

What OLE Is *Really* About

Kraig Brockschmidt
OLE Team, Microsoft Corporation

July 1996 (Copyright © 1996 Microsoft Corporation)

Kraig Brockschmidt is a member of the OLE design team at Microsoft, involved in many aspects of the continuing development and usage of this technology. Prior to holding this position, he was a software engineer in Microsoft's Developer Relations Group for three years, during which time he focused his efforts on OLE, versions 1 and 2, and produced the books Inside OLE 2 and Inside OLE, 2nd Edition for Microsoft Press®. He has worked at Microsoft since 1988.

Abstract

Microsoft has made and is continuing to make heavy investments in OLE-related technologies. OLE itself has been in development for more than seven years, and almost every new technology coming out of Microsoft somehow incorporates elements of OLE. Why does OLE deserve such an investment? And why is OLE significant for the independent software vendor (ISV) and the computer industry as a whole?

The answer, as this paper explores, is that Microsoft created OLE to solve, in an object-oriented manner called *component software*, many practical problems encountered during Microsoft's lengthy experience in operating systems and applications. OLE provides the necessary specifications and the key services that enable component software, which is ultimately a significant gain for the entire computing industry.

A Brief History of OLE

Graphical user interfaces popularized the "clipboard" metaphor—with "copy," "cut," and "paste" operations—that greatly simplified the creation of a "compound document" that included text, graphics, and other types of content. Prior to this invention, you had to print the text and graphics separately, then cut and paste them together with scissors and glue.

Although the clipboard works well for the initial creation of a compound document, Microsoft began to realize its limitations in the late 1980s. Changing the text might require repositioning the graphics. Editing the graphics required many difficult manual steps to get that data back into the original format, if that was possible at all. Microsoft's Applications division then created a complex dynamic data exchange (DDE) protocol to simplify these steps.

Out of the DDE protocol grew OLE version 1.0 (1991), which was made available to all developers as a standard. OLE 1.0 greatly enhanced the creation and management of compound documents: One placed "embedded objects" or "linked objects" in a document that retained the native data used to create them (or a "link" to that data) as well as information about the format. (The acronym "OLE" is an abbreviation of "Object Linking and Embedding.") All the complexity of editing content was reduced to a

double-click of the mouse and presto! the object data was automatically brought back into the original editor.

But the OLE 1.0 designers realized that these "compound document objects" were actually a specific case of *software components*—small elements of software that can be "plugged in" to an application, thereby extending the functionality of that application without requiring changes. In the compound document paradigm, the document editor is a generic "container" that can hold any kind of content object. One can then insert charts, sound, video, pictures, and many other kinds of components into that container without having to update the container.

In the more general sense, component software has much broader applicability than to just compound documents. It is a multi-purpose "plug-in" model that is more powerful and flexible than other means, such as dynamic-link libraries (DLLs), Visual Basic® controls (VBXs), and the like. This was the guiding principle, then, behind the design of OLE version 2.0, released in 1993. Not only did OLE 2.0 improve on the compound document facilities of OLE 1.0, but it built a vast infrastructure to support component software on many levels of complexity.

The core of this infrastructure is a simple, elegant, yet very powerful and extensible architecture called the Component Object Model, or COM. It is within COM that we find the solutions to some of the most perplexing software problems, including those of extensible service architectures, "objects" outside application boundaries, and versioning. Furthermore, these solutions are concerned with *binary* components in a *running system* rather than source-code components in an application. This article explores these problems and the solutions that COM, and therefore the entirety of OLE, provides in the realm of running binary systems.

What makes COM and OLE unique is that the architecture is one of *reusable designs*. There is much talk in the software industry about reusable *code*. Recently, with the new emphasis on *design patterns*, there is increasing interest in being able to reuse a design that still allows flexibility with implementation choices. As we will see, one of the fundamental concepts in COM, that of an "interface," reflects the idea of design patterns.

COM and OLE therefore introduce a programming model based on reusable designs, as well as an implementation that provides fundamental services to make both design and code reuse possible. As a result, Microsoft has been introducing more and more "OLE Technologies" that build on the original architecture of OLE 2.0, including enhancements to COM (Network OLE) and OLE-based technologies built into the operating system (shell extensions). In addition, this architecture has enabled groups outside Microsoft, with or without Microsoft's involvement, to create important technologies in the real-time market data, point-of-sale, health care, and insurance industries.

Of course, people have started to ask about OLE 3.0 and when it would be released, and whether or not the shift from OLE 2.0 would be as major as the shift from OLE 1.0 (where the programming model was completely changed). In fact, there is no OLE 3.0, nor are there plans for such a release. This is because of the *reusable design* idea: OLE's architecture *accommodates* new technologies—regardless of Microsoft's involvement—without requiring modification to the base designs. Thus, the name "OLE" has ceased to be an acronym for "Object Linking and Embedding" and is simply OLE, pronounced "oh-lay" (as opposed to "oh-el-ee").

OLE, then, has moved from being a specific technology for a specific purpose to being a reusable architecture for component software that accommodates new designs and new technology.

So Just What Is OLE?

In the movie *The Gods Must Be Crazy*, an African bushman comes upon an empty Coke bottle that a careless Westerner tossed from an airplane. The bushman believes the bottle is a gift from the gods, especially as he finds more and more uses for it: carrying water, digging holes, pounding stakes, making music, and so on. Are there any limits?

OLE is very much like the Coke bottle (except that at the end of the movie the bushman decides that the bottle has caused too much trouble, so he tosses it off the edge of a cliff. This paper is not suggesting the same for OLE!). Just as the bushman would have a hard time pinning down exactly what the Coke bottle *is*, we (and especially Microsoft's marketing groups) have a hard time pinning down a solid definition of OLE. Throughout its history, OLE has been promoted as many things, from Visual Editing to OLE Automation to OLE Controls to Network OLE, and so on.

OLE might best be understood as a growth curve, like that shown in Figure 1. In other words, as time moves forward, OLE (based on COM) expands to accommodate new technology, never becoming obsolete as an architecture. OLE can be described as an extensible systems object technology whose architecture accommodates new and existing designs.

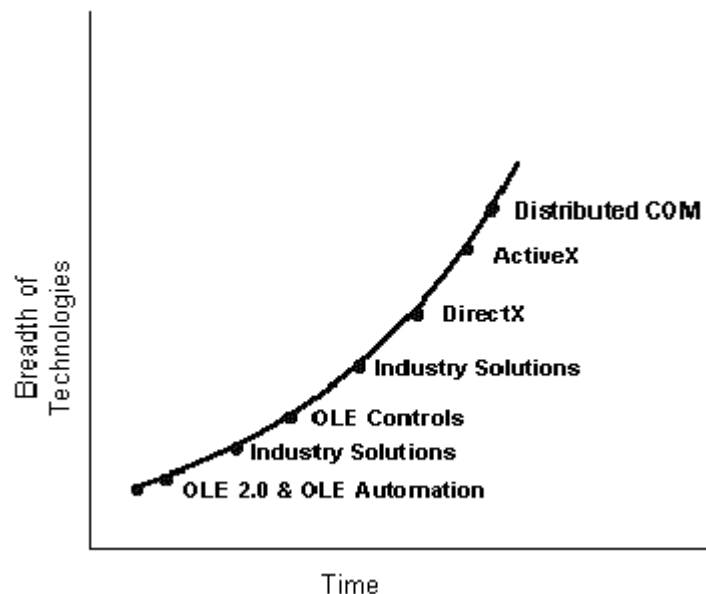


Figure 1. The OLE growth curve

A "systems object technology" means that one can work with object-oriented principles in the context of a running operating system—encapsulated, polymorphic, and reusable components exist and interoperate as binary entities, as opposed to source code definitions. New components, developed by anyone at any time, can be added into the running system, thereby immediately extending the services offered to applications, even if those applications are already running. This is what is meant by an

extensible service architecture, but it poses a number of different problems that are not relevant to source-code programming.

A newly installed operating system offers a basic set of services which developers employ in the creation of applications. COM and OLE provide the ability for anyone to extend the system with new services without requiring a change to the operating system. That is, COM and OLE make it possible for anyone to create new services with which developers create more and more innovative applications. Furthermore, this is accomplished without requiring any kind of central control or coordination between vendors. Microsoft, rather than stifling creative competition has built the infrastructure that allows everyone to create components in isolation but still be able to integrate those components in rich ways.

It is this potential for integration that can lead to significant improvements in how we develop software and, most importantly, the end-user's experience with the computer as a problem-solving tool. The last section of this paper examines that potential in more detail. First, however, let's see how Microsoft's experiences in operating systems and applications led to the designs of COM and OLE.

Problems and Solutions: The COM Architecture

How did Microsoft end up at OLE as an architecture on which it is betting its future? One must understand that the design of COM and OLE were not just dreamed up by a Microsoft architect who drank too much Jolt one night and stayed awake playing Reversi on "expert" level. The fundamental designs in COM are a result of many years of Microsoft's experience in the business of operating systems and applications.

In particular, Microsoft has repeatedly faced the problem of offering new software services for use in applications. The primary concern of an operating systems vendor is *how best to provide these services to applications*, because applications are built on system services and the system will not succeed without applications. Microsoft has had to deal with such issues throughout the lifetimes of MS-DOS® and Windows®.

Some kinds of services, such as device drivers and subsystems (as on Windows NT®), are easy to manage because such services are almost always shipped with the operating system. However, in a component software environment, nearly all new components (that implement a service) are shipped separately—not only separate from the system but separate from each other. Component software requires an architecture through which any developer or vendor can deliver a component at any time and have that component become immediately useful to applications on any given system.

To be used successfully, component software requires that applications always check on what components exist when they need them, instead of assuming there is only a limited set. When a new component is added to the system, it should become instantly available to all applications, even those that are already running. For example, consider a word processor that has a "Check Spelling" menu command whose implementation relies on the existence of a suitable spell-checker component. If a newer and better component is added to the system, that application can immediately take advantage of it the next time the user clicks that menu item.

A system that supports component software must therefore support a generic "service abstraction"—that is, an architecture that defines how all types of components appear and how they are manipulated.

In addition, the architecture must be extensible, so that a new component *category* (as opposed to an *implementation* of an existing type) can be introduced without having to revise the architecture. This is the problem of creating an *extensible service architecture*. For instance, it might be easy to define an architecture that accommodates components that provide content for compound documents, but can that same architecture accommodate later specifications for custom controls? In other words, the architecture must expect that new component types, or categories, will be defined later on.

The other big problem that such an architecture must solve is that of "versioning." It turns out that the first definition of a component type is easy to manage, as is the first implementation of any particular component. The difficulty comes in managing revisions to the designs and the implementations over time. COM and OLE are the results of Microsoft's experience with such problems.

From DLLs to COM

Let's say you have some kind of software service you would like to provide, maybe a library of useful functions. The usual way this is accomplished is to implement a dynamic-link library (DLL) that exports all the functions in the DLL. We call this code module the "server." The server then exposes the "flat API", so named because all the functionality of the service you are providing is described in a flat list of function calls. A "client" application is one that calls these functions to perform various operations.

This flat API setup has been used for quite some time and is considered the *de facto* method of exposing a software service—and it works reasonably well for code modules that are shipped with an operating system. Microsoft has also used the technique for creating VBXes as well as Microsoft® Foundation Classes (MFC) in DLLs. But as we will see in this section, there are many hidden problems that arise when new components are installed separately on a machine over a long period of time. We will see that the basic DLL model is lousy for component software, and that COM was designed to overcome the problems noted above.

As a starting point, let's use the spell-checker component mentioned in the last section. Assume that someone has written a specification that describes what a "spell checker" component does—that is, what it means to be a component that falls into this category. This specification describes the *service* of type "spell-checker" and might be as simple as follows:

A spell-checker component is a DLL that exports one function to check whether or not a particular word exists in the current dictionary. The function prototype is as follows:

```
BOOL WINAPI LookUpWord(wchar_t *pszWord)
```

This function must be exported as ordinal number 10.

Now, anyone can sit down and create an implementation of this service category, producing a *server* that supplies a spell-checker component. Typically there is only one component in any server DLL, and when that DLL is loaded into a process we just call it an "object." This is saying that the relationship between a component/server instance and the service category is just like that between an object instance and a class definition in a language like C++.

Let's say that hypothetical vendor Basic Software creates a server called BASICSP.L.DLL and sells it to the hypothetical ACME Software, Inc. ACME incorporates this DLL into their text application, AcmeNote,

implementing its Tools/Spelling... menu command by parsing each word in the text and passing words one by one into **LookUpWord**. The code for this command might appear as follows:

```
void OnToolsSpelling(void)
{
    wchar_t *pszWord;
    pszWord=GetFirstWord(); //Retrieve first word in text.
    while (NULL!=pszWord)
    {
        if (FALSE==LookUpWord(pszWord)) //Check if word is in dictionary.
            [Alert user]
        pszWord=GetNextWord(pszWord); //Go to next word in text.
    }
    return;
}
```

When an application uses a function provided by some DLL, that application must have an absolutely unique name to identify exactly which function it wants to call. In the basic DLL model on Windows, all absolute function identities include the code module plus the exported function ordinal (or text name). In our example, the **LookUpWord** function (ordinal 10) in BASICSPL.DLL is known as **BASICSPL.DLL:10**.

When the ACME developers compile and link AcmeNote, the call to **LookUpWord** is stored in the .EXE as an "import record" for **BASICSPL.DLL:10**. (An "import library" is what provides the mapping from names like **LookUpWord** as known to the compiler and the import record entry of SPELLCHK.DLL:10. Do a **dumpbin/imports** on some .EXE and you'll see a list of module:ordinal records for all the dynamic links.) An import record is a "dynamic link" to the function in the DLL because the actual address of the function is not known at compile time. Each **module:ordinal** record is essentially an alias for an address that is not known until run time. That is, when the kernel loads the application, it loads the necessary DLLs into memory (always prior to calling **WinMain**), thereby giving each exported function a unique address. The kernel then replaces each import record with a call to that known address.

An application can perform these same steps manually to control when the DLL gets loaded. One reason to do this is to bypass the kernel's default error handling when the DLL or exported function does not exist. If you run AcmeNote and BASICSPL.DLL can't be found, the system will complain with some wonderful message like, "The dynamic link library BASICSPL.DLL cannot be found in the path <list the PATH environment variable>," followed by "The application failed to initialize properly (0xc0000135). Click OK to terminate the application." Uh-huh. Yeah. By having the user perform these steps manually, AcmeNote could inadvertently disable its Tools/Spelling... command if the DLL isn't found.

If the kernel can find the DLL but cannot find the exported function, it says other nasty things. The Windows 3.1 kernel puts up an ugly white system-modal dialog that says, "Call to undefined dynlink." Come again? The Windows NT kernel says, "The application failed to initialize properly (0xc0000139). Click OK to terminate the application." To an end user, *things are definitely NOT OK* and they will probably be calling customer support very soon.

In the case of AcmeNote, we can make the application boot faster if we don't load the DLL at all until we actually execute the Tools/Spelling command. To do this, we explicitly load the DLL with the **LoadLibrary** Win32® API, then retrieve the address of **LookUpWord** (ordinal 10) using the **GetProcAddress** Win32 API. We then call the function through this pointer:

```
//Function pointer type for LookUpWord
```

```

typedef BOOL (WINAPI *PFNLOOKUP)(wchar_t *);
void OnToolsSpelling(void)
{
    HINSTANCE hMod;
    PFNLOOKUP pfnLookup;
    wchar_t *pszWord;
    hMod=LoadLibrary("BASICSPL.DLL");
    if (NULL==hMod)
        [Spell-checker not found, show error]
    //Get address of ordinal 10
    pfnLookup=(PFNLOOKUP)GetProcAddress(hMod, MAKEINTRESOURCE(10));
    if (NULL!=pfnLookup)
    {
        pszWord=GetFirstWord(); //Retrieve first word in text.
        while (NULL!=pszWord)
        {
            if (FALSE==(*pfnLookup)(pszWord)) //Check if word is in dictionary.
                [Alert user]
            pszWord=GetNextWord(pszWord); //Go to next word in text.
        }
    }
    else
        [Export not found, invalid spell checker!]
    FreeLibrary(hMod);
    return;
}

```

For convenience in terminology, we can call the instance identified with *hMod* a "component" or an "object" (the two are not quite equivalent, though often used interchangeably). Describing it as an "object" is interesting because what we see here actually meets the core object-oriented principles. The so-called "object" is *encapsulated* behind the "interface" of **LookUpWord**, because the function is specified outside any implementation. The "object" is also *polymorphic* in that this same client code in AcmeNote would work perfectly fine if we substituted WEBSTERS.DLL in place of BASICSPL.DLL—two implementations of the same specification would be polymorphic. And finally, an instance of BASICSPL.DLL is *reusable*: Someone could implement another spell-checker, such as MEDDICT.DLL, which only recognizes specific terms in the medical field. That DLL could itself load and reuse BASICSPL.DLL in such a way that MEDDICT.DLL checked only medical terms, passing all other unrecognized words to BASICSPL.DLL.

Of course, this begs the question of how on earth AcmeNote would know to load MEDDICT.DLL or WEBSTERS.DLL instead of BASICSPL.DLL. This brings us to the first in a series of problems that arise when DLLs are used as a component software model. In the next few sections, we'll look at a number of problems in detail and the solutions to those problems. These solutions form much of the COM design.

Hard-Coded DLL Names

It should be obvious from the piece of code above that the name "BASICSPL.DLL" is hard-coded into AcmeNote. This happens either way we write the code: If we call **LoadLibrary**, the hard-coded name is explicit. If we use an import library, the name is hard-coded in the import record, which we don't even see. Such hardcoded names lead to the following problems:

Problem 1

The DLL must exist either in the application's directory, in a system directory, or in some directory included in the PATH environment variable.

Problem 1a

If the DLL is stored in the application directory, it cannot be shared between multiple applications that want to use the same service. Therefore multiple copies of the DLL end up on the same machine, wasting disk space. (This problem is amplified when there are dozens or hundreds of different DLLs duplicated on a machine. The result might be 10 megabytes of wasted disk space.)

Problem 1b

If the DLL is stored in a system directory, *there can be only one provider of the service known as BASICSPL.DLL*. Alternatively, each provider must have a different name, but in that case AcmeNote will not work if WESBTERS.DLL exists but not BASICSPL.DLL. You also run the serious risk that someone else will install a BASICSPL.DLL that is written to a different specification, which will lead to the "Call to undefined dynlink" sort of error messages described earlier.

Problem 1c

If the DLL is stored in some other path, the application must add that directory to the PATH environment variable. On Win32 platforms this is a non-issue, but for Windows 3.1 the size of the PATH is severely restricted.

Much to the delight of hard-drive manufacturers, software vendors usually solve this set of problems by installing all DLLs in the same directory as the application. Thus only one application will ever use that given copy of the DLL. What a waste! Yet the only solutions to these problems must occur outside the DLLs themselves.

Part of the solution is to remove all the path dependencies by defining an abstract identifier for each DLL. That is, we might define the number 44,980 as being the abstract ID for BASICSPL.DLL. We then require some way to dynamically map this number to the exact installation point of the DLL. This is the purpose of the system registry (or registration database). In the registry we might create an entry like this:

```
HKEY_CLASSES_ROOT
  ServerIDs
    44980 = c:\libs\44980\basicspl.dll
```

That is, the "ServerIDs" section would list the exact location of many different DLLs. So instead of using "BASICSPL.DLL" directly in the code above, we would write it like this:

```
void OnToolsSpelling(void)
{
    wchar_t    szPath[MAX_PATH];
    [Other locals]
    if (!MapIDToPath(44980, szPath, MAX_PATH)) //Looks up ID in registry and
                                                //returns path.
        [No mapping found, show error]
    hMod=LoadLibrary(szPath);
    [Other code the same]
}
```


This at least allows you to install one BASICSP.L.DLL in a specific directory so that multiple applications can share it. Each application would only hard-code the ID of 44,980 and not depend on the exact location of the DLL.

Of course, this doesn't help the same applications use an alternate *implementation* of the same *service category*, which is a critical requirement for component software. That is, we'd like to write AcmeNote so that it uses whatever implementation of the "Spell Checker" service is around, not relying solely on BASICSP.L.DLL. Stated in object-oriented terms, we'd like to have all implementations of the same service category (specification) be *polymorphic* with one another so that AcmeNote doesn't care which DLL actually provides the **LookUpWord** function.

(Mind you, a critical requirement for component software is not a requirement for things like device drivers. Because there's typically only one piece of a certain type of hardware in a machine at once, there only needs to be one device driver for that type in the system. For example, if you have only one video board, you need only one systemwide device driver.)

Certainly AcmeNote would be shipped with BASICSP.L.DLL as a default component, but if a newer and better implementation—a WEBSTERS.DLL (ID 56791) with more words and a faster algorithm—showed up with some other application, we'd like AcmeNote to automatically benefit.

What we need, then, is an identifier for the *category* along with registry entries that map the category ID to the available implementations. So let's say we assign the value 100,587 to this spell-checker category. We'd then create registry entries that mapped the category to the implementations, which would be mapped elsewhere to the locations of the modules:

```
HKEY_CLASSES_ROOT
Categories
    100587 = Spell Checker
    56791 = Webster's Spell Checker
    44980 = Basic Spell Checker
ServerIDs
    44980 = c:\libs\44980\basicspl.dll
    56791 = c:\libs\56791\websters.dll
```

Now we can write AcmeWord to use the most recently installed spell-checker (which could be an inferior one), or we can add a user interface such as that shown in Figure 2, which allows the user to select the preferred spell-checker from those listed in the registry. We'd then have code as follows:

```
void OnToolsSpelling(void)
{
    DWORD      dwIDServer;
    wchar_t    szPath[MAX_PATH];
    [Other locals]
    dwID=MapCategoryIDToServerID(100587); //Magic to get a server ID, may
                                         //involve UI
    if (!MapIDToPath(dwID, szPath, MAX_PATH)) //Looks up ID in registry and
                                         //returns path
        [No mapping found, show error]
    [Other code the same]
}
```

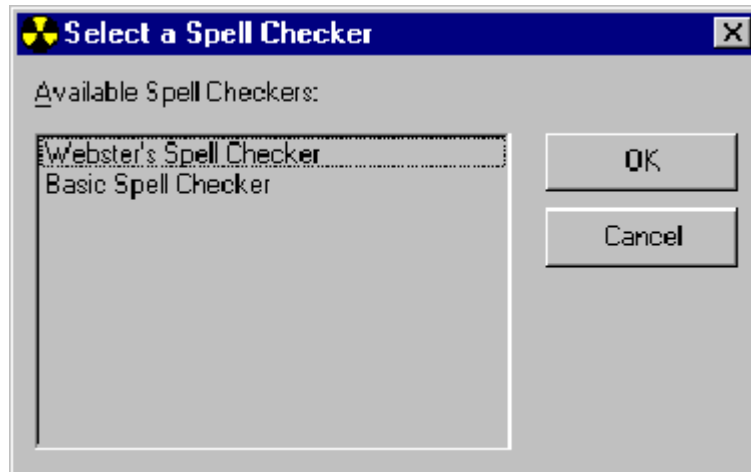


Figure 2. A hypothetical selection dialog based on registry entries

So with category identifiers, server identifiers, and a registry we can solve Problem 1. COM and OLE use the registry in much the same way to bypass exactly these problems. For example, "OLE Documents," the OLE specification for compound documents, describes what are categorized as "insertable" (that is, embeddable) content objects. The various implementations of such objects are described in the registry, and when users want to insert a piece of content into a document, they invoke the standard "Insert Object" dialog, which lists the available object types much like the dialog shown in Figure 2. OLE Controls, defined in another specification, operate in the same manner.

However, the use of abstract identifiers presents another problem itself:

Problem 2

Who defines the category IDs and the server IDs?

There are two possible solutions for this:

- One, some central organization takes responsibility for maintaining a master list of IDs assigned to categories and IDs assigned to DLLs. This means that vendor wishing to either define a category or ship an implementation of some category must get "permission" to do so by obtaining unique identifiers from the organization controlling the master list. Blech.
- Microsoft realized that a centralized approach like the one described above would stifle innovation and be a royal pain to manage. Therefore COM and OLE use *globally unique identifiers* (GUIDs, pronounced goo-ids or gwids)—128-bit values that are generated using an algorithm defined by the Open Systems Foundation to guarantee uniqueness across time and space. Microsoft provides a tool that implements the algorithm so that anyone, anywhere, anytime, can obtain identifiers as needed and still be assured uniqueness.

Of course, it is also wasteful for each application to define its own category IDs and its own registry structure, so Microsoft has defined some standards in these areas. COM and OLE provide APIs to manipulate GUIDs and their various uses in the registry. We'll see more details later in this paper. For

now, simply understand that *COM eliminates dependencies on DLL names, thus allowing multiple polymorphic providers of the same service to coexist.*

Management APIs

You may have noticed that the code in the previous section contained two undefined "magic" functions: **MapCategoryIDToServerID** and **MapServerIDToPath**. The first function looks up the ID of a server for a particular category that may display a UI similar to Figure 2. The second function retrieves the path of the server DLL associated with a particular ID. We might have simply combined all these steps, including the call to **LoadLibrary**, into one function named something like **LoadCategoryServer**.

Such functions are an example of a "management API," which frees an application from the tedium of walking through the registry and so forth. But the API example shown here is actually a pretty advanced one because we have a registry and have defined some kind of standards for the structure of the registry. Therefore the API can be somewhat generic and accommodate many different service types.

However, this was not always the case. For a long time there was no central registry and there were no standards for using any other means to store such information. The result was that whenever someone defined a new service category, they also defined a category-specific management API. The intention was, of course, to make it easier to write an application that used services in that category, and indeed such APIs did for their respective categories. In our example, we might have a specific **LoadSpellChecker** API function that essentially hard-coded the category ID. On the surface, this does simplify the programming model in our code.

Over time, many new categories were defined, each with their own management API. In each case the API initially made working with an individual service type easier, but overall they complicated application programming. Because most applications employ many different services to achieve their desired ends, having a specific API for each service meant a steep learning curve for all the programmers involved. The Win32 API itself, evolved over a decade, is a prime example: There are approximately 70 different "creation" functions, each dealing with a different kind of "object," such as **CreateAcceleratorTable**, **CreateDialog**, **CreateEvent**, **CreateMailslot**, **CreateMutex**, **CreateRectRgn**, and **CreateWindowEx**. Each type of "object" also has a different way to identify it (handle, pointer, structure, and so on) and a different set of APIs to manipulate it. Hence there are hundreds of API functions, each working in different ways, which is what makes programming for Windows so difficult.

We may state the problem as follows:

Problem 3

New service categories introduce new management and manipulation APIs, meaning to simplify the programming model for that category, but tend to complicate the overall programming of an application.

The solution is an obvious one: Create generic abstractions for the common elements found in all these APIs. This is what the **MapCategoryIDToServerID** and **MapServerIDToPath** functions are doing in our example. The process of finding the existing implementations of a category is a common one in nearly all category definitions. By defining a generic API based on unique identifiers, we greatly reduce the overall number of API functions and thereby truly simplify the overall programming model.

This is exactly what COM does—*provide a single generic API for the management of all categories and servers*, based on the standards defined for how these things are listed in the registry. This generic API is extensible in that it can accommodate any new categories and services until the end of time.

Furthermore, this API is made up of only a handful of functions, which can be learned in a matter of hours or days, the most important of which is named **CoCreateInstance**. This function creates an instance of some class (given a CLSID) and returns an interface pointer for it.

A key to this generic API is that object instances are always referenced through an "interface pointer," the exact structure of which we'll see later. For now, suffice it to say that all interface pointers are also polymorphic in a way that allows a generic API to manipulate any kind of interface pointer to any kind of object. This stands in stark contrast to the myriad non-polymorphic handles and references that one finds in the Win32 API.

Shared Server Instances and Lifetime Management

The last section mentioned that many different service types in the Win32 API have different means to identify an "instance" or "object" of the service, component, or object of concern. In the code we've looked at so far, the spell-checker "object" is identified with a module handle through which you can get at function pointers. This is all well and good for the application that loaded the DLL, because the module handle makes sense in that application's process space. However, the same handle is useless to another process. That is, our AcmeNote application could not pass the module handle to another application (via remote procedure call—RPC—or some other inter-process communication mechanism) and expect that other process to make any sense of it. Hence we have another problem (or limitation, one might say):

Problem 4

In general, identifiers for instances or for a server or component cannot be shared with other processes, let alone with processes running on other machines.

On Win32 platforms, process separation means that each application wanting to use a handle-based service must load that service into its own address space, wasting memory (even virtual memory), increasing boot time, and causing an overall performance degradation. The only workaround for this on Windows NT is to create a subsystem type of service where only one instance is running for the whole system. Much of Windows itself (like USER.EXE) works this way, which is why certain handles, like HWNDS, can be passed between processes. But these types of system services are expensive and are not suitable to the general problem of component software.

Even if you could solve the problem for a single machine, it still isn't possible to share instances *across machines*—that is, have multiple applications on multiple machines accessing a service that's running on yet another machine, like a central database. In order to make distributed systems work you must go to full-blown RPC, or named pipes, or another such mechanism, all of which of course have structures and programming models different from anything you might have on a single machine or within a single process. Thus we have another problem:

Problem 5

The mechanisms for working with services differ greatly between the in-process (DLL), local (.EXEs, other processes on the same machine), and remote (other machine) cases. A client application ends up having to understand three or more different programming models to work with all three service locations.

With DLLs in particular, there are other problems that would arise if you *could* share module handles between processes. Let's say AcmeNote has loaded BASICSPL.DLL and then passes the module handle to another application. What happens when AcmeNote terminates and the other application still has the module handle? Because the DLL is loaded into AcmeNote's process, that DLL is unloaded as well. Thus the module handle becomes invalid and we can expect the other application to crash the next time it attempts to use that handle.

There's a similar problem on 16-bit platforms like Windows 3.1, where all applications run in the same address space and one instance of a DLL could be used by multiple applications. The problem there is that abnormal termination of one application might cause the DLL to remain in memory even when all other applications using it are closed normally. Programmers for Windows 3.1 usually become proficient with a little tool called WPS.EXE, which allows you to clean orphaned DLLs out of memory.

In both cases we can state yet another problem:

Problem 6

When instances of a service can be shared between processes, there is usually no robust means of handling abnormal termination of one of the processes, especially the one that first loaded the service. Either the server remains orphaned in memory or other processes using that service may crash. DLLs, for example, have no way of knowing when other processes are referencing them.

COM solves all of these problems by virtue of the "interface pointer" structure (introduced in the previous section) as well as by providing the ability for you to implement servers as *either* .DLLs or .EXEs. Through a process called "marshalling," one can share an interface pointer with other processes on the same machine or on other machines. Of course, you don't pass the exact same pointer value, but the result is that multiple processes can jointly own an interface pointer through which they can call functions in the object attached to that pointer. The structure that makes this possible is borrowed from RPC and also guarantees robustness when the process that loaded the object is terminated. When other processes try to use that object, they receive "disconnected" error codes in return, instead of just crashing. Also, if an object is implemented in its own .EXE and has its own process, it will remain active as long as any client holds a reference to it. (COM automatically releases references held by clients that terminate abnormally, so the server process doesn't get orphaned itself.)

What's more, the *only* way one ever deals with an object instance in COM and OLE is through interface pointers. This applies regardless of the actual location of the object itself. Therefore the programming model is identical for the in-process, local, and remote cases. The technology that makes this work is called "Local/Remote Transparency," as we'll see later.

Multiple Services in a Single DLL

Early in this section we described a *server* and a *component* as pretty much the same thing, and we used the terms *component* and *object* to describe a loaded instance of a server. We could do this because of the base assumption that each server DLL only implements a single service.

However, anyone who has been doing Windows programming for very long knows that a DLL implies a fixed amount of overhead, both in memory footprint and in the time it takes to load the DLL. When performance is a big issue, the first thing you'll want to do as a component vendor is to combine services into one DLL to improve the boot time and working set of applications that use your DLL.

The problem here is that when services are defined as a set of DLL exports, allowing one DLL to offer multiple services requires some kind of central coordination as to which ordinal numbers and/or function names are assigned to which service categories. That is, if two categories both want to use ordinal 10 (or the same text name) for completely different operations, you could not implement both services in one DLL! So once again we're right back to the issue of having some central body (that is, bottleneck) having to grant "permission" to those creating new designs. Yuck.

Previously we saw how COM and OLE use GUIDs to identify both service categories and implementations of a category, which allows anyone to independently define or implement a service without having to talk to anyone else. In order to allow multiple services per DLL, we have to solve the following problem:

Problem 7

A server cannot support multiple services without risking a conflict between the function names or ordinals used to reference those functions. The choice seems to be between central control over ordinal assignments and the inability to combine services together, either stifling innovation or causing performance problems.

The real solution is to get away from referencing functions with names and ordinals altogether. Earlier we saw that the absolute identity of any given function in a DLL-based system is of the form **module:ordinal**, for example, **BASICSPL:10**. In COM and OLE, components and objects expose their functions not through singular exports, but in groups that can be referenced through a single interface pointer. An "interface" is actually defined as a group of related "member functions" and that group is itself assigned a GUID, called an Interface Identifier, or IID, which anyone can generate independently like all other GUIDs.

For example, if we were to redesign the spell-checker specification to use COM, we'd define an interface called **ISpellChecker** as part of the specification (the conventional "I" prefix to the name stands for "interface"), as follows:

A spell checker component is any server supplying an object that implements the ISpellChecker interface which is defined as follows using the Interface Definition Language (IDL):

```
/*
 * IUnknown is the "base" interface for all other interfaces. The value inside
 * the uuid() attribute is the IID, which is given the symbol IID_ISpellChecker.
 * ISpellChecker is the interface type as recognized by compilers.
 */
[uuid(388a05f0-626d-11cf-a231-00aa003d7352), object]
interface ISpellChecker : IUnknown
{
    HRESULT LookUpWord(OLESTR *pszWord);
}
```

Notice how we've eliminated any reference to ordinals as well as any need to stipulate that the spell-checker must be implemented in a DLL. Implementing a server according to this specification is hardly any more work than implementing the DLL according to the old specification.

As we'll see in detail later, an interface pointer (a run-time entity) really points to a table of function pointers. In the case of **ISpellChecker**, the table has four entries—three for the members of the base interface **IUnknown** and one for **LookUpWord**. The table contains the addresses of the implementations of these interface member functions. When a client wants to call a member function through an interface pointer, it actually calls whatever address is at the appropriate offset in the table.

Programming languages make this easy, given the way interfaces are defined in header files. AcmeNote can now use COM's generic **CoCreateInstance** function and the interface to meet this new specification (this is C++ code), as follows:

```
void OnToolsSpelling(void)
{
    ISpellChecker *pSC;
    CLSID          clsID;
    //This hypothetical function uses the registry or invokes the UI.
    GetCLSIDOfSpellChecker(&clsID);
    if (SUCCEEDED(CoCreateInstance(clsID, NULL, CLSCTX_SERVER
        , IID_ISpellChecker, (void **)&pSC)))
        [Spell-checker not found, show error]
    pszWord=GetFirstWord(); //Retrieve first word in text.
    while (NULL!=pszWord)
    {
        if (S_OK!=pSC->LookUpWord(pszWord)) //Check if word is in dictionary.
            [Alert user]
        pszWord=GetNextWord(pszWord); //Go to next word in text
    }
    pSC->Release();
    return;
}
```

Given the interface definition, a compiler will know to generate the right machine code that calls the right offset for the **LookUpWord** member of **ISpellChecker**. To do this reliably requires a standard for the binary structure of an interface—one of the fundamental parts of the COM specification.

Now, putting all this together, we see that any given function—in any given object implementation—is identified with three elements: the object's CLSID, the interface type (IID) through which the client will call member functions, and the offset of the particular member function in the interface. Thus, the absolute function identity in COM and OLE is **CLSID:IID:table_offset**; for example, **CLSID_BasicSpellChecker:IID_ISpellChecker:4**.

Because anyone can generate CLSIDs and IIDs at will, anyone can design a service where the functions describing that service are ultimately identified with two absolutely unique values.

COM allows any given server module to implement as many different classes as it wants, meaning that any server can support as many CLSIDs as desired. When a client wants to access an instance of a class, COM passes the CLSID to the server and the server can decide what to instantiate and which interface pointer to return. That interface pointer points to a table that holds the addresses of the code that is unique to the particular object class. In short, COM removes all barriers to multi-service implementations, regardless of who designed the service.

A small note about the **pSC->Release()** line at the bottom of the code above: **Release** is a member of the **IUnknown** interface that all interfaces share, thus all objects have this member. *Release* is how a client tells the object that it is no longer needed. The object maintains a reference count for all clients using it, so it knows when to free itself from memory. In addition, a server will maintain a count of objects it happens to be servicing, so that when no objects remain it can unload itself (or terminate, if it's an .EXE server).

In this sense, the **Release** function is the universal "delete" abstraction across all of COM—there are only a handful of other special-case "free" or "delete" functions in the whole OLE API. Because all objects are manipulated through interface pointers, freeing the object always means calling its **Release** member. This single yet powerful abstraction significantly simplifies the overall programming model, just like the generic creation function **CoCreateInstance**. In fact, one of the most common patterns in COM/OLE programming is: *(1) Use **CoCreateInstance** to obtain the interface pointer, (2) call functions through that pointer, (3) call **Release** through that pointer when the object is no longer needed.* You see it everywhere, sometimes involving a different creation function, but you still see the pattern.

The Big One: Versioning

We've worked our way now through many of the problems inherent in DLL-based service models and we've introduced many of the solutions found in COM and OLE. But everything we've dealt with so far applies only with the *first* version of a particular service. Not only the first version of the category definition, but also the first version of the server implementation, and the first version of the hypothetical AcmeNote application that uses such a server.

The big question is what happens when we want to (1) change the service specification, (2) update the server, and (3) update the client. You *will* eventually want to do this because innovation is at the heart of the software industry! The primary issue becomes one of compatibility between different versions of the services and the applications that use those services. A new service must be prepared to work with an application that expects an old service. A new application must be prepared to work with an old service.

This is called the "versioning problem," which happens to be the most important problem that COM was designed to solve because it is one that has historically made many lives miserable. This applies even to those services that only require one provider in the system! The problem is this:

Problem 8a

Versioning: Independently changing the specification of a service, the implementation of a server, or the expectations of a client typically results in significant interoperability problems that ultimately prevent innovation and improvements in software.

Or, stated another way:

Problem 8b

When an application talks directly to a service, how does that application dynamically and robustly discover the capabilities supported in that service? How does an application differentiate between different versions of the same service category and between different versions of a server implementation?

It is quite a claim to say that COM fundamentally solves this problem across the board, yet this is true. We must first, however, understand the problem, which we'll do by going back to the original specification for the spell-checker DLL. The specification stipulated that a component, which we'll call "SpellCheck1.0," is a DLL that exports a function called **LookUpWord** as ordinal 10.

Now we've happily created BASICSPL.DLL, which we'll call "BasicSpell1.0" for convenience, and this DLL is used by AcmeNote version 1.0. At this stage, interoperability is a non-issue as long as everyone sticks to the API specifications—this is *always true* with version 1.0 of *any* API. While sticking to the specifications, one is free to make performance improvements that don't affect the interface between client and service. So BASICSPL.DLL can ship a new version with more words and a faster lookup algorithm to compete with WEBSTERS.DLL. No problem.

Problems arise when we want to change the *specification*, usually because we want to innovate, add new features, or otherwise meet additional customer demands. With spell checkers, for example, customers will soon want the ability to add custom words to the dictionary. This will require new functionality in spell-checker components. Working in the DLL world, we would write the SpellCheck2.0 specifications like those shown in Table 1.

Ordinal	Prototype	Description
10	BOOL WINAPI LookUpWord(wchar_t *pszWord)	Checks to see if a particular word exists in the current dictionary.
11	BOOL WINAPI AddWord(wchar_t *pszWord)	Adds a custom word to the dictionary.
12	BOOL WINAPI RemoveWord(wchar_t *pszWord)	Removes a custom word from the dictionary.

Table 1. A Spell-Checker 2.0 Component (DLL) That Exports Three Functions

Note that a SpellCheck2.0 implementation is polymorphic with a SpellCheck1.0 implementation because 1.0 clients will only look for ordinal 10. It seems that we could overwrite the DLL for BasicSpell1.0 with a DLL for BasicSpell2.0, but as we'll soon see, there are many potential problems.

Earlier we saw two different versions of AcmeNote code that used the **LookUpWord** function, one using an import library to resolve the function name, the other retrieving the address of the function directly with **LoadLibrary** and **GetProcAddress**. Both solutions work equally well when there's only one function to worry about. But when we want to write AcmeNote2.0 and add the "Add/Remove Word" feature, the additional functions make it more complicated to program around **GetProcAddress** than to use import libraries. It's still tolerable for AcmeNote2.0 because all these functions are used at the same time and we only need one call to **LoadLibrary**.

But any worthwhile application will not use just one service DLL, nor only a handful of functions in the same place in the code. Most likely there will be dozens of functions strewn all around the application. Calling **LoadLibrary** and **GetProcAddress** for each call introduces many more error conditions, which rapidly increases the complexity of the application and thus the bug count. In the real world, vendors are concerned with shipping the application in order to make money, so the more robust means of dynamic linking quickly gives way to the use of import libraries. After all, having implicit links like this isn't so bad, is it?

Well, let's say we ship AcmeNote2.0 along with BasicSpell2.0, which happens to be in the same module as BasicSpell1.0 (that is, BASICSPL.DLL). When we install AcmeNote2.0, we overwrite the old BASICSPL.DLL with the new one. We know that any application that was using BasicSpell1.0 will still be able to use BasicSpell2.0, so compatibility is ensured.

However, there now exist two different versions of BASICSPL.DLL that, to most people, are indistinguishable from one another. Microsoft has frequently discovered that older versions of a DLL end up overwriting newer versions, no matter what rules you lay down. Some application that uses BasicSpell1.0 might come along and overwrite the newer BASICSPL.DLL with the older one! (The reason for this is that it is difficult to do the right thing—to check the existing DLL version against the one you're installing. This is tedious and it's normally very hard to find a programmer who wants to work on an application's installation program. When crunch mode comes and you want to ship an application, you're much more likely to spend time fixing application bugs than making your installation program perfect. Checking version resources on shared libraries is a likely thing to compromise.)

So what happens when the user now runs AcmeNote2.0? When that application tries to resolve the **AddWord** and **RemoveWord** calls, it will fail. If import libraries are being used, which is most common, the kernel will generate its enigmatic error messages like the ones we saw before—"Call to undefined dynlink," or "The application failed to initialize properly (0xc0000139)." Again, expect customer support calls to come streaming in.

No matter how hard one tries to prevent this situation, it invariably happens when multiple versions of the same service use the same DLL name. There are several ways to prevent the problem that are less than optimal: punt, bloat, stagnate, or suffer. Let's look at these in detail before we see how COM solves the problem for good by making the "right" thing to do an inherent part of the whole programming model.

Punting: Programming to the Least Common Denominator

The easiest way to avoid "Call to undefined dynlink" problems is simply to avoid using a new version of a DLL in the first place! This is often an appealing solution because it takes no work and carries no risk, but it means no innovation.

The result is called the "least common denominator" usage of the feature set of a service DLL, meaning that almost no one uses any functionality that is not guaranteed to be present in *all versions* of that service DLL. What happens then is that the vendor of the DLL is hesitant to add any more features to the DLL because few applications will bother using them—there is little return on investment. Who is going to take the risk first? The DLL vendor or the application vendor? Quite commonly, neither.

Usually the only way to avoid this sort of trouble is never to allow the service DLLs to be distributed by themselves, thereby avoiding the possibility that an old DLL would overwrite a new DLL. This basically means that the new service DLLs must be shipped with a new version of an operating system, and that application should be revised and recompiled to work on the new operating system. Anyone who has been around the Windows development scene for a few years knows how this works: As Microsoft prepares a new version of Windows, it heavily evangelizes independent software vendors (ISVs) to produce new versions of applications that take advantage of the new features in the system, thereby making the old system obsolete, thus encouraging users to upgrade the operating system, the applications, the hardware, and so on.

Microsoft must do this because otherwise the new system will be a complete flop. Microsoft also must invest tremendous resources in making sure that old applications still do work, because without that backwards compatibility the system is also doomed. Not only are evangelism and compatibility testing expensive, but the system ends up bigger and slower because of it. And of course, this doesn't even start to solve the problem for vendors who want to ship their own services!

Bloating: Shipping Multiple DLLs with Redundant Code

Another solution that avoids both the least-common-denominator problem and the versioning problem is always to ship new features in an altogether new DLL. This happens in one of two ways: Either you add new code to the existing DLL code base and recompile the whole mess, or you make a new DLL that contains only the new features. Each solution is basically saying, "DLLs never get versioned at all!"

In the first case, the DLLs themselves are given a filename that reflects the version number, such as BASSPL10.DLL, BASSPL20.DLL, etc. Therefore all the versions of the server DLL can coexist, and there is no chance of overwriting the old version because it has a different name.

This is, in fact, how Microsoft has traditionally solved the versioning problem, as with the Visual Basic run-time libraries in VBRUN100.DLL, VBRUN200.DLL, VBRUN300.DLL, and so on. The trouble, however, is that these things build up on the user's machine. On the machines I'm working on right now, I have the following versions of the 16-bit Visual Basic run-time library: VBRUN200.DLL (348K), VBRUN300.DLL (389K), and VB40016.DLL (16-bit VB 4.0 library, 913K). In all, 1.65MB of Visual Basic run-time libraries!

My question is this: How much code is repeated between these DLLs? I would bet that there is a fair portion of code that exists identically in all of these DLLs. Ideally I should need only the most recent DLL to work with all versions of the applications that use the DLL's services. On my machine, therefore, I should only need the 913K of VB40016.DLL and I should be able to toss 737K of old stuff. Now by itself this isn't bad, but multiply the problem over time with a few added versions and by dozens of DLLs. Sooner or later I might have 50MB of disk space being eaten by worthless DLLs that never get used. Yet I'm afraid to delete them in case one of my applications needs them!

Unless you're a hard-drive manufacturer, this sucks. Users would love to know when they can delete some of these old DLLs to make more disk space, and there are software products on the market that try to help users decide whether or not some old DLL is in use. (Occasionally I'll move some old DLLs to a "DeleteMe" floppy and work on my machine for a while—if nothing uses those DLLs I'll eventually delete them.)

So why does Microsoft put redundant code into each version of these DLLs? Why not just put in VBRUN200.DLL only what is different from VBRUN100.DLL, and do the same for subsequent versions?

Well, there are two reasons. One is that this still bloats the user's machine with a bunch of DLLs, albeit smaller, but still a lot of clutter. The other reason—the big one—is performance. As mentioned before, each DLL has a fixed overhead in load time and memory footprint. Increasing the number of DLLs can exponentially increase application load time. In an industry that rates applications partially on their boot time, this simply is not acceptable.

So in order to sell its applications, Microsoft and all other vendors usually decide to shift the costs to the user, who now must buy more disk space or more powerful machines. Microsoft doesn't like having to do this, but when it's a tradeoff between selling the application and having the user spend more on

hardware, or not selling the application at all, the former is chosen hands-down. Of course, the hardware manufacturers don't mind one bit.

Stagnating: The Service-Level Approach

Another possible solution is to involve a specific management API that allows the application to negotiate with a server to determine if that server has support for a particular "level" of service. We've already seen that having a specific management API is a bad idea; still, this solution has been used in the past.

The idea here is that instead of having just one "functional API" exposed from the server, there are different API "levels," each level consisting of a group of functions that together represent a set of additional features. New features are introduced in new levels. Thus the "Level 2 API" extends the base API, "Level 3" extends "Level 2," and so on. The provider of a Level n service is polymorphic with all other $\leq n$ providers.

Each application that wants to use a service through this API asks the management API to locate those services that implement a certain level of API, using registry entries. Thus the application will guarantee that it will be working with a server that implements the right functionality, or that it can disable certain features in order to work with a down-level server.

The application doesn't actually link to the functions exported from the server itself because it doesn't know the DLL name ahead of time; rather, it links to a "stub" server that provides entry points for all functions in all levels. When the application asks the management API to load a server, that server is loaded and dynamically linked into the stub server (using **GetProcAddress**). Client calls into the stub server are then forwarded to the real server.

A key feature in this kind of architecture is the grouping together of functions as "levels." Each level is a "contract," which means that when a server implements one function in a level it must implement *all* functions in that level. This enables the client application, having successfully located and loaded a server supporting a given level, to trust that all calls to all supported levels will work as expected. The client only has to check for Level 3 support once before it can invoke any function in Level 3.

This "contract" idea also means that the application doesn't have to ask the provider about each function call before invoking that function. Without the manager API, determining "levels" of service in some random DLL would mean calling **GetProcAddress** for each function before trying to call that function. The next section examines the implications of this method.

So why isn't this "level" architecture the right solution? Because a number of the problems discussed earlier are still present in this model:

- First of all, who defines each level? Because a level is defined as a group of exported functions, we still have the problems of centralized control over the specification.
- Second, we have a specific management API with all its inherent problems, including the fact that it is a tedious process to define such an API, which must be done at each level. Solutions like this end up being designed in a committee, which is a notoriously *slow* way to innovate; the result is usually stagnation.

- The other big problem is that this *really doesn't solve the versioning problem*, because the architecture doesn't make any provision for changing the definition of a *level* itself. The only way to add new functionality is to define a new level. But what if you have a new feature that belongs in Level 2 but you have to add it as Level 6? That means that to implement this new feature you must implement Levels 3, 4, and 5, even when you hadn't before. Unless you design each level perfectly from the beginning, the "level" approach here just doesn't help—it complicates everything while solving almost nothing.

Suffering: Programming with GetProcAddress

The fourth way around the "call to undefined dynlink" problem is to return to the **GetProcAddress** programming model to avoid using import libraries altogether. Although this is a difficult programming model to work with, it may be the only choice for an application that really must be robust. Again, the complication is that every call into a DLL turns into two or three calls: **LoadLibrary**, **GetProcAddress**, and the actual call desired. At a minimum, complexity will double, and in our competitive industry this can kill a product (or even a company).

Developers have invented all kinds of interesting strategies to retain the robustness of this programming model while reducing the complexity. One quick optimization is to load any given DLL only once. This means that the application has one global HINSTANCE variable for each DLL it wishes to load in this manner, and each variable is initialized to NULL. Whenever the application needs a certain DLL, it calls an internal function to retrieve the module handle. This function checks if the global variable is NULL, and if so, loads the DLL; otherwise it simply returns the handle that already exists in the global variable. For example, we might have a global array of module handles and symbols for the array index of each DLL we're going to use. Our internal helper function checks the array and loads the module if necessary, as shown here:

```
#define CMODULES      10 //Number of DLLs
#define DLL_ID_SPELLCHECK  0 //Array indices
#define DLL_ID_THESAURUS   1
[other defines here]
typedef BOOL (WINAPI *PFNLOOKUP)(wchar_t *);
HINSTANCE g_rghMods[CMODULES]; //Initialized to NULL on startup.
HINSTANCE ModHandle(UINT id)
{
    ASSERT(id >=0 && id < CMODULES);
    if (NULL==g_rghMods[id])
    {
        g_rghMods[id]=LoadLibrary([name of DLL]);
        if (NULL==g_rghMods[id])
            [Throw exception]
    }
    return g_rghMods[id];
}
void OnToolsSpelling(void)
{
    PFNLOOKUP pfnLookup;
    [Other locals]
    pfnLookup=(PFNLOOKUP)GetProcAddress(ModHandle(DLL_ID_SPELLCHECK),
        MAKEINTRESOURCE(10));
    if (NULL!=pfnLookup)
    {
        [do stuff]
    }
}
```

```

else
    [error]
return;
}

```

The array in **g_rghMods** would be initialized to NULL on startup, and on shutdown the application would call **FreeLibrary** on anything in here that was non-NULL. (Or this might happen when freeing up memory.) Throwing an exception on **LoadLibrary** failure would allow us to write the kind of code shown in **OnToolsSpelling** where we don't have to check for failure of **ModHandle**—we put that code in an exception handler. All of this takes us quite far towards simplifying this programming model throughout the entire application.

Another complication is that the pointers returned from **GetProcAddress** are defined as the type FARPROC, or as a pointer to a function that takes no arguments and has no return value. Unless you define specific types for each function pointer you're going to use, such as PFNLOOKUP in the code above, you'll get no compile-time type-checking. This means that enormously critical run-time errors can arise when the application either calls a function without the requisite arguments or misinterprets a return value. Even if you take the time to write all the necessary typedefs, you'll probably mistype some of them or use the wrong type somewhere in the source code. You'll still have run-time errors, now caused by extremely hard-to-find bugs in header files and sources that otherwise compile without warnings!

The right answer to this problem is for server vendors to provide you with such typedefs in their header files to begin with. Ultimately the types should be defined by those implementing the service, instead of placing the burden on the application. As we'll see, COM enforces this in an even stronger way that allows strong type checking while eliminating the possibility of bugs caused by subtle definition or usage errors.

Now that we have workable solutions for the loading and typing issues, how can we simplify obtaining the function pointers themselves? It complicates the programming model to call **GetProcAddress** each time we need a pointer. A more efficient way to do this would be to cache all the pointers we might need when we first load the DLL, using the same technique as we used with the module handles. That is, for each DLL we define an array of function pointers. We then change our **ModHandle** function to be **FunctionPointer** so we can make calls like this:

```

void OnToolsSpelling(void)
{
    [other code]
    while ([loop on words]
    {
        if (*(PFNLOOKUP)(FunctionPointer(DLL_ID_SPELLCHECK
            , FUNC_ID_LOOKUP))(pszWord))
            [etc.]
    }
    return;
}

```

Again, we can use exceptions so that **OnToolsSpelling** doesn't have to check for a NULL return value from **FunctionPointer**. We might also write macros for each function we wish to use so we don't have to write such ugly code. In the end, we can be writing code like this:

```

#define CALL_SPELLCHECK_LOOKUP(p) (*(PFNLOOKUP)(FunctionPointer(DLL_ID_SPELLCHECK \
    , FUNC_ID_LOOKUP)))(p))
void OnToolsSpelling(void)
{
    [other code]
    while ([loop on words]
    {
        if (!CALL_SPELLCHECK_LOOKUP(pszWord))
            [etc.]
    }
    return;
}

```

Now we've simplified the programming model to the point where it is just as easy to work with as when we were using import libraries, yet we avoid all the import library trouble.

But what happens when we have multiple versions of the server to deal with? We have to expect that some of the functions we want to put in our table of pointers will not actually exist in the server. How should we handle this? A sophisticated application would do something like this:

- Define multiple function tables for each server, where each table represents a group of functions that make up a certain "feature." The application itself would choose the exact groupings. For example, there might be one table for "basic spell checking" that contained only **LookUpWord**, one for "dictionary additions" that contained only **AddWord**, one for "dictionary editing" that contained both **AddWord** and **RemoveWord**.
- Define a flag for each "feature" that is set to TRUE only if all the necessary functions are available for that feature.
- Using these flags, allow other application code to enable or disable certain user commands, depending on which features were available from the various servers. For example, if spell-checking was available, the Tools/Spelling command would be enabled. The spell-check dialog might have Add and Edit buttons inside it, where Add would be enabled only if the "dictionary additions" feature were available, and Edit would be enabled only if the "dictionary editing" *and* "dictionary addition" features were available.

In addition, a *really* sophisticated application would provide its own default code for certain functions that might not exist in all servers, rather than disable the feature altogether. For example, a sophisticated word processor might provide its own backup custom dictionary implementation if a suitable server weren't available. When the application discovered the absence of **AddWord** and **RemoveWord** functions in the server, it would store its own entry points in the table instead. Of course, the application would then install its own proxy implementation of **LookUpWord** to filter out custom entries before calling the server's implementation.

Alternatively, the application might store a pointer to a "do-nothing" function in the table so that the rest of its code could trust that the table is completely full of valid pointers, even if some of the functions don't do anything. That's better than having to check for NULL entries in the table.

In short, a really sophisticated application would define and guarantee its own idea of a "contract" for certain features without complicating its own internal programming model.

Sounds great, doesn't it? Assuming that we can solve all the other problems with the registry, generic management APIs, and so on, the reality is still that (1) it is *very costly* to create an application architecture like this, and (2) such an architecture creates a bloated application that might contain 30 to 50 percent more code than one that just lived with import libraries. Sure, when robustness is the top priority and resources are not an issue, you might do this. But to be honest, 95 percent of all vendors probably cannot afford going to this extent.

But what an awful burden to place on every application! Why should applications have to do this? The type of version control described here is a solution, but it's ridiculous to place this kind of burden on every application for just one DLL, let alone dozens that they might use! Is it not more appropriate to place the definition of "features" on the servers themselves? Is it not appropriate to have the servers provide the do-nothing stub functions? Should not the servers create the tables themselves?

The answer is YES! They should! It is wholly appropriate to have only one implementation of this table-construction code (and instances of the tables themselves!) inside the server rather than have it duplicated across many applications, thus causing code bloat (and happy hardware manufacturers!). It is the *server* that should be responsible for fulfilling the "contract" defined for a certain feature, not the *client*! The client should be able to just ask the server, "Do you support this contract?" and if the answer is yes, invoke any function in that contract without having to check for the existence of individual functions, without having to create tables, without having to provide default implementations, and—most important—without having to do any of the work to define the function types, the macros, and the helper functions.

We are now ready to see that the COM architecture of objects and interfaces, especially *multiple interfaces*, gives us *exactly* what we need to solve versioning without undue burden on clients and without unnecessary complexity for servers.

The Design of the Component Object Model

In the preceding sections I've hinted at how the Component Object Model (a combination of specification and standard system-provided implementation) provides the solutions to a wide range of requirements of a component software architecture. Table 2 below reviews what we have already learned.

Problem	Solution
Path dependencies, multiple providers of a service	Use the registry to map from abstract class identifiers to absolute server locations as well as to map between categories and class identifiers.
Decentralized definition of identifiers	Use GUIDs generated by an algorithm to guarantee uniqueness across time and space, eliminating the need for centralized identifier allocation.
Specific management APIs for each service category	Supply a very simple generic and universal management API that accommodates all service categories, called "Implementation Location"
Sharing object instances across process/machine boundaries	Implement marshalling of interface pointers with Local/Remote Transparency providing the ability to implement a server as an .EXE or .DLL..
Different in-process, local, and remote programming models	Design a single model for all types of client-object connections supported through the Local/Remote Transparency interface structure.

Lifetime management of servers and objects	Institute universal reference counting through the IUnknown base interface that all objects support and from which all other interfaces are derived.
Multiple services per server module	Use the CLSID:IID:table_offset syntax instead of module:ordinal to absolutely identify functions.
Versioning	Support the concept of multiple immutable interfaces as well as interfaces that are strongly typed at both compile time and run time.

Table 2. Summary of Solutions to Component Software Problems

The following sections describe GUIDs, interfaces, implementation location and the registry, local/remote transparency, the **IUnknown** interface, and how all of these elements help solve the versioning problem.

Globally Unique Identifiers

Within the scope of an application or project where one group controls all the source code, it is trivial to avoid naming conflicts between functions, objects, and so on. The compiler will tell you if there's a conflict. In a component software environment, however, especially a distributed one, you have the possibility of many different developers in different organizations having to name functions and modules and classes without running into conflicts. This can either be coordinated through a centralized organization, which is inefficient and expensive, or it can be achieved by a standardized algorithm such as the one created by the Open Systems Foundation (OSF) for their Distributed Computing Environment (DCE).

In this environment, remote procedure calls (RPC) need to travel around large distributed networks knowing exactly what piece of code must be called. The solution is an identifier standard called the Universally Unique Identifier (UUID) generated with an algorithm that uses a machine's Internet protocol (IP) address (unique in *space*, guaranteed unique by network hardware manufacturers) and the current date/time when the algorithm is run (unique in *time*). Uniqueness in both time and space means that no matter when and where the algorithm is run, it will produce a unique value. This is possible because UUIDs are 16-byte values (128 bits), in which you can create a staggering 3.4×10^{38} combinations. Put another way, if everyone in the world (5.5×10^9) generated ten trillion (10^{13}) GUIDs per second, round the clock, it would take *196 million years* to run out of GUIDs. I desperately hope that all this software is vastly obsolete by that time!

COM and OLE use UUIDs for all unique identifications needs, simply calling them "globally unique identifiers," or GUIDs, instead. (Claiming *universal* uniqueness is a bit pretentious—other civilizations on other planets may be using the same IP addresses. So COM calls them *globally* unique, which is more accurate!) According to the DCE standard, GUIDs and UUIDs are spelled out in hexadecimal digits in the following format: {01234567-1234-1234-1234-012345678AB}.

The algorithm itself is commonly implemented in a tool called UUIDGEN, which will spit out one of these UUIDs any time you ask. The Win32® SDK tool supports a command-line switch, *-nXXXX*, to generate XXXX sequential UUIDs. Another Win32 SDK tool, GUIDGEN, generates one GUID at a time in a variety of text formats that are more directly useful in source code.

You use these tools whenever you need a GUID to assign to a category specification, a class implementation, or an interface definition, trusting the absolute uniqueness. And because interfaces are the fundamental extension mechanism in COM and OLE, you can innovate without having to ask permission from anyone.

Interfaces

Interfaces are the point of contact between a client and an object—it is only through an interface that a client and object communicate. In general, an interface is a semantically related group of member functions that carry no implementation (nor any data members). The group of functions acts as a single entity and essentially represents a "feature" or a "design pattern." Objects then expose their features through one or more interfaces, as we'll see shortly.

The definition of an interface is that of an "abstract base class" in C++ parlance. Interfaces are more formally described using the Microsoft COM extensions to the standard DCE Interface Definition Language (IDL). We've seen one example already for a hypothetical **ISpellChecker**, where the "I" is a conventional prefix denoting "interface" (the **object** attribute is the IDL extension that describes a COM interface as opposed to an RPC interface):

```
[uuid(388a05f0-626d-11cf-a231-00aa003d7352), object]
interface ISpellChecker : IUnknown
{
    HRESULT LookUpWord(OLESTR *pszWord);
}
```

The Microsoft IDL (MIDL) compiler can generate header files and various other files from an IDL script like this. At compile time, the interface name (**ISpellChecker** in the example above) is a data type that can be used for type checking. The other important element is the **uuid** attribute, which assigns an IID to the interface. The IID is used as the run-time type of the interface. In other words, the IID identifies the intent of the interface design itself, which means that once defined and assigned an IID, an interface is immutable—any change requires assignment of a new IID because the original intent has now changed.

At run time, an "interface" is always seen as a pointer typed with an IID. The pointer itself points to another pointer that points to a table that holds the addresses of the implementation of each member function in the interface. This binary structure, illustrated in Figure 3, is a core standard of COM, and all of COM and OLE depend upon this standard for interoperability between software components written in arbitrary languages. As long as a compiler can reduce language structures down to this binary standard, it doesn't matter how one programs a component or a client—the point of contact is a run-time binary standard.

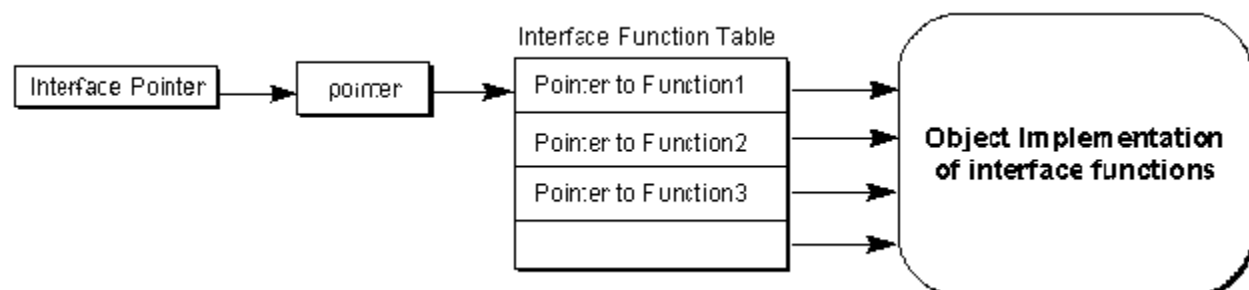


Figure 3. The binary interface structure

It is this exact interface structure that provides the ability to marshall one of these pointers between processes and machines, as described in the section titled "Local/Remote Transparency."

To "implement an interface" on some object means to build this exact binary structure in memory and provide the pointer to the structure. This is what we want instead of having clients do it themselves! You can do this in assembly language if you want, but higher-level languages, especially C++, build the structures automatically. In fact, the interface structure is, by design, identical to that used for C++ virtual functions—programming COM and OLE in C++ is highly convenient.

This is also why COM calls the table portion the "vtable" and the pointer to that table "lpVtbl." A pointer to an interface is a pointer to **lpVtbl**, which points to the vtable. Because this is what C++ expects to see, calling an interface member given an interface pointer is just like calling a C++ object's member function. That is, if I have a pointer to **ISpellChecker** in the variable *pSC*, I can call a member like this:

```
pSC->LookUpWord(pszWord);
```

Because the interface definition describes all the argument types for each interface member function, the compiler will do all the type checking for you. As a client, you never have to define function prototypes for these things yourself.

In C, the same call would look like this, with an explicit indirection through **lpVtbl** and passing the interface pointer as the first argument:

```
pSC->lpVtbl->LookUpWord(pSC, pszWord);
```

This is exactly what C++ does behind the scenes, where the first argument is the *this* pointer. Programming COM and OLE in C++ just saves you a lot of extra typing.

What is also very important about calls to interface member functions is that the address of the call itself is discovered at run time using the value of the interface pointer itself. Therefore the compiler generates *code* that calls whatever address is in the right offset in the vtable pointed to by the interface pointer, which is a value known only at run time. This means that when the kernel loads a client application, there are no "import records" to patch up with absolute addresses. Instead, those addresses are computed at run time. In short, using interfaces is a true form of dynamic linking to a component's functionality, but one that bypasses all the complexity of using **GetProcAddress** and all the risks of using import libraries.

Finally, note that interfaces are considered "contracts" in the sense that when an object implements an interface it must provide at least default or do-nothing code *for every member function in the interface*—every vtable element must contain a valid function pointer. Therefore a client who obtains an interface pointer can call any member function of the interface. Granted, some member functions may just return a "not implemented" error code but the call will always occur. Therefore clients need not check for NULL entries, nor must they ever provide their own default implementations.

Implementation Location and the Registry

When a client wishes to use an object, it always starts by asking COM to locate the object's class server, request the server to create an object, and return an initial interface pointer back to the client. From

that point the client can obtain additional interface pointers from the same object through the **IUnknown::QueryInterface** member function.

For some of its own native classes OLE provides specific creation APIs to streamline the initialization process. But for all custom components (those that are not implemented in OLE itself), COM provides a generic creation API, **CoCreateInstance**, that instantiates an object when given its CLSID. We saw earlier how a client calls this function:

```
if (SUCCEEDED(CoCreateInstance(clsID, NULL, CLSCTX_SERVER, IID_ISpellChecker, (void **)&pSC)))
```

The client passes a CLSID, some flags, and the IID of the initial interface. On output the pointer variable passed by reference in the last argument receives the pointer to the interface.

Internally, COM takes care of mapping the CLSID to the server, loading that server into memory, and asking the server to create the object and return an interface pointer. If the object is in a different process from the client, COM automatically marshalls that pointer to the client's process. The basic process of object instantiation is shown in Figure 4.

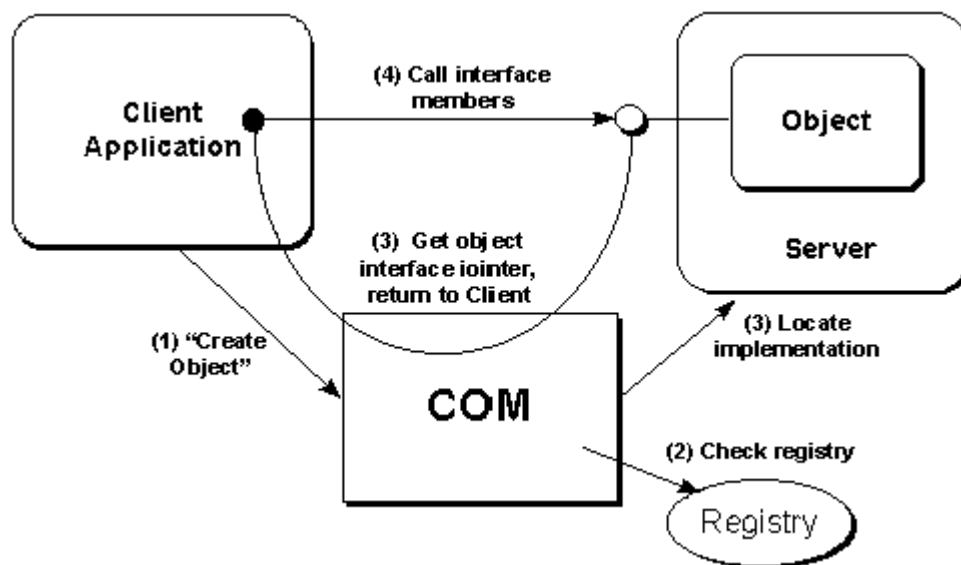


Figure 4. Locating and activating an object

COM looks in the registry to map the CLSID to its server. Servers implemented as Win32 DLLs ("in-process servers") are registered as follows, assuming {01234567-1234-1234-1234-012345678AB} is the CLSID:

```
HKEY_CLASSES_ROOT
  CLSID
    {01234567-1234-1234-1234-012345678AB}
      InprocServer32 = <path to server DLL>
```

To load the DLL into memory, COM needs only to call **LoadLibrary**. A "local server," on the other hand, is implemented when an .EXE uses the **LocalServer32** key instead of **InprocServer32**. For these, COM calls

the **CreateProcess** Win32 API to launch the .EXE, which then initializes COM in its own process space. When this happens, COM in the server's process connects to COM in the client's process.

Entries for a "remote server" (one that runs on another machine) include the machine name as well, and when COM activates such a server it communicates with a resident "service control manager" (SCM) process running on the other machine. The SCM then loads or launches the server on the server machine and sends a marshalled interface pointer back to the client machine and process.

The net result in all three cases is that the client has an interface pointer through which it can begin using the object's services. We're now ready to see how COM's Local/Remote Transparency makes the programming model identical for all three cases.

Local/Remote Transparency: Marshalling and Remoting

When you think of pointers to objects—such as objects implemented in C++—you normally don't even think about passing such a pointer to another process. Ridiculous, right? Well, this is exactly what COM's Local/Remote Transparency allows you to do, if the pointer is an *interface* pointer (which, conveniently, can be generated from a C++ object pointer quite easily, if the C++ class is itself derived from an interface type).

When an in-process object is involved, COM can simply pass the pointer directly from the object to the client, because that pointer is valid in the client's address space. Calls through that pointer end up directly in the object code, as they should, making the in-process case a fast calling model—just as fast as using raw DLLs.

Now, COM cannot, obviously, just pass the object's exact pointer *value* to other processes when local or remote objects are involved. Instead, the mechanism called "marshalling" builds the necessary interprocess communication structures. To "marshall" a pointer means *to create a "marshalling packet" containing the necessary information to connect to the object's process*. This packet is created through the COM API **CoMarshalInterface** function. The packet can then be transported (through any means available) to the client process, where another function, **CoUnmarshalInterface**, turns the packet into an interface pointer (in the client's process), which the client can then use to make calls.

This "marshalling" sequence creates a "proxy" object and a "stub" object that handle the cross-process communication details for that interface. COM creates the "stub" in the object's process and has the stub manage the real interface pointer. COM then creates the "proxy" in the client's process, and connects it to the stub. The proxy then supplies the interface pointer that is given to the client. Those familiar with RPC will recognize this proxy/stub setup as being the same architecture used in raw RPC.

Of course, this proxy does not contain the actual *implementation* of the interface. Instead, each member function packages the arguments it receives into a "remoting" packet and passes that packet to the stub using remote procedure calls. The stub unpacks these arguments, pushes them on the stack, and calls the real object (using the interface pointer the stub is managing). The object executes the function and returns its output. The output is packaged up and sent back to the proxy, which unpacks the output and returns it to the client. This process is illustrated in Figure 5, which shows the differences between in-process, local, and remote cases. The figure also shows that the client only sees in-process objects, which is why we use the term "transparency"—remoting a call is transparent to the client.

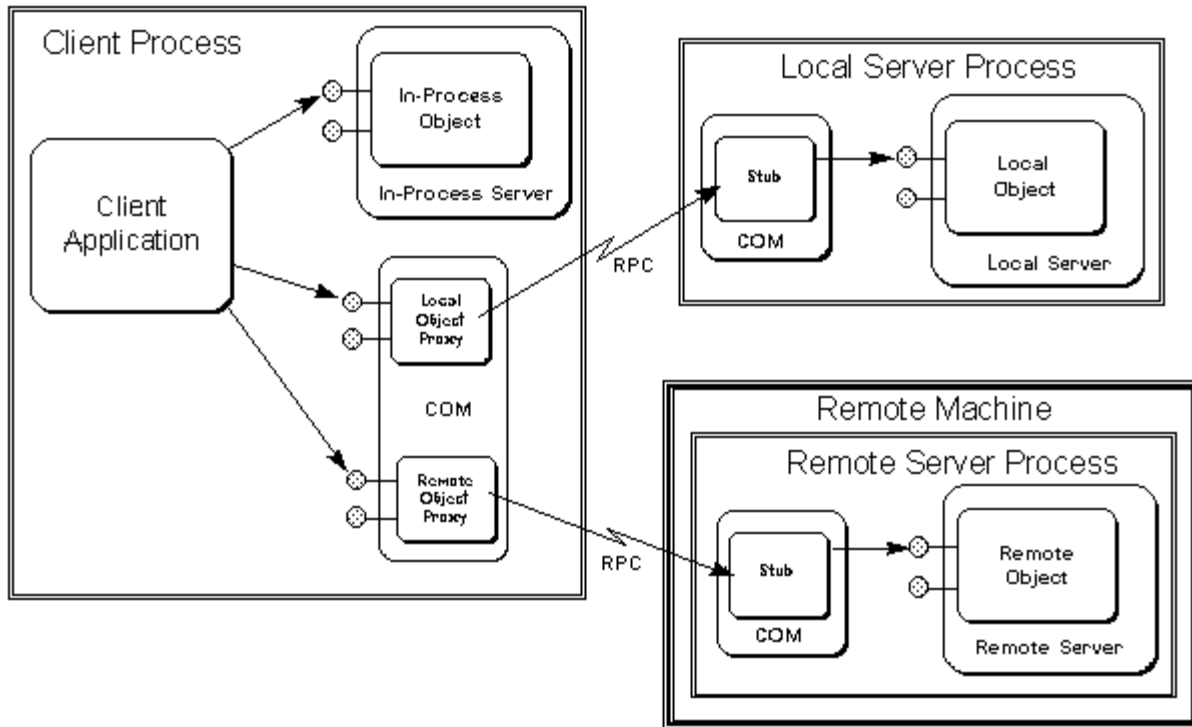


Figure 5. Local/Remote Transparency

An important feature of this architecture is that if the object's process is terminated abnormally, only the stubs for that object are destroyed. The proxies that exist in client processes remain active, but are now "disconnected." Therefore a client can still call the proxy as always and not risk crashing—the proxy simply returns a "disconnected" error code. In addition, if a client process is terminated abnormally, the proxy disappears and its associated stub detects the disconnection. The stub can then clean up any reference counts to the object on behalf of the client that went away.

The necessary proxy and stub code for most "standard" interfaces (those defined by Microsoft) is built into the system. If you define your own "custom" interface, you must supply your own proxy/stub code. Fortunately, the MIDL compiler described earlier will generate this code for you from an IDL file. You need only compile the code into a DLL and you have remoting support for that custom interface.

It is important to note in the remote case that it is not necessary to have an OLE implementation on the server machine in order to interoperate with COM on the client machine. All remote interface calls are transmitted with DCE-compatible RPC such that any DCE-aware system can receive the calls and convert them to fit any object architecture, be it COM, CORBA (the Common Object Request Broker), or any other.

Unknown: Reference Counting and Multiple Interfaces

As discussed earlier, each interface can be thought of as representing a "feature" or a "design pattern." Any given object will be a combination of one or more features or patterns, where the combined functionality of those features is what defines the object's own "component category," described with a Category ID (CATID, a GUID). (Until recently, there were so few categories defined that CATIDs were not

used to identify them; rather, specific registry keys like **Insertable** and **Control** were used. Microsoft has now published a specification for CATID-based categorization.)

This *ability of an object to support multiple interfaces* is a key COM/OLE innovation. No other object model has this concept at its core. Yet it is precisely the idea of *multiple interfaces* that *solves the versioning problem*!

The capacity for multiple interfaces depends on the **IUnknown** interface, which is the core interface of all of COM and OLE. In fact, all other interfaces are derived from **IUnknown** and therefore are polymorphic with **IUnknown**. That is, the three **IUnknown** member functions are universally the first three members of *any* interface:

- **AddRef** increments the object's reference count.
- **Release** decrements the object's reference count, freeing the object if the count becomes zero.
- **QueryInterface** asks the object to return a pointer to another interface given its IID.

The **AddRef** and **Release** members together provide for lifetime management of an object. Every independent external reference to any of the object's interfaces carries a reference count that allows multiple clients to independently use the same instance of an object. Only when all clients have released their reference counts will the object destroy itself and free its resources. As described in the last section, remoting stubs will automatically clean up the reference count for any client process that terminates abnormally without first releasing its references.

QueryInterface is what makes multiple interfaces possible. Once a client obtains the *initial* interface pointer to any object, it obtains another interface pointer to the *same* object through **QueryInterface**. To do a query, pass the IID of the interface you want. In return, you get back the interface pointer, if it is available, or an error code that says, "That interface is not supported." If you get a pointer, you can call the member functions of that interface, calling **Release** when you're through with it. If not, you can't possibly call members of that interface! Thus the object easily protects itself from unexpected calls.

QueryInterface, which is always available through every interface pointer, is how a client asks an object, "Do you support the *feature* identified by this IID?" By asking, a client determines the *greatest common set* of interfaces that both it and the object understand. This completely avoids the "least common denominator" problem described earlier, and also completely solves the versioning problem.

The Solution to Versioning

So how does the idea of multiple interfaces solve the real-world versioning issues? Consider the real issues:

- How does one quantify the difference between any two versions?
- How does a client request a specific version?

Multiple interfaces provide precise answers to these questions. Any given revision of an object will support a certain set of interfaces—between versions, interfaces may be added or removed, or new versions of old interfaces may be introduced (the old interfaces cannot be *changed*, but an object would

support both old and new versions of the interface). **QueryInterface** allows a client to check for whatever "version" or set of interfaces it requires.

A critical feature of the COM programming model is that a client *must* call **QueryInterface** to obtain any interface pointer from the object (creation calls like **CoCreateInstance** by definition have built-in **QueryInterface** calls). In other words, COM forces clients to always check for the presence of an interface before ever trying to invoke members of that interface, resulting in this sort of code:

```
//Assume pObj is an IUnknown * for the object.
if (SUCCEEDED(pObj->QueryInterface(IID_<xxx>, *pInterface)))
{
    //Call members of IID through pInterface, then:
    pInterface->Release();
}
else
{
    //Interface is not available, handle degenerate case.
}
```

A client might also be written to query an object for certain interfaces when the object is first created, thereby enabling or disabling certain features based on what interfaces are available. A client can also cache interface pointers to avoid having to call **QueryInterface**

In contrast, the DLL programming model is either far more complex, costly, or risky. If you program with import libraries, you have no way of checking for functional support in an object before invoking the function, so you risk "Call to undefined dynlink" errors. If you program with **GetProcAddress**, you can *at best* simulate exactly what interfaces provide by default—grouping functions together as "feature" and maintaining flags that say whether or not those features are available.

In short, COM universally simplifies the most sophisticated and difficult programming model that, at great cost and for a single application, one could possibly create with **GetProcAddress**.

A concrete example will illustrate exactly how the **QueryInterface** idea solves the versioning problem. Consider a COM-based specification for the category SpellCheck1.0, which says a spell-checker object like BasicSpell1.0 implements the **ISpellChecker** interface. We also have AcmeNote1.0, the client, which uses BasicSpell1.0 through this interface, as shown in Figure 6. At this point, there is little tangible difference between the COM design and the original DLL design with an exported function.



Figure 6. A version 1.0 object and a client that uses it

But we want to write the SpellCheck2.0 specification to add the custom dictionary feature that includes the **AddWord** and **RemoveWord** functions. (A full design would include a way to enumerate all the words in the dictionary, but a discussion of such a design is beyond the scope of this paper.) With a COM

design, we always introduce these functions in a new interface, but we have great flexibility in how we do this:

- Introduce **ISpellChecker2**, which is derived from **ISpellChecker** but has an entirely distinct IID, as required. Objects would implement this new interface, which would then implement **ISpellChecker**. The advantage here is that the two interfaces end up sharing the same vtable, thereby reducing per-instance memory overhead. This is important when the category being defined typically leads to hundreds or thousands of object instances at run time. By sharing the vtable, you save 16 bytes per instance—the **lpVtbl** pointer plus the **IUnknown** entries in the table. Custom controls are an example of where such savings are important. A spell-checker, on the other hand, doesn't need to worry about this because it will usually have only a single instance per process. So this design isn't advantageous over the others.
- Introduce **ICustomDictionary**, which contains the two new functions (plus the **IUnknown** members). Objects would implement this as another interface alongside **ISpellChecker**. This is the most common design choice.
- Introduce two new interfaces, such as **ICustomDictionaryAdditions** and **ICustomDictionaryEditing**, to separate the features of adding words and editing the contents of the dictionary. In the spell-checker case there is little to gain from this separation; however, in other designs different functions may need different interfaces. One usually wants to avoid going to the extreme of defining several interfaces that each have a single member function. This complicates client programming by forcing two or three calls per function: one to **QueryInterface**, one to the member, one to **Release**, depending on how pointers are cached.

For the purposes of our discussion, we'll choose the second option and define **ICustomDictionary** with its own IID:

```
[uuid(8e47bfb0-633b-11cf-a234-00aa003d7352), object]
interface ICustomDictionary : IUnknown
{
    HRESULT AddWord(OLESTR *pszWord);
    HRESULT RemoveWord(OLESTR *pszWord);
}
```

Now that we have the SpellCheck2.0 specification, the developers at Basic Software and Acme Software will, at some point, each decide to upgrade their respective products to fit this new specification. Earlier we saw that it is very difficult with the plain DLL model to have both companies upgrade their products independently. But this is exactly what COM was designed to support! COM allows either company to add support for the new interface without risking compatibility with the old version of the other.

Let's say that Basic Software chooses to release BasicSpell2.0 first, and the spell-checker object now supports both **ISpellChecker** and **ICustomDictionary**. In making this change, the developers don't have to change anything in the existing **ISpellChecker** code except its **QueryInterface** implementation. The version 1.0 code simply continues to exist as it always has, and COM's multiple interface design encourages this.

Now when AcmeNote1.0 encounters this new object, it sees exactly what it did with BasicSpell1.0, as shown in Figure 7. Because AcmeNote1.0 doesn't know about **ICustomDictionary**, BasicSpell2.0 appears

to be exactly the same as BasicSpell1.0. That is, the two are perfectly polymorphic, and no compatibility problems arise.

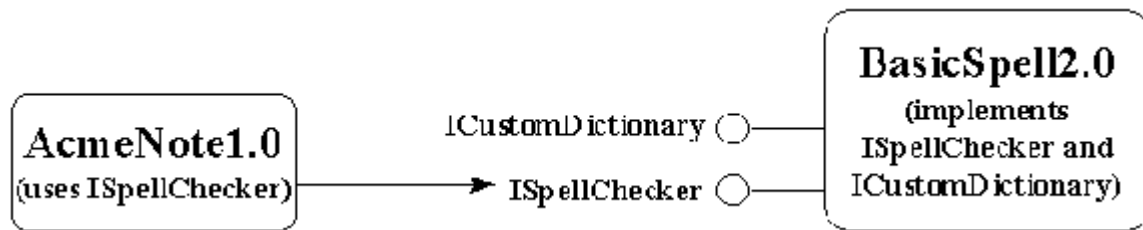


Figure 7. An old client sees a new object exactly as before, because new features show up as new interfaces.

Now let's suppose the reverse happens, that Acme Software releases AcmeNote2.0 first. AcmeNote2.0 is written to enable a custom dictionary feature if the spell checker component it finds supports that feature. At this point, only BasicSpell1.0 exists. But because AcmeNote2.0 must call **QueryInterface** for **ICustomDictionary** in order to use that feature, *it must be coded to expect the absence of the interface!* That is, if a client requests the **ICustomDictionary** interface when initially creating the object, that request will fail and no object is created. A client can specify this interface on creation if it wishes to work *only* with upgraded objects. The client would then provide a meaningful error message, far more meaningful than "Call to undefined Dynlink." The COM programming model forces clients to handle both the presence *and the absence* of a particular interface. Doing this is not difficult, but because you *must* call **QueryInterface**, you must be prepared. So when AcmeNote2.0 encounters BasicSpell1.0, as illustrated in Figure 8, its request for **ICustomDictionary** fails and AcmeNote simply disables its custom dictionary features.

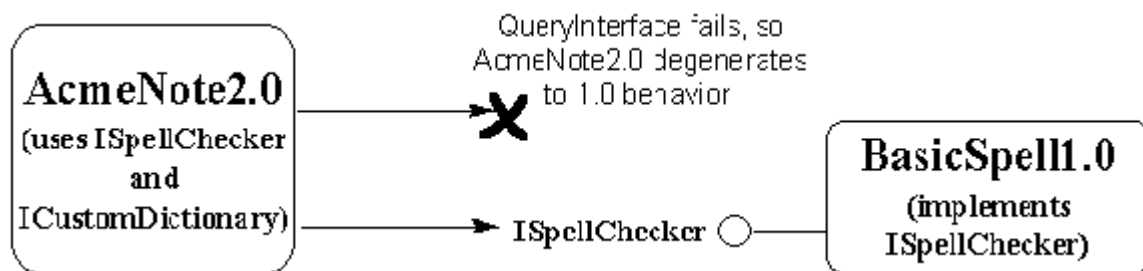


Figure 8. A newer client encounters an older object and degenerates gracefully.

As we saw earlier, the easiest DLL programming model encourages clients to expect the *presence* of exported functions. Catastrophic failure results when those functions are absent, but because it takes extra work to handle such a case robustly, few client programs actually bother. In marked contrast, the COM programming model forces a client to expect the *absence* of interfaces. As a result, clients are written to robustly handle such absences.

What makes the COM model even more powerful is that when a client discovers the *presence* of an interface, it achieves richer integration with the object. That is, expecting the absence of interfaces means the client still works in a degenerate situation (as opposed to failing completely). When interfaces exist on an object, things just get better.

To demonstrate this, let's say that AcmeNote2.0 is written so that it creates an instance of the spell-checker object *only* when the user invokes the Tools/Spelling command. So when users do a spell-check with BasicSpell1.0 installed, they'll get the basic features. Now let's say that *without closing AcmeNote2.0* a user installs BasicSpell2.0. The next time the user invokes Tools/Spelling, AcmeNote2.0 finds that **ICustomDictionary** is present and so it enables additional features that were not available before. That is, once AcmeNote2.0 discovers the availability of the feature, it can enable richer integration between it and the component, as shown in Figure 9. Older clients, like AcmeNote1.0, still continue to work as they always have.

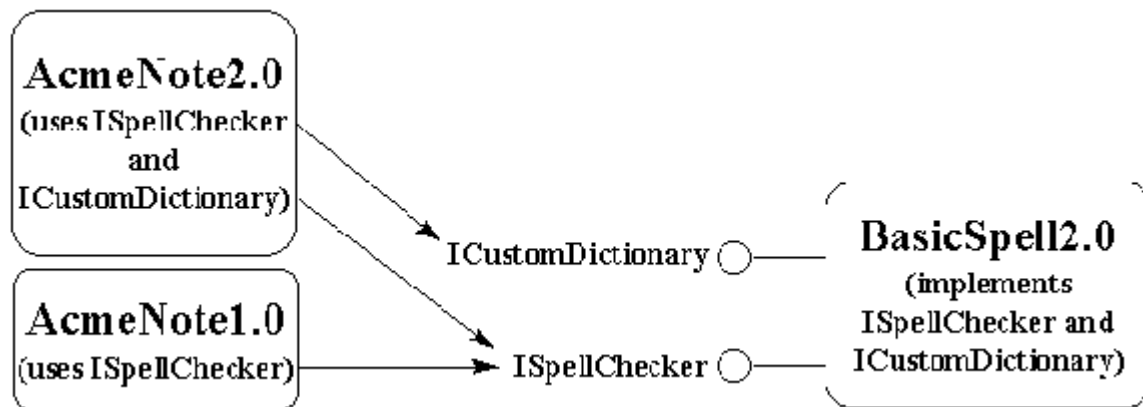


Figure 9. When an upgraded client meets an upgraded component, they integrate on a higher level than before.

This improvement in integration and functionality is *instant* and happens in a *running* system. This is what we mean when we say that COM is meant to solve component software problems in a *binary running* system—without having to turn everything off, we can install a new component and instantly have that component integrate with already running clients, providing new features immediately.

Earlier in this article we saw what happened when an older version of a DLL accidentally overwrote a newer version: Any clients that were implicitly linked to the newer version started to show "Call to undefined dynlink" and wouldn't run at all. We can see here how COM completely solves this problem. Suppose a user installs another application that itself installs BasicSpell1.0 and accidentally overwrites BasicSpell2.0. If we were using raw DLLs and AcmeNote2.0 had implicit links to BasicSpell2.0, AcmeNote2.0 would no longer run at all. But using COM components, AcmeNote2.0 will run just fine, and when the user invokes Tools/Spelling, AcmeNote2.0 will simply not find **ICustomDictionary** and thus will not enable its custom dictionary features. Again, the client expects the *absence* of interfaces and is prepared to work well without them. In addition, because interfaces completely avoid implicit linking, the kernel will never fail to load the application.

It should be obvious now that COM's support for multiple interfaces solves the versioning problem for at least the first revision. But the same solutions also apply over many subsequent versions, as illustrated in Figure 10. Here we have BasicSpell4.0, which contains **ISpellChecker**, **ICustomDictionary**, **ISpecialtyDictionaries**, and maybe even **IThesaurus** and **IGrammar** (extra features, no less!).

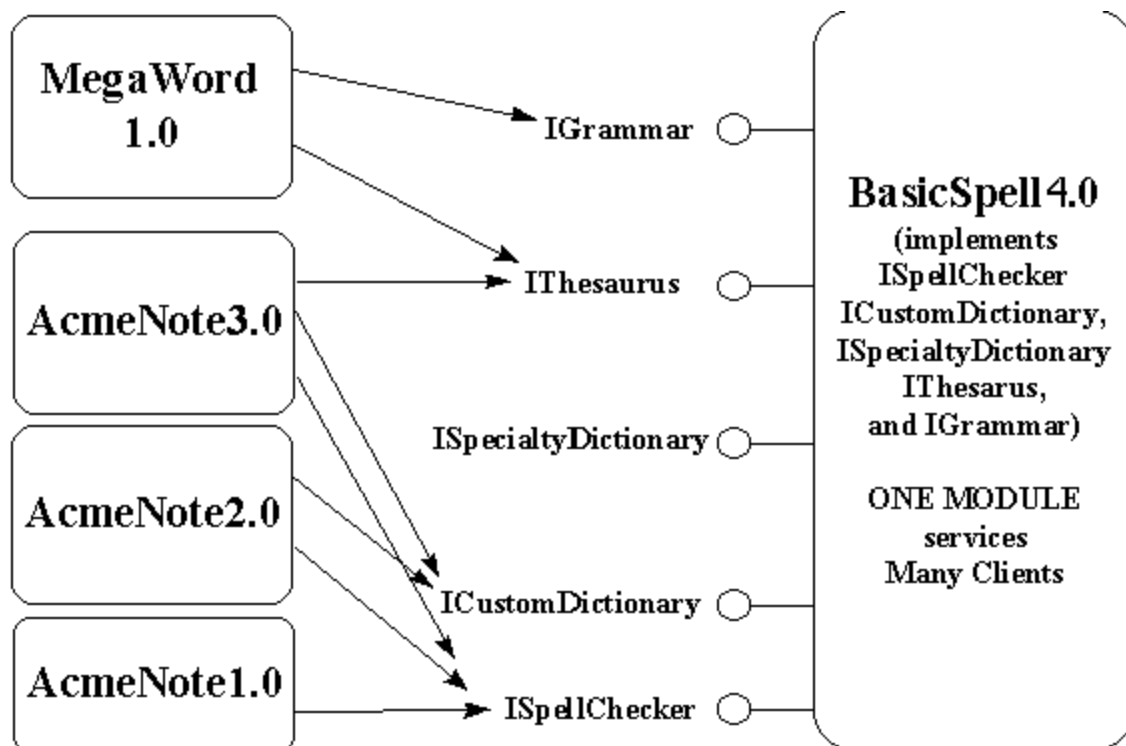


Figure 10. One object can serve many different clients expecting different versions.

This simple version of the object now supports clients of the SpellCheck1.0 specification, clients of the SpellCheck2.0 specification, and those written to specifications that involve the other three added interfaces. Some of those clients may only be using the thesaurus features of this object and never bother with interfaces like **ISpellChecker**. Yet the single object implementation supports them all.

In addition, this new object server can overwrite any old version of the server. Earlier in this paper we saw that building components in the DLL model typically leads to redundant code stored in multiple copies of different DLL versions, like VBRUN100.DLL, VBRUN200.DLL, and so on. COM allows you to factor the differences from version to version in such a way that you can keep the old code available to old clients and provide new code to new clients all within one server module. Hard drive manufacturers may be upset, but users will be happy!

Of course, the newer modules will be larger than the old ones, simply because there's more code for every added interface. However, because all the code is encapsulated behind interfaces, the object can employ many different implementation tricks to (1) reduce the basic module size, (2) improve performance, and (3) minimize memory footprint. In COM, an object is not required to allocate an interface's function table until that interface is requested—each version of an object can thus choose to initially instantiate only those interfaces that are used frequently, leaving others out of memory completely until needed. Code-organization techniques help keep the interface code itself out of memory until it is needed. Such optimizations can minimize the server's load time because only a few code pages are marked "preload." Alternatively, the developer may choose to place the rarely used code in another module altogether, using the object-reuse technique called *aggregation* to instantiate interfaces from that module only when necessary, while still making those interfaces appear as part of the object itself.

Over time, the code for the less-frequently-used interfaces will migrate from being placed in load-on-demand pages in the object module to being separated into another module, to being discarded altogether when those interfaces become obsolete. At such a time, very old clients may no longer run well, but because they were written to expect the absence of interfaces, they'll always degenerate gracefully.

The bottom line is that COM provides a complete solution for versioning problems, at the same time providing a clear and robust migration path for code to gradually move farther out of a working set as it becomes less frequently used. Even when the code is discarded completely, the COM programming model guarantees robust degeneration. COM offers a truly remarkable means for removing obsolete and dead code from a system!

COM vs. "Objects"

There is often a fair bit of debate in the popular trade media about whether or not COM and OLE complement or compete with "objects" or "object-oriented programming." Some say that OLE doesn't support "real" objects, whatever that means. Some think OLE competes with object-oriented programming (OOP) languages like C++, or with frameworks like MFC. All of these claims are untrue, because COM and OLE were designed to complement existing object technologies, as described in the next two sections.

COM/OLE's Problem Space

COM and OLE do not compete with object-oriented languages nor with frameworks because they solve problems in a domain that is not only separate from the problem space for languages and frameworks, but is complementary, as illustrated in Figure 11.

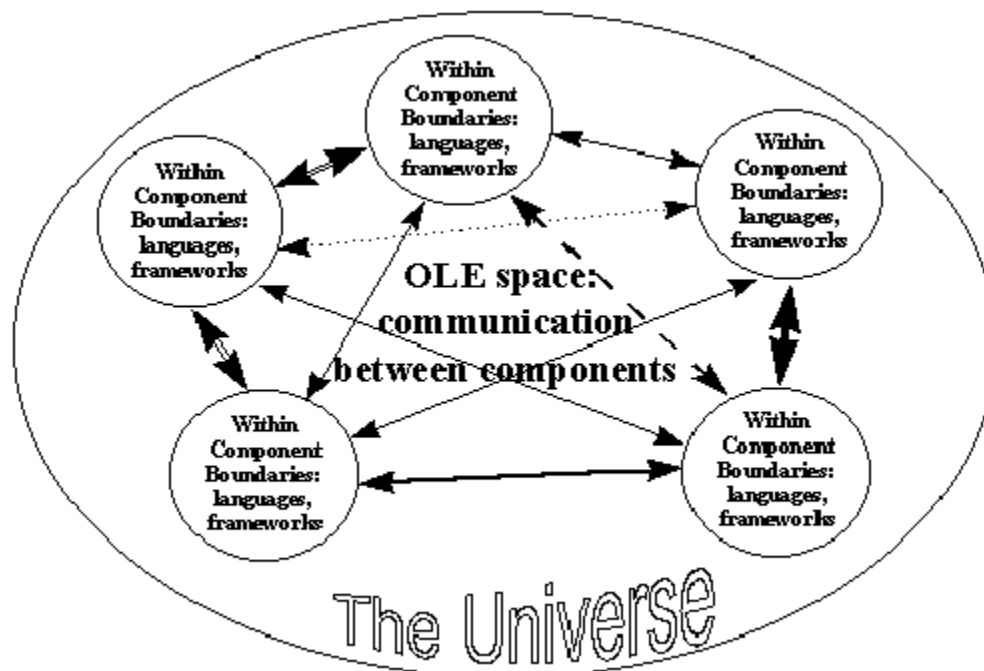


Figure 11. COM and OLE's problem space

Languages and frameworks are absolutely fabulous tools for programming the *internals* of components and objects. However, when you want different objects written by different people to interoperate, especially across time (of deployment) and space (objects running in other processes and running on other machines), COM and OLE provide the means for doing so.

In fact, it is easy to see how one might take a language (such as C++) -based object and create a layer of wrapper code that factors the language object's interface into COM-style interfaces, as shown in Figure 12. By doing so, that language object is suddenly available to use from any other client process through the COM mechanisms. "From CPP to COM" by Markus Horstmann (MSDN Library, Technical Articles) is a good paper on creating this wrapping layer for C++ objects, turning the C++ "interface" into COM interfaces.

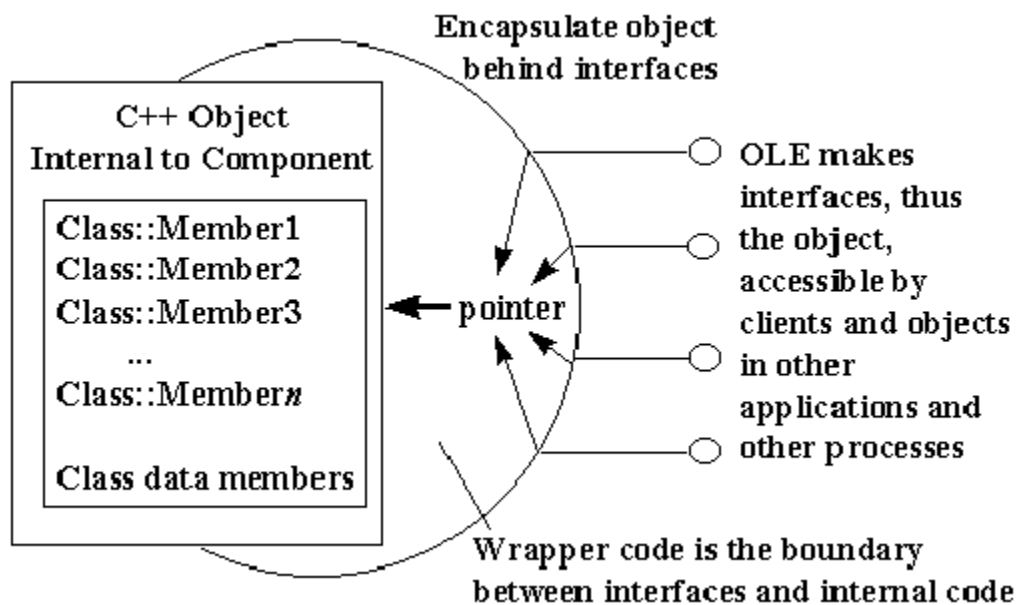


Figure 12. Wrap a language-based object inside interfaces to make the object available to other processes.

To provide OLE support, the framework classes of MFC turn the more raw and complex OLE interfaces for certain high-level features (such as compound documents and OLE Controls) into a simplified C++ class that is easier to work with in application programming. In such cases the interface designs existed first, and MFC developers created a C++ class for them.

So although object-oriented languages are great ways to express object-oriented concepts in source code and great ways to express *implementation details in source code*, OLE is concerned with the *communication* between pieces of *binary code in a running system*. Frameworks like MFC make the boundary between the two worlds simpler by providing default implementation as well as wrapping the communication layers. In this way COM and OLE are to frameworks like MFC as mathematics is to a calculator: A calculator makes mathematical operations effortless but you still have to understand *why* you are using those operations; in the same way, MFC makes certain OLE features very simple to incorporate into an application, but you still have to understand *why* you want those features!

Object Concepts in COM/OLE

If you can take a language-based "object" and wrap it in COM/OLE interfaces, is it still an "object" when seen from a COM/OLE client? Does a client see an entity that is *encapsulated*? Does that entity support *polymorphism*? Does that entity support *reuse*? These are the fundamental concepts behind the "objects" in "object-oriented programming," and COM, at its core, supports all three. Objects in COM/OLE behave as "real" objects in all the ways that matter. The *expression* of the particular concepts is a little different, but the same ends are achieved.

Encapsulation

How does COM support *encapsulation*? Easily. All implementation details are hidden behind the interface structures, exactly as they are in C++ and other languages: The client sees only interfaces and knows nothing about object internals. In fact, COM enforces a stricter encapsulation than many languages because COM interfaces *cannot* expose public data members—all data access must happen through function calls. Certainly language extensions for compilers can let you express public members in *source code*, but on the binary level all data is exchanged through function calls, which is widely recognized as the proper way to do encapsulation.

Polymorphism

How does COM support *polymorphism*? This happens on three different levels:

- First, two interfaces that derive from the same base interface are polymorphic in that base interface, just as they would be in C++. This is because both interfaces have vtable entries that look exactly like those of the base interface. The primary reason for this is that all interfaces are derived from **IUnknown** and thus look exactly alike in the first three entries of the vtable. When you program COM and OLE in C++, you actually use C++ inheritance to express this polymorphic relationship. When programming in C, the interface structures are defined as explicit structures of function pointers where polymorphic interface structures share the same set of initial entries in the structure.
- Second, object classes (and instances) that support the same interface are polymorphic in that interface. That is, if I have two object classes that both support an interface like **IDropTarget**, the exact same client code that manipulates the **IDropTarget** of one object class can be used to manipulate the **IDropTarget** implemented in another class. Such client code exists in the OLE system-level drag-and-drop service, which uses this interface to communicate with any potential recipient (an application, a control, or even the desktop itself on Windows 95) of a drop operation.
- Third, object classes that support the same *set of multiple interfaces* are polymorphic across the entire *set*. Many service categories expect this. The OLE specification for embeddable compound document objects, for example, says that such objects always support the **IOleObject**, **IViewObject2**, **IDataObject**, and **IPersistStorage** interfaces, as well as a few others. A client written to this specification can host any embeddable object regardless of the type of content, be it a chart, video clip, sound bite, table, graphic, text, or other. All OLE Controls are also polymorphic in this manner, supporting a wide range of capabilities behind the same set of polymorphic interfaces.

Reuse

How does COM/OLE support *reuse*? The reader might ask, "Don't you mean *inheritance*?" No—I mean *reuse*. It must be understood that *inheritance* is **not** a fundamental OOP concept; inheritance is how one *expresses* polymorphism in a programming language and how one *achieves* code reuse between classes. *Inheritance is a means to polymorphism and reuse*—inheritance is not an *end* in itself. Many defend inheritance as a core part of OOP, but it is nothing of the sort. Polymorphism and reuse are the things you're really after.

We have already seen how COM supports polymorphism. It supports reuse of one component by another through two mechanisms. The first, *containment*, simply means that one object class uses another class internally for its own implementation. That is, when the "outer" object is instantiated, it internally instantiates an object of the reused "inner" class, just as any other client would do. This is a straightforward client-object relationship, where the inner object doesn't know that its services are being used in the implementation of another class—it just sees some client calling its functions.

The second mechanism is *aggregation*, a more rarely used technique through which the outer object takes interface pointers from the inner object and exposes those pointers through its own **QueryInterface**. This saves the outer object from having any of its own code support an interface and is strictly a *convenient* way to make certain kinds of containment more efficient. This does require some extra coding in the inner object to make its **IUnknown** members behave as the outer object's **IUnknown** members, but this amounts to only a few lines of code. This mechanism even works when multiple levels of inner objects are in use, where the outer object can easily obtain an interface pointer from another object nested down dozens of levels deep.

What OLE Is Made Of

We have seen the reasons why the designs in COM exist and the problems that COM solves. We have seen how COM and OLE complement OOP languages and frameworks, and how COM supports all the fundamental object concepts.

So what, then, is "OLE" itself? Why do people make such a fuss about OLE being "big and slow" or being "difficult to learn"? The reason is that Microsoft has built a large number of additional services and specifications on top of the basic COM architecture and COM services (like Local/Remote Transparency). Taken by itself, COM is simple, elegant, and absolutely no more complex than it has to be. With the right approach, any developer can learn the core of COM in a matter of days and see its real potential.

However, all the layers of OLE that are built on top of COM are overwhelming. There are more than 150 API functions and perhaps 120 interface definitions averaging 6 member functions each. In other words, nearly a thousand bits of functionality that continues to grow! Thus, trying to understand what everything is about is a daunting effort—I spent a few years doing this myself in order to produce the books *Inside OLE 2* (1993) and *Inside OLE, 2nd Edition* (1995) for Microsoft Press. I refer the reader to those titles (primarily the second edition) for a comprehensive and incremental approach to understanding all the pieces of OLE.

For the purposes of this paper, let me provide an overview of this approach to understanding. Step one is to see that *all functionality in OLE falls into three categories*:

1. API functions and interfaces to expose OLE's built-in services, including a fair number of helper functions and helper objects.
2. API functions and interfaces to allow customization of those built-in services.
3. API functions and interfaces that support creation of "custom services" according to various specifications.

The following sections describe each of these in a little more detail.

OLE's Native Services

One reason OLE seems so large is that it offers quite a wide range of native services (we've already discussed some of these):

- Implementation location
- Local/Remote Transparency ("remoting" and "marshalling")
- Task memory allocation
- Structured storage (described in detail below)
- File, Item, Composite, Pointer, and Anti Monikers
- The "running object table"
- Type library and type information creation and management
- Type conversion routines
- Clipboard data exchange
- Drag-and-drop data exchange
- Data caching
- Default "handling" for embedded objects
- Helper and wrapper functions, including support for string and array types.
- Helper objects such as standard dispatch, font, and picture objects, and "advise" holders.

As we've already seen, the services of implementation location, marshalling, and remoting are obviously fundamental to the ability of objects and clients to communicate across any distance. As such, these services are a core part of a COM implementation on any given system.

Other services are not quite as fundamental. Rather, they are implementations of a standard and they support higher-level interoperability. For example, the "Structured Storage" standard describes a "file system within a file" service to standardize how software components can cooperate in sharing an underlying disk file. Because such sharing is essential to cooperation between components, OLE implements the standard in a service called "Compound Files." This eliminates all sorts of interoperability problems that would arise if every component or client application had to implement the standard itself.

To see why a standard like this is important, let us look briefly at how today's file systems themselves came about. When computers first gained mass-storage devices, there weren't these things we call "operating systems." What ran on the computer was *the* application, which did everything and completely owned all system resources, as shown in Figure 13. In such a situation, the application simply wrote its data wherever it wanted on the storage device, itself managing the allocation of the sectors. The application pretty much saw the storage device as one large contiguous byte array in which it could store information however it wanted.

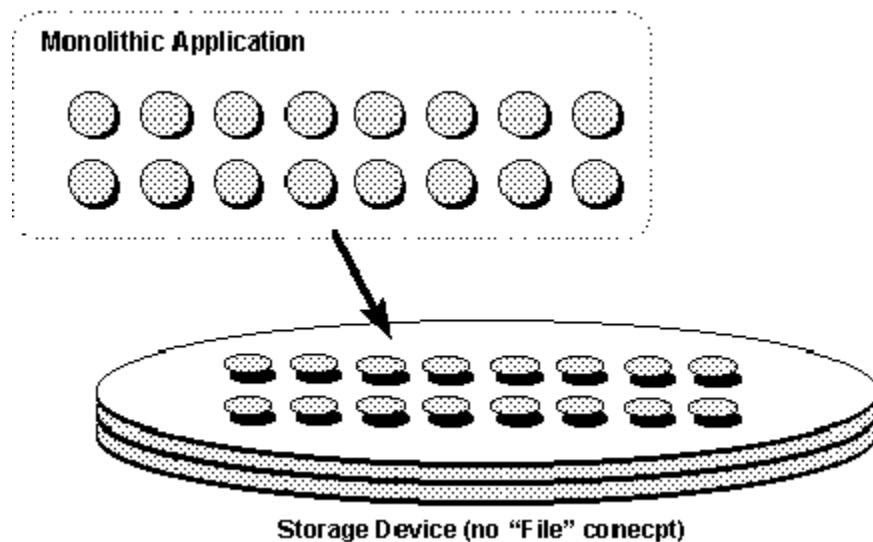


Figure 13. When there was only one "application" on a machine, that application owned all resources, such as the mass-storage device.

Soon it became desirable to run more than one application on the same machine. This required some central agent to *coordinate the sharing* of the machine's resources, such as space on the mass-storage devices. Hence, file systems were born. A file system provides an abstraction layer above the specific allocation of sectors on the device. The abstraction is called the "file" and is made up of usually non-contiguous sectors on the physical device, as shown in Figure 14. The file system maintains a table of the sequence of sectors that make up the "file." When the application wishes to write to the disk, it asks the file system to create a "file," which is to say, "Create a new allocation table and assign it this name." The application then sees the "file" as a continuous byte array and treats it as such, although physically the bytes are not contiguous on disk. In this way the file system controls cooperative resource allocation while still allowing applications to treat the device as a contiguous byte array.

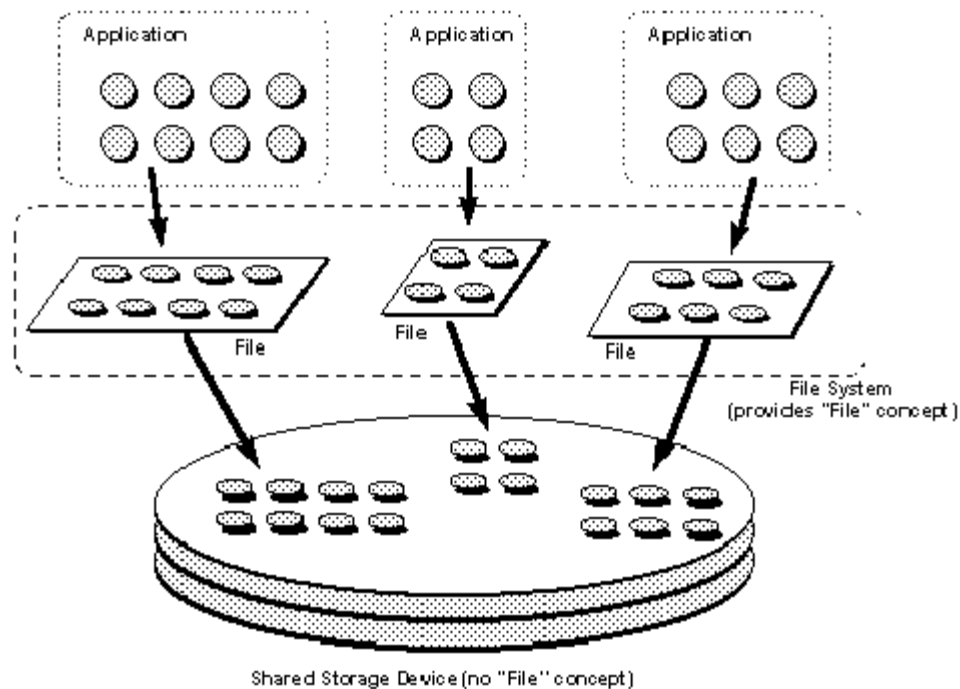


Figure 14. A file system allows multiple applications to share a mass-storage device.

File systems like this work great as long as all parts of the application *itself* are coordinated. In the component software paradigm, this is no longer true! Applications might be built from hundreds of components that were developed at different points in time and space. How, then, can these components work cooperatively in the creation of a single "file" as the user sees it? One solution would have all the components talk to each other through some standardized interfaces so they can negotiate about who gets what part of the file. It should be obvious that such a design would be extremely fragile and horrendously slow.

A better solution is to repeat the file system solution—after all, the problem is the same: different pieces of code must cooperate in using an underlying resource. In this case the underlying resource is a file, not a storage device, although conceptually the two are identical. The OLE Structured Storage specification describes the abstraction layer that turns a single file into a file system itself, in which one can create directory-like elements called "storage objects" (that implement **IStorage**) and file-like elements called "stream objects" (that implement **IStream**), as shown in Figure 15.

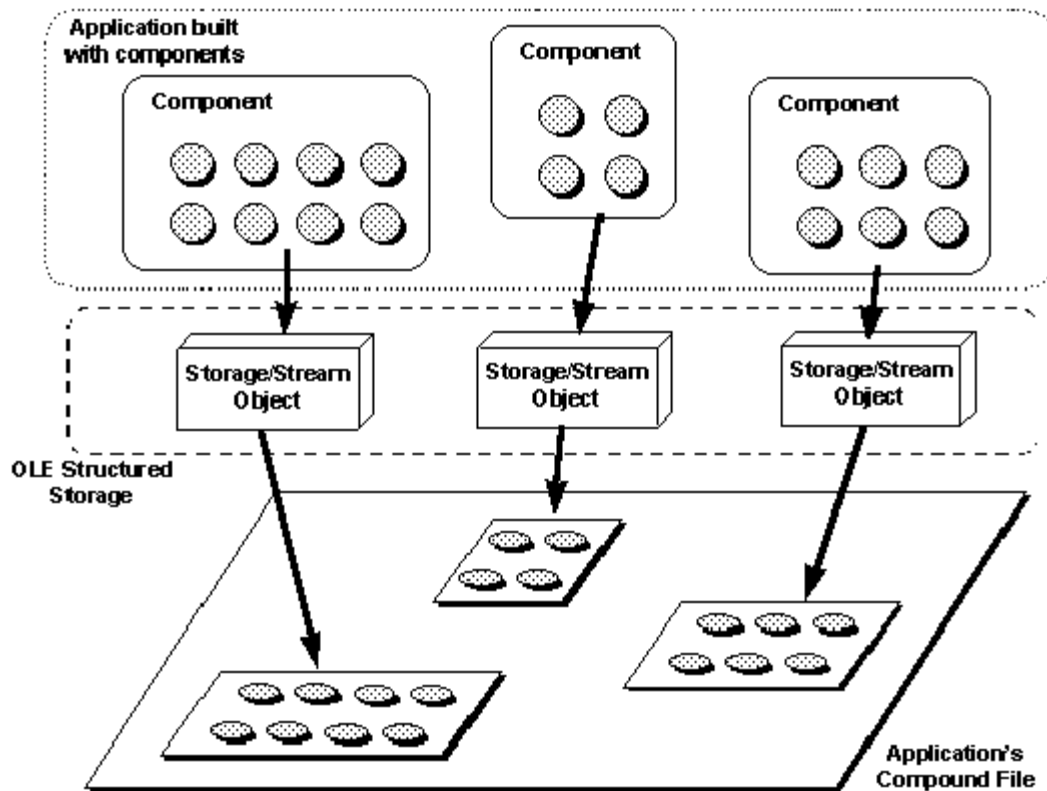


Figure 15. OLE's Structured Storage allows multiple components to share the same file.

OLE's implementation of this standard controls the actual layout of data bits inside the file. The software components themselves, however, see stream objects as contiguous byte arrays once more, so the experience you have working with files translates directly to working with streams. In fact, there is a one-to-one correspondence between typical file-system APIs such as **Read**, **Write**, and **Seek**, and the member functions of the **IStream** interface through which one uses a stream object.

Structured Storage eliminates the need for each component in the application to negotiate its storage requirements; instead, this is coordinated in OLE's implementation. The only necessary bit that the application must communicate to components is to pass each of them an **IStorage** or **IStream** pointer, depending on the object's persistence mechanism. If an object implements, for instance, **IPersistStorage**, it is saying that it would like to have an **IStorage** as a basis for its storage needs. Within the storage object it then creates sub-storages and streams as necessary. If an object implements **IPersistStream**, on the other hand, it says that it only needs a single stream for its storage needs.

Customization of Native Services

Some of the functions in the OLE API and some of the interfaces deal with customizations to OLE's native services. One example will suffice to describe this "customization."

OLE's implementation of the Structured Storage model described in the last section typically works with an underlying file on a file system. However, you can customize this behavior by redirecting the file image to another storage device through what is called a "LockBytes" object (one that implements **ILockBytes**). All this object knows how to do is read or write a certain number of bytes in a certain

location on some storage device. By default, OLE itself uses a LockBytes that talks to a file on the file system. However, you are free to implement your own LockBytes (or use OLE's other implementation that works on a piece of global memory) to redirect data to a database record, a section of another file format, a serial port, or whatever. A few OLE APIs allow you to install your LockBytes underneath the storage implementation. When someone writes data to a stream through **IStream::Write**, for instance, that data ultimately shows up in your **ILockBytes::WriteAt** code, allowing you to place it anywhere you want.

Custom Services

OLE itself can only provide a limited number of native services—specifically those that need to be centralized and standardized. A great deal of OLE's functionality is to support the creation of what I call "custom services" that fit into one or more "service type" specifications.

There are three primary (and sometimes complex) type specifications that Microsoft has defined at present, each of which involves a large number of OLE API functions and interfaces. They are:

- OLE Documents, for the creation and management of compound documents. The specification describes either how you create "embeddable" or "linkable" objects or how you create a "container" that hosts such objects, with in-place activation as an optional feature.
- OLE Controls, for the creation and management of custom controls. The specification describes how to create or host any kind of custom control in an OLE model. Currently there are hundreds, if not thousands, of vendors working on controls or control containers.
- OLE Automation, for programmability and scripting. The specification describes how you create either an automation object or an automation "controller" that can drive that object according to a script of some kind. Visual Basic exploits OLE Automation to a large extent.

Each "type" specification describes the general interaction between an object of the type and a client of such an object. From the specification one can implement a wide assortment of specific objects that fall into these categories. For example, a sound bite, a video clip, a chart, a picture, a presentation, text, and a table are all examples of specific instances of the "embeddable compound document object" type. Buttons, edit boxes, labels, check boxes, tabbed dialogs, list boxes, scrollbars, toolbars, sliders, dials, and status lines are all examples of custom controls, each of which can be implemented as an OLE control. OLE Automation is typically used to create objects that expose the functionality of an entire application, like Microsoft Word, so that a scripting engine like Visual Basic can drive that application to perform certain operations without the need for a user to execute the steps (that is, "automation" in its most generic sense).

In addition to these Microsoft-defined categories, other industry groups have defined service categories for their particular needs. Some examples include the industries of real-time market data, health care, insurance, point-of-sale, and process control. Other groups are actively working on more specifications.

No matter who defines the service, the idea is that many people can implement clients or objects according to the specification, which describes the abstraction that either side (client or object) uses to view the other, so that all clients are polymorphic to the objects and all objects of the category are polymorphic to clients.

What is most important here is that COM is an open architecture in which anyone can create and define new service categories without any input or involvement on Microsoft's part and without any dependency on Microsoft enhancing its operating systems. This independence was intentionally built into the design. With COM, Microsoft has created an extensible service architecture in which anyone—without any centralized control whatsoever—can define new service categories, publish them, implement components and clients around them, and deliver those components to customers.

The true "openness" of an architecture lies in its *design and extensibility*, not just in the *process* through which the design occurs (which at Microsoft is *open*, just not *committee-based*!). COM and OLE enable decentralized and asynchronous innovation, design, development, and deployment of component software, the absolute essential elements of a successful distributed object system.

OLE's Life Purpose

In this paper we have seen why COM and OLE exist, the problems they were designed to solve, the services they provide, and the open flexibility that is inherent to their design. But you're probably asking two questions: *What does OLE mean to me?* and *What does OLE mean to my project?*

If there's one word to describe what OLE is all about—technically speaking—it is **integration**. Think about what integration means to you and the projects you're working on—do you think about integration with the system? How about integration with other applications? Or integration with third-party add-on components? Whatever questions you ask, OLE probably has an answer for you. For example, the Windows 95 and Windows NT 4.0 shells have extension features based on OLE components. Integration between applications, whether through OLE Documents, OLE Automation, Drag and Drop, or other private protocols, almost always involves OLE. And the OLE Component Object Model provides an ideal architecture for plug-ins or add-on models because it solves all those problems that appear across multiple versions and revisions of interfaces.

Everything you need to exploit this integration on a single machine is already in place: COM in its present form shipped nearly three years ago and is already installed in tens of millions of desktops. Since that time Microsoft and other parties have continually added new enhancements that add even more power and more capabilities, such as the introduction of Network OLE in Windows NT 4.0.

All of this leads us to an even greater question: *What is it all for?* We understand what OLE is about, technically, but what is OLE *about*? In what direction is OLE moving? What does OLE mean for the software industry as a whole? What is OLE's "life purpose"? I contend that what OLE is *really* about is much larger than one company such as Microsoft, yet still very personal to everyone involved. OLE's purpose is to make true component software a working reality. Many groups have claimed to have "superior component technology," yet OLE's COM is the only one that solves the right problems for the right reasons.

Why all this fuss about component software? It is a long-standing dream in the software industry to be able to quickly assemble applications from reusable components, allowing us to meet customer demands more quickly. Object-oriented programming to some extent brought us closer to that goal, but it is limited to source code problems. In the real world, problems also need to be solved in a binary running system, because one cannot recompile and redeploy all the software in a system every time one object implementation needs to change. OLE and COM were designed specifically to address the issues

of components in a running system and in running applications. This is why COM is based on a binary standard, not a language or pseudo-language standard. COM is the glue that makes component integration work right.

But we must realize, as an industry, that the purpose of component software is more than just making it easier for *developers* to assemble applications, because no matter what your methodology, developer-written code cannot perfectly match customer needs when they arise. That is, our development methodologies usually involve asking customers what their *problems* are, which assumes they can even articulate them! Once we know the problems, we can try to formulate solutions, and after about a *minimum* three-to-six month period we might be able to deliver a solution. However, the current "application backlog," as it's called, is *three* years, not a matter of months.

The inherent problem here is that by the time we can deliver a prepackaged solution, the customer's problems and needs have changed—often times drastically. We might deliver a solution that solves old problems, but not new ones. So we're caught in a treadmill, always trying to catch up. We try to preach this or that development methodology as "the solution," forgetting that *customers want immediate solutions to their present problems!* Customers care about solutions, not OOP, languages, or "object purity." And if customers care about solutions, then so should developers, because customers are the *single source of income in the computer industry*.

One way we've tried to solve this problem is by creating ever more complex and generic *tools* that we have the nerve to call "applications." These do-it-all applications, with every feature imaginable, generally require the user to *mentally map their problem to the tool*. This is exactly why software seems so hard to use, because this mental mapping is so difficult.

So what's special about component software? It does not just help us run faster in the treadmill. Component software is about enabling end users to quickly and easily construct exactly the applications they need to solve their immediate problems.

This might be done by adding a component to an existing application to provide some new capability that wasn't there before. In a sense, a user who makes a small customization to an existing one has created a new application. Many customer problems can be solved this way, but it requires a strong component architecture to do it. COM is just such an architecture.

As a complement, *users should be able to express their exact problem to the computer and have the computer assemble the application from available components*. If that application isn't quite right, the user can clarify the shortcomings and try again. I believe that with this positive feedback loop—one that developers are very accustomed to—it will be possible for users to express their problem so easily that the computer can construct the right application within a matter of minutes. When the problem is solved, the user throws that application away. After all, the definition of an "application" is "a particular solution to a particular problem" and in the component world, the "application" is nothing more than a list of connections between components.

Component software is about putting more power into the hands of users instead of concentrating it only in the hands of developers. This is precisely what helps users solve their own problems. Of course, there will be many components that are also targeted towards developers as well, for many developers

will be creating higher-level components from lower-level components. Component software serves everyone!

Yet there are two requirements to make this vision of component software really work. One will be the software that allows the user to naturally express a problem, and then turn it into an application. Hopefully this will someday involve natural speech recognition.

The other requirement is a great diversity of components. Science knows that any successful system—biological, ecological, societal, or technological—requires diversity, which is the natural state of the universe. Imagine a world in which the only plant is grass, the only animals are mice, and all the humans are lawyers. Pretty scary. A diversity of components is necessary for a successful component software system, ranging from low, medium, and high-level functions targeted to all people from the low-level hack to the most naive end user.

But where will these components come from? Today there exists a tremendous amount of software functionality inside large, monolithic applications. Unfortunately, because it is part of a monolith, that functionality is only available to that *one* application. Through a process called *componentization*, one breaks up a large application into smaller and smaller components, thereby making that functionality available in many other ways. This won't happen overnight, of course—applications will probably be broken into large pieces first, then those pieces broken down, and those broken down, and so on. It is one of COM and OLE's design points to allow this large-scale componentization to occur *gradually*, while at the same time allowing reintegration of small components into larger ones.

For example, you can first make components available from within a large .EXE that is riddled with legacy code. Then behind the scenes, you can rewrite that code and break it down into smaller and more efficient pieces as you see fit. The COM architecture provides a smooth path through the whole componentization process—from legacy code to small and fast components. You don't have to reengineer everything to receive benefits.

Once a monolith has been broken down into components, you should be able to reassemble the original application. But because there are now possibly hundreds of components, they can also be reassembled in thousands of other ways, thereby creating many other useful applications *from the same code base*. Each of these applications will be far more suited to a specific customer task than was the general-purpose monolith, and each of them will be smaller and easier to use. Customers will be able to solve their problems sooner.

Again, this is what component software is all about! Simply by virtue of programming in a component model that encourages well-factored software designs, we can create thousands of useful components that can be integrated in millions of ways in order to solve customer problems!

So where do you, today's software developer, fit into this world view? There are many opportunities, because Microsoft cannot do it all. Microsoft's role is to provide the "plumbing" on Windows and Macintosh operating systems, and will, of course, be trying to produce some of the best components available. But there are many other opportunities, not only in competing with Microsoft in providing high-value components, but also in providing the plumbing for other operating systems, in designing new services or component types, in integrating components for customers, and in creating the end-

user tools that integrate components automatically. There are likely to be many other opportunities that we cannot even imagine today!

Should you be scared? I think not. When users are empowered to create their own applications and solve their own problems, there will be explosive growth in this industry! The computer industry history shows that each innovation that empowers more people to solve their own problems has generated a surge of growth. This happened with operating systems, with compilers, with graphical user interfaces, and I believe it will happen with component software. Who knows, users might start selling their "applications" to each other, which ultimately would create a larger market with even more users.

This potential for market growth is very exciting, because in an expanding universe, *everyone* has new space in which to grow. Today's software industry (and business in general) is in a phase of serious consolidation, with a new merger announced weekly. When everyone is crowded together in a limited space, the only way for anyone to grow is to consume one another. But in an expanding market, space itself increases and large companies like Microsoft cannot, logistically speaking, grow fast enough to fill in the gaps created in such an expansion. Thus many opportunities will arise for everyone—to grow, to create new businesses, and ultimately to succeed.

What OLE is *really* about is your success.