

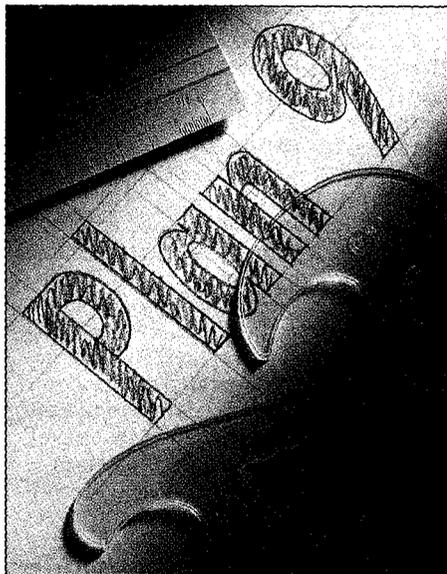
Designing Plan 9

Bell Labs' Plan 9 research project looks to tomorrow

Rob Pike, Dave Presotto, Ken Thompson,
and Howard Trickey

Plan 9 is a distributed computing environment assembled from separate machines acting as CPU servers, file servers, and terminals. The pieces are connected by a single file-oriented protocol and local name space operations. Because the system was built from distinct, specialized components rather than similar general-purpose components, Plan 9 achieves levels of efficiency, security, simplicity, and reliability seldom realized in other distributed systems. This article discusses the building blocks, interconnections, and conventions of Plan 9.

Unhappy with the trends in commercial systems, we began a few years ago to design a system that could adapt well to changes in computing hardware. In particular, we wanted to build a system that could profit from the continuing improvements in personal machines with bitmap graphics, in medium- and high-speed networks, and in high-performance microprocessors. A common approach is to connect a group of small personal timesharing systems — workstations — by a medium-speed network, but this has a number of failings. Because each workstation has private data, each must be administered sepa-



rately; maintenance is difficult to centralize. The machines are replaced every couple of years to take advantage of technological improvements, rendering the hardware obsolete, often before it has been paid for. Most telling, a workstation is a largely self-contained system, not specialized to any particular task; too slow and I/O-bound for fast compilation; too expensive to be used just to run a window system. For our purposes — primarily software development — it seemed that an approach based on distributed specialization rather than compromise would best address issues of cost-effectiveness, maintenance, performance, reliability, and

security. We decided to build a completely new system, including compiler, operating system, networking software, command interpreter, window system, and terminal. This construction would also offer an occasion to rethink, revisit, and perhaps even replace most of the utilities we had accumulated over the years.

Plan 9 is divided along lines of service function. CPU servers concentrate computing power into large (not overloaded) multiprocessors; file servers provide repositories for storage; terminals give each user a dedicated computer with bitmap screen and mouse on which to run a window system. Sharing computing and file storage services provides a sense of community for a group of programmers, amortizes costs, and centralizes and simplifies management and administration.

The pieces communicate by a single protocol, built above a reliable data transport layer offered by an appropriate network, that defines each service as a rooted tree of files. Even for services not usually considered as files, the unified design permits some noteworthy and profitable simplification. Each process has a local filename space that contains attachments to all services the process is using and thereby to the files in those services. One of the most important jobs of a terminal is to support its user's customized view of the entire system as represented by the services visible in the name space.

To be used effectively, the system

The authors are all researchers at AT&T Bell Laboratories, Murray Hill, NJ 07974. This paper was originally delivered at the UKUUG Conference in London in July 1990 and is reprinted here with permission from the UKUUG.

(Continued from page 48)

requires a CPU server and a file server (large machines best housed in an air conditioned machine room with conditioned power) and a terminal. The system is intended to provide service at the level of a departmental computer center or larger, and its strengths stem in part from economies of scale. Accordingly, one of our goals is to unite the computing environment for all of AT&T Bell Laboratories (about 30,000 people) into a single Plan 9 system comprising thousands of CPU and file servers spread throughout, and clustered in, the company's various departments. That is clearly beyond the administrative capacity of workstations on Ethernets.

The following sections describe the basic components of Plan 9, explain the name space and how it is used, and offer examples of unusual services that illustrate how the ideas of Plan 9 can be applied to a variety of problems.

CPU Servers

Several computers provide CPU service for Plan 9. The production CPU server is a Silicon Graphics Power Series machine with four 25-MHz MIPS processors, 128 Mbytes of memory, no disk, and a 20 Mbyte-per-second back-to-back DMA connection to the file server. It also has Datakit and Ethernet controllers to connect to terminals and non-Plan 9 systems. The operating system provides a conventional view of processes, based on *fork* and *exec* system calls, and of files, mostly determined by the remote file server. Once a connection to the CPU server is established, the user may begin typing commands to a command interpreter in a conventional-looking environment.

A multiprocessor CPU server has several advantages. The most important is its ability to absorb load. If the machine is not saturated (which can be economically feasible for a multiprocessor), there is usually a free processor ready to run a new process. This is similar to the notion of free disk blocks in which to store new files on a file system. The comparison extends farther: Just as you might buy a new disk when a file system gets full, you may add processors to a multiprocessor when the system gets busy, without needing to replace or duplicate the entire system. Of course, you may also add new CPU servers and share the file servers.

The CPU server performs compilation, text processing, and other applications. It has no local storage; all the permanent files it accesses are provided by remote servers. Transient parts of

the name space, such as the collected images of active processes or services provided by user processes, may reside locally but these disappear when the CPU server is rebooted. Plan 9 CPU servers are as inter-changeable for their task — computation — as are ordinary terminals for theirs.

*The file server presents
to its clients a file
system rather than,
say, an array of
disks or blocks
or files*

File Servers

The Plan 9 file servers hold all permanent files. The current server is another Silicon Graphics computer with two processors, 64 Mbytes of memory, 600 Mbytes of magnetic disk, and a 300 gigabyte jukebox of write-once optical disk (WORM). (This machine is to be replaced by a MIPS 6280, a single processor with much greater I/O bandwidth.) It connects to Plan 9 CPU servers through 20 Mbyte-per-second DMA links, and to terminals and other machines through conventional networks.

The file server presents to its clients a file system rather than, say, an array of disks or blocks or files. The files are named by slash-separated components that label branches of a tree, and may be addressed for I/O at the byte level. The location of a file in the server is invisible to the client. The true file system resides on the WORM, and is accessed through a two-level cache of magnetic disk and RAM. The contents of recently-used files reside in RAM and are sent to the CPU server rapidly by DMA over a high-speed link, which is much faster than regular disk although not as fast as local memory. The magnetic disk acts as a cache for the WORM and simultaneously as a backup medium for the RAM. With the high-speed links, it is unnecessary for clients to cache data; the file server centralizes the caching for all its clients, avoiding the problems of distributed caches.

The file server actually presents several file systems. One, the "main" system, is used as the file system for most clients. Other systems provide less gen-

erally-used data for private applications. One service is unusual: the backup system. Once a day, the file server freezes activity on the main file system and flushes the data in that system to the WORM. Normal file service continues unaffected, but changes to files are applied to a fresh hierarchy, fabricated on demand, using a copy-on-write scheme. Thus, the file tree is split into two parts: A read-only version representing the system at the time of the dump, and an ordinary system that continues to provide normal service. The roots of these old file trees are available as directories in a file system that may be accessed exactly as any other (read-only) system. For example, the file `/usr/rob/doc/plan9.ms` as it existed on April 1, 1990, can be accessed through the backup file system by the name `/1990/0401/usr/rob/doc/plan9.ms`. This scheme permits recovery or comparison of lost files by traditional commands such as file copy and comparison routines rather than by special utilities in a backup subsystem. Moreover, the backup system is provided by the same file server and the same mechanism as the original files, so permissions in the backup system are identical to those in the main system; you cannot use the backup data to subvert security.

Terminals

The standard terminal for Plan 9 is a Gnot (with silent "G"), a locally-designed machine of which several hundred have been manufactured. The terminal's hardware is reminiscent of a diskless workstation: with 4 or 8 Mbytes of memory, a 25-MHz 68020 processor, a 1024 × 1024 pixel display with 2 bits per pixel, a keyboard, and a mouse. It has no external storage and no expansion bus; it is a terminal, not a workstation. A 2 megabit per second packet-switched distribution network connects the terminals to the CPU and file servers. Although the bandwidth is low for applications such as compilation, it is more than adequate for the terminal's intended purpose: To provide a window system, that is, a multiplexed interface to the rest of Plan 9.

Unlike a workstation, the Gnot does not handle compilation; that is done by the CPU server. The terminal runs a version of the CPU server's operating system, configured for a single, smaller processor with support for bitmap graphics, and uses that to run programs such as a window system and a text editor. Files are provided by the standard file server over the terminal's network connection.

Just like old character terminals, all Gnots are equivalent, as they have no

(continued from page 50)

private storage either locally or on the file server. They are inexpensive enough that every member of our research center can have two — one at work and one at home — and see exactly the same system on both. All the files and computing resources remain at work where they can be shared and maintained effectively.

Networks

Plan 9 has a variety of networks that connect the components. To connect components on a small (computer center or departmental) scale, CPU servers and file servers communicate over back-to-back DMA controllers. More distant machines are connected by traditional networks such as Ethernet or Datakit, which a terminal or CPU server may use completely transparently except for performance considerations. Because our Datakit network spans the country, Plan 9 systems could potentially be assembled on a large scale. (See Figure 1.)

To keep their cost down, Gnats employ an inexpensive network that uses standard telephone wire and a single-chip interface. (The throughput is respectable, about 120 Kbytes per second.) Getting even that bandwidth to home, however, is problematic. Some of us have DS-1 lines at 1.54 megabits per second; others are experimenting with more modest communications equipment. Because the terminal only mediates communication — it instructs the CPU server to connect to the file server but does not participate in the resulting communication — the relatively low bandwidth to the terminal does not affect the overall performance of the system.

Name Spaces

There are two kinds of name space in Plan 9: The global space of the names of the various servers on the network and the local space of files and servers visible to a process. Names of machines and services connected to Datakit are hierarchical: *nj/mb/astro/belix*, for example, roughly defines the area, building, department, and machine. Because the network provides naming for its machines, Plan 9 need not directly handle global naming issues. It does, however, attach network services to the local name space on a per-process basis. This is used to address the issues of customizability, transparency, and heterogeneity.

The protocol for communicating with Plan 9 services is file-oriented; all services, local or remote, are arranged into a set of file-like objects collected into a hierarchy called the name space of the server. For a file server, this is a trivial requirement. Other services must sometimes be more imaginative. For instance, a printing service might be implemented as a directory in which processes create files to be printed. Other examples are described in the following sections. For the moment, consider just a set of ordinary file servers distributed around the network.

When a program calls a Plan 9 service, (using mechanisms inherent in the network and outside Plan 9 itself) the program is connected to the root of the service's name space. Using the protocol, usually as mediated by the local operating system into a set of file-oriented system calls, the program accesses the service by opening, creating, removing, reading, and writing files in the name space.

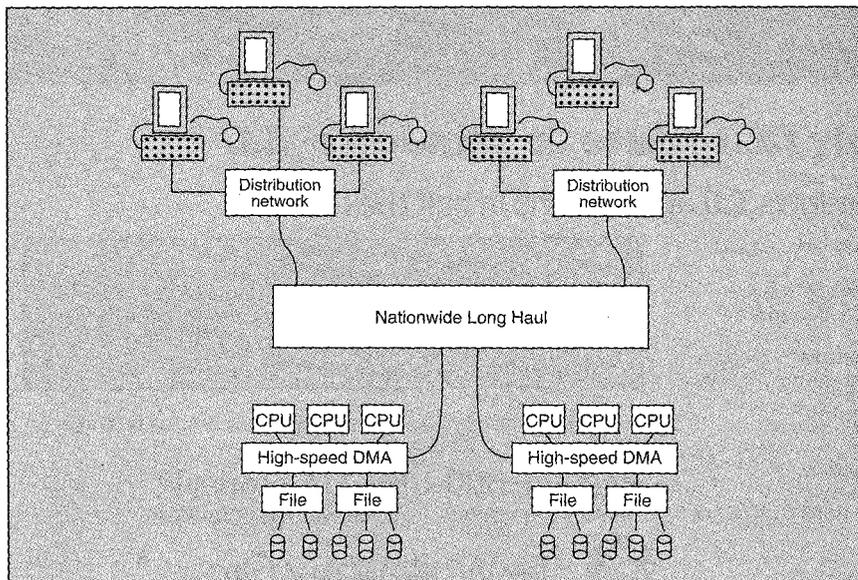


Figure 1: Plan 9 topology

After the user selects desired services (file servers containing personal files, data, or software for a group project, for example), their name spaces are collected and joined to the user's own private name space by a fundamental Plan 9 operator called *attach*. The user's name space is formed by the union of the spaces of the services being used. The local name space is assembled by the local operating system for each user, typically, the terminal. The name space is modifiable on a per-process level, although in practice the name space is assembled at login time and shared by all that user's processes.

To login to the system, the user instructs the terminal which file server to connect to. The terminal calls the server, authenticates the user (described later), and loads the operating system from the server. It then reads a file, called the "profile," in the user's personal directory. The profile contains commands that define what services to use by default, and where in the local name space to attach them. For example, the main file server to be used is attached to the root of the local name space, "/", and the process file system is attached to the directory */proc*. The profile then typically starts the window system.

Within each window, a command interpreter may be used to execute commands locally, using file names interpreted in the name space assembled by the profile. For computation-intensive applications such as compilation, the user runs a command *cpu* that selects (automatically or by name) a CPU server to run commands. After typing *cpu*, the user sees a regular prompt from the command interpreter. But that command interpreter is running on the CPU server in the same name space — even the same current directory — as the *cpu* command itself.

The terminal exports a description of the name space to the CPU server, which then assembles an identical name space, so the customized view of the system assembled by the terminal is the same as that seen on the CPU server. (A description of the name space is used rather than the name space itself so the CPU server may use high-speed links when possible, rather than requiring terminal intervention.) The *cpu* command affects only the performance of subsequent commands; it has nothing to do with the services available or how they are accessed.

The following are a few examples of the usage and possibilities afforded by Plan 9.

(continued on page 54)

(continued from page 52)

The Process File System

An example of a local service is the "process file system," which permits examination and debugging of executing processes through a file-oriented interface. It is related to Killian's process file system but its differences exemplify the way that Plan 9 services are constructed.

The root of the process file system is conventionally attached to the directory */proc*. (Convention is important in Plan 9; many programs have conventional names built in that require the name space to have a certain form. For example, it doesn't matter which server the command interpreter */bin/rc* comes from, but it must have that name to be accessible by the commands that call

on it.) After attachment, the directory */proc* itself contains one subdirectory for each local process in the system, with name equal to the numerical unique identifier of that process. (Processes running on the remote CPU server may also be made visible; this will be discussed shortly.) Each subdirectory contains a set of files that implement the view of that process. For example, */proc/77/mem* contains an image of the virtual memory of process number 77. That file is closely related to the files in Killian's process file system, but unlike Killian's, Plan 9's */proc* implements other functions through other files, rather than through peculiar operations applied to a single file. Table 1 shows a list of the files provided for each process.

The *status* file illustrates how heterogeneity and portability can be handled by a file server model for system functions. The command *cat/proc/*status* presents the status of all processes in the system; in fact, the process status command *ps* is just a reformatting of the ASCII text so gathered. The source for *ps* is a page long and is completely portable across machines. Even when */proc* contains files for processes on several heterogeneous machines, the same implementation works.

Filename	Description
<i>mem</i>	The virtual memory of the process image. Offsets in the file correspond to virtual addresses in the process.
<i>ctl</i>	Control behavior of the processes. Messages sent (by a <i>write</i> system call) to this file cause the process to stop, terminate, resume execution, and so on.
<i>text</i>	The file from which the program originated. This is typically used by a debugger to examine the symbol table of the target process, but is in all respects except name the original file; thus one may type <i>/proc/77/text</i> to the command interpreter to instantiate the program afresh.
<i>note</i>	Any process with suitable permissions may write the <i>note</i> file of another process to send it an asynchronous message for interprocess communication. The system also uses this file to send (poisoned) messages when a process misbehaves, for example, divides by zero.
<i>status</i>	A fixed-format ASCII representation of the status of the process. It includes the name of the file the process was executed from, the CPU time it has consumed, its current state, and so on.

Table 1: Files provided for the "process file system"

The DGIS™ SDK and a TI 34010-based High-Performance Graphics Board for one amazing price.

High performance, high resolution graphics are the wave of the future. With the DGIS Software Developer's Kit™ (SDK), qualified software developers can write for the future today.

The DGIS Developer's Kit provides everything needed to develop applications and drivers for DGIS-compatible 34010 graphics boards—boards from companies such as Compaq, Dell, Hewlett-Packard, NCR, NEC, TI and more than 30 others worldwide. Software developed with this kit can access the full power of the 34010, supporting the greatest number of high resolution graphics boards at the highest levels of performance, resolution and color.

DGIS, the premier and most widely-shipped interface for the TI 340X0 family of graphics coproces-

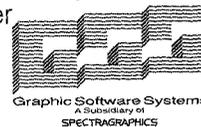
THE POWER OF HIGH RESOLUTION GRAPHICS PROGRAMMING CAN BE REACHED WITH ONE EASY NUMBER:



sors, provides an outstanding feature-rich programming model with 100+ graphics functions. The DGIS SDK includes documentation and language bindings for the DGIS interface, device drivers for Windows 3.0, utilities, and the GSS AT1050™ 1024X768 34010 graphics board (which normally sells for \$1295 alone).

The DGIS SDK is compatible with most C compilers and supports the XMS standard as well as DOS Extenders from Rational and PharLap.

Stepping up to the big screen has never been easier or more attractive. Call today.



**Call (503) 641-2455.
Ask for Dept. DGIS-3.**

All prices subject to change without notice. GSS, DGIS, The DGIS Software Developer's Kit, GSS AT1050 are trademarks of Graphic Software Systems Inc. All other trademarks belong to their respective owners.

CIRCLE NO. 507 ON READER SERVICE CARD

input and output to go through that file, the window system can simulate the expected properties of the file.

The window system serves several files, all conventionally attached to the directory of I/O devices, */dev*. These include *cons*, the port for ASCII I/O; *mouse*, a file that reports the position of the mouse; and *bitblt*, which may be written messages to execute bitmap graphics primitives. Much as the different *cons* files keep separate clients' output in separate windows, the *mouse* and *bitblt* files are implemented by the window system in a way that keeps the various clients independent. For example, when a client process in a window writes a message (to the *bitblt* file) to clear the screen, the window system clears only that window. All graphics sent to partially or totally obscured windows are maintained as bitmap layers, in memory private to the window system. The clients are oblivious of one another.

Because the window system is implemented entirely at user level with file and name space operations, it can be run recursively: It may be a client of itself. The window system functions by opening the files */dev/cons*, */dev/bitblt*, and so forth, as provided by the operating system, and reproduces — multiplexes — their functionality among its clients. Therefore, if a fresh instantiation of the window system is run in a window, it will behave normally, multiplexing its */dev/cons* and other files for its clients. This recursion can be used profitably to debug a new window system in a window or to multiplex the connection to a CPU server. Because the window system has no bitmap graphics code — all its graphics operations are executed by writing standard messages to a file — the window system may be run on any machine that has */dev/bitblt* in its name space, including the CPU server.

The *cpu* Command

The *cpu* command connects from a terminal to a CPU server using a full-duplex network connection and runs a setup process there. The terminal and CPU processes exchange information about the user and name space, and then the terminal-resident process becomes a user-level file server that makes the terminal's private files visible from the CPU server. (At the time of writing, the CPU server builds the name space by reexecuting the user's profile; a version being designed will export the name space using a special terminal-resident server that can be queried to recover the terminal's name space.) The CPU process makes a few adjustments

to the name space, such as making the file */dev/cons* on the CPU server be the same file as on the terminal, and begins a command interpreter. The command interpreter then reads commands from, and prints results on, its file */dev/cons*, which is connected through the terminal process to the appropriate window (for example) on the terminal. Graphics programs such as bitmap editors may also be executed on the CPU server because their definition is entirely based on I/O to files "served" by the terminal for the CPU server. The connection to the CPU server and back again is utterly transparent.

This connection raises the issue of

heterogeneity: The CPU server and the terminal may be, and in the current system are, different types of processors. There are two distinct problems: binary data and executable code. Binary data can be handled two ways: By making it not binary or by strictly defining the format of the data at the byte level. The former is exemplified by the *status* file in */proc*, which enables programs to examine, transparently and portably, the status of remote processes. Another example is the file, provided by the terminal's operating system, */dev/time*. This is a fixed-format ASCII representation of the number of seconds since the epoch that

Power and Craftsmanship

Elegance, Precision, Speed. Greenleaf Functions™ 4.0.

The first general C library ever...and still the best. *So you don't have to do all that grunt work!*

Over 400 powerful, time-tested, well-crafted functions in C and assembler. *Here's a peek under the hood!*

- DOS 4.x and 3.3 support, including disk volumes over 32MB.
- Printer: Buffered Interrupt, DOS Spooler, and via DOS or BIOS.
- File Compression / Expansion.
- Keyboard with (text mode) Mouse.
- Fast Color Text (Direct to video memory.)

- DOS Functions - Incl. Drive Status, File Time/Date - 45 Functions.
- Interrupts incl. Cntrl-Break Handler.
- Graphics & Sound.
- Video Effects, Controls, Cursor.
- Color Print, incl. Centered, Justified.
- String Manipulation, including String Pointer List operations.
- Keyboard, including single-call Enhanced Keyboard, integrated with Ctrl-Break Handler.
- Time, Date, Many Formatting Options.

"Systems Using Greenleaf Libraries are Making Us Rich"

(Ken Baldry - Art & Science Ltd., London)

Call Now for Demo Pack

1-800-523-9830

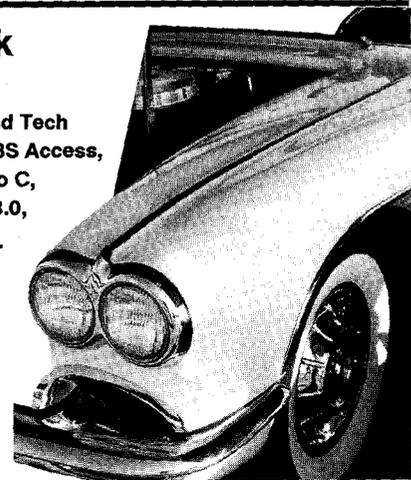
No Royalties. FREE Full Source Code and Tech Support. Professional documentation, BBS Access, Newsletter. Supports Microsoft 6.0, Turbo C, Lattice 6.0, Zortech C++, Watcom 7.0 & 8.0,



JPI TopSpeed C.
All memory models,
of course.

FAX: (214)248-7830

CIRCLE NO. 479 ON READER SERVICE CARD



Powerful Solutions to the User Interface Puzzle

VITAMIN C

Vitamin C gives you the power to create complex applications with a consistent user interface. Applications that are easier to write. Easier to use!

Unrivaled versatility makes you more productive and more competitive. In fact, Vitamin C's open ended design provides virtually limitless customization potential. Even without modifying the source!

Advanced features like programmer defined functions before, on-exit, and after each field, menu, and menu-item complete the picture.

Top it off with comprehensive printed documentation, no royalties, and our 90-day money back guarantee and you have more than a function library. You have solutions. You have Vitamin C!

Complete the picture with proven, reliable resources

Windows: multiple, overlapping, virtual, pop-up, pull-down, moving, scrolling...

Data Entry: forms, fields, validation, formatting, auto-blank, select from list...

Menus: multi-level, Lotus, Mac, separators, blank items, unavailable items...

Help: context-sensitive, pop-up, default, key word highlight...

Keyboard: function keys, jump to function, translation, programmable handler, idle...

Text editor: notepad, word-wrap, justify, search, insert, delete...

Library source included free!

Portable: versions for DOS, OS/2, Unix, Xenix, VAX!

Order Today!

214-416-6447

FAX: 214-418-1915

BBS: 214-418-0059

CREATIVE
PROGRAMMING

Box 112097 Carrollton, Texas 75011-2097

serves as a time base for *make* and other programs. Processes on the CPU server get their time base from the terminal, thereby obviating problems of distributed clocks.

For files that are I/O intensive, such as */dev/bitblt*, the overhead of an ASCII interface can be prohibitive. In Plan 9, therefore, such files accept a binary format in which the byte order is predefined, and programs that access the files use portable libraries that make no assumptions about the order. Thus */dev/bitblt* is usable from any machine, not just the terminal. This principle is used throughout Plan 9. For instance, the format of the compilers' object files and libraries is similarly defined, which means that object files are independent of the type of the CPU that compiled them.

Having different formats of executable binaries is a thornier problem, and Plan 9 solves it adequately if not gracefully. Directories of executable binaries are named appropriately: */mips/bin*, */68020/bin*, and so on, and a program may ascertain, through a special server, what CPU type it is running on. A program, in particular the *cpu* command, may therefore attach the appropriate directory to the conventional name */bin* so that when a program runs, say, */bin/rc*, the appropriate file is found. The various object files and compilers use distinct formats and naming conventions, which makes cross-compilation painless, at least once automated by *make* or a similar program.

Security

Plan 9 does not address security issues directly, but some of its aspects are relevant to the topic. Breaking the file server away from the CPU server enhances security possibilities. Because the file server is a separate machine that can only be accessed over the network by the standard protocol, and therefore can only serve files, it cannot run programs. Many security issues are resolved by the simple observation that the CPU server and file server communicate using a rigorously controlled interface through which it is impossible to gain special privileges.

Of course, certain administrative functions must be performed on the file server, but these are available only through a special command interface accessible only on the console and hence subject to physical security. Moreover, that interface is for administration only. For example, it permits making backups and creating and removing files, but not reading files or changing their permissions. *The contents of a file with read permission for only its owner*

will not be divulged by the file server to any other user, even the administrator.

This begs the question of how a user proves who he or she is. At the moment, we use a simple authentication manager on the Datakit network itself, so that when a user logs in from a terminal, the network assures the authenticity of the maker of calls from the associated terminal. In order to remove the need for trust in our local network, we plan to replace the authentication manager by a Kerberos-like system.

Discussion

A fairly complete version of Plan 9 was built in 1987 and 1988, but development was abandoned. In May of 1989 work was begun on a completely new system, based on the SGI MIPS-based multiprocessors, using the first version as a bootstrap environment. By October, the CPU server could compile all its own software, using the first-draft file server. The SGI file server came on line in February 1990; the true operating system kernel at its core was taken from the CPU server's system, but the file server is otherwise a completely separate program (and computer). The CPU server's system was ported to the 68020 in 13 hours elapsed time in November, 1989. One portability bug was found; the fix affected two lines of code. At the time this article was originally written, work had just begun on a new window system, which has since been implemented. An electronic mail system has also been added, clearing the way for use of Plan 9 on a daily basis by all the authors and 50 to 60 other users. Plan 9 is now up, running, and comfortable to use, although it is certainly too early to pass final judgment.

The multiprocessor operating system for the MIPS-based CPU server has 454 lines of assembly language, more than half of which save and restore registers on interrupts. The kernel proper contains 3647 lines of C plus 774 lines of header files, which includes all process control, virtual memory support, trap handling, and so on. There are 1020 lines of code to interface to the 29 system calls. Much of the functionality of the system is contained in the "drivers" that implement built-in servers such as */proc*; these and the network software add another 9511 lines of code. Most of this code is identical on the 68020 version; for instance, all the code to implement processes, including the process switcher and the *fork* and *exec* system calls, is identical in the two versions; the peculiar properties of each processor are encapsulated in two five-line assembler routines. (The code for

the respective MMUs is quite different, although the page fault handler is substantially the same.) It is only fair to admit, however, that the compilers for the two machines are closely related, and the operating system may depend on properties of the compiler in unknown ways.

The system is efficient. On the four-processor machine connected to the MIPS file server, the 45 source files of the operating system compile in about ten seconds of real time and load in another ten. (The loader runs single-threaded.) Partly due to the register-saving convention of the compiler, the null system call takes only seven microseconds on the MIPS, about half of which is attributed to relatively slow memory on the multiprocessor. A process fork takes 700 microseconds irrespective of the process's size.

Plan 9 does not implement lightweight processes explicitly. We are uneasy about deciding where on the continuum from fine-grained hardware-supported parallelism to the usual timesharing notion of a process we should provide support for user multiprocessing. Existing definitions of threads and lightweight processes seem arbitrary and raise more questions than they resolve. We prefer to have a single kind of process and to permit multiple processes to share their address space. With the ability to share local memory and with efficient process creation and switching, both of which are in Plan 9, we can match the functionality of threads without taking a stand on how users should multiprocess.

Process migration is also deliberately absent from Plan 9. Although Plan 9 makes it easy to instantiate processes where they can most effectively run, it does nothing explicit to make this happen. The compiler, for instance, does not arrange that it run on the CPU server. We prefer to do coarse-grained allocation of computing resources simply by running each new command interpreter on a lightly-loaded CPU server. Reasonable management of computing resources renders process migration unnecessary.

Other aspects of the system lead to other efficiencies. A large single-threaded chess database problem runs about four times as fast on Plan 9 as on the same machine running commercial software because the remote cache on the file server is so large. In general, most file I/O is done by direct DMA from the file server's cache; the file server rarely needs to read from disk at all.

Much of Plan 9 is straightforward. The individual pieces that make Plan 9 up

are relatively ordinary; its unusual aspects lie in their combination. As a case in point, the recent interest in using X terminals connected to timeshared hosts might seem to be similar in spirit to how Plan 9 terminals are used, but that is a mistaken impression. The Gnot, although similar in hardware power to a typical X terminal, serves a much higher-level function in the computing environment. It is a fully programmable computer running a virtual memory operating system that maintains its user's view of the entire Plan 9 system. It off loads from the CPU server all the bookkeeping and I/O intensive chores that a window system must perform. It

is not a workstation either; one would rarely bother to compile on the Gnot, although one would certainly run a text editor there. Like the other pieces of Plan 9, the Gnot's strength derives from careful specialization in concert with other specialized components.

Acknowledgments

Special thanks go to Bart Locanthi, who built the Gnot and encouraged us to program it; Tom Duff, who wrote the command interpreter *rc*; Tom Killian and Ted Kowalski, who cheerfully endured early versions of the software; Dennis Ritchie, who frequently provided us with much-needed wisdom; It

The Test of Time

Few Communication Ideas Make It.
Samuel Morse's Did. CommLib Does.

Greenleaf CommLib is the state of the art in asynchronous communication libraries for C. Telecomputing magazine writes, "The detail of CommLib is exceptional." Its interrupt driven functions provide an astonishing level of communications power for your applications.

- Up to 35 simultaneous channels, up to 115,200 baud.
- XMODEM with CRC, Checksum, G and 1-K options.
- YMODEM with G and Batch options.
- Kermit and ASCII file transfer protocols too.
- XON/XOFF and RTS/CTS Flow Controls.
- Circular buffers for Rx and Tx can be any size up to 64K each.
- Intelligent interrupt service for 8250 and 16450 type UART's with optional Receive Filters.
- "WideTrack Rx" that saves status plus data in buffer for each character.
- Drives IBM PC, XT, AT, PS/2 (COM1..COM8) Digiboard, Arnet, Quadram, StarGate, Contec and Qua Tech multiport boards.
- Forty functions for Hayes compatible modem controls.
- Ctrl-Break and exit() traps for safety.
- Run different baud rates and protocols on different ports at the same time!

CALL 1-800-523-9830

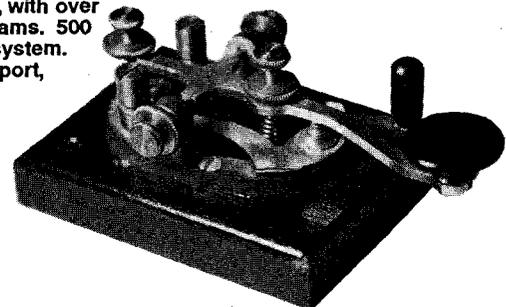
FAX (214) 248-7830

FREE source in C and Assembler. Ask about our Gold Card Support. Available Now for Microsoft C, Quick C, Turbo C, Turbo C++, Lattice C, Zortech C++, TopSpeed C, Watcom C, and JPI TopSpeed C. No Royalties.

All memory models, optimized, with over 170 compilable example programs. 500 page manual and online help system.

FREE unlimited telephone support, FREE BBS, FREE newsletter.

Demo pack available.



CIRCLE NO. 450 ON READER SERVICE CARD

and all those who helped build the system.

References

Accetta, M.J., Robert Baron, William Bolosky, David Golub, Richard Rashid, Avadis Tevanian, and Michael Young. "Mach: A New Kernel Foundation for UNIX Development." In *USENIX Conference Proceedings*. Atlanta, Georgia, 1986.

Duff, T. "Rc — A Shell for Plan 9 and UNIX." In *UNIX Programmer's Manual*. 10th ed. Murray Hill, N.J.: AT&T Bell Laboratories, 1990.

Fraser, A.G. "Datakit — A Modular Network for Synchronous and Asyn-

chronous Traffic." In *Proc. Int. Conf. on Commun.* Boston, Mass., 1980.

Kernighan, Brian W. and Rob Pike. *The UNIX Programming Environment*. Englewood Cliffs, N.J.: Prentice-Hall, 1984.

Killian, T.J. "Processes as Files." In *USENIX Summer Conference Proceedings*. Salt Lake City, Utah, 1984.

Metcalf, R.M. and D.R. Boggs. *The Ethernet Local Network: Three Reports*. Palo Alto, Calif.: Xerox Research Center, 1980.

Miller, S.P., C. Neumann, J.I. Schiller, and J.H. Saltzer. *Kerberos Authentication and Authorization System*. Cambridge, Mass.: MIT Press, 1987.

Pike, R. "Graphics in Overlapping Bitmap Layers." In *Transactions on Graphics*. Vol. 2, No.2, 135-160.

Pike, R. "A Concurrent Window System." In *Computing Systems*. Vol. 2, No. 2, 133-153.

Quinlan, S. "A Cached WORM File System." In *Software—Practice and Experience*. To appear.

Ritchie, D.M. and K. Thompson, "The UNIX Time-Sharing System." In *Comm. Assoc. Comp. Mach.* Vol. 17, 7, 365-375. 1974.

DDJ

Vote for your favorite feature/article. Circle Reader Service **No. 2**.

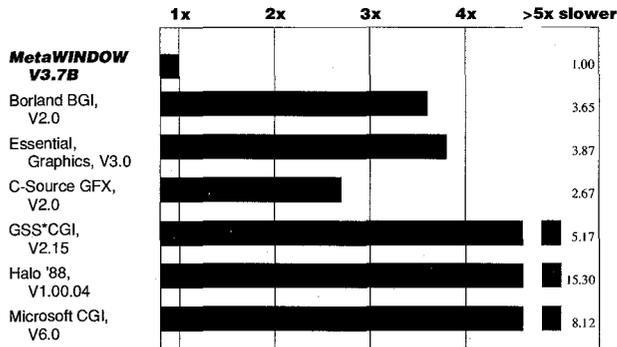
Power Graphics for your PC!

Our apologies. We've created an unfair advantage with the release of the latest version of our graphics development toolkit. We've developed a new version of our award winning MetaWINDOW graphics system that's so blindingly fast and quick - well, we're afraid you won't believe us.

To prove our findings we put together a comprehensive graphics benchmark suite, and then tested MetaWINDOW against every competitive product we could find. The results were horrifying! MetaWINDOW didn't beat the others by merely a factor of 2 or 3, but averaged a whopping 640% performance lead over other graphics tools. Yes, over 6 times faster!! How could we break the news? MetaWINDOW's blistering speed combined with it's already unmatched features is truly an unfair advantage - but prove it yourself!

Obviously, if you're already a Metagraphics customer you're well aware of MetaWINDOW's outstanding capabilities. If you're not using MetaWINDOW, in addition to reinventing a lot of wheels, you may now soon find your competitor's products running over 6 times faster than yours. As we said - our apologies...

Unparalleled Performance!



Benchmarks based on overall line drawing, rectangle fill, oval/circle fill & text bit performance. Timings executed on a standard 6Mhz IBM AT with Video 7 1024i VGA adaptor, 640x480 16-color mode. For further detailed information call, write or fax to request a copy of Metagraphics "white paper" on graphics performance benchmarking. Full benchmark source code downloadable from Metagraphics BBS at 408-438-5368 (1200/2400 baud).

Unmatched Features!

	MW	BGI	EG	GFX	GSS	Halo	MSC
Dynamic recognition & runtime support for over 100+ graphics cards (no clumsy device drivers!)	✓						
Auto-cursor tracking & event queue	✓		*				
Working GUI C source code	✓						
Read/write PCX image files	✓		✓				
Virtual memory/EMS images	✓					✓	
Full clipping	✓		✓		✓		✓
Print graphics screen & virtual memory/EMS images	✓					✓	
Fill/frame rects, rounded-rects, ovals, arcs & polygons	✓						
Full 8-Function RasterOp	✓						
Line join & cap styles	✓						
Rotated/scaled stroked text	✓				✓	✓	
256-color VGA support	✓	✓	✓		✓	✓	✓
Selectable font facing	✓						
Font editor w/extensive fonts	*			*			
No Royalties!	✓	✓	✓	✓			✓
Awarded industry honors for product excellence	✓						

* - Additional cost option.

MetaWINDOW (supports 15 popular C, C++ & Pascal compilers) ... \$ 295
 (single language for Turbo C, Turbo Pascal or MS Quick C)..... \$ 150
 FontWINDOW (font editor plus over 2 megabytes of fonts!)..... \$ 150

30 Day No Risk Money-Back Guarantee!

Order Hotline 800-332-1550



METAGRAPHICS
SOFTWARE CORPORATION

269 Mount Hermon Road, Scotts Valley, CA 95066
 Phone: 408-438-1550 Fax: 408-438-5379

CIRCLE NO. 156 ON READER SERVICE CARD