

USENIX Association

Proceedings of the  
12th USENIX Security Symposium

Washington, D.C., USA  
August 4–8, 2003



© 2003 by The USENIX Association

All Rights Reserved

For more information about the USENIX Association:

Phone: 1 510 528 8649

FAX: 1 510 548 5738

Email: [office@usenix.org](mailto:office@usenix.org)

WWW: <http://www.usenix.org>

Rights to individual papers remain with the author or the author's employer.

Permission is granted for noncommercial reproduction of the work for educational or research purposes.

This copyright notice must be included in the reproduced paper. USENIX acknowledges all trademarks herein.

# Preventing Privilege Escalation

Niels Provos  
*CITI, University of Michigan*

Markus Friedl  
*GeNUA*

Peter Honeyman  
*CITI, University of Michigan*

## Abstract

*Many operating system services require special privilege to execute their tasks. A programming error in a privileged service opens the door to system compromise in the form of unauthorized acquisition of privileges. In the worst case, a remote attacker may obtain superuser privileges. In this paper, we discuss the methodology and design of privilege separation, a generic approach that lets parts of an application run with different levels of privilege. Programming errors occurring in the unprivileged parts can no longer be abused to gain unauthorized privileges. Privilege separation is orthogonal to capability systems or application confinement and enhances the security of such systems even further.*

*Privilege separation is especially useful for system services that authenticate users. These services execute privileged operations depending on internal state not known to an application confinement mechanism. As a concrete example, the concept of privilege separation has been implemented in OpenSSH. However, privilege separation is equally useful for other authenticating services. We illustrate how separation of privileges reduces the amount of OpenSSH code that is executed with special privilege. Privilege separation prevents known security vulnerabilities in prior OpenSSH versions including some that were unknown at the time of its implementation.*

## 1 Introduction

Services running on computers connected to the Internet present a target for adversaries to compromise their security. This can lead to unauthorized access to sensitive data or resources.

Services that require special privilege for their operation are critically sensitive. A programming error here may allow an adversary to obtain and abuse the special privilege.

The degree of the escalation depends on which privileges the adversary is authorized to hold and which privileges can be obtained in a successful attack. For example, a programming error that permits a user to

gain extra privilege after successful authentication limits the degree of escalation because the user is already authorized to hold some privilege. On the other hand, a remote adversary gaining superuser privilege with no authentication presents a greater degree of escalation.

For services that are part of the critical Internet infrastructure it is particularly important to protect against programming errors. Sometimes these services need to retain special privilege throughout their lifetime. For example, in SSH, the SSH daemon needs to know the private host key during re-keying to authenticate the key exchange. The daemon also needs to open new pseudo-terminals when the SSH client so requests. These operations require durable special privilege as they can be requested at any time during the lifetime of a SSH connection. In current SSH implementations, therefore, an exploitable programming error allows an adversary to obtain superuser privilege.

Several approaches to help prevent security problems related to programming errors have been proposed. Among them are type-safe languages [30] and operating system mechanisms such as protection domains [11] or application confinement [18, 21, 28]. However, these solutions do not apply to many existing applications written in C running on generic Unix operating systems. Furthermore, system services that authenticate users are difficult to confine because execution of privileged operations depends on internal state not known to the sandbox.

Instead, this paper discusses the methodology and design of *privilege separation*, a generic approach to limit the scope of programming bugs. The basic principle of privilege separation is to reduce the amount of code that runs with special privilege without affecting or limiting the functionality of the service. This narrows the exposure to bugs in code that is executed with privileges. Ideally, the only consequence of an error in a privilege separated service is denial of service to the adversary himself.

The principle of separating privileges applies to any privileged service on Unix operating systems. It is especially useful for system services that grant authenticated users special privilege. Such services are difficult to confine because the internal state of a service is not

known to an application confinement system and for that reason it cannot restrict operations that the service might perform for authenticated users. As a result, an adversary who gains unauthorized control over the service may execute the same operations as any authenticated user. With privilege separation, the adversary controls only the unprivileged code path and obtains no unauthorized privilege.

Privilege separation also facilitates source code audits by reducing the amount of code that needs to be inspected intensively. While all source code requires auditing, the size of code that is most critical to security decreases.

In Unix, every process runs within its own protection domain, *i.e.*, the operating system protects the address space of a process from manipulation and control by unrelated users. Using this feature, we accomplish privilege separation by spawning unprivileged children from a privileged parent. To execute privileged operations, an unprivileged child asks its privileged parent to execute the operation on behalf of the child. An adversary who gains control over the child is confined in its protection domain and does not gain control over the parent.

In this paper, we use OpenSSH as an example of a service whose privileges can be separated. We show that bugs in OpenSSH that led to system compromise are completely contained by privilege separation. Privilege separation requires small changes to existing code and incurs no noticeable performance penalty.

The rest of the paper is organized as follows. In Section 2, we discuss the principle of least privilege. We introduce the concept of privilege separation in Section 3 and describe a generic implementation for Unix operating system platforms. We explain the implementation of privilege separation in OpenSSH in Section 4. In Section 5, we discuss how privilege separation improves security in OpenSSH. We analyze performance impact in Section 6. Section 7 describes related work. Finally, we conclude in Section 8.

## 2 Least Privilege

We refer to a *privilege* as a security attribute that is required for certain operations. Privileges are not unique and may be held by multiple entities.

The motivation for this effort is the principle of least privilege: every program and every user should operate using the least amount of privilege necessary to complete the job [23]. Applying the principle to application design limits unintended damage resulting from programming errors. Linden [15] suggests three ap-

proaches to application design that help prevent unanticipated consequences from such errors: defensive programming, language enforced protection, and protection mechanisms supported by the operating system.

The latter two approaches are not applicable to many Unix-like operating systems because they are developed in the C language which lacks type-safety or other protection enforcement. Though some systems have started to support non-executable stack pages which prevent many stack overflows from being exploitable, even this simple mechanism is not available for most Unix platforms.

Furthermore, the Unix security model is very coarse grained. Process privileges are organized in a flat tree. At the root of the tree is the superuser. Its leaves are the users of the system. The superuser has access to every process, whereas users may not control processes of other users. Privileges that are related to file system access have finer granularity because the system grants access based on the identity of the user and his group memberships. In general, privileged operations are executed via system calls in the Unix kernel, which differentiates mainly between the superuser and everyone else.

This leaves defensive programming, which attempts to prevent errors by checking the integrity of parameters and data structures at implementation, compile or run time. For example, defensive programming prevents buffer overflows by checking that the buffer is large enough to hold the data that is being copied into it. Improved library interfaces like *strncpy* and *strlcat* help programmers avoid buffer overflows [17].

Nonetheless, for complex applications it is still inevitable that programming errors remain. Furthermore, even the most carefully written application can be affected by third-party libraries and modules that have not been developed with the same stringency. The likelihood of bugs is high, and an adversary will try to use those bugs to gain special privilege. Even if the principle of least privilege has been followed, an adversary may still gain those privileges that are necessary for the application to operate.

## 3 Privilege Separation

This section presents an approach called *privilege separation* that cleaves an application into privileged and unprivileged parts. Its philosophy is similar to the decomposition found in micro-kernels or in Unix command line tools. Privilege separation is orthogonal to other protection mechanisms that an operating system might support, *e.g.*, capabilities or protection domains.

We describe an implementation of privilege separation that does not require special support from the operating system kernel and as such may be implemented on almost any Unix-like operating system.

The goal of privilege separation is to reduce the amount of code that runs with special privilege. We achieve this by splitting an application into parts. One part runs with privileges and the others run without them. We call the privileged part the *monitor* and the unprivileged parts the *slaves*. While there is usually only one slave, it is not a requirement. A slave must ask the monitor to perform any operation that requires privileges. Before serving a request from the slave, the monitor first validates it. If the request is currently permitted, the monitor executes it and communicates the results back to the slave.

In order to separate the privileges in a service, it is necessary to identify the operations that require them. The number of such operations is usually small compared to the operations that can be executed without special privilege. Privilege separation reduces the number of programming errors that occur in a privileged code path. Furthermore, source code audits can focus on code that is executed with special privilege, which can further reduce the incidence of unauthorized privilege escalation.

Although errors in the unprivileged code path cannot result in any immediate privilege escalation, it might still be possible to abuse them for other attacks like resource starvation. Such denial of service attacks are beyond the scope of this paper.

In the remainder of this section, we explain the Unix mechanisms that allow us to implement a privilege separated service. Processes are protection domains in a Unix system. That means that one process cannot control another unrelated process. To achieve privilege separation, we create two entities: a privileged parent process that acts as the monitor and an unprivileged child process that acts as the slave. The privileged parent can be modeled by a finite-state machine (FSM) that monitors the progress of the unprivileged child. The parent accepts requests from the child for actions that require privileges. The set of actions that are permitted changes over time and depends on the current state of the FSM. If the number of actions that require privileges is small, most of the application code is executed by the unprivileged child.

The design of the interface is important as it provides a venue of attack for an adversary who manages to compromise the unprivileged child. For example, the interface should not provide mechanisms that allow the export of sensitive information to the child, like a private signing key. Instead, the interface provides a

request that allows the child to request a digital signature.

A privilege separated service can be in one of two phases:

- Pre-Authentication Phase: A user has contacted a system service but is not yet authenticated. In this case, the unprivileged child has no process privileges and no rights to access the file system.
- Post-Authentication Phase: The user has successfully authenticated to the system. The child has the privileges of the user including file system access, but does not hold any other special privilege. However, special privilege are still required to create new pseudo-terminals or to perform other privileged operations. For those operations, the child must request an action from the privileged parent.

The unprivileged child is created by changing its user identification (UID) and group identification (GID) to otherwise unused IDs. This is achieved by first starting a privileged monitor process. It forks a slave process. To prevent access to the file system, the child changes the root of its file system to an empty directory in which it is not allowed to create any files. Afterwards, the slave changes its UID and GID to lose its process privileges.

To enable slave requests to the monitor, we use inter-process communication (IPC). There are many different ways to allow communication between processes: pipes, shared memory, etc. In our case, we establish a socket between the two processes using the *socket-pair* system call. The file descriptor is inherited by the forked child.

A slave may request different types of privileged operations from the monitor. We classify them depending on the result the slave expects to achieve: *Information*, *Capabilities*, or *Change of Identity*.

A child issues an informational request if acquiring the information requires privileges. The request starts with a 32-bit length field followed by an 8-bit number that determines the request type. In general, the monitor checks every request to see if it is allowed. It may also cache the request and result. In the pre-authentication phase, challenge-response authentication can be handled via informational requests. For example, the child first requests a challenge from the privileged monitor. After receiving the challenge, the child presents it to the user and requests authentication from the monitor by presenting the response to it. In this case, the monitor remembers the challenge that it created and verifies that the response matches. The result is either successful or unsuccessful authentication. In OpenSSH, most privileged operations can

```

cmsg = CMSG_FIRSTHDR(&msg);
cmsg->cmsg_len = CMSG_LEN(sizeof(int));
cmsg->cmsg_level = SOL_SOCKET;
cmsg->cmsg_type = SCM_RIGHTS;
*(int *)CMSG_DATA(cmsg) = fd;

```

Figure 1: File descriptor passing enables us to send a file descriptor to another process using a special control message. With file descriptor passing, the monitor can grant an unprivileged child access to a file that the child is not allowed to open itself.

be implemented with informational requests.

Ordinarily, the only capability available to a process in a Unix operating systems is a file descriptor. When a slave requests a capability, it expects to receive a file descriptor from the privileged monitor that it could not obtain itself. A good example of this is a service that provides a pseudo-terminal to an authenticated user. Creating a pseudo-terminal involves opening a device owned by the superuser and changing its ownership to the authenticated user, which requires special privilege.

Modern Unix operating systems provide a mechanism called *file descriptor passing*. File descriptor passing allows one process to give access to an open file to another process [25]. This is achieved by sending a control message containing the file descriptor to the other process; see Figure 1. When the message is received, the operating system creates a matching file descriptor in the file table of the receiving process that permits access to the sender's file. We implement a capability request by passing a file descriptor over the socket used for the informational requests. The capability request is an informational request in which the slave expects the monitor to answer with a control message containing the passed file descriptor.

The change of identity request is the most difficult to implement. The request is usually issued when a service changes from the pre-authentication to the post-authentication phase. After authentication, the service wants to obtain the privileges of the authenticated user. Unix operating systems provide no portable mechanism to change the user identity<sup>1</sup> a process is associated with unless the process has superuser privilege. However, in our case, the process that wants to change its identity does not have such privilege.

One way to effect a change of identity is to terminate the slave process and ask the monitor to create a new process that can then change its UID and GID to the desired identities. By terminating the child process all

<sup>1</sup>To our knowledge, Solaris is the only Unix operating system to provide such a mechanism.

```

mm_master_t *mm_create(mm_master_t *, size_t);
void mm_destroy(mm_master_t *);
void *mm_malloc(mm_master_t *, size_t);
void mm_free(mm_master_t *, void *);
void mm_share_sync(mm_master_t **, mm_master_t **)

```

Figure 2: These functions represent the interface for shared memory allocation. They allow us to export dynamically allocated data from a child process to its parent without changing address space references contained in opaque data objects.

the state that has been created during its life time is lost. Normally a meaningful continuation of the session is not possible without retaining the state of the slave process. We solve this problem by exporting all state of the unprivileged child process back to the monitor.

Exporting state is messy. For global structures, we use XDR-like [16] data marshaling which allows us to package all data contained in a structure including pointers and send it to the monitor. The data is unpacked by the newly forked child process. This prevents data corruption in the exported data from affecting the privileged monitor in any way.

For structures that are allocated dynamically, *e.g.*, via *malloc*, data export is more difficult. We solve this problem by providing memory allocation from shared memory. As a result, data stored in dynamically allocated memory is also available in the address space of the privileged monitor. Figure 2 shows the interface to the shared memory allocator.

The two functions *mm\_create* and *mm\_share\_sync* are responsible for permitting a complete export of dynamically allocated memory. The *mm\_create* function creates a shared address space of the specified size. There are several ways to implement shared memory, we use anonymous memory maps. The returned value is a pointer to a *mm\_master* structure that keeps track of allocated memory. It is used as parameter in subsequent calls to *mm\_malloc* and *mm\_free*. Every call to those two functions may result in allocation of additional memory for state that keeps track of free or allocated memory in the shared address space. Usually, that memory is allocated with libc's *malloc* function. However, the first argument to the *mm\_create* function may be a pointer to another shared address space. In that case, the memory manager allocates space for additional state from the passed shared address space.

Figure 3 shows an overview of how allocation in the shared address space proceeds. We create two shared address spaces: *back* and *mm*. The address space *mm* uses *back* to allocate state information. When the child

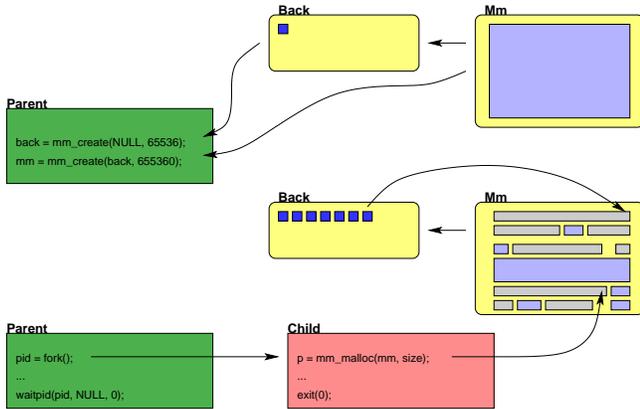


Figure 3: The complete state of a slave process includes dynamically allocated memory. When exporting this state, the dynamically allocated address space in opaque data objects must not change. By employing a shared memory allocator that is backed by another shared address space, we can export state without changing the addresses of dynamically allocated data.

wants to change its identity, it exits and the thread of execution continues in the parent. The parent has access to all the data that was allocated in the child. However, one problem remains. The shared address space *back* uses libc's malloc that allocated memory in the child's address space to keep track of its state. If this information is lost when the child process exits, then subsequent calls to *mm\_malloc* or *mm\_free* fail. To solve the problem, the parent calls the *mm\_share\_sync* function which recreates the state information in the shared address space *back*. Afterwards, freeing and allocating memory proceeds without any problems.

We use shared memory and XDR-like data marshaling to export all state from the child to the parent. After the child process exports its state and terminates, the parent creates a new child process. The new process changes to the desired UID and GID and then imports the exported state. This effects a change of identity in the slave that preserves state information.

## 4 Separating Privileges in OpenSSH

In this section, we show how to use privilege separation in OpenSSH, a free implementation of the SSH protocols. OpenSSH provides secure remote login across the Internet. OpenSSH supports protocol versions one and two; we restrict our explanation of privilege separation to the latter. The procedure is very similar for protocol one and also applies to other services that require authentication.

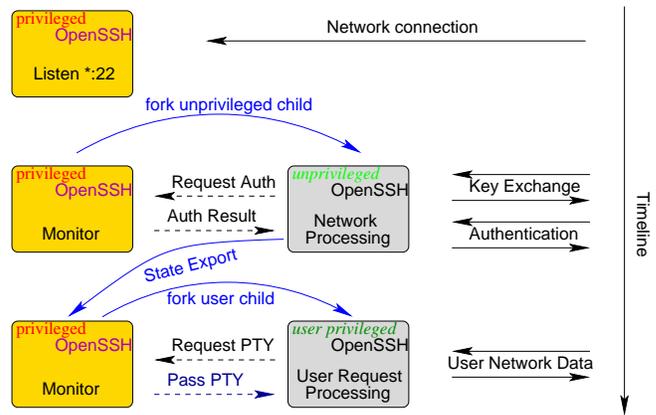


Figure 4: Overview of privilege separation in OpenSSH. An unprivileged slave processes all network communication. It must ask the monitor to perform any operation that requires privileges.

When the SSH daemon starts, it binds a socket to port 22 and waits for new connections. Every new connection is handled by a forked child. The child needs to retain superuser privileges throughout its lifetime to create new pseudo terminals for the user, to authenticate key exchanges when cryptographic keys are replaced with new ones, to clean up pseudo terminals when the SSH session ends, to create a process with the privileges of the authenticated user, etc.

With privilege separation, the forked child acts as the monitor and forks a slave that drops all its privileges and starts accepting data from the established connection. The monitor now waits for requests from the slave; see Figure 4. Requests that are permitted in the pre-authentication phase are shown in Figure 5. If the child issues a request that is not permitted, the monitor terminates.

First, we identify the actions that require special privilege in OpenSSH and show which request types can fulfill them.

### 4.1 Pre-Authentication Phase

In this section, we describe the privileged requests for the pre-authentication phase:

- Key Exchange: SSH v2 supports the Diffie-Hellman Group Exchange which allows the client to request a group of a certain size from the server [10]. To find an appropriate group the server consults the */etc/moduli* file. However, because the slave has no privileges to access the file system, it can not open the file itself, so, it issues an informational request to the monitor. The

```

struct mon_table mon_dispatch_proto20[] = {
    {MONITOR_REQ_MODULI, MON_ONCE, mm_answer_moduli},
    {MONITOR_REQ_SIGN, MON_ONCE, mm_answer_sign},
    {MONITOR_REQ_PWNAM, MON_ONCE, mm_answer_pwnamallow},
    {MONITOR_REQ_AUTHSERV, MON_ONCE, mm_answer_authserv},
    {MONITOR_REQ_AUTHPASSWORD, MON_AUTH, mm_answer_authpassword},
[... ]
    {MONITOR_REQ_KEYALLOWED, MON_ISAUTH, mm_answer_keyallowed},
    {MONITOR_REQ_KEYVERIFY, MON_AUTH, mm_answer_keyverify},
    {0, 0, NULL}
};

```

Figure 5: The table describes valid requests that a slave may send to the monitor in the pre-authentication phase for SSH protocol version two. After authentication, the set of valid requests changes and is described by a separate table.

monitor returns a suitable group after consulting the moduli file. The returned group is used by the slave for the key exchange. As seen in Figure 5, the slave may issue this request only once.

- **Authenticated Key Exchange:** To prevent man-in-the-middle attacks, the key exchange is authenticated. That means that the SSH client requires cryptographic proof of the server identity. At the beginning of the SSH protocol, the server sends its public key to the client for verification. As the public key is public, the slave knows it and no special request is required. However, the slave needs to ask the monitor to authenticate the key exchange by signing a cryptographic hash of all values that have been exchanged between the client and the server. The signature is obtained by an informational request.
- **User Validation:** After successful key exchange, all communication is encrypted and the SSH client informs the server about the identity of the user who wants to authenticate to the system. At this point, the server decides if the user name is valid and allowed to login. If it is invalid, the protocol proceeds but all authentication attempts from the client fail. The slave can not access the password database, so it must issue an informational request to the server. The server caches the user name and reports back to the slave if the name is valid.
- **Password Authentication:** Several methods can be used to authenticate the user. For password authentication, the SSH client needs to send a correct login and password to the server. Once again, the unprivileged slave can not access the password database, so it asks the monitor to verify the password. The monitor informs the slave if the authentication succeeds or fails. If it succeeds, the

pre-authentication phase ends.

- **Public Key Authentication:** Public Key Authentication is similar to password authentication. If it is successful, the pre-authentication phase ends. However, two informational requests are required to use public keys for authentication. The first request allows the slave to determine if a public key presented by the client may be used for authentication. The second request determines if the signature returned by the client is valid and signs the correct data. A valid signature results in successful authentication.

At any time, the number of requests that the slave may issue are limited by the state machine. When the monitor starts, the slave may issue only the first two requests in Figure 5. After the key exchange has finished, the only valid request is for user validation. After validating the user, all authentication requests are permitted. The motivation for keeping the number of valid requests small is to reduce the attack profile available to an intruder who has compromised the slave process.

All requests up to this point have been informational. The pre-authentication phase ends with successful authentication as determined by the monitor. At this point, the slave needs to change its identity to that of the authenticated user. As a result, the slave obtains all privileges of the user, but no other privileges. We achieve this with a change of identity request.

The monitor receives the state of the slave process and waits for it to exit. The state consists of the following: the encryption and authentication algorithms including their secret keys, sequence counters for incoming and outgoing packets, buffered network data and the compression state.

Exporting the cryptographic key material is uncomplcated. The main problem is exporting the compression state. The SSH protocols use the *zlib* compression format [7, 8] which treats network data as a stream instead of sequence of packets. Treating network data as a stream allows *zlib* to improve its dictionary with increasing amount of compressed data. On the other hand, it also means that compression in the server cannot be stopped and then restarted as the client uses a dictionary that depends on all the preceding data. Fortunately, *zlib* provides hooks for user supplied memory management functions. We provide it with functions that use *mm\_malloc* and *mm\_free* as back end. After the child exits, the monitor needs only to call *mm\_share\_sync* to import the compression state.

## 4.2 Post-Authentication Phase

The monitor forks a new process that then changes its process identification to that of the authenticated user. The slave process obtains all the privileges of the authenticated user. At this point, we enter the post-authentication phase which requires only a few privileged operations. They are as follows:

- **Key Exchange:** In SSH protocol version two, it is possible to renew cryptographic keys. This requires a new key exchange, so just as in the pre-authentication phase, the monitor chooses a suitable group for the Diffie-Hellman key exchange and signs for authentication.
- **Pseudo Terminal Creation:** After authentication, the user requires a pseudo terminal whose creation requires superuser privileges. For a Unix application, a pseudo terminal is just a file descriptor. The slave issues a capability request to the monitor. The monitor creates the terminal and passes the corresponding file descriptor to the child process. An informational request suffices when the slave wants to close the pseudo terminal.

## 4.3 Discussion

Observe that the majority of all privileged operations can be implemented with informational requests. In fact, some degree of privilege separation is possible if neither capability nor change of identity requests are available. If the operating system does not support file descriptor passing, privilege separation perforce ends after the pre-authentication phase. To fully support the change of identify request shared memory is required. Without shared memory, the compression state cannot be exported without rewriting *zlib*. Nonetheless, systems that do not support shared memory can

disable compression and still benefit from privilege separation.

Using an alternative design, we can avoid the change of identity request and shared memory. Instead of using only two processes: monitor and slave, we use three processes: one monitor process and two slave processes. The first slave operates similarly to the slave process described in the pre-authentication phase. However, after the user authenticates, the slave continues to run and is responsible for encrypting and decrypting network traffic. The monitor then creates a second slave to execute a shell or remote command with the credentials of the authenticated user. All communication passes via the first child process to the second. This design requires no state export and no shared memory. Although the cryptographic processing is isolated in the first child, it has only a small effect on security. In the original design, a bug in the cryptographic processing may allow an adversary to execute commands with the privilege of the authenticated user. However, after authentication, an adversary can already execute any commands as that user. The three process design may help for environments in which OpenSSH restricts the commands a user is allowed to execute. On the other hand, it adds an additional process, so that every remote login requires three instead of two processes. While removing the state export reduces the complexity of the system, synchronizing three instead of two processes increases it. An additional disadvantage is a decrease in performance because the three process design adds additional data copies and context switches.

For the two process design, the changes to the existing OpenSSH sources are small. About 950 lines of the 44,000 existing lines of source code, or 2%, were changed. Many of the changes are minimal:

```
- authok = auth_password(authctxt, pwd);  
+ authok = PRIVSEP(auth_password(authctxt, pwd));
```

The new code that implements the monitor and the data marshaling amounts to about three thousand lines of source code, or about seven percent increase in the size of the existing sources.

While support for privilege separation increases the source code size, it actually reduces the complexity of the existing code. Privilege separation requires clean and well abstracted subsystem interfaces so that their privileged sections can run in a different process context. During the OpenSSH implementation, the interfaces for several subsystems had to be improved to facilitate their separation. As a result, the source code is better organized, more easily understood and audited, and less complex.

The basic functionality that the monitor provides is independent of OpenSSH. It may be used to enable

privilege separation in other applications. We benefit from reusing security critical source code because it results in more intense security auditing. This idea has been realized in Privman, a library that provides a generic framework for privilege separation [12].

## 5 Security Analysis

To measure the effectiveness of privilege separation in OpenSSH, we discuss attacks that we protect against and analyse how privilege separation would have affected security problems reported in the past. We assume that the employed cryptography is secure, therefore we do not discuss problems of cryptographic primitives.

After privilege separation, two thirds of the source code are executed without privileges as shown in Table 1. The numbers include code from third-party libraries such as *openssl* and *zlib*. For OpenSSH itself, only twenty five percent of the source code require privilege whereas the remaining seventy five percent are executed without special privilege. If we assume that programming errors are distributed fairly uniformly, we can estimate the increase of security by counting the number of source code lines that are now executed without privileges. This back of the envelope analysis suggests that two thirds of newly discovered or introduced programming errors will not result in privilege escalation and that only one third of the source code requires intensive auditing.

We assume that an adversary can exploit a programming error in the slave process to gain complete control over it. Once the adversary compromised the slave process, she can make any system call in the process context of the slave. We assume also that the system call interface to the operating system itself is secure<sup>2</sup>. Still, there are several potential problems that an implementation of privilege separation needs to address:

- The adversary may attempt to signal or ptrace other processes to get further access to the system. This is not possible in our design because the slave processes use their own UID.
- The adversary may attempt to signal or ptrace the slave processes of other SSH sessions. When changing the UID of a process from root to another UID, the operating system marks the process as P\_SUGID so that only root may signal or ptrace it.

<sup>2</sup>This assumption does not always hold. A bug in OpenBSD's select system call allowed an adversary to execute arbitrary code at the kernel-level [5, 20].

Subsystem	Lines of Code	Percentage
<b>Unprivileged</b>	<b>17589</b>	<b>67.70%</b>
<i>OpenSSH</i>	10360	39.88%
Ciphers	267	1.03%
Packet Handling	1093	4.21%
Miscellaneous	7944	30.58%
Privsep Interface	1056	4.06%
<i>OpenSSL</i>	3138	12.08%
Diffie Hellman	369	1.42%
Symmetric Ciphers	2769	10.66%
<i>Zlib</i>	4091	15.75%
<b>Privileged</b>	<b>8391</b>	<b>32.30%</b>
<i>OpenSSH</i>	3403	13.10%
Authentication	803	3.09%
Miscellaneous	1700	6.54%
Monitor	900	3.46%
<i>OpenSSL</i>	4109	15.82%
BigNum/Hash	3178	12.23%
Public Key	931	3.58%
<i>SKey</i>	879	3.17%

Table 1: Number of source code lines that are executed with and without privileges.

As a result, a slave process can not signal another slave.

- She may attempt to use system calls that change the file system, for example to create named pipes for interprocess communication or device nodes. However, as a non-root user the slave process has its file system root set to an empty read-only directory that the adversary can not escape from.
- Using privilege separation, we cannot prevent the adversary from initiating local network connections and potentially abusing trust relations based on IP addresses. However, we may restrict the child's ability to access the system by employing external policy enforcement mechanisms like Sys-trace [21].
- The adversary may attempt to gather information about the system, for example, the system time or PIDs of running processes, that may allow her to compromise a different service. Depending on the operating system, some information is exported only via the file system and can not be accessed by the adversary. A sandbox may help to further restrict the access to system information.

Another way an adversary may try to gain additional privileges is to attack the interface between the

privileged monitor and the slave. The adversary could send badly formatted requests in the hope of exploiting programming errors in the monitor. For that reason, it is important to carefully audit the interface to the monitor. In the current implementation, the monitor imposes strict checks on all requests. Furthermore, the number of valid requests is small and any request detected as invalid causes the privileged monitor to terminate.

Nonetheless, there may be other ways that an adversary might try to harm the system. She might try to starve the resources of the system by forking new processes or by running very time intensive computations. As a result, the system might become unusable. The effect of such an attack can be mitigated by placing process limits on the slave process. For example, we can limit the number of file descriptors the slave may open and the number of processes it is allowed to fork. The monitor may also watch other resource utilization like CPU time and terminate the slave if a certain threshold is reached.

In the following, we discuss how privilege separation would have affected previous programming errors in OpenSSH.

The *SSH-1 Daemon CRC32 Compensation Attack Detector Vulnerability* permits an adversary to gain superuser privileges remotely without authenticating to the systems [31]. The problem is caused by an integer overflow in a function that processes network packets. With privilege separation, this function is executed without any privileges, which makes it impossible for an adversary to directly compromise the system.

Similarly, the off-by-one error in OpenSSH's channel code allows an adversary to gain superuser privileges after authenticating to the system [19]. With privilege separation, the process has only the privileges of the authenticated user. An adversary cannot obtain system privileges in this case either.

A security problem in the external *zlib* compression library was found that might allow a remote adversary to gain superuser privileges without any authentication [3]. This problem occurs in a third-party library, so no audit of the OpenSSH source code itself can find it. Privilege separation prevents a system compromise in this case, too.

At the time of this writing, additional security problems have been found in OpenSSH. A bug in the Kerberos ticket passing functions allowed an authenticated user to gain superuser rights. A more severe problem in code for challenge-response authentication allows a remote adversary to obtain superuser privileges without any authentication [4]. Privilege separation prevents both of these problems and is mentioned in the CERT

advisory as a solution.

The programming errors in the channel code and in the Kerberos ticket passing functions occur in the post-authentication phase. Without privilege separation, these errors allow an authenticated user to gain superuser privilege. The remaining errors occur during pre-authentication and may allow an adversary to gain superuser privilege without any authentication if privilege separation is not used.

These examples demonstrate that privilege separation has the potential to contain security problems yet unknown.

## 6 Performance Analysis

To analyze the performance of privilege separated OpenSSH, we measure the execution time for several different operations in monolithic OpenSSH and the privileged separated version. We conduct the measurements on a 1.13 GHz Pentium III laptop with all data in the memory cache.

Test	Normal	Privsep
Login		
- compressed	0.775s ± 0.0071s	0.777s ± 0.0067s
- uncompressed	0.767s ± 0.0106s	0.774s ± 0.0097s
Data Transfer		
- compressed	4.229s ± 0.0373s	4.243s ± 0.0411s
- uncompressed	1.989s ± 0.0223s	1.994s ± 0.0143s

Table 2: Performance comparison between normal OpenSSH and privilege separated OpenSSH. We measure the overhead in login and data transfer time when employing privilege separation. In both cases, privilege separation imposes no significant performance penalty.

The first test measures the time it takes to login using public key authentication. We measure the time with compression enabled and without compression. The next two tests measure the data transfer time of a 10 MB file filled with random data, with compression enabled, and without compression. The results are shown in Table 2. It is evident that privilege separated OpenSSH does not penalize performance. As the IPC between monitor and slave is never used for moving large amounts of data, this is not surprising.

## 7 Related Work

Confidence in the security of an application starts by source code inspection and auditing. Static analysis

provides methods to automatically analyze a program's source code for security weaknesses. Using static analysis, it is possible to automatically find buffer overrun vulnerabilities [13, 27], format string vulnerabilities [24], etc.

While source code analysis enables us to find some security vulnerabilities, it is even more important to design applications with security in mind. The principle of least privilege is a guideline for developers to secure applications. It states that every program and every user should operate using the least amount of privilege necessary to complete the job [22].

Security mechanisms at the operating system level provide ways to reduce the privileges that applications run with [1, 29, 18, 21]. However, these mechanisms are unaware of an application's internal state. For example, they cannot determine if users authenticate successfully. As a result, they have to allow all operations of authenticated users even when attached by an adversary. Privilege separation remedies this problem because it is built into the application and exposes only an unprivileged child to the adversary. There are several applications that make use of privilege separation as we discuss below. The main difference in this research is the degree and completeness of the separation.

Carson demonstrates how to reduce the number of privileges that are required in the *Sendmail* mail system [2]. His design follows the principle of least privilege. While *Sendmail* is a good example, the degrees of privilege separation demonstrated in *OpenSSH* are much more extensive. For example, we show how to change the effective UID and how to retain privileges securely for the whole duration of the session.

Venema uses semi-resident, mutually-cooperating processes in *Postfix* [26]. He uses the process context as a protection domain similar to our research in privilege separation. However, a mail delivery system does not require the close interaction between privileged and unprivileged processes as necessary for authentication services like *OpenSSH*. For system services that require transitions between different privileges, our approach seems more suitable.

Evans *very secure* FTP daemon uses privilege separation to limit the effect of programming errors [9]. He uses informational and capability requests in his implementation. His work is very similar to the implementation of privilege separation in *OpenSSH*, but not as extensive and less generic.

Solar Designer uses a process approach to switch privileges in his *Owl* Linux distribution [6]. His POP3 daemon *popa3d* forks processes that execute protocol operations with lower privileges and communicate results back to the parent. The interaction between par-

ent and child is based completely on informational requests.

Separating the privileges of an application causes a decomposition into subsystems with well defined functionality. This is similar to the design and functionality of a microkernel where subsystems have to follow the principle of independence and integrity [14]. For a privilege separated application, independence and integrity are realized by multiple processes that have separate address spaces and communicate via IPC.

## 8 Conclusion

Programming errors in privileged services can result in system compromise allowing an adversary to gain unauthorized privileges.

Privilege separation is a concept that allows parts of an application to run without any privileges at all. Programming errors in the unprivileged part of the application cannot lead to privilege escalation.

As a proof of concept, we implemented privilege separation in *OpenSSH* and show that past errors that allowed system compromise would have been contained with privilege separation.

There is no performance penalty when running *OpenSSH* with privilege separation enabled.

## 9 Acknowledgments

We thank Solar Designer, Dug Song, David Wagner and the anonymous reviewers for helpful comments and suggestions.

## References

- [1] Lee Badger, Daniel F. Sterne, David L. Sherman, Kenneth M. Walker, and Sheila A. Haight. A Domain and Type Enforcement UNIX Prototype. In *Proceedings of the 5th USENIX Security Symposium*, pages 127–140, June 1995. 10
- [2] Mark E. Carson. *Sendmail without the Superuser*. In *Proceedings of the 4th USENIX Security Symposium*, October 1993. 10
- [3] CERT/CC. CERT Advisory CA-2002-07 Double Free Bug in zlib Compression Library. <http://www.cert.org/advisories/CA-2002-07.html>, March 2002. 9
- [4] CERT/CC. CERT Advisory CA-2002-18 *OpenSSH Vulnerabilities in Challenge Response Handling*. <http://www.cert.org/advisories/CA-2002-18.html>, June 2002. 9

- [5] Silvio Cesare. FreeBSD Security Advisory FreeBSD-SA-02:38.signed-error. <http://archives.neohapsis.com/archives/freebsd/2002-08/0094.html>, August 2002. 8
- [6] Solar Designer. Design Goals for popa3d. <http://www.openwall.com/popa3d/DESIGN>. 10
- [7] P. Deutsch. DEFLATE Compressed Data Format Specification version 1.3. RFC 1951, 1996. 7
- [8] P. Deutsch and J-L. Gailly. ZLIB Compressed Data Format Specification version 3.3. RFC 1950, 1996. 7
- [9] Chris Evans. Comments on the Overall Architecture of Vsftpd, from a Security Standpoint. <http://vsftpd.beasts.org/>, February 2001. 10
- [10] Markus Friedl, Niels Provos, and William A. Simpson. Diffie-Hellman Group Exchange for the SSH Transport Layer Protocol. Internet Draft, January 2002. Work in progress. 5
- [11] Li Gong, Marianne Mueller, Hemma Prafullchandra, and Roland Schemers. Going Beyond the Sandbox: An Overview of the New Security Architecture in the Java Development Kit 1.2. *USENIX Symposium on Internet Technologies and Systems*, pages 103–112, 1997. 1
- [12] Douglas Kilpatrick. Privman: A Library for Partitioning Applications. In *Proceedings of the USENIX 2003 Annual Technical Conference, FREENIX track*, June 2003. 8
- [13] David Larochelle and David Evans. Statically Detecting Likely Buffer Overflow Vulnerabilities. In *Proceedings of the 10th USENIX Security Symposium*, August 2001. 10
- [14] Jochen Liedtke. On  $\mu$ -Kernel Construction. In *Proceedings of the Symposium on Operating Systems Principles*, pages 237–250, 1995. 10
- [15] Theodore A. Linden. Operating System Structures to Support Security and Reliable Software. *ACM Computing Surveys*, 8(4):409–445, 1976. 2
- [16] Sun Microsystems. XDR: External Data Representation. RFC 1014, June 1987. 4
- [17] Todd C. Miller and Theo de Raadt. strlcpy and strlcat – Consistent, Safe, String Copy and Concatenation. In *Proceedings of the 1999 USENIX Technical Conference, FREENIX track*, June 1999. 2
- [18] David S. Peterson, Matt Bishop, and Raju Pandey. A Flexible Containment Mechanism for Executing Untrusted Code. In *Proceedings of the 11th USENIX Security Symposium*, pages 207–225, August 2002. 1, 10
- [19] Joost Pol. OpenSSH Channel Code Off-By-One Vulnerability. <http://online.securityfocus.com/bid/4241>, March 2002. 9
- [20] Niels Provos. OpenBSD Security Advisory: Select Boundary Condition. <http://monkey.org/openbsd/archive/misc/0208/msg00482.html>, August 2002. 8
- [21] Niels Provos. Improving Host Security with System Call Policies. In *Proceedings of the 12th USENIX Security Symposium*, August 2003. 1, 8, 10
- [22] Jerome H. Saltzer. Protection and the Control of Information in Multics. *Communications of the ACM*, 17(7):388–402, July 1974. 10
- [23] Jerome H. Saltzer and Michael D. Schroeder. The Protection of Information in Computer Systems. In *Proceedings of the IEEE 69*, number 9, pages 1278–1308, September 1975. 2
- [24] Umesh Shankar, Kunal Talwar, Jeffrey S. Foster, and David Wagner. Detecting Format String Vulnerabilities with Type Qualifiers. In *Proceedings of the 10th USENIX Security Symposium*, August 2001. 10
- [25] W. Richard Stevens. *Advanced Programming in the UNIX Environment*. Addison-Wesley, 1992. 4
- [26] Wietse Venema. Postfix Overview. <http://www.postfix.org/motivation.html>. 10
- [27] David Wagner, Jeffrey S. Foster, Eric A. Brewer, and Alexander Aiken. A First Step Towards Automated Detection of Buffer Overrun Vulnerabilities. In *Proceedings of the ISOC Symposium on Network and Distributed System Security*, pages 3–17, February 2000. 10
- [28] David A. Wagner. Janus: an Approach for Confinement of Untrusted Applications. Technical Report CSD-99-1056, University of California, Berkeley, 12, 1999. 1
- [29] Kenneth M. Walker, Daniel F. Sterne, M. Lee Badger, Michael J. Petkac, David L. Shermann, and Karen A. Oostendorp. Confining Root Programs with Domain and Type Enforcement (DTE). In *Proceedings of the 6th Usenix Security Symposium*, July 1996. 10
- [30] Dan S. Wallach, Dirk Balfanz, Drew Dean, and Edward W. Felten. Extensible Security Architectures for Java. *16th Symposium on Operating System Principles*, pages 116–128, 1997. 1
- [31] Michal Zalewski. Remote Vulnerability in SSH Daemon CRC32 Compensation Attack Detector. [http://razor.bindview.com/publish/advisories/adv\\_ssh1crc.html](http://razor.bindview.com/publish/advisories/adv_ssh1crc.html), February 2001. 9