

MobiDesk: Mobile Virtual Desktop Computing

Ricardo A. Baratto, Shaya Potter, Gong Su, and Jason Nieh

Department of Computer Science
Columbia University, New York, NY, USA

{ricardo, spotter, gongsu, nieh}@cs.columbia.edu

ABSTRACT

We present MobiDesk, a mobile virtual desktop computing hosting infrastructure that leverages continued improvements in network speed, cost, and ubiquity to address the complexity, cost, and mobility limitations of today's personal computing infrastructure. MobiDesk transparently virtualizes a user's computing session by abstracting underlying system resources in three key areas: display, operating system, and network. It provides a thin virtualization layer that decouples a user's computing session from any particular end-user device, and moves all application logic to hosting providers. The virtualization layer decouples a user's computing session from the underlying operating system and server instance, enabling high-availability service by transparently migrating sessions from one server to another during server maintenance or upgrades. We have implemented a prototype in Linux that works with existing unmodified applications and operating system kernels. Our experimental results demonstrate that MobiDesk has very low virtualization overhead, can provide a full featured desktop experience including full-motion video support, and is able to migrate users' sessions efficiently and reliably for high-availability, while maintaining existing network connections.

Categories and Subject Descriptors

C.2.4 [Computer-Communication-Networks]: Distributed Systems—*client/server*; D.4.7 [Software-Operating Systems-Organization and Design]: Distributed Systems

General Terms

Design, Experimentation, Measurement, Performance

Keywords

computer utility, virtualization, process migration, thin-client computing, network mobility, on-demand computing

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

MobiCom'04, Sept. 26-Oct. 1, 2004, Philadelphia, Pennsylvania, USA.
Copyright 2004 ACM 1-58113-868-7/04/0009 ...\$5.00.

1. INTRODUCTION

Continuing advances in hardware technology have enabled the proliferation of faster, cheaper, and more portable personal computers to support increasingly mobile users. As personal computers become more ubiquitous in large corporate, government, and academic organizations, the total cost of owning and maintaining them is becoming unmanageable. These computers are increasingly networked which only complicates the management problem. They need to be constantly patched and upgraded to protect them, and their data, from the myriad of viruses and other attacks commonplace in today's networks. Furthermore, as mobile users transport their portable computers from one place to another, it is not uncommon for these machines to be damaged or stolen, resulting in the loss of any important data stored on them. Even in the best case, when such data can be recovered from backup, the time consuming process of reconstituting the state of the lost machine on another device, results in a huge disruption in critical computing service for the user.

We introduce MobiDesk, a mobile virtual desktop computing hosting infrastructure. With wire-speed network technologies scaling at faster Moore's exponents than silicon, MobiDesk leverages rapid improvements in network bandwidth, cost, and ubiquity to address the limitations of the current personal computing model. MobiDesk uses the network to decouple a user's desktop computing session from the end-user device by moving all application logic to hosting providers. In this manner, end-user devices are simply used to transmit user input and display application output, allowing them to be simple stateless clients. MobiDesk also decouples a user's desktop computing session from the underlying operating system and server instance, allowing a user's entire computing environment to be migrated transparently from one server to another. This enables a server to be brought down for maintenance and upgraded in a timely manner with minimal impact on the availability of a user's computing services. Once the original machine has been updated, the user's computing session can be migrated back and continue to execute even though the underlying operating system may have changed. MobiDesk ensures that any network connections associated with the user's computing session are maintained, even as the session is migrated from one machine to another. MobiDesk provides these benefits without modifying, recompiling, or relinking applications or operating system kernels. MobiDesk requires no changes to clients other than being able to execute a simple user-space application to process and display input and output.

MobiDesk provides a mobile virtual desktop computing environment by introducing a thin virtualization layer between a user's computing environment and the underlying system. MobiDesk focuses on virtualizing three key system resources: display, operating system, and network. MobiDesk virtualizes display resources by providing a virtual display driver that efficiently encodes and redirects display updates from the server to an end-user device. MobiDesk virtualizes operating system resources by providing a virtual private namespace for each desktop computing session. The namespace offers a host independent virtualized view of an operating system, enabling the session to be transparently migrated from one server to another. MobiDesk virtualizes network resources by providing virtual address identifiers for connections, and a transport-independent proxy mechanism. Together, they preserve all network connections associated with a user's computing session, even if it is migrated from one server to another inside the MobiDesk server infrastructure.

The MobiDesk hosted desktop computing approach provides a number of important benefits over current computing approaches:

- *High-availability and reliable application services:* Because MobiDesk is designed to work with unmodified legacy applications and commodity operating systems, it offers the potential to bring about more reliable computing without giving up the large investments already made in the existing software base. Furthermore, decoupling from the underlying hardware and operating system allows applications to be moved anywhere, and in particular, migrated off faulty hosts, and before maintenance and upgrades. In contrast to today's long periods of service downtime due to maintenance and upgrades, MobiDesk enables hardware and operating systems to be upgraded in a timely manner with minimal impact on application service availability — by migrating applications to another machine that has already been updated. With MobiDesk, system administrators no longer need to schedule downtime in advance and in cooperation with all the users, thereby closing the vulnerability window of unrepaired systems.
- *Persistence and continuity of business logic:* MobiDesk moves away from the current model of simply backing up file data to secure remote locations, and instead protects entire computing environments by running hosting providers in secure remote locations. This enables academic, business, and government institutions to function much more effectively in times of crisis. Restoring an organization's local computing infrastructure from backup consequent to a crisis is an extremely slow, time consuming process that is increasingly ineffective given the scale of IT infrastructure being deployed today. MobiDesk offers a different, improved model of continuous uptime, especially during a crisis, when infrastructure availability is most crucial.
- *Secure, low-cost global access and transparent user mobility:* MobiDesk client access devices just need to be able to connect to the Internet. They do not need to provide complex computing functionality, making

it unnecessary to continuously upgrade to more powerful desktop machines. Simpler, lower-cost, possibly longer battery life client access devices can be made more readily available for such a service. These devices may come in many shapes and sizes, from desktop machines with megapixel displays to handheld devices with pocket sized screens. Furthermore, because all persistent user state is maintained on the servers, users are able to securely access, and freely move among any client access devices and pick up right where they left off.

- *On-demand access to application and computational resources:* By multiplexing a large pool of shared resources among many users, an individual can gain access to substantially more applications and resources than can be afforded on one's local desktop computer. In terms of applications, MobiDesk can provide a wider range of affordable application services on multiple operating system platforms by amortizing costs over a large number of users. Since not all applications will be in use by all users at one time, statistical multiplexing can serve a larger number of users with fewer software licenses. In terms of resources, a user can be given resource allocations which can be scaled up or down as necessary. Instead of having to throw away their existing local desktop machines every time they need more compute power, users can just ask their service provider to scale up their allocation.
- *Bridging the information gap:* Contrasting with today's environment of information "haves" and "have-nots", MobiDesk can be accessed by low-cost client devices to deliver secure service on a subscription basis, offering a societally better way of providing access to computing as widely as possible. Remote computing and storage services can be paid-for using a subscription-based or use-based economic model that follows current Internet access or telephone pricing policies. A MobiDesk provider could supply basic computing, including a complete suite of desktop productivity tools, for just a few dollars a month. Video-on-demand service, online gaming, and other multimedia services could be obtained from the same or other providers for additional small fees. Furthermore, the more centralized service and management model provides lower total cost of ownership across all applications, making computing more widely affordable. More importantly, MobiDesk enables client devices to be easier to manage and use by removing the local system complexity associated with current software upgrade cycles, making computing more accessible to a wider population of users.

This paper presents the design and implementation of MobiDesk. Section 2 presents the overall MobiDesk system architecture and usage model. Section 3 describes the MobiDesk display virtualization mechanism. Section 4 describes the MobiDesk operating system virtualization mechanism. Section 5 describes the MobiDesk network virtualization mechanism. Section 6 presents experimental results measuring MobiDesk system performance and associated overhead in the context of real desktop computing applications. Section 7 discusses related work. Finally, we present some concluding remarks.

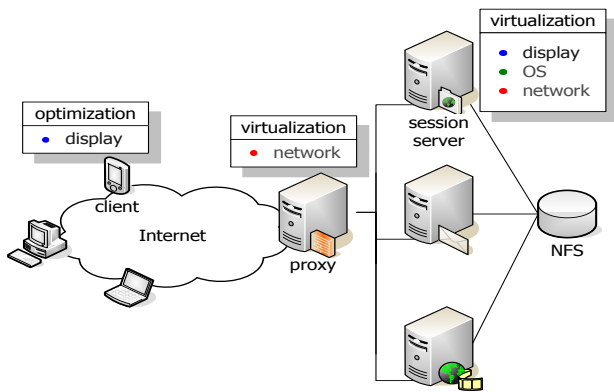


Figure 1: MobiDesk Architecture

2. OVERVIEW OF SYSTEM ARCHITECTURE

MobiDesk is architected as a proxy-based server cluster system, comparable to systems deployed today by application service providers. The overall architecture of the system is depicted in Figure 1. MobiDesk is composed of a proxy, a group of back-end session servers connected in a LAN, a storage server infrastructure, and a number of external, heterogeneous clients through which users access the system.

The proxy acts as a front-end that admits service requests from clients across the Internet, and dispatches the requests to the appropriate back-end application servers. The proxy, operating at layer 7, exposes a single entry point to the clients, and employs suitable admission and service dispatching policies. The back-end compute servers host completely virtualized environments within which the computing sessions of MobiDesk’s users run. The network storage server infrastructure is used for all persistent file storage. The clients are merely *inputting* and *outputting* devices connected to the servers across the Internet.

Users interact with their MobiDesk sessions through a thin-client *session viewer*, a simple device or application that relays the user’s input and the session’s output between the client and the server through a secure channel. Each user in the system is assigned a username and password. Upon the first login, the proxy performs appropriate authentication, and connects the user to a MobiDesk session server. The session server creates a virtual private environment that is populated with a complete set of operating system resources and desktop applications. In contrast to the traditional centralized computing model where users are aware of each others’ presence and activities, MobiDesk’s sessions are isolated from one another and the underlying server environment. To the user, the session appears no different than a private computer, even though the user’s session may coexist with many other sessions on a shared server. When the client disconnects, the session continues to run on the MobiDesk server, unless the user explicitly logs out. On future connections, the session will be in the same state it was when the user last disconnected.

By providing a virtual private environment for each user, MobiDesk is able to dynamically relocate sessions to meet load balancing, system maintenance and/or quality of service requirements. Sessions can be checkpointed and migrated transparently at any point in time. To keep track of the sessions as migration occurs, MobiDesk implements

a *session cookie* mechanism. As new sessions are created, the proxy generates a unique cookie that is passed to the hosting servers and associated with the new session. Whenever a session is migrated, the destination server uses the cookie to inform the proxy of the new location. Finally, the next time the user logs in, the proxy will use the cookie to identify the server where the user’s session is being hosted.

To provide a private and mobile environment for user sessions, MobiDesk virtualizes three key resources: display, operating system, and network. MobiDesk virtualization is designed to work with existing unmodified applications, operating system kernels, and network infrastructure and protocols. The three components work in concert to create a completely virtualized environment for client computing sessions.

MobiDesk virtualizes the server display by providing a virtual display driver. The display driver intercepts drawing commands from user’s applications, and translates the commands into a display protocol between the client and the server. MobiDesk display virtualization focuses on the importance of a latency sensitive design for interactive computing. It also anticipates the kind of high bandwidth network access that is becoming increasingly cost effective and accessible in WAN environments [1]. For example, South Korea is building a nationwide Internet access infrastructure to make speeds up to 100 Mbps available to the home by 2010 [24].

MobiDesk operating system virtualization provides a virtual private namespace for each hosted client computing session. For example, each computing session contains its own host independent view of OS resources, such as PID/GID, IPC, memory, file system, and devices. MobiDesk virtualization operates at a finer granularity than virtual machine approaches such as VMware [52] by virtualizing individual computing sessions instead of complete operating system environments. As a result, computing sessions can be decoupled from the underlying operating system and migrated to other servers. This enables high-availability even in the presence of server hardware and operating system maintenance and upgrades.

MobiDesk network virtualization provides persistent network connections for client computing sessions even as they move among servers in a MobiDesk cluster. All connections operate through the MobiDesk proxy. Similar to operating system virtualization, MobiDesk virtualizes network connections with a virtual private namespace for transport connection identifiers, such as IP address and port number, on both the proxy and the servers. These virtual identifiers remain constant and are simply translated to the underlying physical network identifiers as a session moves amongst servers at different physical network locations. MobiDesk network virtualization provides persistent connections from mobile client computing sessions to outside hosts without running any software on the outside hosts and without any changes to the existing network infrastructure.

3. DISPLAY VIRTUALIZATION

To make MobiDesk a viable replacement to the traditional desktop computing model, it needs to be able to deliver the look and feel of all unmodified desktop applications end-users expect. MobiDesk must work within the framework of existing display systems, intercepting display commands from unmodified applications and redirecting these

Command	Description
RAW	Display raw pixel data at a given location and size
COPY	Copy frame buffer area to specified coordinates
SFILL	Fill an area with a given pixel color value
BITMAP	Fill a region using a bitmap image
PFILL	Tile a pixmap rectangle in a given region

Table 1: MobiDesk Protocol Display Commands

commands to remote clients. To provide good WAN performance, the virtualization must intercept display commands at an appropriate abstraction layer to provide sufficient information to optimize the processing of display commands in a latency sensitive manner. Furthermore, to support transparent user mobility and eliminate client administration complexity, MobiDesk should support the use of thin, stateless clients, by ensuring that all persistent display state is stored in the server infrastructure.

Given that traditional display systems are structured in multiple abstraction layers, there are a number of possible ways in which MobiDesk can interact with existing display systems. We can loosely categorize display system structure into three layers: application, middleware, and hardware. The application layer is the top layer with which applications interact. It presents a high-level model of the overall characteristics of the display system. These characteristics include descriptions of the operation and management of windows, input mechanisms, and display capabilities of the system. These capabilities may range from basic 2D display, to complex operations involving transparency, blends, 3D transformations, and the display of multimedia content. The middleware layer sits between the high-level application layer and the low-level video hardware layer. It is responsible for creating a hardware-independent abstraction of the display hardware to meet the requirements of the display system and its applications. To maintain consistency across hardware with differing abilities, the middleware layer is provisioned with fallback mechanisms and software routines that can implement missing hardware features. The video hardware layer is a low-level, hardware-dependent layer that exposes the video hardware to the display system. It is implemented as a set of device drivers responsible for translating between the middleware’s abstract display operations and the commands understood by the display hardware.

MobiDesk does not intercept at the application layer because that requires a significant amount of application logic and computational power on the client for translating high-level commands. This in turn, would limit the range of target client architectures. Intercepting at the application layer also results in direct synchronization between applications and the client, which can reduce display performance in higher latency WAN environments. MobiDesk also does not intercept at the middleware layer because that would require it to reimplement substantial display system functionality, instead of leveraging continuing advances in existing middleware implementations, such as XFree86.

MobiDesk display virtualization is designed instead as a virtual video device driver that intercepts display commands

at the video hardware layer, by providing a separate virtual video device for each computing session. Rather than sending display commands to local display hardware, the virtual video driver packages up display commands associated with a user’s computing session and sends them over the network to a remote client. For this purpose, MobiDesk implements a simple, low-level, minimum-overhead protocol, as described in Table 1. The protocol mimics the operations most commonly found in display hardware, allowing clients to do little more than forward protocol commands to their local video hardware, thus reducing the latency of display processing. To provide security, all protocol traffic is encrypted using the standard RC4 [44] stream cipher algorithm. MobiDesk’s video hardware layer approach allows it to take full advantage of existing infrastructure and hardware interfaces, while maximizing client resources and requiring minimal computation on the client. Furthermore, new video hardware features can be supported with at most the same amount of work necessary for supporting them in traditional desktop display drivers. While there is some loss of semantic display information at the low-level video device driver interface, our experiments with desktop applications such as web browsers, indicate that the vast majority of application display commands issued can be mapped directly to standard video hardware primitives. Furthermore, we show in Section 6 that the simpler display virtualization used by MobiDesk can provide superior display performance compared to other approaches.

To deliver good performance in WAN environments, the MobiDesk display virtualization architecture couples its virtual video device driver approach with other latency sensitive display mechanisms. In particular, client display hardware resources are exported to the server and leveraged to reduce the latency of display processing. For example, MobiDesk provides direct video support by leveraging alternative YUV video formats natively supported by almost all off-the-shelf video cards available today. Video data can simply be transferred from the server to the client video hardware, which automatically does inexpensive, high speed, color space conversion and scaling. As another example, MobiDesk leverages cursor drawing support available in commodity video cards in use today. MobiDesk uses local cursor drawing in response to mouse movements and maintains local cursor state at the client. Cursor changes still come from the server, but system response time is improved by avoiding network latency with any cursor drawing that does not result in cursor changes. As a final example, MobiDesk transparently supports varying client screen sizes by automatically resizing updates to fit within the client viewport. For mobile devices with small screens, this translates in reduced bandwidth consumption and improved performance.

MobiDesk provides two important server-side mechanisms for improving performance in WAN environments. The first mechanism is the use of a server push model for sending display updates to the client. As soon as display updates are generated on the server, they are delivered to the client. Clients are not required to explicitly request display updates, which add additional network latency to command processing. The second mechanism is the use of display command scheduling to improve the responsiveness of the system. MobiDesk employs a *Shortest-Remaining-Size-First (SRSF)* preemptive command scheduler. SRSF is analogous to Shortest-Remaining-Processing-Time scheduling, which

is known to be optimal for minimizing mean response time in an interactive system. In display applications, short jobs are normally associated with text and general GUI layout components, which are critical to the usability of the system. On the other hand, large jobs are normally lower priority “beautifying” GUI elements, such as image decorations, desktop backgrounds and web page banners.

Finally, MobiDesk supports thin, stateless display clients by storing all session state at the respective session server. Although MobiDesk takes advantage of client resources when available, all client state is considered temporary and destroyed upon client disconnect. When a remote client connects to the MobiDesk infrastructure, the server running the user’s computing session transfers the current session state to the client. For the duration of the connection, the client forwards input events to the server, which in turn forwards display updates back to the client. The client at no point has an intermediate session state differing from the server. Furthermore, if allowed by the user, multiple clients can be connected to the same session at the same time, all of them accessing the same centralized view distributed from the server. When a client eventually disconnects, it leaves no state behind in the local computer.

4. OPERATING SYSTEM VIRTUALIZATION

MobiDesk encapsulates user sessions within a host independent, virtualized view of the operating system. Unlike traditional operating systems, each session is a self contained unit that can be isolated from the system, checkpointed to secondary storage, migrated to another machine, and transparently restarted. This is made possible because each computing session has its own virtual private namespace, which provides the only means for processes to access the underlying operating system. To guarantee correct operation of unmodified applications, MobiDesk session virtualization is completely transparent. This is accomplished by providing a traditional environment with unchanged application interfaces and access to operating system services and resources.

The namespace is private in that only processes within the session can see the namespace, and the namespace in turn masks out resources that are not contained in the session. Processes inside the session appear to one another as normal processes, and they are able to communicate using traditional IPC mechanisms. On the other hand, no IPC interaction is possible across the session’s boundaries, because outside processes are not part of the private namespace. Processes inside a session and those outside of it are only able to communicate over RPC mechanisms, traditionally used to communicate across computers.

The namespace is virtual in that all operating system resources, including processes, user information, files, and devices, are accessed through virtual identifiers. These virtual identifiers are distinct from the host-dependent, physical resource identifiers used by the operating system. The session’s namespace uses the virtual identifiers to provide a host-independent view of the system, which remains consistent throughout a process’s and session’s lifetime. Since the session’s namespace is separate from the underlying namespace, it can preserve naming consistency for its processes, even if the physical namespace changes, as may be the case when sessions are migrated across computers.

Operating system resource identifiers, such as process IDs (PIDs), must remain constant throughout the life of a pro-

cess to ensure its correct operation. However, when a process is moved from one operating system instance to another, there is no guarantee that the destination system will provide the same identifiers to the migrated process; those identifiers may in fact be in use by other processes in the system. The session’s namespace addresses these issues by providing consistent, virtual resource names. Names within a session are trivially assigned in a unique manner in the same way that traditional operating systems assign names, but such names are localized to the session. Since the namespace is virtual, there is no need for it to change when the session is migrated, ensuring that identifiers remain constant throughout the life of the process, as required by applications that use such identifiers. Since the namespace is private to a given session, processes within the session can be migrated as a group, while avoiding resource naming conflicts among processes in different sessions. Finally, the private virtual namespace enables sessions to be securely isolated from each other by providing complete mediation to all operating system resources. Since the only resources within each session are the ones that are accessible to the owner of the session, a compromised session would be unable to harm any other user’s session.

4.1 Session Virtualization

MobiDesk virtualizes sessions by using mechanisms that translate between the session’s virtual resource identifiers and the operating system resource identifiers. For every resource accessed by a process in a session, the virtualization layer associates a *virtual name* to an appropriate operating system *physical name*. When an operating system resource is created for a process in a session, the physical name returned by the system is caught, and a corresponding private virtual name created and returned to the process. Similarly, any time a process passes a virtual name to the operating system, the virtualization layer catches and replaces it with the corresponding physical name. The key virtualization mechanisms used are a system call interposition mechanism and the `chroot` utility with file system stacking for file system resources.

Session virtualization uses system call interposition to virtualize operating system resources, including process identifiers, keys and identifiers for IPC mechanisms such as semaphores, shared memory, and message queues, and network addresses. System call interposition wraps existing system calls to check and replace arguments that take virtual names with the corresponding physical names, before calling the original system call. Similarly, wrappers are used to capture physical name identifiers that the original system calls return, and return corresponding virtual names to the calling process running inside the session. Session virtual names are maintained consistently as a session migrates from one machine to another and are remapped appropriately to underlying physical names that may change as a result of migration. Session system call interposition also masks out processes inside of a session from processes outside of the session to prevent any interprocess host dependencies across the session boundary.

Session virtualization employs the `chroot` utility and file systems stacking to provide each session with its own file system namespace. In the simplest case, where sessions are only checkpointed and restarted on the same computer, the file system can be composed of loopback mounts from the host

computer. In MobiDesk's centralized hosting infrastructure, the session's file system is composed from remote mounts via a network file system such as NFS, which guarantees that the same files can be made consistently available as a session is migrated from one computer to another. More specifically, when a session is created or moved to a host, a private directory is created in the host. This directory, named according to the session identifier, serves as a staging area for the session's virtual file system. Within the directory, the various network accessible directories that the session is configured to access, will be mounted from a network file server. For example, from a Unix-centric viewpoint, this set of directories could include `/etc`, `/lib`, `/bin`, `/usr`, and `/tmp`. The `chroot` system call is then used to set the staging area as the root directory for the session, thereby achieving file system virtualization with negligible performance overhead. This method of file system virtualization provides an easy way to restrict access to files and devices from within a session. This can be done by simply not including file hierarchies and devices within the session's file system namespace. If files and devices are not mounted within the session's virtual file system, they are not accessible to the session's processes.

An alternative MobiDesk usage model to enable disconnected operation without continued access to file servers would be to leverage the use of increasingly popular portable storage devices. In this scenario, the user checkpoints and saves the session and its required file system to a portable device, effectively creating a *pocket-session*, a complete computer environment which can be carried around in a user's pocket. To migrate to another computer, the user only needs to plug the portable device to the destination computer, and restart the session within the file system stored in the device. This functionality may be useful to users running sessions without the need for sensitive data, that plan to move from a high network connectivity location to a low network connectivity location.

Commodity operating systems are not built to support multiple namespaces securely. Therefore, session virtualization must address the fact that there are multiple ways to break out of a chrooted environment, especially when the `chroot` system call is allowed to be used in a session. The session's file system virtualization enforces the chrooted environment and ensures that the session's file system is only accessible to processes within the given session by using a simple form of file system stacking to implement a barrier. This barrier directory prevents processes within the session from traversing it. Since the processes are not allowed to traverse the directory, they are unable to access files outside of the session's file system namespace.

In order for a session to fully replace a regular computer, the session has to allow processes to run as the privileged root user. For instance, programs such as `ping` and `traceroute` that create raw sockets, and `passwd` that is used to change system resources. Because root's UID 0 is treated specially by the operating system kernel, session virtualization also treats UID 0 processes inside of a session specially, to prevent them from breaking the session abstraction. While a session can be configured for administrative reasons to allow full privileged access to the underlying system, we focus on the case of sessions for running application services which do not need to be used in this manner.

While a process is running in user space, the UID it runs as does not have any effect. The UID only has an impact when

the process tries to access the underlying kernel via one of the entry points, specifically devices and system calls. Since a session already provides a virtual file system that includes a virtual `/dev` directory with a limited set of secure devices, the device entry point is already secured. The only system calls of concern are those that could allow a process owned by root to break the session abstraction. However, only a small number of system calls can be used for this purpose. Session virtualization classifies these system calls into three classes which need to be protected.

The first class of system calls are those that only affect the host system and serve no purpose within a session. Examples of these system calls include those that load and unload kernel modules or that reboot the host system. Since these system calls only affect the host, they would break the session security abstraction by allowing processes within it to make system administrative changes to the host. System calls that are part of this class are therefore made inaccessible by default to processes running within a session.

The second class of system calls are those that are forced to run unprivileged. Just like NFS, by default, "squashes" root on a client machine to act as user `nobody`, session virtualization forces privileged processes to act as the `nobody` user when it wants to make use of these system calls. These system calls include those that set resource limits, and `ioctl` system calls. Since system calls such as `setrlimit` and `nice` can allow a privileged process to increase its resource limits beyond predefined limits imposed on session processes, privileged processes are by default treated as unprivileged when executing these system calls within a session. Similarly, the `ioctl` system call is a system call multiplexer that allows any driver on the host to effectively install its own set of system calls. Since the ability to audit the large set of possible system calls is impossible given that sessions may be deployed on a wide range of machine configurations, session virtualization conservatively treats access to this system call as unprivileged by default.

The final class of system calls are calls that are required for regular applications to run, but have options that will give the processes access to underlying host resources, breaking the session abstraction. Since these system calls are required by applications, the session checks all their options to ensure that they are limited to resources that the session has access to, making sure they are not used in a manner that breaks the session abstraction. For example, the `mknod` system call can be used by privileged processes to make named pipes or files. It is therefore desirable to make it available for use within a session. However, it can also be used to create device nodes that provide access to the underlying host resources. To limit how the system call is used, the session system call interposition mechanism checks the options of the system call and only allows it to continue if it is not creating a device.

4.2 Session Migration

MobiDesk provides the ability to maintain session availability in the presence of server downtime due to operating system and hardware upgrades. This is accomplished through a checkpoint-restart mechanism that allows sessions to be migrated between computers with different hardware and operating system kernels. MobiDesk is limited to migrating between machines with a common CPU architecture, and where kernel differences are limited to maintenance and

security patches. These patches often correspond to changes in minor version numbers of the kernel. In particular, the Linux 2.4 kernel has more than 25 minor versions. Migration is limited to these instances because major version changes are allowed to break application compatibility, which may cause running processes to break. However, even with minor versions changes, there can be significant changes in kernel code. For example, during the Linux 2.4 series of kernels, the entire VM subsystem was extensively modified to change the page replacement mechanism. Similarly, migration is limited to scenarios where the application's execution semantics, such as how threads are implemented or dynamic linking is performed, stay constant. On the Linux kernel, this is not an issue as these semantics are enforced by user-space libraries. Since the session's user-space libraries migrate with it, the semantics stay constant. While MobiDesk sessions protect the host from malicious code executed within a session, migrating away from or to a host that has already been compromised is beyond the scope of this paper.

To support migration across different kernels, MobiDesk's checkpoint-restart mechanism employs an intermediate format to represent the state that needs to be saved. On checkpoint, the process image is saved and digitally signed to enable the restart process to verify its integrity. Although the internal state that the kernel maintains on behalf of processes can be different across kernels, the high-level properties of the process are much less likely to change. MobiDesk captures the state of a process in terms of this higher-level semantic information rather than the kernel specific data. For example, part of the state associated with a Unix socket connection consists of the directory entry of the socket, its superblock information, and a hash key. It may be possible to save all of this state in this form and successfully restore on a different machine running the same kernel. But this representation is of limited portability across different kernels. On the other hand, a high-level representation consisting of a four tuple: {virtual source pid, source fd, virtual destination pid, destination fd}, is highly portable. This is because the semantics of a process identifier and a file descriptor are standard across different kernels.

MobiDesk's intermediate representation format is chosen such that it offers the degree of portability needed for migrating between different kernel minor versions. If the representation of state is too high-level, the checkpoint-restart mechanism could become complicated and impose additional overhead. For example, the MobiDesk system saves the address space of a process in terms of discrete memory regions called VM areas. As an alternative, it may be possible to save the contents of a process's address space and denote the characteristics of various portions of it in more abstract terms. However, this would call for an unnecessarily complicated interpretation scheme and make the implementation inefficient. The VM area abstraction is standard even across major Linux kernel revisions. MobiDesk views the VM area abstraction as offering sufficient portability in part because the organization of a process's address space in this manner has been standard across all Linux kernels and has never changed since its inception.

MobiDesk leverages high-level native kernel services in order to transform the intermediate representation of the checkpointed image into the complete internal state required by the target kernel. Continuing with the previous example,

MobiDesk restores a Unix socket connection using high-level kernel functions as follows. First, two new processes are created with virtual PIDs as specified in the four tuple. Then, each one creates a Unix socket with the specified file descriptor and one socket is made to connect to the other. This procedure effectively recreates the original Unix socket connection without depending on many internal kernel details.

This use of high-level functions helps with general portability when using MobiDesk for migration. Security patches and minor version kernel revisions commonly involve modifying the internal details of the kernel while high-level primitives remain unchanged. As such high-level functions are usually made available to kernel modules through exported kernel symbol interface, the MobiDesk system is able to perform cross-kernel migration without requiring modifications to the kernel.

To eliminate possible dependencies on low-level kernel details, MobiDesk's checkpoint-restart mechanism requires processes to be suspended prior to being checkpointed. Suspending processes creates a quiescent state necessary to guarantee the correctness of the checkpointed image, and it also minimizes the amount of information that needs to be saved. As a representative example, consider the case of semaphore wait queues. Although semaphore values can be easily obtained and restored through well known interfaces, saving and restoring the state of the wait queue involves the manipulation of kernel internals. However, by taking advantage of existing semantics which direct the kernel to release a process from a wait queue upon receipt of a signal, MobiDesk is able to empty the wait queues by suspending all processes, and therefore avoid having to save the state of the queue.

Finally, we must ensure that any changes in the system call interfaces are properly accounted for. As MobiDesk has a virtualization layer that uses system call interposition to maintain namespace consistency, a change in the semantics for any system call intercepted could be an issue in migrating across different kernel versions. But such changes usually do not occur, as it would require system libraries to be rewritten. In other words, MobiDesk virtualization is protected from such changes in the same way legacy applications are protected. However, new system calls could be added from time to time. Such system calls could have implications to the encapsulation mechanism. For instance, across all Linux 2.4 kernels, there were two new system calls that used identifiers that needed to be intercepted and virtualized, `gettid` and `tkill`.

5. NETWORK VIRTUALIZATION

Networking support for MobiDesk sessions must address two issues:

- Multiple sessions on the same server may run the same service, e.g. two sessions may both run the apache server, however, only one of them can listen on port 80.
- Ongoing network connections must be preserved when a session is migrated from one server to another.

When all hosting servers are in the same subnet, the two issues can be addressed relatively easily using existing technologies with minor enhancements from MobiDesk. Each session is assigned a unique IP address from a pool maintained by a DHCP server when it is first created. For exam-

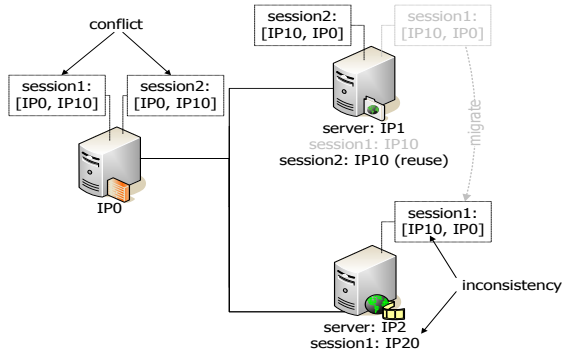


Figure 2: Key Problems of Connection Migration

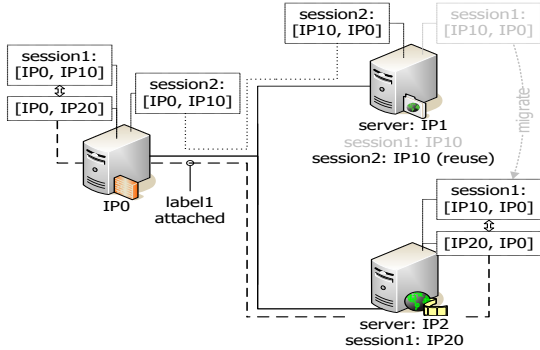


Figure 3: MobiDesk Network Virtualization

ple, the servers may occupy IP address range 192.168.1.2 - 192.168.1.50, and the rest of 192.168.1.5 - 192.168.1.254 may be assigned to MobiDesk sessions. The IP address assigned to a session is created as an alias of the hosting server’s primary IP address. Multiple aliases, each corresponding to a different session, can be created on a server. MobiDesk privatizes the aliases such that a session only sees its own alias, and cannot interfere with traffic of other sessions on the same server.

Since each session has its own IP address, two sessions on the same server can both listen to port 80, bound to their individual private IP address. When a session is migrated from one server to another, the private IP address of the session remains unchanged; it is simply (re)created as an alias of the new hosting servers primary IP address. ARP resolves the MAC address change at the link layer and the migration is transparent to the network layer and above. Ongoing network connections of the session therefore stay intact. Note that ARP caching may cause the migration to be “invisible” to the proxy until its ARP cache times out, which is typically a few minutes. This can be easily addressed by having the migrated session send a “gratuitous” ARP request asking the MAC address of its own IP, which allows the proxy to immediately invalidate its ARP cache for the IP in question. “Gratuitous” ARP requests are commonly used to detect duplicate IP address, and to allow a backup server’s NIC to take over a primary server’s NIC.

While it is possible to have the entire private network behind the proxy to be in a single subnet, it is often desirable to have separate subnets for scalability and management reasons. In this case, when a session is migrated across subnets, its private IP address can no longer persist, since on the destination subnet the address is no longer valid. As

a result, two types of problems can occur, as we illustrate in Figure 2. Note that we omit port numbers for simplicity.

We see that when session1 with IP10 migrates from server IP1 to IP2, its transport connection [IP10, IP0] must persist. However, its IP address IP10 cannot persist because IP2 is on a different subnet. In addition, after session1 with IP10 migrates to server IP2, another session2 may reuse IP10 on server IP1 (or another server) and create another connection [IP10, IP0]; a conflict is created since the proxy will see two identical connections [IP0, IP10]. A potential solution is to use MobileIP. However, MobileIP requires assigning each session a permanent home address that cannot be reused by other sessions (to avoid conflicts as described later in this section). This is difficult since MobiDesk sessions are dynamically created, volatile entities. One can potentially adopt a solution that takes the initial physical address assigned to a session as its home address. However, this still requires additional management infrastructure to (1) assign dynamic address on a per session basis rather than per host basis, and (2) guarantee that the dynamically assigned home address is never reused by any other sessions, even after it has migrated away from its initial subnet. To effectively address these problems without incurring additional management complexity, we have designed an integrated solution of our own.

To address the inconsistency problem on the MobiDesk server, MobiDesk associates each session with two IP addresses, one is a virtual address exposed to the transport layer and above and the other is a physical address seen only at the network layer and below. The virtual address stays constant for the lifetime of the session while the physical address changes whenever the session migrates. The physical address is obtained from a DHCP server, as in the case of a single subnet, but must change when the session is migrated across subnets. There are two ways to assign the virtual address. By default, the virtual address is equal to the initial physical address when the session is created and stays unchanged thereafter. Alternatively, the virtual address can be a predefined value that is constant across all sessions. MobiDesk translates the virtual address to the current physical address (and vice versa) when the session migrates. For example, in Figure 3, after migration, session1’s virtual address IP10 is unchanged while its physical address is assigned by the DHCP server to be IP20 and created as an alias on server IP2. The proxy translates [IP0, IP10] into [IP0, IP20] while the server IP2 translates [IP20, IP0] into [IP10, IP0]. Since the virtual address never changes, the migration is transparent to the transport and above layers.

One potential solution to the conflict problem on the MobiDesk proxy is to require that a physical address, once assigned to a session, is never reused until the session finishes, even after the session has migrated to another subnet. However, this results in undesirable dependency of a session on a trail of addresses if it is migrated many times and new connections are opened between each migration. MobiDesk’s solution is to privatize virtual addresses, i.e., to associate virtual addresses with separate private virtual network interfaces which provide a per-connection address namespace. Instead of having all connections share the same physical interface, each connection is assigned its own private virtual network interface card (VNIC). A VNIC is simply a software emulation of a NIC at the link layer that appears exactly the same as a NIC to the network and above layers. As a re-

sult, two connections using the same virtual IP address due to address reuse can peacefully coexist on the same server, since they are bound to their own private VNIC.

To support per-connection address space, MobiDesk augments the traditional connection tuple with connection labels to identify the VNIC to which a connection is bound. A connection has two labels, independently and uniquely chosen by the MobiDesk proxy and the server at the time when the connection is setup. The two sides also exchange their labels at connection setup time. Before a session is migrated, the labels are not used since the tuple alone is enough to identify the connections of the session. After a session is migrated, both sides will attach its peer’s label learned at connection setup time for all connections between them. The labels allow the connections to be uniquely identified even when a session’s previous physical address is reused. We illustrate these ideas using the example in Figure 3.

In Figure 3, when the connection [IP0, IP10] was setup for session1 while it was on server IP1, the proxy creates a VNIC for the connection and sends its label1 for the connection to server IP1. Similarly, session1 on server IP1 sends its label for the connection to the proxy (the actual label exchange is not shown). Before session1 was migrated, its virtual address, which equals its physical address IP10, cannot be reused by other sessions; therefore the virtual address alone is enough to identify the connection [IP0, IP10] and labels are not used in the absence of migration. However, after session1 is migrated to server IP2, all packets from session1 to the proxy will have label1 attached to them. This allows the proxy to uniquely identify the connection [IP0, IP10] that belongs to session1, instead of another session2 that reuses address IP10. For example, when session2 reuses address IP10 and creates a connection [IP10, IP0] between server IP1 and the proxy after session1 is migrated, the proxy will now have no problem letting both connections coexist. Even though they have exactly the same tuple [IP0, IP10] at the transport layer level, they each are bound to a different VNIC (with the same address IP0), guaranteeing conflict-free operation. Remember that when session2 creates its connection [IP0, IP10], connection labels will be exchanged between server IP1 and the proxy. But since session2 has not been migrated, the labels are not used and the connection is solely identified by the tuple [IP0, IP10]. In other words, packets belonging to session2’s connection will not have a label attached while packets belonging to session1’s connection will have label1 attached. This allows the proxy to correctly identify both connections and perform virtual-physical translation as needed, e.g., session1’s connection is mapped to [IP0, IP20] while no mapping is necessary for session2’s connection.

MobiDesk’s network virtualization and privatization may raise concerns about the scalability of the proxy. We will show in the experimental results section that MobiDesk does not adversely affect the scalability of the proxy, because the mechanisms for supporting virtualization and privatization are all simple functions, such as address translation and interface redirection.

6. EXPERIMENTAL RESULTS

We have implemented a prototype MobiDesk system for serving Linux desktop computing environments. On the server-side, our prototype consists of a virtual display driver module for XFree86 [53] and a loadable kernel module for

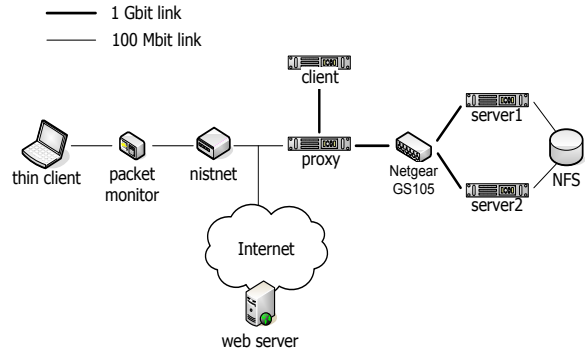


Figure 4: Experimental Testbed

operating system and network virtualization. The display driver runs as part of the display system of the hosting server, and the kernel module is loaded at the hosting server and the proxy. The server-side of our prototype works with unmodified Linux applications and any off-the-shelf Linux 2.4 kernel. On the client-side, our prototype provides a small client application that can be downloaded and run on any unmodified client to provide MobiDesk functionality. We have implemented both Xlib and Java versions of the MobiDesk client application, which can run on both Unix/Linux and Windows clients. We present experimental results using our Linux MobiDesk prototype to quantify its overhead and demonstrate its performance on various desktop computing applications.

Figure 4 shows the isolated network testbed we used for our experiments. The testbed consists of eight IBM Netfinity 4500R machines and a Micron desktop PC. The Netfinity machines each had a 933Mhz Intel Pentium-III CPU and 512MB RAM, and all of them were connected via gigabit Ethernet. The Micron desktop PC had a 450Mhz Intel Pentium-II CPU and 128MB RAM, and was used as the MobiDesk client. Four of the machines served as a MobiDesk server infrastructure consisting of one NFS file server, one proxy server running a delegate 8.9.2 [11] general-purpose application level proxy, and two computing session servers. One machine was connected on the client-side of the MobiDesk proxy and was used as a NISTNet 2.0.12 WAN emulator which could adjust the network characteristics seen by the client. Four machines were connected to the client-side of the WAN emulator, one Micron PC used as a MobiDesk client, a second used as an external web server, a third used as a packet monitor running Ethereal Network Analyzer 0.9.13 for measurement purposes, and the last used as a client for network virtualization overhead measurements. All of the machines ran Debian Linux, with the two computing session servers running Debian Stable with a Linux 2.4.5 kernel and Debian Unstable with a Linux 2.4.18 kernel, respectively. The MobiDesk client machine was installed with a dual boot configuration and also ran Microsoft Windows XP Professional.

6.1 MobiDesk Virtualization Overhead

To measure the cost of MobiDesk’s operating system virtualization, we used a range of micro benchmarks and real application workloads and measured their performance on our prototype and a vanilla Linux system. Table 2 shows the seven microbenchmarks and the four application benchmarks we used to quantify MobiDesk’s operating system

Name	Description	Linux
getpid	average <code>getpid</code> runtime	350 ns
ioctl	average runtime for the FIONREAD <code>ioctl</code>	427ns
shmget-shmctl	IPC Shared memory segment holding an integer is created and removed	3361 ns
semget-semctl	IPC Semaphore variable is created and removed	1370 ns
fork-exit	process forks and waits for child which calls <code>exit</code> immediately	44.7 us
fork-sh	process forks and waits for child to run <code>/bin/sh</code> to run a program that prints "hello world" then exits	3.89 ms
Apache	Runs Apache under load and measures average request time	1.2 ms
Make	Linux Kernel compile with up to 10 process active at one time	224.5s
MySQL	Time per interaction for "TPC-W like" benchmark	8.33s

Table 2: Application Benchmarks

virtualization overhead, as well as the results for a vanilla Linux system. To obtain accurate measurements, we rebooted the system between measurements. Additionally, the system call microbenchmarks directly used the TSC register available Pentium CPUs to record timestamps at the significant measurement events. The average timestamp event cost was 32 ns. The files for the benchmarks were stored on the NFS server. All of these benchmarks were performed in a chrooted environment on the NFS client machine running Debian Unstable with a Linux 2.4.18 kernel. Figure 5 shows the results of running the benchmarks under both configurations, with the vanilla Linux configuration normalized to one. Since all benchmarks measure the time to run the benchmark, a small number is better for all results.

The results in Figure 5 show that the operating system virtualization overhead is small. MobiDesk incurs less than 10% overhead for most of the microbenchmarks and less than 4% overhead for the application workloads. The overhead for the simple system call `getpid` benchmark is only 7% compared to vanilla Linux, reflecting the fact that virtualization for these kinds of system calls only requires an extra procedure call and a hash table lookup. The most expensive benchmarks for MobiDesk is `semget+semctl` which took 51% longer than vanilla Linux. The cost reflects the fact that our unoptimized MobiDesk prototype needs to allocate memory and do a number of namespace translations. The `ioctl` benchmark also has high overhead, because of the 12 separate assignments it does to protect the call against malicious root processes. This is large compared to the simple FIONREAD `ioctl` that just performs a simple dereference. However, since the `ioctl` is simple, we see that it only adds 200 ns of overhead over any `ioctl`. For real applications, the most overhead was only 4% for the Apache workload, where we used the `http_load` benchmark [18] to place a parallel fetch load on the server with 30 clients fetching at the same time. Similarly, we tested MySQL as part of a web commerce scenario outlined by TPC-W with a bookstore servlet running on top of Tomcat with a MySQL back-end [50]. The MobiDesk overhead for this scenario was less than 2% versus vanilla Linux.

To measure the cost of MobiDesk’s network virtualization, we used `netperf 2.2pl4` [28] to measure MobiDesk network I/O overhead versus vanilla Linux in terms of throughput, latency, CPU utilization, and connection setup. We ran

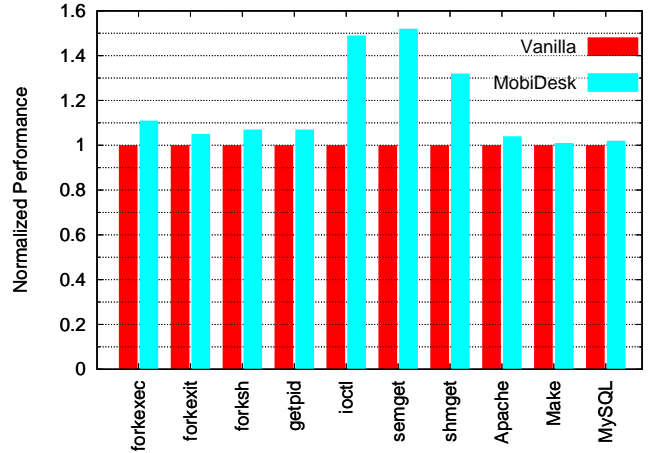


Figure 5: Operating System Virtualization Overhead

the `netperf` client on the Netfinity client and the `netperf` server on the MobiDesk session server. We used the Netfinity client for these experiments instead of the MobiDesk client so that all machines used for the network virtualization measurements were connected via gigabit Ethernet. To ensure that we were accurately measuring the performance overheads of our systems as opposed to raw network link performance, we used gigabit Ethernet for our experiments so that the network link capacity could not be saturated easily. All connections from the `netperf` client to the `netperf` server were made through the delegate proxy. We compared the performance of three different system configurations: Vanilla, MobiDesk1, and MobiDesk2. The Vanilla system is a stock Linux system without MobiDesk loaded into the kernel. The MobiDesk1 and MobiDesk2 are systems with MobiDesk loaded. On MobiDesk1, no connections are migrated and hence only connection virtualization is performed; on MobiDesk2, all connections are migrated and hence both connection virtualization and virtual-physical mapping are performed.

Figures 6 to 8 show the results for running the `netperf` throughput experiment, latency experiment, and connection setup experiment. CPU utilization measurements are omitted due to space constraints, but show similar overhead results. The throughput experiment simply measures the throughput achieved when sending messages of varying sizes as fast as possible from the client to the server. Figure 6 shows the throughput overhead for the three systems we tested. We can see that MobiDesk1 performs very close to Vanilla, with an overhead of about 1.4Mbits/second. MobiDesk2 shows the throughput overhead due to the virtual-physical mapping, which is around 10Mbits/second.

The latency experiment measures the inverse of the transaction rate, where a transaction is the exchange of a request message of size 128 bytes and a reply message of varying sizes between the client and the server over a single connection. Figure 7 shows the latency overhead for the three systems we tested. The results bear the same characteristic as that for the throughput overhead. Performance difference between Vanilla and MobiDesk1 is about 9.4 microseconds, while latency due to the virtual-physical mapping in MobiDesk2 can be observed to be around 40 microseconds. Note that there is a strange drop of latency above a reply

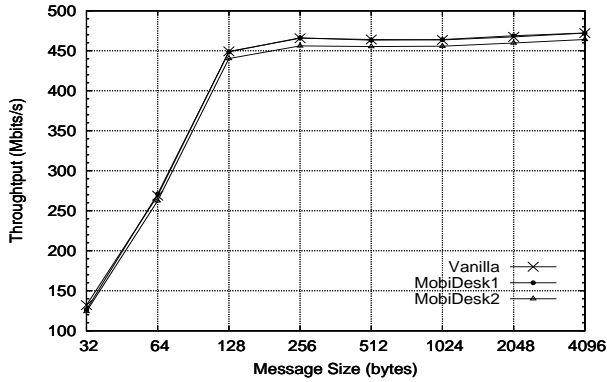


Figure 6: Virtualization Throughput Overhead

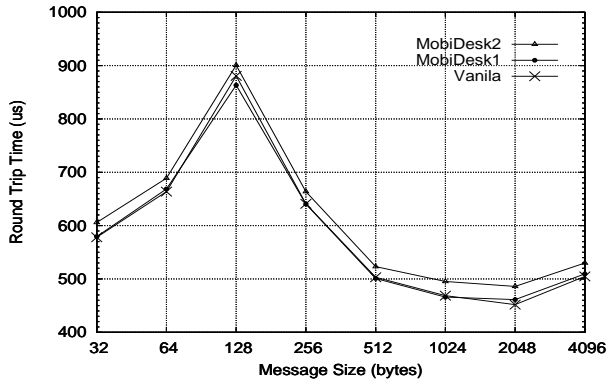


Figure 7: Virtualization Latency Overhead

message size of 128 bytes. We determined that this unusual behavior is due to a problem with the Linux device driver for the Intel Pro/1000 NIC that was used. While the behavior is unusual, it does not affect the key result shown, which is the small relative performance difference between using vanilla Linux and MobiDesk.

The TCP connection setup experiment is the same as the latency experiment except that a new connection is used for every request/response transaction. This experiment simulates the interaction between a client and server in which many short-lived connections are opened and closed. Figure 8 shows the TCP connection setup overhead for Vanilla and MobiDesk1. Note that since connection setup occurs before migration, there is no virtual-physical mapping overhead associated with connection setup, therefore this measurement is not applicable to MobiDesk2. From the figure we can see that the overhead is fewer than 10 transactions per second. Due to the same Linux driver problem in the latency test, we also see a strange increase of connection rate above reply message size of 128 bytes.

To measure the scalability of MobiDesk’s proxy infrastructure, we used `httperf` [17] to stress test the delegate proxy. We ran `httperf` on the Netfinity client and delegate on the Netfinity proxy, which distributes client requests evenly to the two Netfinity servers. The client requests are for a 4KB file which is locally available on both servers. We measured scalability relative to both number of simultaneous connection, and rate of new connections. Due to space constraints, we only present scalability relative to rate of new connections since it is the more stressful test. Figure 9 shows the connection setup latency versus rate of new connections. It shows that while MobiDesk1 has a moderate overhead of

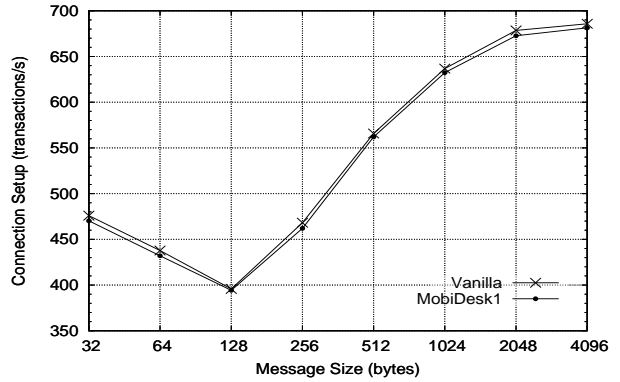


Figure 8: TCP Connection Setup Overhead

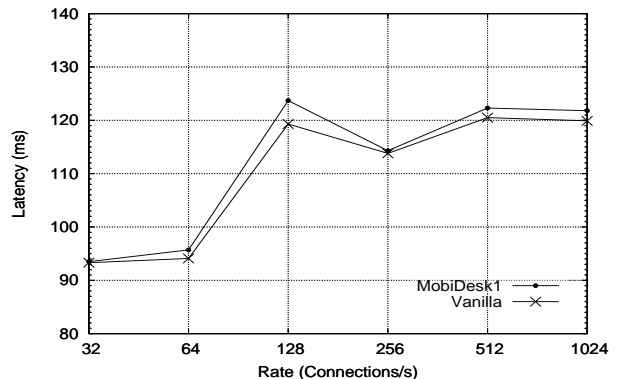


Figure 9: Scalability: Latency vs Rate of Connections

less than 4%, the important thing to note is that the overhead does not increase even after the proxy was overloaded at around 128 connections per second.

6.2 MobiDesk Application Performance

To evaluate MobiDesk performance on real desktop applications, we conducted experiments to measure the display performance of MobiDesk for web and multimedia applications and the migration performance of MobiDesk in moving a user’s desktop computing session from one server to another. To measure display performance, we compared MobiDesk against running applications on a local PC. We also compared MobiDesk running with XFree86 4.3.0 against other popular commercial thin-client systems, including Citrix MetaFrame XP for Windows [9], VNC 3.3.7 for Linux [38], and Sun’s SunRay 2.0 [48]. All of the thin-client systems, except SunRay, used the Micron PC as the client and a Netfinity server as the server. Since SunRay requires Sun hardware to run, we added a SunRay I hardware thin-client and a Solaris 9 v210 server to our experimental testbed since it does not run with the common hardware/software configuration used by the other systems.

We evaluated display performance using two popular desktop application scenarios, web browsing and video playback. Web browsing performance was measured using a Mozilla 1.4 browser to run a benchmark based on the Web Text Page Load test from the Ziff-Davis iBench benchmark suite [19]. The benchmark consists of a sequence of 54 web pages containing a mix of text and graphics. The browser window was set to 1024x768 for all platforms measured. Video playback performance was measured using a video player to

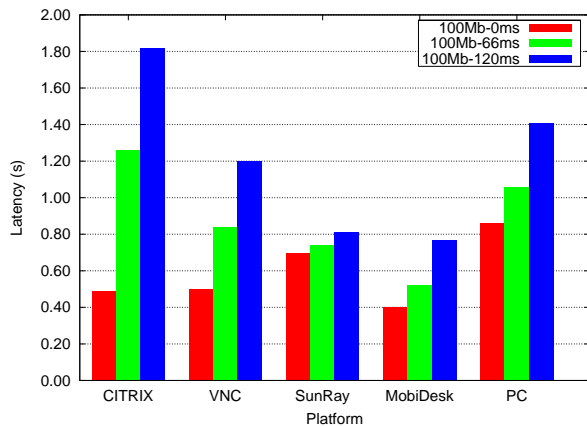


Figure 10: Average Per Web Page latency

play a 34.75 s video clip of original size 352x240 pixels displayed at 1024x768 full screen resolution. In the Unix platforms we used MPlayer 1.0pre3 as the video player, while for the Windows platforms we used the standard Windows Media Player. We used the packet monitor in our testbed to measure benchmark performance on the thin-client systems using slowmotion benchmarking [29], which allows us to quantify system performance in a non-invasive manner by capturing network traffic. The primary measure of web browsing performance was the average page download latency. The primary measure of video playback performance was video quality [29], which accounts for both playback delays and frame drops that degrade playback quality. For example, 100 percent video quality means that all video frames were displayed at real-time speed. On the other hand, 50 percent video quality could mean that half the video frames were dropped when displayed at real-time speed or that the clip took twice as long to play even though all of the video frames were displayed.

For both benchmarks, we measured all systems in three representative network scenarios: LAN, with an available network bandwidth of 100 Mbps and no introduced network latency (100Mb-0ms), and two WAN scenarios, one with 100 Mbps available network bandwidth and 66 ms round-trip network latency (100Mb-66ms), representative of cross-country and transatlantic latencies [15], and another with 100 Mbps available network bandwidth and 120 ms round-trip network latency (100Mb-120ms), representative of typical transpacific latencies [15]. For the WAN tests we increased the default TCP window size for both server and client. SunRay was unaffected by this since it uses UDP.

Figures 10 and 11 show the web browsing performance results in terms of the perceived latency and average per page data transfer, respectively. Figures 12 and 13 show the video playback performance results in terms of the video quality and total data transferred, respectively. Figure 10 shows that MobiDesk has the smallest web page download latencies, thus providing the best overall performance. Figures 10 and 11 together show that MobiDesk outperforms the local PC even though the PC transfers less data. While the PC runs the web browser on the slower PC client, MobiDesk runs the browser on the faster server enabling it to process web pages faster. The worst web browsing platform is Citrix MetaFrame, which adopts a more high-level display approach that results in poor WAN performance because of the tight coupling required between the application running

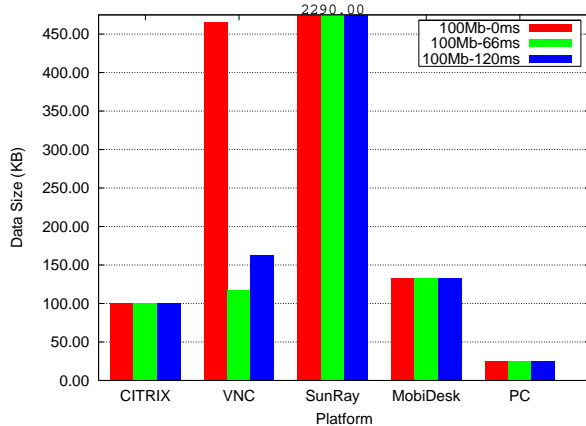


Figure 11: Average Per Web Page Data Transferred

on the server and the Citrix viewer running on the client. VNC has the second worst WAN web browsing performance in part because it relies on a client pull model for sending display updates as opposed to MobiDesk’s server-push model, which avoids roundtrip latencies providing better interactive response time. In addition, as a response to the limited WAN network conditions VNC adaptively uses more efficient compression algorithms, thus reducing its data transfer, but increasing its latency, and worsening its overall web browsing performance.

Our web browsing experiments under WAN conditions show that increased network latency can result in increased web page latencies when using TCP-based thin-client systems. This is due to the fact that TCP implementations reduce the congestion window by half for every roundtrip time that a connection has been idle [16]. As our benchmark mimics traditional web browsing usage by adding delays between the display of each web page, the thin-client connection ends up going idle and goes through slow-start each time a new page starts downloading. As network latency increases, the TCP connection takes longer to recover from its idle state, thus increasing the time it takes for web pages to load. Like other TCP-based thin-client systems, MobiDesk has higher web page latencies for the web benchmark in the presence of transpacific network latencies. However, Figure 10 shows that MobiDesk continues to provide superior sub-second performance over existing systems, even for high latency network connections.

Figure 12 shows that MobiDesk provides perfect video quality in the same manner as the traditional desktop PC, and that all of the other platforms deliver very poor video quality. They suffer from an inability to distinguish video data from normal display updates and apply ineffective compression algorithms on the video data, which are unable to keep up with the stream of updates being generated. In contrast, the results show that MobiDesk’s ability to leverage local client video hardware in delivering video using alternative YUV formats provides substantial performance benefits over other thin-client systems. VNC provides the worst overall performance primarily because of its use of a client pull model instead of MobiDesk’s server push model. In order to display each video frame, the VNC client needs to send an update request to the server. Clearly, video frames are generated faster than the rate at which the client can send requests to the server. Finally, Figure 13 shows that MobiDesk’s 100% video quality does not imply high resource

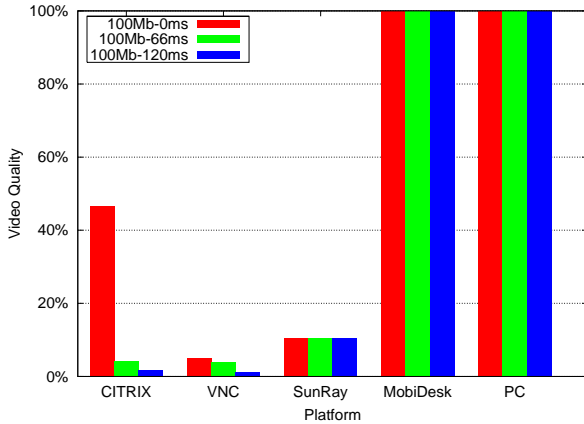


Figure 12: Video Quality

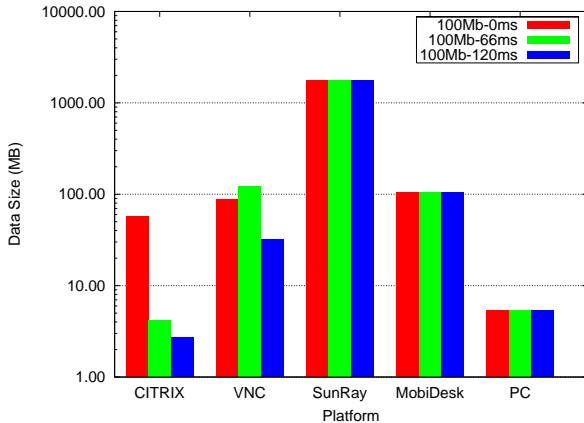


Figure 13: Total Video Data Transferred

utilization. The total data transferred translates into a bandwidth utilization of roughly 24Mbps. While VNC and Citrix consume less bandwidth, their video quality is too low to provide useful video delivery. The PC is the most bandwidth efficient platform overall, since it transfers the compressed video stream for local decoding on the client.

To measure real application performance in terms of the cost of migration, we migrated a complete KDE desktop computing environment from one MobiDesk server to another. The applications running in the KDE computing session when it was migrated are described in Table 3. The KDE session had over 30 different processes running, providing the desktop applications as well as substantial window system infrastructure, in particular, a framework for inter-application sharing. The session also included a rich desktop interface managed by a window manager, and a number of applications running in a panel, such as a clock. To demonstrate the ability to migrate a complete computing session across Linux kernels with different minor versions, we checkpointed the KDE session on the 2.4.5 kernel client machine and restarted it on the 2.4.18 kernel machine. For this experiment, the workloads were checkpointed to and restarted from local disk. The resulting checkpoint and restart times were less than a second, .85 s and .94 s, respectively. The checkpointed image was only 35 MB for a full desktop computing session, which can be easily compressed using bzip2 down to 8.8 MB. Our results show that MobiDesk can be used to provide fast migration of computing sessions among MobiDesk servers with modest checkpoint state.

Application	Description
MobiDesk	Remote display server
KDE	Entire KDE 2.2.2 environment, including window manager, panel and utilities
SSH	openssh 3.4p1 client inside a KDE konsole terminal connected to a remote host
Shell	The Bash 2.05a shell running in a konsole terminal
KGhostView	A PDF viewer with a 450k 16 page PDF file loaded
Konqueror	A modern standards compliant web browser that is part of KDE
KOffice	The KDE word processor and spreadsheet programs

Table 3: Migrated KDE Desktop Computing Session

7. RELATED WORK

MobiDesk provides a utility computing infrastructure for desktop computing. Other utility computing approaches have been proposed, primarily in the context of web services and grid computing [13, 14]. For example, IBM’s Oceano project [2] proposed the use of a pool of web servers that could be reallocated to different customers based on their usage. Web services utilities focus on web applications and grid computing utilities focus on scientific applications and other applications written specifically for grids. These approaches do not support hosting general desktop computing environments, which is the focus of MobiDesk. Plan9 [32] also provides an infrastructure-based approach to desktop computing, but does not provide the same kind of mobility support with unmodified applications and operating system kernels. More generally, IBM’s on demand computing and Hewlett-Packard’s utility computing initiative also demonstrate industry interest and trends toward a utility computing model.

MobiDesk display virtualization builds upon the thin-client computing work of two of the authors [5]. While other thin-client approaches [9, 10, 38, 41, 43] have also been developed, they were not designed for WAN environments with higher network latencies. Approaches such as Citrix MetaFrame operate using higher-level display primitives that are designed more for low bandwidth environments, but can result in worse WAN performance [23]. MobiDesk’s virtual device driver approach is most similar to SunRay [42, 43], but MobiDesk provides a more effective mapping of display commands to protocol primitives to significantly improve performance. MobiDesk also employs additional mechanisms for latency variability, that make its design more suitable for WAN environments. The InfoPad project [27] developed a system to provide access to multimedia content over wireless networks to portable devices. In contrast to MobiDesk’s application transparent approach, InfoPad’s video delivery requires the use of specific formats and playback applications to function properly.

MobiDesk operating system virtualization builds on previous work by three of the authors on Zap [30, 34], a system for transparent migration of unmodified applications inspired in part by [42]. Many other systems have been proposed to support process migration [3, 4, 7, 8, 12, 25, 26, 33, 35, 37, 39], but these systems do not provide general migration of unmodified legacy applications, across independent commodity operating systems. TUI [45] is one of the few systems that provides support for process migration across machines running different operating systems and hardware

architectures. Unlike MobiDesk, TUI has to compile applications on each platform using a special compiler and does not work with unmodified applications. Virtual machine monitors (VMMs) provide an alternative virtualization approach that can be used to migrate an entire operating system environment [21, 40, 51]. Unlike MobiDesk, VMMs decouple processes from the underlying machine hardware, but tie them to an instance of an operating system. As a result, VMMs cannot migrate processes apart from the operating system and cannot continue running those processes if the operating system ever goes down, such as during security upgrades. In contrast, MobiDesk decouples process execution from the underlying operating system allowing it to migrate processes to another computer when an operating system instance is upgraded.

MobiDesk network virtualization provides mobile communication support. Many other approaches have been developed for network mobility [6, 20, 22, 31, 36, 46, 47, 49, 54, 55, 56]. However, with the exception of [55], none of them is designed with process migration integration in mind. In addition, these approaches often require network infrastructure support to address general mobility issues (e.g., locating a mobile host) that do not apply to the application environment of MobiDesk. Transport layer solutions such as [46, 47] can provide fine grain connection migration without additional network infrastructure support. However, they require modifying the transport protocol itself, making them more difficult to deploy and limiting their use. Application layer solutions such as [36, 55, 56] can also provide fine grain connection migration without additional network infrastructure support; neither do they require modifying the transport protocol. However, the migration is emulated by closing the old TCP connection and opening a new one. The emulation requires double buffering at the application layer to account for in-flight data that have been received by TCP but not yet delivered to the application, since these data are lost when the old connection is closed. This results in substantial network I/O overhead [55] even when the connections are not migrated. In contrast, MobiDesk employs a novel, low-overhead virtualization mechanism to provide network mobility without requiring network infrastructure support or transport protocol changes.

8. CONCLUSIONS

We have introduced MobiDesk, an architecture for centralized hosting of desktop computing sessions. MobiDesk hosts computing sessions within virtualized private environments by abstracting three key resources: display, operating system and network. Display virtualization allows MobiDesk to provide fast remote access to sessions across LAN and WAN environments. Operating system virtualization allows MobiDesk to migrate sessions among hosting servers to provide high-availability computing in the presence of server maintenance and upgrades. Network virtualization allows MobiDesk to transparently maintain persistent connections to unmodified outside hosts, even as a session migrates from one server to another.

We have implemented and evaluated the performance of a MobiDesk prototype in Linux. Our implementation demonstrates that MobiDesk can support unmodified applications in hosted computing sessions without any changes to operating system kernels, network infrastructure, or network protocols. Our experimental results with real applications and

hosted desktop computing sessions show that MobiDesk has low virtualization overhead, can migrate computing sessions with subsecond checkpoint/restart times, and provides superior display performance over other remote display systems. MobiDesk is unique in its ability to offer a complete desktop experience remotely with full-motion video support. It can even provide better performance than running a desktop session on a local PC for more resource-constrained clients. Given its performance and centralized hosting model, MobiDesk provides the foundation for a utility computing infrastructure that can dramatically reduce the management complexity and costs of desktop computing.

9. ACKNOWLEDGMENTS

Our shepherd, Ramon Caceres, provided helpful comments on earlier drafts of this paper. Dinesh Subhraveti contributed to the implementation of the session migration and virtualization mechanism. Leo Kim implemented MobiDesk's Java client. This work was supported in part by NSF grants ANI-0117738 and ANI-0240525, a DOE Early Career Award, and an IBM SUR Award.

10. REFERENCES

- [1] 100x100 Project. <http://100x100network.org/>.
- [2] K. Appleby, S. Fakhouri, L. Fong, G. Goldszmidt, and M. Kalantar. Oceano: SLA based Management of a Computing Utility. In *Proceedings of IFIP/IEEE International Symposium on Integrated Network Management*, May 2001.
- [3] Y. Artsy, Y. Chang, and R. Finkel. Interprocess Communication in Charlotte. *IEEE Software*, pages 22–28, Jan 1987.
- [4] A. Barak and R. Wheeler. MOSIX: An Integrated Multiprocessor UNIX. In *Proceedings of the USENIX Winter 1989 Technical Conference*, pages 101–112, San Diego, CA, Feb. 1989.
- [5] R. Baratto, J. Nieh, and L. Kim. THINC: A Remote Display Architecture for Thin-Client Computing. Technical Report CUCS-027-04, Department of Computer Science, Columbia University, July 2004.
- [6] P. Bhagwat and C. Perkins. A Mobile Networking System based on Internet Protocol (IP). In *Proceedings of USENIX Symposium on Mobile and Location Independent Computing*, pages 69–82, Cambridge, MA, Aug. 1993.
- [7] J. Casas, D. Clark, R. Konuru, S. Otto, R. Prouty, and J. Walpole. MPVM: A Migration Transparent Version of PVM. *Computing Systems*, 8(2):171–216, 1995.
- [8] D. Cheriton. The V Distributed System. *Communications of the ACM*, 31(3):314–333, Mar 1988.
- [9] Citrix ICA Technology Brief. Technical White Paper, Boca Research, Boca Raton, FL, 1999.
- [10] B. Cumberland, G. Carius, and A. Muir. *Microsoft Windows NT Server 4.0, Terminal Server Edition: Technical Reference*. Microsoft Press, Redmond, WA, Aug. 1999.
- [11] Delegate. <http://www.delegate.org>.
- [12] F. Douglass and J. Ousterhout. Transparent Process Migration: Design Alternatives and the Sprite Implementation. *Software - Practice and Experience*, 21(8):757–785, Aug. 1991.
- [13] I. Foster, C. Kesselman, J. Nick, and S. Tuecke. The Physiology of the Grid: An Open Grid Services Architecture for Distributed Systems Integration. In *Open Grid Service Infrastructure WG, Global Grid Forum*, June 2002.
- [14] D. Gannon, R. Bramley, G. Fox, S. Smallen, A. Rossi, R. Ananthkrishnan, F. Bertrand, K. Chiu, M. Farrellee,

- M. Govindaraju, S. Krishnan, L. Ramakrishnan, Y. Simmhan, A. Slominski, Y. Ma, C. Olariu, and N. Rey-Cenz. Programming the Grid: Distributed Software Components, P2P and Grid Web Services for Scientific Applications. *Cluster Computing*, 5(3), 2002.
- [15] Global Crossing's IP Network Performance. http://www.globalcrossing.com/xml/network/net_ip_performance.xml.
- [16] M. Handley, J. Padhye, and S. Floyd. TCP Congestion Window Validation, RFC 2861. ACIRI, June 2000.
- [17] httpperf. http://www.hpl.hp.com/personal/David_Mosberger/httpperf.html.
- [18] http.load. http://www.acme.com/software/http_load/.
- [19] i-Bench, Ziff-Davis, Inc. <http://www.veritest.com/benchmarks/i-bench/>.
- [20] J. Ioannidis, D. Duchamp, and G. Q. Maguire. IP-based Protocols for Mobile Internetworking. In *Proceedings of ACM SIGCOMM*, pages 235–245, 1991.
- [21] Internet Suspend/Resume Project. <http://info.pittsburgh.intel-research.net/projects/isr/>.
- [22] D. B. Johnson and C. Perkins. Mobility Support in IPv6. *draft-ietf-mobileip-ipv6-16.txt*, IETF, Mar. 2002.
- [23] A. Lai and J. Nieh. Limits of Wide-Area Thin-Client Computing. In *Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS 2002)*, pages 228–239, Marina del Rey, CA, June 2002.
- [24] D. Legard. Korea to Build 100Mbps Internet System. *Infoworld*, Nov.18 2003. http://www.infoworld.com/article/03/11/18/HNkorea_1.html.
- [25] M. Litzkow, T. Tannenbaum, J. Basney, and M. Livny. Checkpoint and Migration of UNIX Processes in the Condor Distributed Processing System. Technical Report 1346, Department of Computer Science, University of Wisconsin, Madison, Apr. 1997.
- [26] S. J. Mullender, G. v. Rossum, A. S. Tanenbaum, R. v. Renesse, and H. v. Staveren. Amoeba: A Distributed Operating System for the 1990s. *IEEE Computer*, 23(5):44–53, May 1990.
- [27] S. Narayanaswamy, S. Seshan, E. Brewer, R. Brodersen, F. Burghardt, A. Burstein, Y. C. Chang, A. Fox, J. Gilbert, R. Han, R. H. Katz, A. Long, D. Messerschmitt, and J. Rabaey. Application and Network Support for InfoPad. *IEEE Personal Communications*, Mar. 1996.
- [28] Netperf. <http://www.netperf.org/>.
- [29] J. Nieh, S. J. Yang, and N. Novik. Measuring Thin-Client Performance Using Slow-Motion Benchmarking. *ACM Transactions on Computer Systems (TOCS)*, 21(1):87–115, Feb. 2003.
- [30] S. Osman, D. Subhraveti, G. Su, and J. Nieh. The Design and Implementation of Zap: A System for Migrating Computing Environments. In *Proceedings of the Fifth Symposium on Operating Systems Design and Implementation (OSDI 2002)*, Boston, MA, Dec. 2002.
- [31] C. Perkins. IP Mobility Support for IPv4, revised. *draft-ietf-mobileip-rfc2002-bis-08.txt*, Internet Draft, Sept. 2001.
- [32] R. Pike, D. Presotto, S. Dorward, B. Flandrena, K. Thompson, H. Trickey, and P. Winterbottom. Plan 9 from Bell Labs. Technical White Paper, Bell Laboratories, Murray Hill, New Jersey, 1995.
- [33] J. S. Plank, M. Beck, G. Kingsley, and K. Li. Libckpt: Transparent Checkpointing under Unix. In *Proceedings of Usenix Winter 1995 Technical Conference*, pages 213–223, New Orleans, LA, Jan 1995.
- [34] S. Potter, J. Nieh, and D. Subhraveti. Secure Isolation and Migration of Untrusted Legacy Applications. Technical Report CUCS-005-04, Department of Computer Science, Columbia University, Jan. 2004.
- [35] J. Pruyne and M. Livny. Managing Checkpoints for Parallel Programs. In *2nd Workshop on Job Scheduling Strategies for Parallel Processing (In Conjunction with IPPS '96)*, Honolulu, Hawaii, Apr. 1996.
- [36] X. Qu, J. X. Yu, and R. P. Brent. A Mobile TCP Socket. In *International Conference on Software Engineering (SE '97)*, San Francisco, CA, Nov. 1997.
- [37] R. Rashid and G. Robertson. Accent: A Communication Oriented Network Operating System Kernel. In *Proceedings of the 8th Symposium on Operating System Principles*, pages 64–75, Dec 1984.
- [38] T. Richardson, Q. Stafford-Fraser, K. R. Wood, and A. Hopper. Virtual Network Computing. *IEEE Internet Computing*, 2(1):33–38, Jan-Feb 1998.
- [39] M. Rozier, V. Abrossimov, F. Armand, I. Boule, M. Gien, M. Guillemont, F. Herrman, C. Kaiser, S. Langlois, P. Léonard, and W. Neuhauser. Overview of the Chorus Distributed Operating System. In *Workshop on Micro-Kernels and Other Kernel Architectures*, pages 39–70, Seattle, WA, 1992.
- [40] C. P. Sapuntzakis, R. Chandra, B. Pfaff, J. Chow, M. S. Lam, and M. Rosenblum. Optimizing the Migration of Virtual Computers. In *Proceedings of the 5th Symposium on Operating Systems Design and Implementation*, December 2002.
- [41] R. W. Scheifler and J. Gettys. *X Window System*. Digital Press, third edition, 1992.
- [42] B. K. Schmidt. *Supporting Ubiquitous Computing with Stateless Consoles and Computation Caches*. Ph.D Thesis, Computer Science Department, Stanford University, Aug. 2000.
- [43] B. K. Schmidt, M. S. Lam, and J. D. Northcutt. The interactive performance of SLIM: a stateless, thin-client architecture. In *17th ACM Symposium on Operating Systems Principles (SOSP)*, volume 34, pages 32–47, Kiawah Island Resort, SC, December 1999.
- [44] B. Schneier. *Applied Cryptography*. John Wiley and Sons, second edition, 1996.
- [45] P. Smith and N. C. Hutchinson. Heterogeneous Process Migration: The Tui System. *Software – Practice and Experience*, 28(6):611–639, 1998.
- [46] A. C. Snoeren and H. Balakrishnan. An End-to-End Approach to Host Mobility. In *Proceedings of 6th International Conference on Mobile Computing and Networking (MobiCom'00)*, Boston, MA, Aug. 2000.
- [47] F. Sultan, K. Srinivasan, D. Iyer, and L. Iftode. Migratory TCP: Highly Available Internet Services Using Connection Migration. In *Proceedings of ICDCS*, pages 17–26, 2002.
- [48] Sun Ray Integrated Solutions. <http://www.sun.com/products/sunray1/>.
- [49] F. Teraoka, Y. Yokote, and M. Tokoro. A Network Architecture Providing Host Migration Transparency. In *Proceedings of ACM SIGCOMM*, Sept. 1991.
- [50] TPC-W Java Implementation. <http://mitglied.lycos.de/jankiefer/tpcw/>.
- [51] VMware VirtualCenter. <http://www.vmware.com/products/vmanage/>.
- [52] VMware, Inc. <http://www.vmware.com>.
- [53] The XFree86 Project. <http://www.xfree86.org>.
- [54] P. Yalagandula, A. Garg, M. Dahlin, L. Alvisi, and H. Vin. Transparent Mobility with Minimal Infrastructure. In *Technical Report 01-30, University of Texas at Austin*, June 2001.
- [55] V. C. Zandy and B. P. Miller. Reliable Network Connections. In *Proceedings of 8th ACM International Conference on Mobile Computing and Networking (Mobicom '02)*, Atlanta, GA, Sept. 2002.
- [56] Y. Zhang and S. Dao. A Persistent Connection Model for Mobile and Distributed Systems. In *4th International Conference on Computer Communications and Networks (ICCCN)*, Las Vegas, NV, Sept. 1995.