

USENIX Association

Proceedings of the
General Track:
2003 USENIX Annual
Technical Conference

San Antonio, Texas, USA
June 9-14, 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

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.

The Design of the OpenBSD Cryptographic Framework

Angelos D. Keromytis
Columbia University
angelos@cs.columbia.edu

Jason L. Wright
OpenBSD Project
jason@openbsd.org

Theo de Raadt
OpenBSD Project
deraadt@openbsd.org

Abstract

Cryptographic transformations are a fundamental building block in many security applications and protocols. To improve performance, several vendors market hardware accelerator cards. However, until now no operating system provided a mechanism that allowed both *uniform* and *efficient* use of this new type of resource.

We present the OpenBSD Cryptographic Framework (OCF), a service virtualization layer implemented inside the kernel, that provides uniform access to accelerator functionality by hiding card-specific details behind a carefully-designed API. We evaluate the impact of the OCF in a variety of benchmarks, measuring overall system performance, application throughput and latency, and aggregate throughput when multiple applications make use of it.

We conclude that the OCF is extremely efficient in utilizing cryptographic accelerator functionality, attaining 95% of the theoretical peak device performance, and over 800 Mbit/sec aggregate throughput using 3DES. We believe that this validates our decision to opt for ease of use by applications and kernel components through a uniform API, and for seamless support for new accelerators. Furthermore, our evaluation points to several bottlenecks in system and operating system design: data copying between user and kernel modes, PCI bus signaling inefficiency, protocols that use small data units, and single-threaded applications. We offer several suggestions for improvements and directions for future work.

1 Introduction

Today's computing systems are used for applications such as electronic commerce, tele-collaboration of various types, and evolving peer-to-peer systems, often containing sensitive information. Security in these systems depends on several mechanisms that utilize cryptographic primitives as a basic building block. Such cryptographic primitives can be very complex [2] because

the design of these systems is intended to impede simple, brute-force, computational attacks. This complexity drives the belief that strong security is *fundamentally* inimical to good performance.

This belief has led to the common predilection to avoid cryptography in favor of performance [22]. However, the foundation for this belief is often software implementation [8] of algorithms intended for efficient hardware implementation. To address this issue, vendors have been marketing hardware cryptographic accelerators that implement several cryptographic algorithms used by security protocols and applications. However, modern operating systems lack the necessary support to provide *efficient* access to such functionality to applications and the operating system itself through a *uniform* API that abstracts away device details. As a result, accelerators are often used directly through libraries linked with applications, typically requiring device-specific knowledge by the applications, and preventing the operating system itself from easily utilizing such hardware.

We present the OpenBSD Cryptographic Framework (OCF), a service virtualization layer implemented inside the kernel, that provides uniform access to accelerator functionality by hiding device-specific details behind a carefully-designed API. The abstraction introduced allows us to easily support new hardware accelerators and enable applications to use any such accelerator without device-specific knowledge. Furthermore, this intermediate layer does not unduly impact performance, as is common when such abstractions are introduced. The OCF has been in use with OpenBSD [5] for over three years and has proven stable and efficient in practice. It offers features such as load-balancing across multiple accelerators, session migration, and algorithm chaining. We describe the changes we made to the OpenBSD kernel and applications to take advantage of the OCF. In previous work [18] we presented a preliminary analysis of the impact of hardware acceleration on network security protocols, without describing the OCF itself in any detail. Here, we evaluate the impact of the OCF in a variety of micro-benchmarks, measuring overall system perfor-

mance, application throughput and latency, and aggregate throughput when multiple applications use the OCF.

Our evaluation shows that, despite its addition in the system as a device/service virtualization layer, the OCF is extremely efficient in utilizing cryptographic accelerator functionality, attaining 95% of the theoretical peak device performance. In another configuration, we were able to achieve a 3DES aggregate throughput of over 800 Mbps, by employing a multi-threaded application and load-balancing across multiple accelerators. Furthermore, use of hardware accelerators can remove contention for the CPU and thus improve overall system responsiveness and performance for unrelated tasks. Our evaluation allowed us to determine that the limiting factor for high-performance cryptography in modern systems is data copying and the PCI bus. Furthermore, small data-buffers should be processed in software if possible, freeing hardware accelerators to handle larger requests that better amortize the system and PCI transaction costs. On the other hand, multi-threading results in increased utilization of the OCF, improving *aggregate* throughput. We make recommendations for future directions in architectural placement of cryptographic functionality, operating system provisions, and application design, and discuss several improvements and promising directions for future work.

The framework has been in use with IPsec since OpenBSD 2.8, although it continues to evolve in response to new requirements. Public-key support and the `/dev/crypto` interface were introduced in a later version. The OCF has also been ported to FreeBSD and NetBSD, and we are working on Windows and Linux versions.

Paper Organization Section 2 discusses related work. Section 3 describes the OCF's design and implementation, while Section 4 discusses its use by various subsystems and applications. In Section 5, we evaluate the framework's performance, and discuss some of the results and potential improvements and future work in Section 6. Section 7 concludes the paper.

2 Related Work

As interest in security is currently in an upswing, recent work has been examining the overall performance impact of security technologies in real systems. Work by Coarfa, *et al.* [4] has focused on the impact of hardware accelerators in the context of TLS web servers using a trace-based methodology, and concludes that there is some opportunity for acceleration, but given the

choice one might prefer a second processor as it also assists with the substantial (and perhaps dominant) non-cryptographic overheads. [18] provides some basic performance characterizations of IPsec as well as other network security protocols, and the impact acceleration has on throughput. The authors conclude that the relative cost of high-grade cryptography is low enough that it should be the default configuration.

There has been a considerable amount of work on the enhancement of system performance through the addition of cryptographic hardware [2]. This early work was characterized by its focus on the hardware accelerator rather than its implications for overall system performance. [24] began examining cryptographic subsystem issues in the context of securing high-speed networks, and observed that the bus-attached cards would be limited by bus-sharing with a network adapter on systems with a single I/O bus. A second issue pointed out in that time frame [20] was the cost of system calls, and a third [21, 23, 7, 11] the cost of buffer copying. These issues are still with us, and continue to require aggressive design to reduce their impacts.

[25] describes an API to cryptographic functions, the main purpose of which is to separate cryptographic libraries from applications, thus allowing independent development. Our service API is similar at a high level, although several differences were dictated by the need to support actual hardware accelerators and allow it to be used efficiently by protocols such as IPsec and SSL, as we discuss in Section 3. Other work includes the Microsoft CryptoAPI [17], GSS-API [16] and IDUP-GSS-API [1], PKCS #11 [14], SSAPI [26], and the CDSA [19]. These are primarily intended for use by applications that also require authentication, authorization, key management and other higher level security services. Our work focuses on low-level cryptographic operations, providing a simple abstraction layer that does not significantly impact performance, compared to a device-specific approach.

[10] describes an open-source cryptographic coprocessor, focusing on protecting keys and other sensitive information from tampering by unauthorized applications. The author extends the *cryptlib* library to communicate with the co-processor. While he discusses several options for hardware acceleration and identifies some potential performance bottlenecks, it is mostly a qualitative analysis. That work is extended in [9], which presents a comprehensive cryptographic security architecture, again focusing primarily on preserving the confidentiality of users' (and applications') cryptographic keys. We are interested in a much simpler problem: how to accelerate cryptographic operations in a general-

purpose operating system using hardware available in the market and with minimal modifications to the kernel, libraries, and applications.

NetBSD uses the `dmove` facility, which provides an interface to hardware-assisted data movers. This can be used to copy data from one location in memory to another, clear a region of memory, fill a region of memory with a pattern, and perform simple operations on multiple regions of memory, such as XOR, without intervention by the CPU.

3 The Cryptographic Framework

The OpenBSD cryptographic framework (OCF) is an asynchronous service virtualization layer inside the kernel, that provides uniform access to hardware cryptographic accelerator cards. The OCF implements two APIs for use by other kernel subsystems, one for use by *consumers* (other kernel subsystems) and another for use by *producers* (crypto-card device drivers). The OCF supports two classes of algorithms: symmetric (*e.g.*, DES, AES, MD5) and asymmetric (*e.g.*, RSA).

Symmetric-algorithm (*e.g.*, DES, AES, MD5, compression algorithms, *etc.*) operations are built around the concept of the *session*, since such algorithms are typically used for bulk-data processing, and we wanted to take advantage of the session-caching features available in many accelerators. Asymmetric algorithms are implemented as individual operations: no session caching is performed. Session creation and teardown are synchronous operations.

The producer API allows a driver to register with the OCF the various algorithms it supports and any other device characteristics (*e.g.*, support for algorithm chaining, built-in random number generation, *etc.*). The device driver also registers four callback functions that the OCF uses to initialize, use, and teardown symmetric-algorithm sessions, and to issue asymmetric-algorithm requests. The drivers can also selectively de-register algorithms, or remove themselves from the OCF (*e.g.*, a PCMCIA card that is ejected). Any sessions using the defunct driver (or algorithm) are migrated to other cards on-demand (*i.e.*, as the next request for that session arrives). Registration and de-registration can occur at any time; typical device drivers do so at system initialization time. Drivers call the `crypto_done()` and `crypto_kdone()` routines to notify the OCF of completed symmetric and asymmetric requests, respectively. A brief description of the API is given in Appendix A.

In addition to any hardware drivers, a software-crypto pseudo-driver registers a number of symmetric-key algorithms when the system boots. The pseudo-driver acts as a last-resort provider of crypto services; any suitable hardware accelerator will be treated preferably. However, the kernel does not implement asymmetric algorithms in software, for performance reasons; we shall see in Section 4.2 how we handle these. Using a generic API for crypto drivers allows us to easily add support for new cards. We briefly discuss these drivers in Section 3.1.

To use the OCF, consumers first create a session with the OCF using `crypto_newsession()`, specifying the algorithm(s) to use, mode of operation (*e.g.*, CBC, HMAC, *etc.*), cryptographic keys, initialization vectors, and number of rounds (for variable-round algorithms). The OCF supports algorithm-chaining, *i.e.*, performing encryption and integrity-protection in one operation. Such combined operations are used by practically all data-transfer security protocols. At session-creation time, the OCF determines which card to use based on its capabilities and creates a session by calling its `newsession` method, provided at device-registration time. When the session is not needed, `crypto_freesession()` frees any allocated resources.

For the actual encryption/decryption, consumers use `crypto_dispatch()`. The arguments to this include the data to be processed, a copy of the parameters used to initialize the session, consumer-provided opaque data, and a callback function. The data can be provided in the form of *mbufs* (linked lists of data buffers, used by the network subsystem to store packets) or as a collection of potentially non-contiguous memory blocks, called *uio*. The case of a single contiguous data buffer is handled as a *uio*. Although *mbufs* are also a special case of *uio*, we added special support to allow for some processing optimizations when using software cryptography. Furthermore, the issuer of a request can specify whether encryption should be done in-place, or the encrypted data must be returned on a separate buffer. Various offsets indicate where to start and end the encryption, where to place the message authentication code (MAC), and where to find the initialization vector (if already present on the buffer) or where to write it on the output buffer.

The request is queued and `crypto_dispatch()` immediately returns to the consumer. The `crypto` kernel thread is periodically invoked by the scheduler and dispatches all pending requests to the appropriate producers. It also handles all completed requests, by calling the specified callback functions. It then returns to sleep, waiting for more requests. As a result of the OpenBSD kernel architecture (common in most non-SMP kernels), the thread

is not preemptable by user processes, although hardware interrupts are still handled. Currently, the thread must operate at a high priority to avoid synchronization problems. When using the software pseudo-driver, this can cause significant latency in application scheduling and in low-priority kernel operations, although the same problem manifested before the migration to OCF, when encryption was done in-band with IPsec packet processing.

Once the request is processed, the crypto thread calls the consumer-supplied callback routine. If an error has occurred, the callback is responsible for any corrective action. Session migration is implemented by re-creating the session using the initial parameters to *crypto_newsession()*, which accompany all requests as we already mentioned. The error **EAGAIN** is indicated to the callback routine, which re-issues the request after recording the new session number to be used so that subsequent requests are correctly routed. Including the initialization data in each request also allows us to easily integrate cards that do not support the concept of session: the driver simply passes all necessary information (data, algorithm descriptions, and keys) to the card with each request. The opaque data are simply passed back to the consumer unmodified by the OCF; they are used to maintain any additional information for the consumer that is relevant to the request. We shall see an example in Section 4.1.

Asymmetric operations are handled similarly, albeit without support for the concept of session. The parameters in this case include an array of parameters, containing the algorithm-specific big-integers.

When multiple producers implement the same algorithms, the OCF can load-balance sessions across them. This is currently implemented by simply keeping track of the number of sessions active on each producer. At session setup, the OCF picks the producer with the smallest number of active sessions. The software pseudo-driver is currently never used in load-balancing. We evaluate the effectiveness of this simple scheme in Section 5.4. We discuss possible future improvements in Section 6.4.

3.1 Device Drivers

The drivers for the various crypto devices must be able to cope with a wide variety of hardware design decisions (and bugs) made by the manufacturers. These drivers register the algorithms supported by the device and export the appropriate callback functions to the OCF.

The *hifn* driver supports the Hifn 7751, 7811, and 7951 chips and contains around 3,000 lines of code and def-

initions. The driver supports the symmetric operations and hashes available on all these chips. Additionally, it supports the random-number generators available on the 7811 and 7951, but does not support the public key unit on the 7951; the latter was clearly designed for SSL server implementations, as it requires a large amount of CPU-intensive initialization which can be precomputed and used repeatedly on a server but not a client. All these chips support copying-through header and trailer data to the destination buffer, and include full support for scatter-gather I/O. Unfortunately, there is no easy way to coalesce interrupts on this chip, which generates one interrupt per operation, resulting in considerable system overhead. Another important detail is that all of the Hifn symmetric crypto chips poll their descriptor rings in main memory for data to process.

The *nofn* driver supports the Hifn 7814, 7851, and 7854 chips (also known as HIPPI1 packet processors). Currently, there is no support for the symmetric unit on these chips. Fitting these into the current framework is not currently done because they are designed to replace almost all of the IPsec processing (IV generation, MAC checking, replay window handling, *etc.*). In the future, we intend to add support for the IPsec unit by adding a combined-class algorithm and checking for this in IPsec. On the other hand, the public-key unit is almost exactly the same as the Hifn 6500 described below.

The *lofn* driver supports the Hifn 6500 chip, which contains a public-key unit and a random-number generator. This chip is essentially a simple big-number arithmetic logic unit (*i.e.*, it is an ALU capable of performing operations on 1024-bit registers). Unlike all of the other chips, the 6500 is not a bus-master (*i.e.*, has no support for DMA); instead, registers exist within its PCI memory-mapped address space. Because of the expense of modular exponentiations, the somewhat higher overhead of writes to these I/O addresses is still small compared to doing the exponentiation in software.

The *ubsec* driver, which supports the Broadcom 5801, 5802, 5805, 5820, 5821, and 5822 chips, consists of slightly less than 3,000 lines of code and definitions. The symmetric-crypto units on all of the chips are very similar, but the 580x series and 582x series require different formatting for the big numbers on the asymmetric unit. These chips support interrupt coalescing by chaining several commands together, and scatter-gather I/O. Unlike Hifn, these chips do not poll main memory.

We have a variety of other device drivers in various stages of completion. We are aware of other and more modern products from a variety of vendors, but many of them are hesitant to give us the information we need.

4 Use of the OCF in OpenBSD

In this section, we discuss how we extended parts of OpenBSD to make use of the OCF services.

4.1 IPsec

The IP Security Architecture [12], as specified by the Internet Engineering Task Force (IETF), is comprised of a set of protocols that provide data integrity, confidentiality, replay protection, and authentication at the network layer. At the lowest level of the IPsec architecture reside the data encryption/authentication protocols, AH and ESP. These are the “wire protocols,” used for encapsulating the IP packets to be protected. They simply provide a format for the encapsulation; the details of the bit layout are not particularly important for the purposes of this paper. Outgoing packets are authenticated, encrypted, and encapsulated just before being transmitted, and incoming packets are decapsulated, verified, and decrypted immediately upon receipt. These protocols are typically implemented inside the kernel, for performance and security reasons.

IPsec was the first consumer of the OCF services. The original implementation of the OpenBSD IPsec was described in [13]. Here, we give a brief overview and then describe the modifications we had to make to it to allow use of the OCF.

In the OpenBSD kernel, IPsec is implemented as a pair of protocols sitting on top of IP. Thus, incoming IPsec packets destined to the local host are processed by the appropriate IPsec protocol through the protocol switch structure used for all protocols (*e.g.*, TCP and UDP). The selection of the appropriate protocol is based on the protocol number in the IP header. The SA needed to process the packet is found in an in-kernel database using information retrieved from the packet itself. Once the packet has been correctly processed (decrypted, integrity-validated, *etc.*), it is re-queued for further processing by the IP module, accompanied by additional information (such as the fact that it was received under a specific SA) for use by higher-level protocols and the socket layer.

Outgoing packets require somewhat different processing. When a packet is handed to the IP module for transmission (in `ip_output()`), a lookup is made in the Security Policy Database (SPD) to determine whether that packet needs to be processed by IPsec. The decision is made based on the source/destination addresses, transport protocol, and port numbers. If IPsec processing is needed, the lookup will also specify what type of

SA(s) to use for IPsec processing of the packet. If no suitable SA exists, the key-management daemon is notified to acquire one. Otherwise, the packet is processed by IPsec and passed to `ip_output()` again for transmission. The packet also carries an indication as to what IPsec processing has already occurred to it, to avoid processing loops.

In the original IPsec implementation, all cryptographic operations were done in-band with packet processing. This meant that a lot of time was spent performing symmetric-key encryption in the kernel. To make use of the OCF, we split the input and output processing paths. For example, let us consider the case where `ip_output()` determines (by consulting the SPD) that a packet must be IPsec-protected. It then calls `ipsec_process_packet()`, which handles all IPsec outbound-packet processing. After handling encapsulation issues, this routine calls the appropriate “wire protocol” output routine. In the ESP protocol processing, the original `esp_output()` routine was broken up in `esp_output()` and `esp_output_cb()`. `esp_output()` does all the data marshaling and ESP header manipulation, constructs a crypto request, passes it to the OCF and simply returns. Execution returns to `ip_output()` with an indication that the operation was successful.

Once the OCF processes the request, it calls `esp_output_cb()`, a pointer to which is included in the request itself. The callback routine completes the ESP protocol processing by checking for any errors in the crypto processing (re-queuing the request if the OCF indicated so), and calls `ipsec_process_done()`, the second part of the original `ipsec_process_packet()` routine. This routine completes IPsec book-keeping, and calls `ip_output()` with the new packet. `ip_output()` will then perform a new SPD lookup (making sure no IPsec loops occur, by examining the list of SAs that have been already applied to the packet). If necessary, the output processing cycle will occur again. Eventually, `ip_output()` will pass the packet to a network driver for actual transmission.

The cases for output AH and IPcomp processing are similar. Input processing is also similar: `ipsec_common_input()` is called by the network scheduler for all IPsec packets received. It locates the appropriate SA in the kernel SA database and calls `esp_input()`. Similar to the output case, `esp_input()` validates the ESP header fields, constructs a crypto request, passes it to the OCF and returns. Once the request is processed, the OCF will call `esp_input_cb()`, which will verify the packet integrity (by comparing the value on the packet with that computed by the accelerator), remove the ESP header, and pass the packet to `ipsec_common_input_cb()`.

This routine performs further sanity and security checks on the decrypted packet, and re-queues it for further processing by the IP layer. AH and IPcomp input processing is similar, as is the case of IPsec over IPv6.

Input ESP and AH processing offer one example of use of the opaque data passed with each crypto request, discussed in Section 3. All the cryptographic accelerators that support message authentication (MAC) algorithms only offer a “forward-compute” mode. That is, the card can only compute the MAC on the packet, and it is up to the operating system to verify its validity by comparing it with the received value. Thus, we use the opaque data to store the MAC value from the packet and instruct the OCF to write the new MAC value in the appropriate location in the packet — the operation is exactly the same as the output case. In the callbacks, we simply do a byte-wise comparison of the computed value (stored on the packet) and the received value (stored as opaque data in the request itself).

While the code was not very complicated, there were several minor headaches as a result of this asynchronous processing model. For example, one problem was communicating MTU information through arbitrarily-many IPsec SAs to the TCP layer, so as to correctly fragment application data and avoid fragmentation at the IP layer. We could not simply update the appropriate data structures with the correct MTU value after the packet had been encapsulated once, since we could not “peek” inside the encryption. Fortunately, we keep a record of which SAs have been applied to a packet during input and output processing. Thus, on receipt of the appropriate ICMP message, or when the IP layer indicates that the packet is too large to be transmitted without fragmentation, the list of SAs is traversed and each SA is updated with the correct MTU value based on its position in the SA chain (*i.e.*, the first SA on output will advertise a smaller MTU than the last one, the difference being the ESP headers and encryption padding). The next packet that tries to traverse the chain will encounter a correct MTU value.

4.2 /dev/crypto

Building on our experience with the IPsec implementation, we turn our attention to exporting the OCF services to user-level applications. A /dev/crypto device driver exists which abstracts all the OCF functionality and provides a command set that can be used by OpenSSL (or other software that uses /dev/crypto directly). This interface is based on *ioctl()* calls and is thus fully synchronous (*i.e.*, applications can only have one request pending) — in the future, we intend to al-

low processes to issue multiple requests. Both symmetric and asymmetric operations are permitted using this framework; we will first describe the symmetric component.

Similar to the underlying OCF, this uses a session-based model, since the general case assumes that keys will be reused for a sequence of operations. After opening the /dev/crypto device and gaining a file descriptor *fd*, the caller requests that a new session be created with *CIOCGSESSION* for a certain cryptographic operation, and specifies all related parameters (*e.g.*, keys). Similar to the OCF, a single session supports both a cipher and a MAC, as we are simply exporting the same functionality available to the kernel. *CIOCGSESSION* returns a session identifier that can then be reused repeatedly for subsequent operations. When the session is no longer needed, it can be revoked using *CIOCFSESSION*. Many sessions can be requested against a single file descriptor *fd*; all sessions follow a particular *fd* through *fork()* and *exec()* calls, and are not otherwise visible to other processes. Obviously, the last *close()* on *fd* destroys all the sessions.

If the request cannot be satisfied using hardware accelerators, the kernel will return an error of *EINVAL*, so the caller can fall back to a software implementation. We considered adding an *ioctl()* that describes the abilities of the available hardware, allowing an application to determine if the needed algorithm is supported by looking at a list. However, numerous other variables exist (key sizes, block sizes, alignment) which might be difficult to describe. For the time being, we have punted on this issue. However, when first called, the OpenSSL engine will enumerate all OCF-supported algorithms. It does so by issuing a *CIOCGSESSION* request for each algorithm it supports in software, and caches the result. If an algorithm is not provided by the OCF, the library will use its software implementation (in reality, the kernel will admit that it supports cryptographic algorithms that it implements in software, and OpenSSL will make use of them as if they were implemented by hardware, unless a *sysctl* variable is set to prohibit this, which is the default setting).

Once a session is established, blocks can be encrypted or decrypted using the *CIOCCRYPT ioctl()*. Each time this is used, the caller can specify a new IV or MAC information that they wish to fold into the operation. Input and output buffers are specified via separate pointers, but they can point to the same buffer for in-place encryption. Naturally, the data size provided by the caller must be rounded to the default block size of the algorithm being used. A data size limit of 262,140 bytes exists at the moment, to hide a similar limit found in some chipsets.

In the future, we may support larger blocks by splitting operations into smaller chunks.

The user-land data blocks are copied into memory allocated inside the kernel address space. This data is formatted into *uio* blocks as mentioned in Section 3. The OCF is then called to perform the operation using the initialization information stored in the application's `/dev/crypto` session. If the operation is successful, the results are copied back to the application buffers. Obviously, the cost of these two copies is higher for larger block sizes, as we shall see in Section 5.4. In the future, we hope to use page flipping for larger blocks when the kernel memory subsystem supports this.

For asymmetric operations, no session is required. The `CLOCKEY ioctl()` is used in an atomic fashion for each individual operation. Five operations are provided, with `CRK_MOD_EXP` being the most important. Support for the others, `CRK_MOD_EXP_CRT`, `CRK_DSA_SIGN`, `CRK_DSA_VERIFY`, and `CRK_DH_COMPUTE_KEY` has not yet been completed. Each of these has an operation-specific number of input and output parameters, which are always a packed byte array of big integers. The particular format we chose for these parameters makes it easy to interface to OpenSSL “bignums,” and to most of the early hardware we had access to.

Presently, OpenBSD lacks cloning devices. Therefore a cumbersome procedure for opening `/dev/crypto` must be followed. After the initial `open()` call, the caller must use `ioctl()` to retrieve a file descriptor (*fd*) to use, then perform all operations against this replacement *fd*. This replacement *fd* is a unique per-process descriptor, while the initially-opened one would naturally be shared between all callers. Without such semantics, the `fork()` and `_exit()` system calls do not exhibit the expected semantics with respect to file-descriptor inheritance and closing. Just as bad, we would end up with all processes able to see and use each other's keys. When cloning devices are implemented in OpenBSD, we will change the user-level code (mostly OpenSSL) to no longer use this complicated procedure, but the kernel will retain it for backward compatibility. While writing this code, we ran into numerous strange and difficult resource-management issues for session teardown.

It should also be noted that applications using `/dev/crypto` must ensure they use `ioctl()` with the `F_SETFD` command on the crypto descriptor to ensure that the “close-on-exec” flag is set. Otherwise, child processes will inherit unwanted descriptors, which is both a security and a resource-exhaustion concern.

4.2.1 OpenSSL Enhancements

In the past, programmers using OpenSSL (or its predecessor, SSLeay) directly called the generic crypto routines as they existed for each algorithm. More recently, programmers have been encouraged to use the *EVP* layer for dealing with symmetric algorithms. This provides a session-based model much like the `/dev/crypto` layer described in the previous section. Applications like OpenSSH, `mod_ssl` (the Apache SSL module we use), and `sendmail` have matured to use these interfaces.

Newer OpenSSL code-bases contain an “engine” component. This allows asymmetric algorithms to be directed to a hardware driver; a number of stub functions are provided which typically interface with vendor-specific shared libraries to actually do the operation on the vendor's accelerator. Many of these subsystems interact badly and do not consider the effects of `chroot()` or other strange Unix behaviors, resulting in weak security models. Since we run Apache in a `chroot()`'ed environment in which there exists no `/dev/crypto` device, we modified it to perform all necessary initializations prior to being sandboxed. We wrote our own engine modules that interacts directly with `/dev/crypto`, without any of these surprises. Symmetric operations from the *EVP* layer are directly mapped into OCF requests. One major weakness is that the *EVP* layer has no concept of bundling algorithms. Thus, protocols that use encryption and MAC on a message, such as TLS and SSH version 2, sequentially issue two separate requests to `/dev/crypto` through the *EVP* layer, resulting in unnecessary context switches, data copying, and DMA transactions. Thus, the *EVP* layer currently does not pass MAC operations to the OCF.

5 Performance Evaluation

In this section, we analyze the performance of the cryptographic framework. We have ran a series of micro-benchmarks that allowed us to determine the limits of the framework and potential directions for improvement. We use the OCF for simple cryptographic tasks, comparing different cryptographic accelerators with the case of pure-software encryption, and provide a cost breakdown. We also attempt to quantify the benefits to be had by the system at large, when off-loading cryptographic operations to hardware accelerators. Finally, we evaluate the load-balancing feature of OCF, by simultaneously using multiple accelerators on the same machine.

5.1 Testbed

For our tests, we use two identical machines. The machines have 1.4 Ghz Pentium III processors on Tyan Thunder HESL-T motherboards. These motherboards have three independent PCI busses: 32bit/33Mhz/5V, 64bit/66Mhz/5V, and 64bit/66Mhz/3.3V. The boards use 512MB of 133Mhz registered SDRAM and are based on the ServerWorks HESL chipset. We placed the crypto card being tested either on the 64bit/66mhz/3.3V bus or the 32bit/33Mhz/5V bus, as appropriate for the card. The crypto cards we used are:

- Broadcom 5805 reference design board (32bit).
- Broadcom 5820 reference design board (64bit).
- GTGI XL-Crypt (based on the Hifn 7811 chip) (32bit).
- NETSEC 7751 (based on the Hifn 7751 chip) (32bit).
- Hifn 6500 reference design board (32bit).
- Hifn 7814 reference design board (64bit).

The Hifn data-sheet gives a peak performance for the 7751 chip of 62 Mbps for encryption and 110 Mbps decryption, when using IPsec with 3DES/SHA1/LZS (LZS is a data-compression algorithm). When the 3DES engine alone is used, both encryption and decryption throughput are 83 Mbps. Broadcom's web site places the peak performance of the 5820 chip at 310 Mbps of 3DES-SHA1, when used in IPsec. Furthermore, they claim 800 1024-bit RSA signature computations per second.

5.2 OCF Throughput

To determine the raw performance of OCF, we use a single-threaded program that repeatedly encrypts and decrypts a fixed amount of data with various symmetric-key algorithms, using the `/dev/crypto` interface. We run the test against all the hardware accelerators listed in the previous section, as well as using the kernel-resident software implementation of the algorithms. We vary the amount of data to be processed per request across experiments. To measure the overhead of OCF without the cryptographic algorithms, we added to the kernel a *null* algorithm that simply returns the data to the caller without performing any processing. The results can be seen in Figure 1.

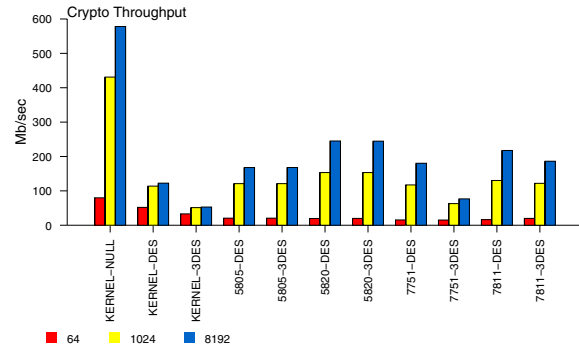


Figure 1: **Crypto-hardware performance.** The *KERNEL-NULL* bar indicates use of the *null* encryption algorithm. The *KERNEL-DES* and *KERNEL-3DES* bars indicate use of the software DES and 3DES implementations in the kernel. The remaining bars indicate use of the various hardware accelerators. The vertical axis unit is Mbits/second.

We can make several observations on this graph. First, even when no actual crypto is done, the ceiling of the throughput is surprisingly low for small-size operations (64 bytes). In this case, the measured cost consists of the overhead of system call invocation, argument validation, and crypto-thread scheduling. As larger buffers are passed to the kernel, the throughput increases dramatically, despite the increasing cost of memory-copying larger buffers in and out of the kernel. When we use 1024-byte buffers, performance in the no-encryption case jumps to 420 Mbps; for 8192-byte buffers, the framework peaks at about 600 Mbps.

Notice however that this peak corresponds to a single process issuing crypto requests. This process is blocked after each request, the scheduler context-switches to the crypto thread (which was blocked waiting for requests), the *null* algorithm executes and the completed request is passed back to the `/dev/crypto` driver, which wakes up the blocked user-level process. If many processes are issuing requests, the crypto thread's request queue will contain multiple requests. When we run multiple processes, each will queue a request (and be blocked by `/dev/crypto`); the crypto thread will process all these requests in a flurry of activity, and cause all processes to wake up in synchrony. The crypto thread will then go back to sleep, while each of the processes will issue another request. This cycle repeats for the duration of the experiment. As a result, more processes using the OCF result in increased aggregate throughput, simultaneously increasing the average processing latency.

These buffer sizes are close to the typical sizes of re-

quests issued by some of the most-commonly used applications:

- *SSH* keyboard input results in many small requests (so we are close to the 64-byte case); responses from the server are larger, but not considerably so. When X forwarding is used, we can occasionally get larger buffers.
- *SCP/SFTP* issue larger requests; OpenSSH, a popular implementation, uses requests of 4 KB.
- *SSL/TLS* also issue large requests. The maximum size of an SSL record is 16 KB, but can be less if (optional) compression is used.
- *IPsec* processes packets at the network layer. Such traffic is trimodal [3]: about 40% of packets are 40–60 bytes (the vast majority of these being TCP acknowledgments), with the remainder split between 576 bytes (TCP MSS when no Path MTU Discovery is used) and 1460 bytes (when Path MTU Discovery is used).

When we use real cryptographic algorithms, we notice that the performance of DES done in software is close to that of no encryption for small packet sizes; even 3DES performance is just half of the no-encryption case. If we use larger buffer sizes, the performance of software crypto done in the kernel (the `KERNEL-*` labeled bars) degrades rapidly. When we use hardware accelerators, we notice two different trends. For small buffers, the performance degrades with respect to the software case. This indicates that the additive costs of system call invocation, OCF processing, and the 2 PCI transactions (to/from the crypto cards) dominate the cost of doing crypto. However, as we move to larger buffer sizes, performance quickly improves as these overheads are amortized over larger buffers, despite the fact that more data has to be copied in and out of the kernel and over the PCI bus. Thus, to improve the performance of the system when applications issue large numbers of small requests, either request-batching should be done, a faster processor should be used, or the number of user/kernel crossings should be minimized. When larger buffers are being processed, it pays off to use some cryptographic accelerators, although not all such cards are equal in terms of performance.

Notice that the performance of DES and 3DES is the same in each of the 5805 and 5820 cards; these cards really implement only 3DES in Encrypt-Decrypt-Encrypt (EDE) mode, and emulate DES by loading the same key in one of the Encrypt and the Decrypt engines (effectively canceling each other out). In contrast, the 7751

seems to implement two separate crypto engines for DES and 3DES, or uses a shortcut in its 3DES engine. The 7811 seems to implement different engines as well, but the performance difference between the two is not as pronounced.

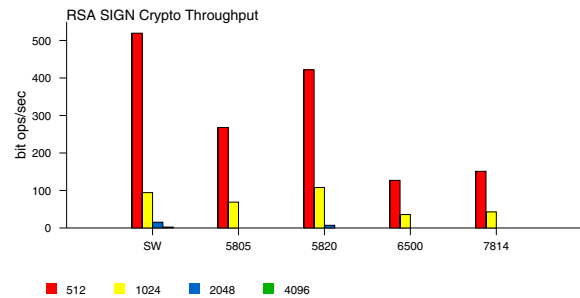


Figure 2: RSA signature generation. The horizontal axis indicates the modulus size, in bits. The vertical axis indicates number of operations per second.

Similarly, we measure the performance of OCF for public-key operations. In this case, there are no kernel-resident software public-key algorithms. We count the number of RSA signature generations and verifications per second, for different accelerators and key sizes (512 to 4096 bits, as supported by the each cards). The results are shown in Figures 2 and 3.

The Hifn 6500 and 7814 are geared more towards slower, embedded applications, so the fact that their performance is considerably worse than software is not surprising. The number of verifications is much larger than the number of signature generations in unit time. This is because, as with most crypto libraries, OpenSSL opts for small values for the public part of the RSA key (typically, $2^{16} + 1$) and correspondingly large values for the private key. This causes the public-key operations (encryption and verification) to be much faster than the private-key operations, even though they are in principle the same operation (modular exponentiation).

Another interesting observation is that the RSA sign throughput is higher in the software case (see Figure 2). This happens because the CPU on the crypto-card is slower than the host CPU and optimized for bit operations, which is as useful for public key cryptography. So the “anomaly” in Figure 2 is actually expected. However, as we mentioned in Section 5.1, Broadcom claims that the 5820 can perform 800 RSA signature operations per second with 1024-bit keys. In our case, we only see slightly over 100. There are two explanations for this. First, we are under-utilizing the 5820: there is only one thread issuing RSA sign operations, which is blocked waiting termination of each request. Once the card com-

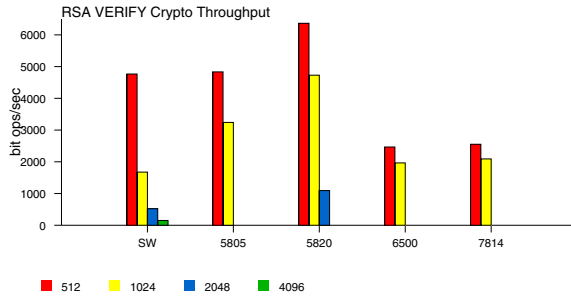


Figure 3: RSA signature verification. The horizontal axis indicates the modulus size, in bits. The vertical axis indicates operations per second.

putes the signature, it has to wait for the crypto framework to wake up the blocked process, then the scheduler to context-switch to it, the process to issue an *ioctl()* call to get the results, and then another *ioctl()* call to issue the next request, which is placed on the crypto thread's queue. Finally, the scheduler has to context-switch to the crypto thread. During all this time, the accelerator is idle, since there is no other process using it. The second reason for the higher vendor-stated performance is that the tests they performed used the CRT parameters for the RSA operations, which make RSA processing considerably faster. However, for implementation reasons, our OpenSSL engine does not use CRT parameters yet.

5.3 System-wide Effects

To determine the system-wide benefits of offloading cryptographic processing, we run multiple threads (up to 24) of the `openssl speed` benchmark with various algorithms, while at the same time we run a simple CPU-intensive job. The CPU "hog" process consists of a small program that performs 2^{32} function calls, each function call performing an integer-multiply operation. The elapsed time for the CPU hog process was recorded for each (*algorithm, number of threads*) tuple. As we see in Figure 4, the crypto accelerators very effectively eliminate contention for the otherwise-shared resource, the CPU, whether the crypto performed is symmetric (DES, 3DES) or asymmetric (DSA with 1024-bit keys). The execution time for the hog process remains constant, regardless of the number of threads of execution.

5.4 Load Balancing

Finally, we wish to determine how well the OCF can load-balance crypto requests when multiple accelerators

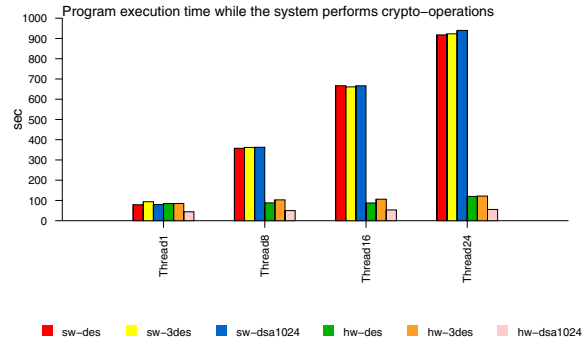


Figure 4: Program execution time while multiple threads perform crypto-operations in parallel. The bars show the elapsed time in seconds for executing the CPU-bound process for different algorithms and numbers of threads.

are available, and the aggregate throughput that can be achieved in that scenario. We use a custom-made card by Avaya that contains four Hifn 7751 chips that can be used as different devices through a PCI bridge resident on the card. We use multiple threads that issue encryption requests for 3DES, and vary the buffer size across different runs. The results are shown in Table 1. As we can see, performance peaks in the case of 32 threads and 16 KB buffers at 320 Mbps, which is over 96% of the maximum rated throughput of four Hifn 7751 chips. The card was installed on the 64bit/66Mhz PCI bus, but because the chip is a 32bit/33Mhz device, the maximum bus transfer rate is 1.056 Gbps. At our peak rate, we use over 640 Mbps of the bus: 320 Mbps for data in each direction (to and from the card), plus the transfer initialization commands and descriptor ring probing, *etc.*, thus utilizing over 60% of the PCI bus. Notice that because the card uses a PCI bridge, a 2-cycle latency is added on each PCI transaction.

The card was installed on the 64bit/66Mhz bus because the system's 32bit/33Mhz bus exhibited surprisingly bad performance, probably because many other system components are found on that bus and likely cause contention: since the machine is operating as it normally would while this test is being run, the scheduler is active, and two clock interrupts are being received at 100 and 128 Hz respectively. Other devices are also generating their own interrupts.

Another possible cause is an artifact of the i386 *spl* protection method: a regular *spl* subsystem disables the interrupts from a certain class of devices at the invocation of an *splX()* call. For instance, calling *splbio()* blocks reception of interrupts from all devices which are in the "bio" class of devices. On the i386, the registers used

Number of threads	16 bytes	64 bytes	256 bytes	1024 bytes	8192 bytes	16384 bytes
1	3.06 Mbps	11.45 Mbps	33.15 Mbps	59.49 Mbps	79.19 Mbps	80.75 Mbps
2	5.53 Mbps	18.40 Mbps	56.07 Mbps	111.60 Mbps	154.18 Mbps	160.02 Mbps
3	6.44 Mbps	23.25 Mbps	71.31 Mbps	152.28 Mbps	229.60 Mbps	238.24 Mbps
4	6.83 Mbps	25.77 Mbps	80.91 Mbps	182.65 Mbps	292.15 Mbps	299.33 Mbps
32	7.37 Mbps	27.51 Mbps	94.05 Mbps	249.17 Mbps	313.79 Mbps	320.19 Mbps

Table 1: Crypto-request load-balancing using a quad-Hifn 7751 card on a PCI 64bit/66Mhz bus.

Number of threads	16 bytes	64 bytes	256 bytes	1024 bytes	8192 bytes	16384 bytes
1	5.42 Mbps	18.88 Mbps	61.94 Mbps	151.95 Mbps	300.88 Mbps	254.79 Mbps
32	9.91 Mbps	37.01 Mbps	120.71 Mbps	410.27 Mbps	758.85 Mbps	801.81 Mbps

Table 2: Crypto-request load-balancing using four 5820 cards on a PCI 64bit/66Mhz bus.

to do interrupt blocking (found on the programmable interrupt controller, also known as the PIC) are located on the 8Mhz ISA bus, which is what OpenBSD uses for interrupt management (as opposed to the APIC).

Worse yet, some operations on this device require a 1 *usec* delay before taking effect. To partially mitigate this extremely high overhead, the i386 kernel interrupt model instead makes the vectors for blocked interrupt routines point to a single-depth queuing function which does the actual interrupt blocking at the time of reception. When the *spl* is lowered again, the original interrupt handler is called. However, the 8Mhz ISA bus still had to be accessed. This has the effect of further reducing the available bandwidth on the PCI bus. One small-buffer benchmark generated over 62,000 interrupts/sec; we believe that the *spl* optimization is failing under such load.

Using four 5820 cards on a 64bit/66Mhz PCI bus allows us to achieve even higher throughput, as shown in Table 2. We show only the 1 and 32-thread tests; the rest of the measurements followed a similar curve as the quad-7751. Performance peaked at over 800 Mbps of crypto throughput. Using the same analysis as before, we are using in excess of 1.6 Gbps of the fast-PCI bus, which has a throughput of 4.22 Gbps, achieving slightly over 38% utilization of the bus. As we mentioned in Section 5.1, the vendor rates this card at 310 Mbps. Thus, the maximum theoretical attainable rate would be 1.24 Gbps. We achieve 64.5% utilization of the four cards in this case. A rough sampling of CPU utilization during these large block benchmarks on both cards showed around 10,000 interrupts/second, which is substantial for a PC.

Investigating further, we determined that all four 5820 cards were sharing *irq* 11. Thus, it is possible that

the culprit is the *spl* optimization previously mentioned, at least for the small buffer sizes: the *vmstat* utility shows us anything from 50,000 to 60,000 interrupts per second when processing buffers of 16 to 1024 bytes. Furthermore, because of a quirk in the processing of shared *irq* handlers, some cards experience slightly worse interrupt-service latency: shared *irq* handlers are placed in a linked list; if multiple cards raise the interrupt at the same time, the list will be traversed from the beginning for each interrupt raised — and each *irq* handler will poll the corresponding card to determine if the interrupt was issued by it. However, fixing this quirk or moving the cards on different *irq*'s did not significantly improve throughput.

When we use 8192-byte buffers, the interrupt count drops to 12,000, which the system can handle. In each of these cases, the system spends approximately 65% of its time inside the kernel. Most of this cost can be attributed to data copying. However, as we move to larger buffer sizes, we find the system spending 89% of its time in the kernel, and only 1.9% in user applications, for the case of 16 KB buffers. The number of interrupts in this case is only 5,600, which the system can easily handle. The problem here is that there is considerable data copyin/copyout between the kernel and the applications; aggravating the situation, while such data copying is in progress no other thread can execute, causing a “convoy” effect: while the kernel is copying a 16 KB buffer to the application buffer, interrupts arrive that cause more completed requests to be placed on the crypto thread’s “completed” queue. The system will not allow the applications to run again before all completed requests are handled, which cause more data copying. Thus, the queue will almost drain before applications will be able to issue requests again and refill it. We intend to further investigate this phenomenon.

Fundamentally, the data copyin/copyout limitation is inherent in the memory subsystem. We measured its write-bandwidth to be approximately 2.4 Gbps. Using the crypto cards, we are in fact doing 3 memory-write operations for each data buffer: one copyin to the kernel, one DMA from the card to main memory, and one copyout to the application. Notice that data DMA'ed in from the card is not resident in the CPU cache, as all such data is considered "suspect" for caching purposes. In addition, there is an equal amount of memory reads (copyin, DMA in from the card, copyout). Each of those transfers represents an aggregate of 800 Mbps. When we ran the same test with three 5820 cards, performance slightly improved to 841.7 Mbps in the case of 16 KB buffers, achieving over 90% utilization of the three cards. In this case, the memory subsystem is still saturated, but the cards can more easily get a PCI-bus grant and perform the DMA.

6 Discussion

6.1 Cryptography in the Kernel

As we saw in the previous section, the influence of multi-threading on performance is strong, which suggests that busy servers can make better use of hardware cryptography than clients. This supports the observations of Dean, *et al.* [6] that it may make sense to make cryptography a shared network service to achieve the best cost/performance in a secure system. Notice that, within the boundaries of one host (operating system instance), this is precisely what the OCF does. We should also mention that use of a threaded model for applications involves an obvious security *vs.* implementation complexity trade-off.

Although the performance of individual applications may not improve drastically when using an accelerator, it appears that the *aggregate* performance of a number of applications (as may be the case in a system with many remote login sessions, a busy web server, or a VPN gateway) does improve, as a result of increased utilization. Furthermore, hardware accelerators can give a performance boost to the rest of the system, as was seen in Figure 4. Very simply, they eliminate contention for the CPU, which is a resource shared by all applications and the operating system itself. Thus, while throughput is not drastically improved (and may in fact degrade in certain scenarios) with use of hardware acceleration, overall system utilization improves because the main CPU is left to perform other tasks.

6.2 System Architecture

As we saw in Section 5.4, data copying and the PCI bus quickly become the limiting factor. In practice, the situation is even worse since cryptography is used in conjunction with either network security protocols, in which case the network interface card (NIC) contents for a slice of the PCI bandwidth, or with filesystem encryption, in which case the storage device claims a portion of the bus. This situation suggests that, for maximum performance, cryptographic support must be provided by the individual devices (*e.g.*, NICs, disk controllers, *etc.*). Alternatively, cryptographic support must be located elsewhere in the system architecture (*e.g.*, attached to the main CPU, the system "north bridge" (as the video subsystem is), or the memory subsystem. Any of these approaches, if implemented correctly, will improve application performance by reducing contention for the PCI bus, but at the same time will create new challenges for operating systems that have to support these new devices, such as session migration and fail-over (which the OCF supports by design, as we discussed in Section 3).

Although the OCF does not directly take advantage of NICs that support IPsec-processing offloading, since they are not general-purpose cryptographic accelerators, we have extended the IPsec stack to use them. The cards of this type we are familiar with are 100 Mbps full-duplex Ethernet, and it seems reasonable to assume that they can achieve that performance, given our results with dedicated cryptographic processors. Unfortunately, at the time this paper was written, we did not have enough information to write a device driver that could take advantage of such features. We are also not aware of any commercially-available hard drive controllers that provide built-in encryption services.

6.3 The Effect of Small Requests

The nature of the challenge for operating systems and their support for cryptography is clear. On every measurement, without exception, small-sized operations fare much worse than those performed on large data buffers. In some cases, buffer size influences performance more than the choice between hardware or software cryptography. This suggests that the per-operation overhead is very high, and this is clear from the larger data sizes, which get close to the throughput advertised by the board manufacturer, which we presume is "best-case". In this respect, our findings confirm those of [15]. Since many cryptographic protocols are transactional in nature rather than bulk transfers, these small data operations will be the common case. Energy should be spent on

reducing the overhead of such cases.

As we mentioned in Section 5.2, there are several possible approaches: request-batching, kernel crossing and/or PCI transaction minimization, or simply use of a faster processor. These are more cost-effective solutions than deploying a hardware accelerator. In situations where bulk data transfer is the norm (as may be the case in the various Storage Area Network technologies currently under consideration), cryptographic accelerators can drastically improve performance, especially for the more “expensive” algorithms such as 3DES. Unfortunately, there were no commercially available hardware accelerators for AES supported by OpenBSD, so we cannot compare the software and hardware cases for that algorithm. However, recent attacks against AES make likely the continued use of 3DES in many environments.

6.4 Other Optimizations and Future Work

Smarter load balancing. The load-balancing currently done in OCF, as discussed in Section 3, is very simple. It performs load-balancing of sessions, by keeping a record of the active sessions per producer and selecting the least-loaded one. However, not all sessions are equivalent in terms of processing requirements: an FTP-over-IPsec session will use the OCF more heavily than a telnet-over-IPsec one. Furthermore, the current scheme does not perform load-balancing for public-key operations. Finally, all producers of crypto services are considered equal, in terms of performance. All these issues point to several potential improvements that can be made to the OCF.

For example, drivers can state their peak performance (experimentally measured, using the vendor-provided numbers, or measured at system boot time), and the OCF can keep a record of the number of operations actively pending on each driver. However, this requires sessions to be simultaneously established on all these cards; as these cards have a limited amount of memory for session caching, this approach is perhaps not optimal for a very busy system. One potential solution is to allow the OCF to do dynamic load-balancing of sessions, replicating and tearing them down on additional cards based on their measured traffic, by maintaining session information internally. Asymmetric operations are easier to load balance, as they do not depend on the concept of the session. An additional benefit of implementing load-balancing in this way is that we can let the software driver handle small requests, reducing latency, and use the hardware producers for larger requests. One complication to this is that many cards (*e.g.*, Hifn) do not export internal state such as IVs or intermediate MAC results,

which makes such session sharing difficult.

Algorithm-chaining across cards. It is possible that an OCF consumer needs to chain together a number of cryptographic algorithms, but no hardware producer implements all these. Currently, this would cause the session to be established on the software pseudo-driver (which implements all algorithms). However, by maintaining session information inside the OCF, it is possible to create “virtual sessions” across multiple (hardware and software) producers. In this case, the OCF will issue multiple sequential requests to the various producers, invoking the consumer-specified callback routine at the end. We have a prototype of this, but we need to further evaluate the performance implications and trade-offs of doing multiple PCI transactions.

Asymmetric Multiprocessing (AMP) support. There is an increasing number of multi-processor systems. Most of these under-utilize the secondary processor, as many modern tasks are I/O-limited. Furthermore, it seems likely that the first version of SMP support for OpenBSD will be very coarse-grained: only one processor (and process) can be inside the kernel at a time. An alternative approach is to designate the secondary processor as a dedicated cryptographic accelerator that registers with the OCF as such. No special support by the OCF is necessary, and we are currently working toward an implementation of this.

OpenSSL support algorithm-chaining with OCF. As we mentioned in Section 4.2, TLS and SSH use the OCF at the granularity of the algorithm. That is, if both an encryption and a message authentication (MAC) algorithm have to be applied on an outgoing message, there will be two distinct calls to the OCF via */dev/crypto*. (The same situation holds for incoming messages.) Since the OCF supports algorithm chaining, there is no reason why OpenSSL cannot take advantage of this to reduce the number of user/kernel crossings. This requires modification of the TLS implementation in OpenSSL and of OpenSSH, to support this algorithm chaining. While this is purely an implementation matter, the complexity of the OpenSSL code is a significant deterrent to progress in this direction.

Minimize number user/kernel crossings and data copying. In most practical uses of the OCF (especially in protocols like TLS or SSH), an application issues one or more crypto requests via */dev/crypto*, followed by a *write()* or *send()* call to transmit the data. Similarly, a

`read()` or `recv()` call is followed by a number of requests to `/dev/crypto`. This implies considerable data copying to and from the kernel, and potentially unnecessary context switching back and forth. An alternative approach is to “link” some crypto context to a socket or file descriptor (when doing application-level file encryption), such that data sent or received on that file descriptor are processed appropriately by the kernel: for example, a TLS implementation might construct a data record and simply `write()` it to the socket (one data copy and kernel crossing), only to have the kernel pass it to the OCF for processing before actually passing it on to TCP for transmission. This requires some discipline by the application, which must set the state on the socket and only `write()` appropriately-formatted record, as well as some support in the kernel to decode incoming TLS or SSH frames for processing by the OCF before passing them on to the application.

Another potential approach is to do “page sharing” of data buffers; when a request is given to `/dev/crypto`, the kernel removes the page from the process’s address space and maps it in its own. When the request is done, the kernel re-maps the page back to the process’s address space, avoiding all data copying. This works well as long as `/dev/crypto` remains a synchronous interface. If processes are allowed to have multiple pending requests, accesses to that page while it is being shared with the kernel must be caught and handled, similar to the way copy-on-write of memory pages is handled. An alternative is to block any process that tries to access such pinned-down pages until the crypto request is completed. Obviously, pages that are shared between processes can cause similar problems even in the current mode of operation. Operations that cross page boundaries also have to be dealt carefully.

7 Conclusions

We presented the OpenBSD Cryptographic Framework (OCF), a service virtualization layer implemented inside the kernel, that provides uniform access to cryptographic hardware accelerator cards by hiding card-specific details behind a carefully designed API. Other kernel subsystems and user-level processes can use the API with symmetric and asymmetric algorithms. The OCF offers several other features, such as load-balancing, session migration, and algorithm-chaining.

Our performance evaluation demonstrated the OCF’s ability to utilize available accelerators to within 95% of their peak performance. This validates our decision to design for ease of use by applications and seamless

support for new accelerators, over a device-specific approach which should be able to fully utilize that device’s capabilities. In addition, we demonstrated aggregate (across several concurrent applications) throughput for 3DES encryption in excess of 800 Mbps. Furthermore, use of hardware accelerators can remove contention for the CPU and thus improve overall system responsiveness and performance for unrelated tasks.

Our evaluation also allowed us to determine that the limiting factor for high-speed cryptography in modern systems is data copying and the PCI bus. Furthermore, small data-buffers should be processed in software, freeing hardware accelerators to handle larger requests that better amortize the system and PCI transaction costs. On the other hand, multi-threading results on increased utilization of the OCF, improving *aggregate* throughput. We made recommendations for future directions in architectural placement of cryptographic functionality, operating system provisions, and application design, and discussed several improvements and promising directions for future work.

Acknowledgements

Bob Beck and Markus Friedl helped with numerous OpenSSL integration issues we faced, since the “engine” code we required was unreleased. Bob also wrote the first working OpenSSL engine interfacing with `/dev/crypto`. Markus helped with regression tests to ensure that `/dev/crypto` operation was correct. Jonathan Smith and Sotiris Ioannidis provided valuable comments and insights. Sam Leffler adapted the OCF to the FreeBSD kernel. We would also like to thank Patrick McDaniel for providing high-quality shepherding of this paper.

We are grateful to Global Technologies Group, Inc. (GTGI) for providing us with two XL-Crypt (Hifn 7811) boards, one Hifn 6500 reference board, and one Hifn 7814 reference board. We are also grateful to Network Security Technologies, Inc. (NETSEC) for providing us with two Hifn 7751 boards, one Broadcom 5820 board, and two Broadcom 5805 boards. In addition, NETSEC funded part of the original development of the device-support software.

Part of this work was supported by DARPA and the Air Force Research Laboratory, Air Force Material Command, USAF, under agreement number F30602-01-2-0537.

References

- [1] C. Adams. Independent Data Unit Protection Generic Security Service Application Program Interface (IDUP-GSS-API). RFC 2479, December 1998.
- [2] A. G. Broscius and J. M. Smith. Exploiting Parallelism in Hardware Implementation of the DES. In *Proceedings of the Crypto Conference*, pages 367–376, August 1991.
- [3] K. Claffy, G. Miller, and K. Thompson. The nature of the beast: Recent traffic measurements from an Internet backbone. In *Proceedings of the ISOC INET Conference*, July 1998.
- [4] C. Coarfa, P. Druschel, and D. Wallach. Performance Analysis of TLS Web Servers. In *Proceedings of the Network and Distributed Systems Security Symposium (NDSS)*, February 2002.
- [5] T. de Raadt, N. Hallqvist, A. Grabowski, A. D. Keromytis, and N. Provos. Cryptography in OpenBSD: An Overview. In *Proceedings of the USENIX Annual Technical Conference, Freenix Track*, pages 93 – 101, June 1999.
- [6] D. Dean, T. Berson, M. Franklin, D. Smetters, and M. Spreitzer. Cryptology as a Network Service. In *Proceedings of the Network and Distributed System Security Symposium (NDSS)*, February 2001.
- [7] P. Druschel, M. B. Abbott, M. A. Pagels, and L. L. Peterson. Network subsystem design. *IEEE Network*, 7(4):8–17, July 1993.
- [8] D. C. Feldmeier and P. R. Karn. UNIX Password Security - Ten Years Later. In *Proceedings of the Crypto Conference*, pages 44–63, August 1990.
- [9] P. Gutmann. The Design of a Cryptographic Security Architecture. In *Proceedings of the 8th USENIX Security Symposium*, August 1999.
- [10] P. Gutmann. An Open-source Cryptographic Coprocessor. In *Proceedings of the 9th USENIX Security Symposium*, August 2000.
- [11] J. Kay and J. Pasquale. The Importance of Non-Data Touching Processing Overheads in TCP/IP. In *Proceedings of the ACM SIGCOMM Conference*, pages 259–269, September 1993.
- [12] S. Kent and R. Atkinson. Security Architecture for the Internet Protocol. RFC 2401, November 1998.
- [13] A. D. Keromytis, J. Ioannidis, and J. M. Smith. Implementing IPsec. In *Proceedings of Global Internet (GlobeCom)*, pages 1948–1952, November 1997.
- [14] RSA Laboratories. *PKCS #11: Cryptographic Token Interface Standard, Version 2.01*, December 1997.
- [15] M. Lindemann and S. W. Smith. Improving DES Coprocessor Throughput for Short Operations. In *Proceedings of the 10th USENIX Security Symposium*, pages 67–81, August 2001.
- [16] J. Linn. Generic Security Service Application Programming Interface. RFC 2078, January 1997.
- [17] Microsoft Corporation. *Microsoft Cryptographic Application Programming Interface (CryptoAPI)*, second edition, December 1998.
- [18] S. Miltchev, S. Ioannidis, and A. D. Keromytis. A Study of the Relative Costs of Network Security Protocols. In *Proceedings of the USENIX Annual Technical Conference, Freenix Track*, pages 41–48, June 2002.
- [19] The Open Group. *Common Data Security Architecture (CDSA)*, second edition, May 1999.
- [20] C. Pu, H. Massalin, J. Ioannidis, and P. Metzger. The Synthesis System. *Computing Systems*, 1(1), 1988.
- [21] C. B. S. and J. M. Smith. Hardware/Software Organization of a High-Performance ATM Host Interface. *IEEE Journal on Selected Areas in Communications (Special Issue on High Speed Computer/Network Interfaces)*, 11(2):240–253, February 1993.
- [22] J. M. Smith. Practical Problems with a Cryptographic Protection Scheme. In *Proceedings of the Crypto Conference*, pages 64–73, August 1990.
- [23] J. M. Smith and C. B. S. Traw. Giving Applications Access to Gb/s Networking. *IEEE Network*, 7(4):44–52, July 1993.
- [24] J. M. Smith, C. B. S. Traw, and D. J. Farber. Cryptographic Support for a Gigabit Network. In *Proceedings of INET*, pages 229–237, June 1992.
- [25] V. Smyslov. Simple Cryptographic Program Interface (Crypto API). RFC 2628, June 1999.
- [26] Cross Organization CAPI Team. *Security Service API: Cryptographic API Recommendation, Updated and Abridged Edition*. National Security Agency, July 1997.

Appendix A: OCF Kernel API

- `int32_t crypto_get_driverid();`
`int crypto_register();`
`int crypto_kregister();`
`int crypto_unregister();`

Used by device drivers to register and unregister symmetric and asymmetric algorithm support with the OCF.

- `void crypto_done();`
`void crypto_kdone();`

Called by device drivers on completion of a request (symmetric and asymmetric, respectively).

- `int crypto_newsession();`

Called by consumers of cryptographic services (such as the IPsec stack) that wish to establish a new session with the framework. On success, the first argument will contain the Session Identifier (SID). The second argument contains all the necessary information for the driver to establish the session (keys, algorithms, offsets, *etc.* The third argument indicates whether only hardware acceleration is acceptable.

- `int crypto_freesession();`

Called to disestablish a previously-established session.

- `int crypto_dispatch();`

Called to process a request, encapsulated in its only argument. The various fields in that structure contain:

- The SID.
- The total length in bytes of the buffer to be processed,
- The total length of the result, which for symmetric crypto operations will be the same as the input length.
- The type of input buffer, as used in the kernel `malloc()` routine. This will be used if the framework needs to allocate a new buffer for the result (or for re-formatting the input).
- The routine that the OCF should invoke upon completion of the request, whether successful or not.
- The error type, if any errors were encountered. If the **EAGAIN** error code is returned, the SID has changed. The consumer should record the new SID and use it in all subsequent requests. In this case, the request may be re-submitted immediately. This mechanism is used by the framework to perform session migration (move a session from one driver to another, because of availability, performance, or other considerations).
- A bitmask of flags associated with this request. Currently, the only flag defined is **CRYPTO_F_IMBUF**, which indicates that the input buffer is an mbuf chain.
- The input and output buffers. The input buffer may be an mbuf chain or a contiguous buffer (as identified by the flags). The output buffer will be of the same type.
- A pointer to opaque data. This is passed through the crypto framework untouched and is intended for the invoking application's use.
- A linked list of operation descriptors, which indicate what operations should be applied, and in what sequence, to the input data. The descriptors indicate where each operation should start, the length of the data to be processed, where on the output buffer should the results be placed, the key material to be used, and various operation-specific flags (*e.g.*, what Initialization Vector to use for CBC-mode encryption).

- `int crypto_kdispatch();`

Similar to `crypto_dispatch()`, for public-key operations.