

Chapter 5

Design with Patterns

"It is not the strongest of the species that survive, nor the most intelligent, but the one most responsive to change."
—Charles Darwin

"Man has a limited biological capacity for change. When this capacity is overwhelmed, the capacity is in future shock."
—Alvin Toffler

Design patterns are convenient solutions for software design problems commonly employed by expert developers. The power of design patterns derives from reusing proven solution "recipes" from similar problems. In other words, patterns are *codifying* practice rather than *prescribing* practice, or, they are capturing the existing best practices, rather than inventing untried procedures. Patterns are used primarily to *improve existing designs or code by rearranging it according to a "pattern."* By reusing a pattern, the developer gains *efficiency*, by avoiding a lengthy process of trials and errors in search of a solution, and *predictability* because this solution is known to work for a given problem.

Design patterns can be of different complexities and for different purposes. In terms of complexity, the design pattern may be as simple as a naming convention for object methods in the JavaBeans specification (see Chapter 7) or can be a complex description of interactions between the multiple classes, some of which will be reviewed in this chapter. In terms of the purpose, a pattern may be intended to facilitate component-based development and reusability, such as in the JavaBeans specification, or its purpose may be to prescribe the rules for responsibility assignment to the objects in a system, as with the design principles described in Section 2.5.

As pointed earlier, finding effective representation(s) is a recurring theme of software engineering. By condensing many structural and behavioral aspects of the design into a few simple concepts, patterns make it easier for team members to discuss the design. As with any symbolic language, one of the greatest benefits of patterns is in *chunking* the design knowledge. Once team members are familiar with the pattern terminology, the use of this terminology shifts

Contents

5.1 Indirect Communication: Publisher-Subscriber	
5.1.1 Control Flow	
5.1.2 Pub-Sub Pattern Initialization	
5.1.3	
5.1.4	
5.1.5	
5.2 More Patterns	
5.2.1 Command	
5.2.2 Decorator	
5.2.3 State	
5.2.4 Proxy	
5.3 Concurrent Programming	
5.3.1 Threads	
5.3.2 Exclusive Resource Access—Exclusion Synchronization	
5.3.3 Cooperation between Threads—Condition Synchronization	
5.3.4	
5.2.3	
5.4 Broker and Distributed Computing	
5.4.1 Broker Pattern	
5.4.2 Java Remote Method Invocation (RMI)	
5.4.3	
5.4.4	
5.5 Information Security	
5.5.1 Symmetric and Public-Key Cryptosystems	
5.5.2 Cryptographic Algorithms	
5.5.3 Authentication	
5.5.4	
5.6 Summary and Bibliographical Notes	
Problems	

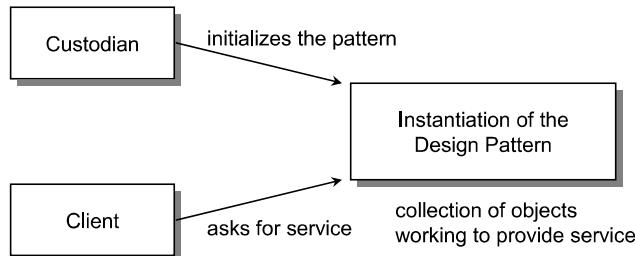


Figure 5-1: The players in a design pattern usage.

the focus to higher-level design concerns. No time is spent in describing the mechanics of the object collaborations because they are condensed into a single pattern name.

This chapter reviews some of the most popular design patterns that will be particularly useful in the rest of the text. What follows is a somewhat broad and liberal interpretation of design patterns. The focus is rather on the techniques of solving specific problems; nonetheless, the “patterns” described below do fit the definition of patterns as recurring solutions. These patterns are conceptual tools that facilitate the development of flexible and adaptive applications as well as reusable software components.

Two important observations are in order. First, finding a name that in one or few words conveys the meaning of a design pattern is very difficult. A similar difficulty is experienced by user interface designers when trying to find graphical icons that convey the meaning of user interface operations. Hence, the reader may find the same or similar software construct under different names by different authors. For example, The *Publisher-Subscriber* design pattern, described in Section 5.1, is most commonly called *Observer* [Gamma et al., 1995], but [Larman, 2005] calls it *Publish-Subscribe*. I prefer the latter because I believe that it conveys better the meaning of the underlying software construct¹. Second, there may be slight variations in what different authors label with the same name. The difference may be due to the particular programming language idiosyncrasies or due to evolution of the pattern over time.

Common players in a design pattern usage are shown in Figure 5-1. A Custodian object assembles and sets up a pattern and cleans up after the pattern’s operation is completed. A client object (can be the same software object as the custodian) needs and uses the services of the pattern. The design patterns reviewed below generally follow this usage “pattern.”

5.1 Indirect Communication: Publisher-Subscriber

¹ The *Publish-Subscribe* moniker has a broader use than presented here and the interested reader should consult [Eugster et al. 2003].

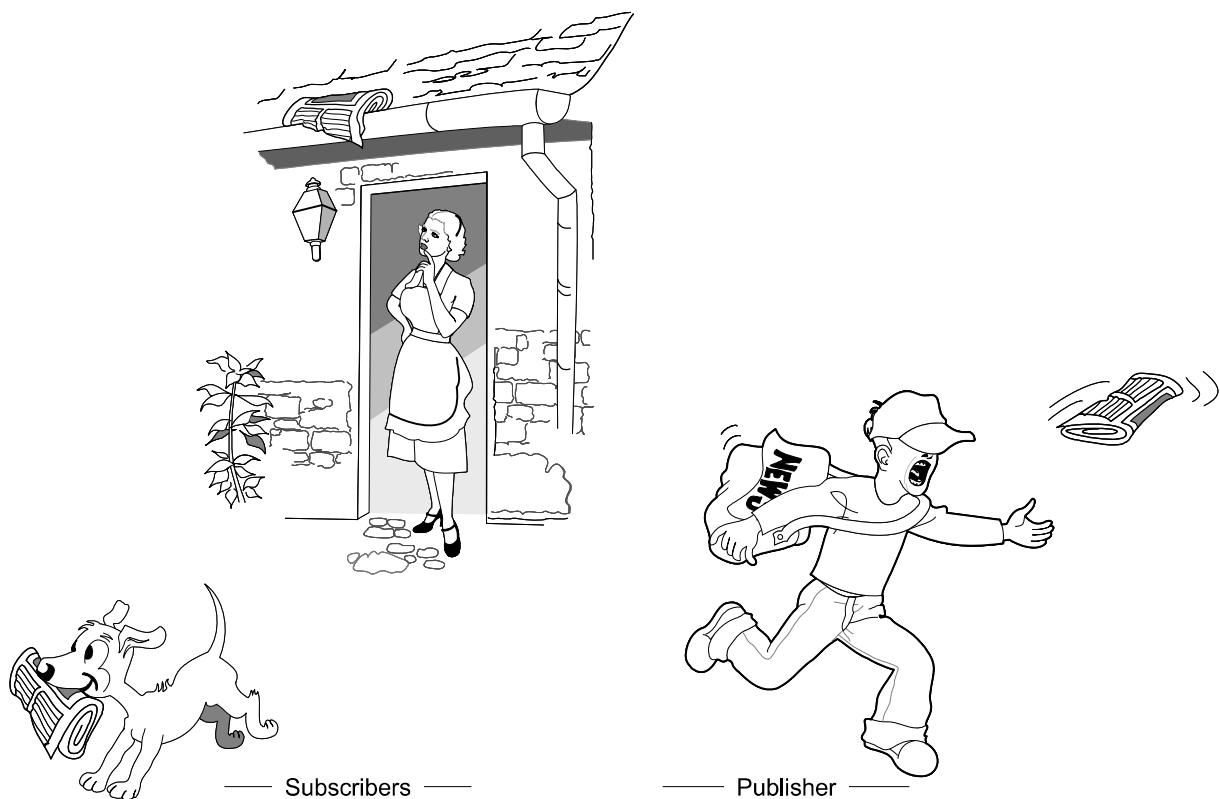


Figure 5-2: The concept of indirect communication in a Publisher/Subscriber system.

"If you find a good solution and become attached to it, the solution may become your next problem."
—Robert Anthony

"More ideas to choose from mean more complexity ... and more opportunities to choose wrongly."
—Vikram Pandit

Publisher-subscriber design pattern (see Figure 5-2) is used to implement *indirect communication* between software objects. Indirect communication is usually used when an object cannot or does not want to know the identity of the object whose method it calls. Another reason may be that it does not want to know what the effect of the call will be. The most popular use of the pub-sub pattern is in building reusable software components.

- 1) Enables building reusable components
- 2) Facilitates separation of the business logic (responsibilities, concerns) of objects

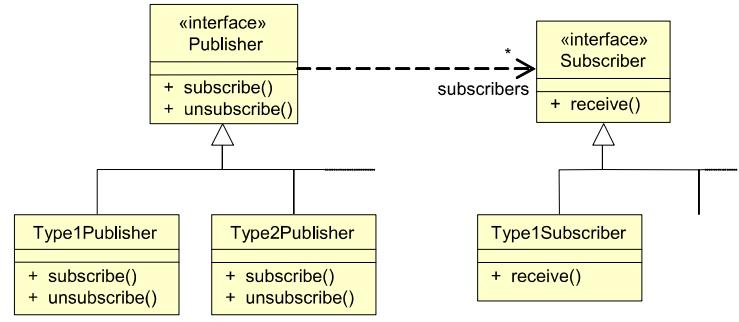
Centralized vs. decentralized execution/program-control method—spreads responsibilities for better balancing. Decentralized control does not necessarily imply concurrent threads of execution.

The problem with building reusable components can be illustrated on our case-study example. Let us assume that we want to reuse the KeyChecker object in an extended version of our case-study application, one that sounds alarm if someone is tampering with the lock. We need to modify the method `unlock()` not only to send message to LockCtrl but also to AlarmCtrl, or to introduce a new method. In either case, we must change the object code, meaning that the object is not reusable as-is.

Publisher
Knowing Responsibilities:
• Knows event source(s)
• Knows interested obj's (subscribers)
Doing Responsibilities:
• Registers/Unregisters subscribers
• Notifies the subscribers of events

Subscriber
Knowing Responsibilities:
• Knows event types of interest
• Knows publisher(s)
Doing Responsibilities:
• Registers/Unregisters with publishers
• Processes received event notifications

(a)



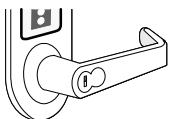
(b)

Figure 5-3: Publisher/Subscriber objects employee cards (a), and the class diagram of their collaborations (b).

Information source acquires information in some way and we assume that this information is important for other objects to do the work they are designed for. Once the source acquires information (becomes “information expert”), it is logical to expect it to pass this information to others and initiate their work. However, this tacitly implies that the source object “knows” what the doer object should do next. This knowledge is encoded in the source object as an “IF-THEN-ELSE” rule and must be modified every time the doer code is modified (as seen earlier in Section 2.5).

Request- vs. event-based communication, Figure 5-4: In the former case, an object makes an explicit request, whereas in the latter, the object expresses interest ahead of time and later gets notified by the information source. In a way, the source is making a method request on the object. Notice also that “request-based” is also synchronous type of communication, whereas event based is asynchronous.

Another way to design the KeyChecker object is to make it become a publisher of events as follows. We need to define two class interfaces: Publisher and Subscriber (see Figure 5-3). The first one, Publisher, allows any object to subscribe for information that it is the source of. The second, Subscriber, has a method, here called `receive()`, to let the Publisher publish the data of interest.



Listing 5-1: Publish-Subscribe class interfaces.

```

public interface Subscriber {
    public void receive(Content content);
}
  
```

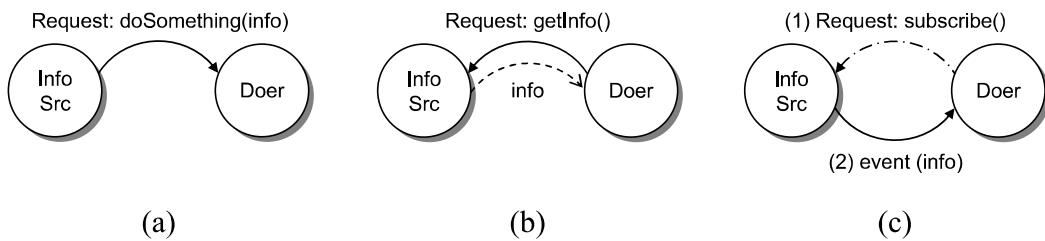


Figure 5-4: Request- vs. event-based communication among objects. (a) Direct request—information source controls the activity of the doer. (b) Direct request—the doer controls its own activity, information source is only for lookup, but doer must know when is the information ready and available. (c) Indirect request—the doer controls its own activity and does not need to worry when the information is ready and available—it gets prompted by the information source.

```

import java.util.ArrayList;
public class Content {
    public Publisher source_;
    public ArrayList data_;

    public Content(Publisher src, ArrayList dat) {
        source_ = src;
        data_ = (ArrayList) dat.clone(); // for write safety...
    } // ...avoid aliasing and create a new copy
}

public interface Publisher {
    public subscribe(Subscriber subscriber);
    public unsubscribe(Subscriber subscriber);
}

```

A `Content` object contains only data, no business logic, and is meant to transfer data from `Publisher` to `Subscriber`. The actual classes then implement those two interfaces. In our example, the key `Checker` object would then implement the `Publisher`, while `DeviceCtrl` would implement the `Subscriber`.

Listing 5-2: Refactored the case-study code of using the Publisher-Subscriber design pattern. Here, the class `DeviceCtrl` implements the `Subscriber` interface and the class `Checker` implements the `Publisher` interface.

```

public class DeviceCtrl implements Subscriber {
    protected LightBulb bulb_;
    protected PhotoSObs sensor_;

    public DeviceCtrl(Publisher keyChecker, PhotoSObs sensor, ...) {
        sensor_ = sensor;
        keyChecker.subscribe(this);
        ...
    }
}

```

```
public void receive(Content content) {
    if (content.source_ instanceof Checker) {
        if ( ((String)content.data_).equals("valid") ) {
            // check the time of day; if daylight, do nothing
            if (!sensor_.isDaylight()) bulb_.setLit(true);
        }
    } else (check for another source of the event ...) {
        ...
    }
}

import java.util.ArrayList;
import java.util.Iterator;

public class Checker implements Publisher {
    protected KeyStorage validKeys_;
    protected ArrayList subscribers_ = new ArrayList();

    public Checker( ... ) { }

    public subscribe(Subscriber subscriber) {
        subscribers_.add(subscriber); // could check whether this
        }                                // subscriber already subscribed

    public unsubscribe(Subscriber subscriber) {
        int idx = subscribers_.indexOf(subscriber);
        if (idx != -1) { subscribers_.remove(idx); }
    }

    public void checkKey(Key user_key) {
        boolean valid = false;
        ... // verify the user key against the "validKeys_" database

        // notify the subscribers
        Content cnt = new Content(this, new ArrayList());

        if (valid) { // authorized user
            cnt.data.add("valid");
        } else { // the lock is being tampered with
            cnt.data.add("invalid");
        }
        cnt.data.add(key);

        for (Iterator e = subscribers_.iterator(); e.hasNext(); ) {
            ((Subscriber) e.next()).receive(cnt);
        }
    }
}
```

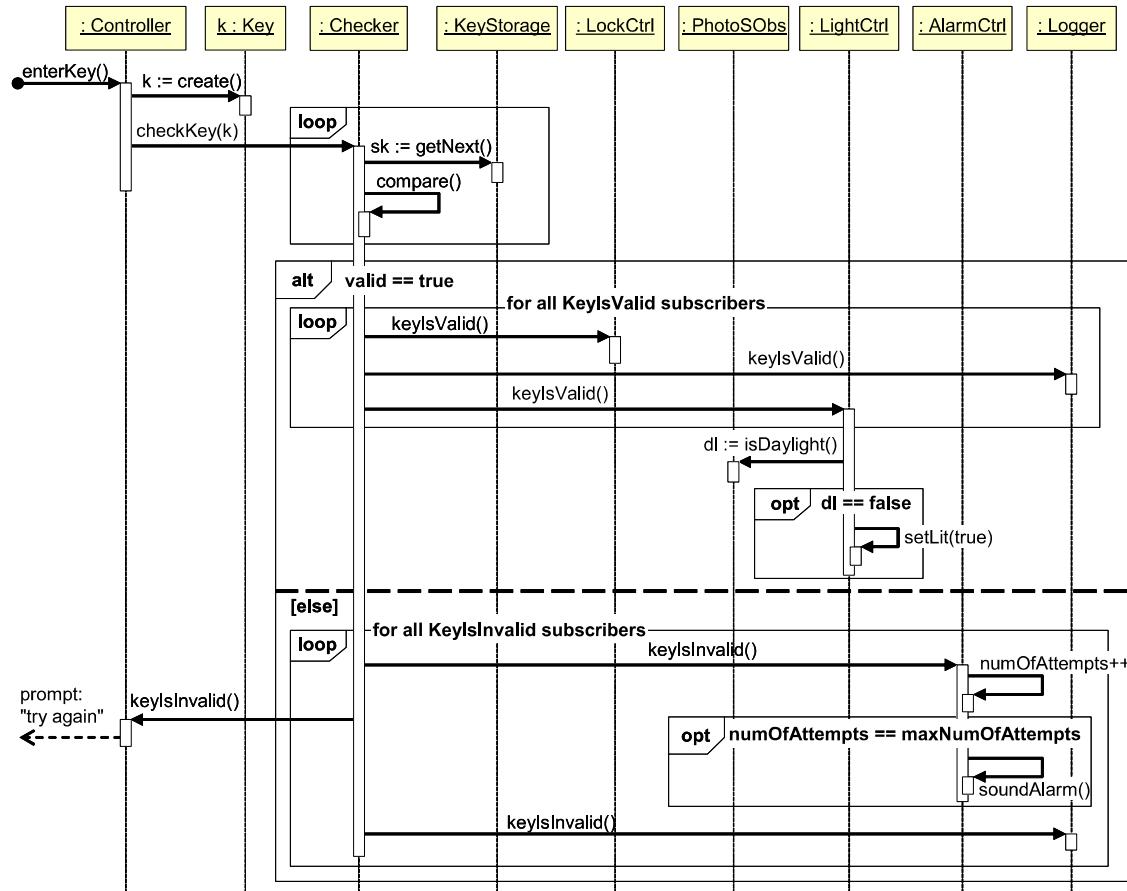


Figure 5-5: Sequence diagram for publish-subscribe version of the use case “Unlock.” Compare this with Figure 2-27.

A Subscriber may be subscribed to several sources of data and each source may provide several types of content. Thus, the Subscriber must determine the source and the content type before it takes any action. If a Subscriber gets subscribed to many sources which publish different content, the Subscriber code may become quite complex and difficult to manage. The Subscriber would contain many `if()` or `switch()` statements to account for different options. A more object-oriented solution for this is to use class *polymorphism*—instead of having one Subscriber, we should have several Subscribers, each specialized for a particular source. The Subscribers may also have more than one `receive()` method, each specialized for a particular data type. Here is an example. We could implement a Switch by inheriting from the generic Subscriber interface defined above, or we can define new interfaces specialized for our problem domain.

Listing 5-3: Subscriber interfaces for “key-is-valid” and “key-is-invalid” events.

```
public interface KeyIsValidSubscriber {
    public void keyIsValid(LockEvent event); // receive() method
}

public interface KeyIsInvalidSubscriber {
    public void keyIsInvalid(LockEvent event); // receive() method
}
```

```
}
```

The new design for the Unlock use case is shown in Figure 5-5, and the corresponding code might look as shown next. Notice that here the attribute numAttempts belongs to the AlarmCtrl, unlike the first implementation in Listing 2-2 (Section 2.7), where it belonged to the Controller. Notice also that the Controller is a KeyIsInvalidSubscriber so it can prompt the user to enter a new key if the previous attempt was unsuccessful.

Listing 5-4: A variation of the Publisher-Subscriber design from Listing 5-2 using the subscriber interfaces from Listing 5-3.

```
public class Checker implements LockPublisher {
    protected KeyStorage validKeys_;
    protected ArrayList keyValidSubscribers_ = new ArrayList();
    protected ArrayList keyInvalidSubscribers_ = new ArrayList();

    public Checker(KeyStorage ks) { validKeys_ = ks; }

    public void subscribeKeyIsValid(KeyIsValidSubscriber sub) {
        keyValidSubscribers_.add(sub);
    }

    public void subscribeKeyIsInvalid(KeyIsInvalidSubscriber sub) {
        keyInvalidSubscribers_.add(sub);
    }

    public void checkKey(Key user_key) {
        boolean valid = false;
        ... // verify the key against the database

        // notify the subscribers
        LockEvent evt = new LockEvent(this, new ArrayList());
        evt.data.add(key);

        if (valid) {
            for (Iterator e = keyValidSubscribers_.iterator();
                 e.hasNext(); ) {
                ((KeyIsValidSubscriber) e.next()).keyIsValid(evt);
            }
        } else { // the lock is being tampered with
            for (Iterator e = keyInvalidSubscribers_.iterator();
                 e.hasNext(); ) {
                ((KeyIsInvalidSubscriber) e.next()).keyIsInvalid(evt);
            }
        }
    }
}

public class DeviceCtrl implements KeyIsValidSubscriber {
    protected LightBulb bulb_;
    protected PhotoSObs photoObserver_;

    public DeviceCtrl(LockPublisher keyChecker, PhotoSObs sensor, .. )
    {
        photoObserver_ = sensor;
    }
}
```

```

        keyChecker.subscribeKeyIsValid(this);
        ...
    }

    public void keyIsValid(LockEvent event) {
        if (!photoObserver_.isDaylight()) bulb_.setLit(true);
    }
}

public class AlarmCtrl implements KeyIsInvalidSubscriber {
    public static final long maxNumOfAttempts_ = 3;
    public static final long interAttemptInterval_ = 300000; //millisec
    protected long numAttempts_ = 0;
    protected long lastTimeAttempt_ = 0;

    public AlarmCtrl(LockPublisher keyChecker, ...) {
        keyChecker.subscribeKeyIsInvalid(this);
        ...
    }

    public void keyIsInvalid(LockEvent event) {
        long currTime = System.currentTimeMillis();
        if ((currTime - lastTimeAttempt_) < interAttemptInterval_) {
            if (++numAttempts_ >= maxNumOfAttempts_) {
                soundAlarm();
                numAttempts_ = 0; // reset for the next user
            }
        } else { // this must be a new user's first mistake ...
            numAttempts_ = 1;
        }
        lastTimeAttempt_ = currTime;
    }
}

```

It is of note that what we just did with the original design for the Unlock use case can be considered refactoring. In software engineering, the term *refactoring* is often used to describe modifying the design and/or implementation of a software module without changing its external behavior, and is sometimes informally referred to as “cleaning it up.” Refactoring is often practiced as part of the software development cycle: developers alternate between adding new tests and functionality and refactoring the code to improve its internal consistency and clarity. In our case, the design from Figure 2-27 has been transformed to the design in Figure 5-5.

There is a tradeoff between the number of `receive()` methods and the `switch()` statements. On one hand, having a long `switch()` statement complicates the Subscriber’s code and makes it difficult to maintain and reuse. On the other hand, having too many `receive()` statements results in a long class interface, difficult to read and represent graphically.

5.1.1 Applications of Publisher-Subscriber

The Publisher-Subscriber design pattern is used in the Java AWT and Swing toolkits for notification of the GUI interface components about user generated events. (This pattern in Java is known as *Source-Listener* or *delegation event model*, see Chapter 7.)

One of the main reasons for software components is easy visualization in integrated development environments (IDEs), so the developer can visually assemble the components. The components are represented as “integrated circuits” in analogy to hardware design, and different `receive()` / `subscribe()` methods represent “pins” on the circuit. If a component has too many pins, it becomes difficult to visualize, and generates too many “wires” in the “blueprint.” The situation is similar to determining the right number of pins on an integrated circuit. (See more about software components in Chapter 7.)

Here I reiterate the key benefits of using the pub-sub design pattern and indirect communication in general:

- The components do not need to know each other’s identity. For example, in the sample code given in Listing 1-1 (Section 1.4.2), `LockCtrl` maintains a reference to a `LightCtrl` object.
- The component’s business logic is contained within the component alone. In the same example, `LockCtrl` explicitly invokes the `LightCtrl`’s method `setLit()`, meaning that it minds `LightCtrl`’s business. In the worst case, even the checking of the time-of-day may be delegated to `LockCtrl` in order to decide when to turn the light on.

Both of the above form the basis for component reusability, because making a component independent of others makes it reusable. The pub-sub pattern is the most basic pattern for reusable software components as will be discussed in Chapter 7.

In the “ideal” case, all objects could be made self-contained and thus reusable by applying the pub-sub design pattern. However, there are penalties to pay. As visible from the examples above, indirect communication requires much more code, which results in increased demand for memory and decreased performance. Thus, if it is not likely that a component will need to be reused or if performance is critical, direct communication should be applied and the pub-sub pattern should be avoided.

When to apply the pub-sub pattern? The answer depends on whether you anticipate that the component is likely to be reused in future projects. If yes, apply pub-sub. You should understand that decoupled objects are independent, therefore reusable and easier to understand, while highly interleaved objects provide fast inter-object communication and compact code. Decoupled objects are better suited for global understanding, whereas interleaved objects are better suited for local understanding. Of course, in a large system, global understanding matters more.

5.1.2 Control Flow

Figure 5-6 highlights the difference in control flow for direct and indirect communication types. In the former case, the control is centralized and all flows emanate from the Controller. In the latter case, the control is decentralized, and it is passed as a token around, cascading from object to object. These diagrams also show the *dynamic (behavioral) architecture* of the system.

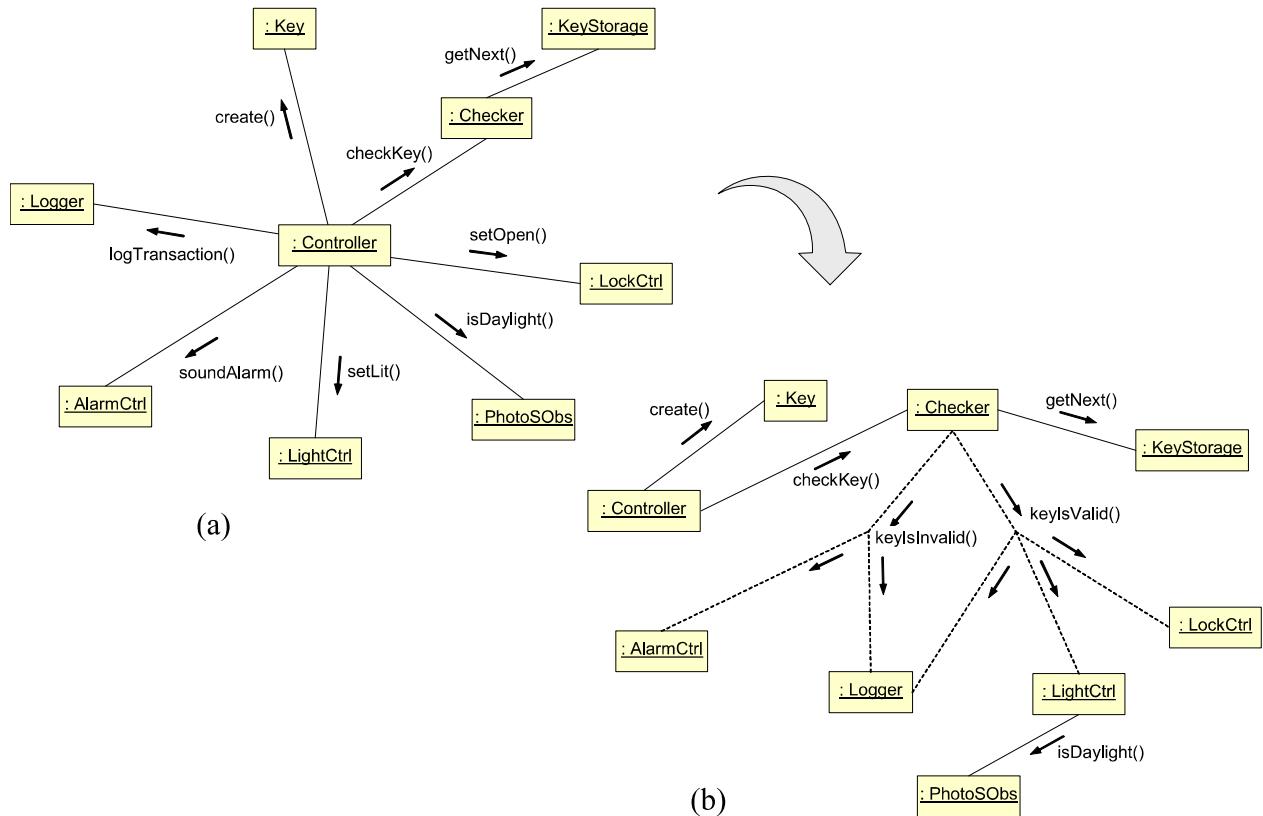


Figure 5-6: Flow control without (a) and with the Pub-Sub pattern (b). Notice that these UML communication diagrams are redrawn from Figure 2-27 and Figure 5-5, respectively.

Although in Figure 5-6(b) it appears as if the Checker plays a central role, this is not so because it is not “aware” of being assigned such a role, i.e., unlike the Controller from Figure 5-6(a), this Checker does not encode the requisite knowledge to play such a role. The outgoing method calls are shown in dashed lines to indicate that these are indirect calls, through the Subscriber interface.

Whatever the rules of behavior are stored in one Controller or distributed (cascading) around in many objects, the *output* (seen from outside of the system) is the same. Organization (internal function) matters only if it simplifies the software maintenance and upgrading.

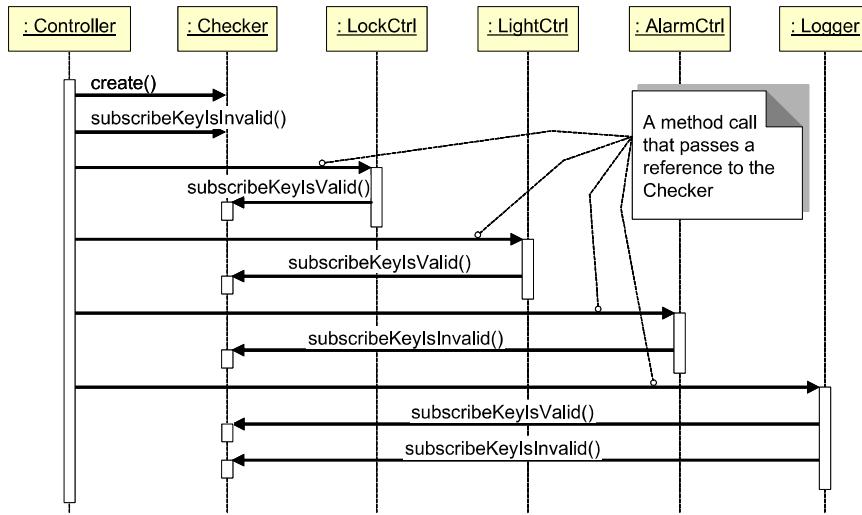


Figure 5-7: Initialization of the pub-sub for the lock control example.

5.1.3 Pub-Sub Pattern Initialization

Note that the “setup” part of the pattern, example shown in Figure 5-7, plays a major, but often ignored, role in the pattern. It essentially represents the master plan of solving the problem using the publish-subscribe pattern and indirect communication.

Most programs are not equipped to split hard problems into parts and then use divide-and-conquer methods. Few programs, too, represent their goals, except perhaps as comments in their source codes. However, a class of programs, called General Problem Solver (GPS), was developed in 1960s by Allen Newell, Herbert Simon, and collaborators, which did have explicit goals and subgoals and solved some significant problems [Newell & Simon, 1962].

I propose that goal representation in object-oriented programs be implemented in the setup part of the program, which then can act at any time during the execution (not only at the initialization) to “rewire” the object relationships.

5.2 More Patterns

Publisher-Subscriber belongs to the category of behavioral design patterns. *Behavioral patterns* separate the interdependent behavior of objects from the objects themselves, or stated differently, they separate functionality from the object to which the functionality applies. This promotes reuse, because different types of functionality can be applied to the same object, as needed. Here I review *Command* as another behavioral pattern.

Another category is structural patterns. An example structural pattern reviewed later is *Proxy*.

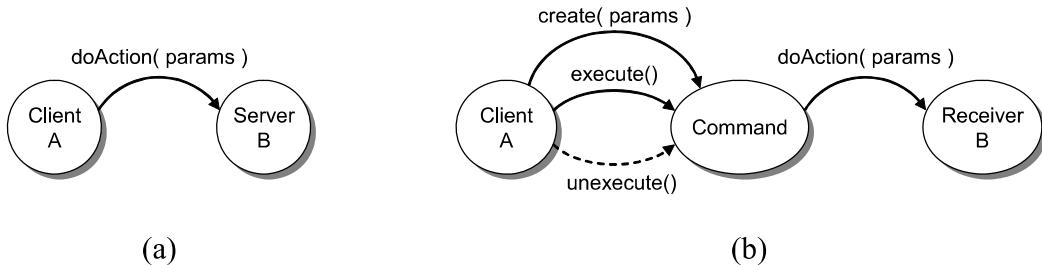


Figure 5-8: Command pattern interposes Command (and other) objects between a client and a server object. Complex actions about rolling back and forward the execution history are delegated to the Command, away from the client object.

A common drawback of design patterns, particularly behavioral patterns, is that we are replacing what would be a single method call with many method calls. This results in performance penalties, which in certain cases may not be acceptable. However, in most cases the benefits of good design outweigh the performance drawbacks.

5.2.1 Command

Objects invoke methods on other objects as depicted in Figure 1-22, which is abstracted in Figure 5-8(a). The need for the Command pattern arises if the invoking object (client) needs to reverse the effect of a previous method invocation. Another reason is the ability to trace the course of the system operation. For example, we may need to keep track of financial transactions for legal or auditing reasons. The purpose of the Command pattern is to delegate the functionality associated with rolling back the server object's state and logging the history of the system operation away from the client object to the Command object, see Figure 5-8(b).

Instead of directly invoking a method on the Receiver (server object), the client object appoints a Command for this task. The Command pattern (Figure 5-9) encapsulates an action or processing task into an object thus increasing flexibility in calling for a service. Command *represents operations as classes* and is used whenever a method call alone is not sufficient. The Command object is the central player in the Command pattern, but as with most patterns, it needs other objects to assist with accomplishing the task. At runtime, a control is passed to the `execute()` method of a non-abstract-class object derived from Command.

Figure 5-9(c) shows a sequence diagram on how to create and execute a command. In addition to executing requests, we may need to be able to trace the course of the system operation. For example, we may need to keep track of financial transactions for legal or auditing reasons. CommandHistory maintains history log of Commands in linear sequence of their execution.

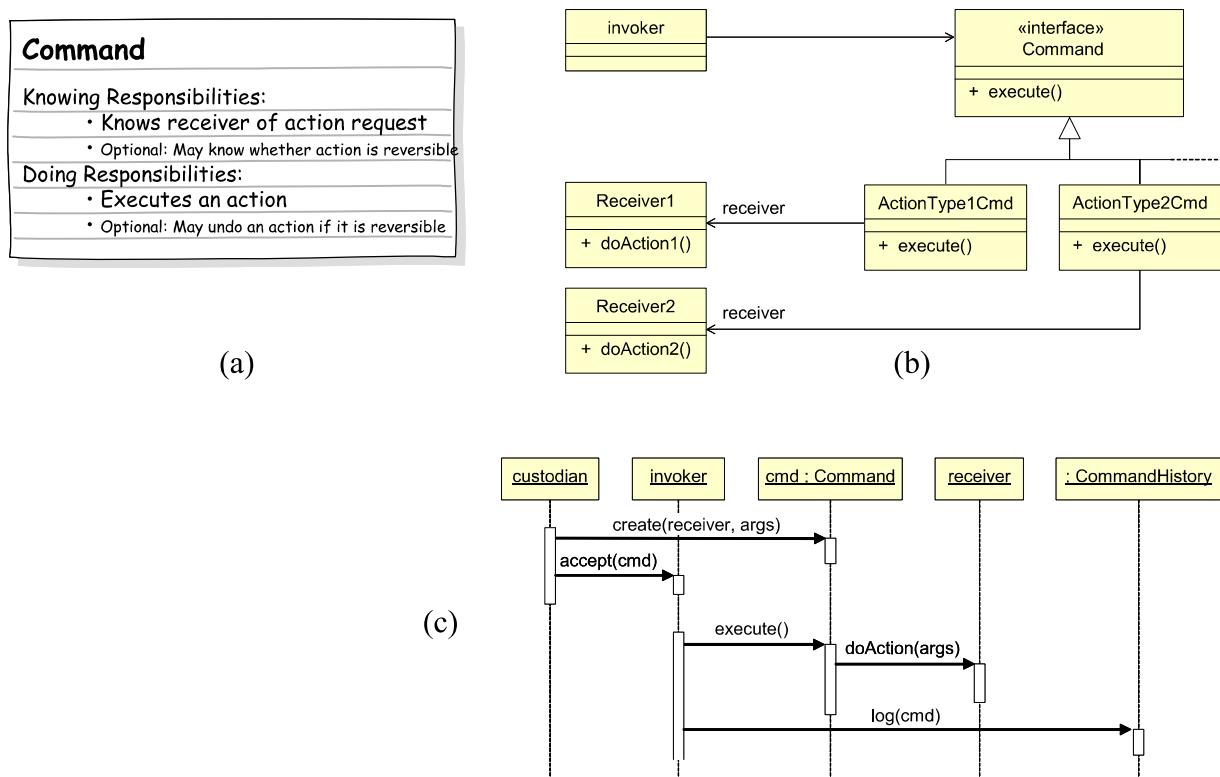


Figure 5-9: (a) Command object employee card. (b) The Command design pattern (class diagram). The base **Command class is an interface implemented by concrete commands. (c) Interaction diagram for creating and executing a command.**

It is common to use Command pattern in operating across the Internet. For example, suppose that client code needs to make a function call on an object of a class residing on a remote server. It is not possible for the client code to make an ordinary method call on this object because the remote object cannot appear in the usual compile-execute process. It is also difficult to employ remote method invocation (Section 5.4.2) here because we often cannot program the client and server at the same time, or they may be programmed by different parties. Instead, the call is made from the client by pointing the browser to the file containing the servlet (a server-side software component). The servlet then calls its method `service(HttpServletRequest, HttpServletResponse)`. The object `HttpServletRequest` includes all the information that a method invocation requires, such as the argument values, obtained from the “environment” variables at standardized global locations. The object `HttpServletResponse` carries the result of invoking `service()`. This technique embodies the basic idea of the Command design pattern. (See also Listing 5-5.)

Web services allow a similar runtime function discovery and invocation, as will be seen in Chapter 8.

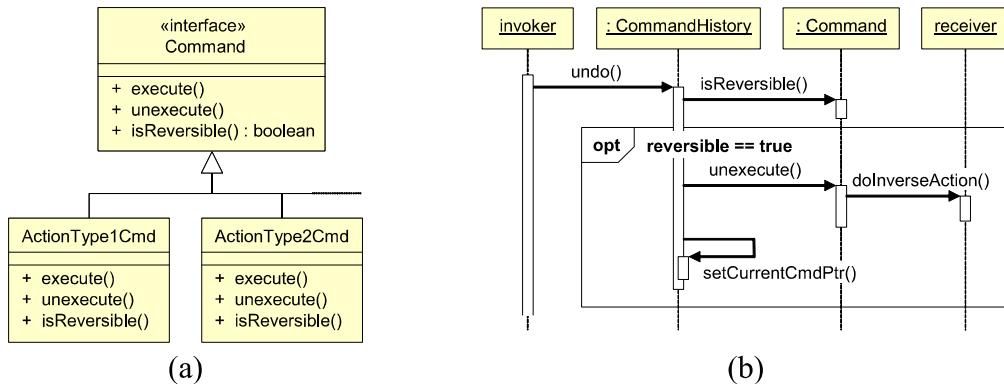


Figure 5-10: (a) Class diagram for commands that can be undone. (b) Interaction diagram for undoing a (reversible) command. Compare to Figure 5-9.

Undo/Redo

The Command pattern may optionally be able to support *rollback* of user's actions in an elegant fashion. Anyone who uses computers appreciates the value of being able to undo their recent actions. Of course, this feature assumes that a command's effect can be reversed. In this case, the Command interface would have two more operations (Figure 5-10(a)): `isReversible()` to allow the invoker to find out whether this command can be undone; and `unexecute()` to undo the effects of a previous `execute()` operation.

Figure 5-10(b) shows a sequence diagram on how to undo/redo a command, assuming that it is undoable. Observe also that `CommandHistory` should decrement its pointer of the current command every time a command is undone and increments it every time a command is redone. An additional requirement on `CommandHistory` is to manage properly the undo/redo caches. For example, if the user backs up along the undo queue and then executes a new command, the whole redo cache should be flushed. Similarly, upon a context switching, both undo/redo caches should be flushed. Obviously, this does not provide for long-term archiving of the commands; if that is required, the archive should be maintained independently of the undo/redo caches.

In physical world, actions are never reversible (because of the laws of thermodynamics). Even an approximate reversibility may not be realistic to expect. Consider a simple light switch. One might think that turning the switch off is exactly opposite of turning it on. Therefore, we could implement a request to turn the switch off as an undo operation of the command to turn the switch on. Unfortunately, this may not be true. For example, beyond the inability to recover the energy lost during the period that the switch was on, it may also happen that the light bulb is burnt. Obviously, this cannot be undone (unless the system has a means of automatically replacing a burnt light bulb with a new one ☺).

In digital world, if the previous state is stored or is easy to compute, then the command can be undone. Even here we need to beware of potential error accumulation. If a number is repeatedly divided and then multiplied by another number, rounding errors or limited number of bits for number representation may yield a different number than the one we started with.

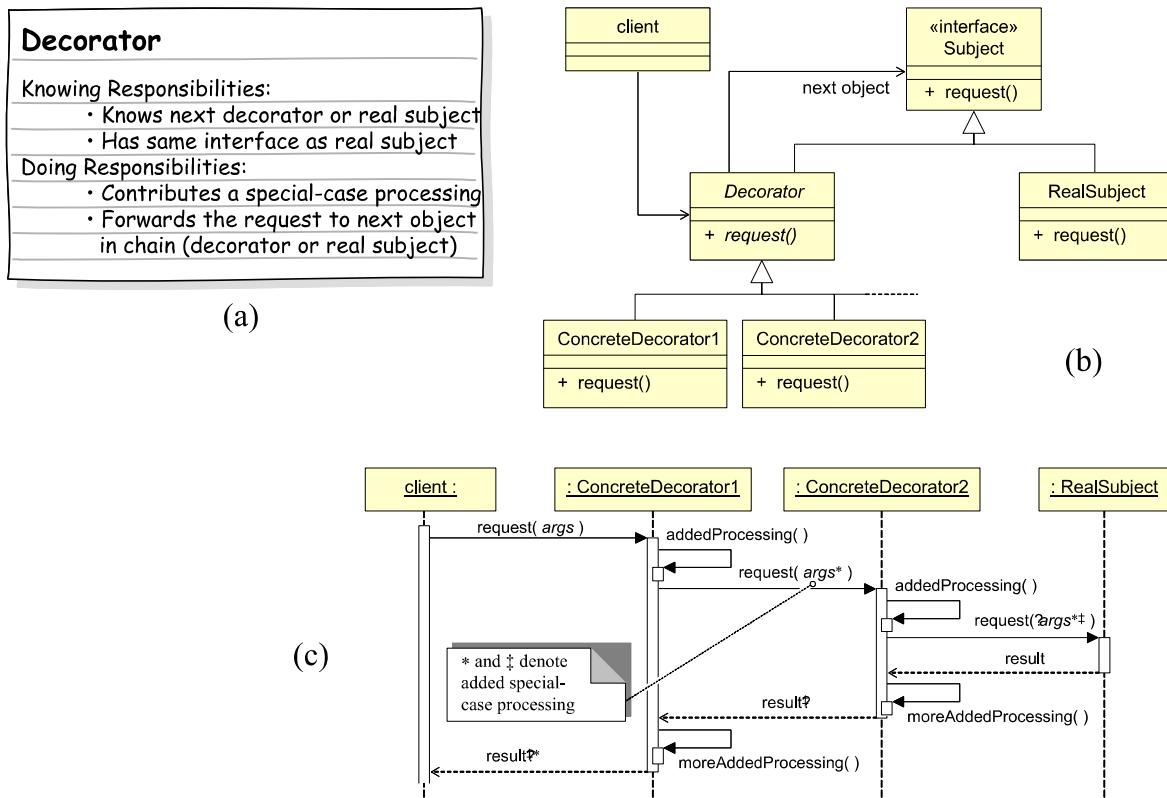


Figure 5-11: (a) Decorator object employee card. (b) The Decorator design pattern (class diagram). (c) Interaction diagram for the Decorator pattern.

5.2.2 Decorator

The Decorator pattern is used to add non-essential behavior to key objects in a software design.

The embellished class (or, *decoratee*) is wrapped up by an arbitrary number of Decorator classes, which provide special-case behaviors (embellishments).

Figure 5-11

Notice that the Decorator is an abstract class (the class and method names are italicized). The reason for this choice is to collect the common things from all different decorators into a base decorator class. In this case, the Decorator class will contain a reference to the next decorator. The decorators are linked in a chain. The client has a reference to the start of the chain and the chain is terminated by the real subject. Figure 5-11(c) illustrates how a request from the client propagates forward through the chain until it reaches the real subject, and how the result propagates back.

To decide whether you need to introduce Decorator, look for special-case behaviors (embellishment logic) in your design.

Consider the following example, where we wish to implement the code that will allow the user to configure the settings for controlling the household devices when the doors are unlocked or locked. The corresponding user interface is shown in Figure 2-2 (Section 2.2). Figure 5-12 and

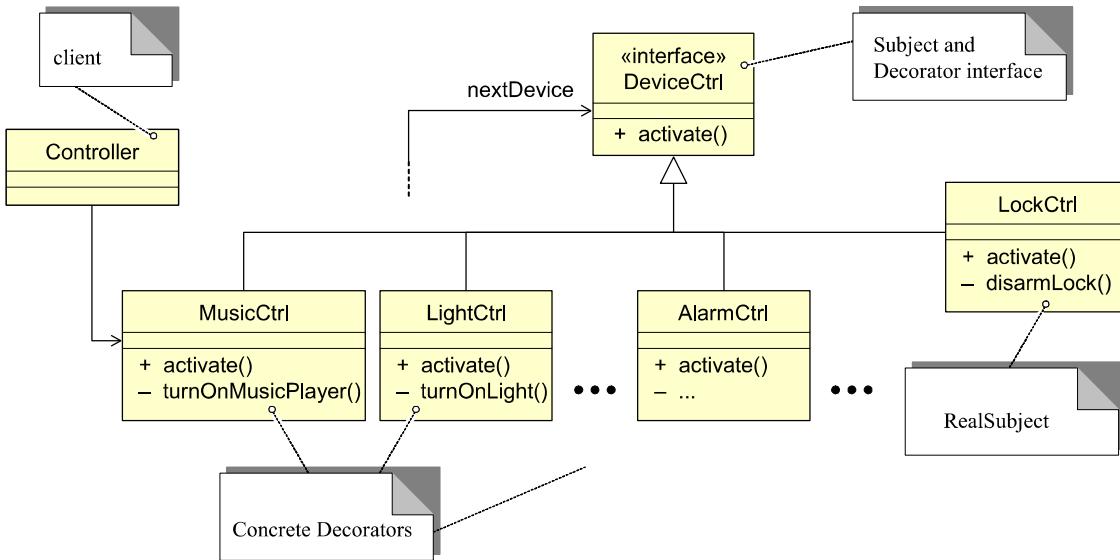


Figure 5-12: Example Decorator class diagram, for implementing the interface in Figure 2-2.

Figure 5-13 show UML diagrams that use the Decorator design pattern in solving this problem. Notice the slight differences in the class diagrams in Figure 5-11(b) and Figure 5-12. As already pointed out, the actual pattern implementation will not always strictly adhere to its generic prototype.

In this example, the decorating functionalities could be added before or after the main function, which is to activate the lock control. For example, in Figure 5-13 the decorating operation `LightCtrl.turnOnLight()` is added before `LockCtrl.activate()`, but `MusicCtrl.turnOnMediaPlayer()` is added after it. In this case all of these operations are commutative and can be executed in any order. This may not always be the case with the decorating functionalities.

5.2.3 State

The State design pattern is usually used when an object's behavior depends on its state in a complex way. In this case, the *state* determines a *mode* of operation. Recall that the *state* of a software object is represented by the current values of its attributes. The State pattern externalizes the relevant attributes into a State object, and this State object has the responsibility of managing the state transitions of the original object. The original object is called “Context” and its attributes are externalized into a State object (Figure 5-14).

A familiar example of object’s state determining its mode of operation includes tools in document editors. Desktop computers normally have only keyboard and mouse as interaction devices. To enable different manipulations of document objects, the document needs to be put in a proper state or mode of operation. That is why we select a proper “tool” in a toolbar before performing a manipulation. The selected tool sets the document state. Consider an example of a graphics editor, such as Microsoft PowerPoint. When the user clicks the mouse pointer on a graphical object and drags the mouse, what will happen depends on the currently selected tool. The default

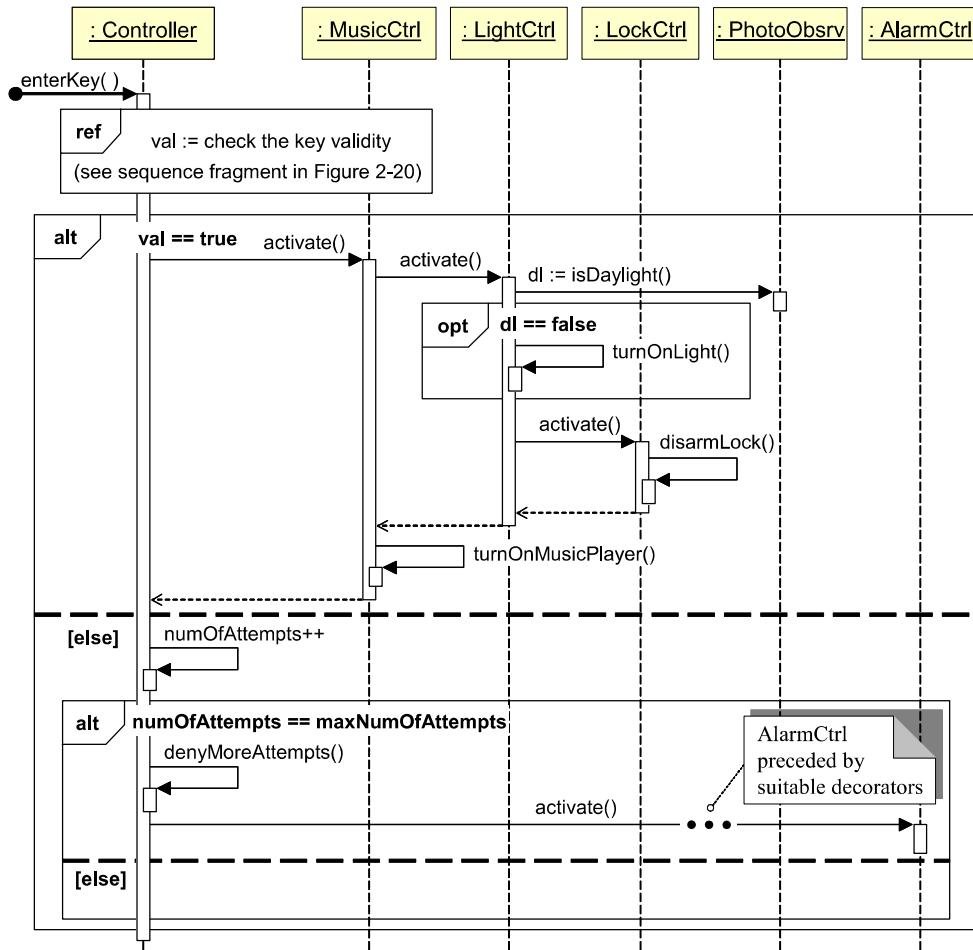


Figure 5-13: Decorator sequence diagram for the class diagram in Figure 5-12.

tool will relocate the object to a new location; the rotation tool will rotate the object for an angle proportional to the distance the mouse is dragged over; etc. Notice that the same action (mouse click and drag) causes different behaviors, depending on the document state (i.e., the currently selected tool).

The State pattern is also useful when an object implements complex conditional logic for changing its state (i.e., the values of this object's attributes). We say that the object is *transitioning from one state* (one set of attribute values) *to another state* (another set of attribute values). To simplify the state transitioning, we define a State interface and different classes that implement this interface correspond to different states of the Context object (Figure 5-14(b)).

Each concrete State class implements the behavior of the Context associated with the state implemented by this State class. The behavior includes calculating the new state of the Context. Because specific attribute values are encapsulated in different concrete states, the current State class just determines the next state and returns it to the Context. Let us assume that the UML state diagram for the Context class is represented by the example in Figure 5-14(c). As shown in Figure 5-14(d), when the Context receives a method call `request()` to handle an event, it calls the method `handle()` on its `currentState`. The current state processes the event and

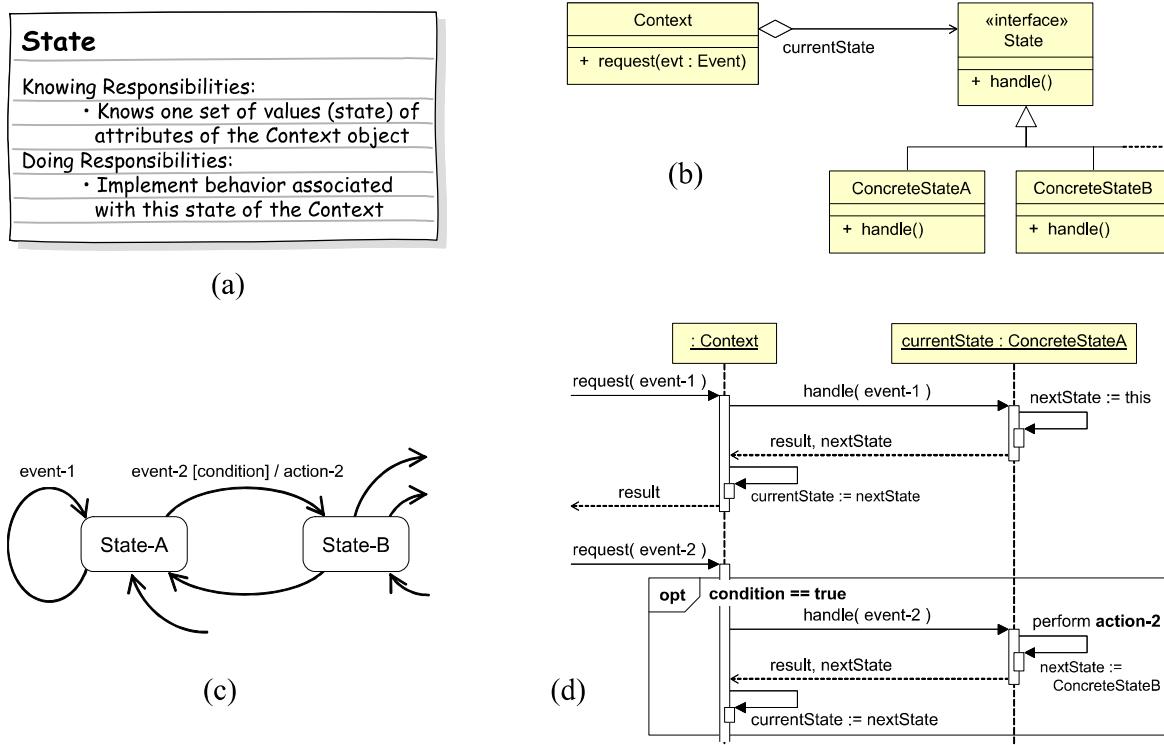


Figure 5-14: (a) State object employee card. (b) The State design pattern (class diagram). (c) Example state diagram for the Context object. (d) Interaction diagram for the state diagram in (c).

performs any action associated with the current state transition. Finally, it returns the next state to the caller Context object. The Context sets this next state as the current state and the next request will be handled by the new current state.

5.2.4 Proxy

The **Proxy pattern** is used to manage or control access to an object. Proxy is needed when the *logistics* of accessing the subject's services is overly complex and comparable or greater in size than that of client's primary responsibility. In such cases, we introduce a helper object (called "proxy") for management of the subject invocation. A Proxy object is a surrogate that acts as a stand-in for the actual subject, and controls or enhances the access to it (Figure 5-15). The proxy object forwards requests to the subject when appropriate, depending on whether the constraint of the proxy is satisfied.

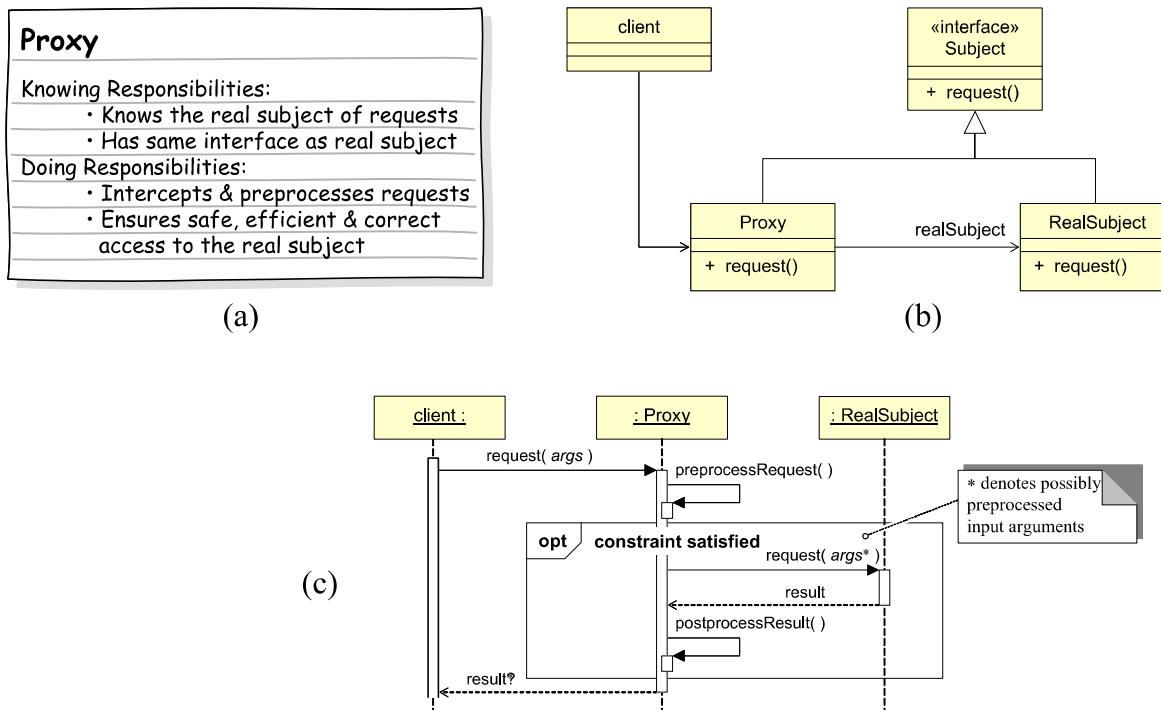


Figure 5-15: (a) Proxy object employee card. (b) The Proxy design pattern (class diagram). (c) Interaction diagram for the Proxy pattern.

The causes of access complexity and the associated constraints include:

- The subject is located in a remote address space, e.g., on a remote host, in which case the invocation (sending messages to it) requires following complex networking protocols.
Solution: use the **Remote Proxy** pattern for crossing the barrier between different memory spaces
- Different access policies constrain the access to the subject. Security policies require that access is provided only to the authorized clients, filtering out others. Safety policies may impose an upper limit on the number of simultaneous accesses to the subject.
Solution: use the **Protection Proxy** pattern for additional housekeeping
- Deferred instantiation of the subject, to speed up the performance (provided that its full functionality may not be immediately necessary). For example, a graphics editor can be started faster if the graphical elements outside the initial view are not loaded until they are needed; only if and when the user changes the viewpoint, the missing graphics will be loaded. Graphical proxies make this process transparent for the rest of the program.
Solution: use the **Virtual Proxy** pattern for optimization in object creation

In essence we could say that proxy allows client objects to cross a barrier to server objects (or, “subjects”). The barrier may be physical (such as network between the client and server computers) or imposed (such as security policies to prevent unauthorized access). As a result, the client cannot or should not access the server by a simple method call as when the barrier does not exist. The additional functionality needed to cross the barrier is extraneous to the client’s business logic. The proxy object abstracts the details of the logistics of accessing the subject’s services

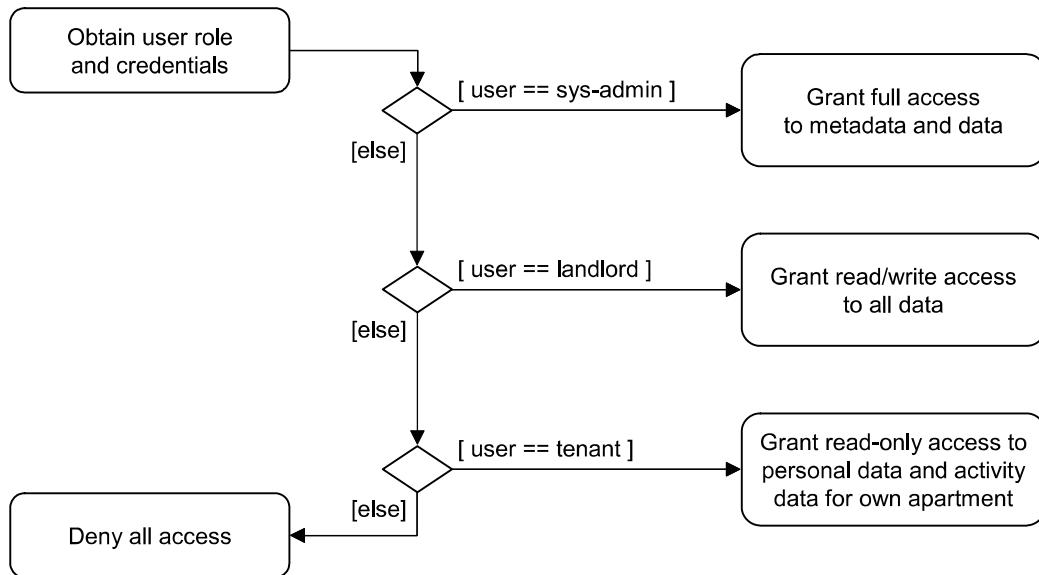


Figure 5-16: Conditional logic for controlling access to the database of the secure home access system.

across different barriers. It does this transparently, so the client has an illusion it is directly communicating with the subject, and does not know that there is a barrier in the middle.

Proxy offers the same interface (set of methods and their signatures) as the real subject and ensures correct access to the real subject. For this reason, the proxy maintains a reference to the real subject (Figure 5-15(b)). Because of the identical interface, the client does not need to change its calling behavior and syntax from that which it would use if there were no barrier involved.

The Remote Proxy pattern will be incorporated into a more complex Broker pattern (Section 5.4). The rest of this section provides more detail on the Protection Proxy.

Protection Proxy

The Protection Proxy pattern can be used to implement different policies to constrain the access to the subject. For example, a security policy may require that a defined service should be seen differently by clients with different privileges. This pattern helps us customize the access, instead of using conditional logic to control the service access. In other words, it is applicable where a subset of capabilities or partial capability should be made available to different actors, based on their roles and privileges.

For example, consider our case study system for secure home access. The sequence diagram for use case UC-5: Inspect Access History is shown in Figure 2-26. Before the Controller calls the method `accessList := retrieve(params : string)` on Database Connection, the system should check that this user is authorized to access the requested data. (This fragment is not shown in Figure 2-26.) Figure 5-16 depicts the Boolean logic for controlling the access to the data in the system database. One way to implement this scheme is to write one large conditional IF-THEN-ELSE statement. This approach would lead to a complex code that is difficult to understand and extend if new policies or roles need to be considered (e.g., the Maintenance

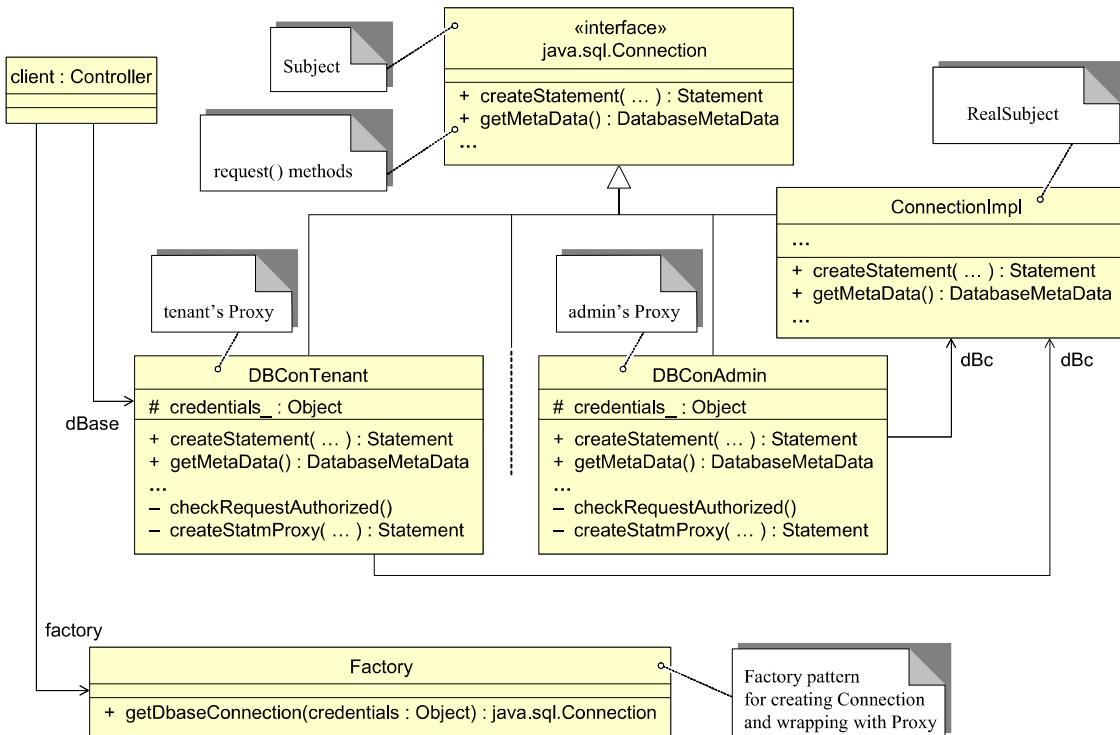


Figure 5-17: Class diagram for the example proxy for enforcing authorized database access. See the interactions in Figure 5-18. (Compare to Figure 5-15(b) for generic Proxy pattern.)

actor). In addition, it serves as a distraction from the main task of the client or server objects. This is where protection proxy enters the picture in and takes on the authorization responsibility.

Figure 5-17 shows how Protection Proxy is implemented in the example of safe database access. Each different proxy type specifies a set of legal messages from client to subject that are appropriate for the current user's access rights. If a message is not legal, the proxy will not forward it to the real subject (the database connection object `ConnectionImpl`); instead, the proxy will send an error message back to the caller (i.e., the client).

In this example, the Factory object acts as a custodian that sets up the Proxy pattern (see Figure 5-17 and Figure 5-18).

It turns out that in this example we need two types of proxies: (a) proxies that implement the database connection interface, such as `java.sql.Connection` if Java is used; and (b) proxies that implement the SQL statement interface, such as `java.sql.Statement` if Java is used. The connection proxy guards access to the database metadata, while the statement proxy guards access to the database data. The partial class diagram in Figure 5-17 shows only the connection-proxy classes, and Figure 5-18 mentions the statement proxy only in the last method call `createStatmProxy()`, by which the database proxy (`DBConTenant`) creates a statement proxy and returns it.

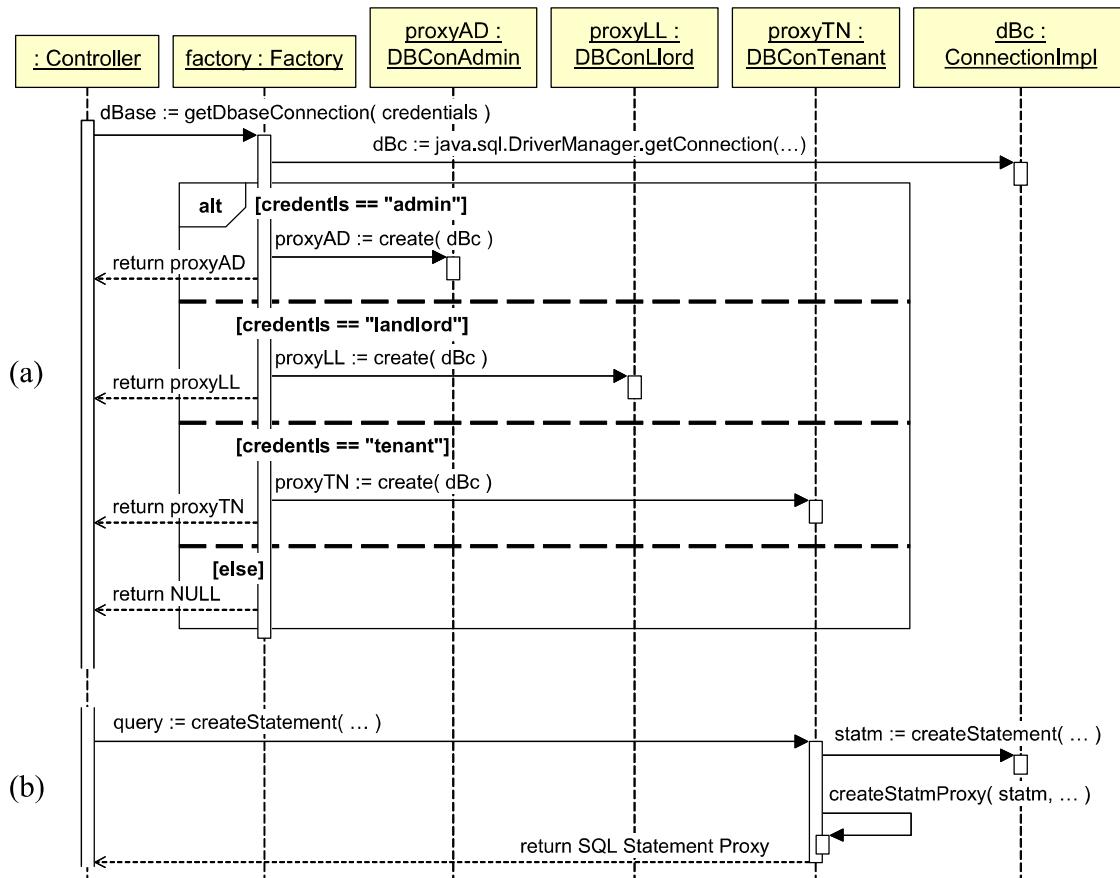


Figure 5-18: Example of Protection Proxy setup (a) and use (b) that solves the access-control problem from Figure 5-16. (See the corresponding class diagram in Figure 5-17.)

Figure 5-18

Listing 5-5: Implementation of the Protection Proxy that provides safe access to the database in the secure home access system.

```

import java.sql.Connection;
import java.sql.DriverManager;
import java.sql.ResultSet;
import java.sql.Statement;
import javax.servlet.ServletConfig;
import javax.servlet.http.HttpServlet;
import javax.servlet.http.HttpServletRequest;
import javax.servlet.http.HttpServletResponse;

public class WebDataAccessServlet extends HttpServlet {
    private String           // database access parameters
        driverClassName = "com.mysql.jdbc.Driver",
        dbURL = "jdbc:mysql://localhost/homeaccessrecords",
        dbUserID = null,
        dbPassword = null;
    private Connection dBase = null;

```

```
public void init(ServletConfig config) throws ServletException {
    super.init(config);
    ...

    dbUserID = config.getInitParameter("userID");
    dbPassword = config.getInitParameter("password");
    Factory factory = new Factory(driverClassName, dbURL);
    dBase = factory.getDbaseConnection(dbUserID, dbPassword);
}

public void service(
    HttpServletRequest req, HttpServletResponse resp
) throws ServletException, java.io.IOException {
    Statement statm = dBase.createStatement();
    // process the request and prepare ...
    String sql = // ... an SQL statement from the user's request
    boolean ok = statm.execute(sql);

    ResultSet result = statm.getResultSet();
    // print the result into the response (resp argument)
}
}

public class Factory {
    protected String dbURL_;
    protected Connection dBc_ = null;

    public Factory(String driverClassName, String dbURL) {
        // load the database driver class (the Driver class creates
        Class.forName(driverClassName); // an instance of itself)
        dbURL_ = dbURL;
    }

    public Connection getDbaseConnection(
        String dbUserID, String dbPassword
    ) {
        dBc_ = DriverManager.getConnection(
            dbURL_, dbUserID, dbPassword
        );

        Connection proxy = null;

        int userType = getUserType(dbUserID, dbPassword);
        switch (userType) {
            case 1:           // dbUserID is a system administrator
                proxy = new DBConAdmin(dBc_, dbUserID, dbPassword);
            case 2:           // dbUserID is a landlord
                proxy = new DBConLord(dBc_, dbUserID, dbPassword);
            case 3:           // dbUserID is a tenant
                proxy = new DBConTenant(dBc_, dbUserID, dbPassword);
            default:          // dbUserID cannot be identified
                proxy = null;
        }
        return proxy;
    }
}
```

```
// Protection Proxy class for the actual java.sql.Connection
public class DBConTenant implements Connection {
    protected Connection dbc_ = null;
    protected String
        dbUserID = null,
        dbPassword = null;;

    public DBConTenant(
        Connection dbc, String dbUserID, String dbPassword
    ) {
        ...
    }

    public Statement createStatement() {
        statm = dbc_.createStatement();

        return createStatmProxy(statm, credentials_);
    }

    private Statement createStatmProxy(
        Statement statm, credentials_
    ) {
        // create a proxy of java.sql.Statement that is appropriate
        // for a user of the type "tenant"
    }
}
```

One may wonder if we should similarly use Protection Proxy to control the access to the locks in a building, so the landlord has access to all apartments and a tenant only to own apartment. When considering the merits of this approach, the developer first needs to compare it to a straightforward conditional statement and see which approach would create a more complex implementation.

The above example illustrated the use of Proxy to implement security policies for authorized data access. Another example involves safety policies to limit the number of simultaneous accesses. For example, to avoid inconsistent reads/writes, the policy may allow at most one client at a time to access the subject, which in effect serializes the access to the subject. This constraint is implemented by passing a token among clients—only the client in possession of the token can access the subject, by presenting the token when requesting access.

The Protection Proxy pattern is structurally identical to the Decorator pattern (compare Figure 5-11 and Figure 5-15). We can also create a chain of Proxies, same as with the Decorators (Figure 5-11(c)). The key difference is in the intent: Protection Proxy protects an object (e.g., from unauthorized access) while Decorator adds special-case behavior to an object.

◆ The reader may have noticed that many design patterns look similar to one another. For example, Proxy is *structurally* almost identical to Decorator. The difference between them is in their *intention*—what they are used for. The intention of Decorator is to *add* functionality, while the intention of Proxy is to *subtract* functionality, particularly for Protection Proxy.

5.3 Concurrent Programming

"The test of a first-rate intelligence is the ability to hold two opposed ideas in mind at the same time and still retain the ability to function." —F. Scott Fitzgerald

The benefits of concurrent programming include better use of multiple processors and easier programming of reactive (event-driven) applications. In event-driven applications, such as graphical user interfaces, the user expects a quick response from the system. If the (single-processor) system processes all requests sequentially, then it will respond with significant delays and most of the requestors will be unhappy. A common technique is to employ *time-sharing* or time slicing—a *single processor dedicates a small amount of time for each task*, so all of them move forward collectively by taking turns on the processor. Although none of the tasks progresses as fast as it would if it were alone, none of them has to wait as long as it could have if the processing were performed sequentially. The task executions are really sequential but *interleaved* with each other, so they virtually appear as concurrent. In the discussion below I ignore the difference between real concurrency, when the system has multiple processors, and virtual concurrency on a single-processor system. From the user's viewpoint, there is no logical or functional difference between these two options—the user would only see difference in the length of execution time.

Computer *process* is, roughly speaking, a task being executed by a processor. A task is defined by a temporally ordered sequence of instructions (program code) for the processor. In general, a process consists of:

- Memory, which contains executable program code and/or associated data
- Operating system resources that are allocated to the process, such as file descriptors (Unix terminology) or handles (Windows terminology)
- Security attributes, such as the identity of process owner and the process's set of privileges
- Processor state, such as the content of registers, physical memory addresses, etc. The state is stored in the actual registers when the process is executing, and in memory otherwise

Threads are similar to processes, in that both represent a single sequence of instructions executed in parallel with other sequences, either by time slicing on a single processor or multiprocessing. A *process* is an entirely independent program, carries considerable state information, and interacts with other processes only through system-provided inter-process communication mechanisms. Conversely, a thread directly shares the state variables with other threads that are part of the same

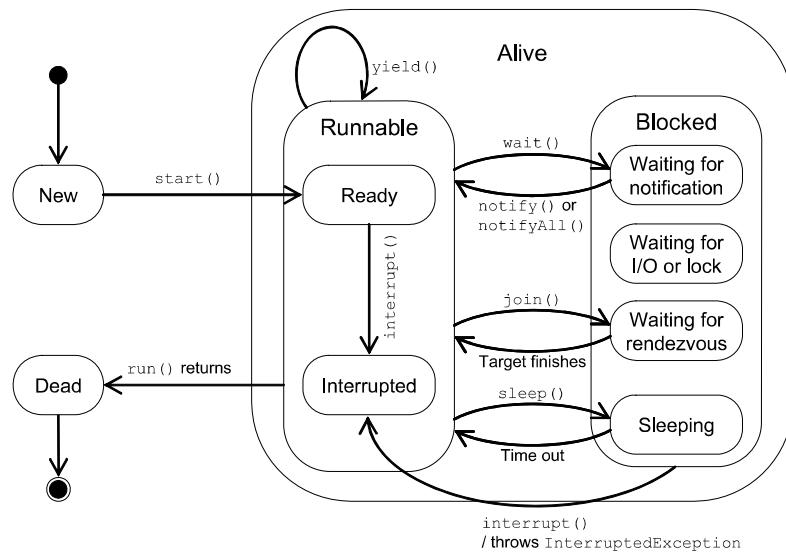


Figure 5-19: State diagram representing the lifecycle of Java threads. (State diagram notation is defined in Section 3.2.2.)

process, as well as memory and other resources. In this section I focus on threads, but many concepts apply to processes as well.

So far, although I promoted the metaphor of an object-based program as a “bucket brigade,” the objects carried their tasks sequentially, one after another, so in effect the whole system consists of a single worker taking the guises one-by-one of different software objects. Threads allow us to introduce true parallelism in the system functioning.

Subdividing a problem to smaller problems (subtasks) is a common strategy in problem solving. It would be all well and easy if the subtasks were always disjoint, clearly partitioned and independent of each other. However, during the execution the subtasks often operate on the same resources or depend on results of other task(s). This is what makes concurrent programming complex: threads (which roughly correspond to subtasks) interact with each other and must coordinate their activities to avoid incorrect results or undesired behaviors.

5.3.1 Threads

A *thread* is a sequence of processor instructions, which can share a single address space with other threads—that is, they can read and write the same program variables and data structures. Threads are a way for a program to split itself into two or more concurrently running tasks. It is a basic unit of program execution. A common use of threads is in reactive applications, having one thread paying attention to the graphical user interface, while others do long calculations in the background. As a result, the application more readily responds to user’s interaction.

Figure 5-19 summarizes different states that a thread may go through in its lifetime. The three main states and their sub-states are:

1. *New*: The thread object has been created, but it has not been started yet, so it cannot run

2. *Alive*: After a thread is started, it becomes alive, at which point it can enter several different sub-states (depending on the method called or actions of other threads within the same process):
 - a. *Runnable*: The thread *can* be run when the time-slicing mechanism has CPU cycles available for the thread. In other words, when there is nothing to prevent it from being run if the scheduler can arrange it
 - b. *Blocked*: The thread could be run, but there is something that prevents it (e.g., another thread is holding the resource needed for this thread to do its work). While a thread is in the blocked state, the scheduler will simply skip over it and not give it any CPU time, so the thread will not perform any operations. As visible from Figure 5-19, a thread can become blocked for the following reasons:
 - i. *Waiting for notification*: Invoking the method `wait()` suspends the thread until the thread gets the `notify()` or `notifyAll()` message
 - ii. *Waiting for I/O or lock*: The thread is waiting for an input or output operation to complete, or it is trying to call a `synchronized` method on a shared object, and that object's lock is not available
 - iii. *Waiting for rendezvous*: Invoking the method `join(target)` suspends the thread until the `target` thread returns from its `run()` method
 - iv. *Sleeping*: Invoking the method `sleep(milliseconds)` suspends the thread for the specified time
3. *Dead*: This normally happens to a thread when it returns from its `run()` method. A dead thread cannot be restarted, i.e., it cannot become alive again

The meaning of the states and the events or method invocations that cause state transitions will become clearer from the example in Section 5.3.4.

A thread object may appear as any other software object, but there are important differences. Threads are not regular objects, so we have to be careful with their interaction with other objects. Most importantly, we cannot just call a method on a thread object, because that would execute the given method from our current thread—neglecting the thread of the method's object—which could lead to conflict. To pass a message from one thread to another, we must use only the methods shown in Figure 5-19. No other methods on thread objects should be invoked.

If two or more threads compete for the same “resource” which can be used by only one at a time, then their access must be serialized, as depicted in Figure 5-20. One of them becomes blocked while the other proceeds. We are all familiar with conflicts arising from people sharing resources. For example, people living in a house/apartment share the same bathroom. Or, many people may be sharing the same public payphone. To avoid conflicts, people follow certain protocols, and threads do similarly.

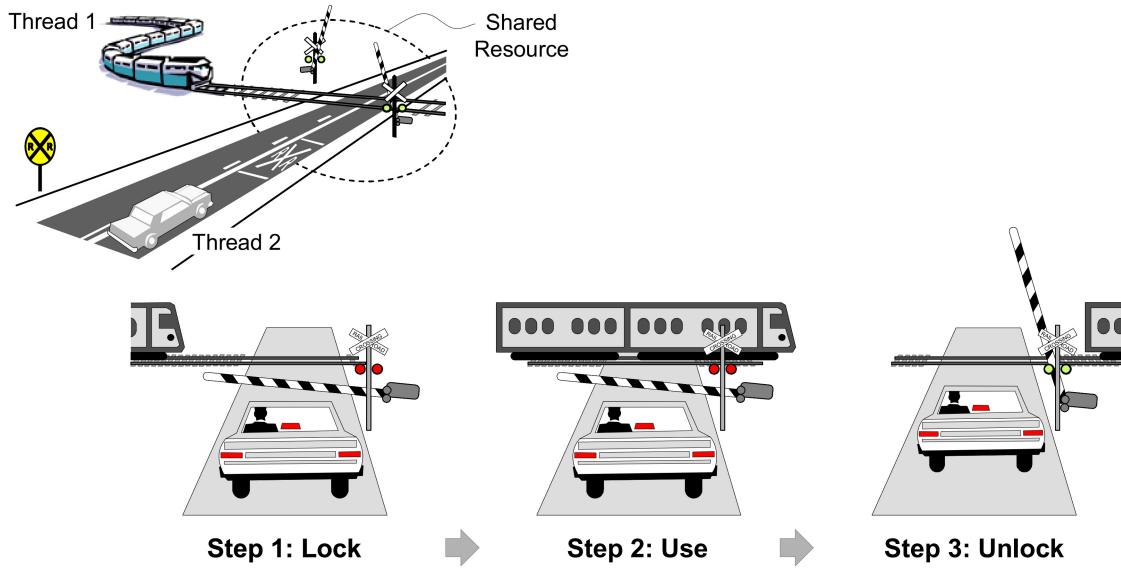


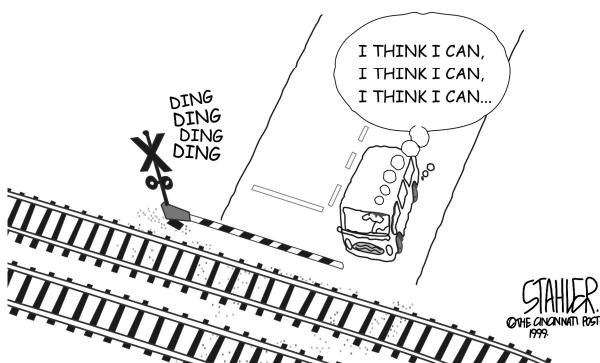
Figure 5-20: Illustration of exclusion synchronization. The lock simply ensures that concurrent accesses to the shared resource are serialized.

5.3.2 Exclusive Resource Access—Exclusion Synchronization

If several threads attempt to access and manipulate the same data concurrently a *race condition* or *race hazard* exists, and the outcome of the execution depends on the particular order in which the access takes place. Consider the following example of two threads simultaneously accessing the same banking account (say, husband and wife interact with the account from different branches):

Thread 1	Thread 2
<pre>oldBalance = account.getBalance(); newBalance = oldBalance + deposit; account.setBalance(newBalance); ... </pre>	<pre>... oldBalance = account.getBalance(); newBalance = oldBalance - withdrawal; account.setBalance(newBalance); </pre>

The final account balance is incorrect and the value depends on the order of access. To avoid race hazards, we need to control access to the common data (shared with other threads) and make the access sequential instead of parallel.



A segment of code in which a thread may modify shared data is known as a *critical section* or *critical region*. The critical-section problem is to design a protocol that threads can use to avoid interfering with each other. *Exclusion synchronization*, or mutual exclusion (mutex), see Figure 5-21, stops different threads from calling methods on the same object at the same time and thereby jeopardizing the integrity of the shared data. If thread is executing in its critical region then no other thread can be

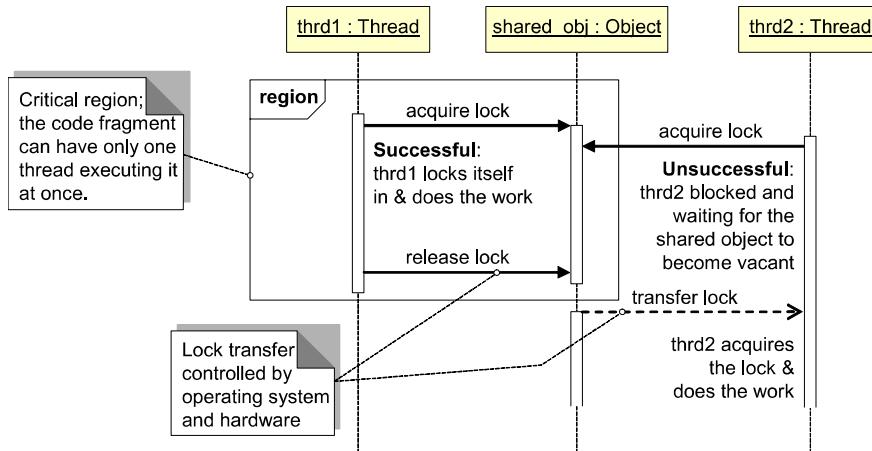


Figure 5-21: Exclusion synchronization pattern for concurrent threads.

executing in its critical region. Only one thread is allowed in a critical region at any moment.

Java provides exclusion synchronization through the keyword `synchronized`, which simply labels a block of code that should be protected by locks. Instead of the programmer explicitly acquiring or releasing the lock, `synchronized` signals to the compiler to do so. As illustrated in Figure 5-22, there are two ways to use the keyword `synchronized`. First technique declares class methods `synchronized`, Figure 5-22(a). If one thread invokes a `synchronized` method on an object, that object is *locked*. Another thread invoking this or another `synchronized` method on that same object will block until the lock is released.

Nesting method invocations are handled in the obvious way: when a `synchronized` method is invoked on an object that is already locked by that same thread, the method returns, but the lock is not released until the outermost `synchronized` method returns.

Second technique designates a statement or a block of code as `synchronized`. The parenthesized *expression* must produce an object to lock—usually, an object reference. In the simplest case, it could be `this` reference to the current object, like so

```
synchronized (this) { /* block of code statements */ }
```

When the lock is obtained, *statement* is executed as if it were `synchronized` method on the locked object. Examples of exclusion synchronization in Java are given in Section 5.3.4.

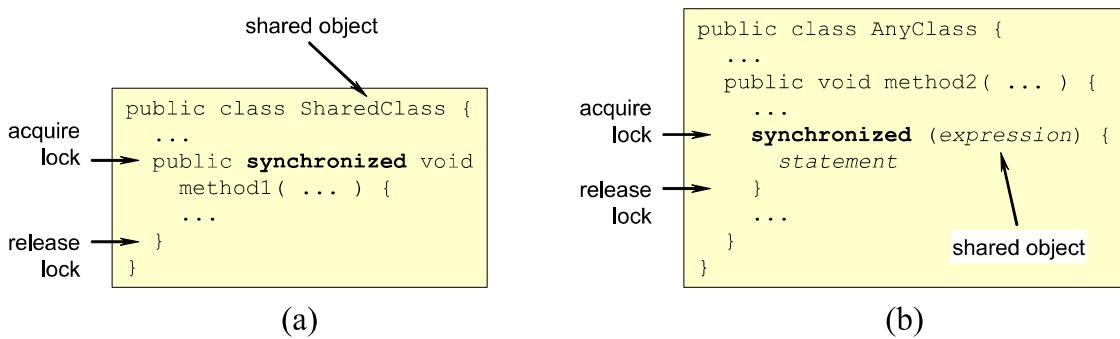


Figure 5-22: Exclusion synchronization in Java: (a) `synchronized` methods, and (b) `synchronized` statements.

5.3.3 Cooperation between Threads—Condition Synchronization

Exclusion synchronization ensures that threads “do not step on each other’s toes,” but other than preventing them from colliding, their activities are completely independent. However, sometimes one thread’s work depends on activities of another thread, so they must cooperate and coordinate. A classic example of cooperation between threads is a `Buffer` object with methods `put()` and `get()`. Producer thread calls `put()` and consumer thread calls `get()`. The producer must wait if the buffer is full, and the consumer must wait if it is empty. In both cases, threads wait for a condition to become fulfilled. Condition synchronization includes no assumption that the wait will be brief; threads could wait indefinitely.

Condition synchronization (illustrated in Figure 5-23) complements exclusion synchronization. In exclusion synchronization, a thread encountering an occupied shared resource becomes blocked and waits until another thread is finished with using the resource. Conversely, in condition synchronization, a thread encountering an unmet condition cannot continue holding the resource on which condition is checked and just wait until the condition is met. If the thread did so, no other thread would have access to the resource and the condition would never change—the resource must be released, so another thread can affect it. The thread in question might release the resource and periodically return to check it, but this would not be an efficient use of processor cycles. Rather, the thread becomes blocked and does nothing while waiting until the condition changes, at which point it must be explicitly notified of such changes.

In the buffer example, a producer thread, t , must first lock the buffer to check if it is full. If it is, t enters the “waiting for notification” state, see Figure 5-19. But while t is waiting for the condition to change, the buffer must remain unlocked so consumers can empty it by calling `get()`. Because the waiting thread is blocked and inactive, it needs to be notified when it is ready to go.

Every software object in Java has the methods `wait()` and `notify()` which makes possible sharing and condition synchronization on every Java object, as explained next. The method `wait()` is used for suspending threads that are waiting for a condition to change. When t finds the buffer full, it calls `wait()`, which *atomically* releases the lock and suspends the thread (see Figure 5-23). Saying that the thread suspension and lock release are *atomic* means that they happen together, indivisibly from the application’s point of view. After some other thread notifies t that the buffer may no longer be full, t regains the lock on the buffer and retests the condition.

The standard Java idiom for condition synchronization is the statement:

```
while (conditionIsNotMet) sharedObject.wait();
```

Such a *wait-loop statement* must be inside a synchronized method or block. Any attempt to invoke the `wait()` or `notify()` methods from outside the synchronized code will throw `IllegalMonitorStateException`. The above idiom states that the condition test should *always* be in a loop—never assume that being woken up means that the condition has been met.

The wait loop blocks the calling thread, t , for as long as the condition is not met. By calling `wait()`, t places itself in the shared object’s wait set and releases all its locks on that object. (It is of note that standard Java implements an unordered “wait set” rather than an ordered “wait queue.” Real-time specification for Java—RTSJ—corrects this somewhat.)

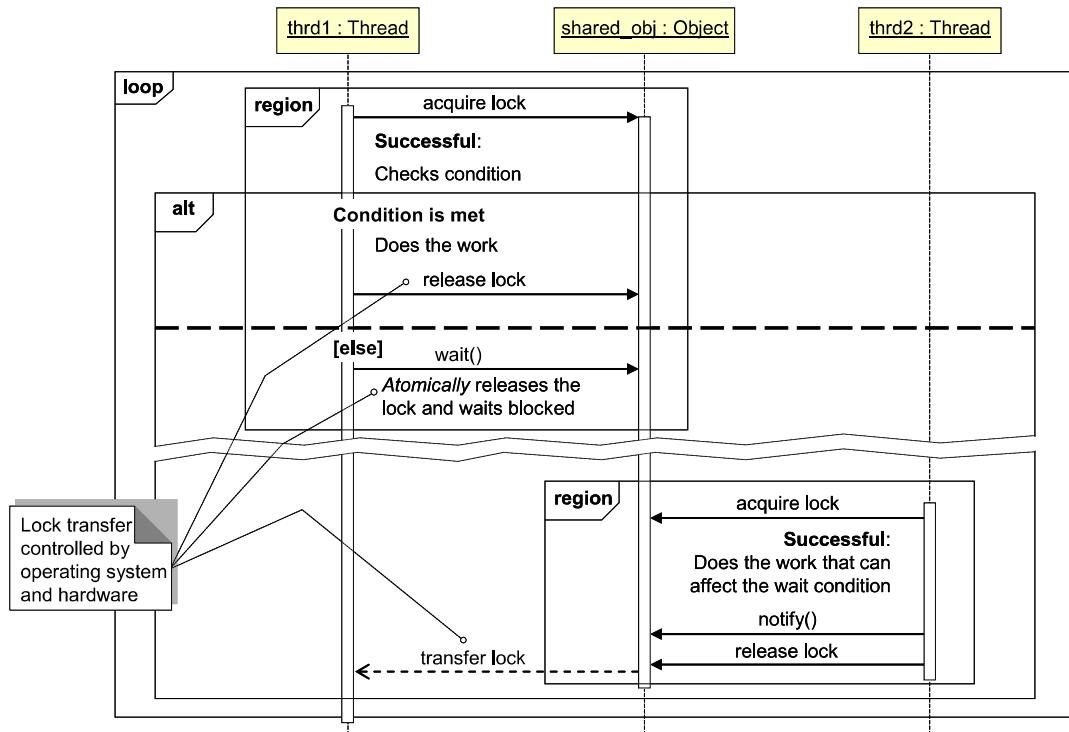


Figure 5-23: Condition synchronization pattern for concurrent threads.

A thread that executes a `synchronized` method on an object, o , and changes a condition that can affect one or more threads in o 's wait set must notify those threads. In standard Java, the call `o.notify()` reactivates one arbitrarily chosen thread, t , in o 's wait set. The reactivated thread then reevaluates the condition and either proceeds into the critical region or reenters the wait set. The call to `o.notifyAll()` releases all threads in the o 's wait set. In standard Java this is the only way to ensure that the highest priority thread is reactivated. This is inefficient, though, because all the threads must attempt access while only one will succeed in acquiring the lock and proceed.

The reader might have noticed resemblance between the above mechanism of wait/notify and the publish/subscribe pattern of Section 5.1. In fact, they are equivalent conceptually, but there are some differences due to concurrent nature of condition synchronization.

5.3.4 Concurrent Programming Example

The following example illustrates cooperation between threads.

Example 5.1 Concurrent Home Access

In our case-study, Figure 1-12 shows lock controls both on front and backyard doors. Suppose two different tenants arrive (almost) simultaneously at the different doors and attempt the access, see Figure 5-24. The single-threaded system designed in Section 2.7 would process them one-by-one, which may cause the second user to wait considerable amount of time. A user unfamiliar with the system intricacies may perceive this as a serious glitch. As shown in Figure 5-24, the processor is *idle* most of the time, such as between individual keystrokes or while the user tries to recall the exact

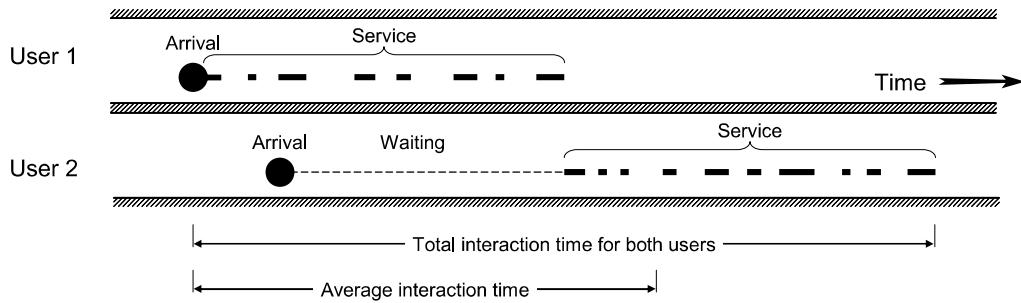


Figure 5-24: Single-threaded, sequential servicing of users in Example 4.1.

password after an unsuccessful attempt. Meanwhile, the second user is needlessly waiting. To improve the user experience, let us design a multithreaded solution.

The solution is given next.

The first-round implementation in Section 2.7, considered the system with a single door lock. We have not yet tackled the architectural issue of running a centralized or a distributed system. In the former case, the main computer runs the system and at the locks we have only embedded processors. We could add an extra serial port, daisy-chained with the other one, and the control would remain as in Section 2.7. In the latter case of a distributed system, each lock would have a proximal embedded computer. The embedded computers would communicate mutually or with the main computer using a local area network. The main computer may even not be necessary, and the embedded processors could coordinate in a “peer-to-peer” mode. Assume for now that we implement the centralized PC solution with multiple serial ports. We also assume a single photosensor and a single light bulb, for the sake of simplicity.

The first question to answer is, how many threads we need and which objects should be turned into threads? Generally, it is not a good idea to add threads indiscriminately, because threads consume system resources, such as computing cycles and memory space.

It may appear attractive to attach a thread to each object that operates physical devices, such as LockCtrl and LightCtrl, but is this the right approach? On the other hand, there are only two users (interacting with the two locks), so perhaps two threads would suffice? Let us roll back, see why we consider introducing threads in the first place. The reason is to improve the user experience, so two users at two different doors can access the home simultaneously, without waiting. Two completely independent threads would work, which would require duplicating all the resources, but this may be wasteful. Here is the list of resources they could share:

- KeyStorage, used to lookup the valid keys
- Serial port(s), to communicate with the devices
- System state, such as the device status or current count of unsuccessful attempts

Sharing KeyStorage seems reasonable—here it is just looked up, not modified. The serial port can also be shared because the communication follows a well-defined RS-232 protocol. However, sharing the system state needs to be carefully examined. Sharing the current count of

unsuccessful attempts seems to make no sense—there must be two counters, each counting accesses for its corresponding door.

There are several observations that guide our design. From the system sequence diagram of Figure 2-15(a), we can observe that the system juggles two distinct tasks: user interaction and internal processing which includes controlling the devices. There are two copies (for the two doors) of each task, which should be able to run in parallel. The natural point of separation between the two tasks is the Controller object, Figure 2-27, which is the entry point of the domain layer of the system. The Controller is a natural candidate for a thread object, so two internal processing tasks can run in parallel, possibly sharing some resources. The threaded controller class, ControllerThd, is defined below. I assume that all objects operating the devices (LockCtrl, LightCtrl, etc.) can be shared as long as the method which writes to the serial port is synchronized. LockCtrl must also know which lock (front or backyard) it currently operates.

Listing 5-6: Concurrent version of the main class for home access control. Compare to Listing 2-1.

```
import javax.comm.*;
import java.io.IOException;
import java.io.InputStream;
import java.util.TooManyListenersException;

public class HomeAccessControlSystem_2x extends Thread
    implements SerialPortEventListener {
    protected ControllerThd contrlFront_; // front door controller
    protected ControllerThd contrlBack_; // back door controller
    protected InputStream inputStream_; // from the serial port
    protected StringBuffer keyFront_ = new StringBuffer();
    protected StringBuffer keyBack_ = new StringBuffer();
    public static final long keyCodeLen_ = 4; // key code of 4 chars

    public HomeAccessControlSystem_2x(
        KeyStorage ks, SerialPort ctrlPort
    ) {
        try {
            inputStream_ = ctrlPort.getInputStream();
        } catch (IOException e) { e.printStackTrace(); }

        LockCtrl lkc = new LockCtrl(ctrlPort);
        LightCtrl lic = new LightCtrl(ctrlPort);
        PhotoObsrv sns = new PhotoObsrv(ctrlPort);
        AlarmCtrl ac = new AlarmCtrl(ctrlPort);

        contrlFront_ = new ControllerThd(
            new KeyChecker(ks), lkc, lic, sns, ac, keyFront_
        );
        contrlBack_ = new ControllerThd(
            new KeyChecker(ks), lkc, lic, sns, ac, keyBack_
        );

        try {
            ctrlPort.addEventListener(this);
        } catch (TooManyListenersException e) {
```

```

        e.printStackTrace(); // limited to one listener per port
    }
    start(); // start the serial-port reader thread
}

/** The first argument is the handle (filename, IP address, ...)
 * of the database of valid keys.
 * The second arg is optional and, if present, names
 * the serial port. */
public static void main(String[] args) {
    ...
    // same as in Listing 2-1 above
}

/** Thread method; does nothing, just waits to be interrupted
 * by input from the serial port. */
public void run() {
    while (true) {
        try { Thread.sleep(100); }
        catch (InterruptedException e) { /* do nothing */ }
    }
}

/** Serial port event handler.
 * Assume that the characters are sent one by one, as typed in.
 * Every character is preceded by a lock identifier (front/back).
 */
public void serialEvent(SerialPortEvent evt) {
    if (evt.getEventType() == SerialPortEvent.DATA_AVAILABLE) {
        byte[] readBuffer = new byte[5]; // just in case, 5 chars

        try {
            while (inputStream_.available() > 0) {
                int numBytes = inputStream_.read(readBuffer);
                // could chk if numBytes == 2 (char + lockId) ...
            }
        } catch (IOException e) { e.printStackTrace(); }

        // Append the new char to a user key, and if the key
        // is complete, awaken the corresponding Controller thread
        if (inputStream_[0] == 'f') { // from the front door
            // If this key is already full, ignore the new chars
            if (keyFront_.length() < keyCodeLen_) {
                synchronized (keyFront_) { // CRITICAL REGION
                    keyFront_.append(new String(readBuffer, 1, 1));
                    // If the key just got completed,
                    // signal the condition to others
                    if (keyFront_.length() >= keyCodeLen_) {
                        // awaken the Front door Controller
                        keyFront_.notify(); //only 1 thrd waiting
                    }
                } // END OF THE CRITICAL REGION
            }
        } else if (inputStream_[0] == 'b') { // from back door
            if (keyBack_.length() < keyCodeLen_) {
                synchronized (keyBack_) { // CRITICAL REGION
                    keyBack_.append(new String(readBuffer, 1, 1));
                }
            }
        }
    }
}

```

```

        if (keyBack_.length() >= keyCodeLen_) {
            // awaken the Back door Controller
            keyBack_.notify();
        }
    } // END OF THE CRITICAL REGION
} // else, exception ?!
}
}
}
```

Each Controller object is a thread, and it synchronizes with HomeAccessControlSystem_2x via the corresponding user key. In the above method `serialEvent()`, the port reader thread fills in the key code until completed; thereafter, it ignores the new keys until the corresponding ControllerThd processes the key and resets it in its `run()` method, shown below. The reader should observe the reuse of a `StringBuffer` to repeatedly build strings, which works here, but in a general case many subtleties of Java `StringBuffers` should be considered.

Listing 5-7: Concurrent version of the `Controller` class. Compare to Listing 2-2.

```
import javax.comm.*;

public class ControllerThd implements Runnable {
    protected KeyChecker checker_;
    protected LockCtrl lockCtrl_;
    protected LightCtrl lightCtrl_;
    protected PhotoObsrv sensor_;
    protected AlarmCtrl alarmCtrl_;
    protected StringBuffer key_;
    public static final long maxNumOfAttempts_ = 3;
    public static final long attemptPeriod_ = 600000; // msec [=10min]
    protected long numOfAttempts_ = 0;

    public ControllerThd(
        KeyChecker kc, LockCtrl lkc, LightCtrl lic,
        PhotoObsrv sns, AlarmCtrl ac, StringBuffer key
    ) {
        checker_ = kc;
        lockCtrl_ = lkc; alarmCtrl_ = ac;
        lightCtrl_ = lic; sensor_ = sns; key_ = key;

        Thread t = new Thread(this, getName());
        t.start();
    }

    public void run() {
        while(true) { // runs forever
            synchronized (key_) { // CRITICAL REGION
                // wait for the key to be completely typed in
                while(key_.length() <
                    HomeAccessControlSystem_2x.keyCodeLen_) {
                    try {
                        key_.wait();
                    } catch(InterruptedException e) {
                        throw new RuntimeException(e);
                    }
                }
            }
        }
    }
}
```

```

        }
    }
} // END OF THE CRITICAL REGION
// Got the key, check its validity:
// First duplicate the key buffer, then release old copy
Key user_key = new Key(new String(key_));
key_.setLength(0); // delete code, so new can be entered
checker_.checkKey(user_key); // assume Publish-Subs. vers.
}
}
}
```

The reader should observe the thread coordination in the above code. We do not want the Controller to grab a half-ready Key and pass it to the Checker for validation. The Controller will do so only when `notify()` is invoked. Once it is done with the key, the Controller resets it to allow the reader to fill it again.

You may wonder how is it that in `ControllerThd.run()` we obtain the lock and then loop until the key is completely typed in—would this not exclude the port reader thread from the access to the `Key` object, so the key would never be completed?! Recall that `wait()` atomically *releases* the lock and suspends the `ControllerThd` thread, leaving `Key` accessible to the port reader thread.

It is interesting to consider the last three lines of code in `ControllerThd.run()`. Copying a `StringBuffer` to a new `String` is a thread-safe operation; so is setting the length of a `StringBuffer`. However, although each of these methods acquires the lock, the lock is released in between and another thread may grab the key object and do something bad to it. In our case this will not happen, because `HomeAccessControlSystem_2x.serialEvent()` checks the length of the key before modifying it, but generally, this is a concern.

Figure 5-25 summarizes the benefit achieved by a multithreaded solution. Notice that there still may be micro periods of waiting for both users and servicing the user who arrived first may take longer than in a single-threaded solution. However, the average service time per user is much shorter, close to the single-user average service time.

Hazards and Performance Penalties

Ideally, we would like that the processor is never idle while there is a task waiting for execution. As seen in Figure 5-25(b), even with threads the processor may be idle while there are users who are waiting for service. The question of “granularity” of the shared resource. Or stated differently, the key issue is how to minimize the length (i.e., processing time) of the critical region.

Solution: Try to narrow down the critical region by lock splitting or using finer-grain locks.

<http://www.cs.panam.edu/~meng/Course/CS6334>Note/master/node49.html>

http://searchstorage.techtarget.com/sDefinition/0,,sid5_gci871100,00.html

[http://en.wikipedia.org/wiki/Thread_\(computer_programming\)](http://en.wikipedia.org/wiki/Thread_(computer_programming))

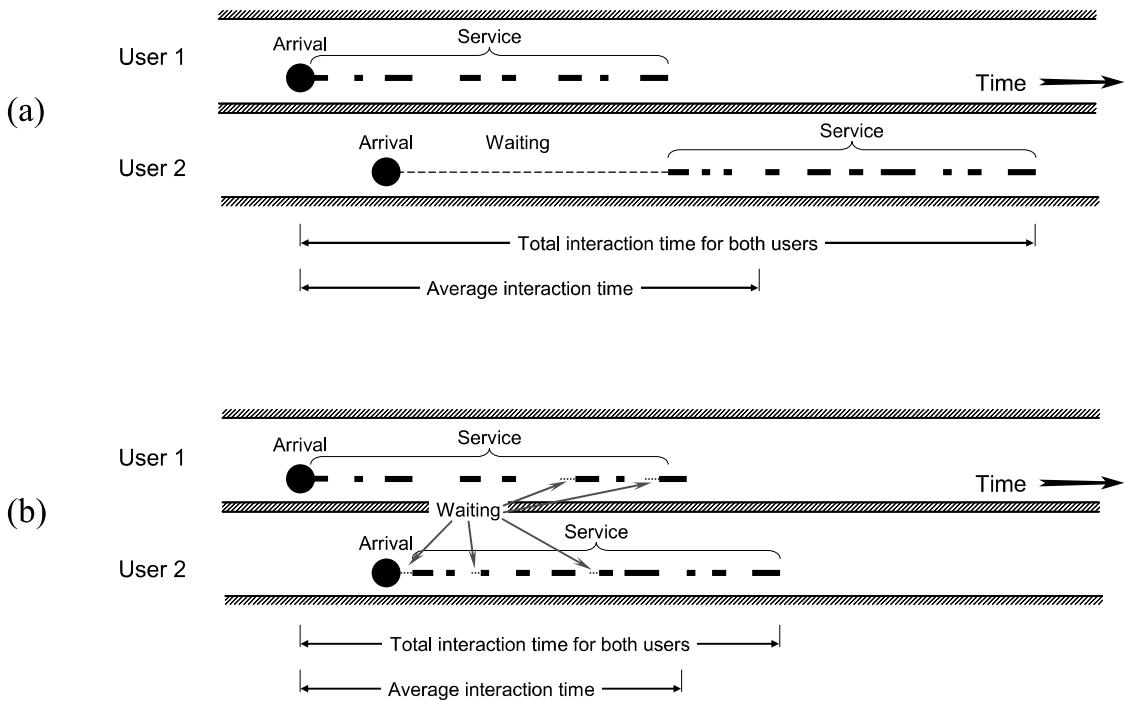


Figure 5-25: Benefit of a multithreaded solution. (a) Sequential servicing, copied from Figure 5-24. (b) Parallel servicing by two threads.

Control of access to shared resources itself can introduce problems, e.g., it can cause deadlock.

5.4 Broker and Distributed Computing

“If computers get too powerful, we can organize them into a committee—that will do them in.”
—Bradley’s Bromide

Let us assume that in our case-study example of home access control the tenants want to remotely download the list of recent accesses. This requires network communication. The most basic network programming uses network sockets, which can call for considerable programming skills (see Appendix B for a brief review). To simplify distributed computing, a set of software techniques have been developed. These generally go under the name of network middleware.

When both client and server objects reside in the same memory space, the communication is carried out by simple method calls on the server object (see Figure 1-22). If the objects reside in different memory spaces or even on different machines, they need to implement the code for interprocess communication, such as opening network connections, sending messages across the network, and dealing with many other aspects of network communication protocols. This significantly increases the complexity of objects. Even worse, in addition to its business logic, objects obtain responsibility for communication logic which is extraneous to their main function.

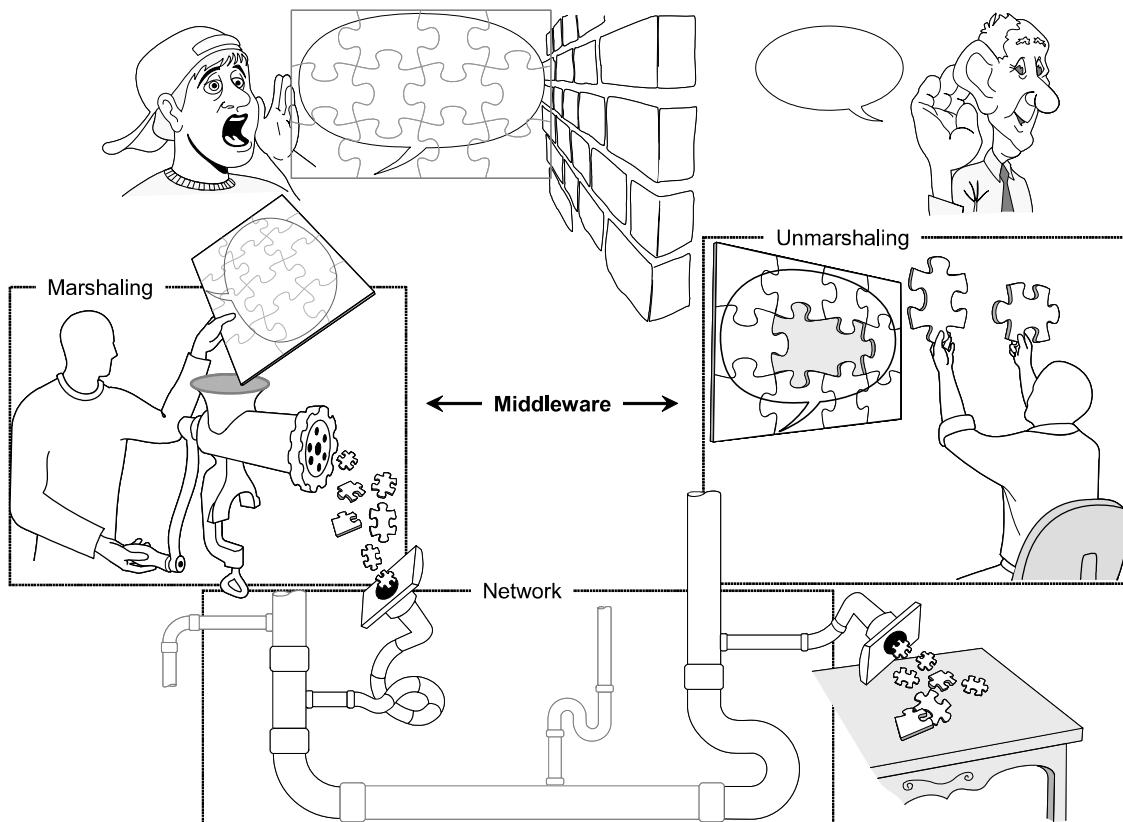


Figure 5-26: Client object invokes a method on a server object. If the message sending becomes too complex, introducing middleware offloads the communication intricacies off the software objects. (Compare with Figure 1-22.)

Employing middleware helps to delegate this complexity away out of the objects, see Figure 5-26. A real world example of middleware is the post service—it deals with the intricacies of delivering letters/packages to arbitrary distances. Another example is the use of different metric systems, currencies, spoken languages, etc., in which case the functionality for conversion between the systems is offloaded to middleware services. Middleware is a good software engineering practice that should be applied any time the communication between objects becomes complex and starts rivaling the object's business logic in terms of the implementation code size.

Middleware is a collection of objects that offer a set of services related to object communication, so that extraneous functionality is offloaded to the middleware. In general, middleware is software used to make diverse applications work together smoothly. The process of deriving middleware is illustrated in Figure 5-27. We start by introducing the proxies for both communicating objects. (The *Proxy* pattern is described in Section 5.2.4.) The proxy object B' of B acts so to provide an illusion for A that it is directly communicating with B . The same is provided to B by A' . Having the proxy objects keeps simple the original objects, because the proxies provide them with an illusion that nothing changed from the original case, where they communicated directly with each other, as in Figure 5-27(a). In other words, proxies provide *location transparency* for the objects—objects remain (almost) the same no matter whether they

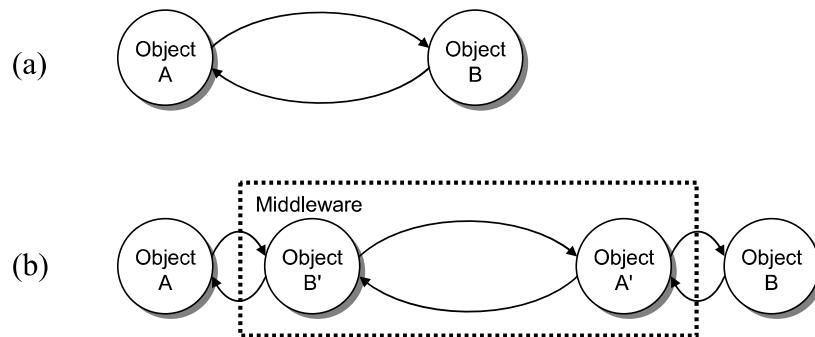


Figure 5-27: Middleware design. Objects A' and B' are the *proxies* of A and B , respectively.

reside in the same address space or in different address spaces / machines. Objects A' and B' comprise the network middleware.

Because it is not likely that we will develop middleware for only two specific objects communicating, further division of responsibilities results in the *Broker* pattern.

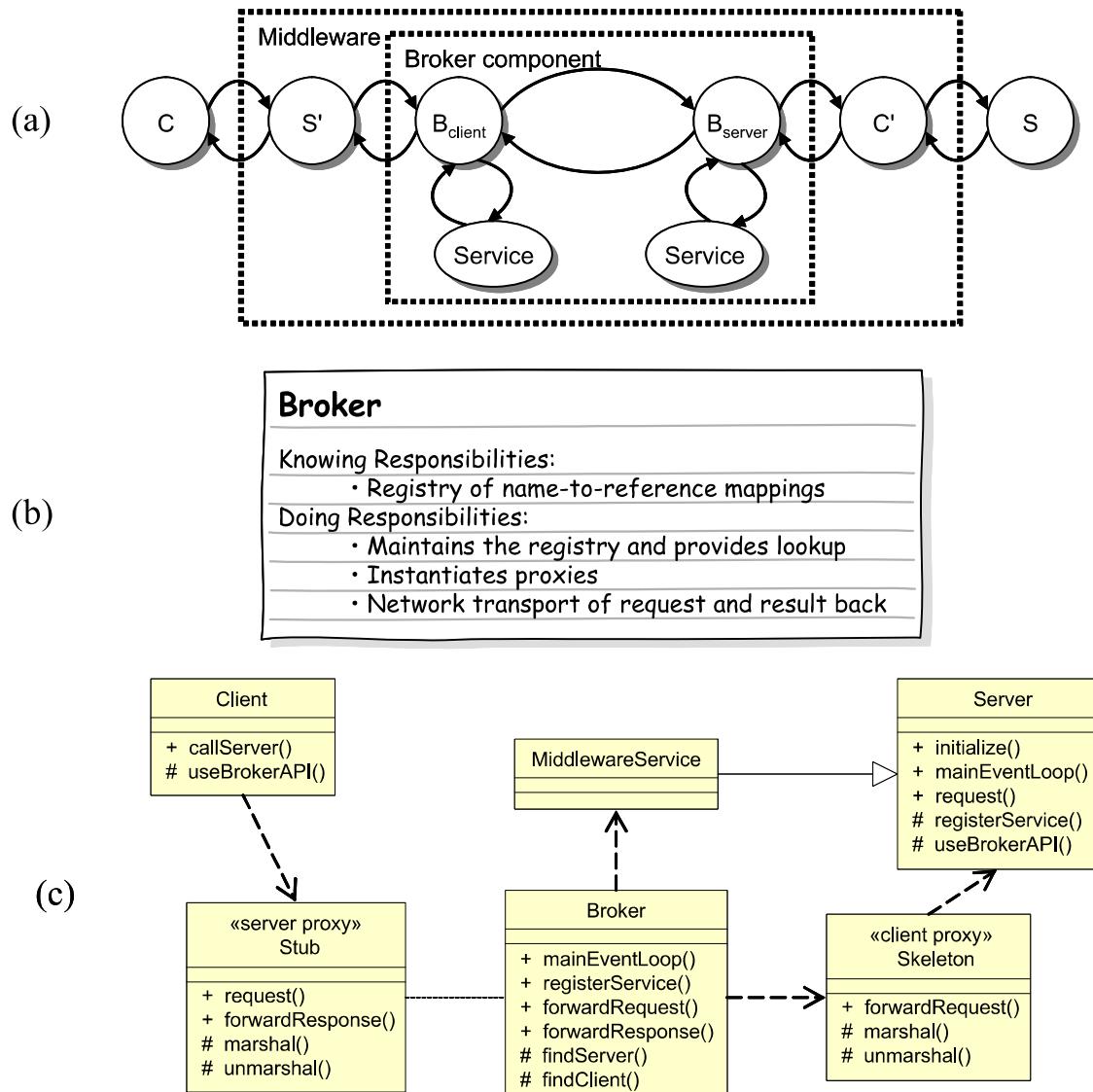


Figure 5-28: (a) The Broker component of middleware represents the intersection of common proxy functions, along with middleware services. (b) Broker's employee card. (c) The Broker pattern class diagram. The server proxy, called *Stub*, resides on the same host as the client and the client proxy, called *Skeleton*, resides on the same host as the server.

5.4.1 Broker Pattern

The Broker pattern is an architectural pattern used to structure distributed software systems with components that interact by remote method calls, see Figure 5-28. The proxies are responsible only for the functions specific to individual objects for which they act as substitutes. In a distributed system the functions that are common to all or most of the proxies are delegated to the *Broker* component, Figure 5-28(b). Figure 5-28(c) shows the objects constituting the Broker pattern and their relationships. The proxies act as representatives of their corresponding objects in

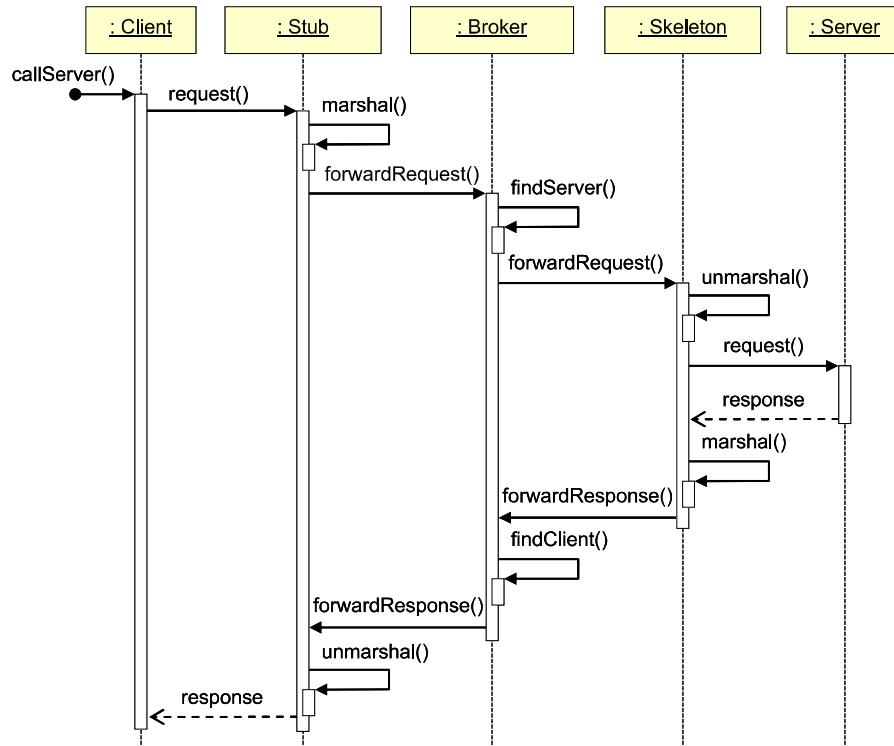


Figure 5-29: Sequence diagram for a client call to the server (remote method invocation).

the foreign address space and contain all the network-related code. The broker component is responsible for coordinating communication and providing links to various middleware services. Although Broker is shown as a single class in Figure 5-28(c), actual brokers consist of many classes.

To use a remote object, a client first *finds* the object through the Broker, which returns a proxy object or *Stub*. As far as the client is concerned, the Stub appears and works like any other local object because they both implement the same interface. But, in reality it only arranges the method call and associated parameters into a stream of bytes using the method `marshal()`. Figure 5-29 shows the sequence diagram for remote method invocation (also called remote procedure call—RPC) via a Broker. The Stub marshals the method call into a stream of bytes and invokes the Broker, which forwards the byte stream to the client's proxy, *Skeleton*, at the remote host. Upon receiving the byte stream, the Skeleton rearranges this stream into the original method call and associated parameters, using the method `unmarshal()`, and invokes this method on the server which contains the actual implementation. Finally, the server performs the requested operation and sends back the return value(s), if any.

The Broker pattern has been proven effective tool in distributed computing, because it leads to greater flexibility, maintainability, and adaptability of the resulting system. By introducing new components and delegating responsibility for communication intricacies, the system becomes potentially distributable and scalable. Java Remote Method Invocation (RMI), which is presented next, is an example of elaborating and implementing the Broker pattern.

5.4.2 Java Remote Method Invocation (RMI)

The above analysis indicates that the Broker component is common to all remotely communicating objects, so it needs to be implemented only once. The proxies are object-specific and need to be implemented for every new object. Fortunately, the process of implementing the proxy objects can be made automatic. It is important to observe that the original object shares the same interface with its Stub and Skeleton (if the object acts as a client, as well). Given the object's interface, there are tools to automatically generate source code for proxy objects. In the general case, the interface is defined in an Interface Definition Language (IDL). Java RMI uses plain Java for interface definition, as well. The basic steps for using RMI are:

1. Define the server object interface
2. Define a class that implements this interface
3. Create the server process
4. Create the client process

Going back to our case-study example of home access control, now I will show how a tenant could remotely connect and download the list of recent accesses.

Step 1: Define the server object interface

The server object will be running on the home computer and currently offers a single method which returns the list of home accesses for the specified time interval:

Listing 5-8: The Informant remote server object interface.

```
/* file Informant.java */
import java.rmi.Remote;
import java.rmi.RemoteException;
import java.util.Vector;

public interface Informant extends Remote {
    public Vector getAccessHistory(long fromTime, long toTime)
        throws RemoteException;
}
```

This interface represents a *contract* between the server and its clients. Our interface must extend `java.rmi.Remote` which is a “tagging” interface that contains no methods. Any parameters of the methods of our interface, such as the return type `Vector` in our case, must be serializable, i.e., implement `java.io.Serializable`.

In the general case, an IDL compiler would generate a Stub and Skeleton files for the `Informant` interface. With Java RMI, we just use the Java compiler, `javac`:

```
% javac Informant.java
```

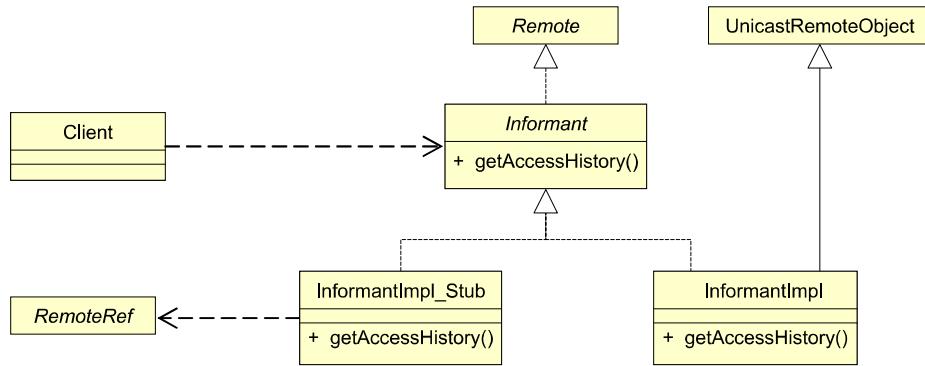


Figure 5-30: Class diagram for the Java RMI example. See text for details.

Step 2: Define a class that implements this interface

The implementation class must extend class `java.rmi.server.RemoteObject` or one of its subclasses. In practice, most implementations extend the subclass `java.rmi.server.UnicastRemoteObject`, because this class supports point-to-point communication using the TCP protocol. The class diagram for this example is shown in Figure 5-30. The implementation class must implement the interface methods and a constructor (even if it has an empty body). I adopt the common convention of adding `Impl` onto the name of our interface to form the implementation class.

Listing 5-9: The class `InformantImpl` implements the actual remote server object.

```

/* file InformantImpl.java (Informant server implementation) */
import java.rmi.RemoteException;
import java.rmi.server.UnicastRemoteObject;
import java.util.Vector;

public class InformantImpl extends UnicastRemoteObject
    implements Informant {
    protected Vector accessHistory_ = new Vector();

    /** Constructor; currently empty. */
    public InformantImpl() throws RemoteException { }

    /** Records all the home access attempts.
     * Called from another class, e.g., from KeyChecker, Listing 2-2
     * @param access Contains the entered key, timestamp, etc.
     */
    public void recordAccess(String access) {
        accessHistory_.add(access);
    }

    /** Implements the "Informant" interface. */
    public Vector getAccessHistory(long fromTime, long toTime)
        throws RemoteException {
        Vector result = new Vector();
        // Extract from accessHistory_ accesses in the
        // interval [fromTime, toTime] into "result"
    }
}
  
```

```

        return result;
    }
}
```

Here, we first use the Java compiler, `javac`, then the RMI compiler, `rmic`:

```
% javac InformantImpl.java
% rmic -v1.2 InformantImpl
```

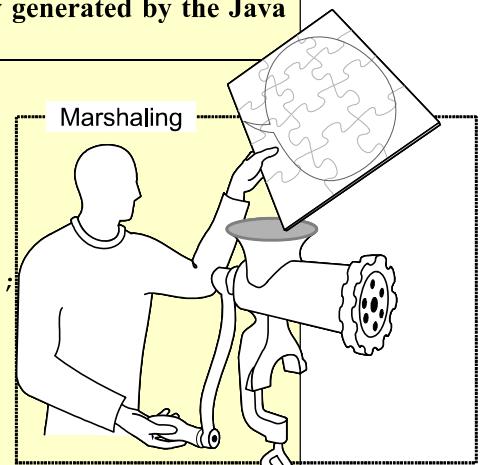
The first statement compiles the Java file and the second one generates the stub and skeleton proxies, called `InformantImpl_Stub.class` and `InformantImpl_Skel.class`, respectively. It is noteworthy, although perhaps it might appear strange, that the RMI compiler operates on a `.class` file, rather than on a source file. In JDK version 1.2 or higher (Java 2), only the stub proxy is used; the skeleton is incorporated into the server itself so that there are no separate entities as skeletons. Server programs now communicate directly with the remote reference layer. This is why the command line option `-v1.2` should be employed (that is, if you are working with JDK 1.2 or higher), so that only the stub file is generated.

As shown in Figure 5-30, the stub is associated with a `RemoteRef`, which is a class in the RMI Broker that represents the handle for a remote object. The stub uses a remote reference to carry out a remote method invocation to a remote object via the Broker. It is instructive to look inside the `InformantImpl_Stub.java`, which is obtained if the RMI compiler is run with the option `-keep`. Here is the stub file (the server proxy resides on client host):

Listing 5-10: Proxy classes (Stub and Skeleton) are automatically generated by the Java `rmic` compiler from the `Informant` interface (Listing 5-8).

```
// Stub class generated by rmic, do not edit.
// Contents subject to change without notice.

1 public final class InformantImpl_Stub
2     extends java.rmi.server.RemoteStub
3     implements Informant, java.rmi.Remote
4 {
5     private static final long serialVersionUID = 2;
6
7     private static java.lang.reflect.Method
$method_getAccessHistory_0;
8
9     static {
10         try {
11             $method_getAccessHistory_0 =
Informant.class.getMethod("getAccessHistory", new java.lang.Class[]
{long.class, long.class});
12         } catch (java.lang.NoSuchMethodException e) {
13             throw new java.lang.NoSuchMethodError(
14                 "stub class initialization failed");
15         }
16     }
17
18     // constructors
19     public InformantImpl_Stub(java.rmi.server.RemoteRef ref) {
20         super(ref);
21     }
22 }
```



```

23     // methods from remote interfaces
24
25     // implementation of getAccessHistory(long, long)
26     public java.util.Vector getAccessHistory(long $param_long_1,
27         long $param_long_2)
28         throws java.rmi.RemoteException
29     {
30         try {
31             Object $result = ref.invoke(this,
32                 $method_getAccessHistory_0, new java.lang.Object[] {new
33                     java.lang.Long($param_long_1), new java.lang.Long($param_long_2)}, -
34                     7208692514216622197L);
35             return ((java.util.Vector) $result);
36         } catch (java.lang.RuntimeException e) {
37             throw e;
38         } catch (java.rmi.RemoteException e) {
39             throw e;
40         } catch (java.lang.Exception e) {
41             throw new java.rmi.UnexpectedException("undclared checked
exception", e);
42         }
43     }
44 }
```

The code description is as follows:

Line 2: shows that our stub extends the `RemoteStub` class, which is the common superclass to client stubs and provides a wide range of remote reference semantics, similar to the broker services in Figure 5-28(a).

Lines 7–15: perform part of the marshaling process of the `getAccessHistory()` method invocation. Computational reflection is employed, which is described in Section 7.3.

Lines 19–21: pass the remote server's reference to the `RemoteStub` superclass.

Line 26: starts the definition of the stub's version of the `getAccessHistory()` method.

Line 30: sends the marshaled arguments to the server and makes the actual call on the remote object. It also gets the result back.

Line 31: returns the result to the client.

The reader should be aware that, in terms of how much of the Broker component is revealed in a stub code, this is only a tip of the iceberg. The Broker component, also known as *Object Request Broker* (ORB), can provide very complex functionality and comprise many software objects.

Step 3: Create the server process

The server process instantiates object(s) of the above implementation class, which accept remote requests. The first problem is, how does a client get handle of such an object, so to be able to invoke a method on it? The solution is for the server to register the implementation objects with a naming service known as *registry*. A naming registry is like a telephone directory. The RMI Registry is a simple name repository which acts as a central management point for Java RMI. The registry must be run before the server and client processes using the following command line:

```
% rmiregistry
```

It can run on any host, including the server's or client's hosts, and there can be several RMI registries running at the same time. (Note: The RMI Registry is an RMI server itself.) For every server object, the registry contains a mapping between the well-known object's name and its reference (usually a globally unique sequence of characters). The process of registration is called ***binding***. The client object can, thus, get handle of a server object by looking up its name in the registry. The lookup is performed by supplying a URL with protocol `rmi`:

```
rmi://host_name:port_number/object_name
```

where `host_name` is the name or IP address of the host on which the RMI registry is running, `port_number` is the port number of the RMI registry, and `object_name` is the name bound to the server implementation object. If the host name is not provided, the default is assumed as `localhost`. The default port number of RMI registry is 1099, although this can be changed as desired. The server object, on the other hand, listens to a port on the server machine. This port is usually an anonymous port that is chosen at runtime by the JVM or the underlying operating system. Or, you can request the server to listen on a specific port on the server machine.

I will use the class `HomeAccessControlSystem`, defined in Listing 2-1, Section 2.7, as the main server class. The class remains the same, except for several modifications:

Listing 5-11: Refactored `HomeAccessControlSystem` class (from Listing 2-1) to instantiate remote server objects of type `Informant`.

```

1 import java.rmi.Naming;
2
3 public class HomeAccessControlSystem extends Thread
4     implements SerialPortEventListener {
5     ...
6     private static final String RMI_REGISTRY_HOST = "localhost";
7
8     public static void main(String[] args) throws Exception {
9         ...
10        InformantImpl temp = new InformantImpl();
11        String rmiObjectName =
12            "rmi://" + RMI_REGISTRY_HOST + "/Informant";
13        Naming.rebind(rmiObjectName, temp);
14        System.out.println("Binding complete...");
15        ...
16    }
17    ...
18 }
```

The code description is as follows:

Lines 3–5: The old `HomeAccessControlSystem` class as defined in Section 2.7.

Line 6: For simplicity's sake, I use `localhost` as the host name, which could be omitted because it is default.

Line 8: The `main()` method now throws `Exception` to account for possible RMI-related exceptions thrown by `Naming.rebind()`.

Line 10: Creates an instance object of the server implementation class.

Lines 11–12: Creates the string URL which includes the object's name, to be bound with its remote reference.

Line 13: Binds the object's name to the remote reference at the RMI naming registry.

Line 14: Displays a message for our information, to know when the binding is completed.

The server is, after compilation, run on the computer located at home by invoking the Java interpreter on the command line:

```
% java HomeAccessControlSystem
```

If Java 2 is used, the skeleton is not necessary; otherwise, the skeleton class, `InformantImpl_Skel.class`, must be located on the server's host because, although not explicitly used by the server, the skeleton will be invoked by Java runtime.

Step 4: Create the client process

The client requests services from the server object. Because the client has no idea on which machine and to which port the server is listening, it looks up the RMI naming registry. What it gets back is a stub object that knows all these, but to the client the stub appears to behave same as the actual server object. The client code is as follows:

Listing 5-12: Client class <code>InformantClient</code> invokes services on a remote <code>Informant</code> object. <pre> 1 import java.rmi.Naming; 2 3 public class InformantClient { 4 private static final String RMI_REGISTRY_HOST = " ... "; 5 6 public static void main(String[] args) { 7 try { 8 Informant grass = (Informant) Naming.lookup(9 "rmi://" + RMI_REGISTRY_HOST + "/Informant" 10); 11 Vector accessHistory = 12 grass.getAccessHistory(fromTime, toTime); 13 14 System.out.println("The retrieved history follows:"); 15 for (Iterator i = accessHistory; i.hasNext();) { 16 String record = (String) i.next(); 17 System.out.println(record); 18 } 19 } catch (ConnectException conEx) { 20 System.err.println("Unable to connect to server!"); 21 System.exit(1); 22 } catch (Exception ex) { 23 ex.printStackTrace(); 24 System.exit(1); 25 } 26 ... 27 } 28 ... 29 }</pre>

The code description is as follows:

Line 4: Specifies the host name on which the RMI naming registry is running.

Line 8: Looks up the RMI naming registry to get handle of the server object. Because the lookup returns a `java.rmi.Remote` reference, this reference must be typecast into an `Informant` reference (not `InformantImpl` reference).

Lines 11–12: Invoke the service method on the stub, which in turn invokes this method on the remote server object. The result is returned as a `java.util.Vector` object named `accessHistory`.

Lines 14–18: Display the retrieved history list.

Lines 19–25: Handle possible RMI-related exceptions.

The client would be run from a remote machine, say from the tenant's notebook or a PDA. The `InformantImpl_Stub.class` file must be located on the client's host because, although not explicitly used by the client, a reference to it will be given to the client by Java runtime. A security-conscious practice is to make the stub files accessible on a website for download. Then, you set the `java.rmi.server.codebase` property equal to the website's URL, in the application which creates the server object (in our example above, this is `HomeAccessControlSystem`). The stubs will be downloaded over the Web, on demand.

The reader should notice that distributed object computing is relatively easy using Java RMI. The developer is required to do just slightly more work, essentially to bind the object with the RMI registry on the server side and to obtain a reference to it on the client side. All the complexity associated with network programming is hidden by the Java RMI tools.

SIDE BAR 5.2: How Transparent Object Distribution?

◆ The reader who experiments with Java RMI, see e.g., Problem 5.22, and tries to implement the same with plain network sockets (see Appendix B), will appreciate how easy it is to work with distributed objects. My recollection from using some CORBA object request brokers was that some provided even greater transparency than Java RMI. Although this certainly is an advantage, there are perils of making object distribution so transparent that it becomes too easy.

The problem is that people tended to forget that there is a network between distributed objects and built applications that relied on fine-grained communication across the network. Too many round-trip communications led to poor performance and reputation problems for CORBA. An interesting discussion of object distribution issues is available in [Waldo *et al.*, 1994] from the same developers who authored Java RMI [Wollrath *et al.*, 1996].

5.5 Information Security

"There is nothing special about security; it's just part of getting the job done." —Rob Short

Information security is a nonfunctional property of the system, it is an *emergent* property. Owing to different types of information use, there are two main security disciplines. *Communication security* is concerned with protecting information when it is being transported between different systems. *Computer security* is related to protecting information within a single system, where it can be stored, accessed, and processed. Although both disciplines must work in accord to successfully protect information, information transport faces greater challenges and so communication security has received greater attention. Accordingly, this review is mainly concerned with communication security. Notice that both communication- and computer security must be complemented with physical (building) security as well as personnel security. Security should be thought of as a chain that is as strong as its weakest link.

The main objectives of information security are:

- *Confidentiality*: ensuring that information is not disclosed or revealed to unauthorized persons
- *Integrity*: ensuring consistency of data, in particular, preventing unauthorized creation, modification, or destruction of data
- *Availability*: ensuring that legitimate users are not unduly denied access to resources, including information resources, computing resources, and communication resources
- *Authorized use*: ensuring that resources are not used by unauthorized persons or in unauthorized ways

To achieve these objectives, we institute various *safeguards*, such as concealing (*encryption*) confidential information so that its meaning is hidden from spying eyes; and *key management* which involves secure distribution and handling of the "keys" used for encryption. Usually, the complexity of one is inversely proportional to that of the other—we can afford relatively simple encryption algorithm with a sophisticated key management method.

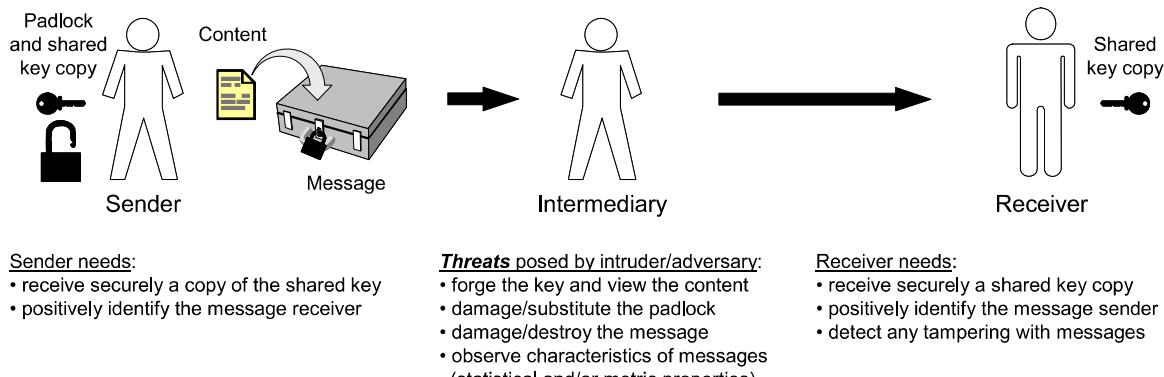


Figure 5-31: Communication security problem: Sender needs to transport a confidential document to Receiver over an untrusted intermediary.

Figure 5-31 illustrates the problem of transmitting a confidential message by analogy with transporting a physical document via untrusted carrier. The figure also lists the security needs of the communicating parties and the potential threats posed by intruders. The sender secures the briefcase, and then sends it on. The receiver must use a correct key in order to unlock the briefcase and read the document. Analogously, a sending computer encrypts the original data using an *encryption algorithm* to make it unintelligible to any intruder. The data in the original form is known as *plaintext* or *cleartext*. The encrypted message is known as *ciphertext*. Without a corresponding “decoder,” the transmitted information (ciphertext) would remain scrambled and be unusable. The receiving computer must regenerate the original plaintext from the ciphertext with the correct *decryption algorithm* in order to read it. This pair of data transformations, encryption and decryption, together forms a *cryptosystem*.

There are two basic types of cryptosystems: (*i*) *symmetric* systems, where both parties use the *same* (secret) key in encryption and decryption transformations; and, (*ii*) *public-key* systems, also known as *asymmetric* systems, where the parties use two related keys, one of which is secret and the other one can be publicly disclosed. I first review the logistics of how the two types of cryptosystems work, while leaving the details of encryption algorithms for the next section.

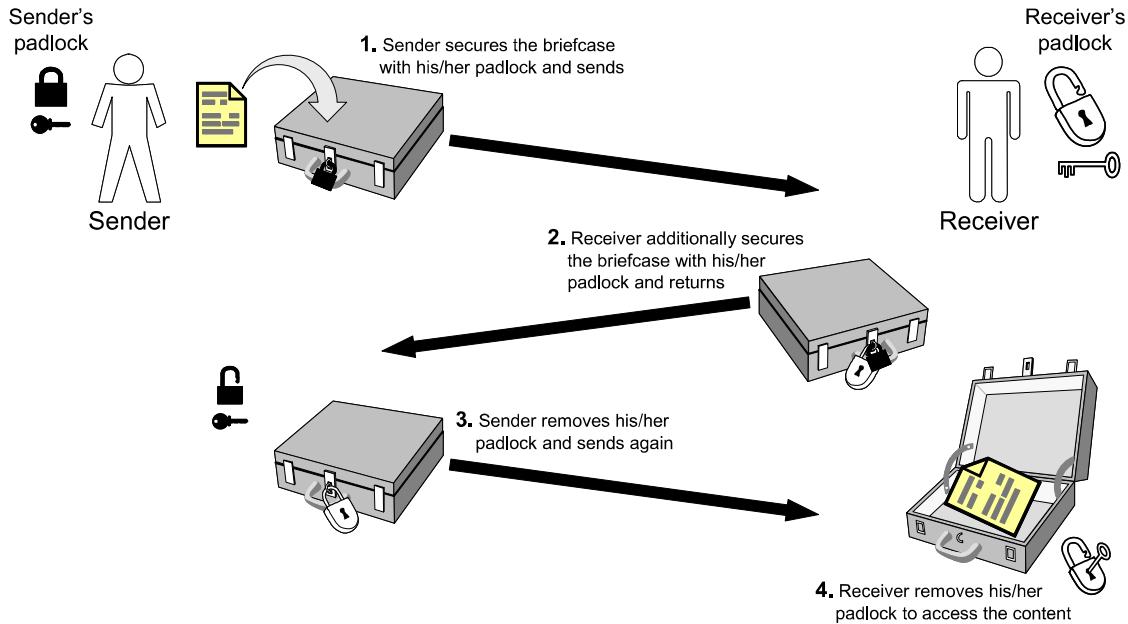


Figure 5-32: Secure transmission via untrustworthy carrier. Note that both sender and receiver keep their own keys with them all the time—the keys are never exchanged.

5.5.1 Symmetric and Public-Key Cryptosystems

In symmetric cryptosystems, both parties use the *same* key in encryption and decryption transformations. The key must remain secret and this, of course, implies trust between the two parties. This is how cryptography traditionally works and prior to the late 1970s, these were the only algorithms available.

The system works as illustrated in Figure 5-31. In order to ensure the secrecy of the shared key, the parties need to meet prior to communication. In this case, the fact that only the parties involved know the secret key implicitly *identifies* one to the other.

Using encryption involves two basic steps: *encoding* a message, and *decoding* it again. More formally, a code takes a message M and produces a coded form $f(M)$. Decoding the message requires an inverse function f^{-1} , such that $f^{-1}(f(M)) = M$. For most codes, anyone who knows how to perform the first step also knows how to perform the second, and it would be unthinkable to release to the adversary the method whereby a message can be turned into code. Merely by “undoing” the encoding procedures, the adversary would be able to break all subsequent coded messages.

In the 1970s Ralph Merkle, Whitfield Diffie, and Martin Hellman realized that this need not be so. The weasel word above is “merely.” Suppose that the encoding procedure is very hard to undo. Then it does no harm to release its details. This led them to the idea of a *trapdoor function*. We call f a trapdoor function if f is easy to compute, but f^{-1} is very hard, indeed impossible for practical purposes. A trapdoor function in this sense is not a very practical code, because the legitimate user finds it just as hard to decode the message as the adversary does. The final twist is

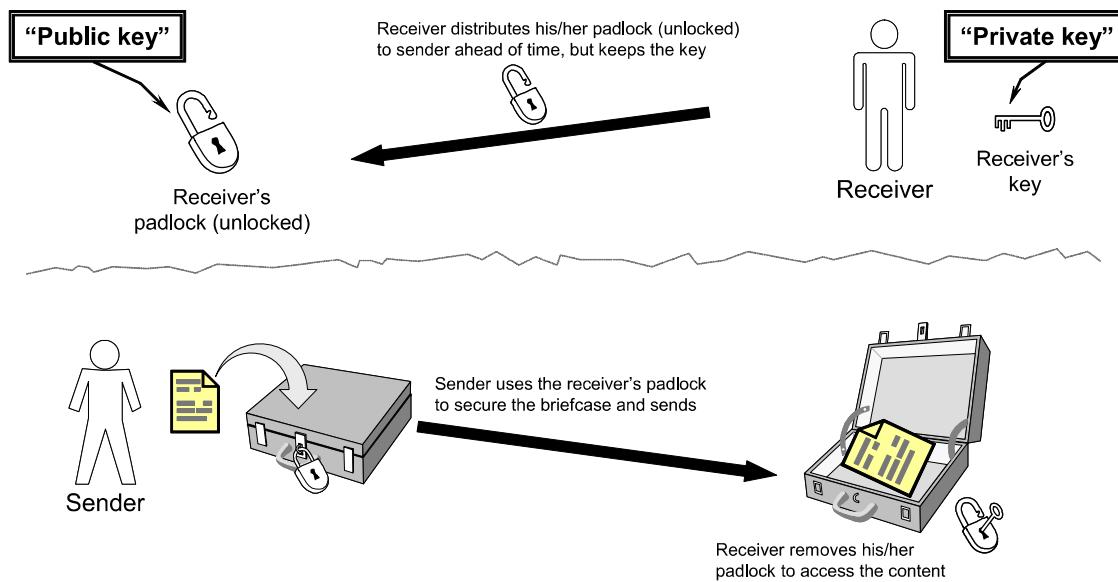


Figure 5-33: Public-key cryptosystem simplifies the procedure from Figure 5-32.

to define f in such a way that a single extra piece of information makes the computation of f^{-1} easy. This is the only bit of information that must be kept secret.

This alternative approach is known as *public-key cryptosystems*. To understand how it works, it is helpful to examine the analogy illustrated in Figure 5-32. The process has three steps. In the first step, the sender secures the briefcase with his or her padlock and sends. Second, upon receiving the briefcase, the receiver secures it additionally with their own padlock and returns. Notice that the briefcase is now doubly secured. Finally, the sender removes his padlock and re-sends. Hence, sender manages to send a confidential document to the receiver without needing the receiver's key or surrendering his or her own key.

There is an inefficiency of sending the briefcase back and forth, which can be avoided as illustrated in Figure 5-33. Here we can skip steps 1 and 2 if the receiver distributed his/her padlock (unlocked, of course!) ahead of time. When the sender needs to send a document, i.e., message, he/she simply uses the receiver's padlock to secure the briefcase and sends. Notice that, once the briefcase is secured, nobody else but receiver can open it, not even the sender. Next I describe how these concepts can be implemented for electronic messages.

5.5.2 Cryptographic Algorithms

Encryption has three aims: keeping data confidential; authenticating who sends data; and, ensuring data has not been tampered with. All cryptographic algorithms involve substituting one thing for another, which means taking a piece of plaintext and then computing and substituting the corresponding ciphertext to create the encrypted message.

Symmetric Cryptography

The Advanced Encryption Standard has a fixed block size of 128 bits and a key size of 128, 192, and 256 bits.

Public-Key Cryptography

As stated above, f is a trapdoor function if f is easy to compute, but f^{-1} is very hard or impossible for practical purposes. An example of such difficulty is factorizing a given number n into prime numbers. An encryption algorithm that depends on this was invented by Ron Rivest, Adi Shamir, and Leonard Adelman (RSA system) in 1978. Another example algorithm, designed by Taher El Gamal in 1984, depends on the difficulty of the discrete logarithm problem.

In the RSA system, the *receiver* does as follows:

1. Randomly select two large prime numbers p and q , which always must be kept secret.
2. Select an integer number E , known as the *public exponent*, such that $(p - 1)$ and E have no common divisors, and $(q - 1)$ and E have no common divisors.
3. Determine the product $n = p \cdot q$, known as *public modulus*.
4. Determine the *private exponent*, D , such that $(E \cdot D - 1)$ is exactly divisible by both $(p - 1)$ and $(q - 1)$. In other words, given E , we choose D such that the integer remainder when $E \cdot D$ is divided by $(p - 1) \cdot (q - 1)$ is 1.
5. Release publicly the **public key**, which is the pair of numbers n and E , $K^+ = (n, E)$. Keep secret the **private key**, $K^- = (n, D)$.

Because a digital message is a sequence of digits, break it into blocks which, when considered as numbers, are each less than n . Then it is enough to encode block by block.

Encryption works so that the sender substitutes each plaintext block B by the ciphertext $C = B^E \% n$, where $\%$ symbolizes the modulus operation. (The *modulus* of two integer numbers x and y , denoted as $x \% y$, is the integer remainder when x is divided by y .)

Then the encrypted message C (ciphertext) is transmitted. To decode, the receiver uses the *decoding key* D , to compute $B = C^D \% n$, that is, to obtain the plaintext B from the ciphertext C .

Example 5.2 RSA cryptographic system

As a simple example of RSA, suppose the receiver chooses $p = 5$ and $q = 7$. Obviously, these are too small numbers to provide any security, but they make the presentation manageable. Next, the receiver chooses $E = 5$, because 5 and $(5 - 1) \cdot (7 - 1)$ have no common factors. Also, $n = p \cdot q = 35$. Finally, the receiver chooses $D = 29$, because $\frac{E \cdot D - 1}{(p-1)(q-1)} = \frac{5 \cdot 29 - 1}{4 \cdot 6} = \frac{144}{24} = 6$, i.e., they are exactly divisible. The receiver's public key is $K^+ = (n, E) = (35, 5)$, which is made public. The private key $K^- = (n, D) = (35, 29)$ is kept secret.

Now, suppose that the sender wants to send the plaintext "hello world." The following table shows the encoding of "hello." I assign to letters a numeric representation, so that $a=1, b=2, \dots, y=25$, and $z=26$, and I assume that block B is one letter long. In an actual implementation, letters are represented as binary numbers, and the blocks B are not necessarily aligned with letters, so that plaintext "l" will not always be represented as ciphertext "12."

Plaintext letter	Plaintext numeric representation	B^E	Ciphertext $B^E \% n$
h	8	$8^5 = 32768$	$8^5 \% 35 = 8$
e	5	$5^5 = 3125$	$5^5 \% 35 = 10$
l	12	$12^5 = 248832$	$12^5 \% 35 = 17$
l	12	248832	17
o	15	$15^5 = 759375$	$15^5 \% 35 = 15$

The sender transmits this ciphertext to the receiver: 8, 10, 17, 17, 15. Upon the receipt, the receiver decrypts the message as shown in the following table.

Ciphertext	C^D	$B = C^D \% n$	Plaintext letter
8	$8^{29} = 154742504910672534362390528$	$8^{29} \% 35 = 8$	h
10	1000000000000000000000000000000000000	5	e
17	481968572106750915091411825223071697	12	l
17	481968572106750915091411825223071697	12	l
15	12783403948858939111232757568359375	15	o

We can see that even this toy example produces some extremely large numbers.

The point is that while the adversary knows n and E , he or she does not know p and q , so they cannot work out $(p - 1) \cdot (q - 1)$ and thereby find D . The designer of the code, on the other hand, knows p and q because those are what he started from. So does any legitimate receiver: the designer will have told them. The security of this system depends on exactly one thing: the notoriously difficulty of factorizing a given number n into primes. For example, given $n = 2^{67} - 1$ it took three years working on Sundays for F. N. Cole to find by hand in 1903 p and q for $n = p \cdot q = 193707721 \times 761838257287$. It would be waste of time (and often combinatorially self-defeating) for the program to grind through all possible options.

Encryption strength is measured by the length of its “key,” which is expressed in bits. The larger the key, the greater the strength of the encryption. Using 112 computers, a graduate student decrypted one 40-bit encrypted message in a little over 7 days. In contrast, data encrypted with a 128-bit key is 309,485,009,821,345,068,724,781,056 times stronger than data encrypted with a 40-bit key. RSA Laboratories recommends that the product of p and q be on the order of 1024 bits long for corporate use and 768 bits for use with less valuable information.

5.5.3 Authentication

5.5.4 Program Security

A virus is malicious code carried from one computer to another by some medium—often an “infected” file. Any operating system that allows third-party programs to run can theoretically run viruses. Some operating systems are more secure than others; earlier versions of Microsoft

Windows did not even provide something as simple as maintain memory space separation. Once on a computer, a virus is executed when its carrier file is “opened” in some meaningful way by software on that system. When the virus executes, it does something unwanted, such as causing software on the host system to send more copies of infected files to other computers over the network, infecting more files, and so on. In other words, a virus typically maximizes its likelihood of being passed on, making itself contagious.

Viral behavior relies on security vulnerabilities that exist in software running on the host system. For example, in the past, viruses have often exploited security vulnerabilities in Microsoft Office macro scripting capabilities. Macro viruses are no longer among the most common virus types. Many viruses take advantage of Trident, the rendering engine behind Internet Explorer and Windows Explorer that is also used by almost every piece of Microsoft software. Windows viruses often take advantage of image-rendering libraries, SQL Server’s underlying database engine, and other components of a complete Windows operating system environment as well.

Viruses are typically addressed by antivirus software vendors. These vendors produce virus definitions used by their antivirus software to recognize viruses on the system. Once a specific virus is detected, the software attempts to quarantine or remove the virus—or at least inform the user of the infection so that some action may be taken to protect the system from the virus.

This method of protection relies on knowledge of the existence of a virus, however, which means that most of the time a virus against which you are protected has, by definition, already infected someone else’s computer and done its damage. The question you should be asking yourself at this point is how long it will be until you are the lucky soul who gets to be the discoverer of a new virus by way of getting infected by it.

It’s worse than that, though. Each virus exploits a vulnerability — but they don’t all have to exploit different vulnerabilities. In fact, it’s common for hundreds or even thousands of viruses to be circulating “in the wild” that, between them, only exploit a handful of vulnerabilities. This is because the vulnerabilities exist in the software and are not addressed by virus definitions produced by antivirus software vendors.

These antivirus software vendors’ definitions match the signature of a given virus — and if they’re really well-designed might even match similar, but slightly altered, variations on the virus design. Sufficiently modified viruses that exploit the same vulnerability are safe from recognition through the use of virus definitions, however. You can have a photo of a known bank robber on the cork bulletin board at the bank so your tellers will be able to recognize him if he comes in — but that won’t change the fact that if his modus operandi is effective, others can use the same tactics to steal a lot of money.

By the same principle, another virus can exploit the same vulnerability without being recognized by a virus definition, as long as the vulnerability itself isn’t addressed by the vendor of the vulnerable software. This is a key difference between open source operating system projects and Microsoft Windows: Microsoft leaves dealing with viruses to the antivirus software vendors, but open source operating system projects generally fix such vulnerabilities immediately when they’re discovered.

Thus, the main reason you don’t tend to need antivirus software on an open source system, unless running a mail server or other software that relays potentially virus-laden files between other systems, isn’t that nobody’s targeting your open source OS; it’s that any time someone targets it,

chances are good that the vulnerability the virus attempts to exploit has been closed up — even if it's a brand-new virus that nobody has ever seen before. Any half-baked script-kiddie has the potential to produce a new virus that will slip past antivirus software vendor virus definitions, but in the open source software world one tends to need to discover a whole new vulnerability to exploit before the “good guys” discover and patch it.

Viruses need not simply be a “fact of life” for anyone using a computer. Antivirus software is basically just a dirty hack used to fill a gap in your system’s defenses left by the negligence of software vendors who are unwilling to invest the resources to correct certain classes of security vulnerabilities.

The truth about viruses is simple, but it’s not pleasant. The truth is that you’re being taken to the cleaners — and until enough software users realize this, and do something about it, the software vendors will continue to leave you in this vulnerable state where additional money must be paid regularly to achieve what protection you can get from a dirty hack that simply isn’t as effective as solving the problem at the source would be.

However, we should not forget that security comes at a cost.

In theory, application programs are supposed to access hardware of the computer only through the interfaces provided by the operating system. But many application programmers who dealt with small computer operating systems of the 1970s and early 1980s often bypassed the OS, particularly in dealing with the video display. Programs that directly wrote bytes into video display memory run faster than programs that didn’t. Indeed, for some applications—such as those that needed to display graphics on the video display—the operating system was totally inadequate.

What many programmers liked most about MS-DOS was that it “stayed out of the way” and let programmers write programs as fast as the hardware allowed. For this reason, popular software that ran on the IBM PC often relied upon idiosyncrasies of the IBM PC hardware.

5.6 Summary and Bibliographical Notes

Design patterns are heuristics for structuring the software modules and their interactions that are proven in practice. They yield in design for change, so the change of the computing environment has as minimal and as local effect on the code as possible.

Key Points:

- Pattern use must be need-driven: use a pattern only when you need it to improve your software design, not because it can be used, or you simply like hitting nails with your new hammer.

- Using the Broker pattern, a client object invokes methods of a remote server object, passing arguments and receiving a return value with each call, using syntax similar to local method calls. Each side requires a proxy that interacts with the system's runtime.

There are many known design patterns and I have reviewed above only few of the major ones. The text that most contributed to the popularity of patterns is [Gamma *et al.*, 1995]. Many books are available, perhaps the best known are [Gamma *et al.*, 1995] and [Buschmann *et al.*, 1996]. The reader can also find a great amount of useful information on the web. In particular, a great deal of information is available in Hillsides.net's Patterns Library: <http://hillside.net/patterns/>.

R. J. Wirfs-Brock, "Refreshing patterns," *IEEE Software*, vol. 23, no. 3, pp. 45-47, May/June 2006.

Section 5.1: Indirect Communication: Publisher-Subscriber

Section 5.2: More Patterns

Section 5.3: Concurrent Programming

Concurrent systems are a large research and practice field and here I provide only the introductory basics. Concurrency methods are usually not covered under design patterns and it is only for the convenience sake that here they appear in the section on software design patterns. I avoided delving into the intricacies of Java threads—by no means is this a reference manual for Java threads. Concurrent programming in Java is extensively covered in [Lea, 2000] and a short review is available in [Sandén, 2004].

[Whiddett, 1987]

Pthreads tutorial: <http://www.cs.nmsu.edu/~jcook/Tools/pthreads/pthreads.html>

Pthreads tutorial from CS6210 (by Phillip Hutto):
http://www.cc.gatech.edu/classes/AY2000/cs6210_spring/pthreads_tutorial.htm

Section 5.4: Broker and Distributed Computing

The *Broker* design pattern is described in [Buschmann *et al.*, 1996; Völter *et al.*, 2005].

Java RMI:

Sun Developer Network (SDN) jGuru: "Remote Method Invocation (RMI)," Sun Microsystems, Inc., Online at: <http://java.sun.com/developer/onlineTraining/rmi/RMI.html>

<http://www.javaworld.com/javaworld/jw-04-2005/jw-0404-rmi.html>

http://www.developer.com/java/ent/article.php/10933_3455311_1

Although Java RMI works only if both client and server processes are coded in the Java programming language, there are other systems, such as CORBA (Common Object Request Broker Architecture), which work with arbitrary programming languages, including Java. A readable appraisal of the state of affairs with CORBA is available in [Henning, 2006].

Section 5.5: Information Security

In an increasingly networked world, all computer users are at risk of having their personally identifying information and private access data intercepted. Even if information is not stolen, computing resources may be misused for criminal activities facilitated by unauthorized access to others' computer systems.

Kerckhoffs' Principle states that a cryptosystem should remain secure even if everything about it other than the key is public knowledge. The security of a system's design is in no way dependent upon the secrecy of the design, in and of itself. Because system designs can be intercepted, stolen, sold, independently derived, reverse engineered by observations of the system's behavior, or just leaked by incompetent custodians, the secrecy of its design can never really be assumed to be secure itself. Hence, the "security through obscurity" security model by attempting to keep system design secret Open source movement even advocates widespread access to the design of a system because more people can review the system's design and detect potential problems. Transparency ensures that the security problems tend to arise more quickly, and to be addressed more quickly. Although an increased likelihood of security provides no guarantees of success, it is beneficial nonetheless.

There is an entire class of software, known as "fuzzers," that is used to quickly detect potential security weaknesses by feeding abusive input at a target application and observing its behavior under that stress. These are the tools that malicious security crackers use all the time to find ways to exploit software systems. Therefore, it is not necessary to have access to software design (or its source code) to be able to detect its security vulnerabilities. This should not be surprising, given that software defects are rarely found by looking at source code. (Recall the software testing techniques from Section 2.7.) Where access to source code becomes much more important is when trying to determine *why* a particular weakness exists, and how to remove it. One might conclude, then, that the open source transparency does not contribute as much to detecting security problems as it does to fixing them.

Cryptography [Menezes *et al.*, 1997], which is available for download, entirely, online at <http://www.cacr.math.uwaterloo.ca/hac/>.

ICS 54: History of Public-Key Cryptography:
<http://www.ics.uci.edu/~ics54/doc/security/pkhistory.html>
<http://www.netip.com/articles/keith/diffie-helman.htm>
<http://www.rsasecurity.com/rsalabs/node.asp?id=2248>
<http://www.scramdisk.clara.net/pgpfaq.html>
<http://postdiluvian.org/~seven/diffie.html>
<http://www.sans.org/rr/whitepapers/vpns/751.php>
<http://www.fors.com/eoug97/papers/0356.htm>

Class iaik.security.dh.DHKeyAgreement

<http://www.cs.utexas.edu/users/chris/cs378/f98/resources/iaikdocs/iaik.security.dh.DHKeyAgreement.html>

Bill Steele, “‘Fabric’ would tighten the weave of online security,” *Cornell Chronicle* (09/30/10): Fabric’s programming language, which is based on Java, builds in security as the program is written. Myers says most of what Fabric does is transparent to the programmer.

<http://www.news.cornell.edu/stories/Sept10/Fabric.html>

P. Dourish, R. E. Grinter, J. Delgado de la Flor, and M. Joseph, “Security in the wild: user strategies for managing security as an everyday, practical problem,” *Personal and Ubiquitous Computing (ACM/Springer)*, vol. 8, no. 6, pp. 391-401, November 2004.

M. J. Ranum, “Security: The root of the problem,” *ACM Queue (Special Issue: Surviving Network Attacks)*, vol. 2, no. 4, pp. 44-49, June 2004.

H. H. Thompson and R. Ford, “Perfect storm: The insider, naivety, and hostility,” *ACM Queue (Special Issue: Surviving Network Attacks)*, vol. 2, no. 4, pp. 58-65, June 2004.
introducing trust and its pervasiveness in information technology

Microsoft offers integrated hardware-level security such as data execution prevention, kernel patch protection and its free Security Essentials software:

http://www.microsoft.com/security_essentials/

Microsoft's 'PassPort' Out, Federation Services In

In 2004 Microsoft issued any official pronouncements on "TrustBridge," its collection of federated identity-management technologies slated to go head-to-head with competing technologies backed by the Liberty Alliance.

<http://www.eweek.com/c/a/Windows/Microsofts-Passport-Out-Federated-Services-In/>

Problems

Problem 5.1

Problem 5.2

Consider the online auction site described in Problem 2.31 (Chapter 2). Suppose you want to employ the Publish-Subscribe (also known as Observer) design pattern in your design solution for Problem 2.31. Which classes should implement the Publisher interface? Which classes should

implement the Subscriber interface? Explain your answer. (Note: You can introduce new classes or additional methods on the existing classes if you feel it necessary for solution.)

Problem 5.3

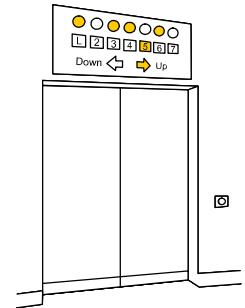
In the patient-monitoring scenario of Problem 2.35 (Chapter 2), assume that multiple recipients must be notified about the patient condition. Suppose that your software is to use the Publish-Subscribe design pattern. Identify the key software objects and draw a UML interaction diagram to represent how software objects in the system could accomplish the notification problem.

Problem 5.4

Problem 5.5

Problem 5.6: Elevator Control

Consider the elevator control problem defined in Problem 3.7 (Chapter 3). Your task is to determine whether the Publisher-Subscriber design pattern can be applied in this design. Explain clearly your answer. If the answer is yes, identify which classes are suitable for the publisher role and which ones are suitable for the subscriber role. Explain your choices, list the events generated by the Publishers, and state explicitly for each Subscriber to which events it is subscribed to.

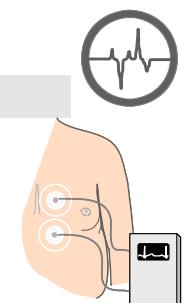


Problem 5.7

Problem 5.8

Problem 5.9

Consider the automatic patient monitoring system described in Problem 2.35. Carefully examine the draft UML sequence diagram in Figure 2-45. Check if the given design already uses some patterns and explain your claim. Identify as many opportunities as you can to improve the design by applying design patterns. Consider how an unnecessary application of some design patterns would make this design worse. Draw UML sequence diagrams or write pseudo-code to describe the proposed design. Always describe your motivation for adopting or rejecting design modifications.



Problem 5.10



Consider the system for inventory management grocery supermarket from Problem 2.15. Suppose you are provided with an initial software design as follows. This design is based on a basic version of the inventory system, but the reader should be aware of extensions that are discussed in the solution of Problem 2.15(c) and Problem 2.16. The software consists of the following classes:

ReaderIface:

This class receives messages from RFID readers that specific tags moved in or out of coverage.

DBaseConn:

This class provides a connection to a relational database that contains data about shelf stock and inventory tasks. The database contains several tables, including ProductsInfo[key = tagID], PendingTasks[key = userID], CompletedTasks, and Statistics[key = infoType] for various information types, such as the count of erroneous messages from RFID readers and the count of reminders sent for individual pending tasks.

Dispatcher:

This class manages inventory tasks by opening new tasks when needed and generates notifications to the concerned store employees.

Monitor:

This class periodically keeps track of potentially overdue tasks. It retrieves the list of pending tasks from the database and generates reminders to the concerned store employees.

Messenger:

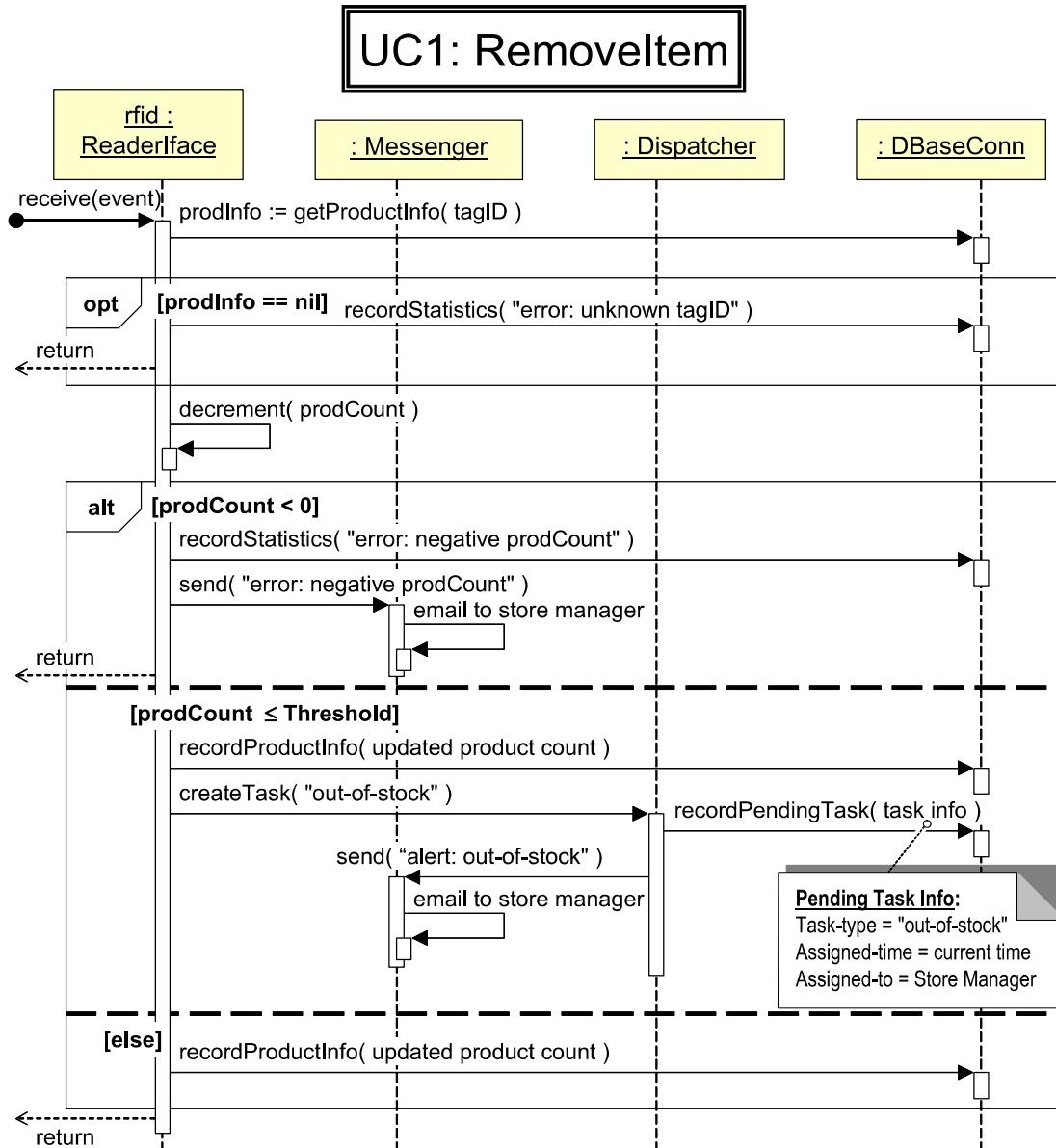
This class sends email notifications to the concerned store employees. (The notifications are generated by other classes.)

Assume that email notifications are used as a supplementary tool, but the system must keep an internal record of sent notifications and pending tasks, so it can take appropriate actions.

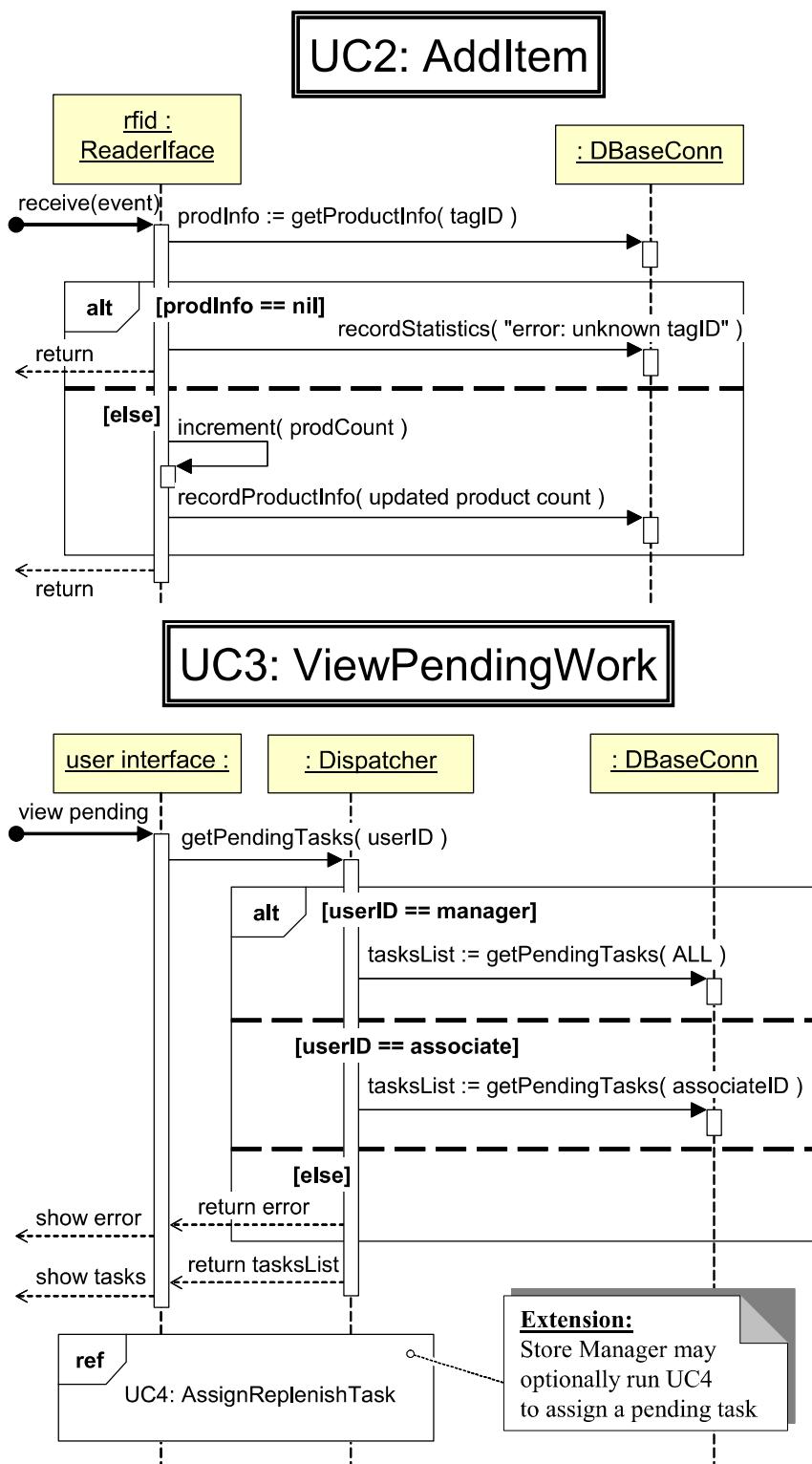
Notice that the current design has a single timer for the whole system. The software designer noticed that sending notifications for overdue tasks does not need to be exactly timed in this system. Delays up to a certain period (e.g., hour or even day) are tolerable. Maintaining many timers would be overkill and would significantly slow down the system. It would not be able to do important activities, such as processing RFID events in a timely fashion. Therefore, the software is designed so that, when a new pending task is created, there is no explicit activation of an associated timer. Instead, the task is simply added to the list of pending tasks. The Monitor object periodically retrieves this list and checks for overdue tasks, as seen below in the design for the use case UC-5 SendReminder.

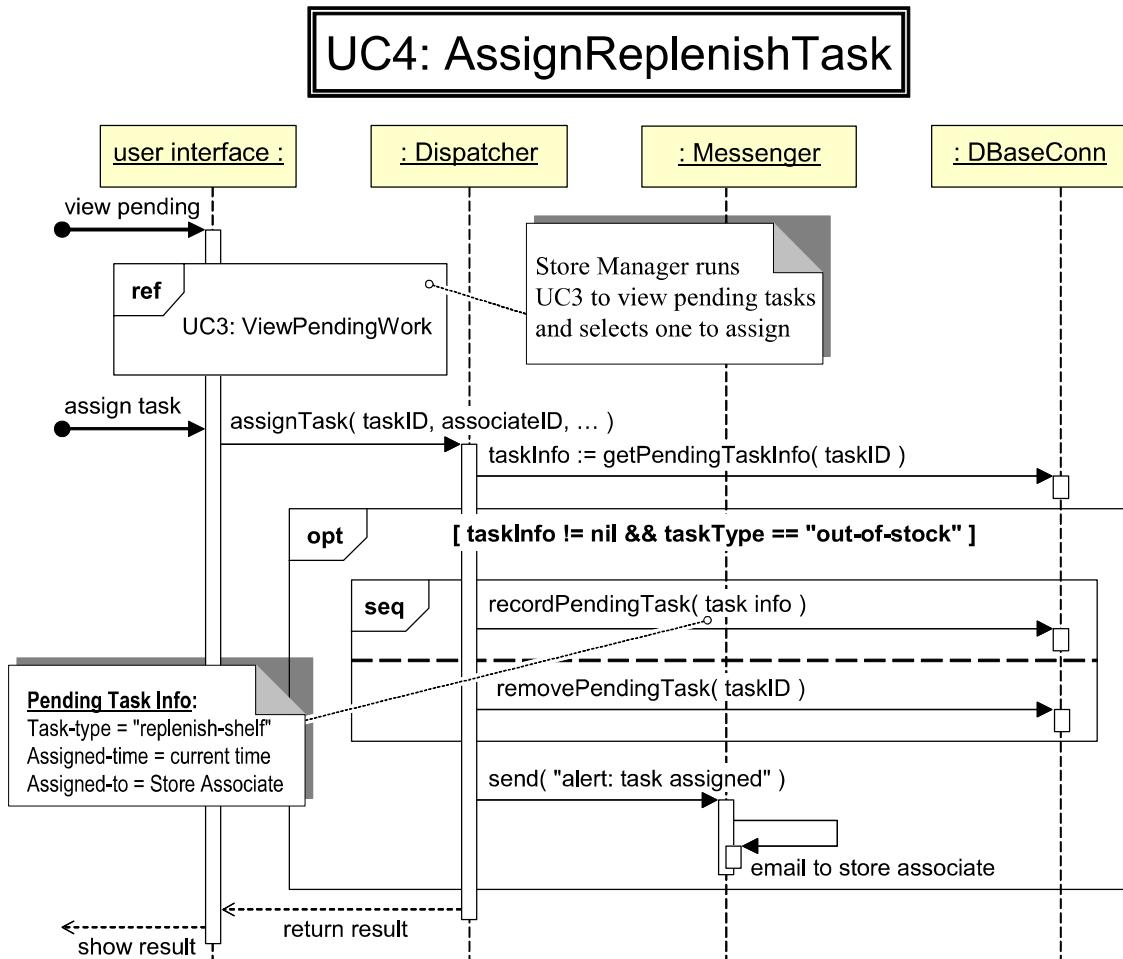
Another simplification is to check only for “out-of-stock” events and not for “low-stock” events. If the customer demands that “low-stock” events be included, then the design of the software-to-be will become somewhat more complex.

The UML sequence diagrams for all the use cases are shown in the following figures. Notice that use cases UC-3, UC-4, and UC-6 «include» UC-7: Login (user authentication), which is not shown to avoid clutter.

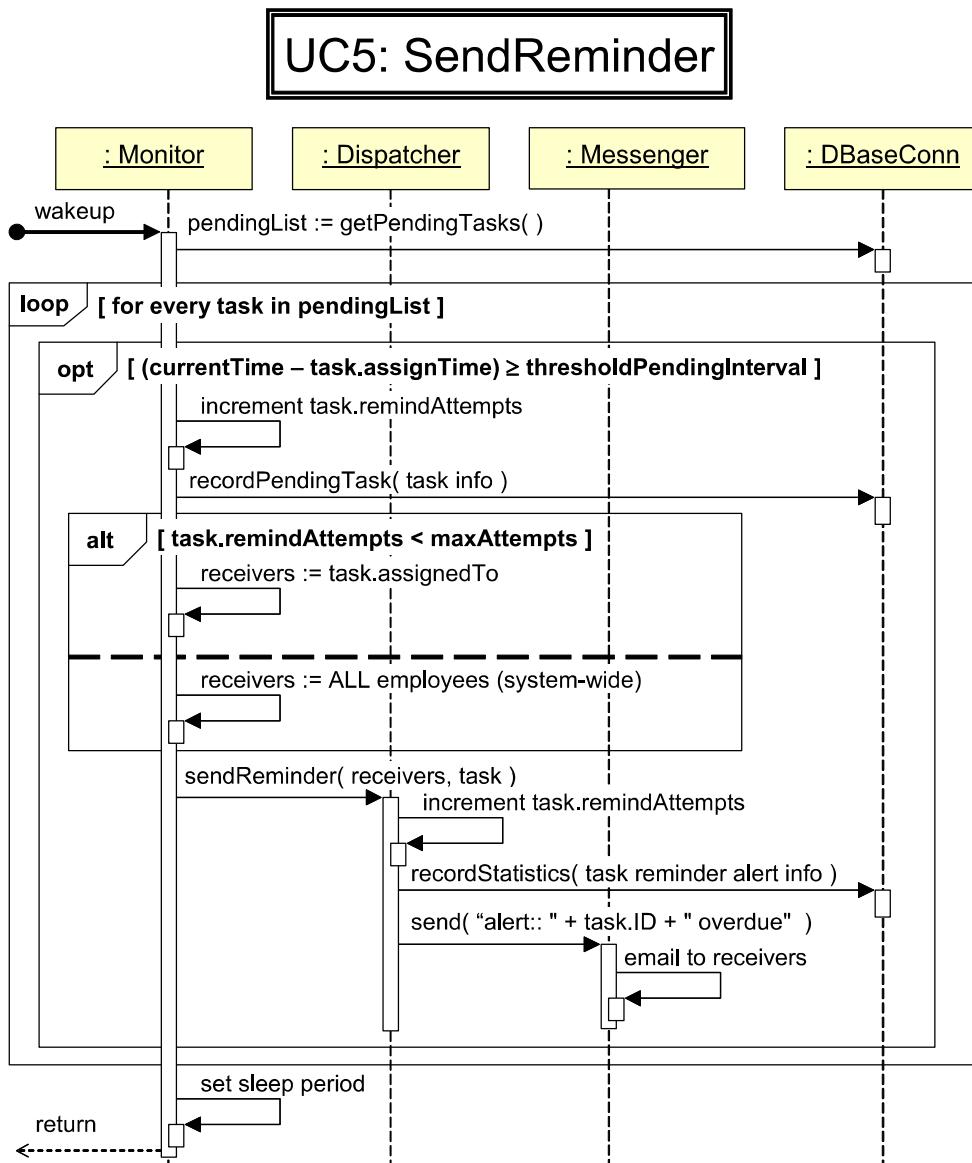


In the design for UC-1, the system may check if a pending task for the given product already exists in the database; if yes, it should not generate a new pending task for the same product.





The [seq] interaction fragment specifies that the interactions contained within the fragment box must occur exactly in the given order. The reason for using this constraint in UC-4 is that the system may crash while the task is being converted from unassigned to pending. If `removePendingTask()` were called first, and the system crashed after it but before `recordPendingTask()`, then all information about this task would be lost! Depending on the database implementation, it may be possible to perform these operations as atomic for they update the same table in the database. To deal with crash scenarios where a task ends up in both tables, the Monitor object in UC-5 SendReminder should be extended to perform a database clean-up after a crash. It should remove those tasks from the PendingTasks table that are marked both as unassigned and pending.

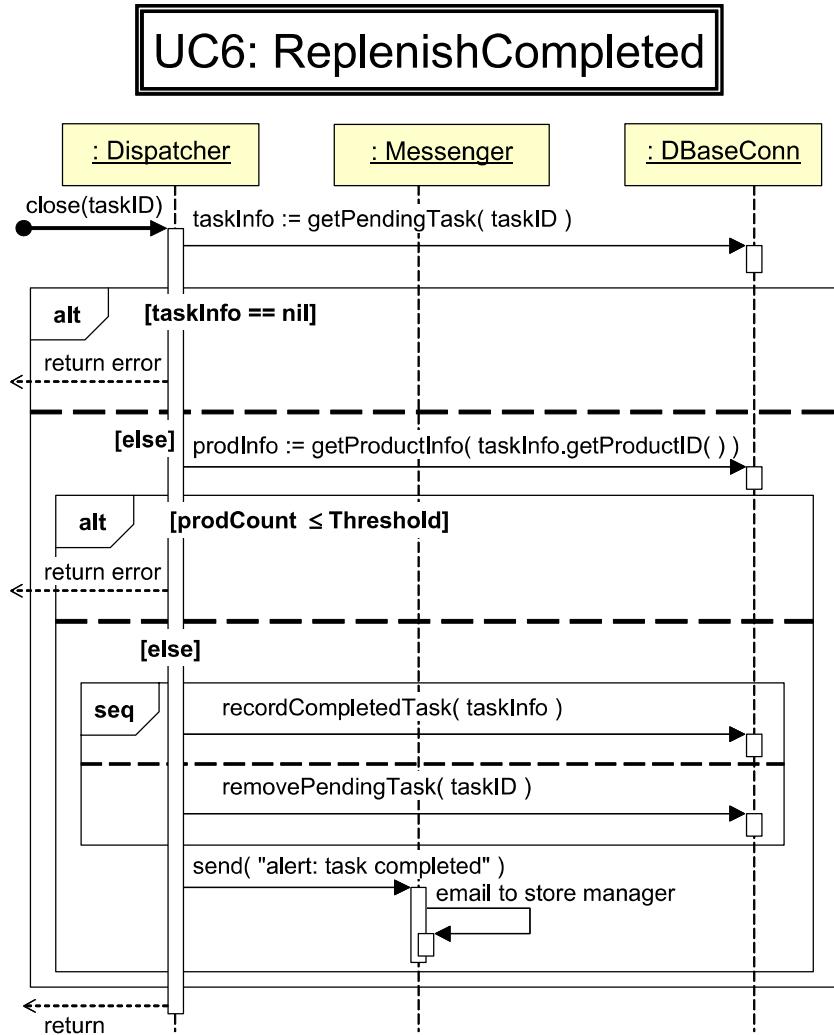


Note that the Monitor discards the list of pending tasks before going to sleep, so it starts every cycle with a fresh list of pending tasks, retrieved from the database, because our assumption is that the database contains the most current information (possibly updated by other objects).

By examining the design for the use case UC-5 SendReminder, we see that the Monitor has to do a lot of subtractions and comparisons every time it wakes up, but this can be done at leisure because seconds or minutes are not critical for this activity. The computing power should be better used for other use cases. Of course, we must ensure that the Monitor still works fast enough not to introduce delays on the order of hours or days during each cycle!

In addition, we need to handle the case where the time values are not incrementing constantly upwards, such as when a full 24 hours passes and a new day starts, the time resets to zero. In Java, using `java.lang.System.currentTimeMillis()` returns the current time in milliseconds as a long integer.

In the solution of Problem 2.16 we discussed using adaptive timeout calculation to adjust the frequency of reminders for busy periods. Another option is to have shorter sleep periods, but during each wakeup, process only part of the list of pending tasks, and leave the rest for subsequent wakeups. Then cycle again from the head of the list. This way, the reminders will be spread over time and not all reminders will be generated at once (avoid generating one “bulk” notification each period).



The logic of UC-6 is that it first retrieves the task, checks if such a task exists, makes sure it is really done, and finally marks it as completed. The `[seq]` interaction fragment specifies that the interactions contained within the fragment box must occur exactly in the given order. Similar to UC-4 AssignReplenishTask, this constraint is needed in case the system crashes while the task is being closed. If `removePendingTask()` were called first, and the system crashed after it but before `recordCompletedTask()`, then all information about this task would be lost! These operations cannot be performed as atomic, because they work on different tables in the database. To deal with crash scenarios where a task ends up in both tables, the Monitor object in UC-5 should be modified to perform a database clean-up after a crash. It should remove those tasks from the PendingTasks table that are already in the CompletedTasks table.

Notice also that the Monitor runs in a separate thread, so while UC-6 is in the process of closing a task, the Monitor may send an unnecessary reminder about this task (in UC-5).

Carefully examine the existing design and identify as many opportunities as you can to improve the design by applying design patterns. Note that the existing design ignores the issue of concurrency, but we will leave the multithreading issue aside for now and focus only on the patterns that improve the quality of software design. (The concurrency issues will be considered later in Problem 5.20.)

- (a) If you introduce a pattern, first provide arguments why the existing design may be problematic.
- (b) Provide as much details as possible about how the pattern will be implemented and how the new design will work (draw UML sequence diagrams or write pseudo-code).
- (c) Explain how the pattern improved the design (i.e., what are the expected benefits compared to the original design).

If considering future evolution and extensions of the system when proposing a modification, then describe explicitly what new features will likely be added and how the existing design would be inadequate to cope with resulting changes. Then introduce a design pattern and explain how the modified version is better.

If you believe that the existing design (or some parts of it) is sufficiently good then explain how the application of some design patterns would make the design worse. Use concrete examples and UML diagrams or pseudo-code to illustrate and refer to specific qualities of software design.

Problem 5.11

Problem 5.12

Problem 5.13

Problem 5.14

Problem 5.15

In Section 5.3, it was stated that the standard Java idiom for condition synchronization is the statement:

```
while (condition) sharedObject.wait();
```

- (a) Is it correct to substitute the `yield()` method call for `wait()`? Explain your answer and discuss any issues arising from the substitution.

(b) Suppose that `if` substitutes for `while`, so we have:

```
if (condition) sharedObject.wait()
```

Is this correct? Explain your answer.

Problem 5.16

Parking lot occupancy monitoring, see Figure 5-34. Consider a parking lot with the total number of spaces equal to `capacity`. There is a single barrier gate with two poles, one for the entrance and the other for the exit. A computer in the barrier gate runs a single program which controls both poles. The program counts the current number of free spaces, denoted by `occupancy`, such that

$$0 \leq \text{occupancy} \leq \text{capacity}$$

When a new car enters, the occupancy is incremented by one; conversely, when a car exits, the occupancy is decremented by one. If `occupancy` equals `capacity`, the red light should turn on to indicate that the parking is full.

In order to be able to serve an entering and an exiting patron in parallel, you should design a system which runs in two threads. `EnterThread` controls the entrance gate and `ExitThread` controls the exit gate. The threads share the `occupancy` counter so to correctly indicate the parking-full state. Complete the UML sequence diagram in Figure 5-35 that shows how the two threads update the shared variable, i.e., `occupancy`.

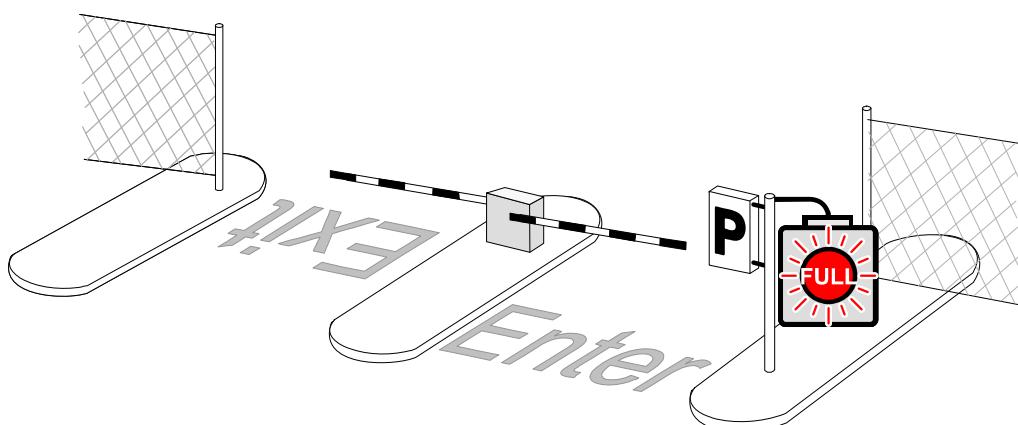


Figure 5-34: Parking lot occupancy monitoring, Problem 5.16.

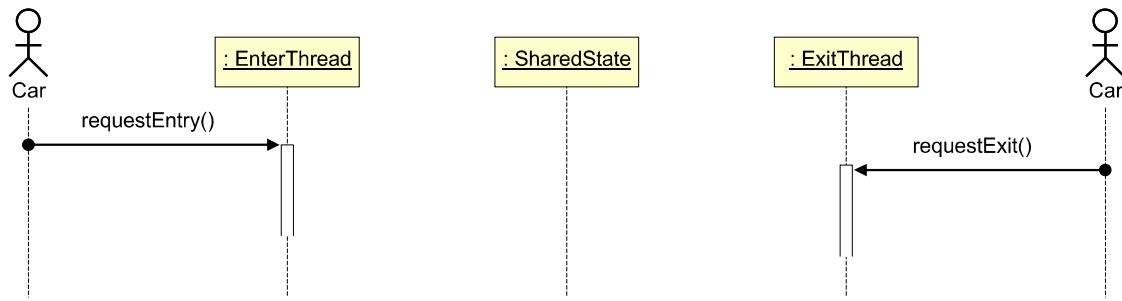


Figure 5-35: UML diagram template for parking lot occupancy monitoring, Problem 5.16.

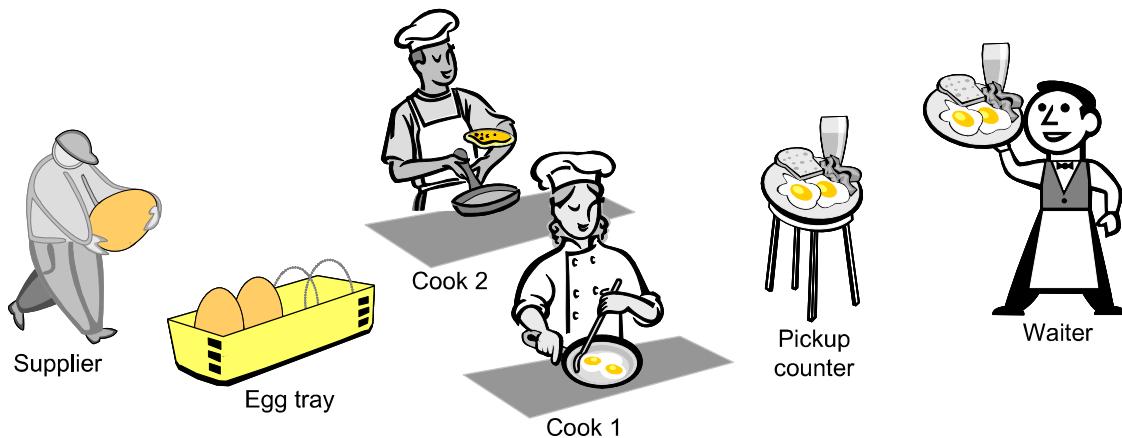


Figure 5-36: Concurrency problem in a restaurant scenario, Problem 5.17.

Hint: Your key concern is to maintain the consistent shared state (occupancy) and indicate when the parking-full sign should be posted. Extraneous actions, such as issuing the ticket for an entering patron and processing the payment for an exiting patron, should not be paid attention—only make a high-level remark where appropriate.

Problem 5.17

Consider a restaurant scenario shown in Figure 5-36. You are to write a simulation in Java such that each person runs in a different thread. Assume that each person takes different amount of time to complete their task. The egg tray and the pickup counter have limited capacities, N_{eggs} and N_{plates} , respectively. The supplier stocks the egg tray but must wait if there are no free slots. Likewise, the cooks must hold the prepared meal if the pickup counter is full.

Problem 5.18

A *priority inversion* occurs when a higher-priority thread is waiting for a lower-priority thread to finish processing of a critical region that is shared by both. Although higher-priority threads normally preempt lower-priority threads, this is not possible when both share the same critical region. While the higher-priority thread is waiting, a third thread, whose priority is between the first two, but it does not share the critical region, preempts the low-priority thread. Now the

higher-priority thread is waiting for more than one lower-priority thread. Search the literature and describe precisely a possible mechanism to avoid priority inversion.

Problem 5.19



Assume that the patient device described in Problem 2.3 (at the end of Chapter 2) runs in a multi-threaded mode, where different threads acquire and process data from different sensors. (See also Problem 2.35 and its solution on the back of this book.) What do you believe is the optimal number of threads? When designing this system, what kind of race conditions or other concurrency issues can you think of? Propose a specific solution for each issue that you identify (draw UML sequence diagrams or write pseudo-code).

Problem 5.20

Consider the supermarket inventory management system from Problem 5.10. A first observation is that the existing design ignores the issue of concurrency—there will be many users simultaneously removing items, and/or several associates may be simultaneously restocking the shelves. Also, it is possible that several employees may simultaneously wish to view pending tasks, assign replenishment tasks, or report replenishment completed. Clearly, it is necessary to introduce multithreading even if the present system will never be extended with new features. Modify the existing design and introduce multithreading.

Problem 5.21

Problem 5.22

Use Java RMI to implement a *distributed* Publisher-Subscriber design pattern.

Requirements: The publisher and subscribers are to be run on different machines. The naming server should be used for rendezvous only; after the first query to the naming server, the publisher should cache the contact information locally.

Handle sudden (unannounced) departures of subscribers by implementing a heartbeat protocol.

Problem 5.23

Suppose that you are designing an online grocery store. The only supported payment method will be using credit cards. The information exchanges between the parties are shown in Figure 5-37. After making the selection of items for purchase, the customer will be prompted to enter information about his/her credit card account. The grocery store (merchant) should obtain this information and relay it to the bank for the transaction authorization.

In order to provide secure communication, you should design a *public-key cryptosystem* as follows. All messages between the involved parties must be encrypted for confidentiality, so that only the appropriate parties can read the messages. Even the information about the purchased items, payment amount, and the outcome of credit-card authorization request should be kept confidential. Only the initial catalog information is not confidential.

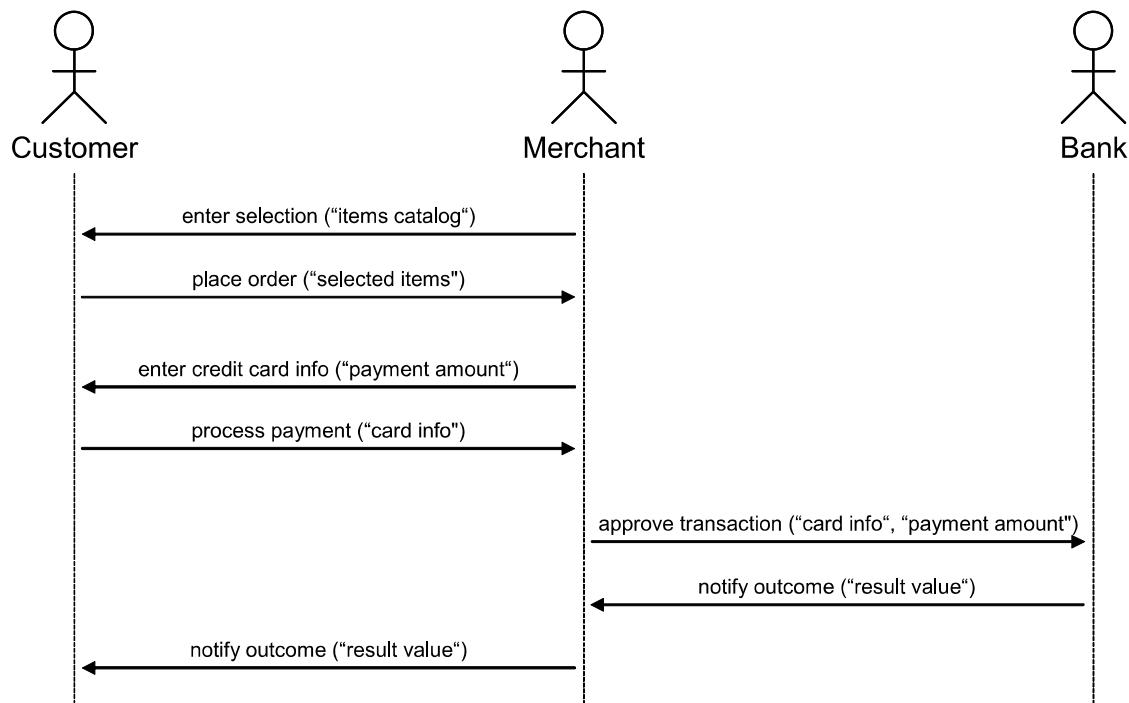


Figure 5-37: Information exchanges between the relevant parties. The quoted variables in the parentheses represent the parameters that are passed on when the operation is invoked.

The credit card information must be encrypted by the customer so that only the bank can read it—the merchant should relay it without being able to view the credit card information. For the sake of simplicity, assume that all credit cards are issued by a single bank.

The message from the bank containing binary decision (“approved” or “rejected”) will be sent to the merchant, who will forward it securely to the customer. *Both* the merchant and customer should be able to read it.

Answer the following questions about the cryptosystem that is to be developed:

- What is the (minimum) total number of public-private key pairs (K_i^+, K_i^-) that must be issued? In other words, which actors need to possess a key pair, or perhaps some actors need more than one pair?
- For each key pair i , specify which actor should issue this pair, to whom the public key K_i^+ should be distributed, and at what time (prior to each shopping session or once for multiple sessions). Provide an explanation for your answer!
- For each key pair i , show which actor holds the public key K_i^+ and which actor holds the private key K_i^- .
- For every message in Figure 5-37, show exactly which key K_i^+/K_i^- should be used in the encryption of the message and which key should be used in its decryption.

Problem 5.24

In the Model-View-Controller design pattern, discuss the merits of having Model subscribe to the Controller using the Publish-Subscribe design pattern? Argue whether Controller should subscribe to the View?