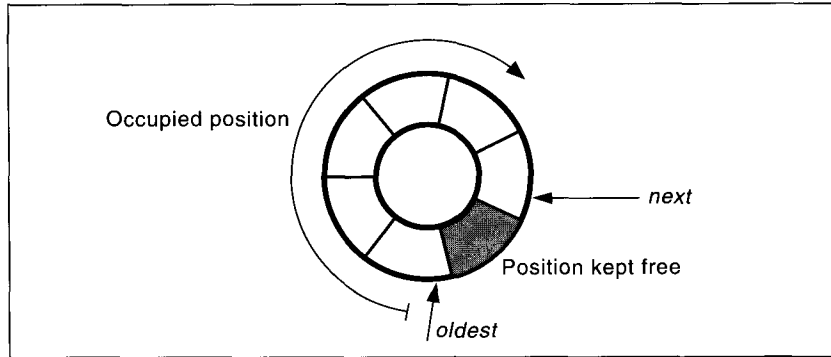
**Adaptability through Libraries**

Many concurrency mechanisms have been proposed over the years, from semaphores and conditional critical regions to Petri nets, monitors and CSP. Each has its partisans, and each may provide the best approach for a certain problem area. It is important that the proposed mechanism should support at least some of these mechanisms. More precisely, the solution must be general enough to allow us to *program* various concurrency constructs in terms of the basic mechanism we will have obtained.

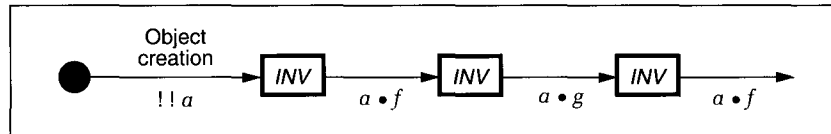
One of the most important aspects of the object-oriented method is that it supports the construction of libraries for widely used schemes. The library construction facilities (classes, assertions, constrained and unconstrained genericity, multiple inheritance, deferred classes and others) should allow us to express many concurrency mechanisms in the form of library components.

**Support for Reuse of Nonconcurrent Software**

A criterion of the desirable (rather than essential) category is the ability to reuse existing, nonconcurrent software, especially libraries of reusable software components. This may not always be achievable, since con-



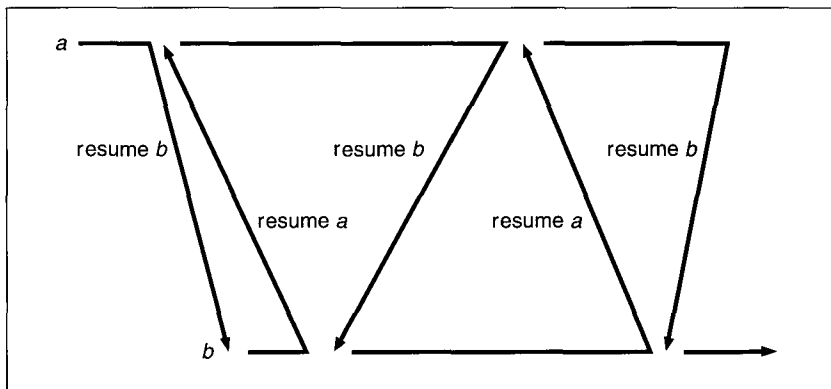
**Figure 1.**



**Figure 2.**

$$\frac{\{ INV \wedge p \in Pre(f) p \} \quad body(f) \quad \{ INV \wedge q \in Post(f) q \}}{\{ p \in Pre(f) t.p' \} \quad t.f \quad \{ q \in Post(f) t.q' \}}$$

**Figure 3.**



**Figure 4.**

currency places new demands on the software structure; for example we must specify the scope of exclusive access segments. But even when existing software cannot be reused exactly "as is" the work involved in making it applicable to concurrent development should be reasonable, involving for example the writing of simple "wrapper" classes encapsulating existing sequential classes. This criterion may be viewed as another form of the minimality requirement, applied here not to the job of the lan-

guage designer but to that of software developers.

**Support for Deadlock Avoidance**

A specific but important question is how the sought mechanism will help solve a difficult problem of concurrent programming: avoiding deadlock. A solution which guarantees deadlock avoidance in all cases would probably be too limited. For example, the requirement to support many different forms of concurrency suggests that our mechanism should



allow writing a class describing semaphores. But as soon as we have semaphores we have the possibility of deadlock. A more realistic requirement, then, is that the mechanism should make it possible to avoid deadlock by observing certain statically enforceable restrictions—which must remain reasonable. Making this possible was an important concern in the design of the mechanism described here. So far, however, the question of deadlock avoidance has not been pursued very far, and no result on this topic is included in the remainder of this article.

### Designing a Solution

Let us now examine how the preceding requirements determine a concurrency mechanism that remains compatible with the letter and spirit of object-oriented software construction.

### No Active-Passive Distinction

The first concern is whether we need a special notion of “active object” and process. Most current proposals for object-oriented concurrency mechanisms include such notions. Two examples among many are the article by Caromel on a concurrent Eiffel extension [7], which includes: “*The first choice faced when designing a concurrent language is process genesis. What language construct and concept permit process definition?*” and a description of the POOL language by America [3], which in an introductory paragraph, states “*each object also has a body, a local process that starts as soon as the object is created and runs in parallel with all the other objects in the system.*” The present discussion differs from such approaches by refusing to include an explicit notion of process or of active object.

The reasoning behind proposals supporting active objects is as follows: In the usual, sequential form of object-oriented computation, objects are “sitting there” waiting for requests addressed to them. Such requests are feature calls, of the form *t.f (...)*, meaning “apply the feature associated with *f* to the object attached to *t*, with the given arguments if any.”

Such a feature call is executed on behalf of a requesting object, the client, and is addressed to a target object, the supplier. In this scheme a supplier object is just a passive repository of features, ready to be triggered at the client’s behest. The client is in the same position with respect to its own clients. In fact, every object except the one created first (the root) is used as supplier of other objects. Execution of a system is started by creating a root object and applying a feature to it, firing off a chain of feature calls.

For concurrent computation, the argument goes, we need active objects with their own computational power—their own agenda. Such objects will correspond to the processes found in most nonobject-oriented models of concurrent computation. Thus we will have objects of two kinds: passive objects, as in sequential object-oriented computation; and active objects. Such a distinction was in fact previewed in Simula [5, 18], where a class could, in addition to its features, include a body—a sequence of instructions describing a behavior associated with the instances of the class. This made instances of such classes the forerunners to active objects. This facility served in particular for the use of Simula classes as coroutines, and for discrete-event simulation.

On closer examination, however, the passive-active distinction appears unjustified and in fact harmful. It is useless to associate a special algorithm with an object when, through the routines of its class, it can have as many as the class author desires.

A typical example of an active object would be the process associated with a printer. The process would describe the algorithm that governs the printer’s life: initialize; repeatedly process user jobs; shutdown. In object-oriented programming, however, this will just be one of the features associated with the printer. The corresponding class, using techniques of behavior encapsulation illustrated previously with class *PRINTER*, could appear in sketched form as:

```
class PRINTER inherit
```

```
PROCESS
```

```
  rename over as off_line, finalize as stop
end
feature
  off_line: BOOLEAN;
  -- over is effected as an attribute,
  -- under the name off_line.
  stop is
    -- Go off-line.
    do off_line := true end;
feature {NONE}
  step is
    -- Execute individual actions
    -- of an iteration step.
    do
      start_job; process_job; finish_job
    end;
  setup is do ... end;
  start_job is do ... end;
  process_job is do ... end;
  finish_job is do ... end
end
```

The procedure *live* describes the process associated with the printer. Were a special notion of active object to be added, this procedure would become the body of class *PRINTER*. But why should we settle for one procedure when we can have as many as we want? By sticking to the scheme shown we retain the ability to have more features than just available to clients, for example *shutoff*, *print\_diagnostics*, *prepare\_for\_maintenance* and many others. This is the starting point of object orientation: the realization that you can almost always do more than one thing with an object. This refusal to consider any one feature as “the main operation” on a class is one of the main tenets of the object-oriented method, and yields some of its key advantages in terms of software extendibility and reusability. By grafting onto object-oriented programming an independent concept of process, in the form of active objects, we would lose this essential property and gain nothing new.

The addition of processes as an independent concept would cause other problems as well. Limiting an active object’s available scripts to just one raises the question of how active objects (processes) request services from each other. The fundamental object-oriented computation mechanism, feature call, would not work any more without some special synchronization mechanism: in the exe-

cution of *t.f.* (...), if the object T\_OBJ attached to *t* is active, it will be busy with its own computation and not ready to handle the call unless special measures are taken. To solve the problem, we would have to add CSP- or Ada-like mechanisms, resulting in a full new language layer. The complexity of the result would be unacceptable.

### Processors

If processes are not the appropriate new basic concept, we must find a better way of expressing the fundamental difference between sequential and concurrent object-oriented programming. The following simple observation may serve as a basis for an answer. Computation, as described by Figure 5, involves three elements: certain processors apply certain actions to certain objects. Object-oriented programming has been quite effective at capturing the last two aspects, by attaching the description of actions (routines) to the description of objects (classes). In ordinary sequential computation, there is only one processor, which is why it tends to remain implicit.

With concurrent computation, however, we have two or more processors, and so we need to make processors explicit. This will be the major result of adding concurrency to the framework of sequential object-oriented computation. For every object T\_OBJ, there must be a processor responsible for executing all calls having T\_OBJ as target. This processor will be said to handle T\_OBJ; the handler of every object is determined when the object is created.

A processor is a separate thread of control capable of supporting the sequential execution of operations on one or more objects. It is important to note that this notion of processor is virtual, not physical. A processor may represent a physical computational device (CPU), for example a computer on a network, but this is not necessarily the case: a processor may just as well be time-shared with other processors on a computer. For example, a Unix task or a lightweight process may be used as processors. The difference between vir-

tual processors and physical CPUs was clearly expressed by Lieberman [10]:

*The number of [processors] need not be bounded in advance, and if there are too many [processors] for the number of real physical [CPUs] you have on your computer system, they are automatically time-shared. Thus the user can pretend that processor resources are practically infinite.*

(Lieberman's terminology and concurrency model are different from those of this article, hence the bracketed words.) To avoid any confusion, the present discussion will employ the term "processor" only in this sense of virtual thread of control; "CPU" is used to refer to an actual computational device.

Because processors are virtual, not physical, the mechanism described here may be used to support distributed processing, in which the processors are physically distinct computers, as well as multiprogramming, in which the processors are supported by operating system processes. An extreme case is the availability of just one CPU: then the mechanism may be used for coroutine programming. The virtual nature of processors has another consequence: although the mechanism as described here does not permit reassigning an object to a new processor, nothing prevents an implementation from offering a way to reassign a processor to a new CPU, which in practice achieves the desired effect.

A process, in this scheme, becomes a trivial notion—an instance of a "process class." A process class has a distinguished procedure, the only one of interest for clients. For the common case in which the distinguished procedure is a loop describing a behavior to be repeated until termination, we can write any process class as an effective descendant of the class *PROCESS* as introduced previously. There is no need for new language constructs.

### Contracts and Concurrency

Moving from sequential to concurrent object-oriented programming, then, will imply making the processors explicit in some way. To find the

appropriate method for doing this, we must understand what having more than one processor implies for the basic scheme of the object-oriented method: design by contract.

As pointed out, a fully defined contract implies a no hidden clause property: clients that "play by the rules," observing the precondition of a call, are guaranteed to obtain the result, as expressed by the postcondition. Unfortunately, this crucial property will not hold in a concurrent context without a change in the semantics.

Consider again the bounded queue example, with the calls to *put* introduced previously as /1/ and /2/. In either case, the condition *not q.empty* will hold prior to the call to *put*. This results from the test in the first case and from the postcondition of *remove* in the second. Another way to express this property is to state that in sequential programming assertions express correctness conditions. If, prior to the call to *put*, the condition *not q.full* is violated, this simply indicates a **bug** in the client containing this call. The only reasonable response is to correct the bug. The presence of a mechanism to check assertions at run time, as exists in the Eiffel environment, does not change anything in this regard: run-time assertion monitoring is a precious aid to quality assurance, in particular testing and debugging; but it is not a technique for treating certain special but expected conditions in a particular way.

Unfortunately, these properties break down for concurrent situations. Assume the processor handling the client is different from the processor handling the object attached to *q*. Then in /1/, as any student having taken an introductory course in concurrent computation knows, one or more processors can sneak in between the test and the call to *put* and call features such as *put* on *q* on behalf of other clients, making the test totally useless. In the same way, between the two instructions of /2/, other processors can invalidate the result (*not q.full*) achieved by the call to *remove*.

In other words, merely ensuring the precondition before a call does



not guarantee correctness any more if the client and the supplier are handled by different processors. This means the sequential contract model does not hold in its original form for concurrent computation.

### Separate Entities

A consequence of the previous discussion is that the sequential semantics of assertions will have to be adapted for cases in which the client and the supplier are handled by different processors. Before we start exploring what the new semantics should be in such cases, and regardless of the eventual answer, we must address an absolute requirement: ensuring that the effect of an operation is clear from the software text.

Any concurrent semantics we choose will imply that  $t.f. (...)$  may have a different effect depending on whether the object attached to  $t$  is handled by the same processor as the client or by a different one. It would be unacceptable to hide this important difference from the reader of the software text.

As a result, we should have a special notation to declare that a certain entity (see "Entities, Types and Values" sidebar for definition of this term) denotes objects that will be handled by a different processor. The syntactic extension is immediate. Instead of the usual declaration

```
--/4/  
x: SOME_TYPE
```

we will declare an entity as

```
--/5/  
y: separate SOME_TYPE
```

to express that  $y$  may become attached to objects handled by a different processor.

With such a declaration, the creation instruction

```
!! y.make(...)
```

has an extra effect. In all cases (*separate* declaration or not) the instruction creates a new object, initializes it to language-defined default values, and applies the creation procedure *make* with the given arguments as a way to override default initialization as needed. If  $y$  is declared as *separate*, the instruction will, in addition,

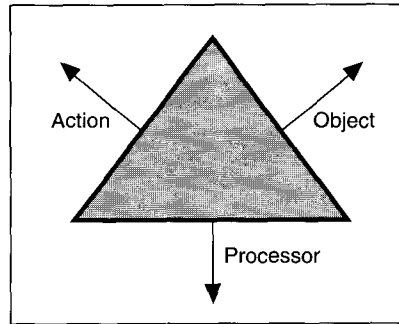


Figure 5.

## Deferred Classes and Features

**A** class is deferred (or abstract) if it has one or more features declared as deferred, that is to say, specified but not implemented. A nondeferred class or feature is said to be effective. A descendant of a deferred class (that is to say, a class which inherits from it directly or indirectly) is effective if it implements (or effects) all deferred features by providing effective forms.

A deferred class describes a general abstraction which may have many different realizations. It may also serve to capture a common set of behaviors by using an effective (nondeferred) routine, such as procedure *live* in class *PROCESS*, which calls deferred ones (*setup*, *over*, *step*). Descendant classes will retain the common behavior and provide the specifics through effective versions of the originally deferred routines.

assign a new processor to handle the newly created object.

The semantics of the language should not specify how the new processor is determined, although it is possible to envision library mechanisms that will give programmers some control in this respect. Another way to obtain a separate object is

through a call to some function of the form

*new\_object: separate T is ...*

Such a function may be declared as *external*, allowing some control from outside the language proper. This makes it possible to assign objects to different threads of control (for example, to various computers on a network or multimicroprocessor system) through some external mechanism, without recompilation of the software. The details of such mechanisms fall beyond the present discussion, but it is important to make sure they are possible.

A possible objection should be examined here. One might be concerned that by declaring  $t$  as *separate* we are giving out too much implementation detail. Should we not just be able to write  $t.f. (...)$  without worrying where the call is executed? But this objection is not justified. What should be hidden from the client in all but special cases is the precise knowledge of *which* processor handles a call. The mechanism described here achieves this objective. But whether the processor in question is the same one as the processor handling the client, or another, is highly relevant information, since the semantics are different. In addition to the change in the interpretation of assertions, there is an even more fundamental difference: execution by the same processor is blocking—the client cannot proceed until execution of  $f$  has been completed—whereas execution by a different processor should not prohibit the client from continuing its own execution until it actually needs the results of the call, if ever. (This property will give rise to the policy of lazy wait, explained later.) So it is indeed necessary to state clearly whether the processor is the same.

Certain classes are meant to be used only as types of separate entities. As a notational convenience, it will be permitted to declare such classes as

```
separate class CLASS_NAME ... the rest as usual ...
```

meaning that an entity declared of type *CLASS\_NAME* will always be

separate. As with the two other existing mechanisms to characterize a class, *deferred* and *expanded*, the separate or nonseparate status of a class is not transmitted through inheritance: a separate class may inherit from a nonseparate one (the common case) and conversely. For example, we may describe bounded buffers through a class declaration consisting simply of

```
separate class
  BOUNDED_BUFFER [G]
inherit
  BOUNDED_QUEUE [G]
end
```

The three specific properties of classes are exclusive: the syntax enables the keyword *class* to be preceded by at most one of *deferred*, *expanded* and *separate*.

The *separate* declaration for entities and classes will turn out to be the only syntactic extension needed by the mechanism described here, although we have yet to examine its full semantic consequences.

It will be convenient to use the term “separate” in various contexts. A class declared as *separate class* is a separate class. An entity declared as *separate SOME\_TYPE*, or as *SOME\_SEPARATE\_TYPE* if *SOME\_SEPARATE\_TYPE* is based on a separate class, is a separate entity. A call  $t.f(\dots)$  is separate if its target  $t$  is separate. An expression  $y$  is separate if it involves at least one separate entity. If  $t$  is a separate entity, and at some point during the execution of the system has a nonvoid value, the attached object, which is handled by a different processor, is said to be a separate object.

For consistency, a validity rule will require that in any assignment of the form  $x := y$ , if the source  $y$  is separate, the target  $x$  must also be separate. Otherwise we would be able to cheat by manipulating a separate object through a nonseparate entity  $x$ , obtaining the wrong semantics. It will be permitted, however, to assign a nonseparate value to a separate entity; calls on the corresponding object will then have the semantics of separate calls, which is harmless.

The same rule applies to actual-formal argument association, as in

the call  $t.f(y)$ , where the corresponding formal argument declared for routine  $f$  is  $x$ . (The semantics of assignment is the same as that of argument passing; the term **attachment** covers both operations.)

### Semantics of Assertions

With the notion of separate entity in place to ensure that every call is clearly identified as intended for the same processor or for another, we need to return to the important question of what assertions mean for separate calls. We have seen that preconditions, or at least those clauses of a precondition that involve separate entities, cannot be taken as correctness conditions. But the requirements expressed by a precondition are still needed for the routine to do its job properly. For example, we cannot write a correct version of *put* without some guarantee that the queue is no-full on entry.

What then should the semantics of a precondition such as *not full* be if the client and supplier are handled by different processors? Only two answers seem to make sense:

- Failure to meet the precondition may mean failure of the call. As noted previously, the contract model provides for such a case: a call may fail if the routine is unable to fulfill its contract; this causes an exception in the client. The client may recover from the exception through a Rescue clause which takes any needed corrective actions and, through a Retry instruction, tries another (or the same) strategy. If no Rescue clause is present, the client itself fails and passes the exception to its own client. This behavior may be called the exception semantics.

- More commonly, failing to meet the precondition may simply mean the conditions are not ripe yet for the routine body to execute, without implying failure. The call in this case should just block the client from progressing, releasing the supplier’s processor for handling requests from other client processors, which may be expected to produce effects that will make the precondition true. When the supplier’s processor completes a call, it will examine the requests from

blocked clients and select one for which the blocking precondition is now true. This behavior may be called the waiting semantics.

Both strategies are compatible with the contract model. Experience with the practice of concurrent programming suggests we should retain the waiting semantics as the default. This is particularly clear in the bounded queue example: handled by a separate processor, the bounded queue becomes a bounded buffer, which various clients use to deposit and withdraw elements. Then if such a client calls *put* with the buffer full, or *item* with the buffer empty, it should normally be made to wait until another client has corrected the situation (through a *remove* in the first case and a *put* in the second).

Thus as a basic semantic rule, a precondition clause involving a separate call should cause the client to wait until the clause becomes satisfied. There may remain a need for the exception semantics when the client expressly wants to treat precondition violation as an exceptional case. A type example is access to a file, a case that occurs in sequential object-oriented programming, since the underlying concurrency is usually hidden. The following extract (using class *FILE* from the basic Eiffel library) would seem safe:

```
f: FILE;
...
if f/= Void and then f.readable then
  f.some_input_routine
  --some_input_routine is any
  --routine that reads data
  --from the file;
  ---its precondition is readable.
end
```

Here the test in the conditional instruction is meant to guarantee that  $f$  is attached to a readable file, meeting the precondition of *some\_input\_routine*. But although this extract does not explicitly show the concurrency, a file is a separate persistent structure, so the client has no way to avoid the case in which an interactive user (or some other software system) will access the file and make it unreadable between the test of *f.readable* and the execution of

*some\_input\_routine*. Here, however, waiting for *readable* to become true again is not the appropriate behavior; instead, the client should probably get an exception (from which it will be able to recover if it has a *Rescue* clause).

Because such cases are unlikely to be the most common, the waiting semantics will be the default. Through calls to routines of the kernel library it will be possible for a client to request the exception semantics.

### Reserving Objects

Considering separate precondition clauses as wait conditions does not yet provide us with enough control, especially if a client needs to reserve an object for the duration of several operations. In the bounded queue example, consider the client extract (again for *q* declared as a queue and *a* as a queue element):

```
--/3/
a := q.item;
... Other instructions (not calling
remove) ...
q.remove
```

The call to *item* accesses a queue element (after a wait if the queue is empty at the time of the call). The call to *remove* is intended to remove that element. But there is no guarantee that the two calls indeed manipulate the same queue element: with the mechanisms introduced so far, we have no way of preventing other separate clients from intruding between the two calls and performing one or more *remove*.

The problem arises even if the calls to *item* and *remove* are consecutive in the client text. Of course, in this case, we could solve the problem by renouncing the command-query distinction and adding a feature *get* to class *BOUNDED\_QUEUE*:

```
get: X is
  --Remove an element,
  --and return it as result.
  require
  not empty
do
  Result := item; remove
end
```

but this would imply a drastic change

in the recommended design style and in any case does not help us for the general case with one or more “Other instructions” between the two calls.

These observations highlight the need for a technique that will allow a client to reserve a separate supplier object for a certain period. Although routine calls similar to P and V operations on semaphores could be envisioned for that purpose, it is preferable, in the interest of deadlock prevention, to look for a linguistic construct: whereas it is practically impossible to ensure that every client that executes a P will later execute a V, a linguistic construct which reserves an object will have a fixed syntactic scope; execution of the construct will automatically release the object at the end of that scope.

Such a construct could have the following form:

```
--Note: uses a form not retained.
--For purposes of discussion only.
hold q then
... Here the client has exclusive
access to the object attached to q ...
end
```

Let *Q\_OBJ* be the object attached to *q*. The execution of such a construct would imply that if *Q\_OBJ* is already reserved the client will wait until *Q\_OBJ* is free again. Unfortunately, a simple *hold* construct as we have defined does not provide us with a flexible enough synchronization mechanism. Often a client will need to wait not just until a certain supplier becomes available, but also until a certain condition becomes true—for example *not q.empty* if the client needs to perform a *q.item* or *q.remove*. This suggests a variant of the *hold* construct, which includes waiting on one or more conditions. For example:

```
--Note: uses a form not retained.
--For purposes of discussion only.
hold q when not q.empty then
a := q.item;
... Other instructions ...
q.remove
end
```

Such a construct could be used as the basis of a workable solution. One advantage of such a solution would

be to give assertions their original semantics: within a *hold*, a precondition is a correctness condition; the proof rule (left for the reader to express) would indicate that, right after the *then*, the condition given by the *when* clause may be assumed.

On further examination, however, this solution is not fully satisfactory. A first problem that will have to be settled by the language designer is whether to allow any separate call (such as *q.item*) outside of a *hold* on the corresponding target (here *q*). The goal of consistency, always so important in language design, suggests prohibiting this. But then client code will be encumbered by many *hold* instructions.

Even if we decide against this strict policy, clients will have to include numerous *hold* to access such separate objects as bounded buffers. Any such situation—in which a clearly identified scheme occurs repeatedly in a certain application area—triggers the object-oriented designer’s basic instinct: encapsulate. Rather than individually wrapping every nontrivial buffer access in its own *local ... end*, the competent object-oriented designer will write a class which encapsulates the corresponding behavior (see Figure 6).

In Figure 6, any class needing that behavior will inherit from *BUFFER\_ACCESS*. Such classes are likely to be needed for every separate data structure. But then we may ask whether the *hold* construct is useful at all. If every significant use of a buffer *q* occurs through a routine with *q* as argument, could we not use argument passing as the mechanism for reserving objects? In other words, we might simply decide that whenever a routine call has a separate argument, any routine call will perform a “hold” operation on the corresponding separate object.

There remains the problem of waiting not just on object availability but also on conditions. Here the idea of preconditions as waiting conditions makes a comeback: since the routine cannot execute its body properly until its precondition becomes true, the precondition again presents itself as offering a natural wait condition. We may note here

that the style that would likely have become the most common for using the *hold* construct, as illustrated by *BUFFER\_ACCESS* in Figure 6, would entail frequent duplications of precondition texts as hold conditions; this is the case in Figure 6 with the use of *put*, *remove* and *item*.

The solution finally retained, then, involves the following elements:

1. No special "holding" construct is needed. *separate* declarations remain the only syntactical extension.
2. No separate call, of the form *t.f* (...) where the target *t* is *separate*, is permitted unless *t* is a formal argument of the routine in which the call appears.
3. If any nonvoid argument to a routine call is *separate*, the call will block until the corresponding object becomes available.
4. If a precondition clause of a routine *f* involves a *separate* call (whose target, because of rule 2, must be a formal argument of *f*), a call to *f* will block until the precondition clause is satisfied.
5. If neither rule 3 nor rule 4 causes a call to block, the call is said to be **satisfiable**. If there are one or more satisfiable calls on available objects handled by a given processor, one of them will proceed. (The semantics does not specify which satisfiable call will be selected if there are two or more; but it does require that one of them will proceed. In other words, a processor may not go on strike.)
6. "Available," as used in rules 3 and 5, is defined as follows: An object is *busy* if some call using it as a target has been started but not yet completed. The object is available if it is not busy and its processor is either idle (not executing any call) or blocked (as per rule 3 or 4) on a separate call executed on behalf of another object.

Rule 2 may seem unduly restrictive. It does bring, however, a much desirable extra safety, avoiding situations such as /3/ mentioned previously where the author of the client code mistakenly believes that two distinct calls using the same *separate* entity as targets, such as the calls *q.item* and *q.remove*, actually apply to

## Entities, Types and Values

The term 'entity', a generalization of the usual notion of variable, covers any name used to denote run-time values. Entities include the following: attributes of classes (see "Basic Terminology" sidebar); formal arguments of routines; local entities of routines (accessible only within the routine's body, and allocated anew for each execution of the routine, as with a local variable in Pascal); and *Result*, a predefined entity used in functions to denote the result to be returned to the client.

Types describe the run-time values that entities denote in the software text. There are two kinds of values: objects, and references to objects. Correspondingly, there are two kinds of types: expanded types, whose values are objects, and reference types, whose values are references to objects.

Expanded types include basic, predefined types such as *INTEGER* and *REAL*. So the value of an entity declared of type *INTEGER* is directly an integer value. A reference type may be obtained from a class which is not declared as *expanded*. For example, with the class *BOUNDED\_QUEUE* introduced earlier, the declaration

```
q: BOUNDED_QUEUE [X]
```

introduces *q* as an entity of reference type. The possible values of such an entity at run time are references: either void (not attached to any object), or attached to an object, here an instance of *BOUNDED\_QUEUE*. "Attached" is the technical term which describes the association between a nonvoid reference and an object.

the same object. Without rule 2, such errors would be likely to occur. They would be difficult to detect and debug, since the corresponding run-time behavior, which depends on how many other calls creep in between the two calls, is nondeterministic. Such errors typify the difficulty of constructing and debugging parallel programs, and we should not lightly forsake an opportunity to avoid them through a statically enforceable rule.

True, rule 2 will require some extra work on the part of the concurrent programmer: we can no longer write *q.remove* freely, but must enclose this call in a routine using *q* (assumed to be of type *separate BOUNDED\_QUEUE [X]* for some *X*) as argument. But this extra requirement seems justified in light of the gain in reliability; in addition, we may expect that many uses of such structures will rely on a class encapsulating the appropriate behavior, such as *BUFFER\_ACCESS* in Figure 6, where the problem is taken care of once and for all. As before, classes needing these facilities will inherit from *BUFFER\_ACCESS*. The

class *BUFFER\_ACCESS* will now be written more simply as shown in Figure 7.

The use of separate actual arguments as ways to reserve objects requires one further comment. If a routine call names two or more such arguments, the call will block until it has got hold of **all** of them, and satisfied the corresponding preconditions. This may prove difficult to implement, especially in a distributed system, which will require consensus between the various processors. It is possible to restrict the mechanism by allowing a routine call to involve at most one separate argument. Then it will be the individual programmer's job to reserve all necessary resources through nested calls, using one of the algorithms described in the literature (e.g., see [4]). It may be preferable, however, to keep the model without restrictions, following the argument that tedious and error-prone tasks should be handled whenever possible, by the programming environment rather than by individual programmers. This reasoning is central in object-oriented programming, where it justifies such important fa-

```

-- Note: uses a form not retained. For purposes of discussion only.
class BUFFER_ACCESS [G] feature
  put (q: separate BOUNDED_QUEUE [G]; x: G) is
    -- Insert x into q, waiting if necessary until there is room.
    do
      hold q when not q.full then q.put (x) end
    end;

  remove (q: separate BOUNDED_QUEUE [G]) is
    -- Remove an element from q, waiting if necessary
    -- until there is such an element.
    do
      hold q when not q.empty then q.remove end
    end;

  item (q: separate BOUNDED_QUEUE [G]): G is
    ... Left to the reader ...
end

```

**Figure 6.**

```

-- Encapsulation of access to bounded buffers
class BUFFER_ACCESS [X] is
  put (q: separate BOUNDED_QUEUE [G]; x: G) is
    -- Insert x into q, waiting if necessary
    -- until there is room.
    require
      not q.full
    do
      q.put (x)
    ensure
      not q.empty
    end;

  remove (q: separate BOUNDED_QUEUE [G]) is
    -- Remove an element from q, waiting if necessary
    -- until there is such an element.
    ... Left to the reader ...

  item (q: separate BOUNDED_QUEUE [G]): G is
    -- Oldest element not yet consumed
    require
      not q.empty
    do
      Result := q.item
    ensure
      not q.full
    end
end

```

**Figure 7.**

$$\frac{\{ INV \wedge p \in Pre(f) \} \text{ body}(f) \{ INV \wedge q \in Post(f) \}}{\{ p \in Nonsep\_pre(f) \} \text{ t.f } \{ q \in Nonsep\_post(f) \}}$$

**Figure 8.**

cilities as garbage collection and compiler-applied static binding.

**A Proof Rule**

In the scheme retained here, waiting on a precondition clause occurs only for a precondition of the form *t.cond*, where *t* is a formal argument of the enclosing routine and is separate. This is why the preconditions on the routines of *BUFFER\_ACCESS* are needed: violation of a precondition of a *BOUNDED\_QUEUE* would be the routines of correctness violation, not a wait condition, since these routines do not have any separate argument. The routines of *BUFFER\_ACCESS*, however, have a separate argument *q*, so any precondition clause of the form *q.some\_condition* in these routines is a waiting condition.

In general, then, in a routine of the form

```

f (... , a: T, ...) is
  require
    clause1;
    clause2;
    ...
  do ... end

```

any of the precondition clauses which does not involve any separate call on a separate formal argument is a correctness condition: any client must ensure that condition prior to any call, otherwise the call is in error. Any precondition clause involving a call of the form *a.some\_condition*, where *a* is a separate formal argument, is a wait condition which will cause calls to block if it is not satisfied. These observations may be expressed as a proof rule, shown in Figure 8 which, for separate computation, replaces the sequential rule given in Figure 3.

In Figure 8, *nonsep\_pre(f)* is the set of clauses in *f*'s precondition which do not involve any separate calls, and similarly for *Nonsep\_post(f)*.

This rule captures in part the essence of parallel computation. To prove a routine correct, we must still prove the same conditions (those above the line) as in the sequential rule. But the consequences on the properties of a call (below the line) are different: the client has fewer properties to ensure before the call,

since trying to ensure the separate part of the precondition would be futile anyway; but we also obtain fewer guarantees on output. The former difference may be considered good news for the client, the latter is bad news.

The separate clauses in preconditions and postconditions thus join the invariant as properties that must be included as part of the internal proof of the routine body, but are not directly usable as properties of the call. The rule also serves to restore the symmetry between preconditions and postconditions, a job that will be completed by considering the lazy wait technique. The discussion so far was essentially based on an analysis of the properties of preconditions.

### Comments on the Use of Preconditions

The idea that assertions, and in particular preconditions, may have two different semantics—sometimes correctness conditions, sometimes wait conditions—is somewhat disturbing. Yet in the design of this mechanism, the idea kept returning each time it was discarded.

One possible objection, however, is unjustified. In Eiffel, run-time assertion checking may be turned on or off as a result of a compilation switch. Is it not dangerous, then, to attach that much semantic importance to preconditions in concurrent object-oriented programming?

Such an objection misses, however, the true nature of assertions. Assertions are not primarily a debugging or run-time checking tool. Instead, one should view assertions as full-fledged components of classes. In the form for *put*, the precondition and postcondition belong to the routine just as much as the *do* clause. They are part of an important property of the routine: its specification. Although this may appear paradoxical, the compilation option that switches run-time assertion checking on or off does *not* affect the semantics of the language. This is because the semantics of any language is defined for correct programs only. But a program whose execution may violate an assertion is incorrect! (The definition of a correct class is precisely that the

*do* clauses of its routines are compatible with the assertions.)

To a practicing programmer, the argument may appear specious, since checking assertions at runtimes may be the best way to determine that a class is incorrect. But in principle it should be possible to prove class correctness statically; run-time monitoring is only an imperfect solution. (This is particularly difficult to explain to C programmers, who when they accept assertions at all, tend to see them just as executable constructs, meant to check certain properties at run time for debugging purposes.)

Assertions, then, are always part of the software, whether or not they are monitored at run time. The difference between their separate and nonseparate clauses is simply that the semantics does not require a nonseparate clause to be evaluated at run time if the software is correct, but does require run-time checking of separate clauses.

### Lazy Wait

One more rule is needed for a full definition of the semantics. This rule will determine how two processors are resynchronized when one needs a result from another. The rule, which may be called lazy wait by analogy with the so-called lazy evaluation of programs written in functional languages such as Lisp, requires a client to wait when it absolutely needs information resulting from a call, but no sooner.

Consider a routine containing a separate call:

```
r(..., t: separate SOME_TYPE, ...) is
do
  ...; t.f(...); --*
  other_instructions
end
```

Assume *r* is executed on an object C\_OBJ and, as part of its execution, executes the call marked \*, with target *t*. Let T\_OBJ be the object attached to *t*. These rules indicate that *t* must be a formal argument of *r*; so C\_OBJ will have obtained hold of T\_OBJ prior to the call.

For the preceding syntax to be legal, *f* must be a procedure. The call to *f* proceeds when it is satisfiable and

selected by the processor of T\_OBJ, which then starts executing it. But C\_OBJ is handled by another processor, and that processor does not need to wait for the call to terminate: it should immediately continue with its further business, *other\_instructions*. This will be the rule according to the lazy wait policy: separate procedure calls should not (except for one special case discussed later) cause the client to wait.

Assume the *other\_instructions* contain more procedure calls on the same separate target T\_OBJ. These calls, like the first one, will not make the client wait. It is essential for consistency, however, that the order in which they will eventually get serviced be the order in which the client, here C\_OBJ, has requested them. The semantic rule must guarantee this property of order preservation.

When then, if ever, must C\_OBJ wait for the termination of the call to *f*? The lazy answer is that an object will wait when it needs the result of a query (function or attribute call) on a separate object. Assume for example that the *other\_instructions* are of the form

```
other_instructions_1; ...
other_instructions_n;
k: = t.some_value
```

where *some\_value* is an integer attribute or function applicable to *t*, and *k* is an auxiliary entity, also of integer type. To assign the value of the last call to *k*, we need access to T\_OBJ; thus the call to *f* must have been completed, as well as any intervening call with T\_OBJ as its target, and if *some\_value* is a function we must also wait for the computation of that function to complete. (Remember that the entire fragment appears in a routine of which *t* is a formal argument, so there is no danger of another client sneaking in and grabbing T\_OBJ between the call to *f* and the call to *some\_value*: C\_OBJ has exclusive separate access to T\_OBJ throughout the execution of this fragment.)

The lazy wait policy requires that a client performing a query on a separate target shall not proceed until all earlier calls have been completed in the order logged, and the query itself

has been executed to completion. Together with the rule that a call proceeds when it is satisfiable, this property is the principal semantic difference between concurrent and sequential object-oriented computation. (It would need to be included in the semantics of calls defined in Chapter 21 of [14].) The lazy policy may be seen as a form of waiting on the postcondition, in the same sense that the basic mechanism causes waiting on the precondition.

A note is in order on the rule that clients do not wait for separate procedure calls to terminate. At first this may be viewed as just an optimization: forcing clients to wait for the termination of such calls before proceeding would simply make less use of the available parallelism. But such a policy would be misguided. If the *other\_instructions* of the preceding pattern are empty, as in

```
r(..., t: separate SOME_TYPE, ...) is
do ... t.f(...) end
```

the execution of *r* will terminate immediately after launching the call to *f*, without waiting for that call to terminate. This may cause the caller of *r* to relinquish the reservation it made on the client of T\_OBJ, which may be crucial if that object is needed by other clients.

A common variant of this scheme, which would not work without the ability to continue execution after starting a separate procedure call, is a loop which gets a number of previously created objects started with their own lives. For example:

```
launch (a: ARRAY [separate X]) is
-- Get every element of a started.
require
-- No element of a is void
local
i: INTEGER
do
from i := 1 until i > a.count loop
  launch_one (a @ i); i := i + 1
end
end
launch_one (p: separate X) is
-- Get p started.
require
p /= Void
do p.live end
```

(The notation *a @ i* denotes the *i*th element of *a*. @ is an infix feature of

the Kernel Library class *ARRAY*. The array must have been previously initialized to hold nonvoid elements, numbered from 1 to *a.count*.) If, as may well be the case, procedure *live* of class *X* describes an infinite process, this scheme only works if each loop iteration proceeds immediately after starting *launch\_one*; otherwise the loop would never get beyond its first iteration.

### Passing Nonseparate References

One technical point remains to be clarified regarding the passing of arguments. In some cases we may wish, in a class *B*, to include a separate call of the form

```
t.f (a)
```

where the actual argument *a* is of a reference type, but not separate. The corresponding routine, in the supplier class *C*, could be of the form

```
f (x:SOME_TYPE) is
do ... y: = x ... end
```

where *y* is some attribute of *C*.

All the entities involved—*a*, *x*, *y*—denote references to objects. The assignment enables *C* to retain a reference to an object A\_OBJ handled by the client processor, so that later, after the termination of *f*, a routine of *C* may execute a call of the form

```
y.some_routine (...)
```

The object to which this call applies is A\_OBJ; but the call should be separate since the processors handling the instances of *B* and *C* are different! Unless *x*, and hence *y* too, are declared as *separate*, the call will have the wrong semantics. A language rule will require that any such formal argument should be declared as *separate*. In the example, *SOME\_TYPE* must be a separate type. The rule on those assignments then requires *y* also to be declared as *separate*.

The case examined here has a consequence on the semantics of calls: because the body of *f* may perform a call with target A\_OBJ, we cannot let the client proceed even if *f* is a procedure. The call to *f* should block the client until *f* terminates. This will be the general semantic rule if we allow separate calls having nonseparate

references as actual arguments: any such call will block the client with no possibility of lazy wait. This rule introduces a small but unpleasant complication. We could avoid it by disallowing nonseparate reference arguments altogether in separate calls. An upcoming example will show, however, that it is convenient (although not strictly indispensable) to keep this possibility open.

The preceding discussion only applies to reference types. If *a* is of an expanded type, its possible values are objects, and attachment (in particular argument passing) implies copying an object onto another rather than assigning a reference, so this problem does not arise. This possibility immediately suggests, however, a new validity constraint: no expanded type may include separate attributes.

### Duels

The preceding discussion concludes the presentation of the basic semantic model. In practice, we need an extra degree of flexibility to have a fully usable mechanism. The major remaining problem is how a client can avoid waiting on a nonavailable separate supplier, or one that does not satisfy a separate precondition clause. As noted in the example of reading a file that may suddenly have become unreadable, the appropriate policy in such a case may be to cause an exception in the client, not to make it wait.

How best should we provide such flexibility? The situation is comparable to the problem of exerting fine control on exception handling in sequential object-oriented computation. The basic Eiffel exception mechanism is extremely simple, providing a Rescue clause enabling a routine to catch exceptions, and a Retry instruction enabling the routine to try again after attempting to correct the cause of the exception. These facilities suffice in most cases. To obtain further facilities for more sophisticated cases, without complicating the language or changing the basic semantics of exception handling, programmers may use the class *EXCEPTIONS* from the Kernel Library, offering such features as the



integer code of the last exception (so that a Rescue clause may discriminate between different causes of exception) and many others. Classes needing these facilities simply inherit from *EXCEPTIONS*.

Here we may use a comparable approach by defining a simple library class *CONCURRENCY*, from which classes may freely inherit. This class will include a procedure *immediate\_service*. A call to this procedure alters the normal response to a nonsatisfiable call: instead of waiting, if the object is not available or a separate condition is not satisfied, it will cause an exception in the client. Procedure *normal\_service* will serve to restore standard behavior.

This can be used immediately in the file-reading example, assuming *FILE* now becomes a separate class:

```
analyze (f: FILE) is
do
  immediate_service;
  if f.readable then
    ... Various read operations of form
    f.input_routine ...
  else
    ... Report impossibility in some way
    ...
  end;
  normal_service
rescue
  ... Take corrective action, or report
  problem ...
end
```

In the contract method of software construction, exceptions should be reserved for truly exceptional cases which would be difficult or impossible to handle otherwise without adding great complication to the text of the software. (This is in marked contrast with, for example, the Ada or CLU mechanisms, which essentially use exceptions as interprocedural control structures.) The preceding use satisfies this requirement, since it applies to a case that is unlikely to occur often—the case in which a file was checked to be readable but becomes unreadable or inexistent before its actual use because of the actions of some independent agent. Although infrequent, this case must still be handled properly.

A similar mechanism may be used to provide a safe form of express messages. Assume we want to allow a

VIP client, the challenger, immediate access to a separate object; if any other client, the holder, currently holds it, it will be interrupted. The only acceptable solution in this case, in light of the Design by Contract principles, is to make the holder's call fail, so the holder will be forced to take corrective action (if it has a Rescue clause) or to fail. Because it has such drastic consequences for the holder, however, this behavior should only occur if the holder has accepted the possibility by calling procedure *yield* of class *CONCURRENCY*. In this case, a challenger that has called *immediate\_service* and finds the desired object nonavailable will not get an exception itself, but will actually get the object, and cause an exception in the holder. (The immediate recipient of the exception is the common supplier, which after executing its Rescue clause to restore its invariant will normally pass on the exception to the holder.) To restore the normal (no-yield) behavior, the holder will call procedure *insist*.

The result of the potential conflict between a holder and a challenger, which we may call a *duel*, is specified in all cases by Table 2. The default behavior for each duelist, as well as the overall default in the absence of any call to the procedures of class *CONCURRENCY* (top-left table entry), are underlined.

### Examples

To illustrate the mechanism, the following examples have been chosen from diverse backgrounds. To conserve space, only the most salient points of the examples are shown; the full texts are included in [16].

#### The Dining Philosophers

We begin with the inevitable dining philosophers. It seems hard to provide a simpler form for describing the philosophers' behavior than the following class, which relies on the general control structure provided by *PROCESS*:

```
separate class PHILOSOPHER
creation
  make
inherit
  PROCESS
```

```
  rename setup as getup end
feature {BUTLER}
step is
  -- Think and eat
do
  think; eat (left, right)
end
feature {NONE}
  -- The two required forks:
left, right: separate FORK;
getup is do ... end;
think is do ... end;
eat (l, r: separate FORK) is
  -- Eat, having grabbed
  -- l and r.
do ... end
feature {NONE}--Creation
make (l, r: separate FORK) is
  --Assign l and r
  --as the required forks.
do
  left := l; right := r
end
end
```

The entire synchronization requirement is embodied by the call to *eat*, which uses arguments *left* and *right* representing the two necessary forks, thus reserving these objects.

Of course the simplicity of this solution comes from the mechanism's ability to reserve several resources through a single call having two or more separate arguments, here *left* and *right*. If we restricted the separate arguments to at most one per call, the solution would use one of the published algorithms (e.g., see [4]) for getting hold of two forks one after the other without causing deadlock.

Thanks to the use of multiple object reservation through arguments, the solution described here does not produce deadlock, but it is not guaranteed to be fair; some of the philosophers can conspire to starve the others. Here too the literature provides various solutions, which may be integrated into the preceding scheme.

The notation *feature* {*BUTLER*} indicates that the following features are exported only to class *BUTLER*, the controller class, and its descendants if any. The notation *feature* {*NONE*}, already encountered in earlier examples, similarly makes features available to *NONE* only. *NONE*,



a Kernel Library class, has no descendants, so this implies making the features secret. Most of the features appearing in the following examples will be secret or selectively exported in this way. Such a situation is much less common in sequential object-oriented programming, but the difference of style is no accident: in concurrent object-oriented programming, many classes are the equivalent of processes in ordinary concurrent programming, which need few or no exported features. Class *FORK* has no feature for this example. The details of class *BUTLER*, used to set up and start sessions, are left to the reader (see [16]).

### Making Full Use of Hardware Parallelism

The example shown in Figure 10, although somewhat academic, illustrates how lazy wait can be used to draw the maximum benefit from any available hardware parallelism—the figure depicts an extract (not involving concurrency) from a class describing binary trees. Function *nodes* makes a standard use of recursion to compute the number of nodes in a tree. The recursion is indirect, through *node\_count*.

In a concurrent environment offering many processors, it is tempting to study how we could offload all the separate node computations to different processors. Declaring the class as *separate*, replacing *nodes* by an attribute and introducing procedures does the job, as shown in Figure 11. The recursive calls to *compute\_nodes* will now be started in parallel. The addition operations wait for these two parallel computations to complete.

If an unbounded number of CPUs (physical processors) are available, this solution seems to make the optimal possible use of the hardware parallelism. If there are fewer CPUs than nodes in the tree, the speedup over sequential computation will depend on how well the implementation allocates CPUs to the (virtual) processors.

The presence of two tests for vacancy of *b* may appear unpleasant. It results, however, from the need to separate the parallelizable part—the

procedure calls, launched concurrently on *left* and *right*—from the additions, which by nature must wait for their operands to become ready.

### Coroutines

One of the requirements stated at the beginning of the article was that the mechanism should support coroutine programming. Figures 12a and 12b show two classes which achieve this goal.

One or more coroutines will share one coroutine controller (set up by a “once” function, the mechanism which allows information sharing while avoiding the dangers of global variables [11, 14]). Each coroutine has an integer identifier. To resume a coroutine of identifier *i*, procedure *resume* will, through *actual\_resume*, set the *next* attribute of the controller to *i*, and then block, waiting on the precondition *next = j*, where *j* is the coroutine’s own identifier. This ensures the desired behavior.

Although this solution looks like normal concurrent software, it is organized in such a way that if all coroutines have different identifiers only one coroutine may proceed at any one time; so it is useless to allocate more than one physical CPU for them.

The recourse to integer identifiers is necessary, since giving *resume* an argument of type *COROUTINE*, a separate type, would cause deadlock. Using Eiffel’s “unique” declaration (similar to Pascal’s enumerated types), programmers do not need to worry about assigning such values manually. This use of integers also has an interesting consequence: if we allow two or more coroutines to have the same identifier, then with a single CPU we obtain a nondeterministic mechanism: a call to *resume* (*i*) will cause the restarting of any coroutine whose identifier has value *i*. With more than one CPU a call *resume* (*i*) will allow all coroutines of identifier *i* to proceed in parallel. This scheme, which for a single CPU provides a coroutine mechanism, doubles up in the case of several CPUs as a mechanism for controlling the maximum number of processes of a certain type which may be simultaneously active.

### Locking and Semaphores

Assume we want to allow a number of clients (the “lockers”) to obtain exclusive access to certain resources (the “lockables”) without having to enclose the exclusive access sections in routines, using a semaphore-like technique. A solution is shown in Figure 13.

Any class describing resources will inherit from *LOCKABLE*. The proper functioning of the mechanism assumes that every locker performs sequences of *grab* and *release* operations, in this order. Other behavior will usually result in deadlock. We can once again rely on the power of object-oriented computation to enforce the required protocol; rather than trusting every locker to behave, we may require lockers to go through a procedure *use* in descendants of a class *LOCKING*, which describes the required behavior. Class *LOCKING* is left to the reader (see [16]); it will inherit from *PROCESS*.

Whether or not we go through class *LOCKING*, a *grab* does not reserve the corresponding lockable for all competing clients: it only excludes other lockers observing the protocol. To exclude any possible client from accessing a resource, you must enclose the operations accessing the resource in a routine to which you pass it as argument.

The reader may have noted that in routine *grab* of class *LOCKER*, the separate call to *set\_holder* passes *Current* as argument. *Current*, a reference to the current object (the call’s client) is a nonseparate reference. Accordingly, the corresponding formal argument of *set\_holder* is declared as *separate*. Without the possibility of passing nonseparate references as arguments to separate calls, we would need to rely, as with the coroutine examples, on a scheme associating an integer with every instance of class *LOCKER*.

### An Elevator Control System

A case in which object-orientation and the mechanism defined in this article can be used to achieve a decentralized event-driven architecture is shown in Figures 14 through 17. The example describes software for an elevator control system, with sev-

eral elevators serving many floors. The design is somewhat fanatically object-oriented in the sense that every significant type of component in the physical system—for example the notion of individual button in an elevator cabin, marked with a floor number—has an associated separate class, so each corresponding object such as a button has its own virtual thread of control (processor). The benefit is that the system is entirely event-driven; in particular, it does not need to include any loop for examining repeatedly the status of objects, for example whether any button has been pressed. The class texts are only sketched, but provide a good idea of what a complete solution would be. The creation procedures, which must perform the necessary initializations, have been left to the reader.

It is convenient to start with class *MOTOR*. This class describes the motor associated with one elevator cabin, and the interface with the mechanical hardware (see Figure 14).

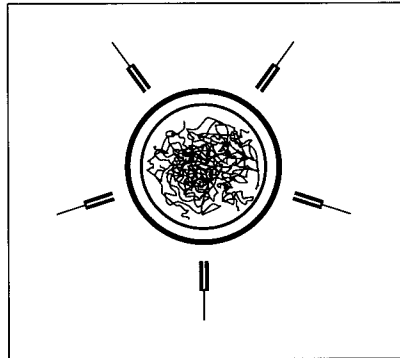
The creation procedure of this class must associate an elevator, *cabin*, with every motor. Class *ELEVATOR* (Figure 15) includes the reverse information through attribute *puller*, indicating the motor pulling the current elevator.

The reason for making an elevator and its motor separate objects is to reduce the grain of locking: once an elevator has sent a *move* request to its motor, it is free again, thanks to the lazy wait policy, to accept requests from buttons either inside or outside the cabin. It will resynchronize with its motor on receipt of a call to procedure *record\_stop*, through *signal\_stopped*. The actual time during which an instance of *ELEVATOR* will be reserved by a call from either a *MOTOR* or *BUTTON* object is very short.

There are two kinds of buttons: floor buttons, which passengers press to call the elevator to a certain floor, and cabin buttons, which are inside a cabin and are pressed by passengers to request a move to a certain floor. There is of course a difference between the corresponding kinds of requests: a request from a cabin but-

**Table 2.** The semantics of duels

Challenger →	<i>normal_service</i>	<i>immediate_service</i>
↓ Holder ↓		
<i>insist</i>	Challenger waits	Exception in challenger
<i>yield</i>	Challenger waits	Exception in both



**Figure 9.**

ton is directed to the cabin to which the button belongs, whereas a floor button request may be handled by any elevator. A request of the latter kind will be sent to a dispatcher object, which will poll the various elevators and select the one that can best handle the request. (The precise selection algorithm is left unimplemented, since it is irrelevant to this discussion; the same applies to the algorithm used by elevators to manage their *pending* queue of requests in class *ELEVATOR* in Figure 15.)

Class *FLOOR\_BUTTON* (Figure 16) assumes that there is only one button on each floor. It is not difficult to update the design to support two buttons, one for up requests and the other for down requests. It is convenient, although not essential, to have a common parent *BUTTON* for the classes representing the two kinds of buttons. Remember that the features exported by *ELEVATOR* to *BUTTON* are, through the standard rules of selective information hiding, also exported to the two descendants of this class.

The class *DISPATCHER* is shown in Figure 17. To refine the algorithm that selects an elevator in procedure *accept*, it will need to access the attri-

butes *position* and *moving* of class *ELEVATOR*, which in the full system should be complemented by a Boolean attribute *going\_up*. Such accesses will not cause any problem as the design ensures that *ELEVATOR* objects never get reserved for a long time.

#### A Watchdog Mechanism

Together with the previous one, the last example helps show the applicability of the mechanism to real-time problems. It also illustrates the concept of duel. We want to enable a class to perform a call to a certain procedure *action*, with the provision that the call will be interrupted, and a Boolean attribute *failed* set to true, if the procedure has not completed its execution after *t* seconds. The only basic timing mechanism available is a procedure *wait(t)*, which will execute for *t* seconds.

Figure 18 depicts the solution, using a duel. A class that needs the mechanism should inherit from class *TIMED* and provide an effective version of the procedure *action* which, in *TIMED*, is deferred. To let *action* execute for at most *t* seconds, it suffices to call *timed\_action(t)*. This procedure sets up a watchdog (an instance of class *WATCHDOG*), which executes *wait(t)* and then interrupts its client. If, however, *action* has been completed in the meantime, it is the client that interrupts the watchdog.

In class *WATCHDOG*, the Boolean attribute *terminated\_before\_time* and the Rescue clause are not actually needed; they have been added for clarity. In a more robust version, the Rescue clauses should test for the type of exception, in order to use the preceding scheme only for exceptions caused by a call from a client having requested *immediate\_service*.

```

class BINARY_TREE [G] feature
  left, right: BINARY_TREE [G]; ... Other features ...

  nodes: INTEGER is
    -- Number of nodes in this tree
    do
      Result := node_count (left) + node_count (right) + 1
    end
  feature {NONE}
    node_count (b: BINARY_TREE [G]): INTEGER is
      -- Number of nodes in b
      do
        if b /= Void then
          Result := b.nodes
        end
      end
end

```

**Figure 10.**

```

separate class BINARY_TREE [G] feature

  left, right: BINARY_TREE [G];
  ... Other features ...

  nodes: INTEGER;

  update_nodes is
    -- Update nodes to reflect the number of nodes in this tree.
    do
      nodes = 1;
      compute_nodes (left); compute_nodes (right);
      adjust_nodes (left); adjust_nodes (right)
    end
  feature {NONE}
    compute_nodes (b: BINARY_TREE [G]) is
      -- Update information about the number of nodes in b.
      do
        if b /= Void then b.update_nodes end
      end;

    adjust_nodes (b: BINARY_TREE [G]) is
      -- Adjust number of nodes from those in b.
      do
        if b /= Void then
          nodes := nodes + b.nodes
        end
      end
end

```

**Figure 11.**

**The Mechanism**

Following is the precise description of the mechanism that results from the preceding discussion. The description consists of three parts: syntax; validity rules (static semantic constraints); semantics. It must be understood as an extension to the Eiffel language specification as given in [14].

**Syntax**

The syntactic extension involves just one new keyword, **separate**. (In preparation for this extension, the

list of keywords in Appendix G of [14] already included **separate**.) A declaration of an entity or function, which normally appears as *x: TYPE*, may now also be of the form

```

--/6/
x: separate TYPE

```

In addition, a class declaration, which normally begins with one of **class C**, **deferred class C** or **expanded class C**, may now be of the form:

```

--/7/
separate class C ...

```

In this case *C* is called a “separate class.” The syntactic convention implies that separate status for a class is incompatible both with expanded status and with deferred status. As in the case of expanded and deferred classes, the property of being separate is not inherited: a class is separate or not according to its own declaration only, regardless of the separateness status of its parents. An entity or function *x* will be said to be separate if either it is declared under form /6/, or its type is based on a separate class (declared under form /7/). It is not an error for both of these to apply: in form /6/, *TYPE* may be based on a separate class, although in this case the use of **separate** in the declaration of *x* is redundant.

**Constraints**

A number of validity rules apply to constructs involving separate entities.

The first rule applies to a declaration of the preceding form /6/:

In a type of the form **separate TYPE**, the base class of *TYPE* must be neither deferred nor expanded.

The next rule governs the combination of separate and nonseparate elements in an attachment. The term “attachment” covers both assignment and argument passing, which have the same semantics. An attachment of *y* to *x* is either an assignment of the form *x* = *y*, or an argument passing in a call of the form *f* (... , *y*, ...) or *t.f* (... , *y*, ...), where the formal argument of the routine *f*, at the position corresponding to *y*, is *x*. Here is the attachment rule:

In an attachment of *y* to *x*, if the source *y* is separate, the target *x* must also be separate.

It is permitted to redefine an entity from separate to nonseparate and conversely (with the corresponding constraints on polymorphic attachment, not detailed here).

The following rules, explained earlier, apply to separate calls (calls whose targets are separate).

- For a separate call to be valid, the call must appear in a routine and its target must be a formal argument of the enclosing routine.

- If an actual argument of a separate call is of a reference type, the corresponding formal argument must be declared as separate.
- If an actual argument of a separate call is of an expanded type, the corresponding base class may not have any attribute of reference type.

The last constraint removes any possibility that the very process of evaluating a precondition clause (as part of an access to a separate object), or more generally an assertion, could lead to further blocking:

If an actual argument *a* to a function call is separate and the call is part of an assertion, then the call must be part of a routine and *a* must be a formal argument of that routine.

This rule (which implies that a function call appearing in a class invariant may not use a separate entity as actual argument) was not introduced in the earlier discussion. Without it, we could have a routine of the form

```
f (x:SOME_TYPE) is
  require
    some_property (separate_attribute)
  do ... end
```

where *separate\_attribute* is a separate attribute of the enclosing class. But then the evaluation of *f*'s precondition, either as part of assertion monitoring for correctness, or as a synchronization condition if the actual argument corresponding to *x* in a call is itself separate, could cause blocking if the attached object is not available! Such behavior is clearly unacceptable.

This concludes the list of validity constraints. The semantics will be discussed in three parts: creation; object states; calls.

### Semantics of Creation

If the target *t* of a creation instruction (which appears after the second exclamation mark in, for example! *t.make (...)*) is nonseparate, the newly created object will be handled by the same processor as the creating object. If, however, *t* is separate, the new object will be allocated to a new processor.

```
separate class COROUTINE creation
  make
  feature {COROUTINE}
    resume (i: INTEGER) is
      -- Wake up coroutine of identifier i and go to sleep.
      do actual_resume (i, controller) end

  feature {NONE}
    controller: COROUTINE_CONTROLLER;

    identifier: INTEGER;

    actual_resume (i: INTEGER; c: COROUTINE_CONTROLLER) is
      -- Wake up coroutine of identifier i and go to sleep.
      -- (Actual work of resume).
      do
        c.set_next (i); request (c)
      end;

    request (c: COROUTINE_CONTROLLER) is
      -- Request eventual re-awakening by c.
      require
        c.is_next (identifier)
      do end

  feature {NONE} -- Creation
    make (i: INTEGER; c: COROUTINE_CONTROLLER) is
      -- Assign i as identifier and c as controller.
      do
        identifier := i; controller := c
      end

end
```

Figure 12a.

```
separate class COROUTINE_CONTROLLER feature {NONE}
  next: INTEGER
  feature {COROUTINE}
    set_next (i: INTEGER) is
      -- Select i as the identifier of the next coroutine to be awakened.
      do next := i end;

    is_next (i: INTEGER): BOOLEAN is
      -- Is i the index of the next coroutine to be awakened?
      do Result := (next = i) end

end
```

Figure 12b.

### Object States

Once it has been created, an object OBJ will be in either of three states:

- Busy: a routine is being executed on OBJ for the benefit of a client object, and is not blocked.
- Blocked: the last routine to be started on OBJ has requested access to an object handled by a different processor, but the processor was not available or a separate precondition was not satisfied.
- Idle: none of the preceding. There is no pending client request on OBJ.

In the object-oriented style of pro-

gramming, a call is always executed on behalf of "someone" else—a client. In sequential computation, the client is always another object, except in the case of the system's root object, whose client is the human or other system that started the execution. With concurrent computation, the only new property is that a system may have more than one root object.

Figure 19 informally shows the transitions between states, as discussed in the following two subsections. A processor is said to be available if and only if every one of the objects that it handles is either idle or blocked.

```

class LOCKER feature
  grab (resource: separate LOCKABLE) is
    -- Request exclusive access to resource.
    require
      not resource.locked
    do resource.set_holder (Current) end

  release (resource: separate LOCKABLE) is
    require
      resource.is_held (Current)
    do resource.release end
end

class LOCKABLE feature {LOCKER}
  set_holder (l: separate LOCKER) is
    -- Designate l as holder.
    require
      l /= Void
    do
      holder := l
    ensure
      locked
    end;

  locked: BOOLEAN is
    -- Is resource reserved by a locker?
    do Result := (holder /= Void) end;

  is_held (l: separate LOCKER): BOOLEAN is
    -- Is resource reserved by l?
    do Result := (holder = l) end;

  release is
    -- Release from current holder.
    do
      holder := Void
    ensure
      not locked
    end

feature {NONE}
  holder: separate LOCKER
end

```

**Figure 13.**

```

separate class MOTOR feature {ELEVATOR}
  move (floor: INTEGER) is
    -- Go to floor, once there, report.
    do
      "Direct the physical device to move to floor";
      signal_s:opped (cabin)
    end;

  signal_stopped (e: ELEVATOR) is
    -- Report that elevator stopped on level e.
    do
      e.record_stop (position)
    end

feature {NONE}
  cabin: ELEVATOR;

  position: INTEGER is
    -- Current floor level
    do
      Result := "The current floor level, read from physical sensors"
    end
end

```

**Figure 14.**

## Semantics of Calls

To study the semantics of calls, we may use the general form  $t.f(\dots, s, \dots)$  and assume that  $f$  is a routine. (If  $f$  is an attribute, we may replace calls to  $f$ , for the purpose of this discussion, by calls to a function whose sole purpose is to return the value of the attribute; this may affect performance but does not change the semantics.)

The introduction of concurrency affects the semantics of a call only if one or more of the elements involved—target and actual arguments—are separate. Let us assume that one or more of the actual arguments are separate, but the target  $t$  is not. (The effect of having the target separate is examined in the next subsection.) The call is executed as part of the execution of a routine on a certain object, say C\_OBJ, which may only be in a busy state at that stage. The property which determines whether the call proceeds or blocks may be defined as follows:

### Definition (satisfiable call)

A call to a routine  $f$  is satisfiable if and only if every separate actual argument  $a$  having a nonvoid value, and hence attached to a separate object A\_OBJ, satisfies the following three conditions: A\_OBJ is idle; the processor handling A\_OBJ is available; every separate clause of the precondition of  $f$ , when evaluated for A\_OBJ and the actual arguments given, has value true. In this definition an assertion clause is separate if and only if it includes a call whose target is separate.

If the call is satisfiable, it proceeds immediately: C\_OBJ remains in the busy state and A\_OBJ enters the busy state, in which it executes  $f$ . When the execution of  $f$  terminates, A\_OBJ returns to the idle state.

If the call is not satisfiable, C\_OBJ enters the blocked state. The call attempt has no immediate effect on its target and actual arguments. The processor handling C\_OBJ becomes available; in particular, if any object also handled by that processor had previously blocked on a call that is now satisfiable, the processor will select one such object, return it to the busy state and execute the corresponding call as described. If there is

more than one selectable candidate, the semantics does not specify which one will be selected, but does require the processor to select one. If a processor is available and there is no immediately satisfiable call, the processor remains available until a state is reached in which there is a satisfiable call for that processor.

A note is in order on exactly when precondition clauses will cause blocking in the definition of satisfiability. A precondition clause is separate (and so may cause blocking) only if it includes a call on a separate target. Other precondition clauses remain correctness conditions, not blocking conditions. For example, a precondition of the form  $i > 1$ , where  $i$  is an integer entity (expanded, and hence nonseparate) is a normal correctness condition, which must be ensured before the call by the client. Failure to ensure this condition at the time of the call indicates a bug in the client, not a run-time waiting condition. Similarly, for separate  $s1$  and  $s2$ , precondition clauses of the form  $s1 \neq \text{Void}$  or  $s1 = s2$  are correctness conditions, since they do not involve any calls. Both are tests for reference equality, which may be performed regardless of the status (idle, busy, waiting) of the objects attached to  $s1$  and  $s2$ , if any. In particular, separate actual arguments may cause blocking only if they are non-void.

An adaptation of the preceding description is applicable to classes inheriting from class *CONCURRENCY*. The effect of calls to procedures *immediate\_service*, *normal\_service*, *insist* and *yield* is defined by Table 2.

### Separate Targets and Lazy Wait

The final semantic change is the effect of having a separate target  $t$  in a call of the form  $t.f(\dots, s, \dots)$ . This appears in a routine  $r$  of which  $t$  must be a formal argument. When the call is executed, the object T\_OBJ attached to  $t$  is busy (it became busy when the latest call to  $r$  was selected) and "reserved" by  $r$ . The previous discussion still applies; the only new effect is the lazy wait policy, which affects the client executing the call.

```

separate class ELEVATOR feature {BUTTON}
  accept (floor: INTEGER) is
    -- Record and process a request to go to floor.
    do
      record (floor);
      if not moving then
        process_request
      end
    end;

  feature {MOTOR}
    record_stop (floor: INTEGER) is
      -- Record information that elevator has stopped on floor.
      do
        moving := false; position := floor;
        process_request
      end;

  feature {DISPATCHER}
    position: INTEGER; moving: BOOLEAN

  feature {NONE}
    puller: MOTOR;

    pending: QUEUE [INTEGER];
      -- The queue of pending requests
      -- (each identified by the number of the destination floor)

    record (floor: INTEGER) is
      -- Record request to go to floor.
      do
        "Algorithm to insert request for floor into pending"
      end;

    process_request is
      -- Handle next pending request, if any.
      local
        floor: INTEGER
      do
        if not pending.empty then
          floor := pending.item;
          actual_process (puller, floor);
          pending.remove
        end
      end;

    actual_process (m: separate MOTOR; floor: INTEGER) is
      -- Direct m to go to floor.
      do
        moving := true; m.move (floor)
      end
  end

```

**Figure 15.**

**Lazy wait policy on satisfiable calls with separate targets:** Once a satisfiable call has been selected for execution:

- If the call is to a procedure, the client continues its execution unless one of its actual arguments is a nonseparate reference.
- If the call is to a query (function or attribute), or includes a nonseparate reference among its actual arguments, the client's execution waits until the completion of the code.

### Conclusion and Open Issues

This presentation has described an approach to concurrent object-oriented computation, and the rationale that led to it. I consider this article only a first step toward a solution of the problem under study. The following points are open to criticism:

- The dual semantics of assertions on separate and nonseparate targets
- The rule that the target of a separate call must be a formal argument

```

separate class BUTTON feature
  target: INTEGER;
end

separate class CABIN_BUTTON inherit BUTTON feature
  ... Left to the reader ...
end

separate class FLOOR_BUTTON inherit BUTTON feature
  controller: DISPATCHER;

  request is
    -- Send to dispatcher a request to stop on level target.
    do
      actual_request (controller)
    end

  actual_request (d: DISPATCHER) is
    -- Send to d a request to stop on level target.
    do
      d.accept (target)
    end
end
  
```

**Figure 16.**

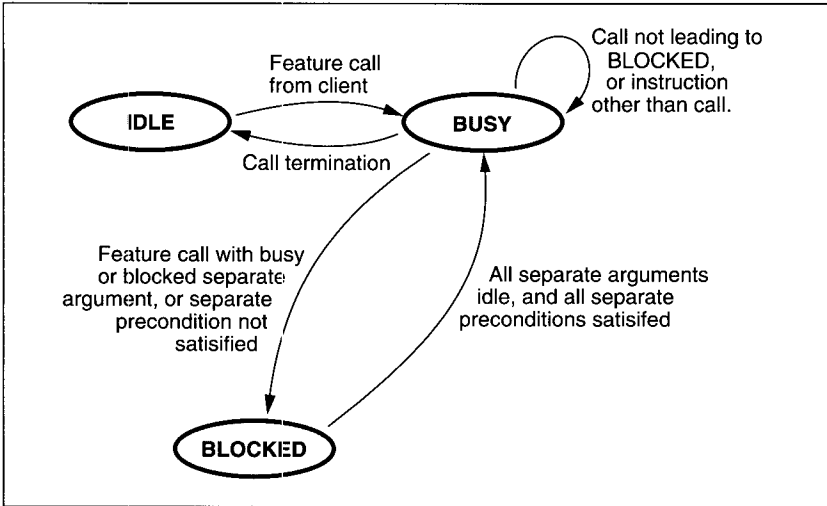
```

separate class DISPATCHER feature {FLOOR_BUTTON}
  accept (floor: INTEGER) is
    -- Handle a request to send an elevator to floor.

    local
      index: INTEGER; chosen: ELEVATOR;
    do
      "A gorithm to determine the elevator
      which should handle the request for floor";
      index := "The index of the chosen elevator";
      chosen := elevators @ index;
      send_request (chosen, floor)
    end
  feature {NONE}
    send_request (e: ELEVATOR; floor: INTEGER) is
      -- Send to e a request to go to floor.
      do e.accept (floor) end;

    elevators: ARRAY [ELEVATOR];
end
  
```

**Figure 17.**



**Figure 19.**

of the enclosing routine, which may require adding apparently superfluous *actual\_xxx* routines


Additionally, much practical and theoretical work remains to be done in the following areas:

- Devising the practical facilities for associating processors, as defined in this article, with physical resources: computers on a network, tasks in an operating system, physical processors in a multiprocessing system.
- Implementing the mechanism in various environments—shared memory, distributed systems and others.
- Exploring the proof rules further and performing proofs of significant systems.
- Studying whether it is possible, by observing certain restrictions in the use of the mechanisms described, to guarantee deadlock avoidance.
- Finding out whether the mechanism can be adapted to ensure fairness, and if so, how (perhaps through library facilities, as was done to achieve the semantics of duels).

One may hope these questions can be answered in a way that preserves the advantages of the concurrency mechanism described here, in particular its minimal departure from the concepts of sequential object-oriented computation and its compatibility with the assertion concepts which are so essential to a proper understanding of this field.

**Acknowledgments**

(Names are included here to acknowledge comments, criticism and help, not to imply endorsement in any way.) Early work leading to this article was presented in March 1990 at the TOOLS EUROPE conference in Paris [12]. The mechanism described here was the topic of a presentation at the TOOLS PACIFIC conference in Sydney in December 1992. I greatly benefited from the work done by John Potter and Ghinwa Jalloul from University of Technology, Sydney (UTS), starting from the original article, and from several discussions with them. (The *hold* construct was their suggestion.) I am also grateful to UTS for the

opportunity to work on the topic of object-oriented concurrency during my stay there from August to October 1992. The design of the lazy wait facility was influenced by the work of Denis Caromel [6]. Important insights were gained at various stages of this work from discussions with Richard Bielak, John Bruno, Steve Cook, Carlo Ghezzi, Peter Löhr, Dino Mandrioli, Jean-Marc Nerson, Robert Switzer and Kim Waldén. 

#### References

1. Agha, G. Concurrent object-oriented programming. *Commun. ACM* 33, 9 (Sept. 1990), 125–141.
2. America, P. Designing an object-oriented programming language with behavioural subtyping. In *Foundations of Object-Oriented Languages*, J.W. de Bakker, W.P. de Roever, and G. Rozenberg, Eds. Springer-Verlag, New York, 1991, pp. 60–90.
3. America, P. and van der Linden, F. A parallel object-oriented language with inheritance and subtyping. In *Proceedings of ECOOP-OOPSLA 90* (Ottawa, Canada, 1990).
4. Ben-Ari, M. *Principles of Concurrent and Distributed Programming*. Prentice-Hall International, Hemel Hempstead, 1991.
5. Birtwistle, G., Dahl, O.-J., Myrhaug, B. and Nygaard, K. *Simula Begin*. Studentlitteratur and Auerbach, 1973.
6. Caromel, D. Concurrency: An object-oriented approach. In *TOOLS 2*, Jean Bézivin et al., Eds., SOL/Angkor, Paris, 1990. 183–197.
7. Caromel, D., Toward a method of object-oriented concurrent programming. *Commun. ACM* 36,9(Sept. 1993).
8. Hailpern, B.T. Verifying concurrent processes using temporal logic. *Lecture Notes in Computer Science 129*, Springer-Verlag, New York, 1982.
9. Hoare, C.A.R. Procedures and parameters: An axiomatic approach. In *Symposium on the Semantics of Programming Languages*, Erwin Engeler, Ed. Springer-Verlag, New York, 1971, pp. 103–116. Reprinted in *Essays in Computing Science*, C.A.R. Hoare and C.B. Jones, Eds. Prentice-Hall International, 1989.
10. Lieberman, H. Concurrent object-oriented programming in Act 1. In *Object-oriented Concurrent Programming*, Akinori Yonezawa and Mario Tokoro, Eds. MIT Press, Cambridge Mass., 1987, pp. 9–36.

```

deferred class TIMED inherit
  CONCURRENCY
feature {NONE}
  failed: BOOLEAN; alarm: WATCHDOG;

  timed_action (t: REAL) is
    -- Execute action, but interrupt after t seconds if not complete.
    -- If interrupted before completion, set failed to true.
    do
      set_alarm (t); action; unset_alarm (t); failed := false
    rescue
      failed := true
    end;

  set_alarm (t: REAL) is
    -- Set alarm to interrupt Current after t seconds.
    do
      -- Create alarm if necessary:
      if alarm = Void then !! alarm end;
      yield; actual_set (alarm, t)
    end;

  unset_alarm (t: REAL) is
    -- Remove the last alarm set.
    do
      insist;
      immediate_service; actual_unset (alarm); normal_service
    end;

  action is
    -- The action to be performed under watchdog control
    deferred
    end

feature {NONE} -- Actual access to watchdog
  ... actual_set and actual_unset, left to the reader ...
feature {WATCHDOG} -- The interrupting operation
  stop is
    -- Empty action to let watchdog interrupt a call to timed_action.
    do end
end

separate class WATCHDOG feature {TIMED}
  set (caller: separate TIMED; t: REAL) is
    -- Set an alarm after t seconds for client caller.
    do
      yield; wait (t);
      immediate_service; caller.stop; normal_service
    rescue
      terminated_before_time := true
    end;

  unset is
    -- Remove alarm.
    do
      terminated_before_time := false
    en

feature {NONE}
  terminated_before_time: BOOLEAN;
end

```

Figure 18.



11. Meyer, B. *Object-oriented software construction*. Prentice-Hall, Englewood Cliffs, N.J., 1988.
12. Meyer, B. Sequential and concurrent object-oriented programming. In *TOOLS 90 (Technology of Object-Oriented Languages and Systems)*, Angkor/SOL, (Paris, June 1990), pp. 17-28.
13. Meyer, B. *Introduction to the Theory of Programming Languages*. Prentice-Hall, Englewood Cliffs, N.J., 1990.
14. Meyer, B. *Eiffel: The Language*. Prentice-Hall, Englewood Cliffs, N.J., 1991.
15. Meyer, B. Applying design by contract. *IEEE Comput.* 40-51, 10 (Oct. 1992).
16. Meyer, B. *Object-Oriented Software Construction, Second ed.* Prentice-Hall, Englewood Cliffs, N.J., 1993.
17. Papatomas, M. Language design rationale and semantic framework for concurrent object-oriented pro-

gramming. Ph.D. dissertation. Université de Genève, 1992.

18. SIS, Standardiseringskommissionen i Sverige (Swedish Standards Institute). Data Processing-Programming Languages—SIMULA, Svensk Standard SS 63 61 14, May 20, 1987.
19. Yonezawa, A. and Tokoro, M., Eds. *Object-oriented Concurrent Programming*. MIT Press, Cambridge, Mass., 1987.

**CR Categories and Subject Descriptors:** D.1.5 [Programming Techniques]: D.1.3 [Software]: Programming Techniques—concurrent programming; Object-Oriented Programming; D.2.2. [Software Engineering]: Tools and techniques—software libraries; H.2.4. [Database Management]: Systems—concurrency

**General Terms:** Design

**Additional Key Words and Phrases:**

Concurrency, Eiffel, object-oriented concurrent programming

**About the Author:**

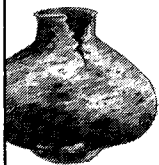
**BERTRAND MEYER** is president of Interactive Software Engineering (Santa Barbara) and SOL (Paris). He is the author of several books devoted to object-oriented software construction, the Eiffel approach, and the theoretical aspects of programming languages. **Author's Present Address:** Interactive Software Engineering, 270 Storke Road, Suite 7, Goleta, CA 93117; email: bertrand@eiffel.com

Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copying is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© ACM 0002-0782/93/0900-056 \$1.50

# IT'S A TRAGEDY BEYOND DESCRIPTION.

96,000 acres of irreplaceable rain forest are being burned every day. These once lush forests are being cleared for grazing and farming. But the tragedy is without the forest this delicate land quickly turns barren.



In the smoldering ashes are the remains of what had taken thousands of years to create. The life-sustaining nutrients of the plants and living matter have been destroyed and the



exposed soil quickly loses its fertility. Wind and rain reap further damage and in as few as five years a land that was teeming with life is turned into a wasteland.

The National Arbor Day Foundation, the world's largest tree-planting environmental organization, has launched Rain Forest Rescue. By joining the Foundation, you will help stop further burning. For the future of our planet, for hungry people everywhere, support Rain Forest Rescue. Call now.



**Call Rain Forest Rescue.  
1-800-255-5500**