

Sanctum: Minimal Architectural Extensions for Isolated Execution

by

Victor Marius Costan

Submitted to the Department of Electrical Engineering and Computer Science

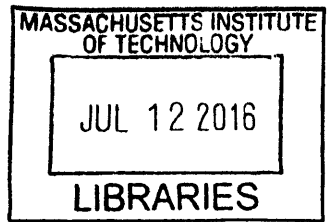
in partial fulfillment of the requirements for the degree of

Doctor of Philosophy in Computer Science and Engineering

at the

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

June 2016



ARCHIVES

© Massachusetts Institute of Technology 2016. All rights reserved.

Signature redacted

Author

Department of Electrical Engineering and Computer Science

February 11, 2016

Signature redacted

Certified by

Srinivas Devadas

Edwin Sibley Webster Professor of EECS

Thesis Supervisor

Signature redacted

Accepted by

Leslie A. Kolodziej

Chairman, Department Committee on Graduate Theses

Sanctum: Minimal Architectural Extensions for Isolated Execution

by

Victor Marius Costan

Submitted to the Department of Electrical Engineering and Computer Science
on February 11, 2016, in partial fulfillment of the
requirements for the degree of
Doctor of Philosophy in Computer Science and Engineering

Abstract

Intel's Software Guard Extensions (SGX) have captured the attention of security practitioners by promising to secure computation performed on a remote computer where all the privileged software is potentially malicious. Unfortunately, an independent analysis of SGX reveals that it is vulnerable to software attacks, and it can only be used by developers licensed by Intel. Furthermore, significant parts of SGX are undocumented, making it impossible for researchers outside of Intel to reason about some of its security properties.

Sanctum offers the same promise as SGX, namely strong provable isolation of software modules running concurrently and sharing resources, but protects against an important class of additional software attacks that infer private information from a program's memory access patterns. Sanctum shuns unnecessary complexity, leading to a simpler security analysis. We follow a principled approach to eliminating entire attack surfaces through isolation, rather than plugging attack-specific privacy leaks. Most of Sanctum's logic is implemented in trusted software, which is easier to analyze than SGX's opaque microcode.

Our prototype targets a Rocket RISC-V core, an open implementation that allows any researcher to reason about its security properties. Sanctum's extensions can be adapted to other RISC cores, because we do not change any major CPU building block. Instead, we add hardware at the interfaces between building blocks, without impacting cycle time.

Sanctum demonstrates that strong software isolation is achievable with a surprisingly small set of minimally invasive hardware changes, and a very reasonable overhead (assuming a software attack model) that is orders of magnitude less than what is incurred by ORAM-enabled processors. Our modifications cause a 2% area increase to the Rocket core. Over a set of benchmarks, Sanctum's worst observed overhead for isolated execution is 15.1% over an idealized insecure baseline, and 2.7% average overhead over a representative insecure baseline.

Thesis Supervisor: Srinivas Devadas
Title: Edwin Sibley Webster Professor of EECS

Acknowledgments

I would like to dedicate this work to my advisor, Prof. Srinivas Devadas. Words cannot express my gratitude for all the trust, patience, and learning opportunities that he has provided. I am a better person thanks to his advice.

I am very grateful to my family, for raising me and for all the sacrifices they made to have me attend MIT.

I owe the determination and inspiration for finishing up my thesis to Ms. Staphany Park. I would have most likely not made it to the finish line without her.

I am thankful to all my teachers. I am the sum of their contributions. Had any of them have been removed, I would not have been here today. I cannot enumerate everyone, but I will mention Ms. Paula Copacel, who taught me how to program, Ms. Victoria Ovanes, who taught me to aim high, and Mr. Tong Chen, who taught me the discipline that I needed to wade through Intel's documentation.

Funding for this research was partially provided by the National Science Foundation under contract number CNS-1413920.

Contents

1	Introduction	27
1.1	The Case for Hardware Isolation	27
1.2	Intel SGX is Not the Answer	28
1.3	Sanctum to the Rescue	30
1.4	Outline	31
2	Related Work	32
2.1	The IBM 4765 Secure Coprocessor	32
2.2	ARM TrustZone	34
2.3	The XOM Architecture	38
2.4	The Trusted Platform Module (TPM)	39
2.5	Intel’s Trusted Execution Technology (TXT)	41
2.6	The Aegis Secure Processor	42
2.7	The Bastion Architecture	43
2.8	Intel SGX	44
2.9	Sanctum in Context	46
2.10	Ascend and Phantom	47
3	Threat Model	48
4	Programming Model Overview	50
5	Hardware Extensions	56
5.1	LLC Address Input Transformation	57

5.2	Page Walker Input	59
5.3	Page Walker Memory Accesses	60
5.4	DMA Transfer Filtering	63
6	Software Design	64
6.1	Attestation Chain of Trust	64
6.1.1	The Measurement Root	64
6.1.2	The Signing Enclave	66
6.2	Security Monitor	68
6.2.1	DRAM Regions	68
6.2.2	Metadata Regions	69
6.2.3	Enclave Lifecycle	70
6.2.4	Enclave Code Execution	72
6.2.5	Mailboxes	73
6.2.6	Multi-Core Concurrency	74
6.3	Enclave Eviction	75
7	Security Argument	77
8	Performance Evaluation	80
8.1	Experiment Design	80
8.2	Cost of Added Hardware	81
8.3	Added Page Walker Latency	82
8.4	Security Monitor Overhead	82
8.5	Overhead of DRAM Region Isolation	83
9	Conclusion	86
A	Computer Architecture Background	87
A.1	Overview	88
A.2	Computational Model	90
A.3	Software Privilege Levels	93

A.4	Address Spaces	94
A.5	Address Translation	96
A.5.1	Address Translation Concepts	96
A.5.2	Address Translation and Virtualization	98
A.5.3	Page Table Attributes	99
A.6	Execution Contexts	100
A.7	Segment Registers	101
A.8	Privilege Level Switching	103
A.8.1	System Calls	104
A.8.2	Faults	105
A.8.3	VMX Privilege Level Switching	106
A.9	A Computer Map	107
A.9.1	The Motherboard	107
A.9.2	The Intel Management Engine (ME)	109
A.9.3	The Processor Die	110
A.9.4	The Core	111
A.10	Out-of-Order and Speculative Execution	111
A.10.1	Out-of-Order Execution	112
A.10.2	Speculative Execution	114
A.11	Cache Memories	115
A.11.1	Caching Principles	116
A.11.2	Cache Organization	117
A.11.3	Cache Coherence	118
A.11.4	Caching and Memory-Mapped Devices	120
A.11.5	Caches and Address Translation	123
A.12	Interrupts	125
A.13	Platform Initialization (Booting)	126
A.13.1	The UEFI Standard	126
A.13.2	SEC on Intel Platforms	128
A.13.3	PEI on Intel Platforms	129

A.14	CPU Microcode	130
A.14.1	The Role of Microcode	131
A.14.2	Microcode Structure	133
A.14.3	Microcode and Address Translation	134
A.14.4	Microcode and Booting	136
A.14.5	Microcode Updates	137
B	Security Background	151
B.1	Cryptographic Primitives	152
B.1.1	Cryptographic Keys	154
B.1.2	Privacy	156
B.1.3	Integrity	159
B.1.4	Freshness	162
B.2	Cryptographic Constructs	164
B.2.1	Certificate Authorities	164
B.2.2	Key Agreement Protocols	168
B.3	Software Attestation Overview	171
B.3.1	Authenticated Key Agreement	171
B.3.2	The Role of Software Measurement	173
B.4	Physical Attacks	174
B.4.1	Port Attacks	175
B.4.2	Bus Tapping Attacks	175
B.4.3	Chip Attacks	176
B.4.4	Power Analysis Attacks	178
B.5	Privileged Software Attacks	179
B.6	Software Attacks on Peripherals	180
B.6.1	PCI Express Attacks	181
B.6.2	DRAM Attacks	181
B.6.3	The Performance Monitoring Side Channel	182
B.6.4	Attacks on the Boot Firmware and Intel ME	182

B.6.5	Accounting for Software Attacks on Peripherals	184
B.7	Address Translation Attacks	184
B.7.1	Passive Attacks	185
B.7.2	Straightforward Active Attacks	185
B.7.3	Active Attacks Using Page Swapping	186
B.7.4	Active Attacks Based on TLBs	187
B.8	Cache Timing Attacks	189
B.8.1	Theory	189
B.8.2	Practical Considerations	190
B.8.3	Known Cache Timing Attacks	191
B.8.4	Defending against Cache Timing Attacks	191
C	Intel SGX Explained	193
C.1	Trusted Hardware	193
C.2	SGX Lightning Tour	196
C.3	Outline and Troubling Findings	198
C.4	SGX Programming Model	199
C.4.1	SGX Physical Memory Organization	199
C.4.2	The Memory Layout of an SGX Enclave	203
C.4.3	The Life Cycle of an SGX Enclave	210
C.4.4	The Life Cycle of an SGX Thread	214
C.4.5	EPC Page Eviction	223
C.4.6	SGX Enclave Measurement	234
C.4.7	SGX Enclave Versioning Support	240
C.4.8	SGX Software Attestation	250
C.4.9	SGX Enclave Launch Control	256
C.5	SGX Analysis	264
C.5.1	SGX Implementation Overview	264
C.5.2	SGX Memory Access Protection	268
C.5.3	SGX Security Check Correctness	274

C.5.4	Tracking TLB Flushes	281
C.5.5	Enclave Signature Verification	283
C.5.6	SGX Security Properties	287
C.6	Conclusion	300

List of Figures

2-1	The IBM 4765 secure coprocessor consists of an entire computer system placed inside an enclosure that can deter and detect physical attacks. The application and the system use separate processors. Sensitive memory can only be accessed by the system code, thanks to access control checks implemented in the system bus' hardware. Dedicated hardware is used to clear the platform's secrets and shut down the system when a physical attack is detected.	34
2-2	Smartphone SoC design based on TrustZone. The red IP blocks are TrustZone-aware. The red connections ignore the TrustZone secure bit in the bus address. Defining the system's security properties requires a complete understanding of all the red elements in this figure.	35
2-3	The measurement stored in a TPM platform configuration register (PCR). The PCR is reset when the system reboots. The software at every boot stage hashes the next boot stage, and sends the hash to the TPM. The PCR's new value incorporates both the old PCR value, and the new software hash. . . .	40
4-1	Software stack of a Sanctum machine	50
4-2	Per-enclave page tables	52
4-3	Enclave layout and data structures	53
5-1	Interesting bit fields in a physical address	57
5-2	Cache address shifting makes DRAM regions contiguous	58
5-3	Cache address shifter that shifts the PPN by 3 bits	59

5-4	A variable shifter that can shift by 2-5 bits can be composed of a fixed shifter by 2 bits and a variable shifter that can shift by 0-3 bits.	59
5-5	Sanctum’s cache address shifter and DMA transfer filter logic in the context of a RISC V-Rocket uncore	60
5-6	Page walker input for per-enclave page tables	61
5-7	Hardware support for per-enclave page tables: transforming the page table entries fetched by the page walker.	62
5-8	Sanctum’s page entry transformation logic in the context of a RISC V Rocket core	63
6-1	Sanctum’s root of trust is a measurement root routine burned into the CPU’s ROM. This code reads the security monitor from flash memory and generates an attestation key and certificate based on the monitor’s hash. Asymmetric key operations, colored in blue, are only performed the first time a monitor is used on a computer.	65
6-2	The certificate chain behind Sanctum’s software attestation signatures . . .	67
6-3	DRAM region allocation states and API calls	68
6-4	Security monitor data structures	69
6-5	Enclave states and enclave management API calls	71
6-6	Enclave thread metadata structure states and thread-related API calls . . .	72
6-7	Mailbox states and security monitor API calls related to inter-enclave communication	73
8-1	Sanctum’s modified page walk has minimal effect on benchmark performance	82
8-2	The impact of DRAM region allocation on the completion time of an enclaved benchmark, relative to an idea insecure baseline	83
8-3	Sanctum’s enclave overheads for one core utilizing 1/4 of the LLC compared against an idealized baseline (non-enclaved app using the entire LLC), and against a representative baseline (non-enclaved app sharing the LLC with concurrent instances)	84

A-1	A computer's core is its processors and memory, which are connected by a system bus. Computers also have I/O devices, such as keyboards, which are also connected to the processor via the system bus.	91
A-2	The memory abstraction	91
A-3	A processor fetches instructions from the memory and executes them. The RIP register holds the address of the instruction to be executed.	139
A-4	The system bus abstraction	139
A-5	The privilege levels in the x86 architecture, and the software that typically runs at each security level.	140
A-6	The four physical address spaces used by an Intel CPU. The registers and MSRs are internal to the CPU, while the memory and I/O address spaces are used to communicate with DRAM and other devices via system buses.	140
A-7	Virtual addresses used by software are translated into physical memory addresses using a mapping defined by the page tables.	140
A-8	The virtual memory abstraction gives each process its own virtual address space. The operating system multiplexes the computer's DRAM between the processes, while application developers build software as if it owns the entire computer's memory.	141
A-9	IA-32e address translation takes in a 48-bit virtual address and outputs a 52-bit physical address.	141
A-10	Address translation can be seen as a mapping between virtual page numbers and physical page numbers.	142
A-11	Virtual addresses used by software are translated into physical memory addresses using a mapping defined by the page tables.	142
A-12	Address translation when hardware virtualization is enabled. The kernel-managed page tables contain guest-physical addresses, so each level in the kernel's page table requires a full walk of the hypervisor's extended page table (EPT). A translation requires up to 20 memory accesses (the bold boxes), assuming the physical address of the kernel's PML4 is cached.	142

A-13 CPU registers in the 64-bit Intel architecture. RSP can be used as a general-purpose register (GPR), e.g., in pointer arithmetic, but it always points to the top of the program's stack. Segment registers are covered in § A.7. . . .	143
A-14 Loading a segment register. The 16-bit value loaded by software is a selector consisting of an index and a ring number. The index selects a GDT entry, which is loaded into the descriptor part of the segment register.	143
A-15 Example address computation process for <code>MOV FS:[RDX], 0</code> . The segment's base address is added to the address in RDX before address translation (§ A.5) takes place.	143
A-16 Modern privilege switching methods in the 64-bit Intel architecture.	144
A-17 The motherboard structures that are most relevant in a system security analysis.	144
A-18 The Intel Management Engine (ME) is an embedded computer hosted in the PCH. The ME has its own execution core, ROM and SRAM. The ME can access the host's DRAM via a memory controller and a DMA controller. The ME is remotely accessible over the network, as it has direct access to an Ethernet PHY via the SMBus.	144
A-19 The major components in a modern CPU package. § A.9.3 gives an uncore overview. § A.9.4 describes execution cores. § A.11.3 takes a deeper look at the uncore.	145
A-20 CPU core with two logical processors. Each logical processor has its own execution context and LAPIC (§ A.12). All the other core resources are shared.	145
A-21 The structures in a CPU core that are relevant to out-of-order and speculative execution. Instructions are decoded into micro-ops, which are scheduled on one of the execution unit's ports. The branch predictor enables speculative execution when a branch is encountered.	146

A-22	The steps taken by a cache memory to resolve an access to a memory address	
	A. A normal memory access (to cacheable DRAM) always triggers a cache lookup. If the access misses the cache, a fill is required, and a write-back might be required.	147
A-23	Cache organization and lookup, for a W -way set-associative cache with 2^l -byte lines and $S = 2^s$ sets. The cache works with n -bit memory addresses. The lowest l address bits point to a specific byte in a cache line, the next s bytes index the set, and the highest $n - s - l$ bits are used to decide if the desired address is in one of the W lines in the indexed set.	148
A-24	The stops on the ring interconnect used for inter-core and core-uncore communication.	148
A-25	The circuit for computing whether a physical address matches a memory type range. Assuming a CPU with 48-bit physical addresses, the circuit uses 36 AND gates and a binary tree of 35 XNOR (equality test) gates. The circuit outputs 1 if the address belongs to the range. The bottom 12 address bits are ignored, because memory type ranges must be aligned to 4 KB page boundaries.	149
A-26	Virtual addresses from the perspective of cache lookup and address translation. The bits used for the L1 set index and line offset are not changed by address translation, so the page tables do not impact L1 cache placement. The page tables do impact L2 and L3 cache placement. Using large pages (2 MB or 1 GB) is not sufficient to make L3 cache placement independent of the page tables, because of the LLC slice hashing function (§ A.11.3).	149
A-27	The phases of the Platform Initialization process in the UEFI specification.	149
A-28	The Firmware Interface Table (FIT) in relation to the firmware's memory map.	150
B-1	In a privacy attack, Eve sees the message sent by Alice to Bob and can understand the information inside it. In this case, Eve can tell that the message is a buy order, and not a sell order.	153

B-2	In an integrity attack, Eve replaces Alice's message with her own. In this case, Eve sends Bob a sell-everything order. In this case, Eve can tell that the message is a buy order, and not a sell order.	154
B-3	In a freshness attack, Eve replaces Alice's message with a message that she sent at an earlier time. In this example, Eve builds a database of labeled messages over time, and is able to send Bob her choice of a BUY or a SELL order.	154
B-4	In symmetric key cryptography, a secret key is shared by the parties that wish to communicate securely.	155
B-5	An asymmetric key generation algorithm produces a private key and an associated public key. The private key is held confidential, while the public key is given to any party who wishes to securely communicate with the private key's holder.	156
B-6	In a symmetric key secure permutation (block cipher), the same secret key must be provided to both the encryption and the decryption algorithm. . . .	156
B-7	In an asymmetric key block cipher, the encryption algorithm operates on a public key, and the decryption algorithm uses the corresponding private key.	157
B-8	Symmetric key block ciphers are combined with operating modes. Most operating modes require a random initialization vector (IV) to be generated for each encrypted message.	158
B-9	Asymmetric key encryption is generally used to bootstrap a symmetric key encryption scheme.	159
B-10	A block hash function operates on fixed-size message blocks and uses a fixed-size internal state.	160
B-11	In the symmetric key setting, integrity is assured by computing a Message Authentication Code (MAC) tag and transmitting it over the network along the message. The receiver feeds the MAC tag into a verification algorithm that checks the message's authenticity.	160

B-12	In the symmetric key setting, integrity is assured by computing a Hash-based Message Authentication Code (HMAC) and transmitting it over the network along the message. The receiver re-computes the HMAC and compares it against the version received from the network.	161
B-13	Signature schemes guarantee integrity in the asymmetric key setting. Signatures are created using the sender's private key, and are verified using the corresponding public key. A cryptographically secure hash function is usually employed to reduce large messages to small hashes, which are then signed.	162
B-14	The RSA signature scheme with PKCS #1 v1.5 padding specified in RFC 3447 combines a secure hash of the signed message with a DER-encoded specification of the secure hash algorithm used by the signature, and a padding string whose bits are all set to 1. Everything except for the secure hash output is considered to be a part of the PKCS #1 v1.5 padding.	163
B-15	Freshness guarantees can be obtained by adding timestamped nonces on top of a system that already offers integrity guarantees. The sender and the receiver use synchronized clocks to timestamp each message and discard unreasonably old messages. The receiver must check the nonce in each new message against a database of the nonces in all the unexpired messages that it has seen.	164
B-16	A certificate is a statement signed by a certificate authority (issuer) binding the identity of a subject to a public key.	165
B-17	A certificate issued by a CA can be validated by any party that has securely obtained the CA's public key. If the certificate is valid, the subject public key contained within can be trusted to belong to the subject identified by the certificate.	166
B-18	A hierarchical CA structure minimizes the usage of the root CA's private key, reducing the opportunities for it to get compromised. The root CA only signs the certificates of intermediate CAs, which sign the end users' certificates.	167

B-19 An ID card is a certificate that binds a subject’s full legal name (identity) to the subject’s physical appearance, which acts as a public key. 168

B-20 In the Diffie-Hellman Key Exchange (DKE) protocol, Alice and Bob agree on a shared secret key $K = g^{AB} \text{ mod } p$. An adversary who observes $g^A \text{ mod } p$ and $g^B \text{ mod } p$ cannot compute K 169

B-21 Any key agreement protocol is vulnerable to a man-in-the-middle (MITM) attack. The active attacker performs key agreements and establishes shared secrets with both parties. The attacker can then forward messages between the victims, in order to observe their communication. The attacker can also send its own messages to either, impersonating the other victim. 170

B-22 The chain of trust in software attestation. The root of trust is a manufacturer key, which produces an endorsement certificate for the secure processor’s attestation key. The processor uses the attestation key to produce the attestation signature, which contains a cryptographic hash of the container and a message produced by the software inside the container. 172

B-23 An example of an active memory mapping attack. The application’s author intends to perform a security check, and only call the procedure that discloses the sensitive information if the check passes. Malicious system software maps the virtual address of the procedure that is called when the check fails, to a DRAM page that contains the disclosing procedure. 186

B-24 An active memory mapping attack where the system software does not modify the page tables. Instead, two pages are evicted from DRAM to a slower storage medium. The malicious system software swaps the two pages’ contents then brings them back into DRAM, building the same incorrect page mapping as the direct attack shown in Figure B-23. This attack defeats protection measures that rely on tracking the virtual and disk addresses for DRAM pages. 187

B-25	An active memory mapping attack where the system software does not invalidate a core's TLBs when it evicts two pages from DRAM and exchanges their locations when reading them back in. The page tables are updated correctly, but the core with stale TLB entries has the same incorrect view of the protected container's code as in Figure B-23.	188
C-1	Secure remote computation. A user relies on a remote computer, owned by an untrusted party, to perform some computation on her data. The user has some assurance of the computation's integrity and privacy.	194
C-2	Trusted computing. The user trusts the manufacturer of a piece of hardware in the remote computer, and entrusts her data to a secure container hosted by the secure hardware.	195
C-3	Software attestation proves to a remote computer that it is communicating with a specific secure container hosted by a trusted platform. The proof is an attestation signature produced by the platform's secret attestation key. The signature covers the container's initial state, a challenge nonce produced by the remote computer, and a message produced by the container.	196
C-4	Enclave data is stored into the EPC, which is a subset of the PRM. The PRM is a contiguous range of DRAM that cannot be accessed by system software or peripherals.	200
C-5	An enclave's EPC pages are accessed using a dedicated region in the enclave's virtual address space, called ELRANGE. The rest of the virtual address space is used to access the memory of the host process. The memory mappings are established using the page tables managed by system software.	204
C-6	A possible layout of an enclave's virtual address space. Each enclave has a SECS, and one TCS per supported concurrent thread. Each TCS points to a sequence of SSAs, and specifies initial values for RIP and for the base addresses of FS and GS.	209

C-7	The SGX enclave life cycle management instructions and state transition diagram	211
C-8	The PAGEINFO structure supplies input data to SGX instructions such as EADD.	212
C-9	The stages of the life cycle of an SGX Thread Control Structure (TCS) that has two State Save Areas (SSAs).	215
C-10	Data flow diagram for a subset of the logic in EENTER. The figure omits the logic for disabling debugging features, such as hardware breakpoints and performance monitoring events.	216
C-11	If a hardware exception occurs during enclave execution, the synchronous execution path is aborted, and an Asynchronous Enclave Exit (AEX) occurs instead.	219
C-12	If a hardware exception occurs during enclave execution, the synchronous execution path is aborted, and an Asynchronous Enclave Exit (AEX) occurs instead.	222
C-13	SGX offers a method for the OS to evict EPC pages into non-PRM DRAM. The OS can then use its standard paging feature to evict the pages out of DRAM.	224
C-14	The VALID and BLOCKED bits in an EPC page's EPCM entry can be in one of three states. EADD and its siblings allocate new EPC pages. EREMOVE permanently deallocates an EPC page. EBLOCK blocks an EPC page so it can be evicted using EWB. ELDB and ELDU load an evicted page back into the EPC.	226
C-15	The EWB instruction outputs the encrypted contents of the evicted EPC page, a subset of the fields in the page's EPCM entry, a MAC tag, and a nonce. All this information is used by the ELDB or ELDU instruction to load the evicted page back into the EPC, with privacy, integrity and freshness guarantees.	229
C-16	The PAGEINFO structure used by the EWB and ELDU / ELDB instructions	230

C-17	The data flow of the EWB instruction that evicts an EPC page. The page's content is encrypted in a non-EPC RAM page. A nonce is created and saved in an empty slot inside a VA page. The page's EPCM metadata and a MAC are saved in a separate area in non-EPC memory.	232
C-18	A version tree formed by evicted VA pages and enclave EPC pages. The enclave pages are leaves, and the VA pages are inner nodes. The OS controls the tree's shape, which impacts the performance of evictions, but not their correctness.	302
C-19	SGX has a certificate-based enclave identity scheme, which can be used to migrate secrets between enclaves that contain different versions of the same software module. Here, enclave A's secrets are migrated to enclave B. . . .	303
C-20	An enclave's Signature Structure (SIGSTRUCT) is intended to be generated by an enclave building toolchain that has access to the enclave author's private RSA key.	304
C-21	EINIT verifies the RSA signature in the enclave's certificate. If the certificate is valid, the information in it is used to populate the SECS fields that make up the enclave's certificate-based identity.	305
C-22	EGETKEY implements a key derivation service that is primarily used by SGX's secret migration feature. The key derivation material is drawn from the SECS of the calling enclave, the information in a Key Request structure, and secure storage inside the CPU's hardware.	306
C-23	Setting up an SGX enclave and undergoing the software attestation process involves the SGX instructions EINIT and EREPORT, and two special enclaves authored by Intel, the SGX Launch Enclave and the SGX Quoting Enclave.	307
C-24	EREPORT data flow	308
C-25	The authenticity of the REPORT structure created by EREPORT can and should be verified by the report's target enclave. The target's code uses EGETKEY to obtain the key used for the MAC tag embedded in the REPORT structure, and then verifies the tag.	309

C-26	SGX’s software attestation is based on two secrets stored in e-fuses inside the processor’s die, and on a key received from Intel’s provisioning service.	310
C-27	When EGETKEY is asked to derive a Provisioning key, it does not use the Seal Secret or OWNEREPOCH. The Provisioning key does, however, depend on MRSIGNER and on the SVN of the SGX implementation.	311
C-28	The derivation material used to produce Provisioning Seal keys does not include the OWNEREPOCH value, so the keys survive computer ownership changes.	312
C-29	The SGX Launch Enclave computes the EINITTOKEN.	313
C-30	SGX adds a few security checks to the PMH. The checks ensure that all the TLB entries created by the address translation unit meet SGX’s memory access restrictions.	314
C-31	The algorithm used to initialize the SECS fields used by the TLB flush tracking method presented in this section.	315
C-32	The algorithm used by ETRACK to activate TLB flush tracking.	315
C-33	The algorithm that updates the TLB flush tracking state when an LP exits an enclave via EEXIT or AEX.	315
C-34	The algorithm that updates the TLB flush tracking state when an LP enters an enclave via EENTER or ERESUME.	316
C-35	The algorithm that ensures that all LPs running an enclave’s code when ETRACK was executed have exited enclave mode at least once.	316
C-36	The algorithm that marks the end of a TLB flushing cycle when EBLOCK is executed.	316
C-37	An RSA signature verification algorithm specialized for the case where the public exponent is 3. s is the RSA signature and m is the RSA key modulus. The algorithm uses two additional inputs, q_1 and q_2 .	317

C-38 A malicious OS can partition a cache between the software running inside an enclave and its own malicious code. Both the OS and the enclave software have cache sets dedicated to them. When allocating DRAM to itself and to the enclave software, the malicious OS is careful to only use DRAM regions that map to the appropriate cache sets. On a system with an Intel CPU, the OS can partition the L2 cache by manipulating the page tables in a way that is completely oblivious to the enclave's software. 318

List of Tables

2.1	Security features overview for the trusted hardware projects related to Intel's SGX	33
A.1	Sample feature-specific Intel architecture registers.	101
A.2	A typical GDT layout in the 64-bit Intel Architecture.	103
A.3	The essential fields of an IDT entry in 64-bit mode. Each entry points to a hardware exception or interrupt handler.	105
A.4	The snapshot pushed on the handler's stack when a hardware exception occurs. IRET restores registers from this snapshot.	106
A.5	Pseudo micro-ops for the out-of-order execution example.	113
A.6	Data written by the renamer into the reorder buffer (ROB), for the micro-ops in Table A.5.	113
A.7	Relevant entries of the register allocation table after the micro-ops in Table A.5 are inserted into the ROB.	113
A.8	Approximate sizes and access times for each level in the memory hierarchy of an Intel processor, from [131]. Memory sizes and access times differ by orders of magnitude across the different levels of the hierarchy. This table does not cover multi-processor systems.	117
A.9	Approximate sizes and access times for each level in the TLB hierarchy, from [4].	123
B.1	Desirable security guarantees and primitives that provide them	152
B.2	Popular cryptographic primitives that are considered to be secure against today's adversaries	153

C.1	The fields in an EPCM entry that track the ownership of pages.	202
C.2	An enclave's attributes are the sub-fields in the ATTRIBUTES field of the enclave's SECS. This table shows a subset of the attributes defined in the SGX documentation.	205
C.3	The fields in an EPCM entry that indicate the enclave's intended virtual memory layout.	207
C.4	64-byte block extended into MRENCLAVE by ECREATE	235
C.5	64-byte block extended into MRENCLAVE by EADD. The ENCLAVEOFFSET is computed by subtracting the BASEADDR in the enclave's SECS from the LINADDR field in the PAGEINFO structure.	237
C.6	64-byte blocks extended into MRENCLAVE by EEXTEND. The ENCLAVE-OFFSET is computed by subtracting the BASEADDR in the enclave's SECS from the LINADDR field in the PAGEINFO structure.	238
C.7	A subset of the metadata fields in a SIGSTRUCT enclave certificate	242
C.8	The format of the RSA signature used in a SIGSTRUCT enclave certificate	242
C.9	A subset of the fields in the KEYREQUEST structure	247
C.10	The fields in an EPCM entry.	267
C.11	The fields in an EPCM entry.	269
C.12	Values of the PT (page type) field in an EPCM entry.	270

Chapter 1

Introduction

Between the Snowden revelations and the seemingly unending series of high-profile hacks of the past few years, the public's confidence in software systems has decreased considerably. At the same time, key initiatives such as cloud computing and the IoT (Internet of Things) require users to trust the systems providing these services. We must therefore develop capabilities to build software systems with better security, and gain back our users' trust.

1.1 The Case for Hardware Isolation

The best known practical method for securing a software system amounts to modularizing the system's code in a way that minimizes the code in the modules responsible for the system's security. Formal verification techniques are then applied to these modules, which make up the system's trusted codebase (TCB). The method assumes that the software modules are isolated, so the TCB automatically includes the mechanism that provides the isolation guarantees.

Today's systems rely on an operating system kernel, or a hypervisor (such as Linux or Xen, respectively) for software isolation. However **each** of the last three years (2012-2014) witnessed over 100 new security vulnerabilities in Linux [8, 35], and over 40 in Xen [10].

One may hope that formal verification methods can produce a secure kernel or hypervisor. Unfortunately, these codebases are far outside our verification capabilities: Linux and Xen have over *17 million*[17] and 150,000[12] lines of code, respectively. In stark contrast, the

seL4 formal verification effort [122] spent *20 man-years* to cover 9,000 lines of code.

Between Linux and Xen’s history of vulnerabilities and dire prospects for formal verification, a prudent system designer cannot include either in a TCB (trusted computing base), and must look elsewhere for a software isolation mechanism.

Fortunately, Intel’s Software Guard Extensions (SGX) [143, 15] has brought attention to the alternative of providing software isolation primitives in the CPU’s hardware. This avenue is appealing because the CPU is an unavoidable TCB component, and processor manufacturers have strong economic incentives to build correct hardware.

1.2 Intel SGX is Not the Answer

Sadly, although the SGX design includes a vast array of defenses against a variety of software and physical attacks, it fails to offer meaningful software isolation guarantees. The SGX threat model protects against all direct attacks, but excludes “side-channel attacks”, even if they can be performed in software.

Alarmingly, **cache timing attacks require only unprivileged software running on the victim’s host computer**, and do not rely on any physical access to the machine. This is particularly concerning in a cloud computing scenario, where gaining software access to the victim’s computer only requires a credit card [161], whereas physical access is a harder prospect, requiring trespass, coercion, or social engineering on the cloud provider’s employees.

Similarly, in many Internet of Things scenarios, the processing units have some amount of physical security, but they run outdated software stacks that have known security vulnerabilities. For example, an attacker might exploit a vulnerability in an IoT lock’s Bluetooth stack and obtain software execution privileges, then mount a cache timing attack on its access-granting process, and obtain the cryptographic key that opens the lock.

Furthermore, a thorough analysis of SGX reveals that it is impossible for anyone outside Intel to reason about the SGX’s security properties, because significant implementation details are not covered by the publicly available documentation. This a meaningful concern because the myriad of security vulnerabilities [199, 197, 200, 47, 169, 198, 195, 53] in

TXT [74], which is Intel's previous attempt at securing remote computation, show that securing the architecture underlying Intel's processors is incredibly challenging, even in the presence of strong economic incentives.

If the SGX successor claimed to protect against cache timing attacks, substantiating such a claim would require an analysis of its hardware and microcode, and ensuring that no implementation detail is vulnerable to cache timing attacks. Barring a highly unlikely shift to open-source hardware from Intel, such analysis will simply never happen.

A concrete example: the SGX documentation [98, 102] does not state where SGX stores the EPCM (enclave page cache map). If the EPCM is stored in cacheable RAM, the page translation verification step is subject to cache timing attacks. Interestingly, this detail is unnecessary for analyzing the security of today's SGX implementation, as we know for certain that SGX uses the operating system's page tables, and therefore page translations are vulnerable to cache timing attacks. The example does, however, demonstrate the fine nature of crucial details that are simply undocumented in today's hardware security implementations.

Adding insult to the injury, SGX includes a **hardware-enforced licensing scheme that requires developers who wish to take advantage of SGX's security to negotiate a business agreement with Intel**. If SGX gains widespread adoption, the licensing scheme built into it would give Intel the power to choose winners and losers in any market where products can benefit from secure remote computation, which includes Software as a Service (SaaS), video games, and Internet of Things (IoT) applications. This is very similar to the Net Neutrality issue in the United States, where a handful of colluding Internet Service Providers are stifling innovation by extorting money from service providers in return for not throttling the traffic between the services and their users.

In summary, while the principles behind SGX have great potential, the SGX design does not offer meaningful isolation guarantees, the SGX implementation is not open enough for independent researchers to be able to analyze its security properties, and the SGX licensing plan threatens to have a chilling effect on security software innovations.

1.3 Sanctum to the Rescue

Our main contribution is a software isolation scheme that addresses the issues raised above. Sanctum’s isolation provably defends against known software side-channel attacks, including cache timing attacks and passive address translation attacks. Sanctum is a co-design that combines **minimal** and **minimally invasive** hardware modifications with a trusted software **security monitor** that is amenable to formal verification.

We achieve minimality by reusing and lightly modifying existing, well-understood mechanisms. For example, our per-enclave page tables implementation uses the core’s existing page walking circuit, and requires very little extra logic. Sanctum is minimally invasive because it does not require modifying any major CPU building block. We only add hardware to the interfaces between blocks, and do not modify any block’s input or output. Our use of conventional building blocks translates into less effort needed to validate a Sanctum implementation.

We demonstrate that memory access pattern attacks can be foiled without incurring unreasonable overheads. Our hardware changes are small enough to present the added circuits, in their entirety, in Figures 5-6 and 5-7. Sanctum cores have the same clock speed as their insecure counterparts, as we do not modify the CPU core critical execution path. Using a straightforward cache partitioning scheme with Sanctum adds a 2.7% execution time overhead, which is orders of magnitude lower than the overheads of the ORAM schemes [69, 178] that are usually employed to conceal memory access patterns.

All the layers of Sanctum’s TCB are open-source and not encumbered by patents, trade secrets, or other similar intellectual property concerns that would disincentivize security researchers from analyzing it. Our prototype targets the Rocket Chip [130], an open-sourced implementation of the RISC-V [194, 192] instruction set architecture, which is an open standard. Sanctum’s software stack is open-sourced under the MIT license, and the authors did not file any patent.

To further encourage analysis, most of our security monitor is written in portable C++ which, once formally verified, can be used across different CPU implementations. Furthermore, even the non-portable assembly code can be reused across different implementations

of the same architecture. In comparison, SGX’s microcode is CPU model-specific, so each micro-architectural revision would require a separate formal verification effort.

1.4 Outline

Chapter 2 reviews the trusted computing hardware projects that are related to Sanctum. Chapter 3 outlines Sanctum’s threat model, which drives our entire design. Chapter 4 describes Sanctum’s programming model and introduces most of the differences between Sanctum and SGX. Chapter 5 describes the hardware modifications required by Sanctum. Remarkably, we are able to our modifications using gate-level circuit figures. Chapter 7 argues that the design presented here meets Sanctum’s threat model, drawing heavily on the security analysis presented in the context of SGX. Last, Chapter 8 evaluates Sanctum’s performance overheads.

We expect that our readers will be more inclined to go through the extensive background material in this thesis if it pertains to a popular, widely deployed architecture. For this reason, all the background material is presented from the context of the Intel architecture and Intel’s SGX, even though Sanctum actually targets the RISC-V architecture.

Appendix A reviews the architectural principles and the details of Intel’s architecture and micro-architecture needed to understand SGX. Most of the principles translate to RISC-V and Sanctum.

Appendix B summarizes the background knowledge needed to reason about an architecture’s security, including a quick overview of cryptographic primitives and constructs, and brief descriptions of the classes of attacks that trusted computing architectures must withstand.

Appendix C describes SGX in painstaking detail and analyzes its security properties. Sanctum uses the same principles as SGX whenever possible, so Appendix C can be used to build the intuition needed to reason about Sanctum’s security properties.

Chapter 2

Related Work

This section describes the broader picture of trusted hardware projects that Sanctum belongs to. Table 2.1 summarizes the security properties of Sanctum and the other trusted hardware presented here.

2.1 The IBM 4765 Secure Coprocessor

Secure coprocessors [207] encapsulate an entire computer system, including a CPU, a cryptographic accelerator, caches, DRAM, and an I/O controller within a tamper-resistant environment. The enclosure includes hardware that deters attacks, such as a Faraday cage, as well as an array of sensors that can detect tampering attempts. The secure coprocessor destroys the secrets that it stores when an attack is detected. This approach has good security properties against physical attacks, but tamper-resistant enclosures are very expensive [16], relatively to the cost of a computer system.

The IBM 4758 [177], and its most current-day successor, the IBM 4765 [2] (shown in Figure 2-1) are representative examples of secure coprocessors. The 4758 was certified to withstand physical attacks to FIPS 140-1 Level 4 [176], and the 4765 meets the rigors of FIPS 140-2 Level 4 [1].

The 4765 relies heavily on physical isolation for its security properties. Its system software is protected from attacks by the application software by virtue of using a dedicated service processor that is completely separate from the application processor. Special-purpose

Table 2.1: Security features overview for the trusted hardware projects related to Intel's SGX

Attack	TrustZone	TPM	TPM+TXT	SGX	XOM	Aegis	Bastion	Ascend, Phantom	Sanctum
Malicious containers (direct probing)	N/A (secure world is trusted)	N/A (The whole computer is one container)	N/A (Does not allow concurrent containers)	Access checks on TLB misses	Identifier tag checks	Security kernel separates containers	Access checks on each memory access	OS separates containers	Access checks on TLB misses
Malicious OS (direct probing)	Access checks on TLB misses	N/A (OS measured and trusted)	Host OS preempted during late launch	Access checks on TLB misses	OS has its own identifier	Security kernel measured and isolated	Memory encryption and HMAC	X	Access checks on TLB misses
Malicious hypervisor (direct probing)	Access checks on TLB misses	N/A (Hypervisor measured and trusted)	Hypervisor preempted during late launch	Access checks on TLB misses	N/A (No hypervisor support)	N/A (No hypervisor support)	Hypervisor measured and trusted	N/A (No hypervisor support)	Access checks on TLB misses
Malicious firmware	N/A (firmware is a part of the secure world)	CPU microcode measures PEI firmware	SINIT ACM signed by Intel key and measured	SMM handler is subject to TLB access checks	N/A (Firmware is not active after booting)	N/A (Firmware is not active after booting)	Hypervisor measured after boot	N/A (Firmware is not active after booting)	Firmware is measured and trusted
Malicious containers (cache timing)	N/A (secure world is trusted)	N/A (Does not allow concurrent containers)	N/A (Does not allow concurrent containers)	X	X	X	X	X	Each enclave gets own cache partition
Malicious OS (page fault recording)	Secure world has own page tables	N/A (OS measured and trusted)	Host OS preempted during late launch	X	N/A (Paging not supported)	X	X	X	Per-enclave page tables
Malicious OS (cache timing)	X	N/A (OS measured and trusted)	Host OS preempted during late launch	X	X	X	X	X	Non-enclave software uses a separate cache partition
DMA from malicious peripheral	On-chip bus bounces secure world accesses	X	IOMMU bounces DMA into TXT memory range	IOMMU bounces DMA into PRM	Equivalent to physical DRAM access	Equivalent to physical DRAM access	Equivalent to physical DRAM access	Equivalent to physical DRAM access	MC bounces DMA outside allowed range
Physical DRAM read	Secure world limited to on-chip SRAM	X	X	Undocumented memory encryption engine	DRAM encryption	DRAM encryption	DRAM encryption	DRAM encryption	X
Physical DRAM write	Secure world limited to on-chip SRAM	X	X	Undocumented memory encryption engine	HMAC of address and data	HMAC of address, data, timestamp	Merkle tree over DRAM	HMAC of address, data, timestamp	X
Physical DRAM rollback write	Secure world limited to on-chip SRAM	X	X	Undocumented memory encryption engine	X	Merkle tree over HMAC timestamps	Merkle tree over DRAM	Merkle tree over HMAC timestamps	X
Physical DRAM address reads	Secure world in on-chip SRAM	X	X	X	X	X	X	ORAM	X
Hardware TCB size	CPU chip package	Motherboard (CPU, TPM, DRAM, buses)	Motherboard (CPU, TPM, DRAM, buses)	CPU chip package	CPU chip package	CPU chip package	CPU chip package	CPU chip package	CPU chip package
Software TCB size	Secure world (firmware, OS, application)	All software on the computer	SINIT ACM + VM (OS, application)	Application module + privileged containers	Application module + hypervisor	Application module + security kernel	Application module + hypervisor	Application process + trusted OS	Application module + security monitor

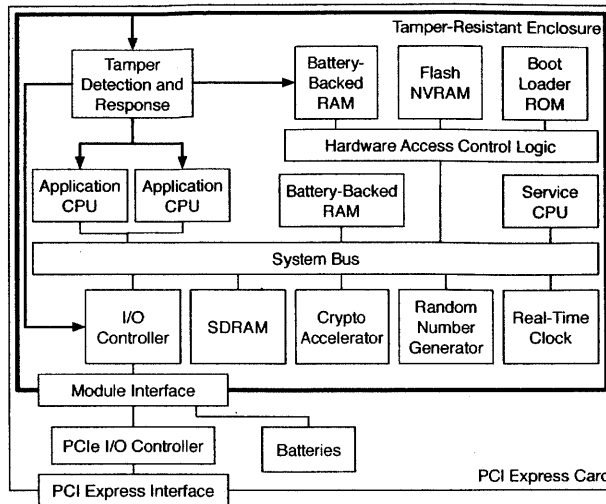


Figure 2-1: The IBM 4765 secure coprocessor consists of an entire computer system placed inside an enclosure that can deter and detect physical attacks. The application and the system use separate processors. Sensitive memory can only be accessed by the system code, thanks to access control checks implemented in the system bus' hardware. Dedicated hardware is used to clear the platform's secrets and shut down the system when a physical attack is detected.

bus logic prevents the application processor from accessing privileged resources, such as the battery-backed memory that stores the system software's secrets.

The 4765 implements software attestation. The coprocessor's attestation key is stored in battery-backed memory that is only accessible to the service processor. Upon reset, the service processor executes a first-stage bootloader stored in ROM, which measures and loads the system software. In turn, the system software measures the application code stored in NVRAM and loads it into the DRAM chip accessible to the application processor. The system software provides attestation services to the application loaded inside the coprocessor.

2.2 ARM TrustZone

ARM's TrustZone [14] is a collection of hardware modules that can be used to conceptually partition a system's resources between a *secure world*, which hosts a secure container, and a *normal world*, which runs an untrusted software stack. The TrustZone documentation [20]

describes semiconductor intellectual property cores (IP blocks) and ways in which they can be combined to achieve certain security properties, reflecting the fact that ARM is an IP core provider, not a chip manufacturer. Therefore, the mere presence of TrustZone IP blocks in a system is not sufficient to determine whether the system is secure under a specific threat model. Figure 2-2 illustrates a design for a smartphone *System-on-Chip* (SoC) design that uses TrustZone IP blocks.

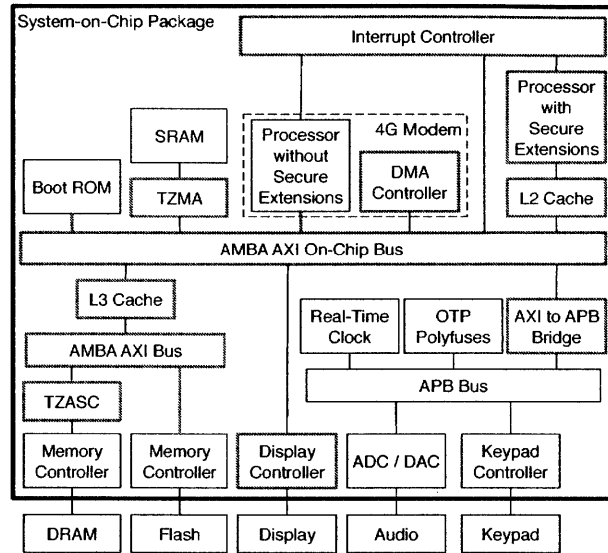


Figure 2-2: Smartphone SoC design based on TrustZone. The red IP blocks are TrustZone-aware. The red connections ignore the TrustZone secure bit in the bus address. Defining the system’s security properties requires a complete understanding of all the red elements in this figure.

TrustZone extends the address lines in the AMBA AXI system bus [19] with one signal that indicates whether an access belongs to the secure or normal (non-secure) world. ARM processor cores that include TrustZone’s “Security Extensions” can switch between the normal world and the secure world when executing code. The address in each bus access executed by a core reflects the world in which the core is currently executing.

The reset circuitry in a TrustZone processor places it in secure mode, and points it to the first-stage bootloader stored in on-chip ROM. TrustZone’s TCB includes this bootloader, which initializes the platform, sets up the TrustZone hardware to protect the secure container from untrusted software, and loads the normal world’s bootloader. The secure container must also implement a monitor that performs the context switches needed to transition an

execution core between the two worlds. The monitor must also handle hardware exceptions, such as interrupts, and route them to the appropriate world.

The TrustZone design gives the secure world's monitor unrestricted access to the normal world, so the monitor can implement inter-process communication (IPC) between the software in the two worlds. Specifically, the monitor can issue bus accesses using both secure and non-secure addresses. In general, the secure world's software can compromise any level in the normal world's software stack. For example, the secure container's software can jump into arbitrary locations in the normal world by flipping a bit in a register. The untrusted software in the normal world can only access the secure world via an instruction that jumps into a well-defined location inside the monitor.

Conceptually, each TrustZone CPU core provides separate address translation units for the secure and normal worlds. This is implemented by two page table base registers, and by having the page walker use the page table base corresponding to the core's current world. The physical addresses in the page table entries are extended to include the values of the secure bit to be issued on the AXI bus. The secure world is protected from untrusted software by having the CPU core force the secure bit in the address translation result to zero for normal world address translations. As the secure container manages its own page tables, its memory accesses cannot be directly observed by the untrusted OS's page fault handler.

TrustZone-aware hardware modules, such as caches, are trusted to use the secure address bit in each bus access to enforce the isolation between worlds. For example, TrustZone's caches store the secure bit in the address tag for each cache line, which effectively provides completely different views of the memory space to the software running in different worlds. This design assumes that memory space is partitioned between the two worlds, so no aliasing can occur.

The TrustZone documentation describes two TLB configurations. If many context switches between worlds are expected, the TLB IP blocks can be configured to include the secure bit in the address tag. Alternatively, the secure bit can be omitted from the TLBs, as long as the monitor flushes the TLBs when switching contexts.

The hardware modules that do not consume TrustZone's address bit are expected to be connected to the AXI bus via IP cores that implement simple partitioning techniques.

For example, the TrustZone Memory Adapter (TZMA) can be used to partition an on-chip ROM or SRAM into a secure region and a normal region, and the TrustZone Address Space Controller (TZASC) partitions the memory space provided by a DRAM controller into secure and normal regions. A TrustZone-aware DMA controller rejects DMA transfers from the normal world that reference secure world addresses.

It follows that analyzing the security properties of a TrustZone system requires a precise understanding of the behavior and configuration of all the hardware modules that are attached to the AXI bus. For example, the caches described in TrustZone's documentation do not enforce a complete separation between worlds, as they allow a world's memory accesses to evict the other world's cache lines. This exposes the secure container software to cache timing attacks from the untrusted software in the normal world. Unfortunately, hardware manufacturers that license the TrustZone IP cores are reluctant to disclose all the details of their designs, making it impossible for security researchers to reason about TrustZone-based hardware.

The TrustZone components do not have any counter-measures for physical attacks. However, a system that follows the recommendations in the TrustZone documentation will not be exposed to physical attacks, under a threat model that trusts the processor chip package. The AXI bus is designed to connect components in an SoC design, so it cannot be tapped by an attacker. The TrustZone documentation recommends having all the code and data in the secure world stored in on-chip SRAM, which is not subject to physical attacks. However, this approach places significant limits on the secure container's functionality, because on-chip SRAM is many orders of magnitude more expensive than a DRAM chip of the same capacity.

TrustZone's documentation does not describe any software attestation implementation. However, it does outline a method for implementing secure boot, which comes down to having the first-stage bootloader verify a signature in the second-stage bootloader against a public key whose cryptographic hash is burned into on-chip *One-Time Programmable* (OTP) polysilicon fuses. A hardware measurement root can be built on top of the same components, by storing a per-chip attestation key in the polyfuses, and having the first-stage bootloader measure the second-stage bootloader and store its hash in an on-chip SRAM

region allocated to the secure world. The polyfuses would be gated by a TZMA IP block that makes them accessible only to the secure world.

2.3 The XOM Architecture

The execute-only memory (XOM) architecture [132] introduced the approach of executing sensitive code and data in isolated containers managed by untrusted host software. XOM outlined the mechanisms needed to isolate a container's data from its untrusted software environment, such as saving the register state to a protected memory area before servicing an interrupt.

XOM supports multiple containers by tagging every cache line with the identifier of the container owning it, and ensures isolation by disallowing memory accesses to cache lines that don't match the current container's identifier. The operating system and the untrusted applications are considered to belong to a container with a null identifier.

XOM also introduced the integration of encryption and HMAC functionality in the processor's memory controller to protect container memory from physical attacks on DRAM. The encryption and HMAC functionality is used for all cache line evictions and fetches, and the ECC bits in DRAM chips are repurposed to store HMAC values.

XOM's design cannot guarantee DRAM freshness, so the software in its containers is vulnerable to physical replay attacks. Furthermore, XOM does not protect a container's memory access patterns, meaning that any piece of malicious software can perform cache timing attacks against the software in a container. Last, XOM containers are destroyed when they encounter hardware exceptions, such as page faults, so XOM does not support paging.

XOM predates the attestation scheme described above, and relies on a modified software distribution scheme instead. Each container's contents are encrypted with a symmetric key, which also serves as the container's identity. The symmetric key, in turn, is encrypted with the public key of each CPU that is trusted to run the container. A container's author can be assured that the container is running on trusted software by embedding a secret into the encrypted container data, and using it to authenticate the container. While conceptually simpler than software attestation, this scheme does not allow the container author to vet the

container's software environment.

2.4 The Trusted Platform Module (TPM)

The Trusted Platform Module (TPM) [75] introduced the software attestation model described at the beginning of this section. The TPM design does not require any hardware modifications to the CPU, and instead relies on an auxiliary tamper-resistant chip. The TPM chip is only used to store the attestation key and to perform software attestation. The TPM was widely deployed on commodity computers, because it does not rely on CPU modifications. Unfortunately, the cost of this approach is that the TPM has very weak security guarantees, as explained below.

The TPM design provides one isolation container, covering all the software running on the computer that has the TPM chip. It follows that the measurement included in an attestation signature covers the entire OS kernel and all the kernel modules, such as device drivers. However, commercial computers use a wide diversity of devices, and their system software is updated at an ever-increasing pace, so it is impossible to maintain a list of acceptable measurement hashes corresponding to a piece of trusted software. Due to this issue, the TPM's software attestation is not used in many security systems, despite its wide deployment.

The TPM design is technically not vulnerable to any software attacks, because it trusts all the software on the computer. However, a TPM-based system is vulnerable to an attacker who has physical access to the machine, as the TPM chip does not provide any isolation for the software on the computer. Furthermore, the TPM chip receives the software measurements from the CPU, so TPM-based systems are vulnerable to attackers who can tap the communication bus between the CPU and the TPM.

Last, the TPM's design relies on the software running on the CPU to report its own cryptographic hash. The TPM chip resets the measurements stored in Platform Configuration Registers (PCRs) when the computer is rebooted. Then, the TPM expects the software at each boot stage to cryptographically hash the software at the next stage, and send the hash to the TPM. The TPM updates the PCRs to incorporate the new hashes it receives, as shown

in Figure 2-3. Most importantly, the PCR value at any point reflects all the software hashes received by the TPM up to that point. This makes it impossible for software that has been measured to “remove” itself from the measurement.

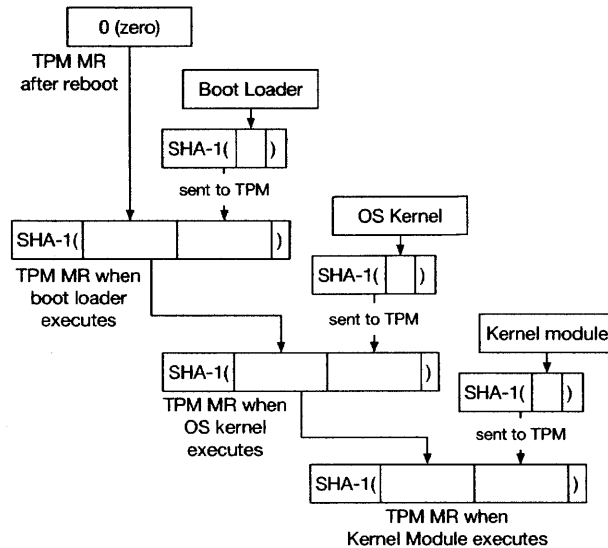


Figure 2-3: The measurement stored in a TPM platform configuration register (PCR). The PCR is reset when the system reboots. The software at every boot stage hashes the next boot stage, and sends the hash to the TPM. The PCR’s new value incorporates both the old PCR value, and the new software hash.

For example, the firmware on most modern computers implements the platform initialization process in the Unified Extensible Firmware Interface (UEFI) specification [186]. Each platform initialization phase is responsible for verifying or measuring the firmware that implements the next phase. The SEC firmware initializes the TPM PCR, and then stores the PEI’s measurement into a measurement register. In turn, the PEI implementation measures the DXE firmware and updates the measurement register that stores the PEI hash to account for the DXE hash. When the OS is booted, the hash in the measurement register accounts for all the firmware that was used to boot the computer.

Unfortunately, the security of the whole measurement scheme hinges on the requirement that the first hash sent to the TPM must reflect the software that runs in the first boot stage. The TPM threat model explicitly acknowledges this issue, and assumes that the firmware responsible for loading the first stage bootloader is securely embedded in the motherboard. However, virtually every TPM-enabled computer stores its firmware in a flash memory

chip that can be re-programmed in software (§ A.9.1), so the TPM’s measurement can be subverted by an attacker who can reflash the computer’s firmware [31].

On very recent Intel processors, the attack described above can be defeated by having the initialization microcode (§ A.14.4) hash the computer’s firmware (specifically, the PEI code in UEFI [186] firmware) and communicate the hash to the TPM chip. This is marketed as the Measured Boot feature of Intel’s Boot Guard [166].

Sadly, most computer manufacturers use Verified Boot (also known as “secure boot”) instead of Measured Boot (also known as “trusted boot”). Verified Boot means that the processor’s microcode only boots into PEI firmware that contains a signature produced by a key burned into the chip’s e-fuses. Verified Boot does not impact the measurements stored on the TPM, so it does not improve the security of software attestation.

2.5 Intel’s Trusted Execution Technology (TXT)

Intel’s Trusted Execution Technology (TXT) [74] uses the TPM’s software attestation model and auxiliary tamper-resistant chip, but reduces the software inside the secure container to a virtual machine (guest operating system and application) hosted by the CPU’s hardware virtualization features (VMX [187]).

TXT isolates the software inside the container from untrusted software by ensuring that the container has exclusive control over the entire computer while it is active. This is accomplished by a secure initialization authenticated code module (SINIT ACM) that effectively performs a warm system reset before starting the container’s VM.

TXT requires a TPM chip with an extended register set. The registers used by the measured boot process described in § 2.4 are considered to make up the platform’s Static Root of Trust Measurement (SRTM). When a TXT VM is initialized, it updates TPM registers that make up the Dynamic Root of Trust Measurement (DRTM). While the TPM’s SRTM registers only reset at the start of a boot cycle, the DRTM registers are reset by the SINIT ACM, every time a TXT VM is launched.

TXT does not implement DRAM encryption or HMACs, and therefore is vulnerable to physical DRAM attacks, just like TPM-based designs. Furthermore, early TXT imple-

mentations were vulnerable to attacks where a malicious operating system would program a device, such as a network card, to perform DMA transfers to the DRAM region used by a TXT container [197, 200]. In recent Intel CPUs, the memory controller is integrated on the CPU die, so the SINIT ACM can securely set up the memory controller to reject DMA transfers targeting TXT memory. An Intel chipset datasheet [108] documents an “Intel TXT DMA Protected Range” IIO configuration register.

Early TXT implementations did not measure the SINIT ACM. Instead, the microcode implementing the TXT launch instruction verified that the code module contained an RSA signature by a hard-coded Intel key. SINIT ACM signatures cannot be revoked if vulnerabilities are found, so TXT’s software attestation had to be revised when SINIT ACM exploits [199] surfaced. Currently, the SINIT ACM’s cryptographic hash is included in the attestation measurement.

Last, the warm reset performed by the SINIT ACM does not include the software running in System Management Mode (SMM). SMM was designed solely for use by firmware, and is stored in a protected memory area (SMRAM) which should not be accessible to non-SMM software. However, the SMM handler was compromised on multiple occasions [47, 169, 198, 195, 53], and an attacker who obtains SMM execution can access the memory used by TXT’s container.

2.6 The Aegis Secure Processor

The Aegis secure processor [180] relies on a security kernel in the operating system to isolate containers, and includes the kernel’s cryptographic hash in the measurement reported by the software attestation signature. [182] argued that Physical Unclonable Functions (PUFs) [60] can be used to endow a secure processor with a tamper-resistant private key, which is required for software attestation. PUFs do not have the fabrication process drawbacks of EEPROM, and are significantly more resilient to physical attacks than e-fuses.

Aegis relies on a trusted security kernel to isolate each container from the other software on the computer by configuring the page tables used in address translation. The security kernel is a subset of a typical OS kernel, and handles virtual memory management, processes,

and hardware exceptions. As the security kernel is a part of the *trusted code base* (TCB), its cryptographic hash is included in the software attestation measurement. The security kernel uses processor features to isolate itself from the untrusted part of the operating system, such as device drivers.

The Aegis memory controller encrypts the cache lines in one memory range, and HMACs the cache lines in one other memory range. The two memory ranges can overlap, and are configurable by the security kernel. Thanks to the two ranges, the memory controller can avoid the latency overhead of cryptographic operations for the DRAM outside containers. Aegis was the first secure processor not vulnerable to physical replay attacks, as it uses a Merkle tree construction [61] to guarantee DRAM freshness. The latency overhead of the Merkle tree is greatly reduced by augmenting the L2 cache with the tree nodes for the cache lines.

Aegis' security kernel allows the OS to page out container memory, but verifies the correctness of the paging operations. The security kernel uses the same encryption and Merkle tree algorithms as the memory controller to guarantee the privacy and integrity of the container pages that are swapped out from DRAM. The OS is free to page out container memory, so it can learn a container's memory access patterns, at page granularity. Aegis containers are also vulnerable to cache timing attacks.

2.7 The Bastion Architecture

The Bastion architecture [33] introduced the use of a trusted hypervisor to provide secure containers to applications running inside unmodified, untrusted operating systems. Bastion's hypervisor ensures that the operating system does not interfere with the secure containers. We only describe Bastion's virtualization extensions to architectures that use nested page tables, like Intel's VMX [187].

The hypervisor enforces the containers' desired memory mappings in the OS page tables, as follows. Each Bastion container has a Security Segment that lists the virtual addresses and permissions of all the container's pages, and the hypervisor maintains a Module State Table that stores an inverted page map, associating each physical memory page to its container and

virtual address. The processor's hardware page walker is modified to invoke the hypervisor on every TLB miss, before updating the TLB with the address translation result. The hypervisor checks that the virtual address used by the translation matches the expected virtual address associated with the physical address in the Module State Table.

Bastion's cache lines are not tagged with container identifiers. Instead, only TLB entries are tagged. The hypervisor's TLB miss handler sets the container identifier for each TLB entry as it is created. Similarly to XOM and Aegis, the secure processor checks the TLB tag against the current container's identifier on every memory access.

Bastion offers the same protection against physical DRAM attacks as Aegis does, without the restriction that a container's data must be stored inside a continuous DRAM range. This is accomplished by extending cache lines and TLB entries with flags that enable memory encryption and HMACing. The hypervisor's TLB miss handler sets the flags on TLB entries, and the flags are propagated to cache lines on memory writes.

The Bastion hypervisor allows the untrusted operating system to evict secure container pages. The evicted pages are encrypted, HMACed, and covered by a Merkle tree maintained by the hypervisor. Thus, the hypervisor ensures the privacy, authenticity, and freshness of the swapped pages. However, the ability to freely evict container pages allows a malicious OS to learn a container's memory accesses with page granularity. Furthermore, Bastion's threat model excludes cache timing attacks.

Bastion does not trust the platform's firmware, and computes the cryptographic hash of the hypervisor after the firmware finishes playing its part in the booting process. The hypervisor's hash is included in the measurement reported by software attestation.

2.8 Intel SGX

Intel's Software Guard Extensions (SGX) [143, 15, 82] implements secure containers for applications without making any modifications to the processor's critical execution path. SGX does not trust any layer in the computer's software stack (firmware, hypervisor, OS). Instead, SGX's TCB consists of the CPU's microcode and a few privileged containers. SGX introduces an approach to solving some of the issues raised by multi-core processors with a

shared, coherent last-level cache.

SGX does not extend caches or TLBs with container identity bits, and does not require any security checks during normal memory accesses. As suggested in the TrustZone documentation, SGX always ensures that a core's TLBs only contain entries for the container that it is executing, which requires flushing the CPU core's TLBs when context-switching between containers and untrusted software.

SGX follows Bastion's approach of having the untrusted OS manage the page tables used by secure containers. The containers' security is preserved by a TLB miss handler that relies on an inverted page map (the EPCM) to reject address translations for memory that does not belong to the current container.

Like Bastion, SGX allows the untrusted operating system to evict secure container pages, in a controlled fashion. After the OS initiates a container page eviction, it must prove to the SGX implementation that it also switched the container out of all cores that were executing its code, effectively performing a very coarse-grained TLB shutdown.

SGX's microcode ensures the privacy, authenticity, and freshness of each container's evicted pages, like Bastion's hypervisor. However, SGX relies on a version-based Merkle tree, inspired by Aegis [180], and adds an innovative twist that allows the operating system to dynamically shape the Merkle tree. SGX also shares Bastion's and Aegis' vulnerability to memory access pattern leaks, namely a malicious OS can directly learn a container's memory accesses at page granularity, and any piece of software can perform cache timing attacks.

SGX's software attestation is implemented using Intel's Enhanced Privacy ID (EPID) group signature scheme [28], which is too complex for a microcode implementation. Therefore, SGX relies on an assortment of privileged containers that receive direct access to the SGX processor's hardware keys. The privileged containers are signed using an Intel private key whose corresponding public key is hard-coded into the SGX microcode, similarly to TXT's SINIT ACM.

As SGX does not protect against cache timing attacks, the privileged enclave's authors cannot use data-dependent memory accesses. For example, cache attacks on the Quoting Enclave, which computes attestation signatures, would provide an attack with a processor's

EPID signing key and completely compromise SGX.

Intel’s documentation states that SGX guarantees DRAM privacy, authentication, and freshness by virtue of a Memory Encryption Engine (MEE). The MEE is informally described in an ISCA 2015 tutorial [106], and appears to lack a formal specification. In the absence of further information, we assume that SGX provides the same protection against physical DRAM attacks that Aegis and Bastion provide.

2.9 Sanctum in Context

Sanctum [40] introduced a straightforward software/hardware co-design that yields the same resilience against software attacks as SGX, and adds protection against memory access pattern leaks, such as page fault monitoring attacks and cache timing attacks.

Sanctum uses a conceptually simple cache partitioning scheme, where a computer’s DRAM is split into equally-sized continuous DRAM regions, and each DRAM region uses distinct sets in the shared last-level cache (LLC). Each DRAM region is allocated to exactly one container, so containers are isolated in both DRAM and the LLC. Containers are isolated in the other caches by flushing on context switches.

Like XOM, Aegis, and Bastion, Sanctum also considers the hypervisor, OS, and the application software to conceptually belong to a separate container. Containers are protected from the untrusted outside software by the same measures that isolate containers from each other.

Sanctum relies on a trusted security monitor, which is the first piece of firmware executed by the processor, and has the same security properties as those of Aegis’ security kernel. The monitor is measured by bootstrap code in the processor’s ROM, and its cryptographic hash is included in the software attestation measurement. The monitor verifies the operating system’s resource allocation decisions. For example, it ensures that no DRAM region is ever accessible to two different containers.

Each Sanctum container manages its own page tables mapping its DRAM regions, and handles its own page faults. It follows that a malicious OS cannot learn the virtual addresses that would cause a page fault in the container. Sanctum’s hardware modifications work

in conjunction with the security monitor to make sure that a container's page tables only reference memory inside the container's DRAM regions.

The Sanctum design focuses completely on software attacks, and does not offer protection from any physical attack. The authors expect Sanctum's hardware modifications to be combined with the physical attack protections in Aegis or Ascend.

2.10 Ascend and Phantom

The Ascend [56] and Phantom [136] secure processors introduced practical implementations of Oblivious RAM [69] techniques in the CPU's memory controller. These processors are resilient to attackers who can probe the DRAM address bus and attempt to learn a container's private information from its DRAM memory access pattern.

Implementing an ORAM scheme in a memory controller is largely orthogonal to the other secure architectures described above. It follows, for example, that Ascend's ORAM implementation can be combined with Aegis' memory encryption and authentication, and with Sanctum's hardware extensions and security monitor, yielding a secure processor that can withstand both software attacks and physical DRAM attacks.

Chapter 3

Threat Model

Like SGX, Sanctum isolates the software inside an **enclave** from any other software on the system, including privileged system software. Programmers are expected to move the sensitive code in their applications into enclaves. In general, an enclave receives encrypted sensitive information from outside, decrypts the information and performs some computation on it, and then returns encrypted results to the outside world.

For example, medical imaging software would use an enclave to decrypt a patient's X-ray and produce a diagnostic via an image processing algorithm. Application code that does not handle sensitive data, such as receiving encrypted X-rays over the network and storing the encrypted images in a database, would execute outside any enclave.

Sanctum protects the integrity and privacy of the code and data inside an enclave against an adversary that can carry out any practical **software** attack. We assume that an attacker can compromise any operating system and hypervisor present on the computer executing the enclave, and can launch rogue enclaves. The attacker knows the target computer's architecture and micro-architecture. The attacker can analyze passively collected data, such as page fault addresses, as well as mount active attacks, such as direct memory probing, memory probing via DMA transfers, and cache timing attacks.

Sanctum also defeats attackers who aim to compromise an operating system or hypervisor by running malicious applications and enclaves. This addresses concerns that enclaves provide new attack vectors for malware [167, 44]. We assume that the benefits of meaningful software isolation outweigh the downside of giving malware authors a new avenue for

frustrating malware detection and reverse engineering [48].

Lastly, Sanctum protects against a malicious computer owner who attempts to lie about the security monitor running on the computer. Specifically, the attacker aims to obtain an attestation stating that the computer is running an un-compromised security monitor, whereas a different supervisor had been loaded in the boot process. The un-compromised security monitor must not have any known vulnerability that causes it to disclose its cryptographic keys. The attacker is assumed to know the target computer's architecture and micro-architecture, and is allowed to run any combination of malicious security monitor, hypervisor, operating system, applications and enclaves.

We do not prevent timing attacks that exploit bottlenecks in the cache coherence directory bandwidth or in the DRAM bandwidth, deferring these to future work.

Sanctum does not protect against denial-of-service (DoS) attacks carried out by compromised system software, as a malicious OS may deny service by refusing to allocate any resources to an enclave. We *do* protect against malicious enclaves attempting to DoS an un-compromised OS.

We assume correct underlying hardware, so we do not protect against software attacks that exploit hardware bugs, such as rowhammer [121, 171] and other fault-injection attacks.

Sanctum's isolation mechanisms exclusively target software attacks. § 2 mentions related work that can harden a Sanctum system against some physical attacks. Furthermore, we consider software attacks that rely on sensor data to be physical attacks. For example, we do not address information leakage due to power variations, because software would require a temperature or current sensor to carry out such an attack.

Chapter 4

Programming Model Overview

By design, Sanctum’s programming model deviates from SGX as little as possible, while providing stronger security guarantees. We expect that application authors will link against a Sanctum-aware runtime that abstracts away most aspects of Sanctum’s programming model. For example, C programs would use a modified implementation of the `libc` standard library. Due to space constraints, we describe the programming model assuming that the reader is familiar with SGX as described in § C.

The software stack on a Sanctum machine, shown in Figure 4-1, is very similar to the stack on a SGX system with one notable exception: SGX’s microcode is replaced by a trusted software component, the **security monitor**, which runs at the highest privilege level (machine level in RISC-V) and therefore is immune to compromised system software.

We follow SGX’s approach of relegating the management of computation resources,

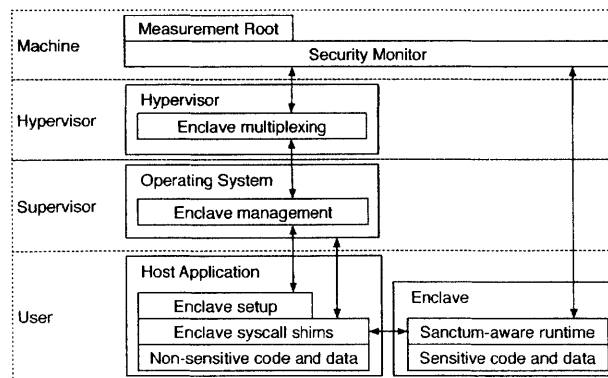


Figure 4-1: Software stack of a Sanctum machine

such as DRAM and execution cores, to untrusted system software. In Sanctum, the security monitor checks the system software's allocation decisions for correctness and commits them into the hardware's configuration. For simplicity, we refer to the software that manages resources as an *OS (operating system)*, even though it may be a combination of a hypervisor and a guest operating system kernel.

Figure 4-1 is representative of today's popular software stacks, where the operating system handles scheduling and demand paging, and the hypervisor multiplexes the computer's CPU cores. Sanctum is easy to integrate in such a stack, because the API calls that make up the security monitor's interface were designed with multiplexing in mind. Furthermore, a security-conscious hypervisor can use Sanctum's cache isolation primitive, the DRAM region, to protect against cross-VM cache timing attacks [18].

Like in SGX, an enclave stores its code and private data in parts of DRAM that have been allocated by the OS exclusively for the enclave's use, which are collectively called **the enclave's memory**. Consequently, we refer to the regions of DRAM that are not allocated to any enclave as **OS memory**. The security monitor tracks DRAM ownership, and ensures that no piece of DRAM is assigned to more than one enclave.

Each Sanctum enclave uses a range of virtual memory addresses (EVRANGE) to access its memory. The enclave's memory is mapped by the enclave's own page tables, which are also stored in the enclave's memory (Figure 4-2). This deviation from SGX is needed because page table dirty and accessed bits reveal memory access patterns at page granularity.

The enclave's virtual address space outside EVRANGE is used to access its host application's memory, via the page tables set up by the OS. Sanctum's hardware extensions implement dual page table lookup (§ 5.2), and make sure that an enclave's page tables can only point into the enclave's memory, while OS page tables can only point into OS memory (§ 5.3).

Like SGX, Sanctum supports multi-threaded enclaves, and enclaves must provision thread state areas for each thread that executes enclave code. Enclave threads, like their SGX cousins, run at the lowest privilege level (user level in RISC-V), so that a malicious enclave cannot compromise the OS. Specifically, enclaves cannot execute privileged instructions, and address translations that use OS page tables will result in page faults when accessing

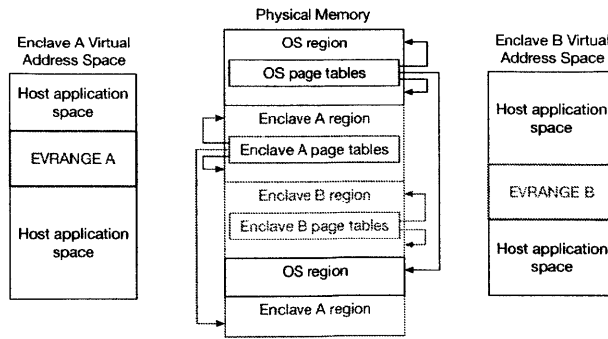


Figure 4-2: Per-enclave page tables

supervisor pages.

In order to prevent timing attacks, the metadata used by the security monitor to operate enclaves is stored in dedicated DRAM regions called **metadata regions**. Each metadata region is managed at the page level by the OS, like SGX's Enclave Page Cache (EPC), and includes a page map that is used by the security monitor to verify the OS' decisions, like SGX's Enclave Page Cache Map (EPCM). Unlike SGX's EPC, the metadata region pages only store enclave and threat metadata. Figure 4-3 shows how these structures are weaved together.

Sanctum, like SGX, considers system software to be untrusted, so it regulates transitions into and out of enclave code. An enclave's host application **enters an enclave** via a security monitor call that locks a thread state area, and transfers control to its entry point. After completing its intended task, the enclave code **exits** by asking the monitor to unlock the thread's state area, and transfer control back to the host application.

Like in SGX, enclaves cannot make system calls directly, as we cannot trust the OS to restore an enclave's execution state. Instead, the enclave's runtime must ask the host application to proxy syscalls such as file system and network I/O requests.

Sanctum's security monitor is the first responder for all hardware exceptions, including interrupts and faults. Like in SGX, an interrupt received during enclave execution causes an *asynchronous enclave exit* (AEX), whereby the monitor saves the core's registers in the current thread's AEX state area, zeroes the registers, exits the enclave, and dispatches the interrupt as if it was received by the code entering the enclave.

Unlike in SGX, resuming enclave execution after an AEX means re-entering the enclave

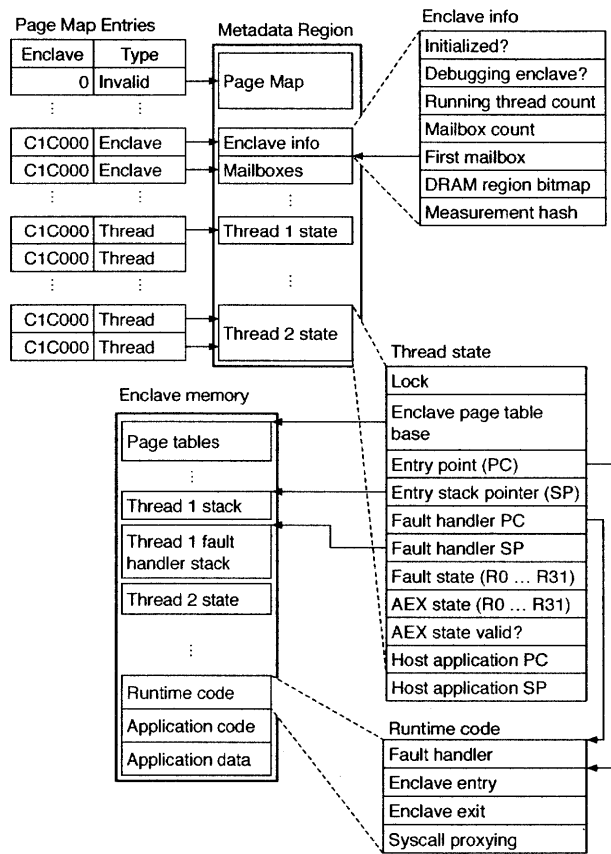


Figure 4-3: Enclave layout and data structures

using its normal entry point, and having the enclave’s code ask the security monitor to restore the pre-AEX execution state. Sanctum enclaves are aware of asynchronous exits so they can implement security policies. For example, an enclave thread that performs time-sensitive work, such as periodic I/O, may terminate itself if it ever gets preempted by an AEX.

Furthermore, whereas SGX dispatches faults (with sanitized information) to the OS, our security monitor dispatches all faults occurring within an enclave to a designated fault handler inside the enclave, which is expected to be implemented by the enclave’s runtime. For example, a `libc` runtime would translate faults into UNIX signals which, by default, would exit the enclave. It is possible, though not advisable for performance reasons (§ 6.3), for a runtime to handle page faults and implement demand paging in a secure manner, without being vulnerable to the attacks described in [204].

Unlike SGX, we isolate each enclave’s data throughout the system’s cache hierarchy. The security monitor flushes per-core caches, such as the L1 cache and the TLB, whenever a core jumps between enclave and non-enclave code. *Last-level cache* (LLC) isolation is achieved by a simple partitioning scheme supported by Sanctum’s hardware extensions (§ 5.1).

Our isolation is also stronger than SGX’s with respect to fault handling. While SGX sanitizes the information that an OS receives during a fault, we achieve full isolation by having the security monitor route the faults that occur inside an enclave to that enclave’s fault handler. This removes all information leaks via the fault timing channel.

Sanctum’s strong isolation yields a simple security model for application developers: *all computation that executes inside an enclave, and only accesses data inside the enclave, is protected from any attack mounted by software outside the enclave.* All communication with the outside world, including accesses to non-enclave memory, is subject to attacks.

We assume that the enclave runtime implements the security measures needed to protect the enclave’s communication with other software modules. For example, any algorithm’s memory access patterns can be protected by ensuring that the algorithm only operates on enclave data. The library can implement this protection simply by copying any input buffer from non-enclave memory into the enclave before computing on it.

The enclave runtime can use Native Client's approach [208] to make sure that the rest of the enclave software only interacts with the host application via the runtime, and mitigate any potential security vulnerabilities in enclave software.

The lifecycle of a Sanctum enclave closely resembles the lifecycle of its SGX equivalent. An enclave is created when its host application performs a system call asking the OS to create an enclave from a dynamically loadable module (.so or .dll file). The OS invokes the security monitor to assign DRAM resources to the enclave, and to load the initial code and data pages into the enclave. Once all the pages are loaded, the enclave is marked as initialized via another security monitor call.

Our software attestation scheme is a simplified version of SGX's scheme, and reuses a subset of its concepts. The data used to initialize an enclave is cryptographically hashed, yielding the enclave's *measurement*. An enclave can invoke a secure inter-enclave messaging service to send a message to a privileged *attestation enclave* that can access the security monitor's attestation key, and produces the attestation signature.

Chapter 5

Hardware Extensions

Sanctum uses an LLC partitioning mechanism that is readily available thanks to the interaction between page tables and direct-mapped or set-associative caches. By manipulating the input to the cache set indexing computation, as described in § 5.1, the computer’s DRAM is divided into equally sized contiguous **DRAM regions** that use disjoint LLC sets. The OS allocates the cache dynamically to its processes and to enclaves, by allocating DRAM regions. Each DRAM region can either be entirely owned by an enclave, or entirely owned by the OS and used by non-enclave processes.

We modify the input to the *page walker* that translates addresses on TLB misses, so it can either use the OS page tables or the current enclave’s page tables, depending on the translated virtual address (§ 5.2). We also add logic between the page walker and the L1 cache, to ensure that the OS page table entries can only point into DRAM regions owned by the OS, and the current enclave’s page table entries can only point into DRAM regions owned by the enclave (§ 5.3).

Lastly, we trust the DMA bus master to reject DMA transfers pointing into DRAM regions allocated to enclaves, to protect against attacks where a malicious OS programs a peripheral to access enclave data. This requires changes similar to the modifications done by SGX and later revisions of TXT to the integrated memory controller on recent Intel chips (§ 5.4).

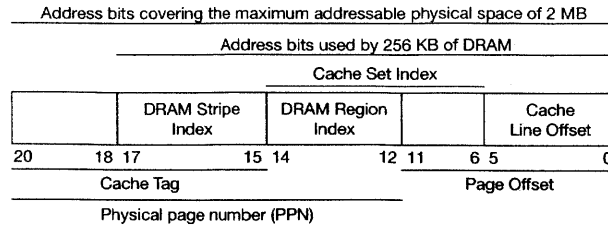


Figure 5-1: Interesting bit fields in a physical address

5.1 LLC Address Input Transformation

Figure 5-1 depicts a physical address in a toy computer with 32-bit virtual addresses and 21-bit physical addresses, 4,096-byte pages, a set-associative LLC with 512 sets and 64-byte lines, and 256 KB of DRAM.

The location where a byte of data is cached in the LLC depends on the low-order bits in the byte’s physical address. The *set index* determines which of the LLC lines can cache the line containing the byte, and the *line offset* locates the byte in its cache line. A virtual address’s low-order bits make up its *page offset*, while the other bits are its *virtual page number* (VPN). Address translation leaves the page offset unchanged, and translates the VPN into a *physical page number* (PPN), based on the mapping specified by the page tables.

We define the **DRAM region index** in a physical address as the intersection between the PPN bits and the cache index bits. This is the maximal set of bits that impact cache placement *and* are determined by privileged software via page tables. We define a **DRAM region** to be the subset of DRAM with addresses having the same DRAM region index. In Figure 5-1, for example, address bits [14 . . . 12] are the DRAM region index, dividing the physical address space into 8 DRAM regions.

DRAM regions are the basis of our cache partitioning because *addresses in a DRAM region do not collide in the LLC with addresses from any other DRAM region*. If programs Alice and Eve use disjoint DRAM regions, they cannot interfere in the LLC, so Eve cannot mount LLC timing attacks on Alice. Furthermore, the OS can place applications in different DRAM regions by manipulating page tables, without having to modify application code.

In a typical system without Sanctum’s hardware extensions, DRAM regions are made up of multiple continuous **DRAM stripes**, where each stripe is exactly one page long. The top

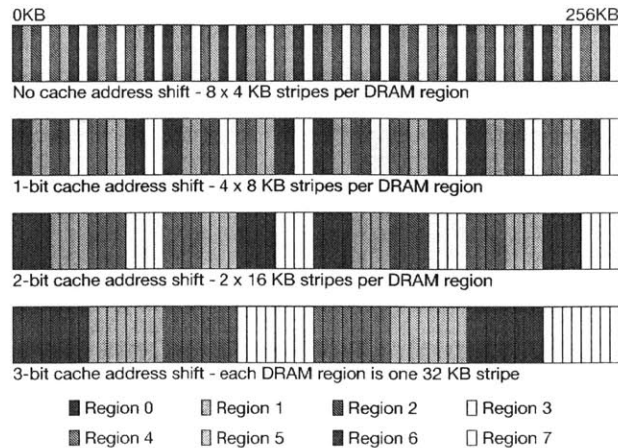


Figure 5-2: Cache address shifting makes DRAM regions contiguous

of Figure 5-2 drives this point home, by showing the partitioning of our toy computer’s 256 KB of DRAM into DRAM regions. The fragmentation of DRAM regions makes it difficult for the OS to allocate contiguous DRAM buffers, which are essential to the efficient DMA transfers used by high performance devices. In our example, if the OS only owns 4 DRAM regions, the largest contiguous DRAM buffer it can allocate is 16 KB.

We observed that, up to a certain point, circularly shifting (rotating) the PPN of a physical address to the right by one bit, before it enters the LLC, doubles the size of each DRAM stripe and halves the number of stripes in a DRAM region, as illustrated in Figure 5-2.

Sanctum takes advantage of this effect by adding a **cache address shifter** that circularly shifts the PPN to the right by a certain amount of bits, as shown in Figures 5-3 and 5-5. In our example, configuring the cache address shifter to rotate the PPN by 3 yields contiguous DRAM regions, so an OS that owns 4 DRAM regions could hypothetically allocate a contiguous DRAM buffer covering half of the machine’s DRAM.

The cache address shifter’s configuration depends on the amount of DRAM present in the system. If our example computer could have 128 KB - 1 MB of DRAM, its cache address shifter must support shift amounts from 2 to 5. Such a shifter can be implemented via a 3-position variable shifter circuit (series of 8-input MUXes), and a fixed shift by 2 (no logic). Alternatively, in systems with known DRAM configuration (embedded, SoC, etc.), the shift amount can be fixed, and implemented with no logic.

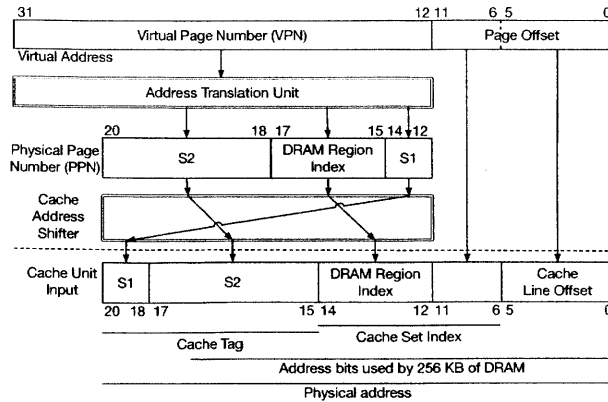


Figure 5-3: Cache address shifter that shifts the PPN by 3 bits

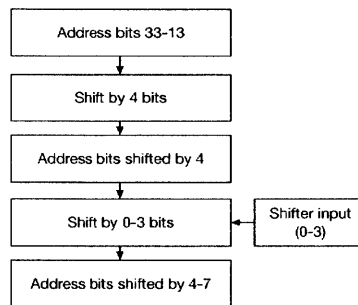


Figure 5-4: A variable shifter that can shift by 2-5 bits can be composed of a fixed shifter by 2 bits and a variable shifter that can shift by 0-3 bits.

5.2 Page Walker Input

Sanctum’s per-enclave page tables require an enclave page table base register `eptbr` that stores the physical address of the currently running enclave’s page tables, and has similar semantics to the page table base register `ptbr` pointing to the operating system-managed page tables. These registers may only be accessed by the Sanctum security monitor, which provides an API call for the OS to modify `ptbr`, and ensures that `eptbr` always points to the current enclave’s page tables.

The circuitry handling TLB misses switches between `ptbr` and `eptbr` based on two registers that indicate the current enclave’s EVRANGE, namely `evbase` (enclave virtual address space base) and `evmask` (enclave virtual address space mask). When a TLB miss occurs, the circuit in Figure 5-6 selects the appropriate page table base by ANDing the faulting virtual address with the mask register and comparing the output against the base

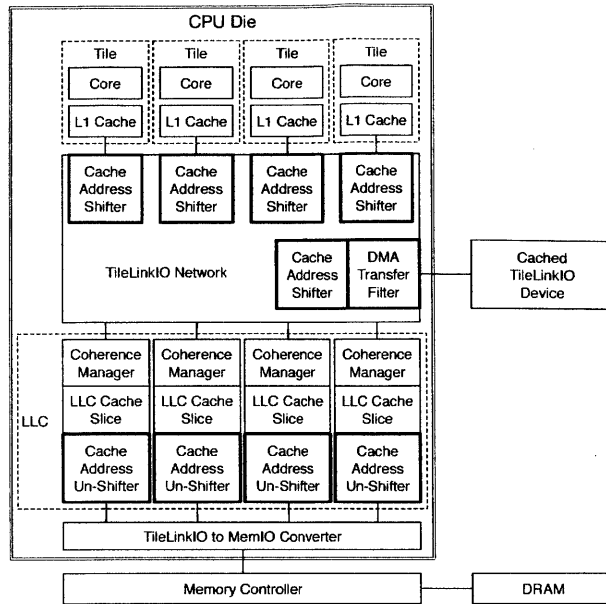


Figure 5-5: Sanctum’s cache address shifter and DMA transfer filter logic in the context of a RISC V-Rocket uncore

register. Depending on the comparison result, either `eptbr` or `ptbr` is forwarded to the page walker as the page table base address.

In addition to the page table base registers, Sanctum uses 4 more pairs of registers that will be described in the next section. On a 64-bit RISC-V computer, the modified FSM input requires 7 extra 51-bit registers (the bottom 13 bits in a 64-bit page-aligned address will always be zero), 2 extra 13-bit registers, 51 AND gates, a 51-bit wide equality comparator (51 XNOR gates and 50 AND gates), and 217 ($51 \times 4 + 13$) 2-bit MUXes.

5.3 Page Walker Memory Accesses

In modern high-speed CPUs, address translation is performed by a hardware **page walker** that traverses the page tables when a TLB miss occurs. The page walker’s latency greatly impacts the CPU’s performance, so it is implemented as a finite-state machine (FSM) that reads page table entries by issuing DRAM read requests using physical addresses, over a dedicated bus to the L1 cache.

Unsurprisingly, page walker modifications require a lot of engineering effort. At the

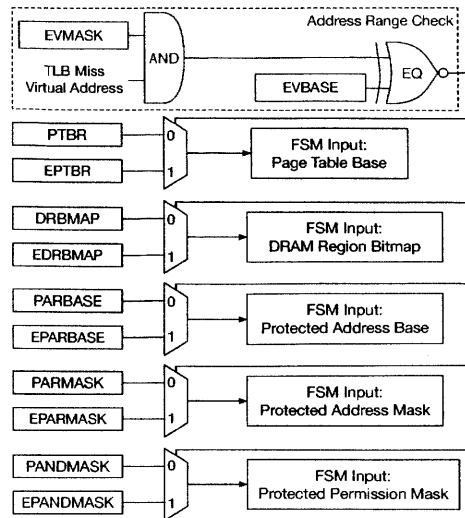


Figure 5-6: Page walker input for per-enclave page tables

same time, Sanctum’s security model demands that the page walker only references enclave memory when traversing the enclave page tables, and only references OS memory when translating the OS page tables. Fortunately, we can satisfy these requirements without modifying the FSM. Instead, the security monitor works in concert with the circuit in Figure 5-7 to ensure that the page tables only point into allowable memory.

Sanctum’s security monitor must guarantee that `ptbr` points into an OS DRAM region, and `eptbr` points into a DRAM region owned by the enclave. This secures the page walker’s initial DRAM read. The circuit in Figure 5-7 receives each page table entry fetched by the FSM, and sanitizes it before it reaches the page walker FSM.

The security monitor configures the set of DRAM regions that page tables may reference by writing to a DRAM region bitmap (`drbmap`) register. The sanitization circuitry extracts the DRAM region index from the address in the page table entry, and looks it up in the DRAM region bitmap. If the address does to belong to an allowable DRAM region, the sanitization logic forces the page table entry’s valid bit to zero, which will cause the page walker FSM to abort the address translation and signal a page fault.

Sanctum’s security monitor must maintain metadata about each enclave, and does so in the enclave’s DRAM regions. For security reasons, the metadata must not be writable by the enclave. Sanctum extends the page table entry transformation described above to implement

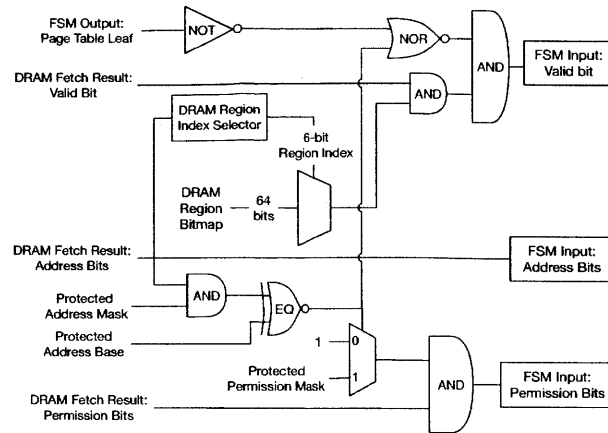


Figure 5-7: Hardware support for per-enclave page tables: transforming the page table entries fetched by the page walker.

per-enclave read-only areas. A protected address range base (`parbase`) register and a protected address range mask (`parmask`) register denote this protected physical address range.

The page table entry sanitization logic in Sanctum’s hardware extensions checks if each page table entry points into the protected address range by ANDing the entry’s address with the protected range mask and comparing the result with the protected range base.

If a leaf page table entry is seen with a protected address, its permission bits are masked with a protected permissions mask (`parpmask`) register. Upon discovering a protected address in an intermediate page table entry, its valid bit is cleared, forcing a page fault.

The above transformation allows the security monitor to set up a read-only range by clearing permission bits (write-enable, for example). Entry invalidation ensures no page table entries are fetched from the protected range, which prevents the page walker FSM from modifying the protected region by setting accessed and dirty bits.

All registers mentioned above come in pairs, as Sanctum maintains separate OS and enclave page tables. The security monitor sets up a protected range in the OS page tables to isolate its own code and data structures (most importantly its private attestation key) from a malicious OS.

Figure 5-8 shows Sanctum’s logic inserted between the page walker and the cache unit that fetches page table entries.

Assuming a 64-bit RISC-V and the example cache above, the logic requires a 64-bit MUX, 65 AND gates, a 51-bit wide equality comparator (51 XNOR gates and 50 AND gates), 13 2-bit MUXs, a 1-bit NOT gate, a 1-bit NOR gate, and a copy of the DRAM region index extraction logic in § 5.1, which could be just wire re-routing if the DRAM configuration is known a priori.

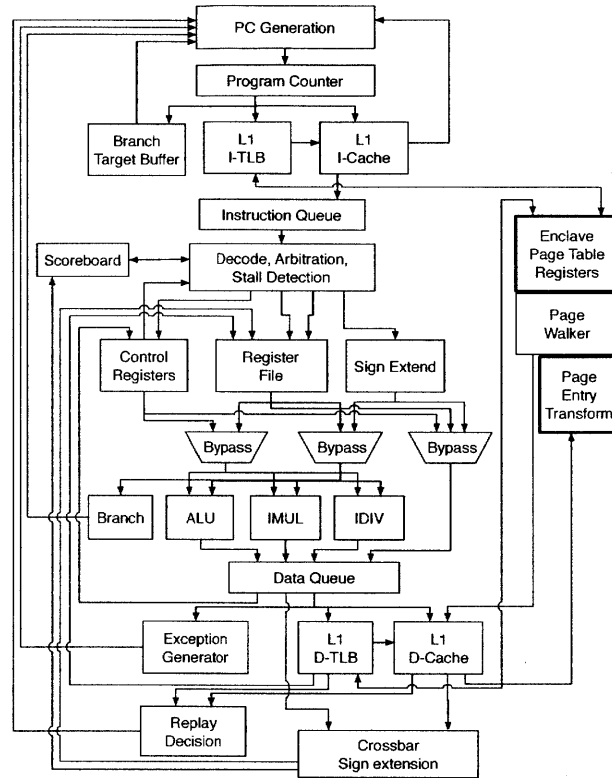


Figure 5-8: Sanctum’s page entry transformation logic in the context of a RISC V Rocket core

5.4 DMA Transfer Filtering

We whitelist a DMA-safe DRAM region instead of following SGX’s blacklist approach. Specifically, Sanctum adds two registers (a base, `dmabase` and an AND mask, `dmarmask`) to the DMA arbiter (memory controller). The range check circuit shown in Figure 5-6 compares each DMA transfer’s start and end addresses against the allowed DRAM range, and the DMA arbiter drops transfers that fail the check.

Chapter 6

Software Design

Sanctum’s chain of trust, discussed in § 6.1, diverges significantly from SGX. We replace SGX’s microcode with a software security monitor that runs at a higher privilege level than the hypervisor and the OS. On RISC-V, the security monitor runs at machine level. Our design only uses one privileged enclave, the signing enclave, which behaves similarly to SGX’s Quoting Enclave.

6.1 Attestation Chain of Trust

Sanctum has three pieces of trusted software: the measurement root, which is burned in on-chip ROM, the security monitor (§ 6.2), which must be stored alongside the computer’s firmware (usually in flash memory), and the signing enclave, which can be stored in any untrusted storage that the OS can access.

We expect the trusted software to be amenable to formal verification: our implementation of a security monitor for Sanctum has fewer than 5 kloc of C++, including a subset of the standard library and the cryptography used for enclave attestation.

6.1.1 The Measurement Root

The measurement root (`mroot`) is stored in a ROM at the top of the physical address space, and covers the reset vector. Its main responsibility is compute a cryptographic hash of the

security monitor and generate a monitor attestation key pair and certificate based on the monitor's hash, as shown in Figure 6-1.

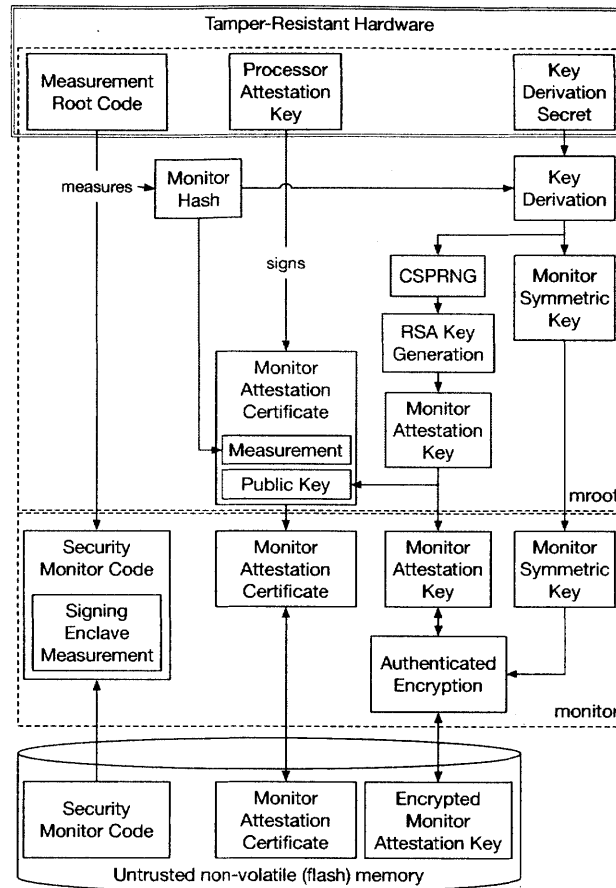


Figure 6-1: Sanctum's root of trust is a measurement root routine burned into the CPU's ROM. This code reads the security monitor from flash memory and generates an attestation key and certificate based on the monitor's hash. Asymmetric key operations, colored in blue, are only performed the first time a monitor is used on a computer.

The security monitor is expected to be stored in non-volatile memory (such as an SPI flash chip) that can respond to memory I/O requests from the CPU, perhaps via a special mapping in the computer's chipset. When `mroot` starts executing, it computes a cryptographic hash over the security monitor. `mroot` then reads the processor's key derivation secret, and derives a symmetric key based on the monitor's hash. `mroot` will eventually hand down the key to the monitor.

The security monitor contains a header that includes the location of an attestation key

existence flag. If the flag is not set, the measurement root generates a monitor attestation key pair, and produces a monitor attestation certificate by signing the monitor's public attestation key with the processor's private attestation key. The monitor attestation certificate includes the monitor's hash.

`mroot` generates a symmetric key for the security monitor so that it can encrypt its private attestation key and store it in the computer's SPI flash memory chip. When writing the key, the monitor also sets the monitor attestation key existence flag, instructing future boot sequences not to re-generate a key. The public attestation key and certificate can be stored unencrypted in any untrusted memory.

Before handing control to the monitor, `mroot` sets a lock that blocks any software from reading the processor's symmetric key derivation seed and private key until a reset occurs. This prevents a malicious security monitor from deriving a different monitor's symmetric key, or from generating a monitor attestation certificate that includes a different monitor's measurement hash.

The symmetric key generated for the monitor is similar in concept to the Seal Keys produced by SGX's key derivation process, as it is used to securely store a secret (the monitor's attestation key) in untrusted memory, in order to avoid an expensive process (asymmetric key attestation and signing). Sanctum's key derivation process is based on the monitor's measurement, so a given monitor is guaranteed to get the same key across power cycles. The cryptographic properties of the key derivation process guarantee that a malicious monitor cannot derive the symmetric key given to another monitor.

6.1.2 The Signing Enclave

In order to avoid timing attacks, the security monitor does not compute attestation signatures directly. Instead, the signing algorithm is executed inside a signing enclave, which is protected by the same isolation guarantees that any other Sanctum enclave enjoys.

The signing enclave receives the monitor's private attestation key via an API call. When the security monitor receives the call, it compares the calling enclave's measurement with a hard-coded value and, upon a successful match, it copies its attestation key into the enclave's

memory. It is worth noting that the monitor only acts on public information (whether the calling enclave is the signing enclave or not), and its memory access pattern does not depend on the private attestation key.

Sanctum's signing enclave authenticates another enclave on the computer and securely receives its attestation data using mailboxes (§ 6.2.5), a simplified version of SGX's local attestation (reporting) mechanism. The enclave's measurement and attestation data are wrapped into a software attestation signature that can be examined by a remote verifier. Figure 6-2 shows the chain of certificates and signatures involved in an instance of software attestation.

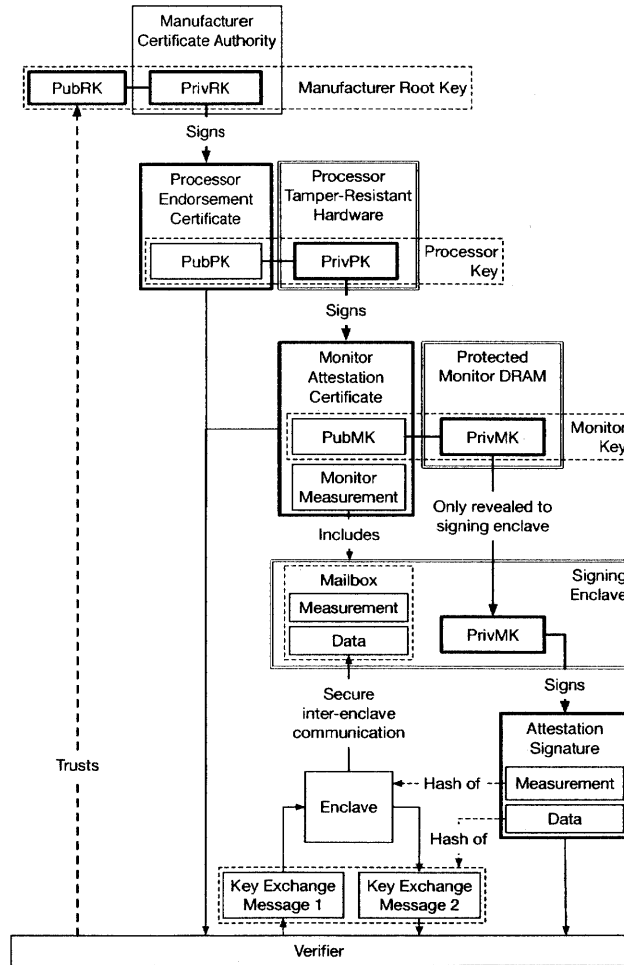


Figure 6-2: The certificate chain behind Sanctum's software attestation signatures



Figure 6-3: DRAM region allocation states and API calls

6.2 Security Monitor

The security monitor receives control after `mroot` finishes setting up the attestation measurement chain.

The monitor provides API calls to the operating system and enclaves for **DRAM region allocation** and **enclave management**. The monitor guards sensitive registers, such as the page table base register (`ptbr`) and the allowed DMA range (`dmabase` and `dmarmask`). The OS can set these registers via monitor calls that ensure the register values are consistent with the current DRAM region allocation.

6.2.1 DRAM Regions

Figure 6-3 shows the DRAM region allocation state transition diagram. After the system boots up, all DRAM regions are allocated to the OS, which can free up DRAM regions so it can re-assign them to enclaves or to itself. A DRAM region can only become free after it is blocked by its owner, which can be the OS or an enclave. While a DRAM region is blocked, any address translations mapping to it cause page faults, so no new TLB entries will be created for that region. Before the OS frees the blocked region, it must flush all the cores' TLBs, to remove any stale entries for the region.

The monitor ensures that the OS performs TLB shutdowns, using a global *block clock*. When a region is blocked, the block clock is incremented, and the current block clock value is stored in the metadata associated with the DRAM region (shown in Figure 4-3). When a core's TLB is flushed, that core's flush time is set to the current block clock value. When the OS asks the monitor to free a blocked DRAM region, the monitor verifies that no core's flush time is lower than the block clock value stored in the region's metadata. As an optimization, freeing a region owned by an enclave only requires TLB flushes on the cores running that enclave's threads. No other core can have TLB entries for the enclave's memory.

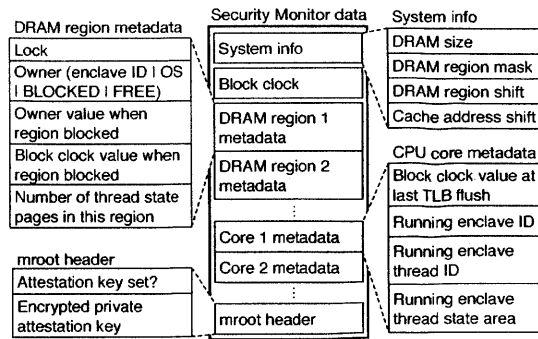


Figure 6-4: Security monitor data structures

The region blocking mechanism guarantees that when a DRAM region is assigned to an enclave or the OS, no stale TLB mappings associated with the DRAM region exist. The monitor uses the MMU extensions described in § 5.2 and § 5.3 to ensure that once a DRAM region is assigned, no software other than the region’s owner may create TLB entries pointing inside the DRAM region. Together, these mechanisms guarantee that the DRAM regions allocated to an enclave cannot be accessed by the operating system or by another enclave.

6.2.2 Metadata Regions

Since the security monitor acts sits between the OS and enclave, and its APIs can be invoked by both sides, it is an easy target for timing attacks. We prevent these attacks with a straightforward policy that states the security monitor is never allowed to access enclave data, and is not allowed to make memory accesses that depend on the attestation key material. The rest of the data handled by the monitor is derived from the OS’ actions, so it is already known to the OS.

A rather obvious consequence of the policy above is that after the security monitor boots the OS, it cannot perform any cryptographic operations that use keys. For example, the security monitor cannot compute an attestation signature directly, and defers that operation to a signing enclave (§ 6.1.2). While it is possible to implement some cryptographic primitives without performing data-dependent accesses, the security and correctness proofs behind these implementations are highly non-trivial. For this reason, Sanctum avoids depending on

any such implementation.

A more subtle aspect of the access policy outlined above is that the metadata structures that the security monitor uses to operate enclaves cannot be stored in DRAM regions owned by enclaves, because that would give the OS an indirect method of accessing the LLC sets that map to enclave's DRAM regions, which could facilitate a cache timing attack.

For this reason, the security monitor requires the OS to set aside at least one DRAM region for enclave metadata before it can create enclaves. The OS has the ability to free up the metadata DRAM region, and regain the LLC sets associated with it, if it predicts that the computer's workload will not involve enclaves.

Each DRAM region that holds enclave metadata is managed independently from the other regions, at page granularity. The first few pages of each region contain a page map that tracks the enclave that tracks the usage of each metadata page, specifically the enclave that it is assigned to, and the data structure that it holds.

Each metadata region is like an EPC region in SGX, with the exception that our metadata regions only hold special pages, like Sanctum's equivalent of SGX's Secure Enclave Control Structure (SECS) and the Thread Control Structure (TCS). These structures will be described in the following sections.

The data structures used to store Sanctum's metadata span multiple pages. When the OS allocates such a structure in a metadata region, it must point the monitor to a sequence of free pages that belong to the same DRAM region. All the pages needed to represent the structure are allocated and released in one API call.

6.2.3 Enclave Lifecycle

The lifecycle of a Sanctum enclave is very similar to that of its SGX counterparts, as shown in Figure 6-5.

The OS creates an enclave by issuing a *create enclave* call that creates the enclave metadata structure, which is Sanctum's equivalent of the SECS. The enclave metadata structure contains an array of mailboxes whose size is established at enclave creation time, so the number of pages required by the structure varies from enclave to enclave. § 6.2.5

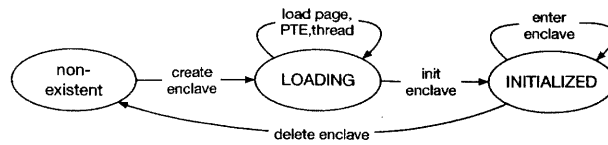


Figure 6-5: Enclave states and enclave management API calls

describes the contents and use of mailboxes.

The *create enclave* API call initializes the enclave metadata fields shown in Figure 4-3, and place the enclave in the **LOADING** state. While the enclave is in this state, the OS sets up the enclave’s initial state via monitor calls that assign DRAM regions to the enclave, create hardware threads and page table entries, and copy code and data into the enclave. The OS then issues a monitor call to transition the enclave to the **INITIALIZED** state, which finalizes its measurement hash. The application hosting the enclave is now free to run enclave threads.

Sanctum stores a measurement hash for each enclave in its metadata area, and updates the measurement to account for every operation performed on an enclave in the **LOADING** state. The policy described in § 6.2.2 does not apply to the secure hash operations used to update enclave’s measurement, because all the data used to compute the hash is already known to the OS.

Enclave metadata is stored in a metadata region (§ 6.2.2), so it can only be accessed by the security monitor. Therefore, the metadata area can safely store public information with integrity requirements, such as the enclave’s measurement hash.

While an OS loads an enclave, it is free to map the enclave’s pages, but the monitor maintains its page tables ensuring all entries point to non-overlapping pages in DRAM owned by the enclave. Once an enclave is initialized, it can inspect its own page tables and abort if the OS created undesirable mappings. Simple enclaves do not require specific mappings. Complex enclaves are expected to communicate their desired mappings to the OS via out-of-band metadata not covered by this work.

Our monitor makes sure that page tables do not overlap by storing the last mapped page’s physical address in the enclave’s metadata. To simplify the monitor, a new mapping is allowed if its physical address is greater than the last mapping’s address, which constrains

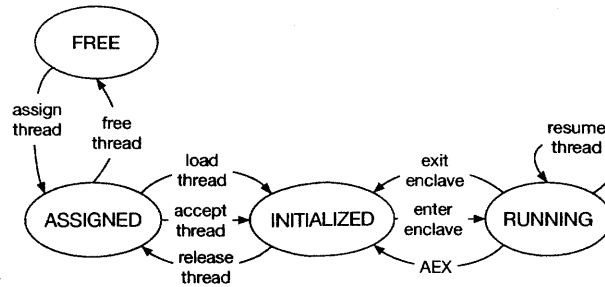


Figure 6-6: Enclave thread metadata structure states and thread-related API calls

the OS to map an enclave’s DRAM pages in monotonically increasing order.

6.2.4 Enclave Code Execution

Sanctum closely follows the threading model of SGX enclaves. Each CPU core that executes enclave code uses a thread metadata structure, which is our equivalent of SGX’s TCS combined with SGX’s State Save Area (SSA). Thread metadata structures are stored in a DRAM region dedicated to enclave metadata in order to prevent a malicious OS from mounting timing attacks against an enclave by causing AEXes on its threads. Figure 6-6 shows the lifecycle of a thread metadata structure.

The OS turns a sequence of free pages in a metadata region into an uninitialized thread structure via an *allocate thread* monitor call. During enclave loading, the OS uses a *load thread* monitor call to initialize the thread structure using data that contributes to the enclave’s measurement. After an enclave is initialized, it can use an *accept thread* monitor call to initialize a thread structure that is allocated to it.

The application hosting an enclave starts executing enclave code by issuing an *enclave enter* API call, which must specify an initialized thread structure. The monitor honors this call by configuring Sanctum’s hardware extensions to allow access to the enclave’s memory, and then by loading the program counter and stack pointer registers from the thread’s metadata structure. The enclave’s code can return control to the hosting application voluntarily, by issuing an *enclave exit* API call, which restores the application’s PC and SP from the thread state area and sets the API call’s return value to `ok`.

When performing an AEX, the security monitor atomically tests-and-sets the *AEX state*

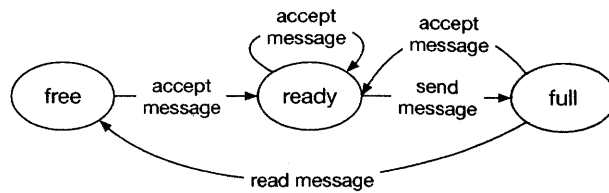


Figure 6-7: Mailbox states and security monitor API calls related to inter-enclave communication

valid flag in the current thread’s metadata. If the flag is clear, the monitor stores the core’s execution state in the thread state’s AEX area. Otherwise, the enclave thread was resuming from an AEX, so the monitor does not change the AEX area. When the host application re-enters the enclave, it will resume from the previous AEX. This reasoning avoids the complexity of SGX’s state stack.

If an interrupt occurs while the enclave code is executing, the security monitor’s exception handler performs an AEX, which sets the API call’s return value to `async_exit`, and invokes the standard interrupt handling code. After the OS handles the interrupt, the enclave’s host application resumes execution, and re-executes the *enter enclave* API call. The enclave’s thread initialization code examines the saved thread state, and seeing that the thread has undergone an AEX, issues a *resume thread* API call. The security monitor restores the enclave’s registers from the thread state area, and clears the AEX flag.

6.2.5 Mailboxes

Sanctum’s software attestation process relies on *mailboxes*, which are a simplified version of SGX’s local attestation mechanism. We could not follow SGX’s approach because it relies on key derivation and MAC algorithms, and our timing attack avoidance policy (§ 6.2.2) states that the security monitor is not allowed to perform cryptographic operations that use keys.

Each enclave’s metadata area contains an array of mailboxes, whose size is specified at enclave creation time, and covered by the enclave’s measurement. Each mailbox goes through the lifecycle shown in Figure 6-7.

An enclave that wishes to receive a message in a mailbox, such as the signing enclave,

declares its intent by performing an *accept message* monitor call. The API call is used to specify the mailbox that will receive the message, and the identity of the enclave that is expected to send the message.

The sending enclave, which is usually the enclave wishing to be authenticated, performs a *send message* call that specifies the identity of the receiving enclave, and a mailbox within that enclave. The monitor only delivers messages to mailboxes that expect them.

When the receiving enclave is notified via an out-of-band mechanism that it has received a message, it issues a *read message* call to the monitor, which moves the message from the mailbox into the enclave's memory. If the API call succeeds, the receiving enclave is assured that the message was sent by the enclave whose identity was specified in the *accept message* call.

Enclave mailboxes are stored in metadata regions (§ 6.2.2), which cannot be accessed by any software other than the security monitor. This guarantees the privacy, integrity, and freshness of the messages sent via the mailbox system.

Our mailbox design has the downside that both the sending and receiving enclave need to be alive in DRAM in order to communicate. By comparison, SGX's local attestation can be done asynchronously. In return, mailboxes do not require any cryptographic operations, and have a much simpler correctness argument.

6.2.6 Multi-Core Concurrency

The security monitor is highly concurrent, with fine-grained locks. API calls targeting two different enclaves may be executed in parallel on different cores. Each DRAM region has a lock guarding that region's metadata. An enclave is guarded by the lock of the DRAM region holding its metadata. Each thread metadata structure also has a lock guarding it, which is acquired when the structure is accessed, but also while the metadata structure is used by a core running enclave code. Thus, the *enter enclave* call acquires a slot lock, which is released by an *enclave exit* call or by an AEX.

We avoid deadlocks by using a form of optimistic locking. Each monitor call attempts to acquire all the locks it needs via atomic test-and-set operations, and errors with a

`concurrent_call` code if any lock is unavailable.

6.3 Enclave Eviction

General-purpose software can be enclaved without source code changes, provided that it is linked against a runtime (e.g., *libc*) modified to work with Sanctum. Any such runtime would be included in the TCB.

The current Sanctum design allows the operating system to over-commit physical memory allocated to enclaves, by paging out to disk DRAM regions from some enclaves. Sanctum does not give the OS visibility into enclave memory accesses, in order to prevent private information leaks, so the OS must decide the enclave whose DRAM regions will be evicted based on other activity, such as network I/O, or based on a business policy, such as Amazon EC2's spot instances.

Once a victim enclave has been decided, the OS asks the enclave to block a DRAM region, giving the enclave an opportunity to rearrange data in its RAM regions. DRAM region management can be transparent to the programmer if handled by the enclave's runtime.

The security monitor does not allow the OS to forcibly reclaim a single DRAM region from an enclave, as doing so would leak memory access patterns. Instead, the OS can delete an enclave, after stopping its threads, and reclaim its DRAM regions. Thus, a small or short-running enclave may well refuse DRAM region management requests from the OS, and expect the OS to delete and re-run it under memory pressure.

To avoid wasted work, large long-running enclaves may elect to use demand paging to overcommit their DRAM, albeit with the understanding that demand paging leaks page-level access patterns to the OS. Securing this mechanism requires the enclave to obfuscate its page faults via periodic I/O using oblivious RAM techniques, as in the Ascend processor [56], applied at page rather than cache line granularity. This carries a high overhead: even with a small chance of paging, an enclave must generate periodic page faults, and access a large set of pages at each period. Using an analytic model, we estimate the overhead to be upwards of 12ms per page per period for a high-end 10K RPM drive, and 27ms for a value

hard drive. Given the number of pages accessed every period grows with an enclave's data size, the costs are easily prohibitive: an enclave accessing pages every second may struggle to make forward progress. While SSDs may alleviate some of this prohibitive overhead, and will be investigated in future work, Sanctum focuses on securing enclaves without demand paging.

Enclaves that perform other data-dependent communication, such as targeted I/O into a large database file, must also use the periodic oblivious I/O to obfuscate their access patterns from the operating system. These techniques are independent of application business logic, and can be provided by libraries such as database access drivers.

Lastly, the presented design requires each enclave to always occupy at least one DRAM region, which contains enclave data structures and the memory management code described above. Evicting all of a live enclave's memory requires an entirely different scheme to be described in future work.

Briefly, the OS can ask the security monitor to *freeze* an enclave, encrypting the enclave's DRAM regions in place, and creating a leaf node in a hash tree. When the monitor *thaws* a frozen enclave, it uses the hash tree leaf to ensure freshness, decrypts the DRAM regions, and *relocates* the enclave, updating its page tables to reflect new owners of relevant DRAM regions. The hash tree is managed by the OS using an approach similar to SGX's version array page eviction.

Chapter 7

Security Argument

Our security argument rests on two pillars, which are the enclave isolation enforced by the security monitor, and the guarantees behind the software attestation signature. This section outlines correctness arguments for each of these pillars.

Sanctum's security monitor isolates enclaves in the processor's caches and in the system's DRAM. This protects each enclave in a system from the (potentially compromised) operating system and from other, potentially malicious, enclaves. We first argue that the monitor effectively isolates enclaves in DRAM, and then reason about the monitor's cache isolation guarantees.

The correctness proof behind our DRAM isolation is very similar to the proof behind SGX's memory access protection scheme presented in § C. Our DRAM regions are equivalent to SGX's EPC pages, and our DRAM region accounting structures, shown in Fig 6-4, are equivalent to SGX's EPCM. While SGX implements its access control in microcode invoked by the PMH, our monitor relies on the hardware implementation described in § 5.3, because RISC cores generally do not have a microcode feature.

Our design avoids passive memory mapping attacks by putting each enclave in charge of mapping its DRAM regions using its own page tables, which removes the need for many of SGX's EPCM-related checks. Instead, Sanctum's memory access checks ensure that an enclave's page tables only point into DRAM regions assigned to it by the OS, and that the OS page tables only point into DRAM regions owned by it. Abstracting this difference away, Sanctum and SGX use very similar methods for access control and provable TLB

shutdown.

Sanctum uses different isolation schemes to address the per-core caches and the shared LLC. The per-core caches are flushed at every transition between enclave and non-enclave mode, which clearly means an enclave will never share a cache with the OS or with another entity. The LLC isolation scheme, presented in § 5.1, ensures that two software modules that do not use the same DRAM region are completely isolated in the LLC. Therefore, LLC isolation for enclaves follows from the DRAM isolation argument presented above.

The OS and its untrusted software is a special case for the LLC isolation scheme, as it is stored in DRAM regions that can be accessed by enclaves. Sanctum does not protect the OS from cache timing attacks. Our security monitor is a special case for both the DRAM and the LLC isolation schemes, as its code and data are stored in DRAM regions that are allocated to the OS.

The monitor protects itself from direct memory probing attacks from the OS by configuring the MMU extensions presented in § 5.2 and § 5.3 to declare a DRAM range that the OS page tables are not allowed to point into. SGX's microcode enjoys the same protection simply by the virtue of being burned into a ROM that is not accessible by software.

Sanctum's security monitor does not protect itself from cache timing attacks. Instead, it follows the memory access policy outlined in § 6.2.2, which ensures that a timing attack on the monitor will not reveal any information that the OS does not already have. By comparison, SGX's microcode does not have this concern, because microcode accesses bypass the caches. However, SGX's architectural enclaves, such as its quoting enclave, must operate under the same regime as Sanctum's monitor, as SGX does not guarantee cache isolation to its enclaves.

Our security monitor stores enclave metadata into dedicated DRAM regions, so the OS can only (indirectly) reach an enclave's LLC sets by performing API calls in very specific situations that require consent from the enclave. For example, the OS can cause evictions in an enclave's LLC sets by having a CPU core execute the enclave's code, via a successful *enter enclave* API call. However, when the API call fails (e.g., because the enclave's thread structure is already in use), the only accesses memory inside the metadata regions, and does not impact the enclave's LLC.

Last, Sanctum enclaves run with the privileges of their host application, like their SGX counterparts. Therefore, all the arguments about OS security in § C apply to Sanctum as well. Specifically, a malicious enclave cannot access the OS memory in a way that would be prohibited by the OS page tables, and cannot carry out a DoS attack against the OS.

Chapter 8

Performance Evaluation

While we propose a high-level set of hardware and software to implement Sanctum, we focus our evaluation on the concrete example of a 4-core RISC-V system generated by Rocket Chip [130]. As Sanctum conveniently isolates concurrent workloads from each other, we can examine its overhead by examining individual applications executing on a single core, discounting the effect of other running software.

8.1 Experiment Design

We use a Rocket-Chip generator modified to model Sanctum’s hardware modifications (§ 5). We generate a 4-core 64-bit RISC-V CPU with private 16KB 4-way set associative instruction and data L1 caches. Using a cycle-accurate simulator for this machine, we produce an LLC access trace and post-process it with a cache emulator for a shared 4MB LLC with and without Sanctum’s DRAM region isolation. We accurately compute the program completion time, in cycles, for each benchmark because Rocket cores have in-order single issue pipelines, and cannot make any progress on a TLB or cache miss. We use a simple model for DRAM accesses and assume unlimited DRAM bandwidth. We also omit an evaluation of the on-chip network and cache coherence overhead, as we do not make any changes that impact any of these subsystems.

Our cache size choices are informed by Intel’s Sandy Bridge [99] desktop models, which have 8 logical CPUs on a 4-core hyper-threaded system with 32KB 8-way L1s, and an

8MB LLC. We do not model Intel’s 256KB per-core L2, because hierarchical caches are not implemented by Rocket. We note, however, that a private L2 would greatly reduce each core’s LLC traffic, which is Sanctum’s main overhead.

We simulate a machine with 4GB of memory that is divided into 64 DRAM regions by Sanctum’s cache address indexing scheme. In our model, an LLC access adds a 12-cycle latency, and a DRAM access costs an additional 100 cycles.

Using the hardware model above, we benchmark the subset of SPECINT 2006 [11] that we could compile using the RISC-V toolchain without additional infrastructure, specifically `bzip2`, `gcc`, `lbm`, `mcf`, `milc`, `sjeng`, and `998.specrand`. This is a mix of memory and compute-bound long-running workloads with diverse access locality.

We choose not to simulate a complete Linux kernel, and instead use the RISC-V proto kernel [193] that provides the few services used by our benchmarks. We schedule each benchmark on Core 0, and run it to completion, while the other cores are idling.

8.2 Cost of Added Hardware

Sanctum’s hardware changes add relatively few gates to the Rocket chip, but do increase its area and power consumption. Like SGX, we avoid modifying the core’s critical path: while our addition to the page walker (as analyzed in the next section) may increase the latency of TLB misses, it does not increase the Rocket core’s clock cycle, which is competitive with an ARM Cortex-A5 [130].

As illustrated at the gate level in Figures 5-6 and 5-7, we estimate Sanctum’s changes to the Rocket hardware to require 1500 (+0.78%) gates and 700 (+1.9%) flip-flops per core, consisting of 50 gates for the cache index calculation, 1000 gates and 700 flip-flops for the extra address page walker configuration, and 400 gates for the page table entry transformations. DMA filtering requires 600 gates (+0.8%) and 128 flip-flops (+1.8%) in the uncore. We do not make any changes to the LLC, and do not include it in the percentages above (the LLC generally accounts for half of chip area).

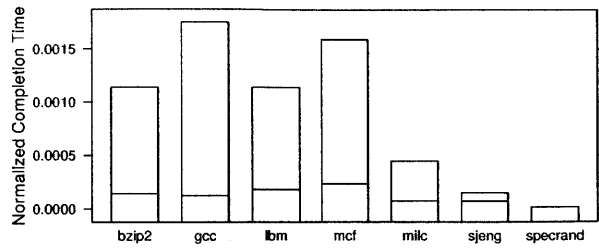


Figure 8-1: Sanctum’s modified page walk has minimal effect on benchmark performance

8.3 Added Page Walker Latency

Sanctum’s page table entry transformation logic is described in § 5.3, and we expect it can be combined with the page walker FSM logic within a single clock cycle.

Nevertheless, in the worst case, the transformation logic would add a pipeline stage between the L1 data cache and the page walker. The transformation logic is small and combinational, significantly simpler than the ALU in the core’s execute stage. In this case, every memory fetch issued by the page walker would experience a 1-cycle latency, which adds 3 cycles of latency to each TLB miss.

Figure 8-1 shows the completion time of selected benchmarks, normalized to the completion time without the extra TLB miss latency. The overheads due to an additional cycle of TLB miss latency are negligible, as quantified by the completion time of SPECINT benchmarks. All overheads are well below 0.01%, relative to the completion time without added TLB latency. This overhead is insignificant relative to the overheads of cache isolation: TLB misses are infrequent and relatively expensive, and a single additional cycle makes little difference.

8.4 Security Monitor Overhead

Invoking Sanctum’s security monitor to load code into an enclave adds a one-time setup cost to each isolated process, when compared against running code without Sanctum’s isolation container. This overhead is amortized by the duration of the computation, so we discount it for long-running workloads.

Entering and exiting enclaves is more expensive than hardware context switches: the

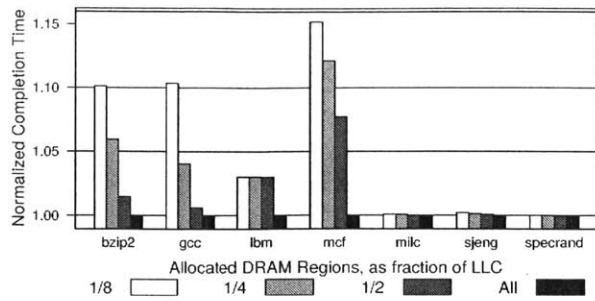


Figure 8-2: The impact of DRAM region allocation on the completion time of an enclaved benchmark, relative to an idea insecure baseline

security monitor must flush TLBs and L1 caches to prevent a privacy leak. However, a sensible OS is expected to minimize the number of context switches by allocating some cores to an enclave and allowing them to execute to completion. We therefore also consider this overhead to be negligible for long-running computations.

8.5 Overhead of DRAM Region Isolation

The crux of Sanctum’s strong isolation is caching DRAM regions in distinct sets. Therefore, when the OS assigns DRAM regions to an enclave, it also confines it to a share of the LLC. An enclaved thread effectively runs on a machine with a smaller LLC, which impacts the enclave’s performance. Note, however, that Sanctum does not partition the per-core caches, so a thread can utilize its core’s entire L1 caches and TLBs.

Figure 8-2 shows the completion times of the SPECINT workloads, each normalized to the completion time of the same benchmark running on an ideal insecure OS that allocates the entire LLC to the benchmark. Sanctum excels at isolating compute-bound workloads operating on sensitive data. Thus, SPECINT’s large, multi-phase workloads heavily exercise the entire memory hierarchy, and therefore paint an accurate picture of a worst case for our system. *mcf*, in particular, is very sensitive to the available LLC size, so it incurs noticeable overheads when being confined to a small subset of the LLC.

We consider *mcf*’s 15.1% decrease in performance when limited to 1/8th of the LLC to be a very pessimistic view of our system’s performance, as it explores the case where the enclave receives 1/4th of the CPU power (a core), but 1/8th of the LLC. For a reasonable

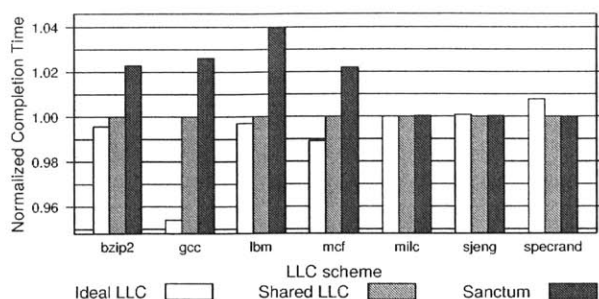


Figure 8-3: Sanctum’s enclave overheads for one core utilizing 1/4 of the LLC compared against an idealized baseline (non-enclaved app using the entire LLC), and against a representative baseline (non-enclaved app sharing the LLC with concurrent instances)

allocation of 1/4 of DRAM regions (in a 4-core system), DRAM regions add a 3-6% overhead to most memory-bound benchmarks (with the exception of `mcf`), and do not impact compute-bound workloads.

We also consider a more representative baseline by considering the performance of a system executing multiple copies of the benchmark concurrently, in different virtual address spaces, effectively normalizing LLC resources available to each instance. Overheads for a reasonable allocation of 1/4th of the LLC shared among 4 instances are shown in Figure 8-3. With this baseline, `mcf` overheads are reduced to 4.6% and 2.2% for allocations of 1/8th and 1/4th of the LLC, respectively. Over the full set of benchmarks, overheads fall below 4%, averaging 1.9%. Memory-bound benchmarks exhibit a 2.7% average overhead over this insecure baseline.

In the LLC, our region-aware cache index translation forces consecutive physical pages in DRAM to map to the same cache sets within a DRAM region, creating interference. We expect the OS memory management implementation to be aware of DRAM regions, and map data structures to pages spanning all available DRAM regions.

The locality of DRAM accesses is also affected: an enclaved process has exclusive access to its DRAM region(s), each a contiguous range of physical addresses. DRAM regions therefore cluster process accesses to physical memory, decreasing the efficiency of bank-level interleaving in a system with multiple DRAM channels. Row or cache line-level interleaving (employed by some Intel processors [99]) of DRAM channels better parallelizes accesses within a DRAM region, but introduces a trade-off in the efficiency of

individual DRAM channels. Considering the low miss rate in a modern cache hierarchy, and multiple concurrent threads, we expect this overhead is small compared to the cost of cache partitioning. We leave a thorough evaluation of DRAM overhead in a multi-channel system for future work.

Chapter 9

Conclusion

We have shown through the design of Sanctum that strong provable isolation of concurrent software modules can be achieved with low overhead. The worst observed overhead across all benchmarks when compared to a representative insecure baseline is 4.6%. This approach provides strong security guarantees against an insidious threat model including cache timing and memory access pattern attacks. With this work, we hope to enable a shift in discourse in secure hardware architecture approaches away from plugging specific security holes to a principled approach to eliminating attack surfaces.

Appendix A

Computer Architecture Background

This section attempts to summarize the general architectural principles behind Intel's most popular computer processors, as well as the peculiarities needed to reason about the security properties of a system running on these processors. Unless specified otherwise, the information here is summarized from Intel's *Software Development Manual (SDM)* [104].

Analyzing the security of a software system requires understanding the interactions between all the parts of the software's execution environment, so this section is quite long. We do refrain from introducing any security concepts here, so readers familiar with x86's intricacies can safely skip this section and refer back to it when necessary.

We use the terms *Intel processor* or *Intel CPU* to refer to the server and desktop versions of Intel's Core line-up. In the interest of space and mental sanity, we ignore Intel's other processors, such as the embedded line of Atom CPUs, or the failed Itanium line. Consequently, the terms *Intel computers* and *Intel systems* refers to computer systems built around Intel's Core processors.

In this paper, the term *Intel architecture* refers to the x86 architecture described in Intel's SDM. The x86 architecture is overly complex, mostly due to the need to support executing legacy software dating back to 1990 directly on the CPU, without the overhead of software interpretation. We only cover the parts of the architecture visible to modern 64-bit software, also in the interest of space and mental sanity.

The 64-bit version of the x86 architecture, covered in this section, was actually invented by Advanced Micro Devices (AMD), and is also known as AMD64, x86_64, and x64.

The term “Intel architecture” highlights our interest in the architecture’s implementation in Intel’s chips, and our desire to understand the mindsets of Intel SGX’s designers.

A.1 Overview

A computer’s main resources (§ A.2) are *memory* and *processors*. On Intel computers, *Dynamic Random-Access Memory* (DRAM) chips (§ A.9.1) provide the memory, and one or more CPU chips expose *logical processors* (§ A.9.4). These resources are managed by *system software*. An Intel computer typically runs two kinds of system software, namely operating systems and hypervisors.

The Intel architecture was designed to support running multiple application software instances, called *processes*. An *operating system* (§ A.3), allocates the computer’s resources to the running processes. Server computers, especially in cloud environments, may run multiple operating system instances at the same time. This is accomplished by having a *hypervisor* (§ A.3) partition the computer’s resources between the operating system instances running on the computer.

System software uses virtualization techniques to isolate each piece of software that it manages (process or operating system) from the rest of the software running on the computer. This isolation is a key tool for keeping software complexity at manageable levels, as it allows application and OS developers to focus on their software, and ignore the interactions with other software that may run on the computer.

A key component of virtualization is address translation (§ A.5), which is used to give software the impression that it owns all the memory on the computer. Address translation provides isolation that prevents a piece of buggy or malicious software from directly damaging other software, by modifying its memory contents.

The other key component of virtualization is the software privilege levels (§ A.3) enforced by the CPU. Hardware privilege separation ensures that a piece of buggy or malicious software cannot damage other software indirectly, by interfering with the system software managing it.

Processes express their computing power requirements by creating execution *threads*,

which are assigned by the operating system to the computer's logical processors. A thread contains an execution context (§ A.6), which is the information necessary to perform a computation. For example, an execution context stores the address of the next instruction that will be executed by the processor.

Operating systems give each process the illusion that it has an infinite amount of logical processors at its disposal, and multiplex the available logical processors between the threads created by each process. Modern operating systems implement *preemptive multithreading*, where the logical processors are rotated between all the threads on a system every few milliseconds. Changing the thread assigned to a logical processor is accomplished by an execution context switch (§ A.6).

Hypervisors expose a fixed number of virtual processors (vCPUs) to each operating system, and also use context switching to multiplex the logical CPUs on a computer between the vCPUs presented to the guest operating systems.

The execution core in a logical processor can execute instructions and consume data at a much faster rate than DRAM can supply them. Many of the complexities in modern computer architectures stem from the need to cover this speed gap. Recent Intel CPUs rely on hyper-threading (§ A.9.4), out-of-order execution (§ A.10), and caching (§ A.11), all of which have security implications.

An Intel processor contains many levels of intermediate memories that are much faster than DRAM, but also orders of magnitude smaller. The fastest intermediate memory is the logical processor's register file (§ A.2, § A.4, § A.6). The other intermediate memories are called caches (§ A.11). The Intel architecture requires application software to explicitly manage the register file, which serves as a high-speed scratch space. At the same time, caches transparently accelerate DRAM requests, and are mostly invisible to software.

Intel computers have multiple logical processors. As a consequence, they also have multiple caches distributed across the CPU chip. On multi-socket systems, the caches are distributed across multiple CPU chips. Therefore, Intel systems use a cache coherence mechanism (§ A.11.3), ensuring that all the caches have the same view of DRAM. Thanks to cache coherence, programmers can build software that is unaware of caching, and still runs correctly in the presence of distributed caches. However, cache coherence does not

cover the dedicated caches used by address translation (§ A.11.5), and system software must take special measures to keep these caches consistent.

CPUs communicate with the outside world via I/O devices (also known as peripherals), such as network interface cards and display adapters (§ A.9). Conceptually, the CPU communicates with the DRAM chips and the I/O devices via a *system bus* that connects all these components.

Software written for the Intel architecture communicates with I/O devices via the I/O address space (§ A.4) and via the memory address space, which is primarily used to access DRAM. System software must configure the CPU's caches (§ A.11.4) to recognize the memory address ranges used by I/O devices. Devices can notify the CPU of the occurrence of events by dispatching interrupts (§ A.12), which cause a logical processor to stop executing its current thread, and invoke a special handler in the system software (§ A.8.2).

Intel systems have a highly complex computer initialization sequence (§ A.13), due to the need to support a large variety of peripherals, as well as a multitude of operating systems targeting different versions of the architecture. The initialization sequence is a challenge to any attempt to secure an Intel computer, and has facilitated many security compromises (§ A.3).

Intel's engineers use the processor's microcode facility (§ A.14) to implement the more complicated aspects of the Intel architecture, which greatly helps manage the hardware's complexity. The microcode is completely invisible to software developers, and its design is mostly undocumented. However, in order to evaluate the feasibility of any architectural change proposals, one must be able to distinguish changes that can be implemented in microcode from changes that can only be accomplished by modifying the hardware.

A.2 Computational Model

This section pieces together a highly simplified model for a computer that implements the Intel architecture, illustrated in Figure A-1. This simplified model is intended to help the reader's intuition process the fundamental concepts used by the rest of the paper. The following sections gradually refine the simplified model into a detailed description of the

Intel architecture.

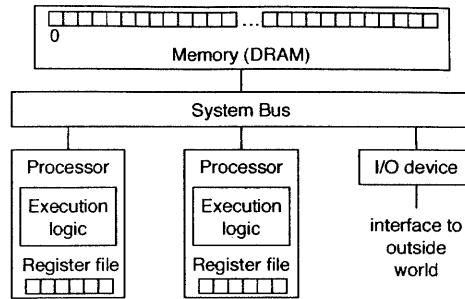


Figure A-1: A computer's core is its processors and memory, which are connected by a system bus. Computers also have I/O devices, such as keyboards, which are also connected to the processor via the system bus.

The building blocks for the model presented here come from [170], which introduces the key abstractions in a computer system, and then focuses on the techniques used to build software systems on top of these abstractions.

The memory is an array of storage cells, addressed using natural numbers starting from 0, and implements the abstraction depicted in Figure A-2. Its salient feature is that the result of reading a memory cell at an address must equal the most value written to that memory cell.

$WRITE(addr, value) \rightarrow \emptyset$ Store <i>value</i> in the storage cell identified by <i>addr</i> .
$READ(addr) \rightarrow value$ Return the <i>value</i> argument to the most recent WRITE call referencing <i>addr</i> .

Figure A-2: The memory abstraction

A logical processor repeatedly reads *instructions* from the computer's memory and executes them, according to the flowchart in Figure A-3.

The processor has an internal memory, referred to as the *register file*. The register file consists of *Static Random Access Memory* (SRAM) cells, generally known as *registers*, which are significantly faster than DRAM cells, but also a lot more expensive.

An instruction performs a simple computation on its inputs and stores the result in an output location. The processor's registers make up an *execution context* that provides the inputs and stores the outputs for most instructions. For example, `ADD RDX, RAX, RBX`

performs an integer addition, where the inputs are the registers RAX and RBX, and the result is stored in the output register RDX.

The registers mentioned in Figure A-3 are the *instruction pointer* (RIP), which stores the memory address of the next instruction to be executed by the processor, and the *stack pointer* (RSP), which stores the memory address of the topmost element in the call stack used by the processor's procedural programming support. The other execution context registers are described in § A.4 and § A.6.

Under normal circumstances, the processor repeatedly reads an instruction from the memory address stored in RIP, executes the instruction, and updates RIP to point to the following instruction. Unlike many RISC architectures, the Intel architecture uses a variable-size instruction encoding, so the size of an instruction is not known until the instruction has been read from memory.

While executing an instruction, the processor may encounter a *fault*, which is a situation where the instruction's preconditions are not met. When a fault occurs, the instruction does not store a result in the output location. Instead, the instruction's result is considered to be the fault that occurred. For example, an integer division instruction `DIV` where the divisor is zero results in a Division Fault (`#DIV`).

When an instruction results in a fault, the processor stops its normal execution flow, and performs the fault handler process documented in § A.8.2. In a nutshell, the processor first looks up the address of the code that will handle the fault, based on the fault's nature, and sets up the execution environment in preparation to execute the fault handler.

The processors are connected to each other and to the memory via a *system bus*, which is a broadcast network that implements the abstraction in Figure A-4.

During each clock cycle, at most one of the devices connected to the system bus can send a message, which is received by all the other devices connected to the bus. Each device attached to the bus decodes the operation codes and addresses of all the messages sent on the bus and ignores the messages that do not require its involvement.

For example, when the processor wishes to read a memory location, it sends a message with the operation code `READ-REQUEST` and the bus address corresponding to the desired memory location. The memory sees the message on the bus and performs the `READ`

operation. At a later time, the memory responds by sending a message with the operation code READ-RESPONSE, the same address as the request, and the data value set to the result of the READ operation.

The computer communicates with the outside world via I/O devices, such as keyboards, displays, and network cards, which are connected to the system bus. Devices mostly respond to requests issued by the processor. However, devices also have the ability to issue *interrupt requests* that notify the processor of outside events, such as the user pressing a key on a keyboard.

Interrupt triggering is discussed in § A.12. On modern systems, devices send interrupt requests by issuing writes to special bus addresses. Interrupts are considered to be *hardware exceptions*, just like faults, and are handled in a similar manner.

A.3 Software Privilege Levels

In an Infrastructure-as-a-Service (IaaS) cloud environment, such as Amazon EC2, commodity CPUs run software at four different privilege levels, shown in Figure A-5.

Each privilege level is strictly more powerful than the ones below it, so a piece of software can freely read and modify the code and data running at less privileged levels. Therefore, a software module can be compromised by any piece of software running at a higher privilege level. It follows that a software module implicitly trusts all the software running at more privileged levels, and a system's security analysis must take into account the software at all privilege levels.

System Management Mode (SMM) is intended for use by the motherboard manufacturers to implement features such as fan control and deep sleep, and/or to emulate missing hardware. Therefore, the bootstrapping software (§ A.13) in the computer's firmware is responsible for setting up a continuous subset of DRAM as *System Management RAM* (SMRAM), and for loading all the code that needs to run in SMM mode into SMRAM. The SMRAM enjoys special hardware protections that prevent less privileged software from accessing the SMM code.

IaaS cloud providers allow their customers to run their operating system of choice in a

virtualized environment. Hardware virtualization [187], called *Virtual Machine Extensions* (VMX) by Intel, adds support for a *hypervisor*, also called a *Virtual Machine Monitor* (VMM) in the Intel documentation. The hypervisor runs at a higher privilege level (VMX root mode) than the operating system, and is responsible for allocating hardware resources across multiple operating systems that share the same physical machine. The hypervisor uses the CPU's hardware virtualization features to make each operating system believe it is running in its own computer, called a *virtual machine* (VM). Hypervisor code generally runs at ring 0 in VMX root mode.

Hypervisors that run in VMX root mode and take advantage of hardware virtualization generally have better performance and a smaller codebase than hypervisors based on binary translation [165].

The systems research literature recommends breaking up an operating system into a small *kernel*, which runs at a high privilege level, known as the *kernel mode* or *supervisor mode* and, in the Intel architecture, as *ring 0*. The kernel allocates the computer's resources to the other system components, such as device drivers and services, which run at lower privilege levels. However, for performance reasons¹, mainstream operating systems have large amounts of code running at ring 0. Their *monolithic kernels* include device drivers, filesystem code, networking stacks, and video rendering functionality.

Application code, such as a Web server or a game client, runs at the lowest privilege level, referred to as *user mode* (*ring 3* in the Intel architecture). In IaaS cloud environments, the virtual machine images provided by customers run in VMX non-root mode, so the kernel runs in VMX non-root ring 0, and the application code runs in VMX non-root ring 3.

A.4 Address Spaces

Software written for the Intel architecture accesses the computer's resources using four distinct physical address spaces, shown in Figure A-6. The address spaces overlap partially, in both purpose and contents, which can lead to confusion. This section gives a high-level overview of the physical address spaces defined by the Intel architecture, with an emphasis

¹Calling a procedure in a different ring is much slower than calling code at the same privilege level.

on their purpose and the methods used to manage them.

The *register* space consists of names that are used to access the CPU's register file, which is the only memory that operates at the CPU's clock frequency and can be used without any latency penalty. The register space is defined by the CPU's architecture, and documented in the SDM.

Some registers, such as the *Control Registers* (CRs) play specific roles in configuring the CPU's operation. For example, CR3 plays a central role in address translation (§ A.5). These registers can only be accessed by system software. The rest of the registers make up an application's *execution context* (§ A.6), which is essentially a high-speed scratch space. These registers can be accessed at all privilege levels, and their allocation is managed by the software's compiler. Many CPU instructions only operate on data in registers, and only place their results in registers.

The *memory* space, generally referred to as *the address space*, or *the physical address space*, consists of 2^{36} (64 GB) - 2^{40} (1 TB) addresses. The memory space is primarily used to access DRAM, but it is also used to communicate with *memory-mapped devices* that read memory requests off a system bus and write replies for the CPU. Some CPU instructions can read their inputs from the memory space, or store the results using the memory space.

A better-known example of memory mapping is that at computer startup, memory addresses 0xFFFFF000 - 0xFFFFFFFF (the 64 KB of memory right below the 4 GB mark) are mapped to a flash memory device that holds the first stage of the code that bootstraps the computer.

The memory space is partitioned between devices and DRAM by the computer's firmware during the bootstrapping process. Sometimes, system software includes motherboard-specific code that modifies the memory space partitioning. The OS kernel relies on address translation, described in § A.5, to control the applications' access to the memory space. The hypervisor relies on the same mechanism to control the guest OSs.

The *input/output (I/O)* space consists of 2^{16} I/O addresses, usually called *ports*. The I/O ports are used exclusively to communicate with devices. The CPU provides specific instructions for reading from and writing to the I/O space. I/O ports are allocated to devices by formal or de-facto standards. For example, ports 0xCF8 and 0xCFC are always used to

access the PCI express (§ A.9.1) configuration space.

The CPU implements a mechanism for system software to provide fine-grained I/O access to applications. However, all modern kernels restrict application software from accessing the I/O space directly, in order to limit the damage potential of application bugs.

The *Model-Specific Register (MSR)* space consists of 2^{32} MSRs, which are used to configure the CPU's operation. The MSR space was initially intended for the use of CPU model-specific firmware, but some MSRs have been promoted to *architectural MSR* status, making their semantics a part of the Intel architecture. For example, architectural MSR 0x10 holds a high-resolution monotonically increasing time-stamp counter.

The CPU provides instructions for reading from and writing to the MSR space. The instructions can only be used by system software. Some MSRs are also exposed by instructions accessible to applications. For example, applications can read the time-stamp counter via the RDTSC and RDTSCP instructions, which are very useful for benchmarking and optimizing software.

A.5 Address Translation

System software relies on the CPU's address translation mechanism for implementing isolation among less privileged pieces of software (applications or operating systems). Virtually all secure architecture designs bring changes to address translation. We summarize the Intel architecture's address translation features that are most relevant when establishing a system's security properties, and refer the reader to [111] for a more general presentation of address translation concepts and its other uses.

A.5.1 Address Translation Concepts

From a systems perspective, address translation is a layer of indirection (shown in Figure A-7) between the *virtual addresses*, which are used by a program's memory load and store instructions, and the *physical addresses*, which reference the physical address space (§ A.4). The mapping between virtual and physical addresses is defined by *page tables*, which are managed by the system software.

Operating systems use address translation to implement the *virtual memory abstraction*, illustrated by Figure A-8. The virtual memory abstraction exposes the same interface as the memory abstraction in § A.2, but each process uses a separate virtual address space that only references the memory allocated to that process. From an application developer standpoint, virtual memory can be modeled by pretending that each process runs on a separate computer and has its own DRAM.

Address translation is used by the operating system to multiplex DRAM among multiple application processes, isolate the processes from each other, and prevent application code from accessing memory-mapped devices directly. The latter two protection measures prevent an application's bugs from impacting other applications or the OS kernel itself. Hypervisors also use address translation, to divide the DRAM among operating systems that run concurrently, and to virtualize memory-mapped devices.

The address translation mode used by 64-bit operating systems, called IA-32e by Intel's documentation, maps 48-bit *virtual addresses* to *physical addresses* of at most 52 bits². The translation process, illustrated in Figure A-9, is carried out by dedicated hardware in the CPU, which is referred to as the *address translation unit* or the *memory management unit* (MMU).

The bottom 12 bits of a virtual address are not changed by the translation. The top 36 bits are grouped into four 9-bit indexes, which are used to index into the page tables. Despite its name, the page tables data structure closely resembles a full 512-ary search tree where nodes have fixed keys. Each node is represented in DRAM as an array of 512 8-byte entries that contain the physical addresses of the next-level children as well as some flags. The physical address of the root node is stored in the CR3 register. The arrays in the last-level nodes contain the physical addresses that are the result of the address translation.

The address translation function, which does not change the bottom bits of addresses, partitions the memory address space into *pages*. A page is the set of all memory locations that only differ in the bottom bits which are not impacted by address translation, so all the memory addresses in a virtual page translate to corresponding addresses in the same physical

²The size of a physical address is CPU-dependent, and is 40 bits for recent desktop CPUs and 44 bits for recent high-end server CPUs.

page. From this perspective, the address translation function can be seen as a mapping between *Virtual Page Numbers* (VPN) and *Physical Page Numbers* (PPN), as shown in Figure A-10.

In addition to isolating application processes, operating systems also use the address translation feature to run applications whose collective memory demands exceed the amount of DRAM installed in the computer. The OS evicts infrequently used memory pages from DRAM to a larger (but slower) memory, such as a hard disk drive (HDD) or solid-state drive (SSD). For historical reason, this slower memory is referred to as the *disk*.

The OS ability to over-commit DRAM is often called *page swapping*, for the following reason. When an application process attempts to access a page that has been evicted, the OS “steps in” and reads the missing page back into DRAM. In order to do this, the OS might have to evict a different page from DRAM, effectively swapping the contents of a DRAM page with a disk page. The details behind this high-level description are covered in the following sections.

The CPU’s address translation is also referred to as “paging”, which is a shorthand for “page swapping”.

A.5.2 Address Translation and Virtualization

Computers that take advantage of hardware virtualization use a hypervisor to run multiple operating systems at the same time. This creates some tension, because each operating system was written under the assumption that it owns the entire computer’s DRAM. The tension is solved by a second layer of address translation, illustrated in Figure A-11.

When a hypervisor is active, the page tables set up by an operating system map between virtual addresses and *guest-physical addresses* in a *guest-physical address space*. The hypervisor multiplexes the computer’s DRAM between the operating systems’ guest-physical address spaces via the second layer of address translations, which uses *extended page tables* (EPT) to map guest-physical addresses to physical addresses.

The EPT uses the same data structure as the page tables, so the process of translating guest-physical addresses to physical addresses follows the same steps as IA-32e address

translation. The main difference is that the physical address of the data structure's root node is stored in the extended page table pointer (EPTP) field in the *Virtual Machine Control Structure* (VMCS) for the guest OS. Figure A-12 illustrates the address translation process in the presence of hardware virtualization.

A.5.3 Page Table Attributes

Each page table entry contains a physical address, as shown in Figure A-9, and some Boolean values that are referred to as *flags* or *attributes*. The following attributes are used to implement page swapping and software isolation.

The *present* (P) flag is set to 0 to indicate unused parts of the address space, which do not have physical memory associated with them. The system software also sets the P flag to 0 for pages that are evicted from DRAM. When the address translation unit encounters a zero P flag, it aborts the translation process and issues a hardware exception, as described in § A.8.2. This hardware exception gives system software an opportunity to step in and bring an evicted page back into DRAM.

The *accessed* (A) flag is set to 1 by the CPU whenever the address translation machinery reads a page table entry, and the *dirty* (D) flag is set to 1 by the CPU when an entry is accessed by a memory write operation. The A and D flags give the hypervisor and kernel insight into application memory access patterns and inform the algorithms that select the pages that get evicted from RAM.

The main attributes supporting software isolation are the *writable* (W) flag, which can be set to 0 to prohibit³ writes to any memory location inside a page, the *disable execution* (XD) flag, which can be set to 1 to prevent instruction fetches from a page, and the *supervisor* (S) flag, which can be set to 1 to prohibit any accesses from application software running at ring 3.

³Writes to non-writable pages result in #GP exceptions (§ A.8.2).

A.6 Execution Contexts

Application software targeting the 64-bit Intel architecture uses a variety of CPU registers to interact with the processor's features, shown in Figure A-13 and Table A.1. The values in these registers make up an application thread's state, or *execution context*.

OS kernels multiplex each logical processor (§ A.9.4) between multiple software threads by *context switching*, namely saving the values of the registers that make up a thread's execution context, and replacing them with another thread's previously saved context. Context switching also plays a part in executing code inside secure containers, so its design has security implications.

Integers and memory addresses are stored in 16 *general-purpose registers* (GPRs). The first 8 GPRs have historical names: RAX, RBX, RCX, RDX, RSI, RDI, RSP, and RBP, because they are extended versions of the 32-bit Intel architecture's GPRs. The other 8 GPRs are simply known as R9-R16. RSP is designated for pointing to the top of the procedure call stack, which is simply referred to as *the stack*. RSP and the stack that it refers to are automatically read and modified by the CPU instructions that implement procedure calls, such as CALL and RET (return), and by specialized stack handling instructions such as PUSH and POP.

All applications also use the RIP register, which contains the address of the currently executing instruction, and the RFLAGS register, whose bits (e.g., the carry flag - CF) are individually used to store comparison results and control various instructions.

Software might use other registers to interact with specific processor features, some of which are shown in Table A.1.

The Intel architecture provides a future-proof method for an OS kernel to save the values of feature-specific registers used by an application. The XSAVE instruction takes in a *requested-feature bitmap* (RFBM), and writes the registers used by the features whose RFBM bits are set to 1 in a memory area. The memory area written by XSAVE can later be used by the XRSTOR instruction to load the saved values back into feature-specific registers. The memory area includes the RFBM given to XSAVE, so XRSTOR does not require an RFBM input.

Feature	Registers	XCR0 bit
FPU	FP0 - FP7, FSW, FTW	0
SSE	MM0 - MM7, XMM0 - XMM15, XMCSR	1
AVX	YMM0 - YMM15	2
MPX	BND0 - BND 3	3
MPX	BNDCFGU, BNDSTATUS	4
AVX-512	K0 - K7	5
AVX-512	ZMM0_H - ZMM15_H	6
AVX-512	ZMM16 - ZMM31	7
PK	PKRU	9

Table A.1: Sample feature-specific Intel architecture registers.

Application software declares the features that it plans to use to the kernel, so the kernel knows what XSAVE bitmap to use when context-switching. When receiving the system call, the kernel sets the XCR0 register to the feature bitmap declared by the application. The CPU generates a fault if application software attempts to use features that are not enabled by XCR0, so applications cannot modify feature-specific registers that the kernel wouldn't take into account when context-switching. The kernel can use the `CPUID` instruction to learn the size of the XSAVE memory area for a given feature bitmap, and compute how much memory it needs to allocate for the context of each of the application's threads.

A.7 Segment Registers

The Intel 64-bit architecture gained widespread adoption thanks to its ability to run software targeting the older 32-bit architecture side-by-side with 64-bit software [174]. This ability comes at the cost of some warts. While most of these warts can be ignored while reasoning about the security of 64-bit software, the segment registers and vestigial segmentation model must be understood.

The semantics of the Intel architecture's instructions include the implicit use of a few segments which are loaded into the processor's *segment registers* shown in Figure A-13. Code fetches use the *code segment* (CS). Instructions that reference the stack implicitly use the *stack segment* (SS). Memory references implicitly use the *data segment* (DS) or the *destination segment* (ES). Via segment override prefixes, instructions can be modified to use

the unnamed segments FS and GS for memory references.

Modern operating systems effectively disable segmentation by covering the entire addressable space with one segment, which is loaded in CS, and one data segment, which is loaded in SS, DS and ES. The FS and GS registers store segments covering *thread-local storage* (TLS).

Due to the Intel architecture's 16-bit origins, segment registers are exposed as 16-bit values, called *segment selectors*. The top 13 bits in a selector are an index in a *descriptor table*, and the bottom 2 bits are the selector's ring number, which is also called requested privilege level (RPL) in the Intel documentation. Also, modern system software only uses rings 0 and 3 (see § A.3).

Each segment register has a hidden *segment descriptor*, which consists of a *base address*, *limit*, and type information, such as whether the descriptor should be used for executable code or data. Figure A-14 shows the effect of loading a 16-bit selector into a segment register. The selector's index is used to read a descriptor from the descriptor table and copy it into the segment register's hidden descriptor.

In 64-bit mode, all segment limits are ignored. The base addresses in most segment registers (CS, DS, ES, SS) are ignored. The base addresses in FS and GS are used, in order to support thread-local storage. Figure A-15 outlines the address computation in this case. The instruction's address, named *logical address* in the Intel documentation, is added to the base address in the segment register's descriptor, yielding the virtual address, also named *linear address*. The virtual address is then translated (§ A.5) to a physical address.

Outside the special case of using FS or GS to reference thread-local storage, the logical and virtual (linear) addresses match. Therefore, most of the time, we can get away with completely ignoring segmentation. In these cases, we use the term “virtual address” to refer to both the virtual and the linear address.

Even though CS is not used for segmentation, 64-bit system software needs to load a valid selector into it. The CPU uses the ring number in the CS selector to track the current privilege level, and uses one of the type bits to know whether it's running 64-bit code, or 32-bit code in compatibility mode.

The DS and ES segment registers are completely ignored, and can have null selectors

loaded in them. The CPU loads a null selector in SS when switching privilege levels, discussed in § A.8.2.

Modern kernels only use one descriptor table, the *Global Descriptor Table* (GDT), whose virtual address is stored in the GDTR register. Table A.2 shows a typical GDT layout that can be used by 64-bit kernels to run both 32-bit and 64-bit applications.

Descriptor	Selector
Null (must be unused)	0
Kernel code	0x08 (index 1, ring 0)
Kernel data	0x10 (index 2, ring 0)
User code	0x1B (index 3, ring 3)
User data	0x1F (index 4, ring 3)
TSS	0x20 (index 5, ring 0)

Table A.2: A typical GDT layout in the 64-bit Intel Architecture.

The last entry in Table A.2 is a descriptor for the *Task State Segment* (TSS), which was designed to implement hardware context switching, named *task switching* in the Intel documentation. The descriptor is stored in the *Task Register* (TR), which behaves like the other segment registers described above.

Task switching was removed from the 64-bit architecture, but the TR segment register was preserved, and it points to a repurposed TSS data structure. The 64-bit TSS contains an *I/O map*, which indicates what parts of the I/O address space can be accessed directly from ring 3, and the *Interrupt Stack Table* (IST), which is used for privilege level switching (§ A.8.2).

Modern operating systems do not allow application software any direct access to the I/O address space, so the kernel sets up a single TSS that is loaded into TR during early initialization, and used to represent all applications running under the OS.

A.8 Privilege Level Switching

Any architecture that has software privilege levels must provide a method for less privileged software to invoke the services of more privileged software. For example, application software needs the OS kernel's assistance to perform network or disk I/O, as that requires

access to privileged memory or to the I/O address space.

At the same time, less privileged software cannot be offered the ability to jump arbitrarily into more privileged code, as that would compromise the privileged software's ability to enforce security and isolation invariants. In our example, when an application wishes to write a file to the disk, the kernel must check if the application's user has access to that file. If the ring 3 code could perform an arbitrary jump in kernel space, it would be able to skip the access check.

For these reasons, the Intel architecture includes privilege-switching mechanisms used to transfer control from less privileged software to well-defined entry points in more privileged software. As suggested above, an architecture's privilege-switching mechanisms have deep implications for the security properties of its software. Furthermore, securely executing the software inside a protected container requires the same security considerations as privilege level switching.

Due to historical factors, the Intel architecture has a vast number of execution modes, and an intimidating amount of transitions between them. We focus on the privilege level switching mechanisms used by modern 64-bit software, summarized in Figure A-16.

A.8.1 System Calls

On modern processors, application software uses the `SYSCALL` instruction to invoke ring 0 code, and the kernel uses `SYSRET` to switch the privilege level back to ring 3. `SYSCALL` jumps into a predefined kernel location, which is specified by writing to a pair of architectural MSRs (§ A.4).

All MSRs can only be read or written by ring 0 code. This is a crucial security property, because it entails that application software cannot modify `SYSCALL`'s MSRs. If that was the case, a rogue application could abuse the `SYSCALL` instruction to execute arbitrary kernel code, potentially bypassing security checks.

The `SYSRET` instruction switches the current privilege level from ring 0 back to ring 3, and jumps to the address in `RCX`, which is set by the `SYSCALL` instruction. The `SYSCALL` / `SYSRET` pair does not perform any memory access, so it out-performs the Intel architecture's

previous privilege switching mechanisms, which saved state on a stack. The design can get away without referencing a stack because kernel calls are not recursive.

A.8.2 Faults

The processor also performs a switch from ring 3 to ring 0 when a *hardware exception* occurs while executing application code. Some exceptions indicate bugs in the application, whereas other exceptions require kernel action.

A *general protection fault* (#GP) occurs when software attempts to perform a disallowed action, such as setting the CR3 register from ring 3.

A *page fault* (#PF) occurs when address translation encounters a page table entry whose P flag is 0, or when the memory inside a page is accessed in way that is inconsistent with the access bits in the page table entry. For example, when ring 3 software accesses the memory inside a page whose S bit is set, the result of the memory access is #PF.

When a hardware exception occurs in application code, the CPU performs a ring switch, and calls the corresponding *exception handler*. For example, the #GP handler typically terminates the application's process, while the #PF handler reads the swapped out page back into RAM and resumes the application's execution.

The exception handlers are a part of the OS kernel, and their locations are specified in the first 32 entries of the Interrupt Descriptor Table (IDT), whose structure is shown in Table A.3. The IDT's physical address is stored in the IDTR register, which can only be accessed by ring 0 code. Kernels protect the IDT memory using page tables, so that ring 3 software cannot access it.

Field	Bits
Handler RIP	64
Handler CS	16
Interrupt Stack Table (IST) index	3

Table A.3: The essential fields of an IDT entry in 64-bit mode. Each entry points to a hardware exception or interrupt handler.

Each IDT entry has a 3-bit index pointing into the Interrupt Stack Table (IST), which is an array of 8 stack pointers stored in the TSS described in § A.7.

When a hardware exception occurs, the execution state may be corrupted, and the current stack cannot be relied on. Therefore, the CPU first uses the handler's IDT entry to set up a known good stack. SS is loaded with a null descriptor, and RSP is set to the IST value to which the IDT entry points. After switching to a reliable stack, the CPU pushes the snapshot in Table A.4 on the stack, then loads the IDT entry's values into the CS and RIP registers, which trigger the execution of the exception handler.

Field	Bits
Exception SS	64
Exception RSP	64
RFLAGS	64
Exception CS	64
Exception RIP	64
Exception code	64

Table A.4: The snapshot pushed on the handler's stack when a hardware exception occurs. IRET restores registers from this snapshot.

After the exception handler completes, it uses the IRET (interrupt return) instruction to load the registers from the on-stack snapshot and switch back to ring 3.

The Intel architecture gives the fault handler complete control over the execution context of the software that incurred the fault. This privilege is necessary for handlers (e.g., #GP) that must perform context switches (§ A.6) as a consequence of terminating a thread that encountered a bug. It follows that all fault handlers must be trusted to not leak or tamper with the information in an application's execution context.

A.8.3 VMX Privilege Level Switching

Intel systems that take advantage of the hardware virtualization support to run multiple operating systems at the same time use a hypervisor that manages the VMs. The hypervisor creates a *Virtual Machine Control Structure* (VMCS) for each operating system instance that it wishes to run, and uses the VMENTER instruction to assign a logical processor to the VM.

When a logical processor encounters a fault that must be handled by the hypervisor, the logical processor performs a VM exit. For example, if the address translation process encounters an EPT entry with the P flag set to 0, the CPU performs a VM exit, and the

hypervisor has an opportunity to bring the page into RAM.

The VMCS shows a great application of the encapsulation principle [134], which is generally used in high-level software, to computer architecture. The Intel architecture specifies that each VMCS resides in DRAM and is 4 KB in size. However, the architecture does not specify the VMCS format, and instead requires the hypervisor to interact with the VMCS via CPU instructions such as `VMREAD` and `VMWRITE`.

This approach allows Intel to add VMX features that require VMCS format changes, without the burden of having to maintain backwards compatibility. This is no small feat, given that huge amounts of complexity in the Intel architecture were introduced due to compatibility requirements.

A.9 A Computer Map

This section outlines the hardware components that make up a computer system based on the Intel architecture.

§ A.9.1 summarizes the structure of a *motherboard*. This is necessary background for reasoning about the cost and impact of physical attacks against a computing system. § A.9.2 describes Intel's Management Engine, which plays a role in the computer's bootstrap process, and has significant security implications.

§ A.9.3 presents the building blocks of an Intel processor, and § A.9.4 models an Intel execution core at a high level. This is the foundation for implementing defenses against physical attacks. Perhaps more importantly, reasoning about software attacks based on information leakage, such as timing attacks, requires understanding how a processor's computing resources are shared and partitioned between mutually distrusting parties.

The information in here is either contained in the SDM or in Intel's Optimization Reference Manual [99].

A.9.1 The Motherboard

A computer's components are connected by a printed circuit board called a *motherboard*, shown in Figure A-17, which consists of *sockets* connected by *buses*. Sockets connect

chip-carrying *packages* to the board. The Intel documentation uses the term “package” to specifically refer to a CPU.

The CPU (described in § A.9.3) hosts the execution cores that run the software stack shown in Figure A-5 and described in § A.3, namely the SMM code, the hypervisor, operating systems, and application processes. The computer’s main memory is provided by *Dynamic Random-Access Memory* (DRAM) chips.

The *Platform Controller Hub* (PCH) houses (relatively) low-speed I/O controllers driving the slower buses in the system, like SATA, used by storage devices, and USB, used by input peripherals. The PCH is also known as the *chipset*. At a first approximation, the *south bridge* term in older documentation can also be considered as a synonym for PCH.

Motherboards also have a non-volatile (flash) memory chip that hosts firmware which implements the *Unified Extensible Firmware Interface* (UEFI) specification [186]. The firmware contains the boot code and the code that executes in System Management Mode (SMM, § A.3).

The components we care about are connected by the following buses: the *Quick-Path Interconnect* (QPI [94]), a network of point-to-point links that connect processors, the *double data rate* (DDR) bus that connects a CPU to DRAM, the *Direct Media Interface* (DMI) bus that connects a CPU to the PCH, the *Peripheral Component Interconnect Express* (PCIe) bus that connects a CPU to peripherals such as a *Network Interface Card* (NIC), and the *Serial Programming Interface* (SPI) used by the PCH to communicate with the flash memory.

The PCIe bus is an extended, point-to-point version of the PCI standard, which provides a method for any peripheral connected to the bus to perform *Direct Memory Access* (DMA), transferring data to and from DRAM without involving an execution core and spending CPU cycles. The PCI standard includes a configuration mechanism that assigns a range of DRAM to each peripheral, but makes no provisions for restricting a peripheral’s DRAM accesses to its assigned range.

Network interfaces consist of a *physical* (PHY) module that converts the analog signals on the network media to and from digital bits, and a *Media Access Control* (MAC) module that implements a network-level protocol. Modern Intel-based motherboards forego a full-fledged NIC, and instead include an Ethernet [87] PHY module.

A.9.2 The Intel Management Engine (ME)

Intel's *Management Engine* (ME) is an embedded computer that was initially designed for remote system management and troubleshooting of server-class systems that are often hosted in data centers. However, all of Intel's recent PCHs contain an ME [83], and it currently plays a crucial role in platform bootstrapping, which is described in detail in § A.13. Most of the information in this section is obtained from an Intel-sponsored book [166].

The ME is part of Intel's *Active Management Technology* (AMT), which is marketed as a convenient way for IT administrators to troubleshoot and fix situations such as failing hardware, or a corrupted OS installation, without having to gain physical access to the impacted computer.

The Intel ME, shown in Figure A-18, remains functional during most hardware failures because it is an entire embedded computer featuring its own execution core, bootstrap ROM, and internal RAM. The ME can be used for troubleshooting effectively thanks to an array of abilities that include overriding the CPU's boot vector and a DMA engine that can access the computer's DRAM. The ME provides remote access to the computer without any CPU support because it can use the *System Management bus* (SMBus) to access the motherboard's Ethernet PHY or an AMT-compatible NIC [103].

The Intel ME is connected to the motherboard's power supply using a power rail that stays active even when the host computer is in the *Soft Off* mode [103], known as ACPI G2/S5, where most of the computer's components are powered off [90], including the CPU and DRAM. For all practical purposes, this means that the ME's execution core is active as long as the power supply is still connected to a power source.

In S5, the ME cannot access the DRAM, but it can still use its own internal memories. The ME can also still communicate with a remote party, as it can access the motherboard's Ethernet PHY via SMBus. This enables applications such as AMT's theft prevention, where a laptop equipped with a cellular modem can be tracked and permanently disabled as long as it has power and signal.

As the ME remains active in deep power-saving modes, its design must rely on low-power components. The execution core is an Argonaut RISC Core (ARC) clocked at

200-400MHz, which is typically used in low-power embedded designs. On a very recent PCH [103], the internal SRAM has 640KB, and is shared with the Integrated Sensor Hub (ISH)'s core. The SMBus runs at 1MHz and, without CPU support, the motherboard's Ethernet PHY runs at 10Mbps.

When the host computer is powered on, the ME's execution core starts running code from the ME's bootstrap ROM. The bootstrap code loads the ME's software stack from the same flash chip that stores the host computer's firmware. The ME accesses the flash memory chip an embedded SPI controller.

A.9.3 The Processor Die

An Intel processor's die, illustrated in Figure A-19, is divided into two broad areas: the *core area* implements the instruction execution pipeline typically associated with CPUs, while the *uncore* provides functions that were traditionally hosted on separate chips, but are currently integrated on the CPU die to reduce latency and power consumption.

At a conceptual level, the uncore of modern processors includes an *integrated memory controller* (iMC) that interfaces with the DDR bus, an *integrated I/O controller* (IIO) that implements PCIe bus lanes and interacts with the DMI bus, and a growing number of integrated peripherals, such as a *Graphics Processing Unit* (GPU). The uncore structure is described in some processor family datasheets [101, 100], and in the overview sections in Intel's uncore performance monitoring documentation [39, 97, 93].

Security extensions to the Intel architecture, such as Trusted Execution Technology (TXT) [74] and Software Guard Extensions (SGX) [143, 15], rely on the fact that the processor die includes the memory and I/O controller, and thus can prevent any device from accessing protected memory areas via *Direct Memory Access* (DMA) transfers. § A.11.3 takes a deeper look at the uncore organization and at the machinery used to prevent unauthorized DMA transfers.

A.9.4 The Core

Virtually all modern Intel processors have core areas consisting of multiple copies of the execution core circuitry, each of which is called a *core*. At the time of this writing, desktop-class Intel CPUs have 4 cores, and server-class CPUs have as many as 18 cores.

Most Intel CPUs feature *hyper-threading*, which means that a core (shown in Figure A-20) has two copies of the register files backing the execution context described in § A.6, and can execute two separate streams of instructions simultaneously. Hyper-threading reduces the impact of memory stalls on the utilization of the fetch, decode and execution units.

A hyper-threaded core is exposed to system software as two *logical processors* (LPs), also named *hardware threads* in the Intel documentation. The logical processor abstraction allows the code used to distribute work across processors in a multi-processor system to function without any change on multi-core hyper-threaded processors.

The high level of resource sharing introduced by hyper-threading introduces a security vulnerability. Software running on one logical processor can use the high-resolution performance counter (RDTSCP, § A.4) [156] to get information about the instructions and memory access patterns of another piece of software that is executed on the other logical processor on the same core.

That being said, the biggest downside of hyper-threading might be the fact that writing about Intel processors in a rigorous manner requires the use of the cumbersome term Logical Processor instead of the shorter and more intuitive “CPU core”, which can often be abbreviated to “core”.

A.10 Out-of-Order and Speculative Execution

CPU cores can execute instructions orders of magnitude faster than DRAM can read data. Computer architects attempt to bridge this gap by using hyper-threading (§ A.9.3), out-of-order and speculative execution, and caching, which is described in § A.11. In CPUs that use out-of-order execution, the order in which the CPU carries out a program’s instructions (*execution order*) is not necessarily the same as the order in which the instructions would be executed by a sequential evaluation system (*program order*).

An analysis of a system's information leakage must take out-of-order execution into consideration. Any CPU actions observed by an attacker match the execution order, so the attacker may learn some information by comparing the observed execution order with a known program order. At the same time, attacks that try to infer a victim's program order based on actions taken by the CPU must account for out-of-order execution as a source of noise.

This section summarizes the out-of-order and speculative execution concepts used when reasoning about a system's security properties. [154] and [79] cover the concepts in great depth, while Intel's optimization manual [99] provides details specific to Intel CPUs.

Figure A-21 provides a more detailed view of the CPU core components involved in out-of-order execution, and omits some less relevant details from Figure A-20.

The Intel architecture defines a *complex instruction set* (CISC). However, virtually all modern CPUs are architected following *reduced instruction set* (RISC) principles. This is accomplished by having the instruction decode stages break down each instruction into *micro-ops*, which resemble RISC instructions. The other stages of the execution pipeline work exclusively with micro-ops.

A.10.1 Out-of-Order Execution

Different types of instructions require different logic circuits, called *functional units*. For example, the arithmetic logic unit (ALU), which performs arithmetic operations, is completely different from the load and store unit, which performs memory operations. Different circuits can be used at the same time, so each CPU core can execute multiple micro-ops in parallel.

The core's out-of-order engine receives decoded micro-ops, identifies the micro-ops that can execute in parallel, assigns them to functional units, and combines the outputs of the units so that the results are equivalent to having the micro-ops executed sequentially in the order in which they come from the decode stages.

For example, consider the sequence of pseudo micro-ops⁴ in Table A.5 below. The OR uses the result of the LOAD, but the ADD does not. Therefore, a good scheduler can have the

⁴The set of micro-ops used by Intel CPUs is not publicly documented. The fictional examples in this section suffice for illustration purposes.

load store unit execute the `LOAD` and the ALU execute the `ADD`, all in the same clock cycle.

#	Micro-op	Meaning
1	<code>LOAD RAX, RSI</code>	$RAX \leftarrow \text{DRAM}[RSI]$
2	<code>OR RDI, RDI, RAX</code>	$RDI \leftarrow RDI \vee RAX$
3	<code>ADD RSI, RSI, RCX</code>	$RSI \leftarrow RSI + RCX$
4	<code>SUB RBX, RSI, RDX</code>	$RBX \leftarrow RSI - RDX$

Table A.5: Pseudo micro-ops for the out-of-order execution example.

The out-of-order engine in recent Intel CPUs works roughly as follows. Micro-ops received from the decode queue are written into a *reorder buffer* (ROB) while they are *in-flight* in the execution unit. The *register allocation table* (RAT) matches each register with the last reorder buffer entry that updates it. The *renamer* uses the RAT to rewrite the source and destination fields of micro-ops when they are written in the ROB, as illustrated in Tables A.6 and A.7. Note that the ROB representation makes it easy to determine the dependencies between micro-ops.

#	Op	Source 1	Source 2	Destination
1	LOAD	RSI	\emptyset	RAX
2	OR	RDI	ROB #1	RSI
3	ADD	RSI	RCX	RSI
4	SUB	ROB # 3	RDX	RBX

Table A.6: Data written by the renamer into the reorder buffer (ROB), for the micro-ops in Table A.5.

Register	RAX	RBX	RCX	RDX	RSI	RDI
ROB #	#1	#4	\emptyset	\emptyset	#3	#2

Table A.7: Relevant entries of the register allocation table after the micro-ops in Table A.5 are inserted into the ROB.

The scheduler decides which micro-ops in the ROB get executed, and places them in the *reservation station*. The reservation station has one port for each functional unit that can execute micro-ops independently. Each reservation station port holds one micro-op from the ROB. The reservation station port waits until the micro-op's dependencies are satisfied and forwards the micro-op to the functional unit. When the functional unit completes executing the micro-op, its result is *written back* to the ROB, and forwarded to any other reservation station port that depends on it.

The ROB stores the results of completed micro-ops until they are *retired*, meaning that the results are *committed* to the register file and the micro-ops are removed from the ROB. Although micro-ops can be executed out-of-order, they must be retired in program order, in order to handle exceptions correctly. When a micro-op causes a hardware exception (§ A.8.2), all the following micro-ops in the ROB are *squashed*, and their results are discarded.

In the example above, the `ADD` can complete before the `LOAD`, because it does not require a memory access. However, the `ADD`'s result cannot be committed before `LOAD` completes. Otherwise, if the `ADD` is committed and the `LOAD` causes a page fault, software will observe an incorrect value for the `RSI` register.

The ROB is tailored for discovering register dependencies between micro-ops. However, micro-ops that execute out-of-order can also have memory dependencies. For this reason, out-of-order engines have a *load buffer* and a *store buffer* that keep track of in-flight memory operations and are used to resolve memory dependencies.

A.10.2 Speculative Execution

Branch instructions, also called *branches*, change the instruction pointer (RIP, § A.6), if a condition is met (*the branch is taken*). They implement conditional statements (`if`) and looping statements, such as `while` and `for`. The most well-known branching instructions in the Intel architecture are in the `jcc` family, such as `je` (jump if equal).

Branches pose a challenge to the decode stage, because the instruction that should be fetched after a branch is not known until the branching condition is evaluated. In order to avoid stalling the decode stage, modern CPU designs include *branch predictors* that use historical information to guess whether a branch will be taken or not.

When the decode stage encounters a branch instruction, it asks the branch predictor for a guess as to whether the branch will be taken or not. The decode stage bundles the branch condition and the predictor's guess into a branch check micro-op, and then continues decoding on the path indicated by the predictor. The micro-ops following the branch check are marked as *speculative*.

When the branch check micro-op is executed, the branch unit checks whether the branch

predictor's guess was correct. If that is the case, the branch check is retired successfully. The scheduler handles *mispredictions* by squashing all the micro-ops following the branch check, and by signaling the instruction decoder to flush the micro-op decode queue and start fetching the instructions that follow the correct branch.

Modern CPUs also attempt to predict memory read patterns, so they can *prefetch* the memory locations that are about to be read into the cache. Prefetching minimizes the latency of successfully predicted read operations, as their data will already be cached. This is accomplished by exposing circuits called prefetchers to memory accesses and cache misses. Each prefetcher can recognize a particular access pattern, such as sequentially reading an array's elements. When memory accesses match the pattern that a prefetcher was built to recognize, the prefetcher loads the cache line corresponding to the next memory access in its pattern.

A.11 Cache Memories

At the time of this writing, CPU cores can process data $\approx 200\times$ faster than DRAM can supply it. This gap is bridged by an hierarchy of cache memories, which are orders of magnitude smaller and an order of magnitude faster than DRAM. While caching is transparent to application software, the system software is responsible for managing and coordinating the caches that store address translation (§ A.5) results.

Caches impact the security of a software system in two ways. First, the Intel architecture relies on system software to manage address translation caches, which becomes an issue in a threat model where the system software is untrusted. Second, caches in the Intel architecture are shared by all the software running on the computer. This opens up the way for *cache timing attacks*, an entire class of software attacks that rely on observing the time differences between accessing a cached memory location and an uncached memory location.

This section summarizes the caching concepts and implementation details needed to reason about both classes of security problems mentioned above. [175], [154] and [79] provide a good background on low-level cache implementation concepts. § B.8 describes cache timing attacks.

A.11.1 Caching Principles

At a high level, caches exploit the high locality in the memory access patterns of most applications to hide the main memory's (relatively) high latency. By *caching* (storing a copy of) the most recently accessed code and data, these relatively small memories can be used to satisfy 90%-99% of an application's memory accesses.

In an Intel processor, the *first-level* (L1) cache consists of a separate data cache (D-cache) and an instruction cache (I-cache). The instruction fetch and decode stage is directly connected to the L1 I-cache, and uses it to read the streams of instructions for the core's logical processors. Micro-ops that read from or write to memory are executed by the memory unit (MEM in Figure A-20), which is connected to the L1 D-cache and forwards memory accesses to it.

Figure A-22 illustrates the steps taken by a cache when it receives a memory access. First, a *cache lookup* uses the memory address to determine if the corresponding data exists in the cache. A *cache hit* occurs when the address is found, and the cache can resolve the memory access quickly. Conversely, if the address is not found, a *cache miss* occurs, and a *cache fill* is required to resolve the memory access. When doing a fill, the cache forwards the memory access to the next level of the memory hierarchy and caches the response. Under most circumstances, a cache fill also triggers a *cache eviction*, in which some data is removed from the cache to make room for the data coming from the fill. If the data that is evicted has been modified since it was loaded in the cache, it must be *written back* to the next level of the memory hierarchy.

Table A.8 shows the key characteristics of the memory hierarchy implemented by modern Intel CPUs. Each core has its own L1 and L2 cache (see Figure A-20), while the L3 cache is in the CPU's uncore (see Figure A-19), and is shared by all the cores in the package.

The numbers in Table A.8 suggest that cache placement can have a large impact on an application's execution time. Because of this, the Intel architecture includes an assortment of instructions that give performance-sensitive applications some control over the caching of their working sets. `PREFETCH` instructs the CPU's prefetcher to cache a specific memory address, in preparation for a future memory access. The memory writes performed by the

Memory	Size	Access Time
Core Registers	1 KB	no latency
L1 D-Cache	32 KB	4 cycles
L2 Cache	256 KB	10 cycles
L3 Cache	8 MB	40-75 cycles
DRAM	16 GB	60 ns

Table A.8: Approximate sizes and access times for each level in the memory hierarchy of an Intel processor, from [131]. Memory sizes and access times differ by orders of magnitude across the different levels of the hierarchy. This table does not cover multi-processor systems.

MOVNT instruction family bypass the cache if a fill would be required. CLFLUSH evicts any cache lines storing a specific address from the entire cache hierarchy.

The methods mentioned above are available to software running at all privilege levels, because they were designed for high-performance workloads with large working sets, which are usually executed at ring 3 (§ A.3). For comparison, the instructions used by system software to manage the address translation caches, described in § A.11.5 below, can only be executed at ring 0.

A.11.2 Cache Organization

In the Intel architecture, caches are completely implemented in hardware, meaning that the software stack has no direct control over the eviction process. However, software can gain some control over which data gets evicted by understanding how the caches are organized, and by cleverly placing its data in memory.

The *cache line* is the atomic unit of cache organization. A cache line has *data*, a copy of a continuous range of DRAM, and a *tag*, identifying the memory address that the data comes from. Fills and evictions operate on entire lines.

The cache line size is the size of the data, and is always a power of two. Assuming n -bit memory addresses and a cache line size of 2^l bytes, the lowest l bits of a memory address are an offset into a cache line, and the highest $n - l$ bits determine the cache line that is used to store the data at the memory location. All recent processors have 64-byte cache lines.

The L1 and L2 caches in recent processors are multi-way set-associative with direct set

indexing, as shown in Figure A-23. A W -way set-associative cache has its memory divided into *sets*, where each set has W lines. A memory location can be cached in any of the w lines in a specific set that is determined by the highest $n - l$ bits of the location's memory address. Direct set indexing means that the S sets in a cache are numbered from 0 to $S - 1$, and the memory location at address A is cached in the set numbered $A_{n-1\dots n-l} \bmod S$.

In the common case where the number of sets in a cache is a power of two, so $S = 2^s$, the lowest l bits in an address make up the cache line offset, the next s bits are the set index. The highest $n - s - l$ bits in an address are not used when selecting where a memory location will be cached. Figure A-23 shows the cache structure and lookup process.

A.11.3 Cache Coherence

The Intel architecture was designed to support application software that was not written with caches in mind. One aspect of this support is the *Total Store Order* (TSO) [151] memory model, which promises that all the logical processors in a computer see the same order of DRAM writes.

The same memory location might be simultaneously cached by different cores' caches, or even by caches on separate chips, so providing the TSO guarantees requires a *cache coherence protocol* that synchronizes all the cache lines in a computer that reference the same memory address.

The cache coherence mechanism is not visible to software, so it is only briefly mentioned in the SDM. Fortunately, Intel's optimization reference [99] and the datasheets referenced in § A.9.3 provide more information. Intel processors use variations of the MESIF [70] protocol, which is implemented in the CPU and in the protocol layer of the QPI bus.

The SDM and the `CPUID` instruction output indicate that the L3 cache, also known as the *last-level cache* (LLC) is *inclusive*, meaning that any location cached by an L1 or L2 cache must also be cached in the LLC. This design decision reduces complexity in many implementation aspects. We estimate that the bulk of the cache coherence implementation is in the CPU's uncore, thanks to the fact that cache synchronization can be achieved without having to communicate to the lower cache levels that are inside execution cores.

The QPI protocol defines *cache agents*, which are connected to the last-level cache in a processor, and *home agents*, which are connected to memory controllers. Cache agents make requests to home agents for cache line data on cache misses, while home agents keep track of cache line ownership, and obtain the cache line data from other cache line agents, or from the memory controller. The QPI routing layer supports multiple agents per socket, and each processor has its own caching agents, and at least one home agent.

Figure A-24 shows that the CPU uncore has a bidirectional ring interconnect, which is used for communication between execution cores and the other uncore components. The execution cores are connected to the ring by *CBoxes*, which route their LLC accesses. The routing is static, as the LLC is divided into same-size slices (common slice sizes are 1.5 MB and 2.5 MB), and an undocumented hashing scheme maps each possible physical address to exactly one LLC slice.

Intel's documentation states that the hashing scheme mapping physical addresses to LLC slices was designed to avoid having a slice become a hotspot, but stops short of providing any technical details. Fortunately, independent researchers have reversed-engineered the hash functions for recent processors [88, 139, 206].

The hashing scheme described above is the reason why the L3 cache is documented as having a “complex” indexing scheme, as opposed to the direct indexing used in the L1 and L2 caches.

The number of LLC slices matches the number of cores in the CPU, and each LLC slice shares a CBox with a core. The CBoxes implement the cache coherence engine, so each CBox acts as the QPI cache agent for its LLC slice. CBoxes use a *Source Address Decoder* (SAD) to route DRAM requests to the appropriate home agents. Conceptually, the SAD takes in a memory address and access type, and outputs a transaction type (coherent, non-coherent, IO) and a node ID. Each CBox contains a SAD replica, and the configurations of all SADs in a package are identical.

The SAD configurations are kept in sync by the *UBox*, which is the uncore configuration controller, and connects the *System agent* to the ring. The UBox is responsible for reading and writing physically distributed registers across the uncore. The UBox also receives interrupts from system and dispatches them to the appropriate core.

On recent Intel processors, the uncore also contains at least one memory controller. Each integrated memory controller (iMC or MBox in Intel’s documentation) is connected to the ring by a *home agent* (HA or *BBox* in Intel’s datasheets). Each home agent contains a *Target Address Decoder* (TAD), which maps each DRAM address to an address suitable for use by the DRAM chips, namely a DRAM channel, bank, rank, and a DIMM address. The mapping in the TAD is not documented by Intel, but it has been reverse-engineered [155].

The integration of the memory controller on the CPU brings the ability to filter DMA transfers. Accesses from a peripheral connected to the PCIe bus are handled by the integrated I/O controller (IIO), placed on the ring interconnect via the UBox, and then reach the iMC. Therefore, on modern systems, DMA transfers go through both the SAD and TAD, which can be configured to abort DMA transfers targeting protected DRAM ranges.

A.11.4 Caching and Memory-Mapped Devices

Caches rely on the assumption that the underlying memory implements the memory abstraction in § A.2. However, the physical addresses that map to memory-mapped I/O devices usually deviate from the memory abstraction. For example, some devices expose command registers that trigger certain operations when written, and always return a zero value. Caching addresses that map to such memory-mapped I/O devices will lead to incorrect behavior.

Furthermore, even when the memory-mapped devices follow the memory abstraction, caching their memory is sometimes undesirable. For example, caching a graphic unit’s framebuffer could lead to visual artifacts on the user’s display, because of the delay between the time when a write is issued and the time when the corresponding cache lines are evicted and written back to memory.

In order to work around these problems, the Intel architecture implements a few caching behaviors, described below, and provides a method for partitioning the memory address space (§ A.4) into regions, and for assigning a desired caching behavior to each region.

Uncacheable (UC) memory has the same semantics as the I/O address space (§ A.4). UC memory is useful when a device’s behavior is dependent on the order of memory reads

and writes, such as in the case of memory-mapped command and data registers for a PCIe NIC (§ A.9.1). The out-of-order execution engine (§ A.10) does not reorder UC memory accesses, and does not issue speculative reads to UC memory.

Write Combining (WC) memory addresses the specific needs of framebuffers. WC memory is similar to UC memory, but the out-of-order engine may reorder memory accesses, and may perform speculative reads. The processor stores writes to WC memory in a write combining buffer, and attempts to group multiple writes into a (more efficient) line write bus transaction.

Write Through (WT) memory is cached, but write misses do not cause cache fills. This is useful for preventing large memory-mapped device memories that are rarely read, such as framebuffers, from taking up cache memory. WT memory is covered by the cache coherence engine, may receive speculative reads, and is subject to operation reordering.

DRAM is represented as *Write Back* (WB) memory, which is optimized under the assumption that all the devices that need to observe the memory operations implement the cache coherence protocol. WB memory is cached as described in § A.11, receives speculative reads, and operations targeting it are subject to reordering.

Write Protected (WP) memory is similar to WB memory, with the exception that every write is propagated to the system bus. It is intended for memory-mapped buffers, where the order of operations does not matter, but the devices that need to observe the writes do not implement the cache coherence protocol, in order to reduce hardware costs.

On recent Intel processors, the cache's behavior is mainly configured by the *Memory Type Range Registers* (MTRRs) and by *Page Attribute Table* (PAT) indices in the page tables (§ A.5). The behavior is also impacted by the Cache Disable (CD) and Not-Write through (NW) bits in Control Register 0 (CR0, § A.4), as well as by equivalent bits in page table entries, namely Page-level Cache Disable (PCD) and Page-level Write-Through (PWT).

The MTRRs were intended to be configured by the computer's firmware during the boot sequence. Fixed MTRRs cover pre-determined ranges of memory, such as the memory areas that had special semantics in the computers using 16-bit Intel processors. The ranges covered by *variable MTRRs* can be configured by system software. The representation used to specify the ranges is described below, as it has some interesting properties that have

proven useful in other systems.

Each variable memory type range is specified using a *range base* and a *range mask*. A memory address belongs to the range if computing a bitwise AND between the address and the range mask results in the range base. This verification has a low-cost hardware implementation, shown in Figure A-25.

Each variable memory type range must have a size that is an integral power of two, and a starting address that is a multiple of its size, so it can be described using the base / mask representation described above. A range's starting address is its base, and the range's size is one plus its mask.

Another advantage of this range representation is that the base and the mask can be easily validated, as shown in Listing A.1. The range is aligned with respect to its size if and only if the bitwise AND between the base and the mask is zero. The range's size is a power of two if and only if the bitwise AND between the mask and one plus the mask is zero. According to the SDM, the MTRRs are not validated, but setting them to invalid values results in undefined behavior.

```
constexpr bool is_valid_range(  
    size_t base, size_t mask) {  
    // Base is aligned to size.  
    return (base & mask) == 0 &&  
        // Size is a power of two.  
        (mask & (mask + 1)) == 0;  
}
```

Listing A.1: The checks that validate the base and mask of a memory-type range can be implemented very easily.

No memory type range can partially cover a 4 KB page, which implies that the range base must be a multiple of 4 KB, and the bottom 12 bits of range mask must be set. This simplifies the interactions between memory type ranges and address translation, described in § A.11.5.

The PAT is intended to allow the operating system or hypervisor to tweak the caching behaviors specified in the MTRRs by the computer's firmware. The PAT has 8 entries that

specify caching behaviors, and is stored in its entirety in a MSR. Each page table entry contains a 3-bit index that points to a PAT entry, so the system software that controls the page tables can specify caching behavior at a very fine granularity.

A.11.5 Caches and Address Translation

Modern system software relies on address translation (§ A.5). This means that all the memory accesses issued by a CPU core use virtual addresses, which must undergo translation. Caches must know the physical address for a memory access, to handle aliasing (multiple virtual addresses pointing to the same physical address) correctly. However, address translation requires up to 20 memory accesses (see Figure A-12), so it is impractical to perform a full address translation for every cache access. Instead, address translation results are cached in the *translation look-aside buffer* (TLB).

Table A.9 shows the levels of the TLB hierarchy. Recent processors have separate L1 TLBs for instructions and data, and a shared L2 TLB. Each core has its own TLBs (see Figure A-20). When a virtual address is not contained in a core’s TLB, the *Page Miss Handler* (PMH) performs a *page walk* (page table / EPT traversal) to translate the virtual address, and the result is stored in the TLB.

Memory	Entries	Access Time
L1 I-TLB	$128 + 8 = 136$	1 cycle
L1 D-TLB	$64 + 32 + 4 = 100$	1 cycle
L2 TLB	$1536 + 8 = 1544$	7 cycles
Page Tables	$2^{36} \approx 6 \cdot 10^{10}$	18 cycles - 200ms

Table A.9: Approximate sizes and access times for each level in the TLB hierarchy, from [4].

In the Intel architecture, the PMH is implemented in hardware, so the TLB is never directly exposed to software and its implementation details are not documented. The SDM does state that each TLB entry contains the physical address associated with a virtual address, and the metadata needed to resolve a memory access. For example, the processor needs to check the writable (W) flag on every write, and issue a General Protection fault (#GP) if the write targets a read-only page. Therefore, the TLB entry for each virtual address caches the

logical-and of all the relevant W flags in the page table structures leading up to the page.

The TLB is transparent to application software. However, kernels and hypervisors must make sure that the TLBs do not get out of sync with the page tables and EPTs. When changing a page table or EPT, the system software must use the INVLPG instruction to invalidate any TLB entries for the virtual address whose translation changed. Some instructions *flush the TLBs*, meaning that they invalidate all the TLB entries, as a side-effect.

TLB entries also cache the desired caching behavior (§ A.11.4) for their pages. This requires system software to flush the corresponding TLB entries when changing MTRRs or page table entries. In return, the processor only needs to compute the desired caching behavior during a TLB miss, as opposed to computing the caching behavior on every memory access.

The TLB is not covered by the cache coherence mechanism described in § A.11.3. Therefore, when modifying a page table or EPT on a multi-core / multi-processor system, the system software is responsible for performing a *TLB shutdown*, which consists of stopping all the logical processors that use the page table / EPT about to be changed, performing the changes, executing TLB-invalidating instructions on the stopped logical processors, and then resuming execution on the stopped logical processors.

Address translation constrains the L1 cache design. On Intel processors, the set index in an L1 cache only uses the address bits that are not impacted by address translation, so that the L1 set lookup can be done in parallel with the TLB lookup. This is critical for achieving a low latency when both the L1 TLB and the L1 cache are hit.

Given a page size $P = 2^p$ bytes, the requirement above translates to $l + s \leq p$. In the Intel architecture, $p = 12$, and all recent processors have 64-byte cache lines ($l = 6$) and 64 sets ($s = 6$) in the L1 caches, as shown in Figure A-26. The L2 and L3 caches are only accessed if the L1 misses, so the physical address for the memory access is known at that time, and can be used for indexing.

A.12 Interrupts

Peripherals use *interrupts* to signal the occurrence of an event that must be handled by system software. For example, a keyboard triggers interrupts when a key is pressed or depressed. System software also relies on interrupts to implement preemptive multi-threading.

Interrupts are a kind of hardware exception (§ A.8.2). Receiving an interrupt causes an execution core to perform a privilege level switch and to start executing the system software's interrupt handling code. Therefore, the security concerns in § A.8.2 also apply to interrupts, with the added twist that interrupts occur independently of the instructions executed by the interrupted code, whereas most faults are triggered by the actions of the application software that incurs them.

Given the importance of interrupts when assessing a system's security, this section outlines the interrupt triggering and handling processes described in the SDM.

Peripherals use bus-specific protocols to signal interrupts. For example, PCIe relies on *Message Signaled Interrupts* (MSI), which are memory writes issued to specially designed memory addresses. The bus-specific interrupt signals are received by the *I/O Advanced Programmable Interrupt Controller* (IOAPIC) in the PCH, shown in Figure A-17.

The IOAPIC routes interrupt signals to one or more *Local Advanced Programmable Interrupt Controllers* (LAPICs). As shown in Figure A-19, each logical CPU has a LAPIC that can receive interrupt signals from the IOAPIC. The IOAPIC routing process assigns each interrupt to an 8-bit *interrupt vector* that is used to identify the interrupt sources, and to a 32-bit *APIC ID* that is used to identify the LAPIC that receives the interrupt.

Each LAPIC uses a 256-bit *Interrupt Request Register* (IRR) to track the unserved interrupts that it has received, based on the interrupt vector number. When the corresponding logical processor is available, the LAPIC copies the highest-priority unserved interrupt vector to the *In-Service Register* (ISR), and invokes the logical processor's interrupt handling process.

At the execution core level, interrupt handling reuses many of the mechanisms of fault handling (§ A.8.2). The interrupt vector number in the LAPIC's ISR is used to locate an interrupt handler in the IDT, and the handler is invoked, possibly after a privilege switch

is performed. The interrupt handler does the processing that the device requires, and then writes the LAPIC's *End Of Interrupt* (EOI) register to signal the fact that it has completed handling the interrupt.

Interrupts are treated like faults, so interrupt handlers have full control over the execution environment of the application being interrupted. This is used to implement pre-emptive multi-threading, which relies on a clock device that generates interrupts periodically, and on an interrupt handler that performs context switches.

System software can cause an interrupt on any logical processor by writing the target processor's APIC ID into the *Interrupt Command Register* (ICR) of the LAPIC associated with the logical processor that the software is running on. These interrupts, called *Inter-Processor Interrupts* (IPI), are needed to implement TLB shoot-downs (§ A.11.5).

A.13 Platform Initialization (Booting)

When a computer is powered up, it undergoes a *bootstrapping* process, also called *booting*, for simplicity. The boot process is a sequence of steps that collectively initialize all the computer's hardware components and load the system software into DRAM. An analysis of a system's security properties must be aware of all the pieces of software executed during the boot process, and must account for the trust relationships that are created when a software module loads another module.

This section outlines the details of the boot process needed to reason about the security of a system based on the Intel architecture. [95] provides a good reference for many of the booting process's low-level details. While some specifics of the boot process depend on the motherboard and components in a computer, this section focuses on the high-level flow described by Intel's documentation.

A.13.1 The UEFI Standard

The firmware in recent computers with Intel processors implements the *Platform Initialization* (PI) process in the *Unified Extensible Firmware Interface* (UEFI) specification [186]. The platform initialization follows the steps shown in Figure A-27 and described below.

The computer powers up, reboots, or resumes from sleep in the *Security phase* (SEC). The SEC implementation is responsible for establishing a temporary memory store and loading the next stage of the firmware into it. As the first piece of software that executes on the computer, the SEC implementation is the system's root of trust, and performs the first steps towards establishing the system's desired security properties.

For example, in a measured boot system (also known as trusted boot), all the software involved in the boot process is measured (cryptographically hashed, and the measurement is made available to third parties, as described in § B.3). In such a system, the SEC implementation takes the first steps in establishing the system's measurement, namely resetting the special register that stores the measurement result, measuring the PEI implementation, and storing the measurement in the special register.

SEC is followed by the *Pre-EFI Initialization phase* (PEI), which initializes the computer's DRAM, copies itself from the temporary memory store into DRAM, and tears down the temporary storage. When the computer is powering up or rebooting, the PEI implementation is also responsible for initializing all the non-volatile storage units that contain UEFI firmware and loading the next stage of the firmware into DRAM.

PEI hands off control to the *Driver eXecution Environment phase* (DXE). In DXE, a loader locates and starts firmware drivers for the various components in the computer. DXE is followed by a Boot Device Selection (BDS) phase, which is followed by a Transient System Load (TSL) phase, where an EFI application loads the operating system selected in the BDS phase. Last, the OS loader passes control to the operating system's kernel, entering the Run Time (RT) phase.

When waking up from sleep, the PEI implementation first initializes the non-volatile storage containing the system snapshot saved while entering the sleep state. The rest of the PEI implementation may use optimized re-initialization processes, based on the snapshot contents. The DXE implementation also uses the snapshot to restore the computer's state, such as the DRAM contents, and then directly executes the operating system's wake-up handler.

A.13.2 SEC on Intel Platforms

Right after a computer is powered up, circuitry in the power supply and on the motherboard starts establishing reference voltages on the power rails in a specific order, documented as “power sequencing” [190] in chipset specifications such as [105]. The rail powering up the Intel ME (§ A.9.2) in the PCH is powered up significantly before the rail that powers the CPU cores.

When the ME is powered up, it starts executing the code in its boot ROM, which sets up the SPI bus connected to the flash memory chip (§ A.9.1) that stores both the UEFI firmware and the ME’s firmware. The ME then loads its firmware from flash memory, which contains the ME’s operating system and applications.

After the Intel ME loads its software, it sets up some of the motherboard’s hardware, such as the PCH bus clocks, and then it kicks off the CPU’s bootstrap sequence. Most of the details of the ME’s involvement in the computer’s boot process are not publicly available, but initializing the clocks is mentioned in a few public documents [110, 5, 45, 7], and is made clear in firmware bringup guides, such as the leaked confidential guide [96] documenting firmware bringup for Intel’s Series 7 chipset.

The beginning of the CPU’s bootstrap sequence is the SEC phase, which is implemented in the processor circuitry. All the logical processors (LPs) on the motherboard undergo *hardware initialization*, which invalidates the caches (§ A.11) and TLBs (§ A.11.5), performs a *Built-In Self Test* (BIST), and sets all the registers (§ A.6) to pre-specified values.

After hardware initialization, the LPs perform the Multi-Processor (MP) initialization algorithm, which results in one LP being selected as the *bootstrap processor* (BSP), and all the other LPs being classified as *application processors* (APs).

According to the SDM, the details of the MP initialization algorithm for recent CPUs depend on the motherboard and firmware. In principle, after completing hardware initialization, all LPs attempt to issue a special no-op transaction on the QPI bus. A single LP will succeed in issuing the no-op, thanks to the QPI arbitration mechanism, and to the UBox (§ A.11.3) in each CPU package, which also serves as a ring arbiter. The arbitration priority of each LP is based on its APIC ID (§ A.12), which is provided by the motherboard when

the system powers up. The LP that issues the no-op becomes the BSP. Upon failing to issue the no-op, the other LPs become APs, and enter the *wait-for-SIPI* state.

Understanding the PEI firmware loading process is unnecessarily complicated by the fact that the SDM describes a legacy process consisting of having the BSP set its RIP register to 0xFFFFFFFF0 (16 bytes below 4 GB), where the firmware is expected to place a instruction that jumps into the PEI implementation.

Recent processors do not support the legacy approach at all [160]. Instead, the BSP reads a word from address 0xFFFFFEE8 (24 bytes below 4 GB) [213, 42], and expects to find the address of a *Firmware Interface Table* (FIT) in the memory address space (§ A.4), as shown in Figure A-28. The BSP is able to read firmware contents from non-volatile memory before the computer is initialized, because the initial SAD (§ A.11.3) and PCH (§ A.9.1) configurations maps a region in the memory address space to the SPI flash chip (§ A.9.1) that stores the computer's firmware.

The FIT [157] was introduced in the context of Intel's Itanium architecture, and its use in Intel's current 64-bit architecture is described in an Intel patent [42] and briefly documented in an obscure piece of TXT-related documentation [92]. The FIT contains *Authenticated Code Modules* (ACMs) that make up the firmware, and other platform-specific information, such as the TPM and TXT configuration [92].

The PEI implementation is stored in an ACM listed in the FIT. The processor loads the PEI ACM, verifies the trustworthiness of the ACM's public key, and ensures that the ACM's contents matches its signature. If the PEI passes the security checks, it is executed. Processors that support Intel TXT only accept Intel-signed ACMs [59, p. 92].

A.13.3 PEI on Intel Platforms

[95] and [37] describe the initialization steps performed by Intel platforms during the PEI phase, from the perspective of a firmware programmer. A few steps provide useful context for reasoning about threat models involving the boot process.

When the BSP starts executing PEI firmware, DRAM is not yet initialized. Therefore the PEI code starts executing in a *Cache-as-RAM* (CAR) mode, which only relies on the

BSP's internal caches, at the expense of imposing severe constraints on the size of the PEI's working set.

One of the first tasks performed by the PEI implementation is enabling DRAM, which requires discovering and initializing the DRAM chips connected to the motherboard, and then configuring the BSP's memory controllers (§ A.11.3) and MTRRs (§ A.11.4). Most firmware implementations use Intel's *Memory Reference Code* (MRC) for this task.

After DRAM becomes available, the PEI code is copied into DRAM and the BSP is taken out of CAR mode. The BSP's LAPIC (§ A.12) is initialized and used to send a broadcast *Startup Inter-Processor Interrupt* (SIPI, § A.12) to wake up the APs. The interrupt vector in a SIPI indicates the memory address of the AP initialization code in the PEI implementation.

The PEI code responsible for initializing APs is executed when the APs receive the SIPI wake-up. The AP PEI code sets up the AP's configuration registers, such as the MTRRs, to match the BSP's configuration. Next, each AP registers itself in a system-wide table, using a memory synchronization primitive, such as a semaphore, to avoid having two APs access the table at the same time. After the AP initialization completes, each AP is suspended again, and waits to receive an INIT Inter-Processor Interrupt from the OS kernel.

The BSP initialization code waits for all APs to register themselves into the system-wide table, and then proceeds to locate, load and execute the firmware module that implements DXE.

A.14 CPU Microcode

The Intel architecture features a large instruction set. Some instructions are used infrequently, and some instructions are very complex, which makes it impractical for an execution core to handle all the instructions in hardware. Intel CPUs use a *microcode* table to break down rare and complex instructions into sequences of simpler instructions. Architectural extensions that only require microcode changes are significantly cheaper to implement and validate than extensions that require changes in the CPU's circuitry.

It follows that a good understanding of what can be done in microcode is crucial to evaluating the cost of security features that rely on architecture extensions. Furthermore, the

limitations of microcode are sometimes the reasoning behind seemingly arbitrary architecture design decisions.

The first sub-section below presents the relevant facts pertaining to microcode in Intel's optimization reference [99] and SDM. The following subsections summarize information gleaned from Intel's patents and other researchers' findings.

A.14.1 The Role of Microcode

The frequently used instructions in the Intel architecture are handled by the core's fast path, which consists of simple decoders (§ A.10) that can emit at most 4 micro-ops per instruction. Infrequently used instructions and instructions that require more than 4 micro-ops use a slower decoding path that relies on a sequencer to read micro-ops from a *microcode store ROM* (MSROM).

The 4 micro-ops limitation can be used to guess intelligently whether an architectural feature is implemented in microcode. For example, it is safe to assume that XSAVE (§ A.6), which takes over 200 micro-ops on recent CPUs [57], is most likely performed in microcode, whereas simple arithmetic and memory accesses are handled directly by hardware.

The core's execution units handle common cases in fast paths implemented in hardware. When an input cannot be handled by the fast paths, the execution unit issues a *microcode assist*, which points the microcode sequencer to a routine in microcode that handles the edge cases. The most common cited example in Intel's documentation is floating point instructions, which issue assists to handle denormalized inputs.

The REP MOVSB family of instructions, also known as *string instructions* because of their use in strcpy-like functions, operate on variable-sized arrays. These instructions can handle small arrays in hardware, and issue microcode assists for larger arrays.

Modern Intel processors implement a microcode update facility. The SDM describes the process of applying microcode updates from the perspective of system software. Each core can be updated independently, and the updates must be reapplied on each boot cycle. A core can be updated multiple times. The latest SDM at the time of this writing states that a microcode update is up to 16 KB in size.

Processor engineers prefer to build new architectural features as microcode extensions, because microcode can be iterated on much faster than hardware, which reduces development cost [202, 203]. The update facility further increases the appeal of microcode, as some classes of bugs can be fixed after a CPU has been released.

Intel patents [142, 112] describing Software Guard Extensions (SGX) disclose that SGX is entirely implemented in microcode, except for the memory encryption engine. A description of SGX's implementation could provide great insights into Intel's microcode, but, unfortunately, the SDM chapters covering SGX do not include such a description. We therefore rely on other public information sources about the role of microcode in the security-sensitive areas covered by previous sections, namely memory management (§ A.5, § A.11.5), the handling of hardware exceptions (§ A.8.2) and interrupts (§ A.12), and platform initialization (§ A.13).

The use of microcode assists can be measured using the *Precise Event Based Sampling* (PEBS) feature in recent Intel processors. PEBS provides counters for the number of micro-ops coming from MSR0M, including complex instructions and assists, counters for the numbers of assists associated with some micro-op classes (SSE and AVX stores and transitions), and a counter for assists generated by all other micro-ops.

The PEBS feature itself is implemented using microcode assists (this is implied in the SDM and confirmed by [123]) when it needs to write the execution context into a PEBS record. Given the wide range of features monitored by PEBS counters, we assume that all execution units in the core can issue microcode assists, which are performed at micro-op retirement. This finding is confirmed by an Intel patent [26], and is supported by the existence of a PEBS counter for the “number of microcode assists invoked by hardware upon micro-op writeback.”

Intel's optimization manual describes one more interesting assist, from a memory system perspective. SIMD masked loads (using `VMSKMOV`) read a series of data elements from memory into a vector register. A mask register decides whether elements are moved or ignored. If the memory address overlaps an invalid page (e.g., the P flag is 0, § A.5), a microcode assist is issued, even if the mask indicates that no element from the invalid page should be read. The microcode checks whether the elements in the invalid page have the

corresponding mask bits set, and either performs the load or issues a page fault.

The description of machine checks in the SDM mentions page assists and page faults in the same context. We assume that the page assists are issued in some cases when a TLB miss occurs (§ A.11.5) and the PMH has to walk the page table. The following section develops this assumption and provides supporting evidence from Intel's assigned patents and published patent applications.

A.14.2 Microcode Structure

According to a 2013 Intel patent [86], the avenues considered for implementing new architectural features are a completely microcode-based implementation, using existing micro-ops, a microcode implementation with hardware support, which would use new micro-ops, and a complete hardware implementation, using finite state machines (FSMs).

The main component of the MSROM is a table of micro-ops [202, 203]. According to an example in a 2012 Intel patent [203], the table contains on the order of 20,000 micro-ops, and a micro-op has about 70 bits. On embedded processors, like the Atom, microcode may be partially compressed [202, 203].

The MSROM also contains an event ROM, which is an array of pointers to event handling code in the micro-ops table [164]. Microcode events are hardware exceptions, assists, and interrupts [26, 153, 38]. The processor described in a 1999 patent [164] has a 64-entry event table, where the first 16 entries point to hardware exception handlers and the other entries are used by assists.

The execution units can issue an assist or signal a fault by associating an event code with the result of a micro-op. When the micro-op is committed (§ A.10), the event code causes the out-of-order scheduler to squash all the micro-ops that are in-flight in the ROB. The event code is forwarded to the microcode sequencer, which reads the micro-ops in the corresponding event handler [26, 153].

The hardware exception handling logic (§ A.8.2) and interrupt handling logic (§ A.12) is implemented entirely in microcode [153]. Therefore, changes to this logic are relatively inexpensive to implement on Intel processors. This is rather fortunate, as the Intel architec-

ture's standard hardware exception handling process requires that the fault handler is trusted by the code that encounters the exception (§ A.8.2), and this assumption cannot be satisfied by a design where the software executing inside a secure container must be isolated from the system software managing the computer's resources.

The execution units in modern Intel processors support microcode procedures, via dedicated microcode call and return micro-ops [38]. The micro-ops manage a hardware data structure that conceptually stores a stack of microcode instruction pointers, and is integrated with out-of-order execution and hardware exceptions, interrupts and assists.

Asides from special micro-ops, microcode also employs special load and store instructions, which turn into special bus cycles, to issue commands to other functional units [163]. The memory addresses in the special loads and stores encode commands and input parameters. For example, stores to a certain range of addresses flush specific TLB sets.

A.14.3 Microcode and Address Translation

Address translation (§ A.5) is configured by CR3, which stores the physical address of the top-level page table, and by various bits in CR0 and CR4, all of which are described in the SDM. Writes to these control registers are implemented in microcode, which stores extra information in microcode-visible registers [66].

When a TLB miss (§ A.11.5) occurs, the memory execution unit forwards the virtual address to the *Page Miss Handler* (PMH), which performs the page walk needed to obtain a physical address. In order to minimize the latency of a page walk, the PMH is implemented as a *Finite-State Machine* (FSM) [81, 158]. Furthermore, the PMH fetches the page table entries from memory by issuing “stuffed loads”, which are special micro-ops that bypass the reorder buffer (ROB) and go straight to the memory execution units (§ A.10), thus avoiding the overhead associated with out-of-order scheduling [67, 163, 81].

The FSM in the PMH handles the fast path of the entire address translation process, which assumes no address translation fault (§ A.8.2) occurs [68, 67, 153, 164], and no page table entry needs to be modified [67].

When the PMH FSM detects the conditions that trigger a Page Fault or a General

Protection Fault, it communicates a microcode event code, corresponding to the detected fault condition, to the execution unit (§ A.10) responsible for memory operations [68, 67, 153, 164]. In turn, the execution unit triggers the fault by associating the event code with the micro-op that caused the address translation, as described in the previous section.

The PMH FSM does not set the Accessed or Dirty attributes (§ A.5.3) in page table entries. When it detects that a page table entry must be modified, the FSM issues a microcode event code for a page walk assist [67]. The microcode handler performs the page walk again, setting the A and D attributes on page table entries when necessary [67]. This finding was indirectly confirmed by the description for a PEBS event in the most recent SDM release.

The patents at the core of our descriptions above [68, 26, 67, 153, 164] were all issued between 1996 and 1999, which raises the concern of obsolescence. As Intel would not be able to file new patents for the same specifications, we cannot present newer patents with the information above. Fortunately, we were able to find newer patents that mention the techniques described above, proving their relevance to newer CPU models.

Two 2014 patents [81, 158] mention that the PMH is executing a FSM which issues stuffing loads to obtain page table entries. A 2009 patent [66] mentions that microcode is invoked after a PMH walk, and that the microcode can prevent the translation result produced by the PMH from being written to the TLB.

A 2013 patent [86] and a 2014 patent [159] on scatter / gather instructions disclose that the newly introduced instructions use a combination of hardware in the execution units that perform memory operations, which include the PMH. The hardware issues microcode assists for slow paths, such as gathering vector elements stored in uncacheable memory (§ A.11.4), and operations that cause Page Faults.

A 2014 patent on APIC (§ A.12) virtualization [173] describes a memory execution unit modification that invokes a microcode assist for certain memory accesses, based on the contents of some range registers. The patent also mentions that the range registers are checked when the TLB miss occurs and the PMH is invoked, in order to decide whether a fast hardware path can be used for APIC virtualization, or a microcode assist must be issued.

The recent patents mentioned above allow us to conclude that the PMH in recent processors still relies on an FSM and stuffed loads, and still uses microcode assists to handle

infrequent and complex operations. This assumption plays a key role in estimating the implementation complexity of architectural modifications targeting the processor's address translation mechanism.

A.14.4 Microcode and Booting

The SDM states that microcode performs the Built-In Self Test (BIST, § A.13.2), but does not provide any details on the rest of the CPU's hardware initialization.

In fact, the entire SEC implementation on Intel platforms is contained in the processor microcode [43, 42, 173]. This implementation has desirable security properties, as it is significantly more expensive for an attacker to tamper with the MSROM circuitry (§ A.14.2) than it is to modify the contents of the flash memory chip that stores the UEFI firmware. § B.4.3 and § B.6 describe the broad classes of attacks that an Intel platform can be subjected to.

The microcode that implements SEC performs MP initialization (§ A.13.2), as suggested in the SDM. The microcode then places the BSP into Cache-as-RAM (CAR) mode, looks up the PEI *Authenticated Code Module* (ACM) in the Firmware Interface Table (FIT), loads the PEI ACM into the cache, and verifies its signature (§ A.13.2) [43, 212, 213, 148, 42]. Given the structure of ACM signatures, we can conclude that Intel's microcode contains implementations of RSA decryption and of a variant of SHA hashing.

The PEI ACM is executed from the CPU's cache, after it is loaded by the microcode [43, 212, 42]. This removes the possibility for an attacker with physical access to the SPI flash chip to change the firmware's contents after the microcode computes its cryptographic hash, but before it is executed.

On motherboards compatible with LaGrande Server Extensions (LT-SX, also known as Intel TXT for servers), the firmware implementing PEI verifies that each CPU connected to motherboard supports LT-SX, and powers off the CPU sockets that don't hold processors that implement LT-SX [148]. This prevents an attacker from tampering with a TXT-protected VM by hot-plugging a CPU in a running computer that is inside TXT mode. When a hot-plugged CPU passes security tests, a hypervisor is notified that a new CPU is available.

The hypervisor updates its internal state, and sends the new CPU a SIPI. The new CPU executes a SIPI handler, inside microcode, that configures the CPU's state to match the state expected by the TXT hypervisor [148]. This implies that the AP initialization described in § A.13.2 is implemented in microcode.

A.14.5 Microcode Updates

The SDM explains that the microcode on Intel CPUs can be updated, and describes the process for applying an update. However, no detail about the contents of an update is provided. Analyzing Intel's microcode updates seems like a promising avenue towards discovering the microcode's structure. Unfortunately, the updates have so far proven to be inscrutable [34].

The microcode updates cannot be easily analyzed because they are encrypted, hashed with a cryptographic hash function like SHA-256, and signed using RSA or elliptic curve cryptography [212]. The update facility is implemented entirely in microcode, including the decryption and signature verification [212].

[78] independently used fault injection and timing analysis to conclude that each recent Intel microcode update is signed with a 2048-bit RSA key and a (possibly non-standard) 256-bit hash algorithm, which agrees with the findings above.

The microcode update implementation places the core's cache into No-Evict Mode (NEM, documented by the SDM) and copies the microcode update into the cache before verifying its signature [212]. The update facility also sets up an MTRR entry to protect the update's contents from modifications via DMA transfers [212] as it is verified and applied.

While Intel publishes the most recent microcode updates for each of its CPU models, the release notes associated with the updates are not publicly available. This is unfortunate, as the release notes could be used to confirm guesses that certain features are implemented in microcode.

However, some information can be inferred by reading through the Errata section in Intel's Specification Updates [91, 107, 109]. The phrase "it is possible for BIOS⁵ to contain

⁵Basic Input/Output System (BIOS) is the predecessor of UEFI-based firmware. Most Intel documentation, including the SDM, still uses the term BIOS to refer to firmware.

a workaround for this erratum” generally means that a microcode update was issued. For example, Errata AH in [91] implies that string instructions (REP MOV) are implemented in microcode, which was confirmed by Intel [13].

Errata AH43 and AH91 in [91], and AAK73 in [107] imply that address translation (§ A.5) is at least partially implemented in microcode. Errata AAK53, AAK63, and AAK70, AAK178 in [107], and BT138, BT210, in [109] imply that VM entries and exits (§ A.8.2) are implemented in microcode, which is confirmed by the APIC virtualization patent [173].

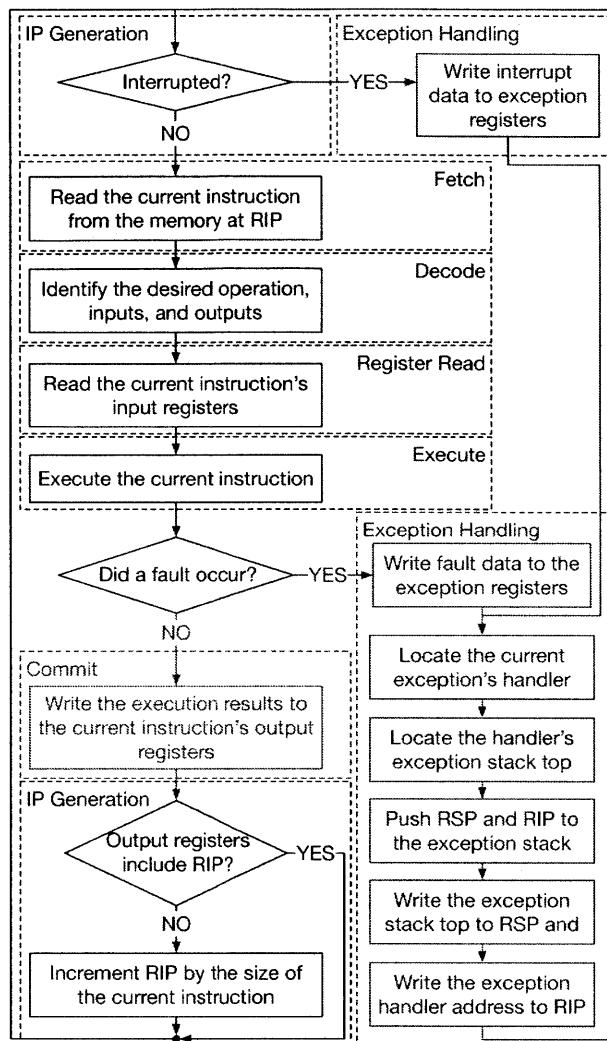


Figure A-3: A processor fetches instructions from the memory and executes them. The RIP register holds the address of the instruction to be executed.

<p>$SEND(op, addr, data) \rightarrow \emptyset$ Place a message containing the operation code op, the bus address $addr$, and the value $data$ on the bus.</p>
<p>$READ() \rightarrow (op, addr, value)$ Return the message that was written on the bus at the beginning of this clock cycle.</p>

Figure A-4: The system bus abstraction

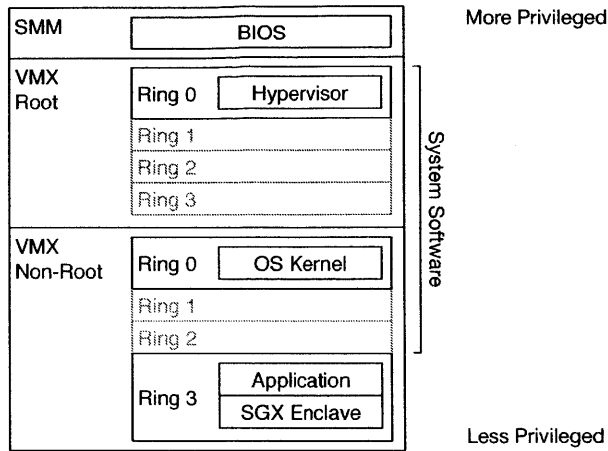


Figure A-5: The privilege levels in the x86 architecture, and the software that typically runs at each security level.

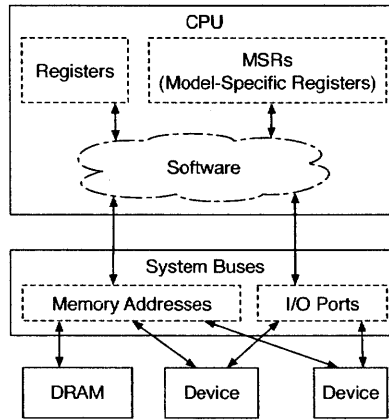


Figure A-6: The four physical address spaces used by an Intel CPU. The registers and MSRs are internal to the CPU, while the memory and I/O address spaces are used to communicate with DRAM and other devices via system buses.

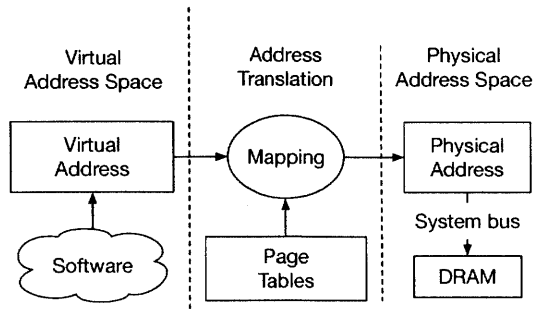


Figure A-7: Virtual addresses used by software are translated into physical memory addresses using a mapping defined by the page tables.

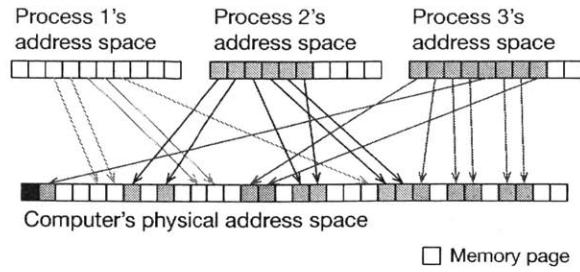


Figure A-8: The virtual memory abstraction gives each process its own virtual address space. The operating system multiplexes the computer's DRAM between the processes, while application developers build software as if it owns the entire computer's memory.

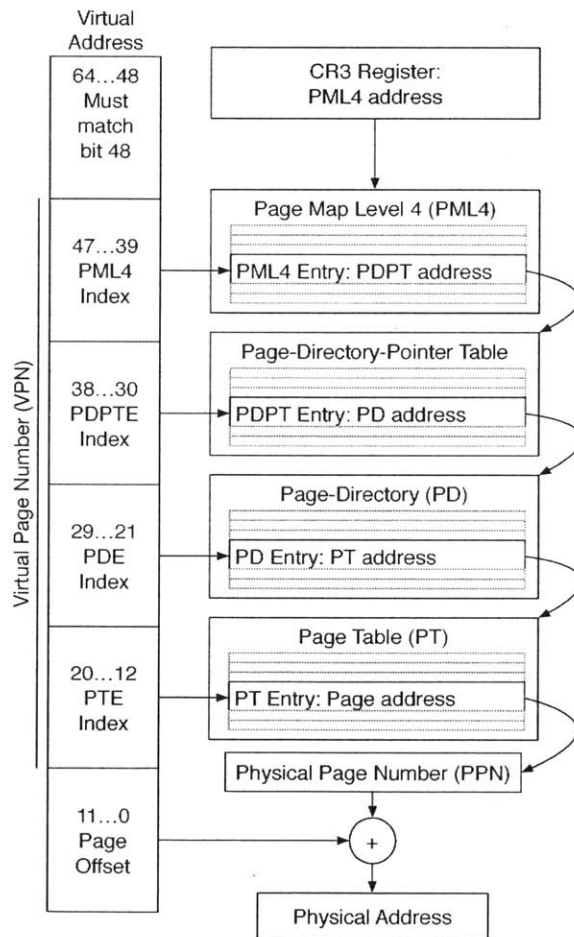


Figure A-9: IA-32e address translation takes in a 48-bit virtual address and outputs a 52-bit physical address.

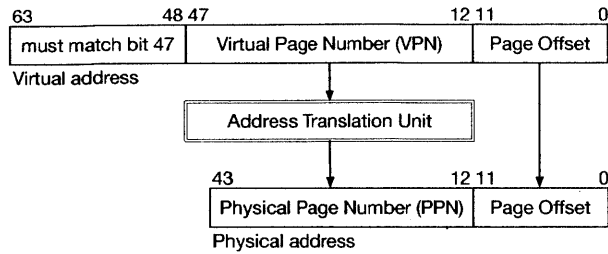


Figure A-10: Address translation can be seen as a mapping between virtual page numbers and physical page numbers.

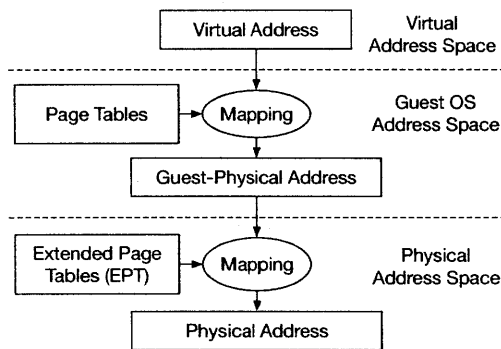


Figure A-11: Virtual addresses used by software are translated into physical memory addresses using a mapping defined by the page tables.

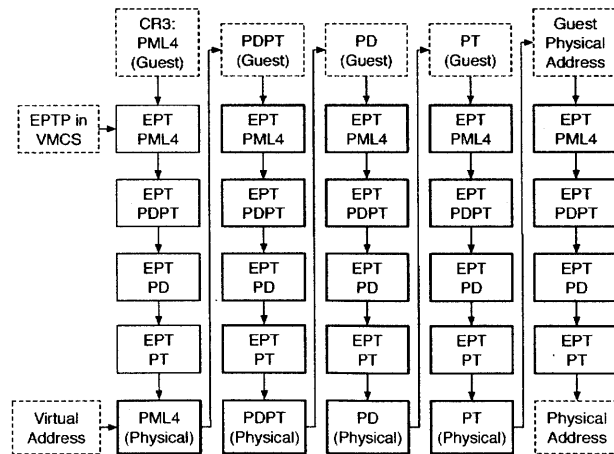


Figure A-12: Address translation when hardware virtualization is enabled. The kernel-managed page tables contain guest-physical addresses, so each level in the kernel's page table requires a full walk of the hypervisor's extended page table (EPT). A translation requires up to 20 memory accesses (the bold boxes), assuming the physical address of the kernel's PML4 is cached.

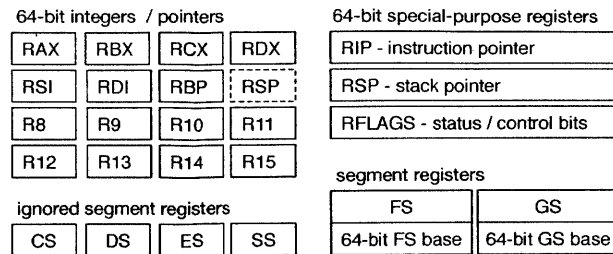


Figure A-13: CPU registers in the 64-bit Intel architecture. RSP can be used as a general-purpose register (GPR), e.g., in pointer arithmetic, but it always points to the top of the program's stack. Segment registers are covered in § A.7.

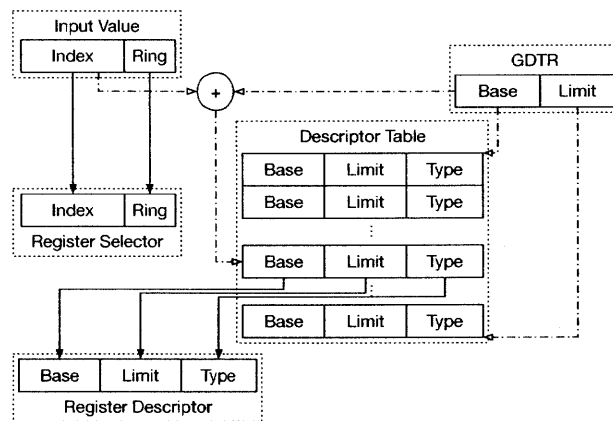


Figure A-14: Loading a segment register. The 16-bit value loaded by software is a selector consisting of an index and a ring number. The index selects a GDT entry, which is loaded into the descriptor part of the segment register.

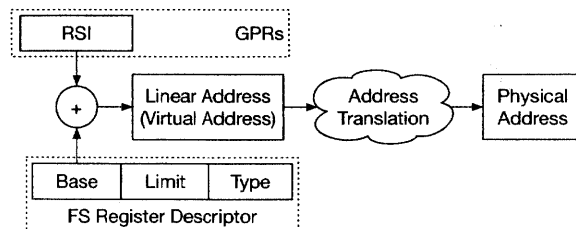


Figure A-15: Example address computation process for `MOV FS:[RDX], 0`. The segment's base address is added to the address in RDX before address translation (§ A.5) takes place.

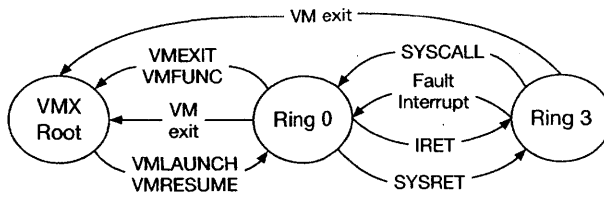


Figure A-16: Modern privilege switching methods in the 64-bit Intel architecture.

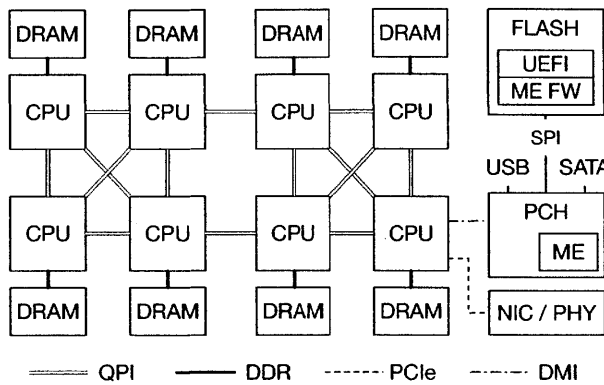


Figure A-17: The motherboard structures that are most relevant in a system security analysis.

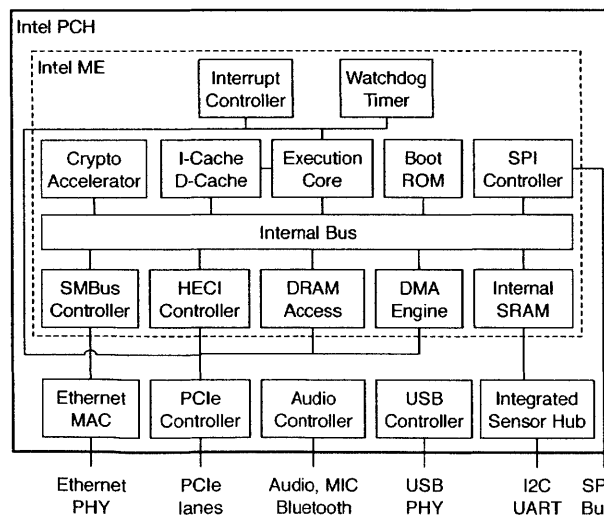


Figure A-18: The Intel Management Engine (ME) is an embedded computer hosted in the PCH. The ME has its own execution core, ROM and SRAM. The ME can access the host's DRAM via a memory controller and a DMA controller. The ME is remotely accessible over the network, as it has direct access to an Ethernet PHY via the SMBus.

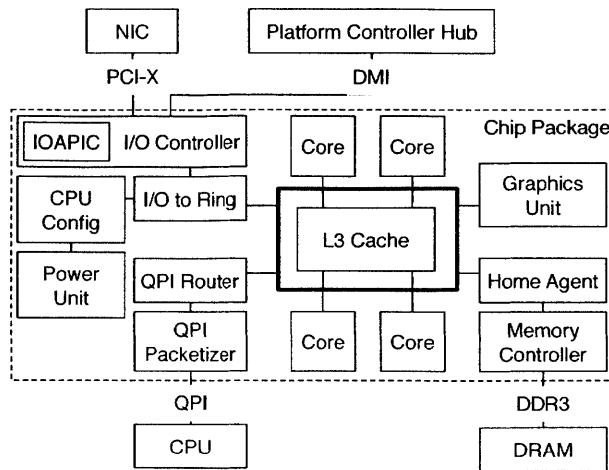


Figure A-19: The major components in a modern CPU package. § A.9.3 gives an uncore overview. § A.9.4 describes execution cores. § A.11.3 takes a deeper look at the uncore.

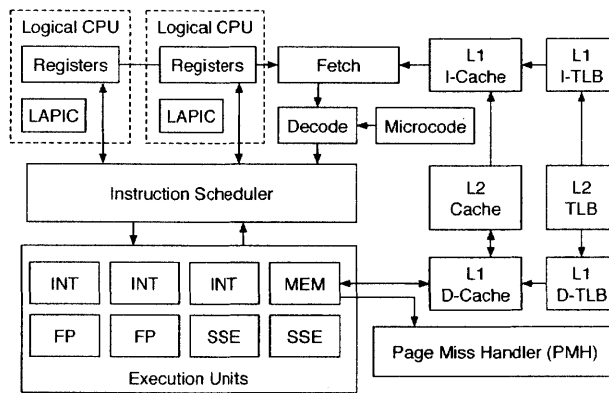


Figure A-20: CPU core with two logical processors. Each logical processor has its own execution context and LAPIC (§ A.12). All the other core resources are shared.

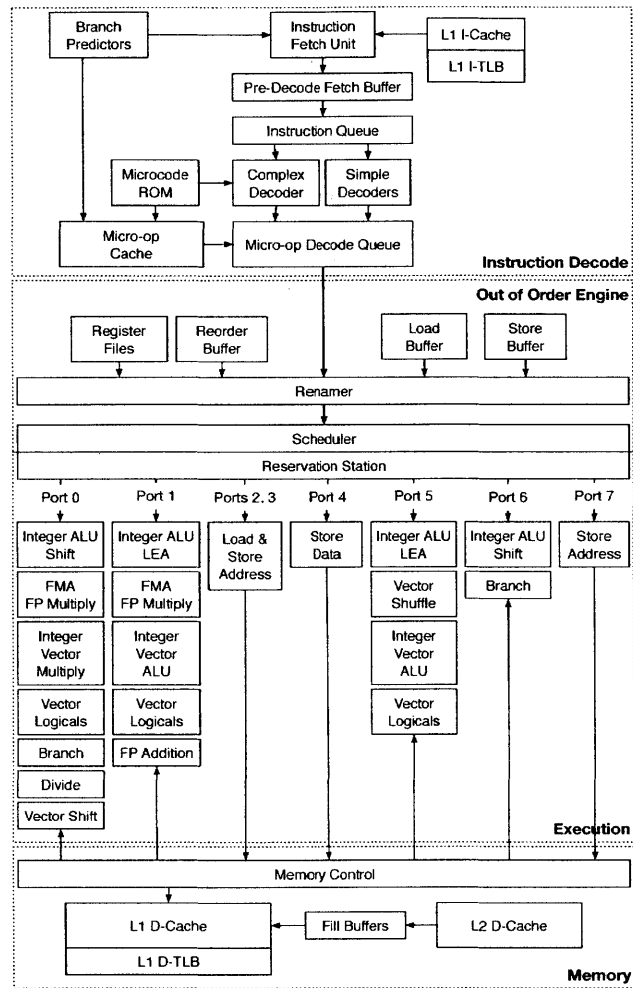


Figure A-21: The structures in a CPU core that are relevant to out-of-order and speculative execution. Instructions are decoded into micro-ops, which are scheduled on one of the execution unit's ports. The branch predictor enables speculative execution when a branch is encountered.

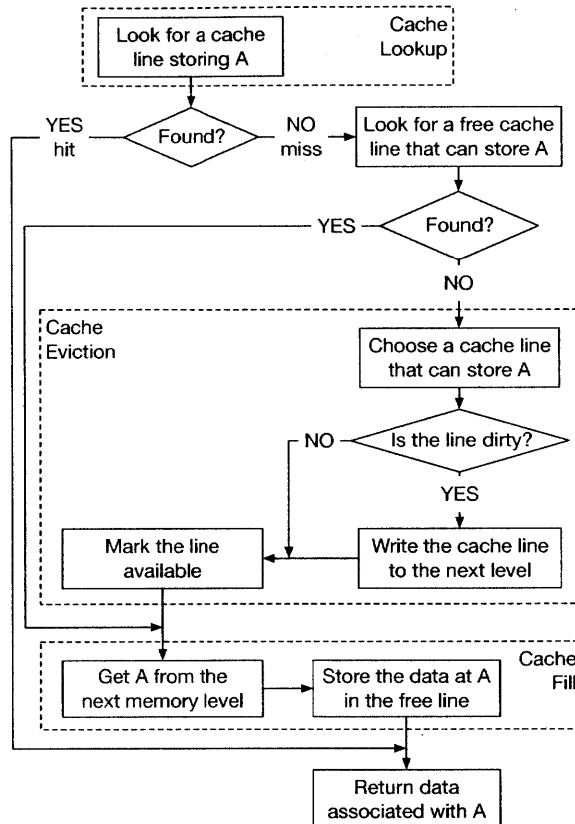


Figure A-22: The steps taken by a cache memory to resolve an access to a memory address A. A normal memory access (to cacheable DRAM) always triggers a cache lookup. If the access misses the cache, a fill is required, and a write-back might be required.

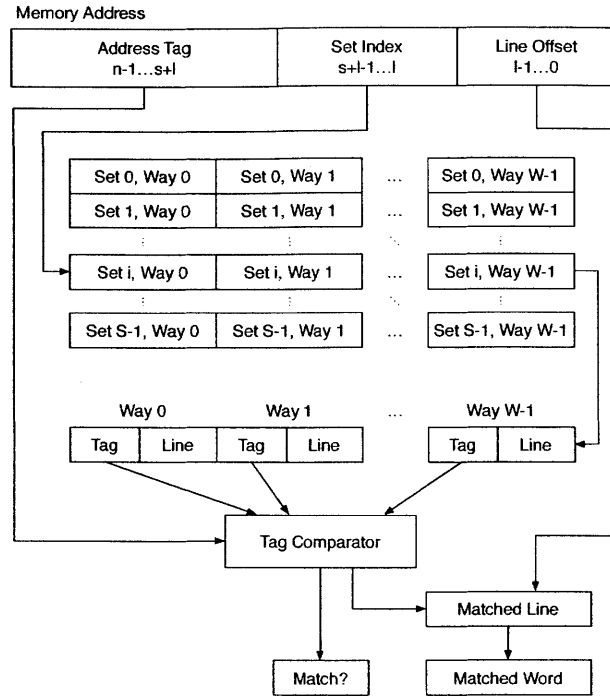


Figure A-23: Cache organization and lookup, for a W -way set-associative cache with 2^l -byte lines and $S = 2^s$ sets. The cache works with n -bit memory addresses. The lowest l address bits point to a specific byte in a cache line, the next s bytes index the set, and the highest $n - s - l$ bits are used to decide if the desired address is in one of the W lines in the indexed set.

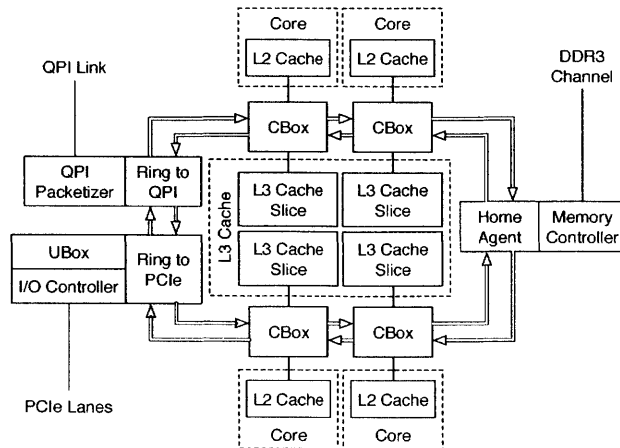


Figure A-24: The stops on the ring interconnect used for inter-core and core-uncore communication.

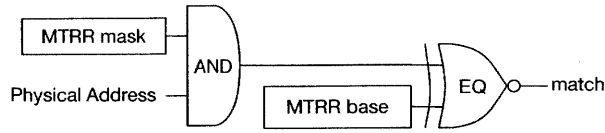


Figure A-25: The circuit for computing whether a physical address matches a memory type range. Assuming a CPU with 48-bit physical addresses, the circuit uses 36 AND gates and a binary tree of 35 XNOR (equality test) gates. The circuit outputs 1 if the address belongs to the range. The bottom 12 address bits are ignored, because memory type ranges must be aligned to 4 KB page boundaries.

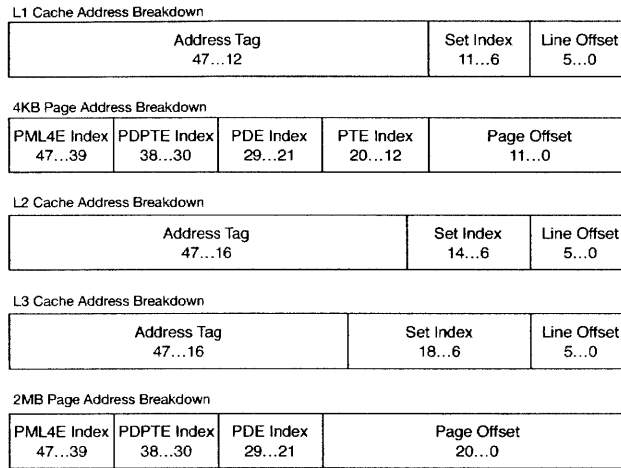


Figure A-26: Virtual addresses from the perspective of cache lookup and address translation. The bits used for the L1 set index and line offset are not changed by address translation, so the page tables do not impact L1 cache placement. The page tables do impact L2 and L3 cache placement. Using large pages (2 MB or 1 GB) is not sufficient to make L3 cache placement independent of the page tables, because of the LLC slice hashing function (§ A.11.3).

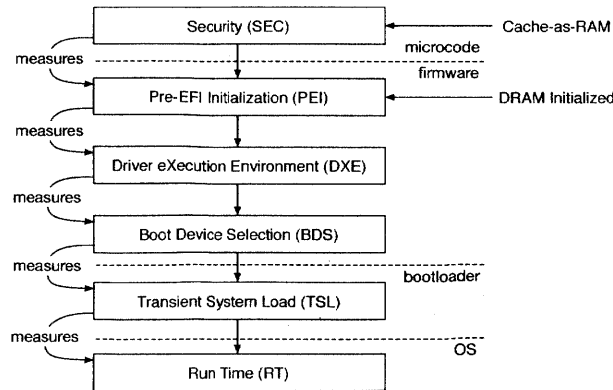


Figure A-27: The phases of the Platform Initialization process in the UEFI specification.

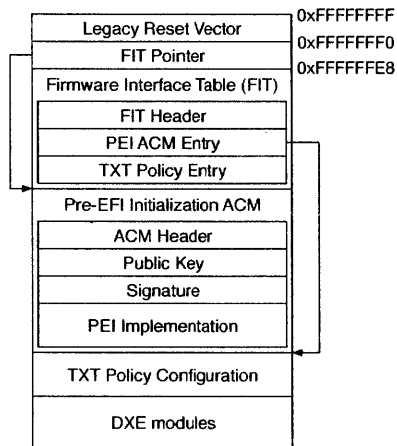


Figure A-28: The Firmware Interface Table (FIT) in relation to the firmware's memory map.

Appendix B

Security Background

Most systems rely on some cryptographic primitives for security. Unfortunately, these primitives have many assumptions, and building a secure system on top of them is a highly non-trivial endeavor. It follows that a system's security analysis should be particularly interested in what cryptographic primitives are used, and how they are integrated into the system.

§ B.1 and § B.2 lay the foundations for such an analysis by summarizing the primitives used by the secure architectures of interest to us, and by describing the most common constructs built using these primitives. § B.3 builds on these concepts and describes software attestation, which is the most popular method for establishing trust in a secure architecture.

Having looked at the cryptographic foundations for building secure systems, we turn our attention to the attacks that secure architectures must withstand. Besides from forming a security checklist for architecture design, these attacks build intuition for the design decisions in the architectures of interest to us.

The attacks that can be performed on a computer system are broadly classified into physical attacks and software attacks. In *physical attacks*, the attacker takes advantage of a system's physical implementation details to perform an operation that bypasses the limitations set by the computer system's software abstraction layers. In contrast, *software attacks* are performed solely by executing software on the victim computer. § B.4 summarizes the main types of physical attacks.

The distinction between software and physical attacks is particularly relevant in cloud

computing scenarios, where gaining software access to the computer running a victim’s software can be accomplished with a credit card backed by modest funds [161], whereas physical access is a more difficult prospect that requires trespass, coercion, or social engineering on the cloud provider’s employees.

However, the distinction between software and physical attacks is blurred by the attacks presented in § B.6, which exploit programmable peripherals connected to the victim computer’s bus in order to carry out actions that are normally associated with physical attacks.

While the vast majority of software attacks exploit a bug in a software component, there are a few attack classes that deserve attention from architecture designers. Memory mapping attacks, described in § B.7, become a possibility on architectures where the system software is not trusted. Cache timing attacks, summarized in § B.8 exploit microarchitectural behaviors that are completely observable in software, but dismissed by the security analyses of most systems.

B.1 Cryptographic Primitives

This section overviews the cryptosystems used by secure architectures. We are interested in cryptographic primitives that guarantee privacy, integrity, and freshness, and we treat these primitives as black boxes, focusing on their use in larger systems. [118] covers the mathematics behind cryptography, while [55] covers the topic of building systems out of cryptographic primitives. Tables B.1 and B.2 summarize the primitives covered in this section.

Guarantee	Primitive
Privacy	<i>Encryption</i>
Integrity	<i>MAC / Signatures</i>
Freshness	<i>Nonces + integrity</i>

Table B.1: Desirable security guarantees and primitives that provide them

A message whose *privacy* is protected can be transmitted over an insecure medium without an adversary being able to obtain the information in the message. When *integrity*

Guarantee	Symmetric Keys	Asymmetric Keys
Privacy	AES-GCM, AES-CTR	RSA with PKCS #1 v2.0
Integrity	HMAC-SHA-2 AES-GCM	DSS-RSA, DSS-ECC

Table B.2: Popular cryptographic primitives that are considered to be secure against today’s adversaries

protection is used, the receiver is guaranteed to either obtain a message that was transmitted by the sender, or to notice that an attacker tampered with the message’s content.

When multiple messages get transmitted over an untrusted medium, a *freshness* guarantee assures the receiver that she will obtain the latest message coming from the sender, or will notice an attack. A freshness guarantee is stronger than the equivalent integrity guarantee, because the latter does not protect against *replay attacks* where the attacker replaces a newer message with an older message coming from the same sender.

The following example further illustrates these concepts. Suppose Alice is a wealthy investor who wishes to either BUY or SELL an item every day. Alice cannot trade directly, and must relay her orders to her broker, Bob, over a network connection owned by Eve.

A communication system with privacy guarantees would prevent Eve from distinguishing between a BUY and a SELL order, as illustrated in Figure B-1. Without privacy, Eve would know Alice’s order before it is placed by Bob, so Eve would presumably gain a financial advantage at Alice’s expense.

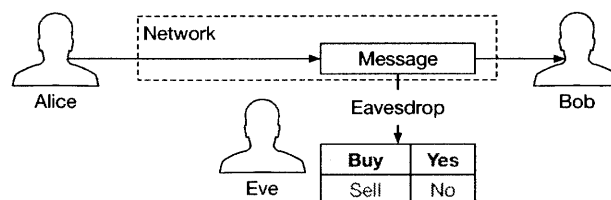


Figure B-1: In a privacy attack, Eve sees the message sent by Alice to Bob and can understand the information inside it. In this case, Eve can tell that the message is a **buy** order, and not a **sell** order.

A system with integrity guarantees would prevent Eve from replacing Alice’s message with a false order, as shown in Figure B-2. In this example, without integrity guarantees,

Eve could replace Alice's message with a SELL-EVERYTHING order, and buy Alice's assets at a very low price.

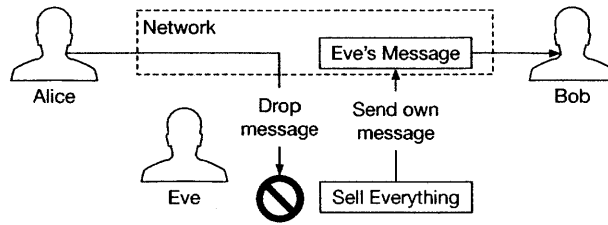


Figure B-2: In an integrity attack, Eve replaces Alice's message with her own. In this case, Eve sends Bob a **sell-everything** order. In this case, Eve can tell that the message is a **buy** order, and not a **sell** order.

Last, a communication system that guarantees freshness would ensure that Eve cannot perform the replay attack pictured in Figure B-3, where she would replace Alice's message with an older message. Without freshness guarantees, Eve could mount the following attack, which bypasses both privacy and integrity guarantees. Over a few days, Eve would copy and store Alice's messages from the network. When an order would reach Bob, Eve would observe the market and determine if the order was BUY or SELL. After building up a database of messages labeled BUY or SELL, Eve would replace Alice's message with an old message of her choice.

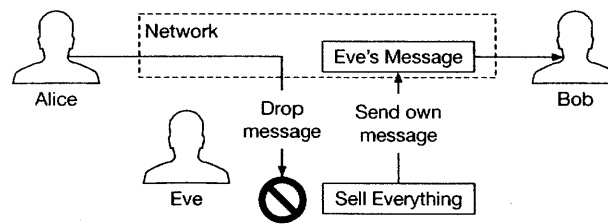


Figure B-3: In a freshness attack, Eve replaces Alice's message with a message that she sent at an earlier time. In this example, Eve builds a database of labeled messages over time, and is able to send Bob her choice of a BUY or a SELL order.

B.1.1 Cryptographic Keys

All cryptographic primitives that we describe here rely on *keys*, which are small pieces of information that must only be disclosed according to specific rules. A large part of a system's

security analysis focuses on ensuring that the keys used by the underlying cryptographic primitives are produced and handled according to the primitives' assumptions.

Each cryptographic primitive has an associated *key generation algorithm* that uses random data to produce a unique key. The random data is produced by a *cryptographically strong pseudo-random number generator* (CSPRNG) that expands a small amount of *random seed* data into a much larger amount of data, which is computationally indistinguishable from true random data. The random seed must be obtained from a true source of randomness whose output cannot be predicted by an adversary, such as the least significant bits of the temperature readings coming from a hardware sensor.

Symmetric key cryptography requires that all the parties in the system establish a shared *secret key*, which is usually referred to as “the key”. Typically, one party executes the key generation algorithm and securely transmits the resulting key to the other parties, as illustrated in Figure B-4. The channel used to distribute the key must provide privacy and integrity guarantees, which is a non-trivial logistical burden. The symmetric key primitives mentioned here do not make any assumption about the key, so the key generation algorithm simply grabs a fixed number of bits from the CSPRNG.

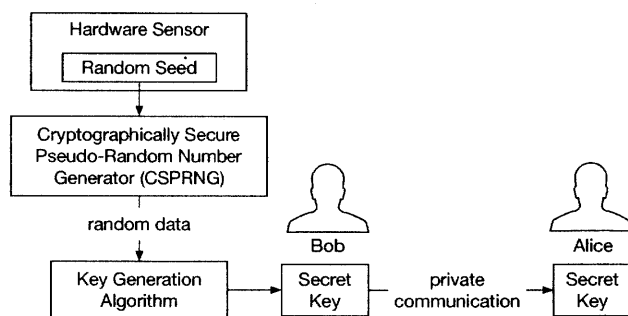


Figure B-4: In symmetric key cryptography, a secret key is shared by the parties that wish to communicate securely.

The defining feature of *asymmetric key* cryptography is that it does not require a private channel for key distribution. Each party executes the key generation algorithm, which produces a *private key* and a *public key* that are mathematically related. Each party's public key is distributed to the other parties over a channel with integrity guarantees, as shown in Figure B-5. Asymmetric key primitives are more flexible than their symmetric counterparts,

but are more complicated and consume more computational resources.

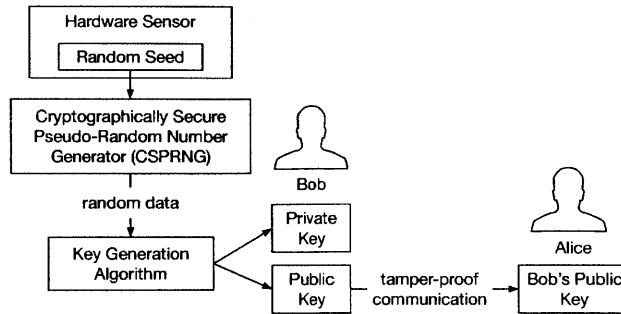


Figure B-5: An asymmetric key generation algorithm produces a private key and an associated public key. The private key is held confidential, while the public key is given to any party who wishes to securely communicate with the private key's holder.

B.1.2 Privacy

Many cryptosystems that provide integrity guarantees are built upon *block ciphers* that operate on fixed-size message blocks. The sender transforms a block using an *encryption* algorithm, and the receiver inverts the transformation using a *decryption* algorithm. The encryption algorithms in block ciphers obfuscate the message block's content in the output, so that an adversary who does not have the decryption key cannot obtain the original message block from the encrypted output.

Symmetric key encryption algorithms use the same secret key for encryption and decryption, as shown in Figure B-6, while asymmetric key block ciphers use the public key for encryption, and the corresponding private key for decryption, as shown in Figure B-7.

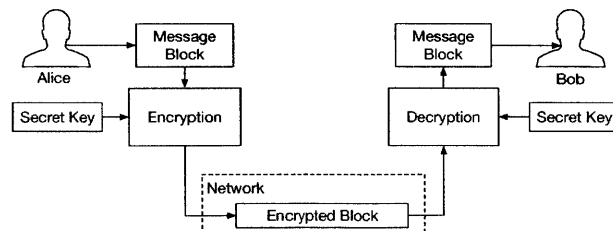


Figure B-6: In a symmetric key secure permutation (block cipher), the same secret key must be provided to both the encryption and the decryption algorithm.

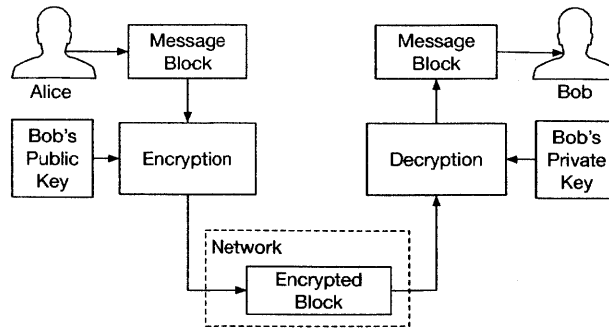


Figure B-7: In an asymmetric key block cipher, the encryption algorithm operates on a public key, and the decryption algorithm uses the corresponding private key.

The most popular block cipher based on symmetric keys at the time of this writing is the *American Encryption Standard* (AES) [41, 145], with two variants that operate on 128-bit blocks using 128-bit keys or 256-bit keys. AES is a *secure permutation* function, as it can transform any 128-bit block into another 128-bit block. Recently, the United States *National Security Agency* (NSA) required the use of 256-bit AES keys for protecting sensitive information [147].

The most deployed asymmetric key block cipher is the *Rivest-Shamir-Adelman* (RSA) [162] algorithm. RSA has variable key sizes, and 3072-bit key pairs are considered to provide the same security as 128-bit AES keys [22].

A block cipher does not necessarily guarantee privacy, when used on its own. A noticeable issue is that in our previous example, a block cipher would generate the same encrypted output for any of Alice's BUY orders, as they all have the same content. Furthermore, each block cipher has its own assumptions that can lead to subtle vulnerabilities if the cipher is used directly.

Symmetric key block ciphers are combined with operating modes to form symmetric encryption schemes. Most operating modes require a random *initialization vector* (IV) to be used for each message, as shown in Figure B-8. When analyzing the security of systems based on these cryptosystems, an understanding of the IV generation process is as important as ensuring the confidentiality of the encryption key.

Counter (CTR) and Cipher Block Chaining (CBC) are examples of operating modes recommended [49] by the United States *National Institute of Standards and Technology* (NIST),

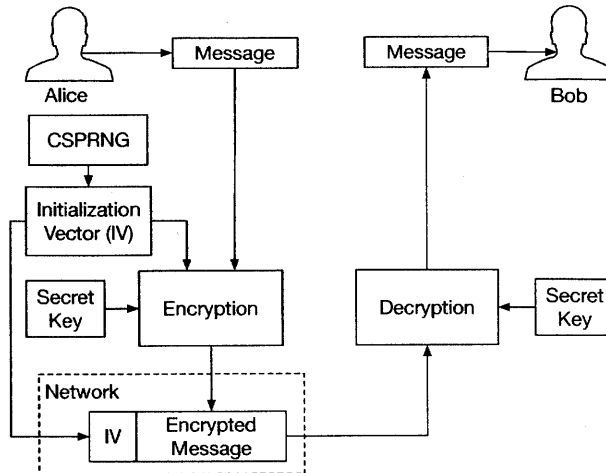


Figure B-8: Symmetric key block ciphers are combined with operating modes. Most operating modes require a random initialization vector (IV) to be generated for each encrypted message.

which informs the NSA's requirements. Combining a block cipher, such as AES, with an operating mode, such as CTR, results in an encryption method, such as AES-CTR, which can be used to add privacy guarantees.

In the asymmetric key setting, there is no concept equivalent to operating modes. Each block cipher has its own assumptions, and requires a specialized scheme for general-purpose usage.

The RSA algorithm is used in conjunction with *padding methods*, the most popular of which are the methods described in the *Public-Key Cryptography Standard (PKCS) #1* versions 1.5 [114] and 2.0 [115]. A security analysis of a system that uses RSA-based encryption must take the padding method into consideration. For example, the padding in PKCS #1 v1.5 can leak the private key under certain circumstances [25]. While PKCS #1 v2.0 solves this issue, it is complex enough that some implementations have their own security issues [138].

Asymmetric encryption algorithms have much higher computational requirements than symmetric encryption algorithms. Therefore, when non-trivial quantities of data is encrypted, the sender generates a single-use secret key that is used to encrypt the data, and encrypts the secret key with the receiver's public key, as shown in Figure B-9.

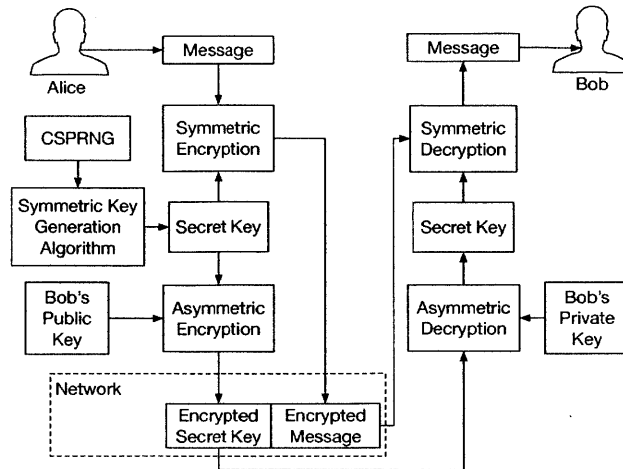


Figure B-9: Asymmetric key encryption is generally used to bootstrap a symmetric key encryption scheme.

B.1.3 Integrity

Many cryptosystems that provide integrity guarantees are built upon *secure hashing* functions. These hash functions operate on an unbounded amount of input data and produce a small fixed-size output. Secure hash functions have a few guarantees, such as *pre-image resistance*, which states that an adversary cannot produce input data corresponding to a given hash output.

At the time of this writing, the most popular secure hashing function is the *Secure Hashing Algorithm* (SHA) [52]. However, due to security issues in SHA-1 [179], new software is recommended to use at least 256-bit SHA-2 [23] for secure hashing.

The SHA hash functions are members of a large family of *block hash functions* that consume their input in fixed-size message blocks, and use a fixed-size internal state. A block hash function is used as shown in Figure B-10. An INITIALIZE algorithm is first invoked to set the internal state to its initial values. An EXTEND algorithm is executed for each message block in the input. After the entire input is consumed, a FINALIZE algorithm produces the hash output from the internal state.

In the symmetric key setting, integrity guarantees are obtained using a *Message Authentication Code* (MAC) cryptosystem, illustrated in Figure B-11. The sender uses a MAC algorithm that reads in a symmetric key and a variable-length message, and produces a

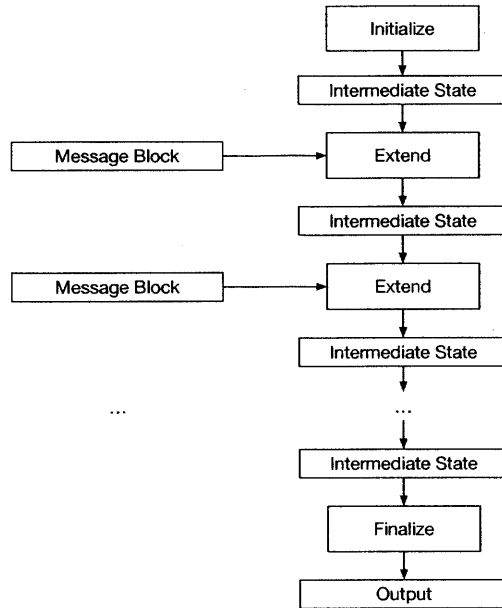


Figure B-10: A block hash function operates on fixed-size message blocks and uses a fixed-size internal state.

fixed-length, short *MAC tag*. The receiver provides the original message, the symmetric key, and the MAC tag to a MAC verification algorithm that checks the authenticity of the message.

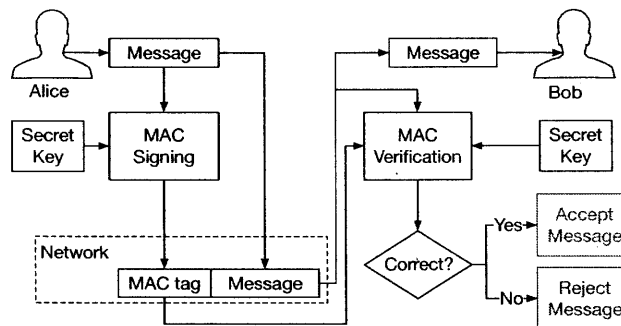


Figure B-11: In the symmetric key setting, integrity is assured by computing a Message Authentication Code (MAC) tag and transmitting it over the network along the message. The receiver feeds the MAC tag into a verification algorithm that checks the message's authenticity.

The key property of MAC cryptosystems is that an adversary cannot produce a MAC tag that will validate a message without the secret key.

Many MAC cryptosystems do not have a separate MAC verification algorithm. Instead,

the receiver checks the authenticity of the MAC tag by running the same algorithm as the sender to compute the expected MAC tag for the received message, and compares the output with the MAC tag received from the network.

This is the case for the *Hash Message Authentication Code (HMAC)* [127] generic construction, whose operation is illustrated in Figure B-12. HMAC can use any secure hash function, such as SHA, to build a MAC cryptosystem.

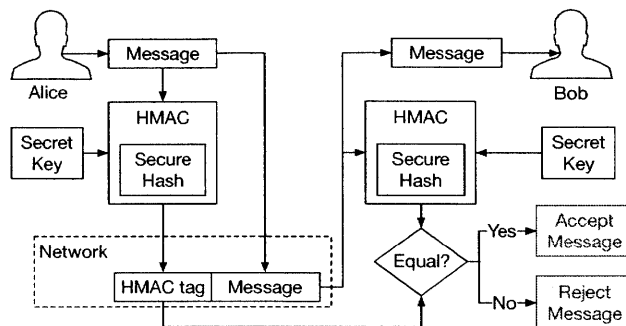


Figure B-12: In the symmetric key setting, integrity is assured by computing a Hash-based Message Authentication Code (HMAC) and transmitting it over the network along the message. The receiver re-computes the HMAC and compares it against the version received from the network.

Asymmetric key primitives that provide integrity guarantees are known as *signatures*. The message sender provides her private key to a *signing* algorithm, and transmits the output signature along with the message, as shown in Figure B-13. The message receiver feeds the sender's public key and the signature to a *signature verification* algorithm, which returns TRUE if the message matches the signature, and FALSE if the message has been tampered with.

Signing algorithms can only operate on small messages and are computationally expensive. Therefore, in practice, the message to be transmitted is first ran through a cryptographically strong hash function, and the hash is provided as the input to the signing algorithm.

At the time of this writing, the most popular choice for guaranteeing integrity in shared secret settings is HMAC-SHA, an HMAC function that uses SHA for hashing.

Authenticated encryption, which combines a block cipher with an operating mode that offers both privacy and integrity guarantees, is often an attractive alternative to HMAC.

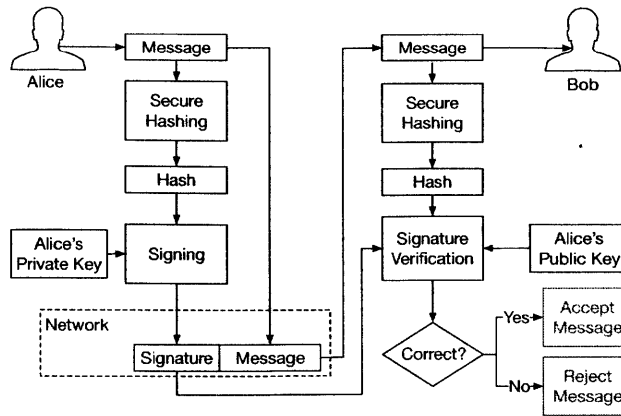


Figure B-13: Signature schemes guarantee integrity in the asymmetric key setting. Signatures are created using the sender’s private key, and are verified using the corresponding public key. A cryptographically secure hash function is usually employed to reduce large messages to small hashes, which are then signed.

The most popular authenticated encryption operating mode is *Galois/Counter operation mode* (GCM) [141], which has earned NIST’s recommendation [51] when combined with AES to form AES-GCM.

The most popular signature scheme combines the RSA encryption algorithms with a padding schemes specified in PKCS #1, as illustrated in Figure B-14. Recently, elliptic curve cryptography (ECC) [124] has gained a surge in popularity, thanks to its smaller key sizes. For example, a 384-bit ECC key is considered to be as secure as a 3072-bit RSA key [22, 147]. The NSA requires the Digital Signature Standard (DSS)[146], which specifies schemes based on RSA and ECC.

B.1.4 Freshness

Freshness guarantees are typically built on top of a system that already offers integrity guarantees, by adding a unique piece of information to each message. The main challenge in freshness schemes comes down to economically maintaining the state needed to generate the unique pieces of information on the sender side, and verify their uniqueness on the receiver side.

A popular solution for gaining freshness guarantees relies on *nonces*, single-use random numbers. Nonces are attractive because the sender does not need to maintain any state; the

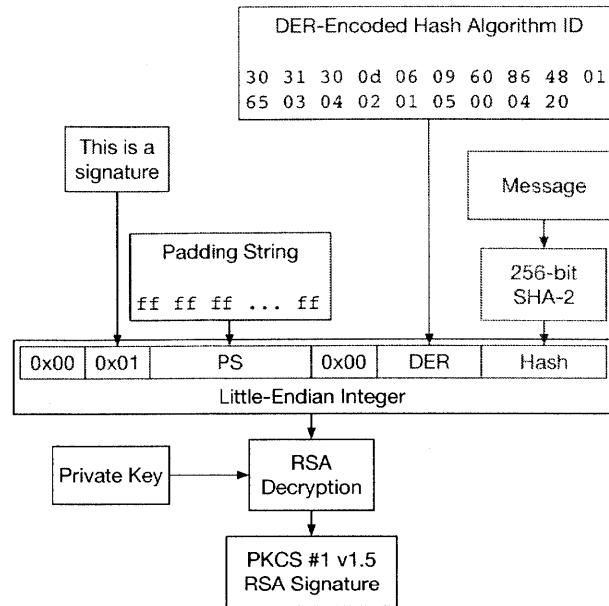


Figure B-14: The RSA signature scheme with PKCS #1 v1.5 padding specified in RFC 3447 combines a secure hash of the signed message with a DER-encoded specification of the secure hash algorithm used by the signature, and a padding string whose bits are all set to 1. Everything except for the secure hash output is considered to be a part of the PKCS #1 v1.5 padding.

receiver, however, must store the nonces of all received messages.

Nonces are often combined with a message timestamping and expiration scheme, as shown in Figure B-15. An expiration can greatly reduce the receiver's storage requirement, as the nonces for expired messages can be safely discarded. However, the scheme depends on the sender and receiver having synchronized clocks. The message expiration time is a compromise between the desire to reduce storage costs, and the need to tolerate clock skew and delays in message transmission and processing.

Alternatively, nonces can be used in challenge-response protocols, in a manner that removes the storage overhead concerns. The challenger generates a nonce and embeds it in the challenge message. The response to the challenge includes an acknowledgement of the embedded nonce, so the challenger can distinguish between a fresh response and a replay attack. The nonce is only stored by the challenger, and is small in comparison to the rest of the state needed to validate the response.

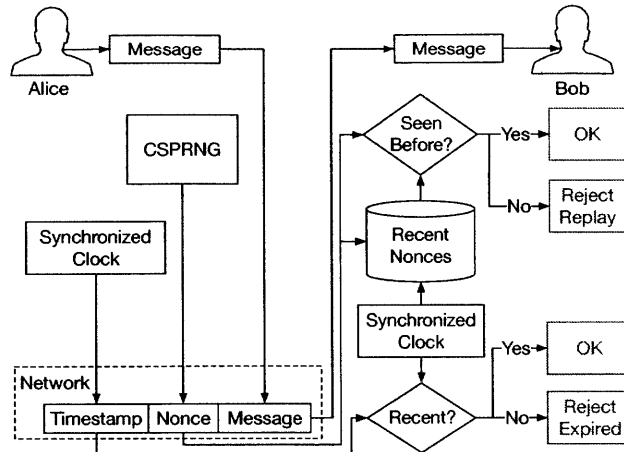


Figure B-15: Freshness guarantees can be obtained by adding timestamped nonces on top of a system that already offers integrity guarantees. The sender and the receiver use synchronized clocks to timestamp each message and discard unreasonably old messages. The receiver must check the nonce in each new message against a database of the nonces in all the unexpired messages that it has seen.

B.2 Cryptographic Constructs

This section summarizes two constructs that are built on the cryptographic primitives described in § B.1, and are used in the rest of this work.

B.2.1 Certificate Authorities

Asymmetric key cryptographic primitives assume that each party has the correct public keys for the other parties. This assumption is critical, as the entire security argument of an asymmetric key system rests on the fact that certain operations can only be performed by the owners of the private keys corresponding to the public keys. More concretely, if Eve can convince Bob that her own public key belongs to Alice, Eve can produce message signatures that seem to come from Alice.

The introductory material in § B.1 assumed that each party transmits their public key over a channel with integrity guarantees. In practice, this is not a reasonable assumption, and the secure distribution of public keys is still an open research problem.

The most widespread solution to the public key distribution problem is the Certificate Authority (CA) system, which assumes the existence of a trusted authority whose public key

is securely transmitted to all the other parties in the system.

The CA is responsible for securely obtaining the public key of each party, and for issuing a *certificate* that binds a party's identity (e.g., "Alice") to its public key, as shown in Figure B-16.

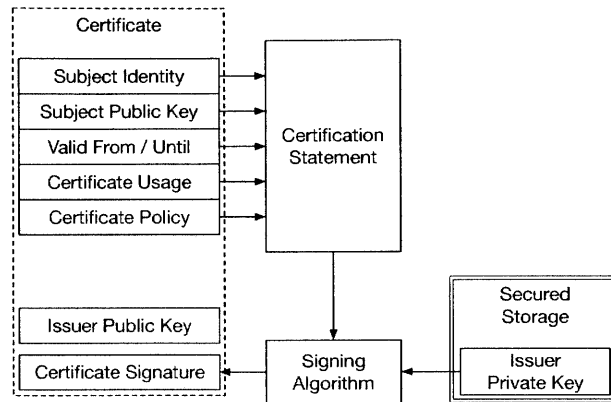


Figure B-16: A certificate is a statement signed by a certificate authority (issuer) binding the identity of a subject to a public key.

A certificate is essentially a cryptographic signature produced by the private key of the certificate's *issuer*, who is generally a CA. The message signed by the issuer states that a public key belongs to a *subject*. The certificate message generally contains identifiers that state the intended use of the certificate, such as "the key in this certificate can only be used to sign e-mail messages". The certificate message usually also includes an identifier for the issuer's *certification policy*, which summarizes the means taken by the issuer to ensure the authenticity of the subject's public key.

A major issue in a CA system is that there is no obvious way to revoke a certificate. A revocation mechanism is desirable to handle situations where a party's private key is accidentally exposed, to avoid having an attacker use the certificate to impersonate the compromised party. While advanced systems for certificate revocation have been developed, the first line of defense against key compromise is adding expiration dates to certificates.

In a CA system, each party presents its certificate along with its public key. Any party that trusts the CA and has obtained the CA's public key securely can verify any certificate using the process illustrated in Figure B-17.

One of the main drawbacks of the CA system is that the CA's private key becomes a very

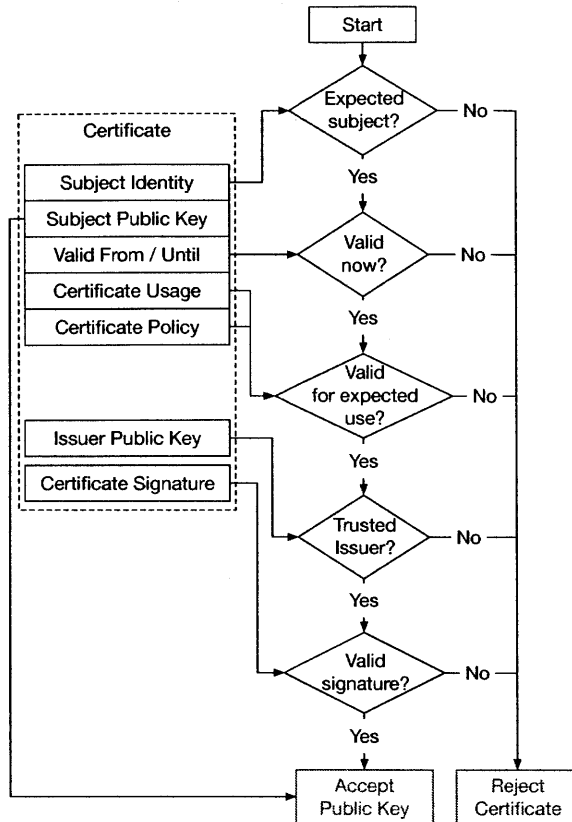


Figure B-17: A certificate issued by a CA can be validated by any party that has securely obtained the CA’s public key. If the certificate is valid, the subject public key contained within can be trusted to belong to the subject identified by the certificate.

attractive attack target. This issue is somewhat mitigated by minimizing the use of the CA’s private key, which reduces the opportunities for its compromise. The authority described above becomes the *root CA*, and their private key is only used to produce certificates for the *intermediate CAs* who, in turn, are responsible for generating certificates for the other parties in the system, as shown in Figure B-18.

In hierarchical CA systems, the only public key that gets distributed securely to all the parties is the root CA’s public key. Therefore, when two parties wish to interact, each party must present their own certificate, as well as the certificate of the issuing CA. For example, given the hierarchy in Figure B-18, Alice would prove the authenticity of her public key to Bob by presenting her certificate, as well as the certificate of Intermediate CA 1. Bob would first use the steps in Figure B-17 to validate Intermediate CA 1’s certificate against

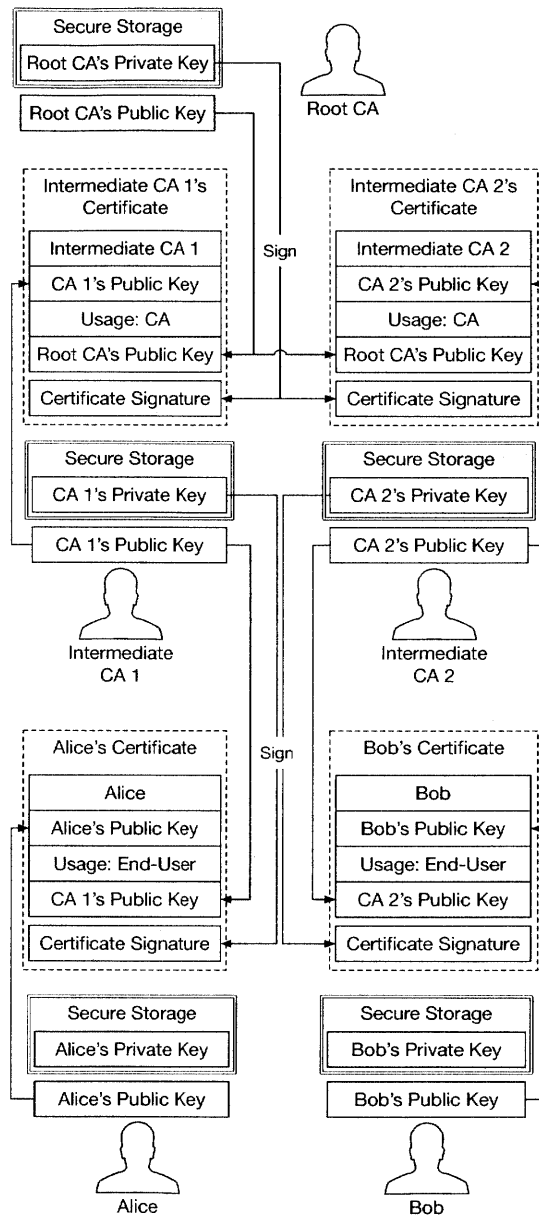


Figure B-18: A hierarchical CA structure minimizes the usage of the root CA's private key, reducing the opportunities for it to get compromised. The root CA only signs the certificates of intermediate CAs, which sign the end users' certificates.

the root CA's public key, which would assure him of the authenticity of Intermediate CA 1's public key. Bob would then validate Alice's certificate using Intermediate CA 1's public key, which he now trusts.

In most countries, the government issues ID cards for its citizens, and therefore acts as a certificate authority. An ID card, shown in Figure B-19, is a certificate that binds a subject's identity, which is a full legal name, to the subject's physical appearance, which is used as a public key.

The CA system is very similar to the identity document (ID card) systems used to establish a person's identity, and a comparison between the two may help further the reader's understanding of the concepts in the CA system.

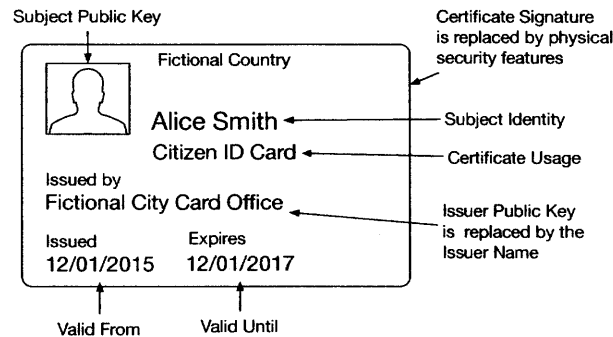


Figure B-19: An ID card is a certificate that binds a subject's full legal name (identity) to the subject's physical appearance, which acts as a public key.

Each government's ID card issuing operations are regulated by laws, so an ID card's issue date can be used to track down the laws that make up its certification policy. Last, the security of ID cards does not (yet) rely on cryptographic primitives. Instead, ID cards include physical security measures designed to deter tampering and prevent counterfeiting.

B.2.2 Key Agreement Protocols

The initial design of symmetric key primitives, introduced in § B.1, assumed that when two parties wish to interact, one party generates a secret key and shares it with the other party using a communication channel with privacy and integrity guarantees. In practice, a pre-existing secure communication channel is rarely available.

Key agreement protocols are used by two parties to establish a shared secret key, and only require a communication channel with integrity guarantees. Figure B-20 outlines the Diffie-Hellman Key Exchange (DKE) [46] protocol, which should give the reader an intuition for how key agreement protocols work.

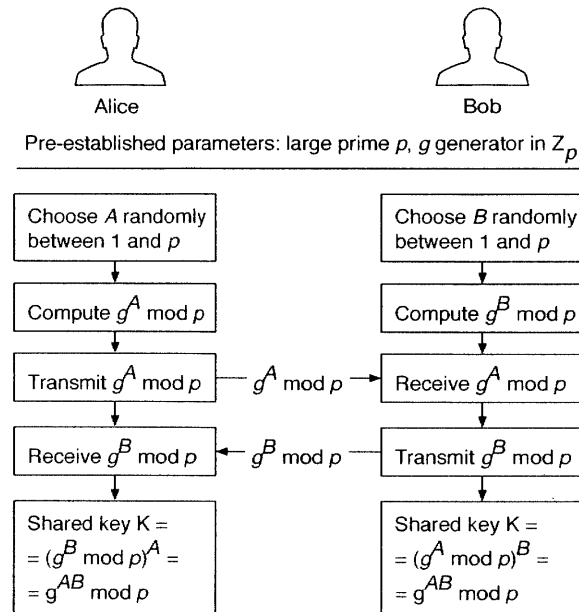


Figure B-20: In the Diffie-Hellman Key Exchange (DKE) protocol, Alice and Bob agree on a shared secret key $K = g^{AB} \bmod p$. An adversary who observes $g^A \bmod p$ and $g^B \bmod p$ cannot compute K .

This work is interested in using key agreement protocols to build larger systems, so we will neither explain the mathematic details in DKE, nor prove its correctness. We note that both Alice and Bob derive the same shared secret key, $K = g^{AB} \bmod p$, without ever transmitting K . Furthermore, the messages transmitted in DKE, namely $g^A \bmod p$ and $g^B \bmod p$, are not sufficient for an eavesdropper Eve to determine K , because efficiently solving for x in $g^x \bmod p$ is an open problem assumed to be very difficult.

Key agreement protocols require a communication channel with integrity guarantees. If an active adversary Eve can tamper with the messages transmitted by Alice and Bob, she can perform a *man-in-the-middle* (MITM) attack, as illustrated in Figure B-21.

In a MITM attack, Eve intercepts Alice's first key exchange message, and sends Bob her own message. Eve then intercepts Bob's response and replaces it with her own, which

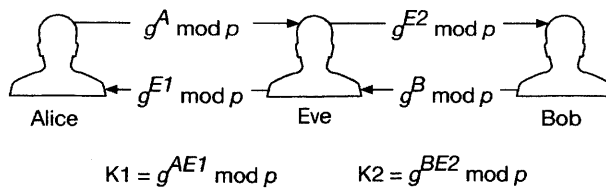


Figure B-21: Any key agreement protocol is vulnerable to a man-in-the-middle (MITM) attack. The active attacker performs key agreements and establishes shared secrets with both parties. The attacker can then forward messages between the victims, in order to observe their communication. The attacker can also send its own messages to either, impersonating the other victim.

she sends to Alice. Eve effectively performs key exchanges with both Alice and Bob, establishing a shared secret with each of them, with neither Bob nor Alice being aware of her presence.

After establishing shared keys with both Alice and Bob, Eve can choose to observe the communication between Alice and Bob, by forwarding messages between them. For example, when Alice transmits a message, Eve can decrypt it using K1, the shared key between herself and Alice. Eve can then encrypt the message with K2, the key established between Bob and herself. While Bob still receives Alice's message, Eve has been able to see its contents.

Furthermore, Eve can impersonate either party in the communication. For example, Eve can create a message, encrypt it with K2, and then send it to Bob. As Bob thinks that K2 is a shared secret key established between himself and Alice, he will believe that Eve's message comes from Alice.

MITM attacks on key agreement protocols can be foiled by authenticating the party who sends the last message in the protocol (in our examples, Bob) and having them sign the key agreement messages. When a CA system is in place, Bob uses his public key to sign the messages in the key agreement and also sends Alice his certificate, along with the certificates for any intermediate CAs. Alice validates Bob's certificate, ensures that the subject identified by the certificate is whom she expects (Bob), and verifies that the key agreement messages exchanged between herself and Bob match the signature provided by Bob.

In conclusion, a key agreement protocol can be used to bootstrap symmetric key primitives from an asymmetric key signing scheme, where only one party needs to be able to sign messages.

B.3 Software Attestation Overview

The security of systems that employ trusted processors hinges on *software attestation*. The software running inside an *isolated container* established by trusted hardware can ask the hardware to sign (§ B.1.3) a small piece of *attestation data*, producing an *attestation signature*. Besides from the attestation data, the signed message includes a *measurement* that uniquely identifies the software inside the container. Therefore, an attestation signature can be used to convince a *verifier* that the attestation data was produced by a specific piece of software, which is hosted inside a container that is isolated by trusted hardware from outside interference.

Each hardware platform discussed in this section uses a slightly different software attestation scheme. Platforms differ by the amount of software that executes inside an isolated container, by the isolation guarantees provided to the software inside a container, and by the process used to obtain a container's measurement. The threat model and security properties of each trusted hardware platform follow directly from the design choices outlined above, so a good understanding of attestation is a prerequisite to discussing the differences between existing platforms.

B.3.1 Authenticated Key Agreement

Software attestation can be combined with a key agreement protocol (§ B.2.2), as software attestation provides the authentication required by the key agreement protocol. The resulting protocol can assure a verifier that it has established a shared secret with a specific piece of software, hosted inside an isolated container created by trusted hardware. The next paragraph outlines the augmented protocol, using Diffie-Hellman Key Exchange (DKE) [46] as an example of the key exchange protocol.

The verifier starts executing the key exchange protocol, and sends the first message, g^A ,

to the software inside the secure container. The software inside the container produces the second key exchange message, g^B , and asks the trusted hardware to attest the cryptographic hash of both key exchange messages, $h(g^A||g^B)$. The verifier receives the second key exchange and attestation signature, and authenticates the software inside the secure container by checking all the signatures along the *attestation chain* of trust shown in Figure B-22.

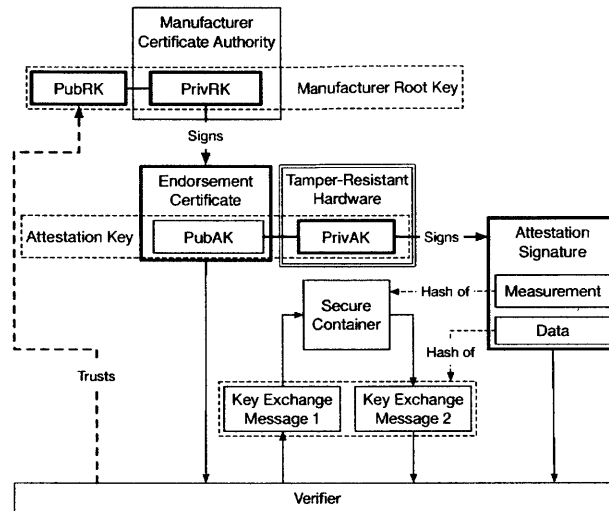


Figure B-22: The chain of trust in software attestation. The root of trust is a manufacturer key, which produces an endorsement certificate for the secure processor’s attestation key. The processor uses the attestation key to produce the attestation signature, which contains a cryptographic hash of the container and a message produced by the software inside the container.

The chain of trust used in software attestation is rooted at a signing key owned by the hardware manufacturer, which must be trusted by the verifier. The manufacturer acts as a Certificate Authority (CA, § B.2.1), and provisions each secure processor that it produces with a unique *attestation key*, which is used to produce *attestation signatures*. The manufacturer also issues an *endorsement certificate* for each secure processor’s attestation key. The certificate indicates that the key is meant to be used for software attestation. The certification policy generally states that, at the very least, the private part of the attestation key be stored in tamper-resistant hardware, and only be used to produce attestation signatures.

A secure processor identifies each isolated container by storing a cryptographic hash of the code and data loaded inside the container. When the processor is asked to sign a piece of attestation data, it uses the cryptographic hash associated with the container

as the measurement in the attestation signature. After a verifier validates the processor's attestation key using its endorsement certificate, the verifier ensures that the signature is valid, and that the measurement in the signature belongs to the software with which it expects to communicate. Having checked all the links in the attestation chain, the verifier has authenticated the other party in the key exchange, and is assured that it now shares a secret with the software that it expects, running in an isolated container on hardware that it trusts.

B.3.2 The Role of Software Measurement

The measurement that identifies the software inside a secure container is always computed using a secure hashing algorithm (§ B.1.3). Trusted hardware designs differ in their secure hash function choices, and in the data provided to the hash function. However, all the designs share the principle that each step taken to build a secure container contributes data to its measurement hash.

The philosophy behind software attestation is that the computer's owner can load any software she wishes in a secure container. However, the computer owner is assumed to have an incentive to participate in a distributed system where the secure container she built is authenticated via software attestation. Without the requirement to undergo software attestation, the computer owner can build any container without constraints, which would make it impossible to reason about the security properties of the software inside the container.

By the argument above, a trusted hardware design based on software attestation must assume that each container is involved in software attestation, and that the remote party will refuse to interact with a container whose reported measurement does not match the expected value set by the distributed system's author.

For example, a cloud infrastructure provider should be able to use the secure containers provided by trusted hardware to run any software she wishes on her computers. However, the provider makes money by renting her infrastructure to customers. If security savvy customers are only willing to rent containers provided by trusted hardware, and use software attestation to authenticate the containers that they use, the cloud provider will have a strong

financial incentive to build the customers' containers according to their specifications, so that the containers pass the software attestation.

A container's measurement is computed using a secure hashing algorithm, so the only method of building a container that matches an expected measurement is to follow the exact sequence of steps specified by the distributed system's author. The cryptographic properties of the secure hash function guarantee that if the computer's owner strays in any way from the prescribed sequence of steps, the measurement of the created container will not match the value expected by the distributed system's author, so the container will be rejected by the software attestation process.

Therefore, it makes sense to state that a trusted hardware design's measurement scheme guarantees that a property has a certain value in a secure container. The precise meaning of this phrase is that the property's value determines the data used to compute the container's measurement, so an expected measurement hash effectively specifies an expected value for the property. All containers in a distributed system that correctly uses software attestation will have the desired value for the given property.

For example, the measuring scheme used by trusted hardware designed for cloud infrastructure should guarantee that the container's memory was initialized using the customer's content, often referred to as an image.

B.4 Physical Attacks

Physical attacks are generally classified according to their cost, which factors in the equipment needed to carry out the attack and the attack's complexity. Joe Grand's DefCon presentation [73] provides a good overview with a large number of intuition-building figures and photos.

The simplest type of physical attack is a denial of service attack performed by disconnecting the victim computer's power supply or network cable. The threat models of most secure architectures ignore this attack, because denial of service can also be achieved by software attacks that compromise system software such as the hypervisor.

B.4.1 Port Attacks

Slightly more involved attacks rely on connecting a device to an existing port on the victim computer's case or motherboard (§ A.9.1). A simple example is a *cold boot attack*, where the attacker plugs in a USB flash drive into the victim's case and causes the computer to boot from the flash drive, whose malicious system software receives unrestricted access to the computer's peripherals.

More expensive physical attacks that still require relatively little effort target the debug ports of various peripherals. The cost of these attacks is generally dominated by the expense of acquiring the development kits needed to connect to the debug ports. For example, recent Intel processors include the Generic Debug eXternal Connection (GDXC) [209, 129], which collects and filters the data transferred by the uncore's ring bus (§ A.11.3), and reports it to an external debugger.

The threat models of secure architectures generally ignore debug port attacks, under the assumption that devices sold for general consumption have their debug ports irreversibly disabled. In practice, manufacturers have strong incentives to preserve debugging ports in production hardware, as this facilitates the diagnosis and repair of defective units. Due to insufficient documentation on this topic, we ignore the possibility of GDXC-based attacks.

B.4.2 Bus Tapping Attacks

More complex physical attacks consist of installing a device that taps a bus on the computer's motherboard (§ A.9.1). *Passive attacks* are limited to monitoring the bus traffic, whereas *active attacks* can modify the traffic, or even place new commands on the bus. *Replay attacks* are a notoriously challenging class of active attacks, where the attacker first records the bus traffic, and then selectively replays a subset of the traffic. Replay attacks bypass systems that rely on static signatures or HMACs, and generally aim to double-spend a limited resource.

The cost of bus tapping attacks is generally dominated by the cost of the equipment used to tap the bus, which increases with bus speed and complexity. For example, the flash chip that stores the computer's firmware is connected to the PCH via an SPI bus (§ A.9.1), which is simpler and much slower than the DDR bus connecting DRAM to the CPU. Consequently,

tapping the SPI bus is much cheaper than tapping the DDR bus. For this reason, systems whose security relies on a cryptographic hash of the firmware will first copy the firmware into DRAM, hash the DRAM copy of the firmware, and then execute the firmware from DRAM.

Although the speed of the DDR bus makes tapping very difficult, there are well-publicized records of successful attempts. The original Xbox console's booting process was reverse-engineered, thanks to a passive tap on the DRAM bus [85], which showed that the firmware used to boot the console was partially stored in its southbridge. The protection mechanisms of the PlayStation 3 hypervisor were subverted by an active tap on its memory bus [84] that targeted the hypervisor's page tables.

The Ascend secure processor (§ 2.10) shows that concealing the addresses of the DRAM cells accessed by a program is orders of magnitude more expensive than protecting the memory's contents. Therefore, we are interested in analyzing attacks that tap the DRAM bus, but only use the information on the address lines. These attacks use the same equipment as normal DRAM bus tapping attacks, but require a significantly more involved analysis to learn useful information. One of the difficulties of such attacks is that the memory addresses observed on the DRAM bus are generally very different from the application's memory access patterns, because of the extensive cache hierarchies in modern processors (§ A.11).

We are not aware of any successful attack based on tapping the address lines of a DRAM bus and analyzing the sequence of memory addresses.

B.4.3 Chip Attacks

The most equipment-intensive physical attacks involve removing a chip's packaging and directly interacting with its electrical circuits. These attacks generally take advantage of equipment and techniques that were originally developed to diagnose design and manufacturing defects in chips. [24] covers these techniques in depth.

The cost of chip attacks is dominated by the required equipment, although the reverse-engineering involved is also non-trivial. This cost grows very rapidly as the circuit components shrink. At the time of this writing, the latest Intel CPUs have a 14nm feature size,

which requires ion beam microscopy.

The least expensive classes of chip attacks are destructive, and only require imaging the chip's circuitry. These attacks rely on a microscope capable of capturing the necessary details in each layer, and equipment for mechanically removing each layer and exposing the layer below it to the microscope.

Imaging attacks generally target global secrets shared by all the chips in a family, such as ROM masks that store global encryption keys or secret boot code. They are also used to reverse-engineer undocumented functionality, such as debugging backdoors. E-fuses and polyfuses are particularly vulnerable to imaging attacks, because of their relatively large sizes.

Non-destructive passive chip attacks require measuring the voltages across a module at specific times, while the chip is operating. These attacks are orders of magnitude more expensive than imaging attacks, because the attacker must maintain the integrity of the chip's circuitry, and therefore cannot de-layer the chip.

The simplest active attacks on a chip create or destroy an electric connection between two components. For example, the debugging functionality in many chips is disabled by "blowing" an e-fuse. Once this e-fuse is located, an attacker can reconnect its two ends, effectively undoing the "blowing" operation. More expensive attacks involve changing voltages across a component as the chip is operating, and are typically used to reverse-engineer complex circuits.

Surprisingly, active attacks are not significantly more expensive to carry out than passive non-destructive attacks. This is because the tools used to measure the voltage across specific components are not very different from the tools that can tamper with the chip's electric circuits. Therefore, once an attacker develops a process for accessing a module without destroying the chip's circuitry, the attacker can use the same process for both passive and active attacks.

At the architectural level, we cannot address physical attacks against the CPU's chip package. Active attacks on the CPU change the computer's execution semantics, leaving us without any hardware that can be trusted to make security decisions. Passive attacks can read the private data that the CPU is processing. Therefore, many secure computing

architectures assume that the processor chip package is invulnerable to physical attacks.

Thankfully, physical attacks can be deterred by reducing the value that an attacker obtains by compromising an individual chip. As long as this value is below the cost of carrying out the physical attack, a system's designer can hope that the processor's chip package will not be targeted by the physical attacks.

Architects can reduce the value of compromising an individual system by avoiding shared secrets, such as global encryption keys. Chip designers can increase the cost of a physical attack by not storing a platform's secrets in hardware that is vulnerable to destructive attacks, such as e-fuses.

B.4.4 Power Analysis Attacks

An entirely different approach to physical attacks consists of indirectly measuring the power consumption of a computer system or its components. The attacker takes advantage of a known correlation between power consumption and the computed data, and learns some property of the data from the observed power consumption.

The earliest power analysis attacks have directly measured the processor chip's power consumption. For example, [125] describes a simple power analysis (SPA) attack that exploits the correlation between the power consumed by a smart card chip's CPU and the type of instruction it executed, and learned a DSA key that the smart card was supposed to safeguard.

While direct power analysis attacks necessitate some equipment, their costs are dominated by the complexity of the analysis required to learn the desired information from the observed power trace which, in turn, is determined by the complexity of the processor's circuitry. Today's smart cards contain special circuitry [185] and use hardened algorithms [80] designed to frustrate power analysis attacks.

Recent work demonstrated successful power analysis attacks against full-blown out-of-order Intel processors using inexpensive off-the-shelf sensor equipment. [64] extracts an RSA key from GnuPG running on a laptop using a microphone that measures its acoustic emissions. [63] and [62] extract RSA keys from power analysis-resistant implementations

using a voltage meter and a radio. All these attacks can be performed quite easily by a disgruntled data center employee.

Unfortunately, power analysis attacks can be extended to displays and human input devices, which cannot be secured in any reasonable manner. For example, [188] documented a very early attack that measures the radiation emitted by a CRT display's ion beam to reconstitute the image on a computer screen in a different room. [128] extended the attack to modern LCD displays. [211] used a directional microphone to measure the sound emitted by a keyboard and learn the password that its operator typed. [152] applied similar techniques to learn a user's input on a smartphone's on-screen keyboard, based on data from the device's accelerometer.

In general, power attacks cannot be addressed at the architectural level, as they rely on implementation details that are decided during the manufacturing process. Therefore, it is unsurprising that the secure computing architectures described in § 2 do not protect against power analysis attacks.

B.5 Privileged Software Attacks

The rest of this section points to successful exploits that execute at each of the privilege levels described in § A.3, motivating the SGX design decision to assume that all the privileged software on the computer is malicious. [168] describes all the programmable hardware inside Intel computers, and outlines the security implications of compromising the software running it.

SMM, the most privileged execution level, is only used to handle a specific kind of interrupts (§ A.12), namely *System Management Interrupts* (SMI). SMIs were initially designed exclusively for hardware use, and were only triggered by asserting a dedicated pin (SMI#) in the CPU's chip package. However, in modern systems, system software can generate an SMI by using the LAPIC's IPI mechanism. This opens up the avenue for SMM-based software exploits.

The SMM handler is stored in *System Management RAM* (SMRAM) which, in theory, is not accessible when the processor isn't running in SMM. However, its protection mecha-

nisms were bypassed multiple times [47, 169, 198, 116], and SMM-based rootkits [195, 53] have been demonstrated. Compromising the SMM grants an attacker access to all the software on the computer, as SMM is the most privileged execution mode.

Xen [210] is a very popular representative of the family of hypervisors that run in VMX root mode and use hardware virtualization. At 150,000 lines of code [12], Xen's codebase is relatively small, especially when compared to a kernel. However, Xen still has had over 40 security vulnerabilities patched in **each** of the last three years (2012-2014) [10].

[140] proposes using a very small hypervisor together with Intel TXT's dynamic root of trust for measurement (DRTM) to implement trusted execution. [189] argues that a dynamic root of trust mechanism, like Intel TXT, is necessary to ensure a hypervisor's integrity. Unfortunately, the TXT design requires an implementation complex enough that exploitable security vulnerabilities have crept in [200, 199]. Furthermore, any SMM attack can be used to compromise TXT [197].

The monolithic kernel design leads to many opportunities for security vulnerabilities in kernel code. Linux is by far the most popular kernel for IaaS cloud environments. Linux has *17 million* lines of code [17], and has had over 100 security vulnerabilities patched in **each** of the last three years (2012-2014) [8, 35].

B.6 Software Attacks on Peripherals

Threat models for secure architectures generally only consider software attacks that directly target other components in the software stack running on the CPU. This assumption results in security arguments with the very desirable property of not depending on implementation details, such as the structure of the motherboard hosting the processor chip.

The threat models mentioned above must classify attacks from other motherboard components as physical attacks. Unfortunately, these models would mis-classify all the attacks described in this section, which can be carried out solely by executing software on the victim processor. The incorrect classification matters in cloud computing scenarios, where physical attacks are significantly more expensive than software attacks.

B.6.1 PCI Express Attacks

The PCIe bus (§ A.9.1) allows any device connected to the bus to perform *Direct Memory Access* (DMA), reading from and writing to the computer's DRAM without the involvement of a CPU core. Each device is assigned a range of DRAM addresses via a standard PCI configuration mechanism, but can perform DMA on DRAM addresses outside of that range.

Without any additional protection mechanism, an attacker who compromises system software can take advantage of programmable devices to access any DRAM region, yielding capabilities that were traditionally associated with a DRAM bus tap. For example, an early implementation of Intel TXT [74] was compromised by programming a PCIe NIC to read TXT-reserved DRAM via DMA transfers [199]. Recent versions have addressed this attack by adding extra security checks in the DMA bus arbiter. § 2.5 provides a more detailed description of Intel TXT.

B.6.2 DRAM Attacks

The rowhammer DRAM bit-flipping attack [121, 171, 76] is an example of a different class of software attacks that exploit design defects in the computer's hardware. Rowhammer took advantage of the fact that some mobile DRAM chips (§ A.9.1) refreshed the DRAM's contents slowly enough that repeatedly changing the contents of a memory cell could impact the charge stored in a neighboring cell, which resulted in changing the bit value obtained from reading the cell. By carefully targeting specific memory addresses, the attackers caused bit flips in the page tables used by the CPU's address translation (§ A.5) mechanism, and in other data structures used to make security decisions.

The defect exploited by the rowhammer attack most likely stems from an incorrect design assumption. The DRAM engineers probably only thought of non-malicious software and assumed that an individual DRAM cell cannot be accessed too often, as repeated accesses to the same memory address would be absorbed by the CPU's caches (§ A.11). However, malicious software can take advantage of the `CLFLUSH` instruction, which flushes the cache line that contains a given DRAM address. `CLFLUSH` is intended as a method for applications to extract more performance out of the cache hierarchy, and is therefore available to software

running at all privilege levels. Rowhammer exploited the combination of CLEFLUSH's availability and the DRAM engineers' invalid assumptions, to obtain capabilities that are normally associated with an active DRAM bus attack.

B.6.3 The Performance Monitoring Side Channel

Intel's *Software Development Manual* (SDM) [104] and *Optimization Reference Manual* [99] describe a vast array of performance monitoring events exposed by recent Intel processors, such as branch mispredictions (§ A.10). The SDM also describes digital temperature sensors embedded in each CPU core, whose readings are exposed using Model-Specific Registers (MSRs) (§ A.4) that can be read by system software.

An attacker who compromises a computer's system software and gains access to the performance monitoring events or the temperature sensors can obtain the information needed to carry out a power analysis attack, which normally requires physical access to the victim computer and specialized equipment.

B.6.4 Attacks on the Boot Firmware and Intel ME

Virtually all motherboards store the firmware used to boot the computer in a flash memory chip (§ A.9.1) that can be written by system software. This implementation strategy provides an inexpensive avenue for deploying firmware bug fixes. At the same time, an attack that compromises the system software can subvert the firmware update mechanism to inject malicious code into the firmware. The malicious code can be used to carry out a cold boot attack, which is typically considered a physical attack. Furthermore, malicious firmware can run code at the highest software privilege level, System Management Mode (SMM, § A.3). Last, malicious firmware can modify the system software as it is loaded during the boot process. These avenues give the attacker capabilities that have traditionally been associated with DRAM bus tapping attacks.

The Intel Management Engine (ME) [166] loads its firmware from the same flash memory chip as the main computer, which opens up the possibility of compromising its firmware. Due to its vast management capabilities (§ A.9.2), a compromised ME would

leak most of the powers that come with installing active probes on the DRAM bus, the PCI bus, and the System Management bus (SMBus), as well as power consumption meters. Thanks to its direct access to the motherboard's Ethernet PHY, the probe would be able to communicate with the attacker while the computer is in the Soft-Off state, also known as S5, where the computer is mostly powered off, but is still connected to a power source. The ME has significantly less computational power than probe equipment, however, as it uses low-power embedded components, such as a 200-400MHz execution core, and about 600KB of internal RAM.

The computer and ME firmware are protected by a few security measures. The first line of defense is a security check in the firmware's update service, which only accepts firmware updates that have been digitally signed by a manufacturer key that is hard-coded in the firmware. This protection can be circumvented with relative ease by foregoing the firmware's update services, and instead accessing the flash memory chip directly, via the PCH's SPI bus controller.

The deeper, more powerful, lines of defense against firmware attacks are rooted in the CPU and ME's hardware. The bootloader in the ME's ROM will only load flash firmware that contains a correct signature generated by a specific Intel RSA key. The ME's boot ROM contains the SHA-256 cryptographic hash of the RSA public key, and uses it to validate the full Intel public key stored in the signature. Similarly, the microcode bootstrap process in recent CPUs will only execute firmware in an Authenticated Code Module (ACM, § A.13.2) signed by an Intel key whose SHA-256 hash is hard-coded in the microcode ROM.

However, both the computer firmware security checks [201, 58] and the ME security checks [184] have been subverted in the past. While the approaches described above are theoretically sound, the intricate details and complex interactions in Intel-based systems make it very likely that security vulnerabilities will creep into implementations. Further proving this point, a security analysis [191] found that early versions of Intel's Active Management Technology (AMT), the flagship ME application, contained an assortment of security issues that allowed an attacker to completely take over a computer whose ME firmware contained the AMT application.

B.6.5 Accounting for Software Attacks on Peripherals

The attacks described in this section show that a system whose threat model assumes no software attacks must be designed with an understanding of all the system's buses, and the programmable devices that may be attached to them. The system's security analysis must argue that the devices will not be used in physical-like attacks. The argument will rely on barriers that prevent untrusted software running on the CPU from communicating with other programmable devices, and on barriers that prevent compromised programmable devices from tampering with sensitive buses or DRAM.

Unfortunately, the ME, PCH and DMI are Intel-proprietary and largely undocumented, so we cannot assess the security of the measures set in place to protect the ME from being compromised, and we cannot reason about the impact of a compromised ME that runs malicious software.

B.7 Address Translation Attacks

§ B.5 argues that today's system software is virtually guaranteed to have security vulnerabilities. This suggests that a cautious secure architecture should avoid having the system software in the TCB.

However, removing the system software from the TCB requires the architecture to provide a method for isolating sensitive application code from the untrusted system software. This is typically accomplished by designing a mechanism for loading application code in isolated containers whose contents can be certified via software attestation (§ B.3). One of the more difficult problems these designs face is that application software relies on the memory management services provided by the system software, which is now untrusted.

Intel's SGX [143, 15], leaves the system software in charge of setting up the page tables (§ A.5) used by address translation, inspired by Bastion [33], but instantiates access checks that prevent the system software from directly accessing the isolated container's memory.

This section discusses some attacks that become relevant when the application software does not trust the system software, which is in charge of the page tables. Understanding these attacks is a prerequisite to reasoning about the security properties of architectures

with this threat model. For example, many of the mechanisms in SGX target a subset of the attacks described here.

B.7.1 Passive Attacks

System software uses the CPU's address translation feature (§ A.5) to implement page swapping, where infrequently used memory pages are evicted from DRAM to a slower storage medium. Page swapping relies the accessed (A) and dirty (D) page table entry attributes (§ A.5.3) to identify the DRAM pages to be evicted, and on a page fault handler (§ A.8.2) to bring evicted pages back into DRAM when they are accessed.

Unfortunately, the features that support efficient page swapping turn into a security liability, when the system software managing the page tables is not trusted by the application software using the page tables. The system software can be prevented from reading the application's memory directly by placing the application in an isolated container. However, potentially malicious system software can still infer partial information about the application's memory access patterns, by observing the application's page faults and page table attributes.

We consider this class of attacks to be passive attacks that exploit the CPU's address translation feature. It may seem that the page-level memory access patterns provided by these attacks are not very useful. However, [204] describes how this attack can be carried out against Intel's SGX, and implements the attack in a few practical settings. In one scenario, which is particularly concerning for medical image processing, the outline of a JPEG image is inferred while the image is decompressed inside a container protected by SGX's isolation guarantees.

B.7.2 Straightforward Active Attacks

We define active address translation attacks to be the class of attacks where malicious system software modifies the page tables used by an application in a way that breaks the virtual memory abstraction (§ A.5). Memory mapping attacks do not include scenarios where the system software breaks the memory abstraction by directly writing to the application's

memory pages.

We begin with an example of a straight-forward active attack. In this example, the application inside a protected container performs a security check to decide whether to disclose some sensitive information. Depending on the security check's outcome, the enclave code either calls a `errorOut` procedure, or a `disclose` procedure. The simplest version of the attack assumes that each procedure's code starts at a page boundary, and takes up less than a page. These assumptions are relaxed in more complex versions of the attack.

In the most straightforward setting, the malicious system software directly modifies the page tables of the application inside the container, as shown in Figure B-23, so the virtual address intended to store the `errorOut` procedure is actually mapped to a DRAM page that contains the `disclose` procedure. Without any security measures in place, when the application's code jumps to the virtual address of the `errorOut` procedure, the CPU will execute the code of the `disclose` procedure instead.

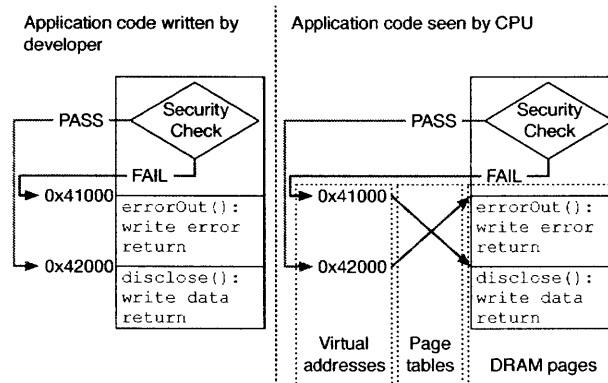


Figure B-23: An example of an active memory mapping attack. The application's author intends to perform a security check, and only call the procedure that discloses the sensitive information if the check passes. Malicious system software maps the virtual address of the procedure that is called when the check fails, to a DRAM page that contains the disclosing procedure.

B.7.3 Active Attacks Using Page Swapping

The most obvious active attacks on memory mapping can be defeated by tracking the correct virtual address for each DRAM page that belongs to a protected container. However, a naive protection measure based on address tracking can be defeated by a more subtle active attack

that relies on the architectural support for page swapping. Figure B-24 illustrates an attack that does not modify the application’s page tables, but produces the same corrupted CPU view of the application as the straight-forward attack described above.

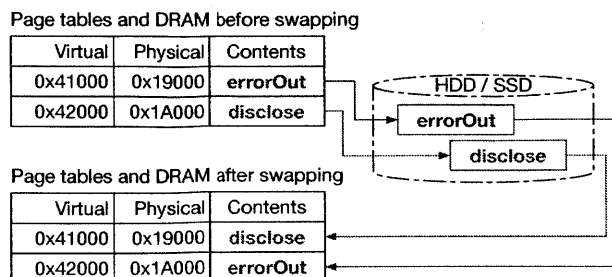


Figure B-24: An active memory mapping attack where the system software does not modify the page tables. Instead, two pages are evicted from DRAM to a slower storage medium. The malicious system software swaps the two pages’ contents then brings them back into DRAM, building the same incorrect page mapping as the direct attack shown in Figure B-23. This attack defeats protection measures that rely on tracking the virtual and disk addresses for DRAM pages.

In the swapping attack, malicious system software evicts the pages that contain the `errorOut` and `disclose` procedures from DRAM to a slower medium, such as a hard disk. The system software exchanges the hard disk bytes storing the two pages, and then brings the two pages back into DRAM. Remarkably, all the steps taken by this attack are indistinguishable from legitimate page swapping activity, with the exception of the I/O operations that exchange the disk bytes storing evicted pages.

The subtle attack described in this section can be defeated by cryptographically binding the contents of each page that is evicted from DRAM to the virtual address to which the page should be mapped. The cryptographic primitive (§ B.1) used to perform the binding must obviously guarantee integrity. Furthermore, it must also guarantee freshness, in order to foil replay attacks where the system software “undoes” an application’s writes by evicting one of its DRAM pages to disk and bringing in an older version of the same page.

B.7.4 Active Attacks Based on TLBs

Today’s multi-core architectures can be subjected to an even more subtle active attack, illustrated in Figure B-25, which can bypass any protection measures that solely focus on

the integrity of the page tables.

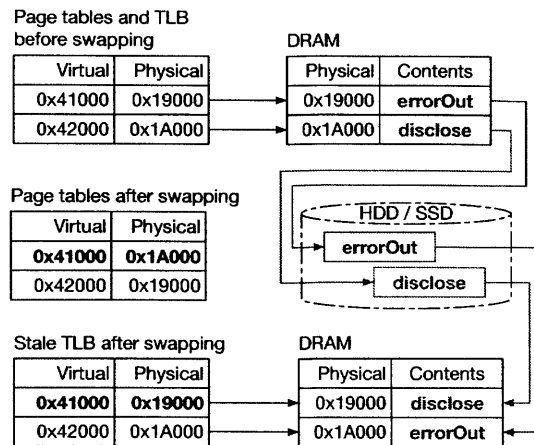


Figure B-25: An active memory mapping attack where the system software does not invalidate a core's TLBs when it evicts two pages from DRAM and exchanges their locations when reading them back in. The page tables are updated correctly, but the core with stale TLB entries has the same incorrect view of the protected container's code as in Figure B-23.

For performance reasons, each execution core caches address translation results in its own translation look-aside buffer (TLB, § A.11.5). For simplicity, the TLBs are not covered by the cache coherence protocol that synchronizes data caches across cores. Instead, the system software is responsible for invalidating TLB entries across all the cores when it modifies the page tables.

Malicious system software can take advantage of the design decisions explained above by carrying out the following attack. While the same software used in the previous examples is executing on a core, the system software executes on a different core and evicts the `errorOut` and `disclose` pages from DRAM. As in the previous attack, the system software loads the `disclose` code in the DRAM page that previously held `errorOut`. In this attack, however, the system software also updates the page tables.

The core where the system software executed sees the code that the application developer intended. Therefore, the attack will pass any security checks that rely upon cryptographic associations between page contents and page table data, as long as the checks are performed by the core used to load pages back into DRAM. However, the core that executes the protected container's code still uses the old page table data, because the system software did not invalidate its TLB entries. Assuming the TLBs are not subjected to any additional security

checks, this attack causes the same private information leak as the previous examples.

In order to avoid the attack described in this section, the trusted software or hardware that implements protected containers must also ensure that the system software invalidates the relevant TLB entries on all the cores when it evicts a page from a protected container to DRAM.

B.8 Cache Timing Attacks

Cache timing attacks [21] are a powerful class of software attacks that can be mounted entirely by application code running at ring 3 (§ A.3). Cache timing attacks do not learn information by reading the victim's memory, so they bypass the address translation-based isolation measures (§ A.5) implemented in today's kernels and hypervisors.

B.8.1 Theory

Cache timing attacks exploit the unfortunate dependency between the location of a memory access and the time it takes to perform the access. A cache miss requires at least one memory access to the next level cache, and might require a second memory access if a write-back occurs. On the Intel architecture, the latency between a cache hit and a miss can be easily measured by the `RDTSC` and `RDTSCP` instructions (§ A.4), which read a high-resolution time-stamp counter. These instructions have been designed for benchmarking and optimizing software, so they are available to ring 3 software.

The fundamental tool of a cache timing attack is an attacker process that measures the latency of accesses to carefully designated memory locations in its own address space. The memory locations are chosen so that they map to the same cache lines as those of some interesting memory locations in a victim process, in a cache that is shared between the attacker and the victim. This requires in-depth knowledge of the shared cache's organization (§ A.11.2).

Armed with the knowledge of the cache's organization, the attacker process sets up the attack by accessing its own memory in such a way that it fills up all the cache sets that would hold the victim's interesting memory locations. After the targeted cache sets are

full, the attacker allows the victim process to execute. When the victim process accesses an interesting memory location in its own address space, the shared cache must evict one of the cache lines holding the attacker's memory locations.

As the victim is executing, the attacker process repeatedly times accesses to its own memory locations. When the access times indicate that a location was evicted from the cache, the attacker can conclude that the victim accessed an interesting memory location in its own cache. Over time, the attacker collects the results of many measurements and learns a subset of the victim's memory access pattern. If the victim processes sensitive information using data-dependent memory fetches, the attacker may be able to deduce the sensitive information from the learned memory access pattern.

B.8.2 Practical Considerations

Cache timing attacks require control over a software process that shares a cache memory with the victim process. Therefore, a cache timing attack that targets the L2 cache would have to rely on the system software to schedule a software thread on a logical processor in the same core as the target software, whereas an attack on the L3 cache can be performed using any logical processor on the same CPU. The latter attack relies on the fact that the L3 cache is inclusive, which greatly simplifies the processor's cache coherence implementation (§ A.11.3).

The cache sharing requirement implies that L3 cache attacks are feasible in an IaaS environment, whereas L2 cache attacks become a significant concern when running sensitive software on a user's desktop.

Out-of-order execution (§ A.10) can introduce noise in cache timing attacks. First, memory accesses may not be performed in program order, which can impact the lines selected by the cache eviction algorithms. Second, out-of-order execution may result in cache fills that do not correspond to executed instructions. For example, a load that follows a faulting instruction may be scheduled and executed before the fault is detected.

Cache timing attacks must account for speculative execution, as mispredicted memory accesses can still cause cache fills. Therefore, the attacker may observe cache fills that don't

correspond to instructions that were actually executed by the victim software. Memory prefetching adds further noise to cache timing attacks, as the attacker may observe cache fills that don't correspond to instructions in the victim code, even when accounting for speculative execution.

B.8.3 Known Cache Timing Attacks

Despite these difficulties, cache timing attacks are known to retrieve cryptographic keys used by AES [150, 27], RSA [30], Diffie-Hellman [126], and elliptic-curve cryptography [29].

Early attacks required access to the victim's CPU core, but more sophisticated recent attacks [205, 135] are able to use the L3 cache, which is shared by all the cores on a CPU die. L3-based attacks can be particularly devastating in cloud computing scenarios, where running software on the same computer as a victim application only requires modest statistical analysis skills and a small amount of money [161]. Furthermore, cache timing attacks were recently demonstrated using JavaScript code in a page visited by a Web browser [149].

Given this pattern of vulnerabilities, ignoring cache timing attacks is dangerously similar to ignoring the string of demonstrated attacks which led to the deprecation of SHA-1 [9, 6, 3].

B.8.4 Defending against Cache Timing Attacks

Fortunately, invalidating any of the preconditions for cache timing attacks is sufficient for defending against them. The easiest precondition to focus on is that the attacker must have access to memory locations that map to the same sets in a cache as the victim's memory. This assumption can be invalidated by the judicious use of a cache partitioning scheme.

Performance concerns aside, the main difficulty associated with cache partitioning schemes is that they must be implemented by a trusted party. When the system software is trusted, it can (for example) use the principles behind page coloring [183, 119] to partition the caches [133] between mutually distrusting parties. This comes down to setting up the page tables in such a way that no two mutually distrusting software module are stored in physical pages that map to the same sets in any cache memory. However, if the system

software is not trusted, the cache partitioning scheme must be implemented in hardware.

The other interesting precondition is that the victim must access its memory in a data-dependent fashion that allows the attacker to infer private information from the observed memory access pattern. It becomes tempting to think that cache timing attacks can be prevented by eliminating data-dependent memory accesses from all the code handling sensitive data.

However, removing data-dependent memory accesses is difficult to accomplish in practice because instruction fetches must also be taken into consideration. [117] gives an idea of the level of effort required to remove data-dependent accesses from AES, which is a relatively simple data processing algorithm. At the time of this writing, we are not aware of any approach that scales to large pieces of software.

While the focus of this section is cache timing attacks, we would like to point out that any shared resource can lead to information leakage. A worrying example is hyper-threading (§ A.9.4), where each CPU core is represented as two logical processors, and the threads executing on these two processors share execution units. An attacker who can run a process on a logical processor sharing a core with a victim process can use `RDTSCP` [156] to learn which execution units are in use, and infer what instructions are executed by the victim process.

Appendix C

Intel SGX Explained

Intel's Software Guard Extensions (SGX) is a set of extensions to the Intel architecture that aims to provide integrity and privacy guarantees to security-sensitive computation performed on a computer where all the privileged software (kernel, hypervisor, etc) is potentially malicious.

This paper analyzes Intel SGX, based on the 3 papers [143, 15, 82] that introduced it, on the Intel Software Developer's Manual [104] (which supersedes the SGX manuals [98, 102]), on an ISCA 2015 tutorial [106], and on two patents [142, 112]. We use the papers, reference manuals, and tutorial as primary data sources, and only draw on the patents to fill in missing information.

This section's contributions are a detailed and structured presentation of the publicly available information on SGX, a series of intelligent guesses about some important but undocumented aspects of SGX, and an analysis of SGX's security properties.

C.1 Trusted Hardware

Secure remote computation (Figure C-1) is the problem of executing software on a remote computer **owned and maintained by an untrusted party**, with some integrity and privacy guarantees. In the general setting, secure remote computation is an unsolved problem. Fully Homomorphic Encryption [65] solves the problem for a limited family of computations, but has an impractical performance overhead [144].

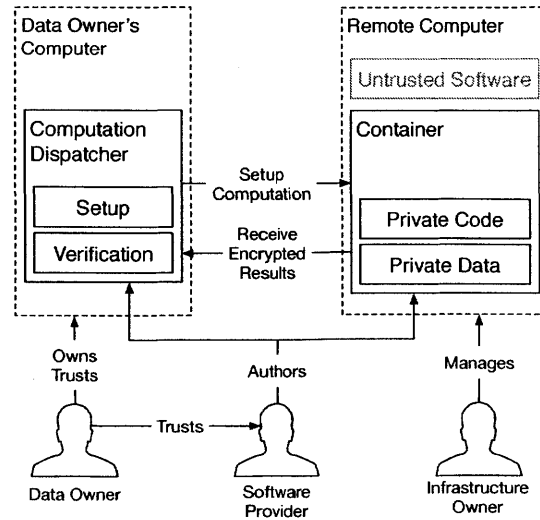


Figure C-1: Secure remote computation. A user relies on a remote computer, owned by an untrusted party, to perform some computation on her data. The user has some assurance of the computation's integrity and privacy.

Intel's Software Guard Extensions (SGX) is the latest iteration in a long line of trusted computing (Figure C-2) designs, which aim to solve the secure remote computation problem by leveraging trusted hardware in the remote computer. The trusted hardware establishes a secure container, and the remote computation service user uploads the desired computation and data into the secure container. The trusted hardware protects the data's privacy and integrity while the computation is being performed on it.

SGX relies on *software attestation*, like its predecessors, the TPM [75] and TXT [74]. Attestation (Figure C-3) proves to a user that she is communicating with a specific piece of software running in a secure container hosted by the trusted hardware. The proof is a cryptographic signature that certifies the hash of the secure container's contents. It follows that the remote computer's owner can load any software in a secure container, but the remote computation service user will refuse to load her data into a secure container whose contents' hash does not match the expected value.

The remote computation service user verifies the *attestation key* used to produce the signature against an *endorsement certificate* created by the trusted hardware's manufacturer. The certificate states that the attestation key is only known to the trusted hardware, and only used for the purpose of attestation.

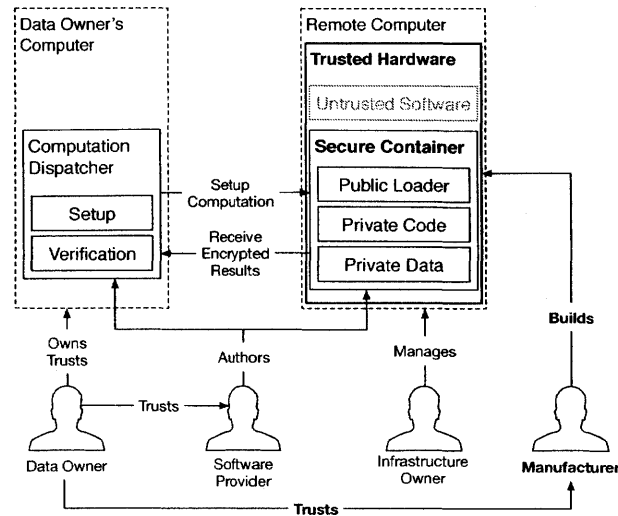


Figure C-2: Trusted computing. The user trusts the manufacturer of a piece of hardware in the remote computer, and entrusts her data to a secure container hosted by the secure hardware.

SGX stands out from its predecessors by the amount of code covered by the attestation, which is in the Trusted Computing Base (TCB) for the system using hardware protection. The attestations produced by the original TPM design covered all the software running on a computer, and TXT attestations covered the code inside a VMX [187] virtual machine. In SGX, an *enclave* (secure container) only contains the private data in a computation, and the code that operates on it.

For example, a cloud service that performs image processing on confidential medical images could be implemented by having users upload encrypted images. The users would send the encryption keys to software running inside an enclave. The enclave would contain the code for decrypting images, the image processing algorithm, and the code for encrypting the results. The code that receives the uploaded encrypted images and stores them would be left outside the enclave.

An SGX-enabled processor protects the integrity and privacy of the computation inside an enclave by isolating the enclave's code and data from the outside environment, including the operating system and hypervisor, and hardware devices attached to the system bus. At the same time, the SGX model remains compatible with the traditional software layering in the Intel architecture, where the OS kernel and hypervisor manage the computer's resources.

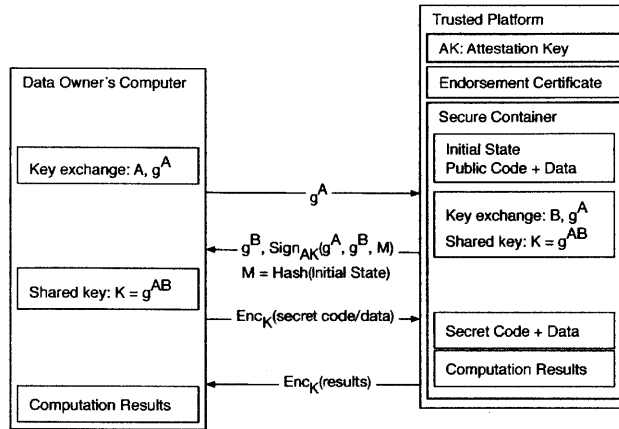


Figure C-3: Software attestation proves to a remote computer that it is communicating with a specific secure container hosted by a trusted platform. The proof is an attestation signature produced by the platform’s secret attestation key. The signature covers the container’s initial state, a challenge nonce produced by the remote computer, and a message produced by the container.

This work discusses the original version of SGX, also referred to as SGX 1. While SGX 2 brings very useful improvements for enclave authors, it is a small incremental improvement, from a design and implementation standpoint. After understanding the principles behind SGX 1 and its security properties, the reader should be well equipped to face Intel’s reference documentation and learn about the changes brought by SGX 2.

C.2 SGX Lightning Tour

SGX sets aside a memory region, called the *Processor Reserved Memory* (PRM, § C.4.1). The CPU protects the PRM from all non-enclave memory accesses, including kernel, hypervisor and SMM (§ A.3) accesses, and DMA accesses (§ A.9.1) from peripherals.

The PRM holds the *Enclave Page Cache* (EPC, § C.4.1), which consists of 4 KB pages that store enclave code and data. The system software, which is untrusted, is in charge of assigning EPC pages to enclaves. The CPU tracks each EPC page’s state in the *Enclave Page Cache Metadata* (EPCM, § C.4.1), to ensure that each EPC page belongs to exactly one enclave.

The initial code and data in an enclave is loaded by untrusted system software. During

the loading stage (§ C.4.3), the system software asks the CPU to copy data from unprotected memory (outside PRM) into EPC pages, and assigns the pages to the enclave being setup (§ C.4.1). It follows that the initial enclave state is known to the system software.

After all the enclave's pages are loaded into EPC, the system software asks the CPU to mark the enclave as initialized (§ C.4.3), at which point application software can run the code inside the enclave. After an enclave is initialized, the loading method described above is disabled.

While an enclave is loaded, its contents is cryptographically hashed by the CPU. When the enclave is initialized, the hash is finalized, and becomes the enclave's *measurement hash* (§ C.4.6).

A remote party can undergo a *software attestation* process (§ C.4.8) to convince itself that it is communicating with an enclave that has a specific measurement hash, and is running in a secure environment.

Execution flow can only enter an enclave via special CPU instructions (§ C.4.4), which are similar to the mechanism for switching from user mode to kernel mode. Enclave execution always happens in protected mode, at ring 3, and uses the address translation set up by the OS kernel and hypervisor.

To avoid leaking private data, a CPU that is executing enclave code does not directly service an interrupt, fault (e.g., a page fault) or VM exit. Instead, the CPU first performs an Asynchronous Enclave Exit (§ C.4.4) to switch from enclave code to ring 3 code, and then services the interrupt, fault, or VM exit. The CPU performs an AEX by saving the CPU state into a predefined area inside the enclave and transfers control to a pre-specified instruction outside the enclave, replacing CPU registers with synthetic values.

The allocation of EPC pages to enclaves is delegated to the OS kernel (or hypervisor). The OS communicates its allocation decisions to the SGX implementation via special ring 0 CPU instructions (§ C.4.3). The OS can also evict EPC pages into untrusted DRAM and later load them back, using dedicated CPU instructions. SGX uses cryptographic protections to assure the privacy, integrity and freshness of the evicted EPC pages while they are stored in untrusted memory.

C.3 Outline and Troubling Findings

Reasoning about the security properties of Intel's SGX requires a significant amount of background information that is currently scattered across many sources. For this reason, a significant portion of this work is dedicated to summarizing this prerequisite knowledge.

Section A summarizes the relevant subset of the Intel architecture and the micro-architectural properties of recent Intel processors. Section B outlines the security landscape around trusted hardware system, including cryptographic tools and relevant attack classes. Last, section 2 briefly describes the trusted hardware systems that make up the context in which SGX was created.

After having reviewed the background information, section C.4 provides a (sometimes painstakingly) detailed description of SGX's programming model, mostly based on Intel's Software Development Manual.

Section C.5 analyzes other public sources of information, such as Intel's SGX-related patents, to fill in some of the missing details in the SGX description. The section culminates in a detailed review of SGX's security properties that draws on information presented in the rest of the paper. This review outlines some troubling gaps in SGX's security guarantees, as well as some areas where no conclusions can be drawn without additional information from Intel.

That being said, perhaps the most troubling finding in our security analysis is that Intel added a launch control feature to SGX that forces each computer's owner to gain approval from a third party (which is currently Intel) for any enclave that the owner wishes to use on the computer. § C.4.9 explains that the only publicly documented intended use for this launch control feature is a licensing mechanism that requires software developers to enter a (yet unspecified) business agreement with Intel to be able to author software that takes advantage of SGX's protections. All the official documentation carefully sidesteps this issue, and has a minimal amount of hints that lead to the Intel's patents on SGX. Only these patents disclose the existence of licensing plans.

The licensing issue might not bear much relevance right now, because our security analysis reveals that the limitations in SGX's guarantees mean that a security-conscious

software developer cannot in good conscience rely on SGX for secure remote computation. At the same time, should SGX ever develop better security properties, the licensing scheme described above becomes a major problem, given Intel's near-monopoly market share of desktop and server CPUs. Specifically, the licensing limitations effectively give Intel the power to choose winners and losers in industries that rely on cloud computing.

C.4 SGX Programming Model

The central concept of SGX is the *enclave*, a protected environment that contains the code and data pertaining to a security-sensitive computation.

SGX-enabled processors provide trusted computing by isolating each enclave's environment from the untrusted software outside the enclave, and by implementing a software attestation scheme that allows a remote party to authenticate the software running inside an enclave. SGX's isolation mechanisms are intended to protect the privacy and integrity of the computation performed inside an enclave from attacks coming from malicious software executing on the same computer, as well as from a limited set of physical attacks.

This section summarizes the SGX concepts that make up a mental model which is sufficient for programmers to author SGX enclaves and to add SGX support to existing system software. All the information in this section is backed up by Intel's Software Developer Manual (SDM). The following section builds on the concepts introduced here to fill in some of the missing pieces in the manual, and analyzes some of SGX's security properties.

C.4.1 SGX Physical Memory Organization

The enclaves' code and data is stored in *Processor Reserved Memory* (PRM), which is a subset of DRAM that cannot be directly accessed by other software, including system software and SMM code. The CPU's integrated memory controllers (§ A.9.3) also reject DMA transfers targeting the PRM, thus protecting it from access by other peripherals.

The PRM is a continuous range of memory whose bounds are configured using a base and a mask register with the same semantics as a variable memory type range (§ A.11.4).

Therefore, the PRM's size must be an integer power of two, and its start address must be aligned to the same power of two. Due to these restrictions, checking if an address belongs to the PRM can be done very cheaply in hardware, using the circuit outlined in § A.11.4.

The SDM does not describe the PRM and the PRM range registers (PRMRR). These concepts are documented in the SGX manuals [98, 102] and in one of the SGX papers [143]. Therefore, the PRM is a micro-architectural detail that might change in future implementations of SGX. Our security analysis of SGX relies on implementation details surrounding the PRM, and will have to be re-evaluated for SGX future implementations.

The Enclave Page Cache (EPC)

The contents of enclaves and the associated data structures are stored in the *Enclave Page Cache* (EPC), which is a subset of the PRM, as shown in Figure C-4.

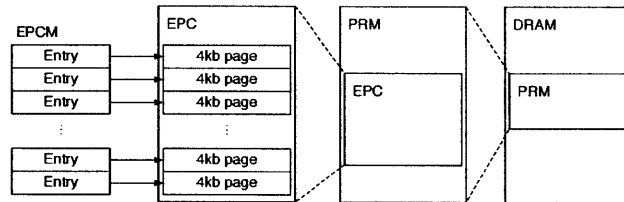


Figure C-4: Enclave data is stored into the EPC, which is a subset of the PRM. The PRM is a contiguous range of DRAM that cannot be accessed by system software or peripherals.

The SGX design supports having multiple enclaves on a system at the same time, which is a necessity in multi-process environments. This is achieved by having the EPC split into 4 KB pages that can be assigned to different enclaves. The EPC uses the same page size as the architecture's address translation feature (§ A.5). This is not a coincidence, as future sections will reveal that the SGX implementation is tightly coupled with the address translation implementation.

The EPC is managed by the same system software that manages the rest of the computer's physical memory. The system software, which can be a hypervisor or an OS kernel, uses SGX instructions to allocate unused pages to enclaves, and to free previously allocated EPC pages. The system software is expected to expose enclave creation and management services to application software.

Non-enclave software cannot directly access the EPC, as it is contained in the PRM. This restriction plays a key role in SGX's enclave isolation guarantees, but creates an obstacle when the system software needs to load the initial code and data into a newly created enclave. The SGX design solves this problem by having the instructions that allocate an EPC page to an enclave also initialize the page. Most EPC pages are initialized by copying data from a non-PRM memory page.

The Enclave Page Cache Map (EPCM)

The SGX design expects the system software to allocate the EPC pages to enclaves. However, as the system software is not trusted, SGX processors check the correctness of the system software's allocation decisions, and refuse to perform any action that would compromise SGX's security guarantees. For example, if the system software attempts to allocate the same EPC page to two enclaves, the SGX instruction used to perform the allocation will fail.

In order to perform its security checks, SGX records some information about the system software's allocation decisions for each EPC page in the *Enclave Page Cache Map* (EPCM). The EPCM is an array with one entry per EPC page, so computing the address of a page's EPCM entry only requires a bitwise shift operation and an addition.

The EPCM's contents is only used by SGX's security checks. Under normal operation, the EPCM does not generate any software-visible behavior, and enclave authors and system software developers can mostly ignore it. Therefore, the SDM only describes the EPCM at a very high level, listing the information contained within and noting that the EPCM is "trusted memory". The SDM does not disclose the storage medium or memory layout used by the EPCM.

The EPCM uses the information in Table C.1 to track the ownership of each EPC page. We defer a full discussion of the EPCM to a later section, because its contents is intimately coupled with all of SGX's features, which will be described over the next few sections.

The SGX instructions that allocate an EPC page set the VALID bit of the corresponding EPCM entry to 1, and refuse to operate on EPC pages whose VALID bit is already set.

The instruction used to allocate an EPC page also determines the page's intended usage, which is recorded in the *page type* (PT) field of the corresponding EPCM entry. The pages

Field	Bits	Description
VALID	1	0 for un-allocated EPC pages
PT	8	page type
ENCLAVESECS		identifies the enclave owning the page

Table C.1: The fields in an EPCM entry that track the ownership of pages.

that store an enclave’s code and data are considered to have a *regular* type (PT_REG in the SDM). The pages dedicated to the storage of SGX’s supporting data structures are tagged with special types. For example, the PT_SECS type identifies pages that hold SGX Enclave Control Structures, which will be described in the following section. The other EPC page types will be described in future sections.

Last, a page’s EPCM entry also identifies the enclave that owns the EPC page. This information is used by the mechanisms that enforce SGX’s isolation guarantees to prevent an enclave from accessing another enclave’s private information. As the EPCM identifies a single owning enclave for each EPC page, it is impossible for enclaves to communicate via shared memory using EPC pages. Fortunately, enclaves can share untrusted non-EPC memory, as will be discussed in § C.4.2.

The SGX Enclave Control Structure (SECS)

SGX stores per-enclave metadata in a *SGX Enclave Control Structure* (SECS) associated with each enclave. Each SECS is stored in a dedicated EPC page with the page type PT_SECS. These pages are not intended to be mapped into any enclave’s address space, and are exclusively used by the CPU’s SGX implementation.

An enclave’s identity is almost synonymous to its SECS. The first step in bringing an enclave to life allocates an EPC page to serve as the enclave’s SECS, and the last step in destroying an enclave deallocates the page holding its SECS. The EPCM entry field identifying the enclave that owns an EPC page points to the enclave’s SECS. The system software uses the virtual address of an enclave’s SECS to identify the enclave when invoking SGX instructions.

All SGX instructions take virtual addresses as their inputs. Given that SGX instructions use SECS addresses to identify enclaves, the system software must create entries in its page

tables pointing to the SECS of the enclaves it manages. However, the system software cannot access any SECS page, as these pages are stored in the PRM. SECS pages are not intended to be mapped inside their enclaves' virtual address spaces, and SGX-enabled processors explicitly prevent enclave code from accessing SECS pages.

This seemingly arbitrary limitation is in place so that the SGX implementation can store sensitive information in the SECS, and be able to assume that no potentially malicious software will access that information. For example, the SDM states that each enclave's measurement is stored in its SECS. If software would be able to modify an enclave's measurement, SGX's software attestation scheme would provide no security assurances.

The SECS is strongly coupled with many of SGX's features. Therefore, the pieces of information that make up the SECS will be gradually introduced as the different aspects of SGX are described.

C.4.2 The Memory Layout of an SGX Enclave

SGX was designed to minimize the effort required to convert application code to take advantage of enclaves. History suggests this is a wise decision, as a large factor in the continued dominance of the Intel architecture is its ability to maintain backward compatibility. To this end, SGX enclaves were designed to be conceptually similar to the leading software modularization construct, dynamically loaded libraries, which are packaged as `.so` files on Unix, and `.dll` files on Windows.

For simplicity, we describe the interaction between enclaves and non-enclave software assuming that each enclave is used by exactly one application process, which we shall refer to as the enclave's *host process*. We do note, however, that the SGX design does not explicitly prohibit multiple application processes from sharing an enclave.

The Enclave Linear Address Range (ELRANGE)

Each enclave designates an area in its virtual address space, called the *enclave linear address range* (ELRANGE), which is used to map the code and the sensitive data stored in the enclave's EPC pages. The virtual address space outside ELRANGE is mapped to access

non-EPC memory via the same virtual addresses as the enclave’s host process, as shown in Figure C-5.

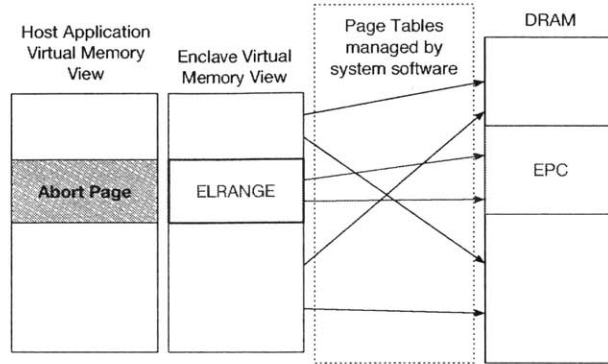


Figure C-5: An enclave’s EPC pages are accessed using a dedicated region in the enclave’s virtual address space, called ELRANGE. The rest of the virtual address space is used to access the memory of the host process. The memory mappings are established using the page tables managed by system software.

The SGX design guarantees that the enclave’s memory accesses inside ELRANGE obey the virtual memory abstraction (§ A.5.1), while memory accesses outside ELRANGE receive no guarantees. Therefore, enclaves must store all their code and private data inside ELRANGE, and must consider the memory outside ELRANGE to be an untrusted interface to the outside world.

The word “linear” in ELRANGE references the linear addresses produced by the vestigial segmentation feature (§ A.7) in the 64-bit Intel architecture. For most purposes, “linear” can be treated as a synonym for “virtual”.

ELRANGE is specified using a base (the BASEADDR field) and a size (the SIZE) in the enclave’s SECS (§ C.4.1). ELRANGE must meet the same constraints as a variable memory type range (§ A.11.4) and as the PRM range (§ C.4.1), namely the size must be a power of 2, and the base must be aligned to the size. These restrictions are in place so that the SGX implementation can inexpensively check whether an address belongs to an enclave’s ELRANGE, in either hardware (§ A.11.4) or software.

When an enclave represents a dynamic library, it is natural to set ELRANGE to the memory range reserved for the library by the loader. The ability to access non-enclave memory from enclave code makes it easy to reuse existing library code that expects to work

with pointers to memory buffers managed by code in the host process.

Non-enclave software cannot access PRM memory. A memory access that resolves inside the PRM results in an aborted transaction, which is undefined at an architectural level. On current processors, aborted writes are ignored, and aborted reads return a value whose bits are all set to 1. This comes into play in the scenario described above, where an enclave is loaded into a host application process as a dynamically loaded library. The system software maps the enclave’s code and data in ELRANGE into EPC pages. If application software attempts to access memory inside ELRANGE, it will experience the abort transaction semantics. The current semantics do not cause the application to crash (e.g., due to a Page Fault), but also guarantee that the host application will not be able to tamper with the enclave or read its private information.

SGX Enclave Attributes

The execution environment of an enclave is heavily influenced by the value of the ATTRIBUTES field in the enclave’s SECS (§ C.4.1). The rest of this work will refer to the field’s sub-fields, shown in Table C.2, as *enclave attributes*.

Field	Bits	Description
DEBUG	1	Opts into enclave debugging features.
XFRM	64	The value of XCR0 (§ A.6) while this enclave’s code is executed.
MODE64BIT	1	Set for 64-bit enclaves.

Table C.2: An enclave’s attributes are the sub-fields in the ATTRIBUTES field of the enclave’s SECS. This table shows a subset of the attributes defined in the SGX documentation.

The most important attribute, from a security perspective, is the DEBUG flag. When this flag is set, it enables the use of SGX’s debugging features for this enclave. These debugging features include the ability to read and modify most of the enclave’s memory. Therefore, DEBUG should only be set in a development environment, as it causes the enclave to lose all the SGX security guarantees.

SGX guarantees that enclave code will always run with the XCR0 register (§ A.6) set to the value indicated by *extended features request mask* (XFRM). Enclave authors are expected to use XFRM to specify the set of architectural extensions enabled by the compiler

used to produce the enclave's code. Having XFRM be explicitly specified allows Intel to design new architectural extensions that change the semantics of existing instructions, such as Memory Protection Extensions (MPX), without having to worry about the security implications on enclave code that was developed without an awareness of the new features.

The MODE64BIT flag is set to true for enclaves that use the 64-bit Intel architecture. From a security standpoint, this flag should not even exist, as supporting a secondary architecture adds unnecessary complexity to the SGX implementation, and increases the probability that security vulnerabilities will creep in. It is very likely that the 32-bit architecture support was included due to Intel's strategy of offering extensive backwards compatibility, which has paid off quite well so far.

In the interest of mental sanity, this work does not analyze the behavior of SGX for enclaves whose MODE64BIT flag is cleared. However, a security researcher who wishes to find vulnerabilities in SGX might study this area.

Last, the INIT flag is always false when the enclave's SECS is created. The flag is set to true at a certain point in the enclave lifecycle, which will be summarized in § C.4.3.

Address Translation for SGX Enclaves

Under SGX, the operating system and hypervisor are still in full control of the page tables and EPTs, and each enclave's code uses the same address translation process and page tables (§ A.5) as its host application. This minimizes the amount of changes required to add SGX support to existing system software. At the same time, having the page tables managed by untrusted system software opens SGX up to the address translation attacks described in § B.7. As future sections will reveal, a good amount of the complexity in SGX's design can be attributed to the need to prevent these attacks.

SGX's active memory mapping attacks defense mechanisms revolve around ensuring that each EPC page can only be mapped at a specific virtual address (§ A.7). When an EPC page is allocated, its intended virtual address is recorded in the EPCM entry for the page, in the ADDRESS field.

When an address translation (§ A.5) result is the physical address of an EPC page, the

CPU ensures¹ that the virtual address given to the address translation process matches the expected virtual address recorded in the page's EPCM entry.

SGX also protects against some passive memory mapping attacks and fault injection attacks by ensuring that the access permissions of each EPC page always match the enclave author's intentions. The access permissions for each EPC page are specified when the page is allocated, and recorded in the *readable* (R), *writable* (W), and *executable* (X) fields in the page's EPCM entry, shown in Table C.3.

Field	Bits	Description
ADDRESS	48	the virtual address used to access this page
R	1	allow reads by enclave code
W	1	allow writes by enclave code
X	1	allow execution of code inside the page, inside enclave

Table C.3: The fields in an EPCM entry that indicate the enclave's intended virtual memory layout.

When an address translation (§ A.5) resolves into an EPC page, the corresponding EPCM entry's fields override the access permission attributes (§ A.5.3) specified in the page tables. For example, the W field in the EPCM entry overrides the writable (W) attribute, and the X field overrides the disable execution (XD) attribute.

It follows that an enclave author must include memory layout information along with the enclave, in such a way that the system software loading the enclave will know the expected virtual memory address and access permissions for each enclave page. In return, the SGX design guarantees to the enclave authors that the system software, which manages the page tables and EPT, will not be able to set up an enclave's virtual address space in a manner that is inconsistent with the author's expectations.

The `.so` and `.dll` file formats, which are SGX's intended enclave delivery vehicles, already have provisions for specifying the virtual addresses that a software module was designed to use, as well as the desired access permissions for each of the module's memory areas.

Last, a SGX-enabled CPU will ensure that the virtual memory inside ELRANGE (§ C.4.2) is mapped to EPC pages. This prevents the system software from carrying out an address

¹A mismatch triggers a general protection fault (#GP, § A.8.2).

translation attack where it maps the enclave's entire virtual address space to DRAM pages outside the PRM, which do not trigger any of the checks above, and can be directly accessed by the system software.

The Thread Control Structure (TCS)

The SGX design fully embraces multi-core processors. It is possible for multiple logical processors (§ A.9.3) to concurrently execute the same enclave's code at the same time, via different threads.

The SGX implementation uses a *Thread Control Structure* (TCS) for each logical processor that executes an enclave's code. It follows that an enclave's author must provision at least as many TCS instances as the maximum number of concurrent threads that the enclave is intended to support.

Each TCS is stored in a dedicated EPC page whose EPCM entry type is PT_TCS. The SDM describes the first few fields in the TCS. These fields are considered to belong to the architectural part of the structure, and therefore are guaranteed to have the same semantics on all the processors that support SGX. The rest of the TCS is not documented.

The contents of an EPC page that holds a TCS cannot be directly accessed, even by the code of the enclave that owns the TCS. This restriction is similar to the restriction on accessing EPC pages holding SECS instances. However, the architectural fields in a TCS can be read by enclave debugging instructions.

The architectural fields in the TCS lay out the context switches (§ A.6) performed by a logical processor when it transitions between executing non-enclave and enclave code.

For example, the OENTRY field specifies the value loaded in the instruction pointer (RIP) when the TCS is used to start executing enclave code, so the enclave author has strict control over the entry points available to enclave's host application. Furthermore, the OFSBASGX and OFSBASGX fields specify the base addresses loaded in the FS and GS segment registers (§ A.7), which typically point to Thread Local Storage (TLS).

The State Save Area (SSA)

When the processor encounters a hardware exception (§ A.8.2), such as an interrupt (§ A.12), while executing the code inside an enclave, it performs a privilege level switch (§ A.8.2) and invokes a hardware exception handler provided by the system software. Before executing the exception handler, however, the processor needs a secure area to store the enclave code's execution context (§ A.6), so that the information in the execution context is not revealed to the untrusted system software.

In the SGX design, the area used to store an enclave thread's execution context while a hardware exception is handled is called a State Save Area (SSA), illustrated in Figure C-6. Each TCS references a contiguous sequence of SSAs. The *offset of the SSA array* (OSSA) field specifies the location of the first SSA in the enclave's virtual address space. The *number of SSAs* (NSSA) field indicates the number of available SSAs.

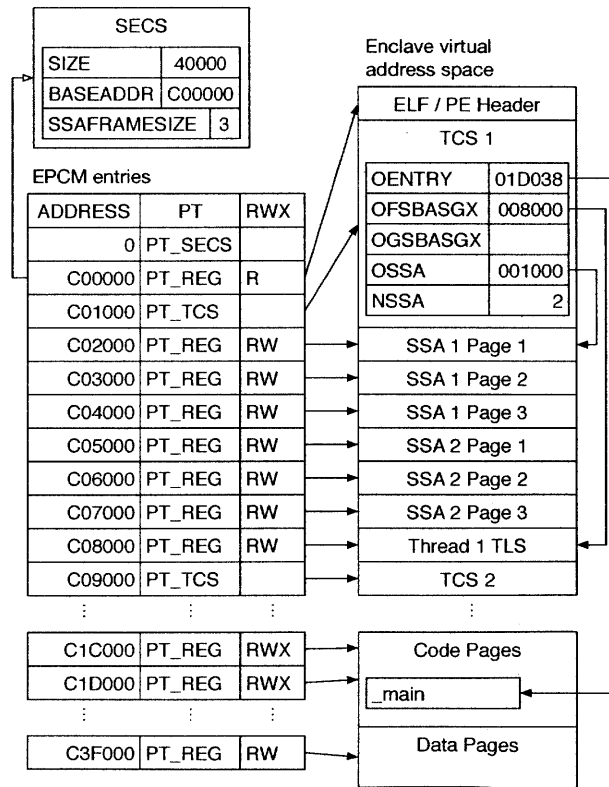


Figure C-6: A possible layout of an enclave's virtual address space. Each enclave has a SECS, and one TCS per supported concurrent thread. Each TCS points to a sequence of SSAs, and specifies initial values for RIP and for the base addresses of FS and GS.

Each SSA starts at the beginning of an EPC page, and uses up the number of EPC pages that is specified in the SSAFRAMESIZE field of the enclave's SECS. These alignment and size restrictions most likely simplify the SGX implementation by reducing the number of special cases that it needs to handle.

An enclave thread's execution context consists of the general-purpose registers (GPRs) and the result of the XSAVE instruction (§ A.6). Therefore, the size of the execution context depends on the requested-feature bitmap (RFBM) used by XSAVE. All the code in an enclave uses the same RFBM, which is declared in the XFRM enclave attribute (§ C.4.2). The number of EPC pages reserved for each SSA, specified in SSAFRAMESIZE, must² be large enough to fit the XSAVE output for the feature bitmap specified by XFRM.

SSAs are stored in regular EPC pages, whose EPCM page type is PT_REG. Therefore, the SSA contents is accessible to enclave software. The SSA layout is architectural, and is completely documented in the SDM. This opens up possibilities for an enclave exception handler that is invoked by the host application after a hardware exception occurs, and acts upon the information in a SSA.

C.4.3 The Life Cycle of an SGX Enclave

An enclave's life cycle is deeply intertwined with resource management, specifically the allocation of EPC pages. Therefore, the instructions that transition between different life cycle states can only be executed by the system software. The system software is expected to expose the SGX instructions described below as enclave loading and teardown services.

The following subsections describe the major steps in an enclave's lifecycle, which is illustrated by Figure C-7.

Creation

An enclave is born when the system software issues the ECREATE instruction, which turns a free EPC page into the SECS (§ C.4.1) for the new enclave.

ECREATE initializes the newly created SECS using the information in a non-EPC page

²ECREATE (§ C.4.3) fails if SSAFRAMESIZE is too small.

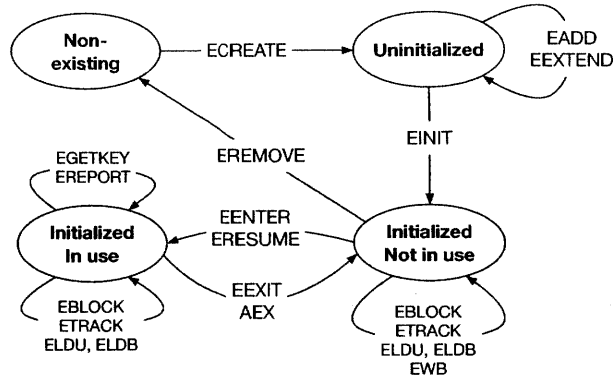


Figure C-7: The SGX enclave life cycle management instructions and state transition diagram

owned by the system software. This page specifies the values for all the SECS fields defined in the SDM, such as `BASEADDR` and `SIZE`, using an architectural layout that is guaranteed to be preserved by future implementations.

While it is very likely that the actual SECS layout used by initial SGX implementations matches the architectural layout quite closely, future implementations are free to deviate from this layout, as long as they maintain the ability to initialize the SECS using the architectural layout. Software cannot access an EPC page that holds a SECS, so it cannot become dependent on an internal SECS layout. This is a stronger version of the encapsulation used in the Virtual Machine Control Structure (VMCS, § A.8.3).

`ECREATE` validates the information used to initialize the SECS, and results in a page fault (`#PF`, § A.8.2) or general protection fault (`#GP`, § A.8.2) if the information is not valid. For example, if the `SIZE` field is not a power of two, `ECREATE` results in `#GP`. This validation, combined with the fact that the SECS is not accessible by software, simplifies the implementation of the other SGX instructions, which can assume that the information inside the SECS is valid.

Last, `ECREATE` initializes the enclave's `INIT` attribute (sub-field of the `ATTRIBUTES` field in the enclave's SECS, § C.4.2) to the false value. The enclave's code cannot be executed until the `INIT` attribute is set to true, which happens in the initialization stage that will be described in § C.4.3.

Loading

ECREATE marks the newly created SECS as *uninitialized*. While an enclave's SECS is in this state, the system software can use EADD instructions to load the initial code and data into the enclave. EADD is used to create both TCS pages (§ C.4.2) and regular pages.

EADD reads its input data from a *Page Information* (PAGEINFO) structure, illustrated in Figure C-8. The structure's contents are only used to communicate information to the SGX implementation, so it is entirely architectural and documented in the SDM.

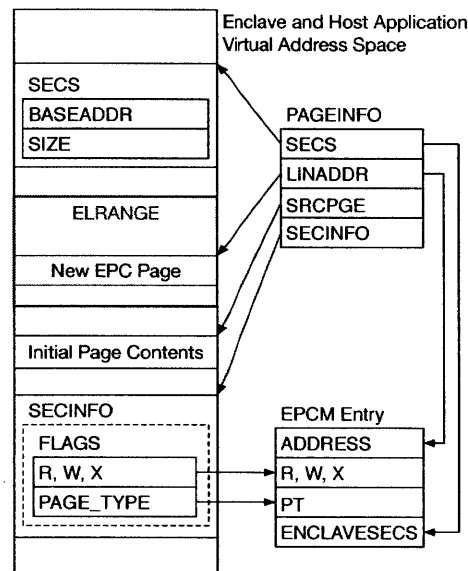


Figure C-8: The PAGEINFO structure supplies input data to SGX instructions such as EADD.

Currently, the PAGEINFO structure contains the virtual address of the EPC page that will be allocated (LINADDR), the virtual address of the non-EPC page whose contents will be copied into the newly allocated EPC page (SRCPGE), a virtual address that resolves to the SECS of the enclave that will own the page (SECS), and values for some of the fields of the EPCM entry associated with the newly allocated EPC page (SECINFO).

The SECINFO field in the PAGEINFO structure is actually a virtual memory address, and points to a *Security Information* (SECINFO) structure, some of which is also illustrated in Figure C-8. The SECINFO structure contains the newly allocated EPC page's access permissions (R, W, X) and its EPCM page type (PT_REG or PT_TCS). Like PAGEINFO,

the SECINFO structure is solely used to communicate data to the SGX implementation, so its contents are also entirely architectural. However, most of the structure's 64 bytes are reserved for future use.

Both the PAGEINFO and the SECINFO structures are prepared by the system software that invokes the EADD instruction, and therefore must be contained in non-EPC pages. Both structures must be aligned to their sizes – PAGEINFO is 32 bytes long, so each PAGEINFO instance must be 32-byte aligned, while SECINFO has 64 bytes, and therefore each SECINFO instance must be 64-byte aligned. The alignment requirements likely simplify the SGX implementation by reducing the number of special cases that must be handled.

EADD validates its inputs before modifying the newly allocated EPC page or its EPCM entry. Most importantly, attempting to EADD a page to an enclave whose SECS is in the initialized state will result in a #GP. Furthermore, attempting to EADD an EPC page that is already allocated (the VALID field in its EPCM entry is 1) results in a #PF. EADD also ensures that the page's virtual address falls within the enclave's ELRANGE, and that all the reserved fields in SECINFO are set to zero.

While loading an enclave, the system software will also use the EEXTEND instruction, which updates the enclave's measurement used in the software attestation process. Software attestation is discussed in § C.4.8.

Initialization

After loading the initial code and data pages into the enclave, the system software must use a *Launch Enclave* (LE) to obtain an EINIT Token Structure, via an under-documented process that will be described in more detail in § C.4.9. The token is then provided to the EINIT instruction, which marks the enclave's SECS as *initialized*.

The LE is a privileged enclave provided by Intel, and **is a prerequisite for the use of enclaves authored by parties other than Intel**. The LE is an SGX enclave, so it must be created, loaded and initialized using the processes described in this section. However, the LE is cryptographically signed (§ B.1.3) with a special Intel key that is hard-coded into the SGX implementation, and that causes EINIT to initialize the LE without checking for a

valid EINIT Token Structure.

When EINIT completes successfully, it sets the enclave's INIT attribute to true. This opens the way for ring 3 (§ A.3) application software to execute the enclave's code, using the SGX instructions described in § C.4.4. On the other hand, once INIT is set to true, EADD cannot be invoked on that enclave anymore, so the system software must load all the pages that make up the enclave's initial state before executing the EINIT instruction.

Teardown

After the enclave has done the computation it was designed to perform, the system software executes the EREMOVE instruction to deallocate the EPC pages used by the enclave.

EREMOVE marks an EPC page as available by setting the VALID field of the page's EPCM entry to 0 (zero). Before freeing up the page, EREMOVE makes sure that there is no logical processor executing code inside the enclave that owns the page to be removed.

An enclave is completely destroyed when the EPC page holding its SECS is freed. EREMOVE refuses to deallocate a SECS page if it is referenced by any other EPCM entry's ENCLAVESECS field, so an enclave's SECS page can only be deallocated after all the enclave's pages have been deallocated.

C.4.4 The Life Cycle of an SGX Thread

Between the time when an enclave is initialized (§ C.4.3) and the time when it is torn down (§ C.4.3), the enclave's code can be executed by any application process that has the enclave's EPC pages mapped into its virtual address space.

When executing the code inside an enclave, a logical processor is said to be *in enclave mode*, and the code that it executes can access the regular (PT_REG, § C.4.1) EPC pages that belong to the currently executing enclave. When a logical process is outside enclave mode, it bounces any memory accesses inside the Processor Reserved Memory range (PRM, § C.4.1), which includes the EPC.

Each logical processor that executes enclave code uses a Thread Control Structure (TCS, § C.4.2). When a TCS is used by a logical processor, it is said to be *busy*, and it cannot be used by

any other logical processor. Figure C-9 illustrates the instructions used by a host process to execute enclave code and their interactions with the TCS that they target.

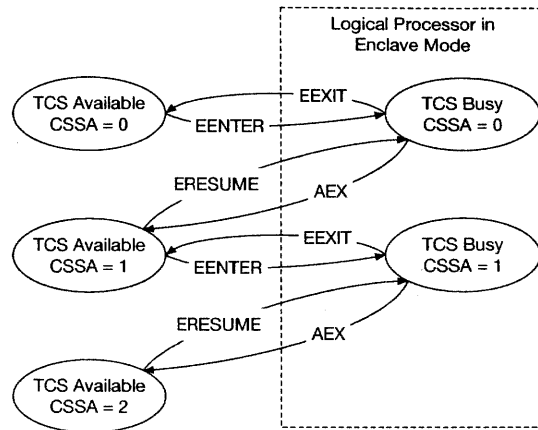


Figure C-9: The stages of the life cycle of an SGX Thread Control Structure (TCS) that has two State Save Areas (SSAs).

Assuming that no hardware exception occurs, an enclave’s host process uses the `EENTER` instruction, described in § C.4.4, to execute enclave code. When the enclave code finishes performing its task, it uses the `EEXIT` instruction, covered in § C.4.4, to return the execution control to the host process that invoked the enclave.

If a hardware exception occurs while a logical processor is in enclave mode, the processor is taken out of enclave mode using an *Asynchronous Enclave Exit* (AEX), summarized in § C.4.4, before the system software’s exception handler is invoked. After the system software’s handler is invoked, the enclave’s host process can use the `ERESUME` instruction, described in § C.4.4, to re-enter the enclave and resume the computation that it was performing.

Synchronous Enclave Entry

At a high level, `EENTER` performs a controlled jump into enclave code, while performing the processor configuration that is needed by SGX’s security guarantees. Going through all the configuration steps is a tedious exercise, but it is a necessary prerequisite to understanding how all data structures used by SGX work together. For this reason, `EENTER` and its siblings are described in much more detail than the other SGX instructions.

EENTER, illustrated in Figure C-10 can only be executed by unprivileged application software running at ring 3 (§ A.3), and results in an undefined instruction (#UD) fault if is executed by system software.

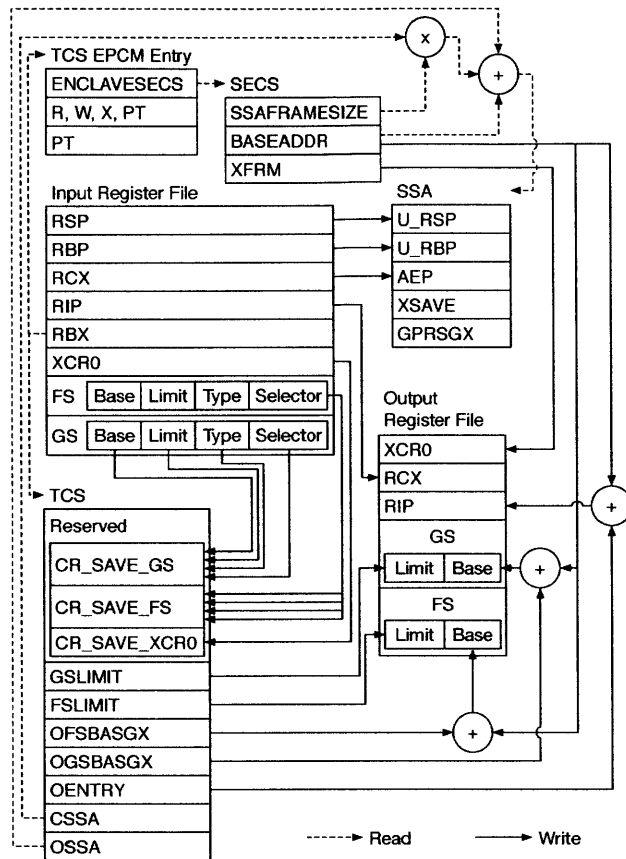


Figure C-10: Data flow diagram for a subset of the logic in EENTER. The figure omits the logic for disabling debugging features, such as hardware breakpoints and performance monitoring events.

EENTER switches the logical processor to enclave mode, but does not perform a privilege level switch (§ A.8.2). Therefore, enclave code always executes at ring 3, with the same privileges as the application code that calls it. This makes it possible for an infrastructure owner to allow user-supplied software to create and use enclaves, while having the assurance that the OS kernel and hypervisor can still protect the infrastructure from buggy or malicious software.

EENTER takes the virtual address of a TCS as its input, and requires that the TCS is *available* (not busy), and that at least one State Save Area (SSA, § C.4.2) is available in the

TCS. The latter check is implemented by making sure that the *current SSA index* (CSSA) field in the TCS is less than the number of SSAs (NSSA) field. The SSA indicated by the CSSA, which shall be called the *current SSA*, is used in the event that a hardware exception occurs while enclave code is executed.

EENTER transitions the logical processor into enclave mode, and sets the instruction pointer (RIP) to the value indicated by the *entry point offset* (OENTRY) field in the TCS that it receives. EENTER is used by an untrusted caller to execute code in a protected environment, and therefore has the same security considerations as SYSCALL (§ A.8), which is used to call into system software. Setting RIP to the value indicated by OENTRY guarantees to the enclave author that the enclave code will only be invoked at well defined points, and prevents a malicious host application from bypassing any security checks that the enclave author may perform.

EENTER also sets XCR0 (§ A.6), the register that controls which extended architectural features are in use, to the value of the XFRM enclave attribute (§ C.4.2). Ensuring that XCR0 is set according to the enclave author's intentions prevents a malicious operating system from bypassing an enclave's security by enabling architectural features that the enclave is not prepared to handle.

Furthermore, EENTER loads the bases of the segment registers (§ A.7) FS and GS using values specified in the TCS. The segments' selectors and types are hard-coded to safe values for ring 3 data segments. This aspect of the SGX design makes it easy to implement per-thread Thread Local Storage (TLS). For 64-bit enclaves, this is a convenience feature rather than a security measure, as enclave code can securely load new bases into FS and GS using the WRFSBASE and WRGSBASE instructions.

The EENTER implementation backs up the old values of the registers that it modifies, so they can be restored when the enclave finishes its computation. Just like SYSCALL, EENTER saves the address of the following instruction in the RCX register.

Interestingly, the SDM states that the old values of the XCR0, FS, and GS registers are saved in new registers dedicated to the SGX implementation. However, given that they will only be used on an enclave exit, we expect that the registers are saved in DRAM, in the reserved area in the TCS.

Like `SYSCALL`, `EENTER` does not modify the stack pointer register (`RSP`). To avoid any security exploits, enclave code should set `RSP` to point to a stack area that is entirely contained in EPC pages. Multi-threaded enclaves can easily implement per-thread stack areas by setting up each thread's TLS area to include a pointer to the thread's stack, and by setting `RSP` to the value obtained by reading the TLS area at which the FS or GS segment points.

Last, when `EENTER` enters enclave mode, it suspends some of the processor's debugging features, such as hardware breakpoints and Precise Event Based Sampling (PEBS). Conceptually, a debugger attached to the host process sees the enclave's execution as one single processor instruction.

Synchronous Enclave Exit

`EEXIT` can only be executed while the logical processor is in enclave mode, and results in a `(#UD)` if executed in any other circumstances. In a nutshell, the instruction returns the processor to ring 3 outside enclave mode and restores the registers saved by `EENTER`, which were described above.

Unlike `SYSRET`, `EEXIT` sets `RIP` to the value read from `RBX`, after exiting enclave mode. This is inconsistent with `EENTER`, which saves the `RIP` value to `RCX`. Unless this inconsistency stems from an error in the SDM, enclave code must be sure to note the difference.

The SDM explicitly states that `EEXIT` does not modify most registers, so enclave authors must make sure to clear any secrets stored in the processor's registers before returning control to the host process. Furthermore, enclave software will most likely cause a fault in its caller if it doesn't restore the stack pointer `RSP` and the stack frame base pointer `RBP` to the values that they had when `EENTER` was called.

It may seem unfortunate that enclave code can induce faults in its caller. For better or for worse, this perfectly matches the case where an application calls into a dynamically loaded module. More specifically, the module's code is also responsible for preserving stack-related registers, and a buggy module might jump anywhere in the application code of the host process.

This section describes the `EENTER` behavior for 64-bit enclaves. The `EENTER` implementation for 32-bit enclaves is significantly more complex, due to the extra special cases introduced by the full-fledged segmentation model that is still present in the 32-bit Intel architecture. As stated in the introduction, we are not interested in such legacy aspects.

Asynchronous Enclave Exit (AEX)

If a hardware exception, like a fault (§ A.8.2) or an interrupt (§ A.12), occurs while a logical processor is executing an enclave’s code, the processor performs an *Asynchronous Enclave Exit (AEX)* before invoking the system software’s exception handler, as shown in Figure C-11.

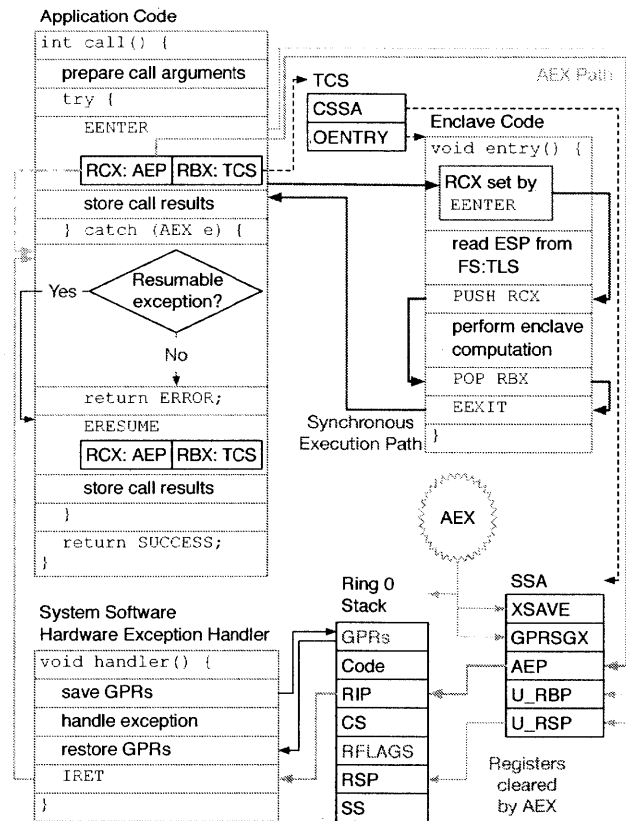


Figure C-11: If a hardware exception occurs during enclave execution, the synchronous execution path is aborted, and an Asynchronous Enclave Exit (AEX) occurs instead.

The AEX saves the enclave code’s execution context (§ A.6), restores the state saved by `EENTER`, and sets up the processor registers so that the system software’s hardware excep-

tion handler will return to an *asynchronous exit handler* in the enclave's host process. The exit handler is expected to use the `ERESUME` instruction to resume the enclave computation that was interrupted by the hardware exception.

Asides from the behavior described in § C.4.4, `EENTER` also writes some information to the current SSA, which is only used if an AEX occurs. As shown in Figure C-10, `EENTER` stores the stack pointer register `RSP` and the stack frame base pointer register `RBP` into the `U_RSP` and `U_RBP` fields in the current SSA. Last, `EENTER` stores the value in `RCX` in the *Asynchronous Exit handler Pointer (AEP)* field in the current SSA.

When a hardware exception occurs in enclave mode, the SGX implementation performs a sequence of steps that takes the logical processor out of enclave mode and invokes the hardware exception handler in the system software. Conceptually, the SGX implementation first performs an AEX to take the logical processor out of enclave mode, and then the hardware exception is handled using the standard Intel architecture's behavior described in § A.8.2. Actual Intel processors may interleave the AEX implementation with the exception handling implementation. However, for simplicity, this work describes AEX as a separate process that is performed before any exception handling steps are taken.

In the Intel architecture, if a hardware exception occurs, the application code's execution context can be read and modified by the system software's exception handler (§ A.8.2). This is acceptable when the system software is trusted by the application software. However, under SGX's threat model, the system software is not trusted by enclaves. Therefore, the AEX step erases any secrets that may exist in the execution state by resetting all its registers to predefined values.

Before the enclave's execution state is reset, it is backed up inside the current SSA. Specifically, an AEX backs up the general purpose registers (GPRs, § A.6) in the `GPRSGX` area in the SSA, and then performs an `XSAVE` (§ A.6) using the requested-feature bitmap (RFBM) specified in the `XFRM` field in the enclave's `SECS`. As each SSA is entirely stored in EPC pages allocated to the enclave, the system software cannot read or tamper with the backed up execution state. When an SSA receives the enclave's execution state, it is marked as used by incrementing the `CSSA` field in the current `TCS`.

After clearing the execution context, the AEX process sets `RSP` and `RBP` to the values

saved by `EENTER` in the current SSA, and sets `RIP` to the value in the current SSA's `AEP` field. This way, when the system software's hardware exception handler completes, the processor will execute the asynchronous exit handler code in the enclave's host process. The SGX design makes it easy to set up the asynchronous handler code as an exception handler in the routine that contains the `EENTER` instruction, because the `RSP` and `RBP` registers will have the same values as they had when `EENTER` was executed.

Many of the actions taken by `AEX` to get the logical processor outside of enclave mode match `EEXIT`. The segment registers `FS` and `GS` are restored to the values saved by `EENTER`, and all the debugging facilities that were suppressed by `EENTER` are restored to their previous states.

Recovering from an Asynchronous Exit

When a hardware exception occurs inside enclave mode, the processor performs an `AEX` before invoking the exception's handler set up by the system software. The `AEX` sets up the execution context in such a way that when the system software finishes processing the exception, it returns into an asynchronous exit handler in the enclave's host process. The asynchronous exception handler usually executes the `ERESUME` instruction, which causes the logical processor to go back into enclave mode and continue the computation that was interrupted by the hardware exception.

`ERESUME` shares much of its functionality with `EENTER`. This is best illustrated by the similarity between Figures C-12 and C-11.

`EENTER` and `ERESUME` receive the same inputs, namely a pointer to a TCS, described in § C.4.4, and an `AEP`, described in § C.4.4. The most common application design will pair each `EENTER` instance with an asynchronous exit handler that invokes `ERESUME` with exactly the same arguments.

The main difference between `ERESUME` and `EENTER` is that the former uses an SSA that was "filled out" by an `AEX` (§ C.4.4), whereas the latter uses an empty SSA. Therefore, `ERESUME` results in a `#GP` fault if the `CSSA` field in the provided TCS is 0 (zero), whereas `EENTER` fails if `CSSA` is greater than or equal to `NSSA`.

When successful, `ERESUME` decrements the `CSSA` field of the TCS, and restores the

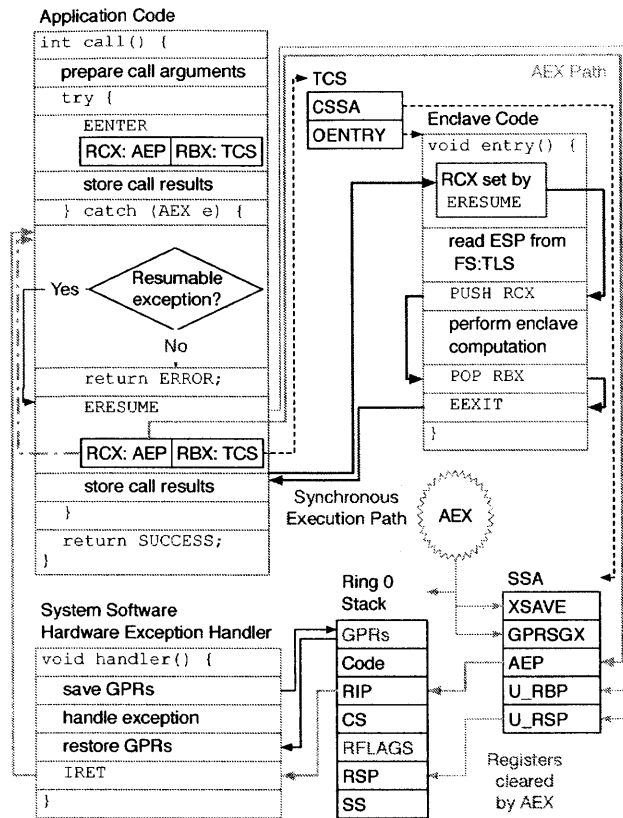


Figure C-12: If a hardware exception occurs during enclave execution, the synchronous execution path is aborted, and an Asynchronous Enclave Exit (AEX) occurs instead.

execution context backed up in the SSA pointed to by the CSSA field in the TCS. Specifically, the ERESUME implementation restores the GPRs (§ A.6) from the GPRSGX field in the SSA, and performs an XRSTOR (§ A.6) to load the execution state associated with the extended architectural features used by the enclave.

ERESUME shares the following behavior with EENTER (§ C.4.4). Both instructions write the U_RSP, U_RBP, and AEP fields in the current SSA. Both instructions follow the same process for backing up XCR0 and the FS and GS segment registers, and set them to the same values, based on the current TCS and its enclave’s SECS. Last, both instructions disable the same subset of the logical processor’s debugging features.

An interesting edge case that ERESUME handles correctly is that it sets XCR0 to the XFRM enclave attribute **before** performing an XRSTOR. It follows that ERESUME fails if the requested feature bitmap (RFBM) in the SSA is not a subset of XFRM. This matters

because, while an AEX will always use the XFRM value as the RFBM, enclave code executing on another thread is free to modify the SSA contents before `ERESUME` is called.

The correct sequencing of actions in the `ERESUME` implementation prevents a malicious application from using an enclave to modify registers associated with extended architectural features that are not declared in `XFRM`. This would break the system software's ability to provide thread-level execution context isolation.

C.4.5 EPC Page Eviction

Modern OS kernels take advantage of address translation (§ A.5) to implement page swapping, also referred to as paging (§ A.5). In a nutshell, paging allows the OS kernel to over-commit the computer's DRAM by evicting rarely used memory pages to a slower storage medium called the disk.

Paging is a key contributor to utilizing a computer's resources effectively. For example, a desktop system whose user runs multiple programs concurrently can evict memory pages allocated to inactive applications without a significant degradation in user experience.

Unfortunately, the OS cannot be allowed to evict an enclave's EPC pages via the same methods that are used to implement page swapping for DRAM memory outside the PRM range. In the SGX threat model, enclaves do not trust the system software, so the SGX design offers an EPC page eviction method that can defend against a malicious OS that attempts any of the active address translation attacks described in § B.7.

The price of the security afforded by SGX is that an OS kernel that supports evicting EPC pages must use a modified page swapping implementation that interacts with the SGX mechanisms. Enclave authors can mostly ignore EPC evictions, similarly to how today's application developers can ignore the OS kernel's paging implementation.

As illustrated in Figure C-13, SGX supports evicting EPC pages to DRAM pages outside the PRM range. The system software is expected to use its existing page swapping implementation to evict the contents of these pages out of DRAM and onto a disk.

SGX's eviction feature revolves around the `EWB` instruction, described in detail in § C.4.5. Essentially, `EWB` evicts an EPC page into a DRAM page outside the EPC and marks the

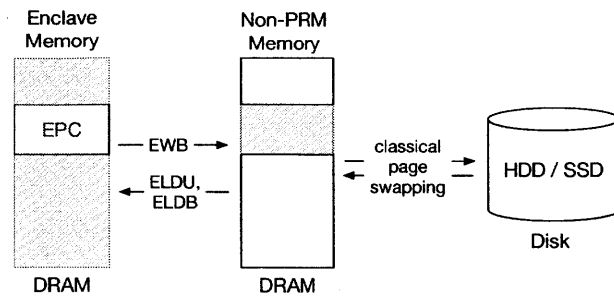


Figure C-13: SGX offers a method for the OS to evict EPC pages into non-PRM DRAM. The OS can then use its standard paging feature to evict the pages out of DRAM.

EPC page as available, by zeroing the VALID field in the page’s EPCM entry.

The SGX design relies on symmetric key cryptography B.1.1 to guarantee the privacy and integrity of the evicted EPC pages, and on nonces (§ B.1.4) to guarantee the freshness of the pages brought back into the EPC. These nonces are stored in Version Arrays (VAs), covered in § C.4.5, which are EPC pages dedicated to nonce storage.

Before an EPC page is evicted and freed up for use by other enclaves, the SGX implementation must ensure that no TLB has address translations associated with the evicted page, in order to avoid the TLB-based address translation attack described in § B.7.4.

As explained in § C.4.1, SGX leaves the system software in charge of managing the EPC. It naturally follows that the SGX instructions described in this section, which are used to implement EPC paging, are only available to system software, which runs at ring 0 § A.3.

In today’s software stacks (§ A.3), only the OS kernel implements page swapping in order to support the over-committing of DRAM. The hypervisor is only used to partition the computer’s physical resources between operating systems. Therefore, this section is written with the expectation that the OS kernel will also take on the responsibility of EPC page swapping. For simplicity, we often use the term “OS kernel” instead of “system software”. The reader should be aware that the SGX design does not preclude a system where the hypervisor implements its own EPC page swapping. Therefore, “OS kernel” should really be read as “the system software that performs EPC paging”.

Page Eviction and the TLBs

One of the least promoted accomplishments of SGX is that it does not add any security checks to the memory execution units (§ A.9.4, § A.10). Instead, SGX's access control checks occur after an address translation (§ A.5) is performed, right before the translation result is written into the TLBs (§ A.11.5). This aspect is generally downplayed throughout the SDM, but it becomes visible when explaining SGX's EPC page eviction mechanism.

A full discussion of SGX's memory access protections checks merits its own section, and is deferred to § C.5.2. The EPC page eviction mechanisms can be explained using only two requirements from SGX's security model. First, when a logical processor exits an enclave, either via `EEXIT` (§ C.4.4) or via an `AEX` (§ C.4.4), its TLBs are flushed. Second, when an EPC page is deallocated from an enclave, all logical processors executing that enclave's code must be directed to exit the enclave. This is sufficient to guarantee the removal of any TLB entry targeting the deallocated EPC.

System software can cause a logical processor to exit an enclave by sending it an Inter-Processor Interrupt (IPI, § A.12), which will trigger an `AEX` when received. Essentially, this is a very coarse-grained TLB shutdown.

SGX does not trust system software. Therefore, before marking an EPC page's EPCM entry as free, the SGX implementation must ensure that the OS kernel has flushed all the TLBs that might contain translations for the page. Furthermore, performing IPIs and TLB flushes for each page eviction would add a significant overhead to a paging implementation, so the SGX design allows a batch of pages to be evicted using a single IPI / TLB flush sequence.

The TLB flush verification logic relies on a 1-bit EPCM entry field called `BLOCKED`. As shown in Figure C-14, the `VALID` and `BLOCKED` fields yield three possible EPC page states. A page is *free* when both bits are zero, *in use* when `VALID` is zero and `BLOCKED` is one, and *blocked* when both bits are one.

Blocked pages are not considered accessible to enclaves. If an address translation results in a blocked EPC page, the SGX implementation causes the translation to result in a Page Fault (`#PF`, § A.8.2). This guarantees that once a page is blocked, the CPU will not create

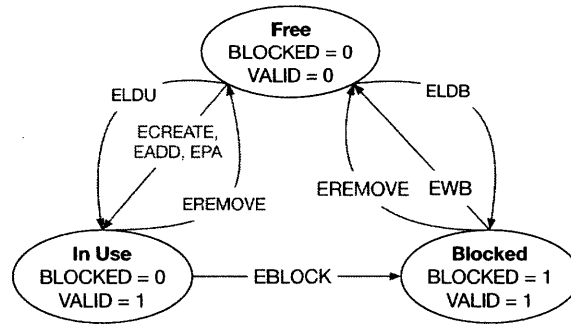


Figure C-14: The VALID and BLOCKED bits in an EPC page’s EPCM entry can be in one of three states. EADD and its siblings allocate new EPC pages. EREMOVE permanently deallocates an EPC page. EBLOCK blocks an EPC page so it can be evicted using EWB. ELDB and ELDU load an evicted page back into the EPC.

any new TLB entries pointing to it.

Furthermore, every SGX instruction makes sure that the EPC pages on which it operates are not blocked. For example, EENTER ensures that the TCS it is given is not blocked, that its enclave’s SECS is not blocked, and that every page in the current SSA is not blocked.

In order to evict a batch of EPC pages, the OS kernel must first issue EBLOCK instructions targeting them. The OS is also expected to remove the EPC page’s mapping from page tables, but is not trusted to do so.

After all the desired pages have been blocked, the OS kernel must execute an ETRACK instruction, which directs the SGX implementation to keep track of which logical processors have had their TLBs flushed. ETRACK requires the virtual address of an enclave’s SECS (§ C.4.1). If the OS wishes to evict a batch of EPC pages belonging to multiple enclaves, it must issue an ETRACK for each enclave.

Following the ETRACK instructions, the OS kernel must induce enclave exits on all the logical processors that are executing code inside the enclaves that have been ETRACKED. The SGX design expects that the OS will use IPIs to cause AEXs in the logical processors whose TLBs must be flushed.

The EPC page eviction process is completed when the OS executes an EWB instruction for each EPC page to be evicted. This instruction, which will be fully described in § C.4.5, writes an encrypted version of the EPC page to be evicted into DRAM, and then frees the page by clearing the VALID and BLOCKED bits in its EPCM entry. Before carrying out its

tasks, EWB ensures that the EPC page that it targets has been blocked, and checks the state set up by ETRACK to make sure that all the relevant TLBs have been flushed.

An evicted page can be loaded back into the EPC via the ELDU and ELDB instructions. Both instructions start up with a free EPC page and a DRAM page that has the evicted contents of an EPC page, decrypt the DRAM page's contents into the EPC page, and restore the corresponding EPCM entry. The only difference between ELDU and ELDB is that the latter sets the BLOCKED bit in the page's EPCM entry, whereas the former leaves it cleared.

ELDU and ELDB resemble ECREATE and EADD, in the sense that they populate a free EPC page. Since the page that they operate on was free, the SGX security model predicates that no TLB entries can possibly target it. Therefore, these instructions do not require a mechanism similar to EBLOCK or ETRACK.

The Version Array (VA)

When EWB evicts the contents of an EPC, it creates an 8-byte nonce (§ B.1.4) that Intel's documentation calls a *page version*. SGX's freshness guarantees are built on the assumption that nonces are stored securely, so EWB stores the nonce that it creates inside a *Version Array (VA)*.

Version Arrays are EPC pages that are dedicated to storing nonces generated by EWB. Each VA is divided into slots, and each slot is exactly large enough to store one nonce. Given that the size of an EPC page is 4KB, and each nonce occupies 8 bytes, it follows that each VA has 512 slots.

VA pages are allocated using the EPA instruction, which takes in the virtual address of a free EPC page, and turns it into a Version Array with empty slots. VA pages are identified by the PT_VA type in their EPCM entries. Like SECS pages, VA pages have the ENCLAVEADDRESS fields in their EPCM entries set to zero, and cannot be accessed directly by any software, including enclaves.

Unlike the other page types discussed so far, VA pages are not associated with any enclave. This means they can be deallocated via EREMOVE without any restriction. However, freeing up a VA page whose slots are in use effectively discards the nonces in those slots, which results in losing the ability to load the corresponding evicted pages back into the EPC.

Therefore, it is unlikely that a correct OS implementation will ever call `EREMOVE` on a VA with non-free slots.

According to the pseudo-code for `EPA` and `EWB` in the SDM, SGX uses the zero value to represent the free slots in a VA, implying that all the generated nonces have to be non-zero. This also means that `EPA` initializes a VA simply by zeroing the underlying EPC page. However, since software cannot access a VA's contents, neither the use of a special value, nor the value itself is architectural.

Enclave IDs

The `EWB` and `ELDU / ELDB` instructions use an *enclave ID* (EID) to identify the enclave that owns an evicted page. The EID has the same purpose as the `ENCLAVESECS` (§ C.4.1) field in an EPCM entry, which is also used to identify the enclave that owns an EPC page. This section explains the need for having two values represent the same concept by comparing the two values and their uses.

The SDM states that `ENCLAVESECS` field in an EPCM entry is used to identify the SECS of the enclave owning the associated EPC page, but stops short of describing its format. In theory, the `ENCLAVESECS` field can change its representation between SGX implementations since SGX instructions never expose its value to software.

However, we will later argue that the most plausible representation of the `ENCLAVESECS` field is the physical address of the enclave's SECS. Therefore, the `ENCLAVESECS` value associated with a given enclave will change if the enclave's SECS is evicted from the EPC and loaded back at a different location. It follows that the `ENCLAVESECS` value is only suitable for identifying an enclave while its SECS remains in the EPC.

According to the SDM, the EID field is a 64-bit field stored in an enclave's SECS. `ECREATE`'s pseudocode in the SDM reveals that an enclave's ID is generated when the SECS is allocated, by atomically incrementing a global counter. Assuming that the counter does not roll over³, this process guarantees that every enclave created during a power cycle has a unique EID.

Although the SDM does not specifically guarantee this, the EID field in an enclave's

³A 64-bit counter incremented at 4Ghz would roll over in slightly more than 136 years

SECS does not appear to be modified by any instruction. This makes the EID's value suitable for identifying an enclave throughout its lifetime, even across evictions of its SECS page from the EPC.

Evicting an EPC Page

The system software evicts an EPC page using the EWB instruction, which produces all the data needed to restore the evicted page at a later time via the ELDU instruction, as shown in Figure C-15.

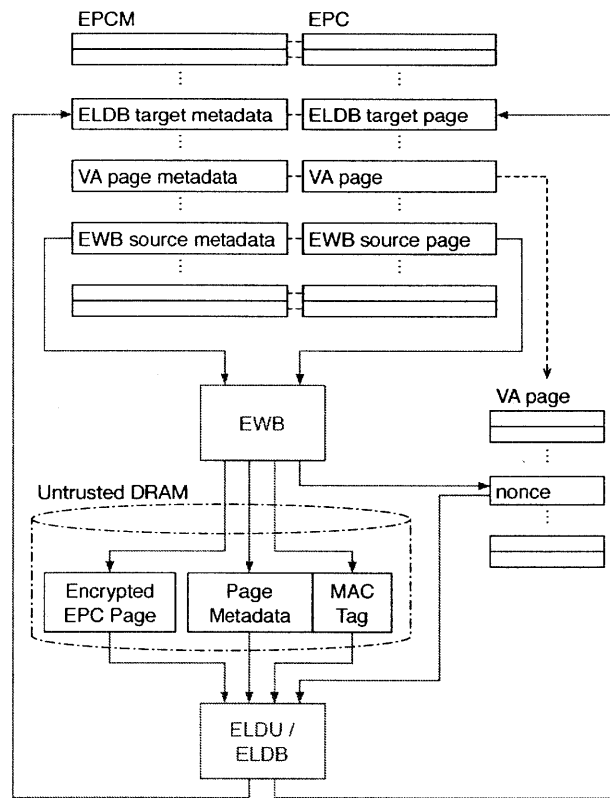


Figure C-15: The EWB instruction outputs the encrypted contents of the evicted EPC page, a subset of the fields in the page's EPCM entry, a MAC tag, and a nonce. All this information is used by the ELDB or ELDU instruction to load the evicted page back into the EPC, with privacy, integrity and freshness guarantees.

EWB's output consists of an encrypted version of the evicted EPC page's contents, a subset of the fields in the EPCM entry corresponding to the page, the nonce discussed in § C.4.5, and a message authentication code (MAC, § B.1.3) tag. With the exception of the

nonce, EWB writes its output in DRAM outside the PRM area, so the system software can choose to further evict it to disk.

The EPC page contents is encrypted, to protect the privacy of the enclave's data while the page is stored in the untrusted DRAM outside the PRM range. Without the use of encryption, the system software could learn the contents of an EPC page by evicting it from the EPC.

The page metadata is stored in a *Page Information* (PAGEINFO) structure, illustrated in Figure C-16. This structure is similar to the PAGEINFO structure described in § C.4.3 and depicted in Figure C-8, except that the SECINFO field has been replaced by a PCMD field, which contains the virtual address of a *Page Crypto Metadata* (PCMD) structure.

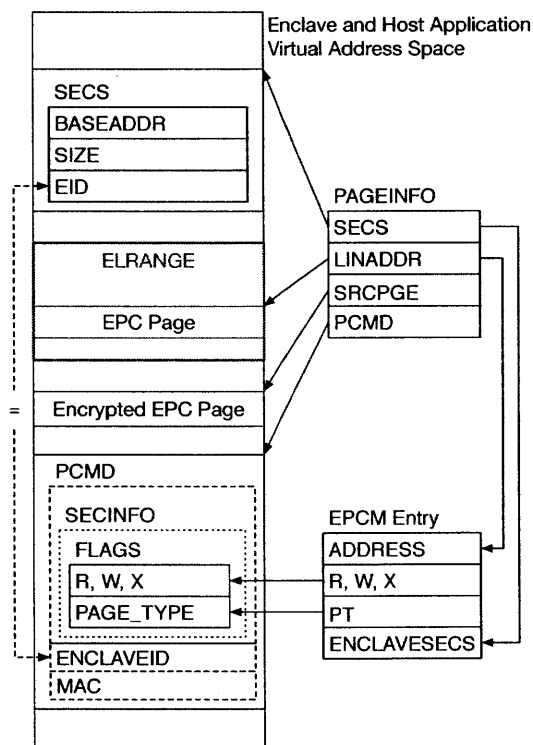


Figure C-16: The PAGEINFO structure used by the EWB and ELDU / ELDB instructions

The LINADDR field in the PAGEINFO structure is used to store the ADDRESS field in the EPCM entry, which indicates the virtual address intended for accessing the page. The PCMD structure embeds the *Security Information* (SECINFO) described in § C.4.3, which is used to store the page type (PT) and the access permission flags (R, W, X) in the EPCM

entry. The PCMD structure also stores the enclave's ID (EID, § C.4.5). These fields are later used by ELDU or ELDB to populate the EPCM entry for the EPC page that is reloaded.

The metadata described above is stored unencrypted, so the OS has the option of using the information inside as-is for its own bookkeeping. This has no negative impact on security, because the metadata is not confidential. In fact, with the exception of the enclave ID, all the metadata fields are specified by the system software when ECREATE is called. The enclave ID is only useful for identifying the enclave that the EPC page belongs to, and the system software already has this information as well.

Asides from the metadata described above, the PCMD structure also stores the MAC tag generated by EWB. The MAC tag covers the authenticity of the EPC page contents, the metadata, and the nonce. The MAC tag is checked by ELDU and ELDB, which will only load an evicted page back into the EPC if the MAC verification confirms the authenticity of the page data, metadata, and nonce. This security check protects against the page swapping attacks described in § B.7.3.

Similarly to EREMOVE, EWB will only evict the EPC page holding an enclave's SECS if there is no other EPCM entry whose ENCLAVESECS field references the SECS. At the same time, as an optimization, the SGX implementation does not perform ETRACK-related checks when evicting a SECS. This is safe because a SECS is only evicted if the EPC has no pages belonging to the SECS' enclave, which implies that there isn't any TCS belonging to the enclave in the EPC, so no processor can be executing enclave code.

The pages holding Version Arrays can be evicted, just like any other EPC page. VA pages are never accessible by software, so they can't have any TLB entries pointing to them. Therefore, EWB evicts VA pages without performing any ETRACK-related checks. The ability to evict VA pages has profound implications that will be discussed in § C.4.5.

EWB's data flow, shown in detail in Figure C-17, has an aspect that can be confusing to OS developers. The instruction reads the virtual address of the EPC page to be evicted from a register (RBX) and writes it to the LINADDR field of the PAGEINFO structure that it is provided. The separate input (RBX) could have been removed by providing the EPC page's address in the LINADDR field.

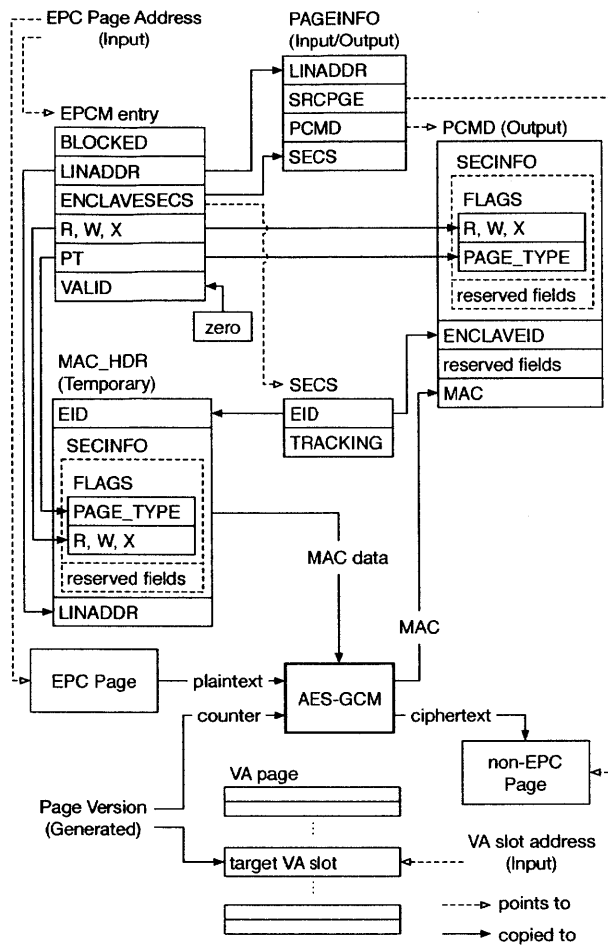


Figure C-17: The data flow of the EWB instruction that evicts an EPC page. The page’s content is encrypted in a non-EPC RAM page. A nonce is created and saved in an empty slot inside a VA page. The page’s EPCM metadata and a MAC are saved in a separate area in non-EPC memory.

Loading an Evicted Page Back into EPC

After an EPC page belonging to an enclave is evicted, any attempt to access the page from enclave code will result in a Page Fault (#PF, § A.8.2). The #PF will cause the logical processor to exit enclave mode via AEX (§ C.4.4), and then invoke the OS kernel’s page fault handler.

Page faults receive special handling from the AEX process. While leaving the enclave, the AEX logic specifically checks if the hardware exception that triggered the AEX was #PF. If that is the case, the AEX implementation clears the least significant 12 bits of the CR2

register, which stores the virtual address whose translation caused a page fault.

In general, the OS kernel's page handler needs to be able to extract the virtual page number (VPN, § A.5.1) from CR2, so that it knows which memory page needs to be loaded back into DRAM. The OS kernel may also be able to use the 12 least significant address bits, which are not part of the VPN, to better predict the application software's memory access patterns. However, unlike the bits that make up the VPN, the bottom 12 bits are not absolutely necessary for the fault handler to carry out its job. Therefore, SGX's AEX implementation clears these 12 bits, in order to limit the amount of information that is learned by the page fault handler.

When the OS page fault handler examines the address in the CR2 register and determines that the faulting address is inside the EPC, it is generally expected to use the ELDU or ELDB instruction to load the evicted page back into the EPC. If the outputs of EWB have been evicted from DRAM to a slower storage medium, the OS kernel will have to read the outputs back into DRAM before invoking ELDU / ELDB.

ELDU and ELDB verify the MAC tag produced by EWB, described in § C.4.5. This prevents the OS kernel from performing the page swapping-based active address translation attack described in § B.7.3.

Eviction Trees

The SGX design allows VA pages to be evicted from the EPC, just like enclave pages. When a VA page is evicted from EPC, all the nonces stored by the VA slots become inaccessible to the processor. Therefore, the evicted pages associated with these nonces cannot be restored by ELDB until the OS loads the VA page back into the EPC.

In other words, an evicted page depends on the VA page storing its nonce, and cannot be loaded back into the EPC until the VA page is reloaded as well. The dependency graph created by this relationship is a forest of `eviction trees`. An eviction tree, shown in Figure C-18, has enclave EPC pages as leaves, and VA pages as inner nodes. A page's parent is the VA page that holds its nonce. Since EWB always outputs a nonce in a VA page, the root node of each eviction tree is always a VA page in the EPC.

A straightforward inductive argument shows that when an OS wishes to load an evicted

enclave page back into the EPC, it needs to load all the VA pages on the path from the eviction tree's root to the leaf corresponding to the enclave page. Therefore, the number of page loads required to satisfy a page fault inside the EPC depends on the shape of the eviction tree that contains the page.

The SGX design leaves the OS in complete control of the shape of the eviction trees. This has no negative impact on security, as the tree shape only impacts the performance of the eviction scheme, and not its correctness.

C.4.6 SGX Enclave Measurement

SGX implements a software attestation scheme that follows the general principles outlined in § B.3. For the purposes of this section, the most relevant principle is that a remote party authenticates an enclave based on its measurement, which is intended to identify the software that is executing inside the enclave. The remote party compares the enclave measurement reported by the trusted hardware with an expected measurement, and only proceeds if the two values match.

§ C.4.3 explains that an SGX enclave is built using the `ECREATE` (§ C.4.3), `EADD` (§ C.4.3) and `EEXTEND` instructions. After the enclave is initialized via `EINIT` (§ C.4.3), the instructions mentioned above cannot be used anymore. As the SGX measurement scheme follows the principles outlined in § B.3.2, the measurement of an SGX enclave is obtained by computing a secure hash (§ B.1.3) over the inputs to the `ECREATE`, `EADD` and `EEXTEND` instructions used to create the enclave and load the initial code and data into its memory. `EINIT` finalizes the hash that represents the enclave's measurement.

Along with the enclave's contents, the enclave author is expected to specify the sequence of instructions that should be used in order to create an enclave whose measurement will match the expected value used by the remote party in the software attestation process. The `.so` and `.dll` dynamically loaded library file formats, which are SGX's intended enclave delivery methods, already include informal specifications for loading algorithms. We expect the informal loading specifications to serve as the starting points for specifications that prescribe the exact sequences of SGX instructions that should be used to create enclaves

from .so and .dll files.

As argued in § B.3.2, an enclave’s measurement is computed using a secure hashing algorithm, so the system software can only build an enclave that matches an expected measurement by following the exact sequence of instructions specified by the enclave’s author.

The SGX design uses the 256-bit SHA-2 [23] secure hash function to compute its measurements. SHA-2 is a block hash function (§ B.1.3) that operates on 64-byte blocks, uses a 32-byte internal state, and produces a 32-byte output. Each enclave’s measurement is stored in the MRENCLAVE field of the enclave’s SECS. The 32-byte field stores the internal state and final output of the 256-bit SHA-2 secure hash function.

Measuring ECREATE

The ECREATE instruction, overviewed in § C.4.3, first initializes the MRENCLAVE field in the newly created SECS using the 256-bit SHA-2 initialization algorithm, and then extends the hash with the 64-byte block depicted in Table C.4.

Offset	Size	Description
0	8	"ECREATE\0"
8	8	SECS.SSAFRAMESIZE (§ C.4.2)
16	8	SECS.SIZE (§ C.4.2)
32	8	32 zero (0) bytes

Table C.4: 64-byte block extended into MRENCLAVE by ECREATE

The enclave’s measurement does not include the BASEADDR field. The omission is intentional, as it allows the system software to load an enclave at any virtual address inside a host process that satisfies the ELRANGE restrictions (§ C.4.2), without changing the enclave’s measurement. This feature can be combined with a compiler that generates position-independent enclave code to obtain relocatable enclaves.

The enclave’s measurement includes the SSAFRAMESIZE field, which guarantees that the SSAs (§ C.4.2) created by AEX and used by EENTER (§ C.4.4) and ERESUME (§ C.4.4) have the size that is expected by the enclave’s author. Leaving this field out of an enclave’s measurement would allow a malicious enclave loader to attempt to attack the enclave’s

security checks by specifying a bigger SSAFRAMESIZE than the enclave's author intended, which could cause the SSA contents written by an AEX to overwrite the enclave's code or data.

Measuring Enclave Attributes

The enclave's measurement does not include the enclave attributes (§ C.4.2), which are specified in the ATTRIBUTES field in the SECS. Instead, it is included directly in the information that is covered by the attestation signature, which will be discussed in § C.4.8.

The SGX software attestation definitely needs to cover the enclave attributes. For example, if XFRM (§ C.4.2, § C.4.2) would not be covered, a malicious enclave loader could attempt to subvert an enclave's security checks by setting XFRM to a value that enables architectural extensions that change the semantics of instructions used by the enclave, but still produces an XSAVE output that fits in SSAFRAMESIZE.

The special treatment applied to the ATTRIBUTES SECS field seems questionable from a security standpoint, as it adds extra complexity to the software attestation verifier, which translates into more opportunities for exploitable bugs. This decision also adds complexity to the SGX software attestation design, which is described in § C.4.8.

The most likely reason why the SGX design decided to go this route, despite the concerns described above, is the wish to be able to use a single measurement to represent an enclave that can take advantage of some architectural extensions, but can also perform its task without them.

Consider, for example, an enclave that performs image processing using a library such as OpenCV, which has routines optimized for SSE and AVX, but also includes generic fallbacks for processors that do not have these features. The enclave's author will likely wish to allow an enclave loader to set bits 1 (SSE) and 2 (AVX) to either true or false. If ATTRIBUTES (and, by extension, XFRM) was a part of the enclave's measurement, the enclave author would have to specify that the enclave has 4 valid measurements. In general, allowing n architectural extensions to be used independently will result in 2^n valid measurements.

Measuring EADD

The *EADD* instruction, described in § C.4.3, extends the SHA-2 hash in MRENCLAVE with the 64-byte block shown in Table C.5.

Offset	Size	Description
0	8	"EADD\0\0\0\0"
8	8	ENCLAVEOFFSET
16	48	SECINFO (first 48 bytes)

Table C.5: 64-byte block extended into MRENCLAVE by *EADD*. The ENCLAVEOFFSET is computed by subtracting the BASEADDR in the enclave's SECS from the LINADDR field in the PAGEINFO structure.

The address included in the measurement is the address where the *EADD*ed page is expected to be mapped in the enclave's virtual address space. This ensures that the system software sets up the enclave's virtual memory layout according to the enclave author's specifications. If a malicious enclave loader attempts to set up the enclave's layout incorrectly, perhaps in order to mount an active address translation attack (§ B.7.2), the loaded enclave's measurement will differ from the measurement expected by the enclave's author.

The virtual address of the newly created page is measured relatively to the start of the enclave's ELRANGE. In other words, the value included in the measurement is LINADDR - BASEADDR. This makes the enclave's measurement invariant to BASEADDR changes, which is desirable for relocatable enclaves. Measuring the relative addresses still preserves all the information about the memory layout inside ELRANGE, and therefore has no negative security impact.

EADD also measures the first 48 bytes of the SECINFO structure (§ C.4.3) provided to *EADD*, which contain the page type (PT) and access permissions (R, W, X) field values used to initialize the page's EPCM entry. By the same argument as above, including these values in the measurement guarantees that the memory layout built by the system software loading the enclave matches the specifications of the enclave author.

The EPCM field values mentioned above take up less than one byte in the SECINFO structure, and the rest of the bytes are reserved and expected to be initialized to zero. This leaves plenty of expansion room for future SGX features.

The most notable omission from Table C.5 is the data used to initialize the newly created EPC page. Therefore, the measurement data contributed by EADD guarantees that the enclave's memory layout will have pages allocated with prescribed access permissions at the desired virtual addresses. However, the measurements don't cover the code or data loaded in these pages.

For example, EADD's measurement data guarantees that an enclave's memory layout consists of three executable pages followed by five writable data pages, but it does not guarantee that any of the code pages contains the code supplied by the enclave's author.

Measuring EEXTEND

The EEXTEND instruction exists solely for the reason of measuring data loaded inside the enclave's EPC pages. The instruction reads in a virtual address, and extends the enclave's measurement hash with the five 64-byte blocks in Table C.6, which effectively guarantee the contents of a 256-byte chunk of data in the enclave's memory.

Offset	Size	Description
0	8	"EEXTEND\0"
8	8	ENCLAVEOFFSET
16	48	48 zero (0) bytes
64	64	bytes 0 - 64 in the chunk
128	64	bytes 64 - 128 in the chunk
192	64	bytes 128 - 192 in the chunk
256	64	bytes 192 - 256 in the chunk

Table C.6: 64-byte blocks extended into MRENCLAVE by EEXTEND. The ENCLAVE-OFFSET is computed by subtracting the BASEADDR in the enclave's SECS from the LINADDR field in the PAGEINFO structure.

Before examining the details of EEXTEND, we note that SGX's security guarantees only hold when the contents of the enclave's key pages is measured. For example, EENTER (§ C.4.4) is only guaranteed to perform controlled jumps inside an enclave's code if the contents of all the Thread Control Structure (TCS, § C.4.2) pages are measured. Otherwise, a malicious enclave loader can change the OENTRY field (§ C.4.2, § C.4.4) in a TCS while building the enclave, and then a malicious OS can use the TCS to perform an arbitrary jump inside en-

clave code. By the same argument, all the enclave's code should be measured by `EEXTEND`. Any code fragment that is not measured can be replaced by a malicious enclave loader.

Given these pitfalls, it is surprising that the SGX design opted to decouple the virtual address space layout measurements done by `EADD` from the memory content measurements done by `EEXTEND`.

At a first pass, it appears that the decoupling only has one benefit, which is the ability to load un-measured user input into an enclave while it is being built. However, this benefit only translates into a small performance improvement, because enclaves can alternatively be designed to copy the user input from untrusted DRAM after being initialized. At the same time, the decoupling opens up the possibility of relying on an enclave that provides no meaningful security guarantees, due to not measuring all the important data via `EEXTEND` calls.

However, the real reason behind the `EADD` / `EEXTEND` separation is hinted at by the `EINIT` pseudo-code in the SDM, which states that the instruction opens an interrupt (§ A.12) window while it performs a computationally intensive RSA signature check. If an interrupt occurs during the check, `EINIT` fails with an error code, and the interrupt is serviced. This very unusual approach for a processor instruction suggests that the SGX implementation was constrained in respect to how much latency its instructions were allowed to add to the interrupt handling process.

In light of the concerns above, it is reasonable to conclude that `EEXTEND` was introduced because measuring an entire page using 256-bit SHA-2 is quite time-consuming, and doing it in `EADD` would have caused the instruction to exceed SGX's latency budget. The need to hit a certain latency goal is a reasonable explanation for the seemingly arbitrary 256-byte chunk size.

The `EADD` / `EEXTEND` separation will not cause security issues if enclaves are authored using the same tools that build today's dynamically loaded modules, which appears to be the workflow targeted by the SGX design. In this workflow, the tools that build enclaves can easily identify the enclave data that needs to be measured.

It is correct and meaningful, from a security perspective, to have the message blocks provided by `EEXTEND` to the hash function include the address of the 256-byte chunk, in

addition to the contents of the data. If the address were not included, a malicious enclave loader could mount the memory mapping attack described in § B.7.2 and illustrated in Figure B-23.

More specifically, the malicious loader would EADD the errorOut page contents at the virtual address intended for disclose, EADD the disclose page contents at the virtual address intended for errorOut, and then EEXTEND the pages in the wrong order. If EEXTEND would not include the address of the data chunk that is measured, the steps above would yield the same measurement as the correctly constructed enclave.

The last aspect of EEXTEND worth analyzing is its support for relocating enclaves. Similarly to EADD, the virtual address measured by EEXTEND is relative to the enclave's BASEADDR. Furthermore, the only SGX structure whose content is expected to be measured by EEXTEND is the TCS. The SGX design has carefully used relative addresses for all the TCS fields that represent enclave addresses, which are OENTRY, OFSBASGX and OGSBASGX.

Measuring EINIT

The EINIT instruction (§ C.4.3) concludes the enclave building process. After EINIT is successfully invoked on an enclave, the enclave's contents are “sealed”, meaning that the system software cannot use the EADD instruction to load code and data into the enclave, and cannot use the EEXTEND instruction to update the enclave's measurement.

EINIT uses the SHA-2 finalization algorithm (§ B.1.3) on the MRENCLAVE field of the enclave's SECS. After EINIT, the field no longer stores the intermediate state of the SHA-2 algorithm, and instead stores the final output of the secure hash function. This value remains constant after EINIT completes, and is included in the attestation signature produced by the SGX software attestation process.

C.4.7 SGX Enclave Versioning Support

The software attestation model (§ B.3) introduced by the Trusted Platform Module (§ 2.4) relies on a measurement (§ C.4.6), which is essentially a content hash, to identify the

software inside a container. The downside of using content hashes for identity is that there is no relation between the identities of containers that hold different versions of the same software.

In practice, it is highly desirable for systems based on secure containers to handle software updates without having access to the remote party in the initial software attestation process. This entails having the ability to migrate secrets between the container that has the old version of the software and the container that has the updated version. This requirement translates into a need for a separate identity system that can recognize the relationship between two versions of the same software.

SGX supports the migration of secrets between enclaves that represent different versions of the same software, as shown in Figure C-19.

The secret migration feature relies on a one-level certificate hierarchy (§ B.2.1), where each enclave author is a Certificate Authority, and each enclave receives a certificate from its author. These certificates must be formatted as Signature Structures (SIGSTRUCT), which are described in § C.4.7. The information in these certificates is the basis for an enclave identity scheme, presented in § C.4.7, which can recognize the relationship between different versions of the same software.

The `EINIT` instruction (§ C.4.3) examines the target enclave's certificate and uses the information in it to populate the SECS (§ C.4.1) fields that describe the enclave's certificate-based identity. This process is summarized in § C.4.7.

Last, the actual secret migration process is based on the key derivation service implemented by the `EGETKEY` instruction, which is described in § C.4.7. The sending enclave uses the `EGETKEY` instruction to obtain a symmetric key (§ B.1.1) based on its identity, encrypts its secrets with the key, and hands off the encrypted secrets to the untrusted system software. The receiving enclave passes the sending enclave's identity to `EGETKEY`, obtains the same symmetric key as above, and uses the key to decrypt the secrets received from system software.

The symmetric key obtained from `EGETKEY` can be used in conjunction with cryptographic primitives that protect the privacy (§ B.1.2) and integrity (§ B.1.3) of an enclave's secrets while they are migrated to another enclave by the untrusted system software. How-

ever, symmetric keys alone cannot be used to provide freshness guarantees (§ B.1), so secret migration is subject to replay attacks. This is acceptable when the secrets being migrated are immutable, such as when the secrets are encryption keys obtained via software attestation

Enclave Certificates

The SGX design requires each enclave to have a certificate issued by its author. This requirement is enforced by `EINIT` (§ C.4.3), which refuses to operate on enclaves without valid certificates.

The SGX implementation consumes certificates formatted as *Signature Structures* (SIGSTRUCT), which are intended to be generated by an enclave building toolchain, as shown in Figure C-20.

A SIGSTRUCT certificate consists of metadata fields, the most interesting of which are presented in Table C.7, and an RSA signature that guarantees the authenticity of the metadata, formatted as shown in Table C.8. The semantics of the fields will be revealed in the following sections.

Field	Bytes	Description
ENCLAVEHASH	32	Must equal the enclave's measurement (§ C.4.6).
ISVPRODID	32	Differentiates modules signed by the same public key.
ISVSVN	32	Differentiates versions of the same module.
VENDOR	4	Differentiates Intel enclaves.
ATTRIBUTES	16	Constrains the enclave's attributes.
ATTRIBUTEMASK	16	Constrains the enclave's attributes.

Table C.7: A subset of the metadata fields in a SIGSTRUCT enclave certificate

Field	Bytes	Description
MODULUS	384	RSA key modulus
EXPONENT	4	RSA key public exponent
SIGNATURE	384	RSA signature (See § C.5.5)
Q1	384	Simplifies RSA signature verification. (See § C.5.5)
Q2	384	Simplifies RSA signature verification. (See § C.5.5)

Table C.8: The format of the RSA signature used in a SIGSTRUCT enclave certificate

The enclave certificates must be signed by RSA signatures (§ B.1.3) that follow the method described in RFC 3447 [113], using 256-bit SHA-2 [23] as the hash function that

reduces the input size, and the padding method described in PKCS #1 v1.5 [114], which is illustrated in Figure B-14.

The SGX implementation only supports 3072-bit RSA keys whose public exponent is 3. The key size is likely chosen to meet FIPS' recommendation [22], which makes SGX eligible for use in U.S. government applications. The public exponent 3 affords a simplified signature verification algorithm, which is discussed in § C.5.5. The simplified algorithm also requires the fields Q1 and Q2 in the RSA signature, which are also described in § C.5.5.

Certificate-Based Enclave Identity

An enclave's identity is determined by three fields in its certificate (§ C.4.7): the modulus of the RSA key used to sign the certificate (MODULUS), the enclave's product ID (ISVPRODID) and the security version number (ISVSVN).

The public RSA key used to issue a certificate identifies the enclave's author. All RSA keys used to issue enclave certificates must have the public exponent set to 3, so they are only differentiated by their moduli. SGX does not use the entire modulus of a key, but rather a 256-bit SHA-2 hash of the modulus. This is called a *signer measurement* (MRSIGNER), to parallel the name of *enclave measurement* (MRENCLAVE) for the SHA-2 hash that identifies an enclave's contents.

The SGX implementation relies on a hard-coded MRSIGNER value to recognize certificates issued by Intel. Enclaves that have an Intel-issued certificate can receive additional privileges, which are discussed in § C.4.8.

An enclave author can use the same RSA key to issue certificates for enclaves that represent different software modules. Each module is identified by a unique Product ID (ISVPRODID) value. Conversely, all the enclaves whose certificates have the same ISVPRODID and are issued by the same RSA key (and therefore have the same MRENCLAVE) are assumed to represent different versions of the same software module. Enclaves whose certificates are signed by different keys are always assumed to contain different software modules.

Enclaves that represent different versions of a module can have different security version numbers (SVN). The SGX design disallows the migration of secrets from an enclave with a

higher SVN to an enclave with a lower SVN. This restriction is intended to assist with the distribution of security patches, as follows.

If a security vulnerability is discovered in an enclave, the author can release a fixed version with a higher SVN. As users upgrade, SGX will facilitate the migration of secrets from the vulnerable version of the enclave to the fixed version. Once a user's secrets have migrated, the SVN restrictions in SGX will deflect any attack based on building the vulnerable enclave version and using it to read the migrated secrets.

Software upgrades that add functionality should not be accompanied by an SVN increase, as SGX allows secrets to be migrated freely between enclaves with matching SVN values. As explained above, a software module's SVN should only be incremented when a security vulnerability is found. SIGSTRUCT only allocates 2 bytes to the ISVSVN field, which translates to 65,536 possible SVN values. This space can be exhausted if a large team (incorrectly) sets up a continuous build system to allocate a new SVN for every software build that it produces, and each code change triggers a build.

CPU Security Version Numbers

The SGX implementation itself has a security version number (CPUSVN), which is used in the key derivation process implemented [142] by EGETKEY, in addition to the enclave's identity information. CPUSVN is a 128-bit value that, according to the SDM, reflects the processor's microcode update version.

The SDM does not describe the structure of CPUSVN, but it states that comparing CPUSVN values using integer comparison is not meaningful, and that only some CPUSVN values are valid. Furthermore, CPUSVNs admit an ordering relationship that has the same semantics as the ordering relationship between enclave SVNs. Specifically, an SGX implementation will consider all SGX implementations with lower SVNs to be compromised due to security vulnerabilities, and will not trust them.

An SGX patent [142] discloses that CPUSVN is a concatenation of small integers representing the SVNs of the various components that make up SGX's implementation. This structure is consistent with all the statements made in the SDM.

Establishing an Enclave's Identity

When the `EINIT` (§ C.4.3) instruction prepares an enclave for code execution, it also sets the `SECS` (§ C.4.1) fields that make up the enclave's certificate-based identity, as shown in Figure C-21.

`EINIT` requires the virtual address of the `SIGSTRUCT` certificate issued to the enclave, and uses the information in the certificate to initialize the certificate-based identity information in the enclave's `SECS`. Before using the information in the certificate, `EINIT` first verifies its RSA signature. The `SIGSTRUCT` fields `Q1` and `Q2`, along with the RSA exponent `3`, facilitate a simplified verification algorithm, which is discussed in § C.5.5.

If the `SIGSTRUCT` certificate is found to be properly signed, `EINIT` follows the steps discussed in the following few paragraphs to ensure that the certificate was issued to the enclave that is being initialized. Once the checks have completed, `EINIT` computes `MRSIGNER`, the 256-bit SHA-2 hash of the `MODULUS` field in the `SIGSTRUCT`, and writes it into the enclave's `SECS`. `EINIT` also copies the `ISVPRODID` and `ISVSVN` fields from `SIGSTRUCT` into the enclave's `SECS`. As explained in § C.4.7, these fields make up the enclave's certificate-based identity.

After verifying the RSA signature in `SIGSTRUCT`, `EINIT` copies the signature's padding into the `PADDING` field in the enclave's `SECS`. The PKCS #1 v1.5 padding scheme, outlined in Figure B-14, does not involve randomness, so `PADDING` should have the same value for all enclaves.

`EINIT` performs a few checks to make sure that the enclave undergoing initialization was indeed authorized by the provided `SIGSTRUCT` certificate. The most obvious check involves making sure that the `MRENCLAVE` value in `SIGSTRUCT` equals the enclave's measurement, which is stored in the `MRENCLAVE` field in the enclave's `SECS`.

However, `MRENCLAVE` does not cover the enclave's attributes, which are stored in the `ATTRIBUTES` field of the `SECS`. As discussed in § C.4.6, omitting `ATTRIBUTES` from `MRENCLAVE` facilitates writing enclaves that have optimized implementations that can use architectural extensions when present, and also have fallback implementations that work on CPUs without the extensions. Such enclaves can execute correctly when built with

a variety of values in the XFRM (§ C.4.2, § C.4.2) attribute. At the same time, allowing system software to use arbitrary values in the ATTRIBUTES field would compromise SGX's security guarantees.

When an enclave uses software attestation (§ B.3) to gain access to secrets, the ATTRIBUTES value used to build it is included in the SGX attestation signature (§ C.4.8). This gives the remote party in the attestation process the opportunity to reject an enclave built with an undesirable ATTRIBUTES value. However, when secrets are obtained using the migration process facilitated by certificate-based identities, there is no remote party that can check the enclave's attributes.

The SGX design solves this problem by having enclave authors convey the set of acceptable attribute values for an enclave in the ATTRIBUTES and ATTRIBUTESMASK fields of the SIGSTRUCT certificate issued for the enclave. EINIT will refuse to initialize an enclave using a SIGSTRUCT if the bitwise AND between the ATTRIBUTES field in the enclave's SECS and the ATTRIBUTESMASK field in the SIGSTRUCT does not equal the SIGSTRUCT's ATTRIBUTES field. This check prevents enclaves with undesirable attributes from obtaining and potentially leaking secrets using the migration process.

Any enclave author can use SIGSTRUCT to request any of the bits in an enclave's ATTRIBUTES field to be zero. However, certain bits can only be set to one for enclaves that are signed by Intel. EINIT has a mask of restricted ATTRIBUTES bits, discussed in § C.4.8. The EINIT implementation contains a hard-coded MRSIGNER value that is used to identify Intel's privileged enclaves, and only allows privileged enclaves to be built with an ATTRIBUTES value that matches any of the bits in the restricted mask. This check is essential to the security of the SGX software attestation process, which is described in § C.4.8.

Last, EINIT also inspects the VENDOR field in SIGSTRUCT. The SDM description of the VENDOR field in the section dedicated to SIGSTRUCT suggests that the field is essentially used to distinguish between special enclaves signed by Intel, which use a VENDOR value of 0x8086, and everyone else's enclaves, which should use a VENDOR value of zero. However, the EINIT pseudocode seems to imply that the SGX implementation only checks that VENDOR is either zero or 0x8086.

Enclave Key Derivation

SGX's secret migration mechanism is based on the symmetric key derivation service that is offered to enclaves by the `EGETKEY` instruction, illustrated in Figure C-22.

The keys produced by `EGETKEY` are derived based on the identity information in the current enclave's SECS and on two secrets stored in secure hardware inside the SGX-enabled processor. One of the secrets is the input to a largely undocumented series of transformations that yields the symmetric key for the cryptographic primitive underlying the key derivation process. The other secret, referred to as the `CR_SEAL_FUSES` in the SDM, is one of the pieces of information used in the key derivation material.

The SDM does not specify the key derivation algorithm, but the SGX patents [142, 112] disclose that the keys are derived using the method described in FIPS SP 800-108 [36] using AES-CMAC [50] as a Pseudo-Random Function (PRF). The same patents state that the secrets used for key derivation are stored in the CPU's e-fuses, which is confirmed by the ISCA 2015 SGX tutorial [106].

This additional information implies that all `EGETKEY` invocations that use the same key derivation material will result in the same key, even across CPU power cycles. Furthermore, it is impossible for an adversary to obtain the key produced from a specific key derivation material without access to the secret stored in the CPU's e-fuses. SGX's key hierarchy is further described in § C.4.8.

The following paragraphs discuss the pieces of data used in the key derivation material, which are selected by the Key Request (`KEYREQUEST`) structure shown in in Table C.9,

Field	Bytes	Description
<code>KEYNAME</code>	2	The desired key type; secret migration uses Seal keys
<code>KEYPOLICY</code>	2	The identity information (<code>MRENCLAVE</code> and/or <code>MR-SIGNER</code>)
<code>ISVSVN</code>	2	The enclave SVN used in derivation
<code>CPUSVN</code>	16	SGX implementation SVN used in derivation
<code>ATTRIBUTEMASK</code>	16	Selects enclave attributes
<code>KEYID</code>	32	Random bytes

Table C.9: A subset of the fields in the `KEYREQUEST` structure

The `KEYNAME` field in `KEYREQUEST` always participates in the key generation

material. It indicates the type of the key to be generated. While the SGX design defines a few key types, the secret migration feature always uses Seal keys. The other key types are used by the SGX software attestation process, which will be outlined in § C.4.8.

The `KEYPOLICY` field in `KEYREQUEST` has two flags that indicate if the `MRENCLAVE` and `MRSIGNER` fields in the enclave's `SECS` will be used for key derivation. Although the field admits 4 values, only two seem to make sense, as argued below.

Setting the `MRENCLAVE` flag in `KEYPOLICY` ties the derived key to the current enclave's measurement, which reflects its contents. No other enclave will be able to obtain the same key. This is useful when the derived key is used to encrypt enclave secrets so they can be stored by system software in non-volatile memory, and thus survive power cycles.

If the `MRSIGNER` flag in `KEYPOLICY` is set, the derived key is tied to the public RSA key that issued the enclave's certificate. Therefore, other enclaves issued by the same author may be able to obtain the same key, subject to the restrictions below. This is the only `KEYPOLICY` value that allows for secret migration.

It makes little sense to have no flag set in `KEYPOLICY`. In this case, the derived key has no useful security property, as it can be obtained by other enclaves that are completely unrelated to the enclave invoking `EGETKEY`. Conversely, setting both flags is redundant, as setting `MRENCLAVE` alone will cause the derived key to be tied to the current enclave, which is the strictest possible policy.

The `KEYREQUEST` structure specifies the enclave SVN (`ISVSVN`, § C.4.7) and SGX implementation SVN (`CPUSVN`, § C.4.7) that will be used in the key derivation process. However, `EGETKEY` will reject the derivation request and produce an error code if the desired enclave SVN is greater than the current enclave's SVN, or if the desired SGX implementation's SVN is greater than the current implementation's SVN.

The SVN restrictions prevent the migration of secrets from enclaves with higher SVNs to enclaves with lower SVNs, or from SGX implementations with higher SVNs to implementations with lower SVNs. § C.4.7 argues that the SVN restrictions can reduce the impact of security vulnerabilities in enclaves and in SGX's implementation.

`EGETKEY` always uses the `ISVPRODID` value from the current enclave's `SECS` for key derivation. It follows that secrets can never flow between enclaves whose `SIGSTRUCT`

certificates assign them different Product IDs.

Similarly, the key derivation material always includes the value of an 128-bit Owner Epoch (OWNEREPOCH) SGX configuration register. This register is intended to be set by the computer's firmware to a secret generated once and stored in non-volatile memory. Before the computer changes ownership, the old owner can clear the OWNEREPOCH from non-volatile memory, making it impossible for the new owner to decrypt any enclave secrets that may be left on the computer.

Due to the cryptographic properties of the key derivation process, outside observers cannot correlate keys derived using different OWNEREPOCH values. This makes it impossible for software developers to use the EGETKEY-derived keys described in this section to track a processor as it changes owners.

The EGETKEY derivation material also includes a 256-bit value supplied by the enclave, in the KEYID field. This makes it possible for an enclave to generate a collection of keys from EGETKEY, instead of a single key. The SDM states that KEYID should be populated with a random number, and is intended to help prevent key wear-out.

Last, the key derivation material includes the bitwise AND of the ATTRIBUTES (§ C.4.2) field in the enclave's SECS and the ATTRIBUTESMASK field in the KEYREQUEST structure. The mask has the effect of removing some of the ATTRIBUTES bits from the key derivation material, making it possible to migrate secrets between enclaves with different attributes. § C.4.6 and § C.4.7 explain the need for this feature, as well as its security implications.

Before adding the masked attributes value to the key generation material, the EGETKEY implementation forces the mask bits corresponding to the INIT and DEBUG attributes (§ C.4.2) to be set. From a practical standpoint, this means that secrets will never be migrated between enclaves that support debugging and production enclaves.

Without this restriction, it would be unsafe for an enclave author to use the same RSA key to issue certificates to both debugging and production enclaves. Debugging enclaves receive no integrity guarantees from SGX, so it is possible for an attacker to modify the code inside a debugging enclave in a way that causes it to disclose any secrets that it has access to.

C.4.8 SGX Software Attestation

The software attestation scheme implemented by SGX follows the principles outlined in § B.3. An SGX-enabled processor computes a measurement of the code and data that is loaded in each enclave, which is similar to the measurement computed by the TPM (§ 2.4). The software inside an enclave can start a process that results in an SGX attestation signature, which includes the enclave's measurement and an enclave message.

The cryptographic primitive used in SGX's attestation signature is too complex to be implemented in hardware, so the signing process is performed by a privileged *Quoting Enclave*, which is issued by Intel, and can access the SGX attestation key. This enclave is discussed in § C.4.8.

Pushing the signing functionality into the Quoting Enclave creates the need for a secure communication path between an enclave undergoing software attestation and the Quoting Enclave. The SGX design solves this problem with a local attestation mechanism that can be used by an enclave to prove its identity to any other enclave hosted by the same SGX-enabled CPU. This scheme, described in § C.4.8, is implemented by the EREPORT instruction.

The SGX attestation key used by the Quoting Enclave does not exist at the time SGX-enabled processors leave the factory. The attestation key is provisioned later, using a largely undocumented process that is known to involve at least one other enclave issued by Intel, and two special EGETKEY (§ C.4.7) key types. The publicly available details of this process are summarized in § C.4.8.

The SGX Launch Enclave and EINITTOKEN structure will be discussed in § C.4.9.

Local Attestation

An enclave proves its identity to another *target enclave* via the EREPORT instruction shown in Figure C-24. The SGX instruction produces an attestation *Report* (REPORT) that cryptographically binds a message supplied by the enclave with the enclave's measurement-based (§ C.4.6) and certificate-based (§ C.4.7) identities. The cryptographic binding is accomplished by a MAC tag (§ B.1.3) computed using a symmetric key that is only shared

between the target enclave and the SGX implementation.

The EREPORT instruction reads the current enclave's identity information from the enclave's SECS (§ C.4.1), and uses it to populate the REPORT structure. Specifically, EREPORT copies the SECS fields indicating the enclave's measurement (MRENCLAVE), certificate-based identity (MRSIGNER, ISVPRODID, ISVSVN), and attributes (ATTRIBUTES). The attestation report also includes the SVN of the SGX implementation (CPUSVN) and a 64-byte (512-bit) message supplied by the enclave.

The target enclave that receives the attestation report can convince itself of the report's authenticity as shown in Figure C-25. The report's authenticity proof is its MAC tag. The key required to verify the MAC can only be obtained by the target enclave, by asking EGETKEY (§ C.4.7) to derive a Report key. The SDM states that the MAC tag is computed using a block cipher-based MAC (CMAC, [50]), but stops short of specifying the underlying cipher. One of the SGX papers [15] states that the CMAC is based on 128-bit AES.

The Report key returned by EGETKEY is derived from a secret embedded in the processor (§ C.4.7), and the key material includes the target enclave's measurement. The target enclave can be assured that the MAC tag in the report was produced by the SGX implementation, for the following reasons. The cryptographic properties of the underlying key derivation and MAC algorithms ensure that only the SGX implementation can produce the MAC tag, as it is the only entity that can access the processor's secret, and it would be impossible for an attacker to derive the Report key without knowing the processor's secret. The SGX design guarantees that the key produced by EGETKEY depends on the calling enclave's measurement, so only the target enclave can obtain the key used to produce the MAC tag in the report.

EREPORT uses the same key derivation process as EGETKEY does when invoked with KEYNAME set to the value associated with Report keys. For this reason, EREPORT requires the virtual address of a *Report Target Info* (TARGETINFO) structure that contains the measurement-based identity and attributes of the target enclave.

When deriving a Report key, EGETKEY behaves slightly differently than it does in the case of seal keys, as shown in Figure C-25. The key generation material never includes the fields corresponding to the enclave's certificate-based identity (MRSIGNER, ISVPRODID,

ISVSVN), and the KEYPOLICY field in the KEYREQUEST structure is ignored. It follows that the report can only be verified by the target enclave.

Furthermore, the SGX implementation's SVN (CPUSVN) value used for key generation is determined by the current CPUSVN, instead of being read from the Key Request structure. Therefore, SGX implementation upgrades that increase the CPUSVN invalidate all outstanding reports. Given that CPUSVN increases are associated with security fixes, the argument in § C.4.7 suggests that this restriction may reduce the impact of vulnerabilities in the SGX implementation.

Last, EREPORT sets the KEYID field in the key generation material to the contents of an SGX configuration register (CR_REPORT_KEYID) that is initialized with a random value when SGX is initialized. The KEYID value is also saved in the attestation report, but it is not covered by the MAC tag.

Remote Attestation

The SDM paints a complete picture of the local attestation mechanism that was described in § C.4.8. In comparison, the remote attestation process, which includes the Quoting Enclave and the underlying keys, is shrouded in mystery. This section presents the information that can be gleaned from the SDM, from one [15] of the SGX papers, and from the ISCA 2015 SGX tutorial [106].

SGX's software attestation scheme, which is illustrated in Figure C-26, relies on a key generation facility and on a provisioning service, both operated by Intel.

During the manufacturing process, an SGX-enabled processor communicates with Intel's key generation facility, and has two secrets burned into e-fuses, which are a one-time programmable storage medium that can be economically included on a high-performance chip's die. We shall refer to the secrets stored in e-fuses as the *Provisioning Secret* and the *Seal Secret*.

The Provisioning Secret is the main input to a largely undocumented process that outputs the SGX master derivation key used by EGETKEY, which was referenced in Figures C-22, C-23, C-24, and C-25.

The Seal Secret is not exposed to software by any of the architectural mechanisms

documented in the SDM. The secret is only accessed when it is included in the material used by the key derivation process implemented by EGETKEY (§ C.4.7). The pseudocode in the SDM uses the CR_SEAL_FUSES register name to refer to the Seal Secret.

The names “Seal Secret” and “Provisioning Secret” deviate from Intel’s official documents, which confusingly use the “Seal Key” and “Provisioning Key” names to refer to both secrets stored in e-fuses and keys derived by EGETKEY.

The SDM briefly describes the keys produced by EGETKEY, but no official documentation explicitly describes the secrets in e-fuses. The description below is the only interpretation of all the public information sources that is consistent with all the SDM’s statements about key derivation.

The Provisioning Secret is generated at the key generation facility, where it is burned into the processor’s e-fuses and stored in the database used by Intel’s provisioning service. The Seal Secret is generated inside the processor chip, and therefore is not known to Intel. This approach has the benefit that an attacker who compromises Intel’s facilities cannot derive most keys produced by EGETKEY, even if the attacker also compromises a victim’s firmware and obtains the OWNEREPOCH (§ C.4.7) value. These keys include the Seal keys (§ C.4.7) and Report keys (§ C.4.8) introduced in previous sections.

The only documented exception to the reasoning above is the *Provisioning key*, which is effectively a shared secret between the SGX-enabled processor and Intel’s provisioning service. Intel has to be able to derive this key, so the derivation material does not include the Seal Secret or the OWNEREPOCH value, as shown in Figure C-27.

EGETKEY derives the Provisioning key using the current enclave’s certificate-based identity (MRSIGNER, ISVPRODID, ISVSVN) and the SGX implementation’s SVN (CPUSVN). This approach has a few desirable security properties. First, Intel’s provisioning service can be assured that it is authenticating a Provisioning Enclave signed by Intel. Second, the provisioning service can use the CPUSVN value to reject SGX implementations with known security vulnerabilities. Third, this design admits multiple mutually distrusting provisioning services.

EGETKEY only derives Provisioning keys for enclaves whose PROVISIONKEY attribute is set to true. § C.4.9 argues that this mechanism is sufficient to protect the computer owner

from a malicious software provider that attempts to use Provisioning keys to track a CPU chip across OWNEREPOCH changes.

After the Provisioning Enclave obtains a Provisioning key, it uses the key to authenticate itself to Intel's provisioning service. Once the provisioning service is convinced that it is communicating to a trusted Provisioning enclave in the secure environment provided by a SGX-enabled processor, the service generates an *Attestation Key* and sends it to the Provisioning Enclave. The enclave then encrypts the *Attestation Key* using a *Provisioning Seal key*, and hands off the encrypted key to the system software for storage.

Provisioning Seal keys, are the last publicly documented type of special keys derived by EGETKEY, using the process illustrated in Figure C-28. As their name suggests, Provisioning Seal keys are conceptually similar to the Seal Keys (§ C.4.7) used to migrate secrets between enclaves.

The defining feature of Provisioning Seal keys is that they are not based on the OWNEREPOCH value, so they survive computer ownership changes. Since Provisioning Seal keys can be used to track a CPU chip, their use is gated on the PROVISIONKEY attribute, which has the same semantics as for Provisioning keys.

Like Provisioning keys, Seal keys are based on the current enclave's certificate-based identity (MRSIGNER, ISVPROD, ISVSVN), so the *Attestation Key* encrypted by Intel's Provisioning Enclave can only be decrypted by another enclave signed with the same Intel RSA key. However, unlike Provisioning keys, the Provisioning Seal keys are based on the Seal Secret in the processor's e-fuses, so they cannot be derived by Intel.

When considered independently from the rest of the SGX design, Provisioning Seal keys have desirable security properties. The main benefit of these keys is that when a computer with an SGX-enabled processor exchanges owners, it does not need to undergo the provisioning process again, so Intel does not need to be aware of the ownership change. The privacy issue that stems from not using OWNEREPOCH was already introduced by Provisioning keys, and is mitigated using the access control scheme based on the PROVISIONKEY attribute that will be discussed in § C.4.9.

Similarly to the Seal key derivation process, both the Provisioning and Provisioning Seal keys depend on the bitwise AND of the ATTRIBUTES (§ C.4.2) field in the enclave's SECS

and the ATTRIBUTESMASK field in the KEYREQUEST structure. While most attributes can be masked away, the DEBUG and INIT attributes are always used for key derivation.

This dependency makes it safe for Intel to use its production RSA key to issue certificates for Provisioning or Quoting Enclaves with debugging features enabled. Without the forced dependency on the DEBUG attribute, using the production Intel signing key on a single debug Provisioning or Quoting Enclave could invalidate SGX's security guarantees on all the CPU chips whose attestation-related enclaves are signed by the same key. Concretely, if the issued SIGSTRUCT would be leaked, any attacker could build a debugging Provisioning or Quoting enclave, use the SGX debugging features to modify the code inside it, and extract the 128-bit Provisioning key used to authenticated the CPU to Intel's provisioning service.

After the provisioning steps above have been completed, the Quoting Enclave can be invoked to perform SGX's software attestation. This enclave receives local attestation reports (§ C.4.8) and verifies them using the Report keys generated by EGETKEY. The Quoting Enclave then obtains the Provisioning Seal Key from EGETKEY and uses it to decrypt the Attestation Key, which is received from system software. Last, the enclave replaces the MAC in the local attestation report with an *Attestation Signature* produced with the Attestation Key.

The SGX patents state that the name "Quoting Enclave" was chosen as a reference to the TPM (§ 2.4)'s quoting feature, which is used to perform software attestation on a TPM-based system.

The Attestation Key uses Intel's *Enhanced Privacy ID* (EPID) cryptosystem [28], which is a group signature scheme that is intended to preserve the anonymity of the signers. Intel's key provisioning service is the issuer in the EPID scheme, so it publishes the Group Public Key, while securely storing the Master Issuing Key. After a Provisioning Enclave authenticates itself to the provisioning service, it generates an EPID Member Private Key, which serves as the Attestation Key, and executes the EPID Join protocol to join the group. Later, the Quoting Enclave uses the EPID Member Private Key to produce Attestation Signatures.

The Provisioning Secret stored in the e-fuses of each SGX-enabled processor can be used by Intel to trace individual chips when a Provisioning Enclave authenticates itself to

the provisioning service. However, if the EPID Join protocol is blinded, Intel's provisioning service cannot trace an Attestation Signature to a specific Attestation Key, so Intel cannot trace Attestation Signatures to individual chips.

Of course, the security properties of the description above hinge on the correctness of the proofs behind the EPID scheme. Analyzing the correctness of such cryptographic schemes is beyond the scope of this work, so we defer the analysis of EPID to the crypto research community.

C.4.9 SGX Enclave Launch Control

The SGX design includes a launch control process, which introduces an unnecessary approval step that is required before running most enclaves on a computer. The approval decision is made by the *Launch Enclave* (LE), which is an enclave issued by Intel that gets to approve every other enclave before it is initialized by *EINIT* (§ C.4.3). The officially documented information about this approval process is discussed in § C.4.9.

The SGX patents [142, 112] disclose in no uncertain terms that the Launch Enclave was introduced to ensure that each enclave's author has a business relationship with Intel, and implements a software licensing system. § C.4.9 briefly discusses the implications, should this turn out to be true.

The remainder of the section argues that the Launch Enclave should be removed from the SGX design. § C.4.9 explains that the LE is not required to enforce the computer owner's launch control policy, and concludes that the LE is only meaningful if it enforces a policy that is detrimental to the computer owner. § C.4.9 debunks the myth that an enclave can host malware, which is likely to be used to justify the LE. § C.4.9 argues that Anti-Virus (AV) software is not fundamentally incompatible with enclaves, further disproving the theory that Intel needs to actively police the software that runs inside enclaves.

Enclave Attributes Access Control

The SGX design requires that all enclaves be vetted by a Launch Enclave (LE), which is only briefly mentioned in Intel's official documentation. Neither its behavior nor its interface

with the system software is specified. We speculate that Intel has not been forthcoming about the LE because of its role in enforcing software licensing, which will be discussed in § C.4.9. This section abstracts away the licensing aspect and assumes that the LE enforces a black-box Launch Control Policy.

The LE approves an enclave by issuing an *EINIT Token* (EINITTOKEN), using the process illustrated in Figure C-29. The EINITTOKEN structure contains the approved enclave's measurement-based (§ C.4.6) and certificate-based (§ C.4.7) identities, just like a local attestation REPORT (§ C.4.8). This token is inspected by EINIT (§ C.4.3), which refuses to initialize enclaves with incorrect tokens.

While an EINIT token is handled by untrusted system software, its integrity is protected by a MAC tag (§ B.1.3) that is computed using a *Launch Key* obtained from EGETKEY. The EINIT implementation follows the same key derivation process as EGETKEY to convince itself that the EINITTOKEN provided to it was indeed generated by an LE that had access to the Launch Key.

The SDM does not document the MAC algorithm used to confer integrity guarantees to the EINITTOKEN structure. However, the EINIT pseudocode verifies the token's MAC tag using the same function that the *EREPORT* pseudocode uses to create the REPORT structure's MAC tag. It follows that the reasoning in § C.4.8 can be reused to conclude that EINITTOKEN structures are MACed using AES-CMAC with 128-bit keys.

The EGETKEY instruction only derives the Launch Key for enclaves that have the LAUNCHKEY attribute set to true. The Launch Key is derived using the same process as the Seal Key (§ C.4.7). The derivation material includes the current enclave's versioning information (ISVPRODID and ISVSVN) but it does not include the main fields that convey an enclave's identity, which are MRSIGNER and MRENCLAVE. The rest of the derivation material follows the same rules as the material used for Seal Keys.

The EINITTTOKEN structure contains the identities of the approved enclave (MRENCLAVE and MRSIGNER) and the approved enclave attributes (ATTRIBUTES). The token also includes the information used for the Launch Key derivation, which includes the LE's Product ID (ISVPRODIDLE), SVN (ISVSVNLE), and the bitwise AND between the LE's ATTRIBUTES and the ATTRIBUTEMASK used in the KEYREQUEST (MASKEDAT-

TRIBUTESLE).

The EINITOKEN information used to derive the Launch Key can also be used by EINIT for damage control, e.g. to reject tokens issued by Launch Enclaves with known security vulnerabilities. The reference pseudocode supplied in the SDM states that EINIT checks the DEBUG bit in the MASKEDATTRIBUTESLE field, and will not initialize a production enclave using a token issued by a debugging LE. It is worth noting that MASKEDATTRIBUTESLE is guaranteed to include the LE's DEBUG attribute, because EGETKEY forces the DEBUG attribute's bit in the attributes mask to 1 (§ C.4.7).

The check described above make it safe for Intel to supply SGX enclave developers with a debugging LE that has its DEBUG attribute set, and performs minimal or no security checks before issuing an EINITOKEN. The DEBUG attribute disables SGX's integrity protection, so the only purpose of the security checks performed in the debug LE would be to help enclave development by mimicking its production counterpart. The debugging LE can only be used to launch any enclave with the DEBUG attribute set, so it does not undermining Intel's ability to enforce a Launch Control Policy on production enclaves.

The enclave attributes access control system described above relies on the LE to reject initialization requests that set privileged attributes such as PROVISIONKEY on unauthorized enclaves. However, the LE cannot vet itself, as there will be no LE available when the LE itself needs to be initialized. Therefore, the Launch Key access restrictions are implemented in hardware.

EINIT accepts an EINITOKEN whose VALID bit is set to zero, if the enclave's MRSIGNER (§ C.4.7) equals a hard-coded value that corresponds to an Intel public key. For all other enclave authors, an invalid EINIT token causes EINIT to reject the enclave and produce an error code.

This exemption to the token verification policy provides a way to bootstrap the enclave attributes access control system, namely using a zeroed out EINITOKEN to initialize the Launch Enclave. At the same time, the cryptographic primitives behind the MRSIGNER check guarantee that only Intel-provided enclaves will be able to bypass the attribute checks. This does not change SGX's security properties because Intel is already a trusted party, as it is responsible for generating the Provisioning Keys and Attestation Keys used by software

attestation (§ C.4.8).

Curiously, the `EINIT` pseudocode in the SDM states that the instruction enforces an additional restriction, which is that all enclaves with the `LAUNCHKEY` attribute must have its certificate issued by the same Intel public key that is used to bypass the `EINITTTOKEN` checks. This restriction appears to be redundant, as the same restriction could be enforced in the Launch Enclave.

Licensing

The SGX patents [142, 112] disclose that `EINIT` Tokens and the Launch Enclave (§ C.4.9) were introduced to verify that the `SIGSTRUCT` certificates associated with production enclaves are issued by enclave authors who have a business relationship with Intel. In other words, the Launch Enclave is intended to be **an enclave licensing mechanism that allows Intel to force itself as an intermediary in the distribution of all enclave software.**

The SGX patents are likely to represent an early version of the SGX design, due to the lengthy timelines associated with patent application approval. In light of this consideration, we cannot make any claims about Intel's current plans. However, given that we know for sure that Intel considered enclave licensing at some point, we briefly discuss the implications of implementing such a licensing plan.

Intel has a near-monopoly on desktop and server-class processors, and being able to decide which software vendors are allowed to use SGX can effectively put Intel in a position to decide winners and losers in many software markets.

Assuming SGX reaches widespread adoption, this issue is the software security equivalent to the Net Neutrality debates that have pitted the software industry against telecommunication giants. Given that virtually all competent software development companies have argued that losing Net Neutrality will stifle innovation, it is fairly safe to assume that Intel's ability to regulate access to SGX will also stifle innovation.

Furthermore, from a historical perspective, the enclave licensing scheme described in the SGX patents is very similar to Verified Boot, which was briefly discussed in § 2.4. Verified Boot has mostly received negative reactions from software developers, so it is likely that an enclave licensing scheme would meet the same fate, should the developer community

become aware of it.

System Software Can Enforce a Launch Policy

§ C.4.3 explains that the SGX instructions used to load and initialize enclaves (ECREATE, EADD, EINIT) can only be issued by privileged system software, because they manage the EPC, which is a system resource.

A consequence on the restriction that only privileged software can issue ECREATE and EADD instructions is that the system software is able to track all the public contents that is loaded into each enclave. The privilege requirements of EINIT mean that the system software can also examine each enclave's SIGSTRUCT. It follows that the system software has access to a superset of the information that the Launch Enclave might use.

Furthermore, EINIT's privileged instruction status means that the system software can perform its own policy checks before allowing application software to initialize an enclave. So, the system software can enforce a Launch Control Policy set by the computer's owner. For example, an IaaS cloud service provider may use its hypervisor to implement a Launch Control Policy that limits what enclaves its customers are allowed to execute.

Given that the system software has access to a superset of the information that the Launch Enclave might use, it is easy to see that the set of policies that can be enforced by system software is a superset of the policies that can be supported by an LE. Therefore, the only rational explanation for the existence of the LE is that it was designed to implement a Launch Control Policy that is not beneficial to the computer owner.

As an illustration of this argument, we consider the case of restricting access to EGETKEY's Provisioning keys (§ C.4.8). The derivation material for Provisioning keys does not include OWNEREPOCH, so malicious enclaves can potentially use these keys to track a CPU chip package as it exchanges owners. For this reason, the SGX design includes a simple access control mechanism that can be used by system software to limiting enclave access to Provisioning keys. EGETKEY refuses to derive Provisioning keys for enclaves whose PROVISIONKEY attribute is not set to true.

It follows that a reasonable Launch Control Policy would only allow the PROVISIONKEY attribute to be set for the enclaves that implement software attestation, such as Intel's

Provisioning Enclave and Quoting Enclave. This policy can easily be implemented by system software, given its exclusive access to the `EINIT` instruction.

The only concern with the approach outlined above is that a malicious system software might abuse the `PROVISIONKEY` attribute to generate a unique identifier for the hardware that it runs on, similar to the much maligned Intel Processor Serial Number [89]. We dismiss this concern by pointing out that system software has access to many unique identifiers, such as the Media Access Control (MAC) address of the Ethernet adapter integrated into the motherboard's chipset (§ A.9.1).

Enclaves Cannot Damage the Host Computer

SGX enclaves execute at the lowest privilege level (user mode / ring 3), so they are subject to the same security checks as their host application. For example, modern operating systems set up the I/O maps (§ A.7) to prevent application software from directly accessing the I/O address space (§ A.4), and use the supervisor (S) page table attribute (§ A.5.3) to deny application software direct access to memory-mapped devices (§ A.4) and to the DRAM that stores the system software. Enclave software is subject to I/O privilege checks and address translation checks, so a malicious enclave cannot directly interact with the computer's devices, and cannot tamper the system software.

It follows that software running in an enclave has the same means to compromise the system software as its host application, which come down to exploiting a security vulnerability. The same solutions used to mitigate vulnerabilities exploited by application software (e.g., `seccomp/bpf` [120]) apply to enclaves.

The only remaining concern is that an enclave can perform a denial of service (DoS) attack against the system software. The rest of this section addresses the concern.

The SGX design provides system software the tools it needs to protect itself from enclaves that engage in CPU hogging and DRAM hogging. As enclaves cannot perform I/O directly, these are the only two classes of DoS attacks available to them.

An enclave that attempts to hog an LP assigned to it can be preempted by the system software via an Inter-Processor Interrupt (IPI, § A.12) issued from another processor. This method is available as long as the system software reserves at least one LP for non-enclave

computation.

Furthermore, most OS kernels use tick schedulers, which use a real-time clock (RTC) configured to issue periodical interrupts (ticks) to all cores. The RTC interrupt handler invokes the kernel's scheduler, which chooses the thread that will get to use the logical processor until the next RTC interrupt is received. Therefore, kernels that use tick schedulers always have the opportunity to de-schedule enclave threads, and don't need to rely on the ability to send IPIs.

In SGX, the system software can always evict an enclave's EPC pages to non-EPC memory, and then to disk. The system software can also outright deallocate an enclave's EPC pages, though this will probably cause the enclave code to encounter page faults that cannot be resolved. The only catch is that the EPC pages that hold metadata for running enclave threads cannot be evicted or removed. However, this can easily be resolved, as the system software can always preempt enclave threads, using one of the methods described above.

Interaction with Anti-Virus Software

Today's anti-virus (AV) systems are glorified pattern matchers. AV software simply scans all the executable files on the system and the memory of running processes, looking for bit patterns that are thought to only occur in malicious software. These patterns are somewhat pompously called "virus signatures".

SGX (and TXT, to some extent) provides a method for executing code in an isolated container that we refer to as an enclave. Enclaves are isolated from all the other software on the computer, including any AV software that might be installed.

The isolation afforded by SGX opens up the possibility for bad actors to structure their attacks as a generic loader that would end up executing a malicious payload without tripping the AV's pattern matcher. More specifically, the attack would create an enclave and initialize it with a generic loader that looks innocent to an AV. The loader inside the enclave would obtain an encrypted malicious payload, and would undergo software attestation with an Internet server to obtain the payload's encryption key. The loader would then decrypt the malicious payload and execute it inside the enclave.

In the scheme suggested here, the malicious payload only exists in a decrypted form inside an enclave's memory, which cannot be accessed by the AV. Therefore, the AV's pattern matcher will not trip.

This issue does not have a solution that maintains the status-quo for the AV vendors. The attack described above would be called a protection scheme if the payload would be a proprietary image processing algorithm, or a DRM scheme.

On a brighter note, enclaves do not bring the complete extinction of AV, they merely require a change in approach. Enclave code always executes at the lowest privilege mode (ring 3 / user mode), so it cannot perform any I/O without invoking the services of system software. For all intents and purposes, this effectively means that enclave software cannot perform any malicious action without the complicity of system software. Therefore, enclaves can be policed effectively by intelligent AV software that records and filters the I/O performed by software, and detects malicious software according to the actions that it performs, rather than according to bit patterns in its code.

Furthermore, SGX's enclave loading model allows the possibility of performing static analysis on the enclave's software. For simplicity, assume the existence of a standardized static analysis framework. The initial enclave contents is not encrypted, so the system software can easily perform static analysis on it. Dynamically loaded code or Just-In-Time code generation (JIT) can be handled by requiring that all enclaves that use these techniques embed the static analysis framework and use it to analyze any dynamically loaded code before it is executed. The system software can use static verification to ensure that enclaves follow these rules, and refuse to initialize any enclaves that fail verification.

In conclusion, enclaves in and of themselves don't introduce new attack vectors for malware. However, the enclave isolation mechanism is fundamentally incompatible with the approach employed by today's AV solutions. Fortunately, it is possible (though non-trivial) to develop more intelligent AV software for enclave software.

C.5 SGX Analysis

C.5.1 SGX Implementation Overview

An under-documented and overlooked feat achieved by the SGX design is that implementing it on an Intel processor has a very low impact on the chip's hardware design. SGX's modifications to the processor's execution cores (§ A.9.4) are either very small or completely inexistent. The CPU's uncore (§ A.9.3, § A.11.3) receives a new module, the Memory Encryption Engine, which appears to be fairly self-contained.

The bulk of the SGX implementation is relegated to the processor's microcode (§ A.14), which supports a much higher development speed than the chip's electrical circuitry.

Execution Core Modifications

At a minimum, the SGX design requires a very small modification to the processor's execution cores (§ A.9.4), in the Page Miss Handler (PMH, § A.11.5).

The PMH resolves TLB misses, and consists of a fast path that relies on an FSM page walker, and a microcode assist fallback that handles the edge cases (§ A.14.3). The bulk of SGX's memory access checks, which are discussed in § C.5.2, can be implemented in the microcode assist.

The only modification to the PMH hardware that is absolutely necessary to implement SGX is developing an ability to trigger the microcode assist for all address translations when a logical processor (§ A.9.4) is in enclave mode (§ C.4.4), or when the physical address produced by the page walker FSM matches the Processor Reserved Memory (PRM, § C.4.1) range.

The PRM range is configured by the PRM Range Registers (§ C.4.1), which have exactly the same semantics as the Memory Type Range Registers (MTRRs, § A.11.4) used to configure a variable memory range. The page walker FSM in the PMH is already configured to issue a microcode assist when the page tables are in uncacheable memory (§ A.11.4). Therefore, the PRMRR can be represented as an extra MTRR pair.

Uncore Modifications

The SDM states that DMA transactions (§ A.9.1) that target the PRM range are aborted by the processor. The SGX patents disclose that the PRMRR protection against unauthorized DMA is implemented by having the SGX microcode set up entries in the Source Address Decoder (SAD) in the uncore CBoxes and in the Target Address Decoder (TAD) in the integrated Memory Controller (MC).

§ A.11.3 mentions that Intel's Trusted Execution Technology (TXT) [74] already takes advantage of the integrated MC to protect a DRAM range from DMA. It is highly likely that the SGX implementation reuses the mechanisms brought by TXT, and only requires the extension of the SADs and TADs by one entry.

SGX's major hardware modification is the Memory Encryption Engine (MEE) that is added to the processor's uncore (§ A.9.3, § A.11.3) to protect SGX's Enclave Page Cache (EPC, § C.4.1) against physical attacks.

The MEE is briefly described in the ISCA 2015 SGX tutorial [106]. According to the information presented there, the MEE roughly follows the approach introduced by Aegis [180] [182], which relies on a variation of Merkle trees to provide the EPC with privacy, integrity, and freshness guarantees (§ B.1). Unlike Aegis, the MEE uses non-standard cryptographic primitives that include a slightly modified AES operating mode (§ B.1.2) and a Carter-Wegman [32, 196] MAC (§ B.1.3) construction.

Both the ISCA SGX tutorial and the patents state that the MEE is connected to the Memory Controller (MC) integrated in the CPU's uncore. However, all sources are completely silent on further implementation details. The MEE overview slide states that "the Memory Controller detects [the] address belongs to the MEE region, and routes transaction to MEE", which suggests that the MEE is fairly self-contained and has a narrow interface to the rest of the MC.

Intel's SGX patents use the name Crypto Memory Aperture (CMA) to refer to the MEE. The CMA description matches the MEE and PRM concepts, as follows. According to the patents, the CMA is used to securely store the EPC, relies on crypto controllers in the MC, and loses its keys during deep sleep. These details align perfectly with the SDM's statements

regarding the MEE and PRM.

The Intel patents also disclose that the EPCM (§ C.4.1) and other structures used by the SGX implementation are also stored in the PRM. This rules out the possibility that the EPCM requires on-chip memory resembling the last-level cache (§ A.11, § A.11.3).

Last, the SGX patents shine a bit of light on an area that the official Intel documentation is completely silent about, namely the implementation concerns brought by computer systems with multiple processor chips. The patents state that the MEE also protects the Quick-Path Interconnect (QPI, § A.9.1) traffic using link-layer encryption.

Microcode Modifications

According to the SGX patents, all the SGX instructions are implemented in microcode. This can also be deduced by reading the SDM's pseudocode for all the instructions, and realizing that it is highly unlikely that any SGX instruction can be implemented in 4 or fewer micro-ops (§ A.10), which is the most that can be handled by the simple decoders used in the hardware fast paths (S A.14.1).

The Asynchronous Enclave Exit (AEX, § C.4.4) behavior is also implemented in microcode. § A.14.2 draws on an assortment of Intel patents to conclude that hardware exceptions (§ A.8.2), including both faults and interrupts, trigger microcode events (§ A.14.2). It follows that the SGX implementation can implement AEX by modifying the hardware exception handlers in the microcode.

The SGX initialization sequence is also implemented in microcode. SGX is initialized in two phases. First, it is very likely that the boot sequence in microcode (§ A.14.4) was modified to initialize the registers associated with the SGX microcode. The ISCA SGX tutorial states that the MEE' keys are initialized during the boot process. Second, SGX instructions are enabled by setting a bit in a Model-Specific Register (MSR, § A.4). This second phase involves enabling the MEE and configuring the SAD and TAD to protect the PRM range. Both tasks are amenable to a microcode implementation.

The SGX description in the SDM implies that the SGX implementation uses a significant number of new registers, which are only exposed to microcode. However, the SGX patents reveal that most of these registers are actually stored in DRAM.

For example, the patents state that each TCS (§ C.4.2) has two fields that receive the values of the DR7 and IA32_DEBUGCTL registers when the processor enters enclave mode (§ C.4.4), and are used to restore the original register values during enclave exit (§ C.4.4). The SDM documents these fields as “internal CREGs” (CR_SAVE_DR7 and CR_SAVE_DEBUGCTL), which are stated to be “hardware specific registers”.

The SGX patents document a small subset of the CREGs described in the SDM, summarized in Table C.10, as microcode registers. While in general we trust official documentation over patents, in this case we use the CREG descriptions provided by the patents, because they appear to be more suitable for implementation purposes.

SDM Name	Bits	Scope	Description
CSR_SGX_OWNEREPOCH	128	CPU Chip Package	Used by EGETKEY (§ C.4.7)
CR_ENCLAVE_MODE	1	Logical Processor	1 when executing code inside an enclave
CR_ACTIVE_SECS	16	Logical Processor	The index of the EPC page storing the current enclave’s SECS
CR_TCS_LA	64	Logical Processor	The virtual address of the TCS (§ C.4.2) used to enter (§ C.4.4) the current enclave
CR_TCS_PH	16	Logical Processor	The index of the EPC page storing the TCS used to enter the current enclave
CR_XSAVE_PAGE_0	16	Logical Processor	The index of the EPC page storing the first page of the current SSA (§ C.4.2)

Table C.10: The fields in an EPCM entry.

From a cost-performance standpoint, the cost of register memory only seems to be justified for the state used by the PMH to implement SGX’s memory access checks, which will be discussed in § C.5.2). The other pieces of state listed as CREGs are accessed so infrequently that storing them in dedicated SRAM would make very little sense.

The SGX patents state that SGX requires very few hardware changes, and most of the implementation is in microcode, as a positive fact. We therefore suspect that minimizing hardware changes was a high priority in the SGX design, and that any SGX modification

proposals need to be aware of this priority.

C.5.2 SGX Memory Access Protection

SGX guarantees that the software inside an enclave is isolated from all the software outside the enclave, including the software running in other enclaves. This isolation guarantee is at the core of SGX's security model.

It is tempting to assume that the main protection mechanism in SGX is the Memory Encryption Engine (MEE) described in § C.5.1, as it encrypts and MACs the DRAM's contents. However, the MEE sits in the processor's memory controller, which is at the edge of the on-chip memory hierarchy, below the caches (§ A.11). Therefore, the MEE cannot protect an enclave's memory from software attacks.

The root of SGX's protections against software attacks is a series of memory access checks which prevents the currently running software from accessing memory that does not belong to it. Specifically, non-enclave software is only allowed to access memory outside the PRM range, while the code inside an enclave is allowed to access non-PRM memory, and the EPC pages owned by the enclave.

Although it is believed [54] that SGX's access checks are performed on every memory access check, Intel's patents disclose that the checks are performed in the Page Miss Handler (PMH, § A.11.5), which only handles TLB misses.

Functional Description

The intuition behind SGX's memory access protections can be built by considering what it would take to implement the same protections in a trusted operating system or hypervisor, solely by using the page tables that direct the CPU's address translation feature (§ A.5).

The hypothetical trusted software proposed above can implement enclave entry (§ C.4.4) as a system call § A.8.1 that creates page table entries mapping the enclave's memory. Enclave exit (§ C.4.4) can be a symmetric system call that removes the page table entries created during enclave entry. When modifying the page tables, the system software has to consider TLB coherence issues (§ A.11.5) and perform TLB shutdowns when appropriate.

SGX leaves page table management under the system software’s control, but it cannot trust the software to set up the page tables in any particular way. Therefore, the hypothetical design described above cannot be used by SGX as-is. Instead, at a conceptual level, the SGX implementation approximates the effect of having the page tables set up correctly by inspecting every address translation that comes out of the Page Miss Handler (PMH, § A.11.5). The address translations that do not obey SGX’s access control restrictions are rejected before they reach the TLBs.

SGX’s approach relies on the fact that software always references memory using virtual addresses, so all the micro-ops (§ A.10) that reach the memory execution units (§ A.10.1) use virtual addresses that must be resolved using the TLBs before the actual memory accesses are carried out. By contrast, the processor’s microcode (§ A.14) has the ability to issue physical memory accesses, which bypass the TLBs. Conveniently, SGX instructions are implemented in microcode (§ C.5.1), so they can bypass the TLBs and access memory that is off limits to software, such as the EPC page holding an enclave’s SECS (§ C.4.1).

The SGX address translation checks use the information in the Enclave Page Cache Map (EPCM, § C.4.1), which is effectively an inverted page table that covers the entire EPC. This means that each EPC page is accounted for by an EPCM entry, using the structure is summarized in Table C.11. The EPCM fields were described in detail in § C.4.1, § C.4.2, § C.4.2, § C.4.5, and § C.4.5.

Field	Bits	Description
VALID	1	0 for un-allocated EPC pages
BLOCKED	1	page is being evicted
R	1	enclave code can read
W	1	enclave code can write
X	1	enclave code can execute
PT	8	page type (Table C.12)
ADDRESS	48	the virtual address used to access this page
ENCLAVESECS		the EPC slot number for the SECS of the enclave owning the page

Table C.11: The fields in an EPCM entry.

Conceptually, SGX adds the access control logic illustrated in Figure C-30 to the PMH. SGX’s security checks are performed after the page table attributes-based checks (§ A.5.3)

Type	Allocated by	Contents
PT_REG	EADD	enclave code and data
PT_SECS	ECREATE	SECS (§ C.4.1)
PT_TCS	EADD	TCS (§ C.4.2)
PT_VA	EPA	VA (§ C.4.5)

Table C.12: Values of the PT (page type) field in an EPCM entry.

defined by the Intel architecture. It follows that SGX’s access control logic has access to the physical address produced by the page walker FSM.

SGX’s security checks depend on whether the logical processor (§ A.9.4) is in enclave mode (§ C.4.4) or not. While the processor is outside enclave mode, the PMH allows any address translation that does not target the PRM range (§ C.4.1). When the processor is inside enclave mode, the PMH performs the checks described below, which provide the security guarantees described in § C.4.2.

First, virtual addresses inside the enclave’s virtual memory range (ELRANGE, § C.4.2) must always translate into physical addresses inside the EPC. This way, an enclave is assured that all the code and data stored in ELRANGE receives SGX’s privacy, integrity, and freshness guarantees. Since the memory outside ELRANGE does not enjoy these guarantees, the SGX design disallows having enclave code outside ELRANGE. This is most likely accomplished by setting the disable execution (XD, § A.5.3) attribute on the TLB entry.

Second, an EPC page must only be accessed by the code of the enclave who owns the page. For the purpose of this check, each enclave is identified by the index of the EPC page that stores the enclave’s SECS (§ C.4.1). The current enclave’s identifier is stored in the CR_ACTIVE_SECS microcode register during enclave entry. This register is compared against the enclave identifier stored in the EPCM entry corresponding to the EPC page targeted by the address translation.

Third, some EPC pages cannot be accessed by software. Pages that hold SGX internal structures, such as a SECS, a TCS (§ C.4.2), or a VA (§ C.4.5) must only be accessed by SGX’s microcode, which uses physical addresses and bypasses the address translation unit, including the PMH. Therefore, the PMH rejects address translations targeting these pages.

Blocked (§ C.4.5) EPC pages are in the process of being evicted (§ C.4.5), so the PMH

must not create new TLB entries targeting them.

Next, an enclave's EPC pages must always be accessed using the virtual addresses associated with them when they were allocated to the enclave. Regular EPC pages, which can be accessed by software, are allocated to enclaves using the `EADD` (§ C.4.3) instruction, which reads in the page's address in the enclave's virtual address space. This address is stored in the `LINADDR` field in the corresponding EPCM entry. Therefore, all the PMH has to do is to ensure that `LINADDR` in the address translation's target EPCM entry equals the virtual address that caused the TLB miss which invoked the PMH.

At this point, the PMH's security checks have completed, and the address translation result will definitely be added to the TLB. Before that happens, however, the SGX extensions to the PMH apply the access restrictions in the EPCM entry for the page to the address translation result. While the public SGX documentation we found did not describe this process, there is a straightforward implementation that fulfills SGX's security requirements. Specifically, the TLB entry bits `P`, `W`, and `XD` can be AND-ed with the EPCM entry bits `R`, `W`, and `X`.

EPCM Entry Representation

Most EPCM entry fields have obvious representations. The exception is the `LINADDR` and `ENCLAVESECS` fields, described below. These representations explain SGX's seemingly arbitrary limit on the size of an enclave's virtual address range (`ELRANGE`).

The SGX patents disclose that the `LINADDR` field in an EPCM entry stores the virtual page number (`VPN`, § A.5.1) of the corresponding EPC page's expected virtual address, relative to the `ELRANGE` base of the enclave that owns the page.

The representation described above reduces the number of bits needed to store `LINADDR`, assuming that the maximum `ELRANGE` size is significantly smaller than the virtual address size supported by the CPU. This desire to save EPCM entry bits is the most likely motivation for specifying a processor model-specific `ELRANGE` size, which is reported by the `CPUID` instruction.

The SDM states that the `ENCLAVESECS` field of an EPCM entry corresponding to an EPC page indicates the `SECS` of belonging to the enclave that owns the page. Intel's

patents reveal that the SECS address in ENCLAVESECS is represented as a physical page number (PPN, § A.5.1) relative to the start of the EPC. Effectively, this relative PPN is the 0-based EPC page index.

The EPC page index representation saves bits in the ECPM entry, assuming that the EPCM size is significantly smaller than the physical address space supported by the CPU. The ISCA 2015 SGX tutorial slides mention an EPC size of 96MB, which is significantly smaller than the physical addressable space on today's typical processors, which is $2^{36} - 2^{40}$ bytes.

PMH Hardware Modifications

The SDM describes the memory access checks performed after SGX is enabled, but does not provide any insight into their implementation. Intel's patents hint at three possible implementations that make different cost-performance tradeoffs. This section summarizes the three approaches and argues in favor of the implementation that requires the fewest hardware modifications to the PMH.

All implementations of SGX's security checks entail adding a pair of memory type range registers (MTRRs, § A.11.4) to the PMH. These registers are named the *Secure Enclave Range Registers* (SERR) in Intel's patents. Enabling SGX on a logical processor initializes the SERR to the values of the Protected Memory Range Registers (PMRR, § C.4.1).

Furthermore, all implementations have the same behavior when a logical processor is outside enclave mode. The memory type range described by the SERR is enabled, causing a microcode assist to trigger for every address translation that resolves inside the PRM. SGX's implementation uses the microcode assist to replace the address translation result with an address that causes memory access transactions to be aborted.

The three implementations differ in their behavior when the processor enters enclave mode (§ C.4.4) and starts executing enclave code.

The alternative that requires the least amount of hardware changes sets up the PMH to trigger a microcode assist for every address translation. This can be done by setting the SERR to cover all the physical memory (e.g., by setting both the base and the mask to zero). In this approach, the microcode assist implements all the enclave mode security checks

illustrated in Figure C-30.

A speedier alternative adds a pair of registers to the PMH that represents the current enclave's ELRANGE and modifies the PMH so that, in addition to checking physical addresses against the SERR, it also checks the virtual addresses going into address translations against ELRANGE. When either check is true, the PMH invokes the microcode assist used by SGX to implement its memory access checks. Assuming the ELRANGE registers use the same base / mask representation as variable MTRRs, enclave exists can clear ELRANGE by zeroing both the base and the mask. This approach uses the same microcode assist implementation, minus the ELRANGE check that moves into the PMH hardware.

The second alternative described above has the benefit that the microcode assist is not invoked for enclave mode accesses outside ELRANGE. However, § C.4.2 argues that an enclave should treat all the virtual memory addresses outside ELRANGE as untrusted storage, and only use that memory to communicate with software outside the enclave. Taking this into consideration, well-designed enclaves would spend relatively little time performing memory accesses outside ELRANGE. Therefore, this second alternative is unlikely to obtain performance gains that are worth its cost.

The last and most performant alternative would entail implementing all the access checks shown in Figure C-30 in hardware. Similarly to the address translation FSM, the hardware would only invoke a microcode assist when a security check fails and a Page Fault needs to be handled.

The high-performance implementation described above avoids the cost of microcode assists for all TLB misses, assuming well-behaved system software. In this association, a microcode assist results in a Page Fault, which triggers an Asynchronous Enclave Exit (AEX, § C.4.4). The cost of the AEX dominates the performance overhead of the microcode assist.

While this last implementation looks attractive, one needs to realize that TLB misses occur quite infrequently, so a large improvement in the TLB miss speed translates into a much less impressive improvement in overall enclave code execution performance. Taking this into consideration, it seems unwise to commit to extensive hardware modifications in the PMH before SGX gains adoption.

C.5.3 SGX Security Check Correctness

In § C.5.2, we argued that SGX’s security guarantees can be obtained by modifying the Page Miss Handler (PMH, § A.11.5) to block undesirable address translations from reaching the TLB. This section builds on the result above and outlines a correctness proof for SGX’s memory access protection.

Specifically, we outline a proof for the following invariant. **At all times, all the TLB entries in every logical processor will be consistent with SGX’s security guarantees.** By the argument in § C.5.2, the invariant translates into an assurance that all the memory accesses performed by software obey SGX’s security model. The high-level proof structure is presented because it helps understand how the SGX security checks come together. By contrast, a detailed proof would be incredibly tedious, and would do very little to boost the reader’s understanding of SGX.

Top-Level Invariant Breakdown

We first break down the above invariant into specific cases based on whether a logical processor (LP) is executing enclave code or not, and on whether the TLB entries translate virtual addresses in the current enclave’s ELRANGE (§ C.4.2). When the processor is outside enclave mode, ELRANGE can be considered to be empty. This reasoning yields the three cases outlined below.

1. At all times when an LP is outside enclave mode, its TLB may only contain physical addresses belonging to DRAM pages outside the PRM.
2. At all times when an LP is inside enclave mode, the TLB entries for virtual addresses outside the current enclave’s ELRANGE must contain physical addresses belonging to DRAM pages outside the PRM.
3. At all times when an LP is in enclave mode, the TLB entries for virtual addresses inside the current enclave’s ELRANGE must match the virtual memory layout specified by the enclave author.

The first two invariant cases can be easily proven independently for each LP, by induction over the sequence of instructions executed by the LP. For simplicity, the reader can assume that instructions are executed in program mode. While the assumption is not true on processors with out-of-order execution (§ A.10), the arguments presented here also hold when the executed instruction sequence is considered in retirement order, for reasons that will be described below.

An LP will only transition between enclave mode and non-enclave mode at a few well-defined points, which are `EENTER` (§ C.4.4), `ERESUME` (§ C.4.4), `EEXIT` (§ C.4.4), and Asynchronous Enclave Exits (AEX, § C.4.4). According to the SDM, all the transition points flush the TLBs and the out-of-order execution pipeline. In other words, the TLBs are guaranteed to be empty after every transition between enclave mode and non-enclave mode, so we can consider all these transitions to be trivial base cases for our induction proofs.

While SGX initialization is not thoroughly discussed, the SDM mentions that loading some Model-Specific Registers (MSRs, § A.4) triggers TLB flushes, and that system software should flush TLBs when modifying Memory Type Range Registers (MTRRs, § A.11.4). Given that all the possible SGX implementations described in § C.5.2 entail adding a MTRR, it is safe to assume that enabling SGX mode also results in a TLB flush and out-of-order pipeline flush, and can be used by our induction proof as well.

All the base cases in the induction proofs are serialization points for out-of-order execution, as the pipeline is flushed during both enclave mode transitions and SGX initialization. This makes the proofs below hold when the program order instruction sequence is replaced with the retirement order sequence.

The first invariant case holds because while the LP is outside enclave mode, the SGX security checks added to the PMH (§ C.5.2, Figure C-30) reject any address translation that would point into the PRM before it reaches the TLBs. A key observation for proving the induction step of this invariant case is that the PRM never changes after SGX is enabled on an LP.

The second invariant case can be proved using a similar argument. While an LP is executing an enclave's code, the SGX memory access checks added to the PMH reject any address translation that resolves to a physical address inside the PRM, if the translated

virtual address falls outside the current enclave's ELRANGE. The induction step for this invariant case can be proven by observing that a change in an LP's current ELRANGE is always accompanied by a TLB flush, which results in an empty TLB that trivially satisfies the invariant. This follows from the constraint that an enclave's ELRANGE never changes after it is established, and from the observation that the LP's current enclave can only be changed by an enclave entry, which must be preceded by an enclave exit, which triggers a TLB flush.

The third invariant case is best handled by recognizing that the Enclave Page Cache Map (EPCM, § C.4.1) is an intermediate representation for the virtual memory layout specified by the enclave authors. This suggests breaking down the case into smaller sub-invariants centered around the EPCM, which will be proven in the sub-sections below.

1. At all times, each EPCM entry for a page that is allocated to an enclave matches the virtual memory layout desired by the enclave's author.
2. Assuming that the EPCM contents is constant, at all times when an LP is in enclave mode, the TLB entries for virtual addresses inside the current enclave's ELRANGE must match EPCM entries that belong to the enclave.
3. An EPCM entry is only modified when there is no mapping for it in any LP's TLB.

The second and third invariant combined prove that all the TLBs in an SGX-enabled computer always reflect the contents of the EPCM, as the third invariant essentially covers the gaps in the second invariant. This result, in combination with the first invariant, shows that the EPCM is a bridge between the memory layout specifications of the enclave authors and the TLB entries that regulate what memory can be accessed by software executing on the LPs. When further combined with the reasoning in § C.5.2, the whole proof outlined here results in an end-to-end argument for the correctness of SGX's memory protection scheme.

EPCM Entries Reflect Enclave Author Design

This sub-section outlines the proof for the following invariant. **At all times, each EPCM entry for a page that is allocated to an enclave matches the virtual memory layout**

desired by the enclave's author.

A key observation, backed by the SDM pseudocode for SGX instructions, is that all the instructions that modify the EPCM pages allocated to an enclave are synchronized using a lock in the enclave's SECS. This entails the existence of a time ordering of the EPCM modifications associated with an enclave. We prove the invariant stated above using a proof by induction over this sequence of EPCM modifications.

EPCM entries allocated to an enclave are created by instructions that can only be issued before the enclave is initialized, specifically `ECREATE` (§ C.4.3) and `EADD` (§ C.4.3). The contents of the EPCM entries created by these instructions contributes to the enclave's measurement (§ C.4.6), together with the initial data loaded into the corresponding EPC pages.

§ B.3.2 argues that we can assume that enclaves with incorrect measurements do not exist, as they will be rejected by software attestation. Therefore, we can assume that the attributes used to initialize EPCM pages match the enclave authors' memory layout specifications.

EPCM entries can be evicted to untrusted DRAM, together with their corresponding EPC pages, by the `EWB` (§ C.4.5) instruction. The `ELDU / ELDB` (§ C.4.5) instructions re-load evicted page contents and metadata back into the EPC and EPCM. By induction, we can assume that an EPCM entry matches the enclave author's specification when it is evicted. Therefore, if we can prove that the EPCM entry that is reloaded from DRAM is equivalent to the entry that was evicted, we can conclude that the reloaded entry matches the author's specification.

A detailed analysis of the cryptographic primitives used by the SGX design to protect the evicted EPC page contents and its associated metadata is outside the scope of this work. Summarizing the description in § C.4.5, the contents of evicted pages is encrypted using AES-GMAC (§ B.1.3), which is an authenticated encryption mechanism. The MAC tag produced by AES-GMAC covers the EPCM metadata as well as the page data, and includes a 64-bit version that is stored in a version tree whose nodes are Version Array (VA, (§ C.4.5) pages.

Assuming no cryptographic weaknesses, SGX's scheme does appear to guarantee the privacy, integrity, and freshness of the EPC page contents and associated metadata while it

is evicted in untrusted memory. It follows that EWB will only reload an EPCM entry if the contents is equivalent to the contents of an evicted entry.

The equivalence notion invoked here is slightly different from perfect equality, in order to account for the allowable operation of evicting an EPC page and its associated EPCM entry, and then reloading the page contents to a different EPC page and a different EPCM entry, as illustrated in Figure C-13. Loading the contents of an EPC page at a different physical address than it had before does not break the virtual memory abstraction, as long as the contents is mapped at the same virtual address (the LINEARADDRESS EPCM field), and has the same access control attributes (R, W, X, PT EPCM fields) as it had when it was evicted.

The rest of this section enumerates the address translation attacks prevented by the MAC verification that occurs in ELDU / ELDB. This is intended to help the reader develop some intuition for the reasoning behind using the page data and all the EPCM fields to compute and verify the MAC tag.

The most obvious attack is prevented by having the MAC cover the contents of the evicted EPC page, so the untrusted OS cannot modify the data in the page while it is stored in untrusted DRAM. The MAC also covers the metadata that makes up the EPCM entry, which prevents the more subtle attacks described below.

The enclave ID (EID) field is covered by the MAC tag, so the OS cannot evict an EPC page belonging to one enclave, and assign the page to a different enclave when it is loaded back into the EPC. If EID was not covered by authenticity guarantees, a malicious OS could read any enclave's data by evicting an EPC page belonging to the victim enclave, and loading it into a malicious enclave that would copy the page's contents to untrusted DRAM.

The virtual address (LINADDR) field is covered by the MAC tag, so the OS cannot modify the virtual memory layout of an enclave by evicting an EPC page and specifying a different LINADDR when loading it back. If LINADDR was not covered by authenticity guarantees, a malicious OS could perform the exact attack shown in Figure B-24 and described in § B.7.3.

The page access permission flags (R, W, X) are also covered by the MAC tag. This prevents the OS from changing the access permission bits in a page's EPCM entry by

evicting the page and loading it back in. If the permission flags were not covered by authenticity guarantees, the OS could use the ability to change EPCM access permissions to facilitate exploiting vulnerabilities in enclave code. For example, exploiting a stack overflow vulnerability is generally easier if OS can make the stack pages executable.

The nonce stored in the VA slot is also covered by the MAC. This prevents the OS from mounting a replay attack that reverts the contents of an EPC page to an older version. If the nonce would not be covered by integrity guarantees, the OS could evict the target EPC page at different times t_1 and t_2 in the enclave's life, and then provide the EWB outputs at t_1 to the ELDU / ELDB instruction. Without the MAC verification, this attack would successfully revert the contents of the EPC page to its version at t_1 .

While replay attacks look relatively benign, they can be quite devastating when used to facilitate double spending.

TLB Entries for ELRANGE Reflect EPCM Contents

This sub-section sketches a proof for the following invariant. **At all times when an LP is in enclave mode, the TLB entries for virtual addresses inside the current enclave's ELRANGE must match EPCM entries that belong to the enclave.** The argument makes the assumption that **the EPCM contents is constant**, which will be justified in the following sub-section.

The invariant can be proven by induction over the sequence of TLB insertions that occur in the LP. This sequence is well-defined because an LP has a single PMH, so the address translation requests triggered by TLB misses must be serialized to be processed by the PMH.

The proof's induction step depends on the fact that the TLB on hyper-threaded cores (§ A.9.4) is dynamically partitioned between the two LPs that share the core, and no TLB entry is shared between the LPs. This allows our proof to consider the TLB insertions associated with one LP independently from the other LP's insertions, which means we don't have to worry about the state (e.g., enclave mode) of the other LP on the core.

The proof is further simplified by observing that when an LP exits enclave mode, both its TLB and its out-of-order instruction pipeline are flushed. Therefore, the enclave mode and current enclave register values used by address translations are guaranteed to match the

values obtained by performing the translations in program order.

Having eliminated all the complexities associated with hyper-threaded (§ A.9.4) out-of-order (§ A.10) execution cores, it is easy to see that the security checks outlined in Figure C-30 and § C.5.2 ensure that TLB entries that target EPC pages are guaranteed to reflect the constraints in the corresponding EPCM entries.

Last, the SGX access checks implemented in the PMH reject any address translation for a virtual address in ELRANGE that does not resolve to an EPC page. It follows that memory addresses inside ELRANGE can only map to EPC pages which, by the argument above, must follow the constraints of the corresponding EPCM entries.

EPCM Entries are Not In TLBs When Modified

In this sub-section, we outline a proof that **an EPCM entry is only modified when there is no mapping for it in any LP's TLB.** This proof analyzes each of the instructions that modify EPCM entries.

For the purposes of this proof, we consider that setting the BLOCKED attribute does not count as a modification to an EPCM entry, as it does not change the EPC page that the entry is associated with, or the memory layout specification associated with the page.

The instructions that modify EPCM entries in such a way that the resulting EPCM entries have the VALID field set to true require that the EPCM entries were invalid before they were modified. These instructions are ECREATE (§ C.4.3), EADD (§ C.4.3), EPA (§ C.4.5), and ELDU / ELDB (§ C.4.5). The EPCM entry targeted by any these instructions must have had its VALID field set to false, so the invariant proved in the previous sub-section implies that the EPCM entry had no TLB entry associated with it.

Conversely, the instructions that modify EPCM entries and result in entries whose VALID field is false start out with valid entries. These instructions are EREMOVE (§ C.4.3) and EWB (§ C.4.5).

The EPCM entries associated with EPC pages that store Version Arrays (VA, § C.4.5) represent a special case for both instructions mentioned above, as these pages are not associated with any enclave. As these pages can only be accessed by the microcode used to implement SGX, they never have TLB entries representing them. Therefore, both EREMOVE

and EWB can invalidate EPCM entries for VA pages without additional checks.

EREMOVE only invalidates an EPCM entry associated with an enclave when there is no LP executing in enclave mode using a TCS associated with the same enclave. An EPCM entry can only result in TLB translations when an LP is executing code from the entry's enclave, and the TLB translations are flushed when the LP exits enclave mode. Therefore, when EREMOVE invalidates an EPCM entry, any associated TLB entry is guaranteed to have been flushed.

EWB's correctness argument is more complex, as it relies on the EBLOCK / ETRACK sequence described in § C.4.5 to ensure that any TLB entry that might have been created for an EPCM entry is flushed before the EPCM entry is invalidated.

Unfortunately, the SDM pseudocode for the instructions mentioned above leaves out the algorithm used to verify that the relevant TLB entries have been flushed. Thus, we must base our proof on the assumption that the SGX implementation produced by Intel's engineers matches the claims in the SDM. In § C.5.4, we propose a method for ensuring that EWB will only succeed when all the LPs executing an enclave's code at the time when ETRACK is called have exited enclave mode at least once between the ETRACK call and the EWB call. Having proven the existence of a correct algorithm by construction, we can only hope that the SGX implementation uses our algorithm, or a better algorithm that is still correct.

C.5.4 Tracking TLB Flushes

This section proposes a straightforward method that the SGX implementation can use to verify that the system software plays its part correctly in the EPC page eviction (§ C.4.5) process. Our method meets the SDM's specification for EBLOCK (§ C.4.5), ETRACK (§ C.4.5) and EWB (§ C.4.5).

The motivation behind this section is that, at least at the time of this writing, there is no official SGX documentation that contains a description of the mechanism used by EWB to ensure that all the Logical Processors (LPs, § A.9.4) running an enclave's code exit enclave mode (§ C.4.4) between an ETRACK invocation and a EWB invocation. Knowing that there exists a correct mechanism that has the same interface as the SGX instructions described in

the SDM gives us a reason to hope that the SGX implementation is also correct.

Our method relies on the fact that an enclave's SECS (§ C.4.1) is not accessible by software, and is already used to store information used by the SGX microcode implementation (§ C.5.1). We store the following fields in the SECS. *tracking* and *done-tracking* are Boolean variables. *tracked-threads* and *active-threads* are non-negative integers that start at zero and must store numbers up to the number of LPs in the computer. *lp-mask* is an array of Boolean flags that has one member per LP in the computer. The fields are initialized as shown in Figure C-31.

The *active-threads* SECS field tracks the number of LPs that are currently executing the code of the enclave who owns the SECS. The field is atomically incremented by EENTER (§ C.4.4) and ERESUME (§ C.4.4) and is atomically decremented by EEXIT (§ C.4.4) and Asynchronous Enclave Exits (AEXs, § C.4.4). Besides from helping track TLB flushes, this field can also be used by EREMOVE (§ C.4.3) to decide when it is safe to free an EPC page that belongs to an enclave.

As specified in the SDM, ETRACK activates TLB flush tracking for an enclave. In our method, this is accomplished by setting the *tracking* field to TRUE and the *done-tracking* field to FALSE.

When tracking is enabled, *tracked-threads* is the number of LPs that were executing the enclave's code when the ETRACK instruction was issued, and have not yet exited enclave mode. Therefore, executing ETRACK atomically reads *active-threads* and writes the result into *tracked-threads*. Also, *lp-mask* keeps track of the LPs that have exited the current enclave after the ETRACK instruction was issued. Therefore, the ETRACK implementation atomically zeroes *lp-mask*. The full ETRACK algorithm is listed in Figure C-32.

When an LP exits an enclave that has TLB flush tracking activated, we atomically test and set the current LP's flag in *lp-mask*. If the flag was not previously set, it means that an LP that was executing the enclave's code when ETRACK was invoked just exited enclave mode for the first time, and we atomically decrement *tracked-threads* to reflect this fact. In other words, *lp-mask* prevents us from double-counting an LP when it exits the same enclave while TLB flush tracking is active.

Once *active-threads* reaches zero, we are assured that all the LPs running the enclave's

code when ETRACK was issued have exited enclave mode at least once, and can set the *done-tracking* flag. Figure C-33 enumerates all the steps taken on enclave exit.

Without any compensating measure, the method above will incorrectly decrement *tracked-threads*, if the LP exiting the enclave had entered it after ETRACK was issued. We compensate for this with the following trick. When an LP starts executing code inside an enclave that has TLB flush tracking activated, we set its corresponding flag in *lp-mask*. This is sufficient to avoid counting the LP when it exits the enclave. Figure C-34 lists the steps required by our method when an LP enters an enclave.

With these algorithms in place, EWB can simply verify that both *tracking* and *done-tracking* are TRUE. This ensures that the system software has triggered enclave exits on all the LPs that were running the enclave's code when ETRACK was executed. Figure C-35 lists the algorithm used by the EWB tracking verification step.

Last, EBLOCK marks the end of a TLB flush tracking cycle by clearing the *tracking* flag. This ensures that system software must go through another cycle of ETRACK and enclave exits before being able to use EWB on the page whose BLOCKED EPCM field was just set to TRUE by EBLOCK. Figure C-36 shows the details.

Our method's correctness can be easily proven by arguing that each SECS field introduced in this section has its intended value throughout enclave entries and exits.

C.5.5 Enclave Signature Verification

Let m be the public modulus in the enclave author's RSA key, and s be the enclave signature. Since the SGX design fixes the value of the public exponent e to 3, verifying the RSA signature amounts to computing the signed message $M = s^3 \bmod m$, checking that the value meets the PKCS v1.5 padding requirements, and comparing the 256-bit SHA-2 hash inside the message with the value obtained by hashing the relevant fields in the SIGSTRUCT supplied with the enclave.

This section describes an algorithm for computing the signed message while only using subtraction and multiplication on large non-negative integers. The algorithm admits a significantly simpler implementation than the typical RSA signature verification algorithm,

by avoiding the use of long division and negative numbers. The description here is essentially the idea in [77], specialized for $e = 3$.

The algorithm provided here requires the signer to compute the q_1 and q_2 values shown below. The values can be computed from the public information in the signature, so they do not leak any additional information about the private signing key. Furthermore, the algorithm verifies the correctness of the values, so it does not open up the possibility for an attack that relies on supplying incorrect values for q_1 and q_2 .

$$q_1 = \left\lfloor \frac{s^2}{m} \right\rfloor$$
$$q_2 = \left\lfloor \frac{s^3 - q_1 \times s \times m}{m} \right\rfloor$$

Due to the desirable properties mentioned above, it is very likely that the algorithm described here is used by the SGX implementation to verify the RSA signature in an enclave's SIGSTRUCT (§ C.4.7).

The algorithm in Figure C-37 computes the signed message $M = s^3 \bmod m$, while also verifying that the given values of q_1 and q_2 are correct. The latter is necessary because the SGX implementation of signature verification must handle the case where an attacker attempts to exploit the signature verification implementation by supplying invalid values for q_1 and q_2 .

The rest of this section proves the correctness of the algorithm in Figure C-37.

Analysis of Steps 1 - 4

Steps 1 – 4 in the algorithm check the correctness of q_1 and use it to compute $s^2 \bmod m$. The key observation to understanding these steps is recognizing that q_1 is the quotient of the integer division s^2/m .

Having made this observation, we can use elementary division properties to prove that the supplied value for q_1 is correct if and only if the following property holds.

$$0 \leq s^2 - q_1 \times m < m$$

We observe that the first comparison, $0 \leq s^2 - q_1 \times m$, is equivalent to $q_1 \times m \leq s^2$, which is precisely the check performed by step 2. We can also see that the second comparison, $s^2 - q_1 \times m < m$ corresponds to the condition verified by step 4. Therefore, if the algorithm passes step 4, it must be the case that the value supplied for q_1 is correct.

We can also plug s^2 , q_1 and m into the integer division remainder definition to obtain the identity $s^2 \bmod m = s^2 - q_1 \times m$. However, according to the computations performed in steps 1 and 3, $w = s^2 - q_1 \times m$. Therefore, we can conclude that $w = s^2 \bmod m$.

Analysis of Steps 5 - 8

Similarly, steps 5 – 8 in the algorithm check the correctness of q_2 and use it to compute $w \times s \bmod m$. The key observation here is that q_2 is the quotient of the integer division $(w \times s)/m$.

We can convince ourselves of the truth of this observation by using the fact that $w = s^2 \bmod m$, which was proven above, by plugging in the definition of the remainder in integer division, and by taking advantage of the distributivity of integer multiplication with respect to addition.

$$\begin{aligned}
 \left\lfloor \frac{w \times s}{m} \right\rfloor &= \left\lfloor \frac{(s^2 \bmod m) \times s}{m} \right\rfloor \\
 &= \left\lfloor \frac{(s^2 - \lfloor \frac{s^2}{m} \rfloor \times m) \times s}{m} \right\rfloor \\
 &= \left\lfloor \frac{s^3 - \lfloor \frac{s^2}{m} \rfloor \times m \times s}{m} \right\rfloor \\
 &= \left\lfloor \frac{s^3 - q_1 \times m \times s}{m} \right\rfloor \\
 &= \left\lfloor \frac{s^3 - q_1 \times s \times m}{m} \right\rfloor \\
 &= q_2
 \end{aligned}$$

By the same argument used to analyze steps 1 – 4, we use elementary division properties to prove that q_2 is correct if and only if the equation below is correct.

$$0 \leq w \times s - q_2 \times m < m$$

The equation's first comparison, $0 \leq w \times s - q_2 \times m$, is equivalent to $q_2 \times m \leq w \times s$, which corresponds to the check performed by step 6. The second comparison, $w \times s - q_2 \times m < m$, matches the condition verified by step 8. It follows that, if the algorithm passes step 8, it must be the case that the value supplied for q_2 is correct.

By plugging $w \times s$, q_2 and m into the integer division remainder definition, we obtain the identity $w \times s \bmod m = w \times s - q_2 \times m$. Trivial substitution reveals that the computations in steps 5 and 7 result in $z = w \times s - q_2 \times m$, which allows us to conclude that $z = w \times s \bmod m$.

In the analysis for steps 1 – 4, we have proven that $w = s^2 \bmod m$. By substituting this into the above identity, we obtain the proof that the algorithm's output is indeed the desired signed message.

$$\begin{aligned} z &= w \times s \bmod m \\ &= (s^2 \bmod m) \times s \bmod m \\ &= s^2 \times s \bmod m \\ &= s^3 \bmod m \end{aligned}$$

Implementation Requirements

The main advantage of the algorithm in Figure C-37 is that it relies on the implementation of very few arithmetic operations on large integers. The maximum integer size that needs to be handled is twice the size of the modulus in the RSA key used to generate the signature.

Steps 1 and 5 use large integer multiplication. Steps 3 and 7 use integer subtraction. Steps 2, 4, 6, and 8 use large integer comparison. The checks in steps 2 and 6 guarantee that the results of the subtractions performed in steps 3 and 7 will be non-negative. It follows that the algorithm will never encounter negative numbers.

C.5.6 SGX Security Properties

We have summarized SGX's programming model and the implementation details that are publicly documented in Intel's official documentation and published patents. We are now ready to bring this the information together in an analysis of SGX's security properties. We start the analysis by restating SGX's security guarantees, and spend the bulk of this section discussing how SGX fares when pitted against the attacks described in § B. We conclude the analysis with some troubling implications of SGX's lack of resistance to software side-channel attacks.

Overview

Intel's Software Guard Extensions (SGX) is Intel's latest iteration of a trusted hardware solution to the secure remote computation problem. The SGX design is centered around the ability to create an isolated container whose contents receives special hardware protections that are intended to translate into privacy, integrity, and freshness guarantees.

An enclave's initial contents is loaded by the system software on the computer, and therefore cannot contain secrets in plain text. Once initialized, an enclave is expected to participate in a software attestation process, where it authenticates itself to a remote server. Upon successful authentication, the remote server is expected to disclose some secrets to an enclave over a secure communication channel. The SGX design attempts to guarantee that the measurement presented during software attestation accurately represents the contents loaded into the enclave.

SGX also offers a certificate-based identity system that can be used to migrate secrets between enclaves that have certificates issued by the same authority. The migration process involves securing the secrets via authenticated encryption before handing them off to the untrusted system software, which passes them to another enclave that can decrypt them.

The same mechanism used for secret migration can also be used to cache the secrets obtained via software attestation in an untrusted storage medium managed by system software. This caching can reduce the number of times that the software attestation process needs to be performed in a distributed system. In fact, SGX's software attestation process

is implemented by enclaves with special privileges that use the certificate-based identity system to securely store the CPU's attestation key in untrusted memory.

Physical Attacks

We begin by discussing SGX's resilience to the physical attacks described in § B.4. Unfortunately, this section is set to disappoint readers expecting definitive statements. The lack of publicly available details around the hardware implementation aspects of SGX precludes any rigorous analysis. However, we do know enough about SGX's implementation to point out a few avenues for future exploration.

Due to insufficient documentation, one can only hope that the SGX security model is not trivially circumvented by a port attack (§ B.4.1). We are particularly concerned about the Generic Debug eXternal Connection (GDXC) [209, 129], which collects and filters the data transferred by the uncore's ring bus (§ A.11.3), and reports it to an external debugger.

The SGX memory protection measures are implemented at the core level, in the Page Miss Handler (PMH, § A.11.5) (§ C.5.2) and at the chip die level, in the memory controller (§ C.5.1). Therefore, the code and data inside enclaves is stored in plaintext in on-chip caches (§ A.11), which entails that the enclave contents travels without any cryptographic protection on the uncore's ring bus (§ A.11.3).

Fortunately, a recent Intel patent [172] indicates that Intel engineers are tackling at least some classes of attacks targeting debugging ports.

The SDM and SGX papers discuss the most obvious class of bus tapping attacks (§ B.4.2), which is the DRAM bus tapping attack. SGX's threat model considers DRAM and the bus connecting it to the CPU chip to be untrusted. Therefore, SGX's Memory Encryption Engine (MEE, § C.5.1) provides privacy, integrity and freshness guarantees to the Enclave Page Cache (EPC, § C.4.1) data while it is stored in DRAM.

However, both the SGX papers and the ISCA 2015 tutorial on SGX admit that the MEE does not protect the addresses of the DRAM locations accessed when cache lines holding EPC data are evicted or loaded. This provides an opportunity for a malicious computer owner to observe an enclave's memory access patterns by combining a DRAM address line bus tap with carefully crafted system software that creates artificial pressure on the last-level

cache (LLC, § A.11) lines that hold the enclave's EPC pages.

On a brighter note, as mentioned in § B.4.2, we are not aware of any successful DRAM address line bus tapping attack. Furthermore, SGX is vulnerable to cache timing attacks that can be carried out completely in software, so malicious computer owners do not need to bother setting up a physical attack to obtain an enclave's memory access patterns.

While the SGX documentation addresses DRAM bus tapping attacks, it makes no mention of the System Management bus (SMBus, § A.9.2) that connects the Intel Management Engine (ME, § A.9.2) to various components on the computer's motherboard.

In § C.5.6, we will explain that the ME needs to be taken into account when evaluating SGX's memory protection guarantees. This makes us concerned about the possibility of an attack that taps the SMBus to reach into the Intel ME. The SMBus is much more accessible than the DRAM bus, as it has fewer wires that operate at a significantly lower speed. Unfortunately, without more information about the role that the Intel ME plays in a computer, we cannot move beyond speculation on this topic.

The threat model stated by the SGX design excludes physical attacks targeting the CPU chip (§ B.4.3). Fortunately, Intel's patents disclose an array of countermeasures aimed at increasing the cost of chip attacks.

For example, the original SGX patents [142, 112] disclose that the Fused Seal Key and the Provisioning Key, which are stored in e-fuses (§ C.4.8), are encrypted with a *global wrapping logic key* (GWK). The GWK is a 128-bit AES key that is hard-coded in the processor's circuitry, and serves to increase the cost of extracting the keys from an SGX-enabled processor.

As explained in § B.4.3, e-fuses have a large feature size, which makes them relatively easy to "read" using a high-resolution microscope. In comparison, the circuitry on the latest Intel processors has a significantly smaller feature size, and is more difficult to reverse engineer. Unfortunately, the GWK is shared among all the chip dies created from the same mask, so it has all the drawbacks of global secrets explained in § B.4.3.

Newer Intel patents [71, 72] describe SGX-enabled processors that employ a *Physical Unclonable Function* (PUF), e.g., [181], [137], which generates a symmetric key that is used during the provisioning process.

Specifically, at an early provisioning stage, the PUF key is encrypted with the GWK and transmitted to the key generation server. At a later stage, the key generation server encrypts the key material that will be burned into the processor chip's e-fuses with the PUF key, and transmits the encrypted material to the chip. The PUF key increases the cost of obtaining a chip's fuse key material, as an attacker must compromise both provisioning stages in order to be able to decrypt the fuse key material.

As mentioned in previous sections, patents reveal design possibilities considered by the SGX engineers. However, due to the length of timelines involved in patent applications, patents necessarily describe earlier versions of the SGX implementation plans, which might not match the shipping implementation. We expect this might be the case with the PUF provisioning patents, as it makes little sense to include a PUF in a chip die and rely on e-fuses and a GWK to store SGX's root keys. Deriving the root keys from the PUF would be more resilient to chip imaging attacks.

SGX's threat model excludes power analysis attacks (§ B.4.4) and other side-channel attacks. This is understandable, as power attacks cannot be addressed at the architectural level. Defending against power attacks requires expensive countermeasures at the lowest levels of hardware implementation, which can only be designed by engineers who have deep expertise in both system security and Intel's manufacturing process. It follows that defending against power analysis attacks has a very high cost-to-benefit ratio.

Privileged Software Attacks

The SGX threat model considers system software to be untrusted. This is a prerequisite for SGX to qualify as a solution to the secure remote computation problem encountered by software developers who wish to take advantage of Infrastructure-as-a-Service (IaaS) cloud computing.

SGX's approach is also an acknowledgement of the realities of today's software landscape, where the system software that runs at high privilege levels (§ A.3) is so complex that security researchers constantly find vulnerabilities in it (§ B.5).

The SGX design prevents malicious software from directly reading or from modifying the EPC pages that store an enclave's code and data. This security property relies on two

pillars in the SGX design.

First, the SGX implementation (§ C.5.1) runs in the processor's microcode (§ A.14), which is effectively a higher privilege level that system software does not have access to. Along the same lines, SGX's security checks (§ C.5.2) are the last step performed by the PMH, so they cannot be bypassed by any other architectural feature.

This implementation detail is only briefly mentioned in SGX's official documentation, but has a large impact on security. For context, Intel's Trusted Execution Technology (TXT, [74]), which is the predecessor of SGX, relied on Intel's Virtual Machine Extensions (VMX) for isolation. The approach was unsound, because software running in System Management Mode (SMM, § A.3) could bypass the restrictions used by VMX to provide isolation.

The security properties of SGX's memory protection mechanisms are discussed in detail in § C.5.6.

Second, SGX's microcode is always involved when a CPU transitions between enclave code and non-enclave code (§ C.4.4), and therefore regulates all interactions between system software and an enclave's environment.

On enclave entry (§ C.4.4), the SGX implementation sets up the registers (§ A.2) that make up the execution state (§ A.6) of the logical processor (LP § A.9.4), so a malicious OS or hypervisor cannot induce faults in the enclave's software by tampering with its execution environment.

When an LP transitions away from an enclave's code due to a hardware exception (§ A.8.2), the SGX implementation stashes the LP's execution state into a State Save Area (SSA, § C.4.2) area inside the enclave and scrubs it, so the system software's exception handler cannot access any enclave secrets that may be stored in the execution state.

The protections described above apply to all the levels of privileged software. SGX's transitions between an enclave's code and non-enclave code place SMM software on the same footing as the system software at lower privilege levels. System Management Interrupts (SMI, § A.12, § B.5), which cause the processor to execute SMM code, are handled using the same Asynchronous Enclave Exit (AEX, § C.4.4) process as all other hardware exceptions.

Reasoning about the security properties of SGX's transitions between enclave mode and non-enclave mode is very difficult. A correctness proof would have to take into account all the CPU's features that expose registers. Difficulty aside, such a proof would be very short-lived, because every generation of Intel CPUs tends to introduce new architectural features. The paragraph below gives a taste of what such a proof would look like.

`EENTER` (§ C.4.4) stores the `RSP` and `RBP` register values in the `SSA` used to enter the enclave, but stores `XCR0` (§ A.6), `FS` and `GS` (§ A.7) in the non-architectural area of the `TCS` (§ C.5.1). At first glance, it may seem elegant to remove this inconsistency and have `EENTER` store the contents of the `XCR0`, `FS`, and `GS` registers in the current `SSA`, along with `RSP` and `RBP`. However, this approach would break the Intel architecture's guarantees that only system software can modify `XCR0`, and application software can only load segment registers using selectors that index into the `GDT` or `LDT` set up by system software. Specifically, a malicious application could modify these privileged registers by creating an enclave that writes the desired values to the current `SSA` locations backing up the registers, and then executes `EEXIT` (§ C.4.4).

Unfortunately, the following sections will reveal that while SGX offers rather thorough guarantees against straightforward attacks on enclaves, its guarantees are almost non-existent when it comes to more sophisticated attacks, such as side-channel attacks. This section concludes by describing what might be the most egregious side-channel vulnerability in SGX.

Most modern Intel processors feature hyper-threading. On these CPUs, the execution units (§ A.10) and caches (§ A.11) on a core (§ A.9.4) are shared by two LPs, each of which has its own execution state. SGX does not prevent hyper-threading, so malicious system software can schedule a thread executing the code of a victim enclave on an LP that shares the core with an LP executing a snooping thread. This snooping thread can use the processor's high-resolution performance counter [156], in conjunction with microarchitectural knowledge of the CPU's execution units and out-of-order scheduler, to learn the instructions executed by the victim enclave, as well as its memory access patterns.

This vulnerability can be fixed using two approaches. The straightforward solution is to require cloud computing providers to disable hyper-threading when offering SGX. The SGX

enclave measurement would have to be extended to include the computer's hyper-threading configuration, so the remote parties in the software attestation process can be assured that their enclaves are hosted by a secure environment.

A more complex approach to fixing the hyper-threading vulnerability would entail having the SGX implementation guarantee that when an LP is executing an enclave's code, the other LP sharing its core is either inactive, or is executing the same enclave's code. While this approach is possible, its design would likely be quite cumbersome.

Memory Mapping Attacks

§ C.4.4 explained that the code running inside an enclave uses the same address translation process (§ A.5) and page tables as its host application. While this design approach makes it easy to retrofit SGX support into existing codebases, it also enables the address translation attacks described in § B.7.

The SGX design protects the code inside enclaves against the active attacks described in § B.7. These protections have been extensively discussed in prior sections, so we limit ourselves to pointing out SGX's answer to each active attack. We also explain the lack of protections against passive attacks, which can be used to learn an enclave's memory access pattern at 4KB page granularity.

SGX uses the Enclave Page Cache Map (EPCM, § C.4.1) to store each EPC page's position in its enclave's virtual address space. The EPCM is consulted by SGX's extensions to the Page Miss Handler (PMH, § C.5.2), which prevent straightforward active address translation attacks (§ B.7.2) by rejecting undesirable address translations before they reach the TLB (§ A.11.5).

SGX allows system software to evict (§ C.4.5) EPC pages into untrusted DRAM, so that the EPC can be over-subscribed. The contents of the evicted pages and the associated EPCM metadata are protected by cryptographic primitives that offer privacy, integrity and freshness guarantees. This protects against the active attacks using page swapping described in § B.7.3.

When system software wishes to evict EPC pages, it must follow the process described in § C.4.5, which guarantees to the SGX implementation that all the LPs have invalidated

any TLB entry associated with pages that will be evicted. This defeats the active attacks based on stale TLB entries described in § B.7.4.

§ C.5.3 outlines a correctness proof for the memory protection measures described above.

Unfortunately, SGX does not protect against passive address translation attacks (§ B.7.1), which can be used to learn an enclave's memory access pattern at page granularity. While this appears benign, recent work [204] demonstrates the use of these passive attacks in a few practical settings, which are immediately concerning for image processing applications.

The rest of this section describes the theory behind planning a passive attack against an SGX enclave. The reader is directed to [204] for a fully working system.

Passive address translation attacks rely on the fact that memory accesses issued by SGX enclaves go through the Intel architecture's address translation process (§ A.5), including delivering page faults (§ A.8.2) and setting the accessed (A) and dirty (D) attributes (§ A.5.3) on page table entries.

A malicious OS kernel or hypervisor can obtain the page-level trace of an application executing inside an enclave by setting the present (P) attribute to 0 on all the enclave's pages before starting enclave execution. While an enclave executes, the malicious system software maintains exactly one instruction page and one data page present in the enclave's address space.

When a page fault is generated, CR2 contains the virtual address of a page accessed by enclave, and the error code indicates whether the memory access was a read or a write (bit 1) and whether the memory access is a data access or an instruction fetch access (bit 4). On a data access, the kernel tracing the enclave code's memory access pattern would set the P flag of the desired page to 1, and set the P flag of the previously accessed data page to 0. Instruction accesses can be handled in a similar manner.

For a slightly more detailed trace, the kernel can set a desired page's writable (W) attribute to 0 if the page fault's error code indicates a read access, and only set it to 1 for write accesses. Also, applications that use a page as both code and data (self-modifying code and just-in-time compiling VMs) can be handled by setting a page's disable execution (XD) flag to 0 for a data access, and by carefully accounting for the case where the last

accessed data page is the same as the last accessed code page.

Leaving an enclave via an Asynchronous Enclave Exit (AEX, § C.4.4) and re-entering the enclave via `ERESUME` (§ C.4.4) causes the CPU to flush TLB entries that contain enclave addresses, so a tracing kernel would not need to worry about flushing the TLB. The tracing kernel does not need to flush the caches either, because the CPU needs to perform address translation even for cached data.

A straightforward way to reduce this attack's power is to increase the page size, so the trace contains less information. However, the attack cannot be completely prevented without removing the kernel's ability to oversubscribe the EPC, which is a major benefit of paging.

Software Attacks on Peripherals

Since the SGX design does not trust the system software, it must be prepared to withstand the attacks described in § B.6, which can be carried out by the system software thanks to its ability to control peripheral devices on the computer's motherboard (§ A.9.1). This section summarizes the security properties of SGX when faced with these attacks, based on publicly available information.

When SGX is enabled on an LP, it configures the memory controller (MC, § A.11.3) integrated on the CPU chip die to reject any DMA transfer that falls within the Processor Reserved Memory (PRM, § C.4.1) range. The PRM includes the EPC, so the enclaves' contents is protected from the PCI Express attacks described in § B.6.1. This protection guarantee relies on the fact that the MC is integrated on the processor's chip die, so the MC configuration commands issued by SGX's microcode implementation (§ C.5.1) are transmitted over a communication path that never leaves the CPU die, and therefore can be trusted.

SGX regards DRAM as an untrusted storage medium, and uses cryptographic primitives implemented in the MEE to guarantee the privacy, integrity and freshness of the EPC contents that is stored into DRAM. This protects against software attacks on DRAM's integrity, like the rowhammer attack described in § B.6.2.

The SDM describes an array of measures that SGX takes to disable processor features intended for debugging when a LP starts executing an enclave's code. For example, enclave

entry (§ C.4.4) disables Precise Event Based Sampling (PEBS) for the LP, as well as any hardware breakpoints placed inside the enclave's virtual address range (ELRANGE, § C.4.2). This addresses some of the attacks described in § B.6.3, which take advantage of performance monitoring features to get information that typically requires access to hardware probes.

At the same time, the SDM does not mention anything about uncore PEBS counters, which can be used to learn about an enclave's LLC activity. Furthermore, the ISCA 2015 tutorial slides mention that **SGX does not protect against software side-channel attacks** that rely on performance counters.

This limitation in SGX's threat model leaves security-conscious enclave authors in a rather terrible situation. These authors know that SGX does not protect their enclaves against a class of software attacks. At the same time, they cannot even contemplate attempting to defeat these attacks on their own, due to lack of information. Specifically, the documentation that is publicly available from Intel does not provide enough information to model the information leakage due to performance counters.

For example, Intel does not document the mapping implemented in CBoxes (§ A.11.3) between physical DRAM addresses and the LLC slices used to cache the addresses. This mapping impacts several uncore performance counters, and the impact is strong enough to allow security researches to reverse-engineer the mapping [88, 139, 206]. Therefore, it is safe to assume that a malicious computer owner who knows the CBox mapping can use the uncore performance counters to learn about an enclave's memory access patterns.

The SGX papers mention that SGX's threat model includes attacks that overwrite the flash memory chip that stores the computer's firmware, which result in malicious code running in SMM. However, all the official SGX documentation is silent about the implications of an attack that compromises the firmware executed by the Intel ME.

§ B.6.4 states that the ME's firmware is stored in the same flash memory as the boot firmware, and enumerates some of ME's special privileges that enable it to help system administrators remotely diagnose and fix hardware and software issues. Given that the SGX design is concerned about the possibility of malicious computer firmware, it is reasonable to be concerned about malicious ME firmware.

§ B.6.4 argues that an attacker who compromises the ME can carry out actions that are

usually classified as physical attacks. An optimistic security researcher can observe that the most scary attack vector afforded by an ME takeover appears to be direct DRAM access, and SGX already assumes that the DRAM is untrusted. Therefore, an ME compromise would be equivalent to the DRAM attacks analyzed in § C.5.6.

However, we are troubled by the lack of documentation on the ME's implementation, as certain details are critical to SGX's security analysis. For example, the ME is involved in the computer's boot process (§ A.13, § A.14.4), so it is unclear if it plays any part in the SGX initialization sequence. Furthermore, during the security boot stage (SEC, § A.13.2), the bootstrap LP (BSP) is placed in Cache-As-Ram (CAR) mode so that the PEI firmware can be stored securely while it is measured. This suggests that it would be convenient for the ME to receive direct access to the CPU's caches, so that the ME's TPM implementation can measure the firmware directly. At the same time, a special access path from the ME to the CPU's caches might sidestep the MEE, allowing an attacker who has achieved ME code execution to directly read the EPC's contents.

Cache Timing Attacks

The SGX threat model excludes the cache timing attacks described in § B.8. The SGX documentation bundles these attacks together with other side-channel attacks and summarily dismisses them as complex physical attacks. However, cache timing attacks can be mounted entirely by unprivileged software running at ring 3. This section describes the implications of SGX's environment and threat model on cache timing attacks.

The main difference between SGX and a standard architecture is that SGX's threat model considers the system software to be untrusted. As explained earlier, this accurately captures the situation in remote computation scenarios, such as cloud computing. SGX's threat model implies that the system software can be carrying out a cache timing attack on the software inside an enclave.

A malicious system software translates into significantly more powerful cache timing attacks, compared to those described in § B.8. The system software is in charge of scheduling threads on LPs, and also in charge of setting up the page tables used by address translation (§ A.5), which control cache placement (§ A.11.5).

For example, the malicious kernel set out to trace an enclave's memory access patterns described in § C.5.6 can improve the accuracy of a cache timing attack by using page coloring [119] principles to partition [133] the cache targeted by the attack. In a nutshell, the kernel divides the cache's sets (§ A.11.2) into two regions, as shown in Figure C-38.

The system software stores all the victim enclave's code and data in DRAM addresses that map to the cache sets in one of the regions, and stores its own code and data in DRAM addresses that map to the other region's cache sets. The snooping thread's code is assumed to be a part of the OS. For example, in a typical 256 KB (per-core) L2 cache organized as 512 8-way sets of 64-byte lines, the tracing kernel could allocate lines 0-63 for the enclave's code page, lines 64-127 for the enclave's data page, and use lines 128-511 for its own pages.

To the best of our knowledge, there is no minor modification to SGX that would provably defend against cache timing attacks. However, the SGX design could take a few steps to increase the cost of cache timing attacks. For example, SGX's enclave entry implementation could flush the core's private caches, which would prevent cache timing attacks from targeting them. This measure would defeat the cache timing attacks described below, and would only be vulnerable to more sophisticated attacks that target the shared LLC, such as [205, 135]. The description above assumes that multi-threading has been disabled, for the reasons explained in § C.5.6.

Barring the additional protection measures described above, a tracing kernel can extend the attack described in § C.5.6 with the steps outlined below to take advantage of cache timing and narrow down the addresses in an application's memory access trace to cache line granularity.

Right before entering an enclave via `EENTER` or `ERESUME`, the kernel would issue `CLFLUSH` instructions to flush the enclave's code page and data page from the cache. The enclave could have accessed a single code page and a single data page, so flushing the cache should be reasonably efficient. The tracing kernel then uses 16 bogus pages (8 for the enclave's code page, and 8 for the enclave's data page) to load all the 8 ways in the 128 cache sets allocated by enclave pages. After an `AEX` gives control back to the tracing kernel, it can read the 16 bogus pages, and exploit the time difference between an L2 cache hit and a miss to see which cache lines were evicted and replaced by the enclave's memory accesses.

An extreme approach that can provably defeat cache timing attacks is disabling caching for the PRM range, which contains the EPC. The SDM is almost completely silent about the PRM, but the SGX manuals that it is based on state that the allowable caching behaviors (§ A.11.4) for the PRM range are uncacheable (UC) and write-back (WB). This could become useful if the SGX implementation would make sure that the PRM's caching behavior cannot be changed while SGX is enabled, and if the selected behavior would be captured by the enclave's measurement (§ C.4.6).

Software Side-Channel Attacks and SGX

The SGX design reuses a few terms from the Trusted Platform Module (TPM, § 2.4) design. This helps software developers familiar with TPM understand SGX faster. At the same time, the term reuse invites the assumption that SGX's software attestation is implemented in tamper-resistant hardware, similarly to the TPM design.

§ C.4.8 explains that, in fact, the SGX design delegates the creation of attestation signatures to software that runs inside a Quoting Enclave with special privileges that allows it to access the processor's attestation key. Re-stated, SGX includes an enclave whose software reads the attestation key and produces attestation signatures.

Creating the Quoting Enclave is a very elegant way of reducing the complexity of the hardware implementation of SGX, assuming that the isolation guarantees provided by SGX are sufficient to protect the attestation key. However, the security analysis in § C.5.6 reveals that enclaves are vulnerable to a vast array of software side-channel attacks, which have been demonstrated effective in extracting a variety of secrets from isolated environments.

The gaps in the security guarantees provided to enclaves place a large amount of pressure on Intel's software developers, as they must attempt to implement the EPID signing scheme used by software attestation without leaking any information. Intel's ISCA 2015 SGX tutorial slides suggest that the SGX designers will advise developers to write their code in a way that avoids data-dependent memory accesses, as suggested in § B.8.4, and perhaps provide analysis tools that detect code that performs data-dependent memory accesses.

The main drawback of the approach described above is that it is extremely cumbersome. § B.8.4 describes that, while it may be possible to write simple pieces of software in such a

way that they do not require data-dependent memory accesses, there is no known process that can scale this to large software systems. For example, each virtual method call in an object-oriented language results in data-dependent code fetches.

The ISCA 2015 SGX tutorial slides also suggest that the efforts of removing data-dependent memory accesses should focus on cryptographic algorithm implementations, in order to protect the keys that they handle. This is a terribly misguided suggestion, because cryptographic key material has no intrinsic value. Attackers derive benefits from obtaining the data that is protected by the keys, such as medical and financial records.

Some security researchers focus on protecting cryptographic keys because they are the target of today's attacks. Unfortunately, it is easy to lose track of the fact that keys are being attacked simply because they are the lowest hanging fruit. A system that can only protect the keys will have a very small positive impact, as the attackers will simply shift their focus on the algorithms that process the valuable information, and use the same software side-channel attacks to obtain that information directly.

The second drawback of the approach described towards the beginning of this section is that while eliminating data-dependent memory accesses should thwart the attacks described in § C.5.6 and § C.5.6, the measure may not be sufficient to prevent the hyper-threading attacks described in § C.5.6. The level of sharing between the two logical processors (LP, § A.9.4) on the same CPU core is so high that it is possible that a snooping LP can learn more than the memory access pattern from the other LP on the same core.

For example, if the number of cycles taken by an integer ALU to execute a multiplication or division micro-op (§ A.10) depends on its inputs, the snooping LP could learn some information about the numbers multiplied or divided by the other LP. While this may be a simple example, it is safe to assume that the Quoting Enclave will be studied by many motivated attackers, and that any information leak will be exploited.

C.6 Conclusion

Shortly after we learned about Intel's Software Guard Extensions (SGX) initiative, we set out to study it in the hope of finding a practical solution to its vulnerability to cache timing

attacks. After reading the official SGX manuals, we were left with more questions than when we started. The SGX patents filled some of the gaps in the official documentation, but also revealed Intel's enclave licensing scheme, which has troubling implications.

After learning about the SGX implementation and inferring its design constraints, we discarded our draft proposals for defending enclave software against cache timing attacks. We concluded that it would be impossible to claim to provide this kind of guarantee given the design constraints and all the unknowns surrounding the SGX implementation. Instead, we applied the knowledge that we gained to design Sanctum [40], which is briefly described in § 2.9.

This paper describes our findings while studying SGX. We hope that it will help fellow researchers understand the breadth of issues that need to be considered before accepting a trusted hardware design as secure. We also hope that our work will prompt the research community to expect more openness from the vendors who ask us to trust their hardware.

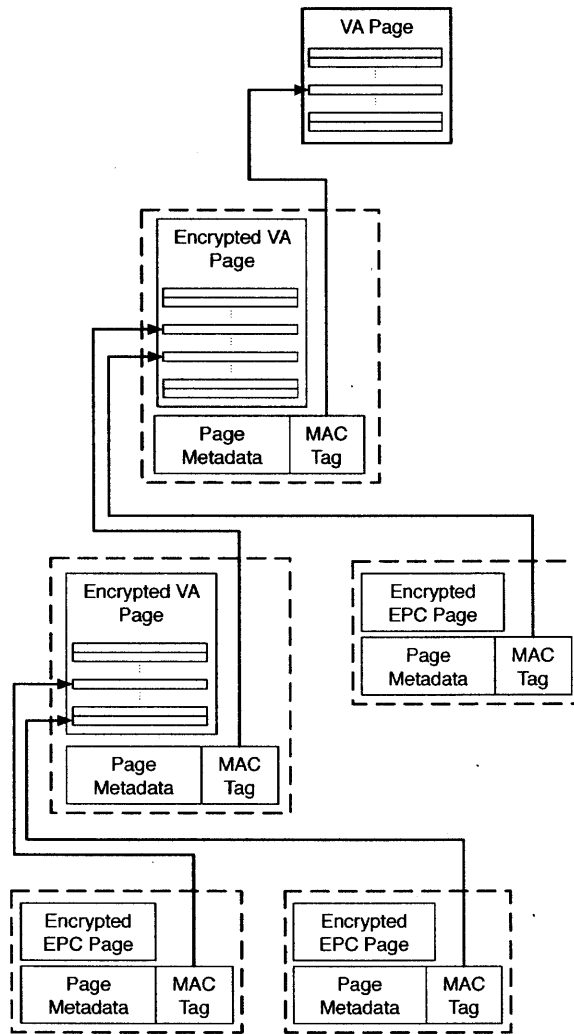


Figure C-18: A version tree formed by evicted VA pages and enclave EPC pages. The enclave pages are leaves, and the VA pages are inner nodes. The OS controls the tree's shape, which impacts the performance of evictions, but not their correctness.

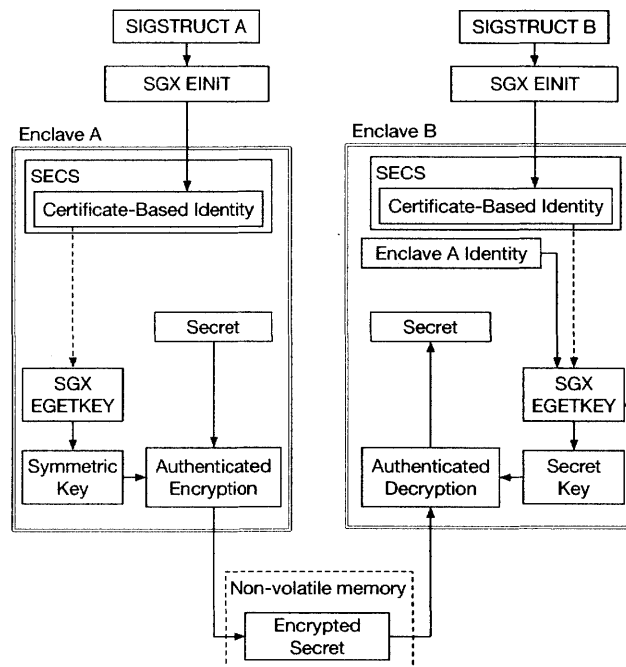


Figure C-19: SGX has a certificate-based enclave identity scheme, which can be used to migrate secrets between enclaves that contain different versions of the same software module. Here, enclave A's secrets are migrated to enclave B.

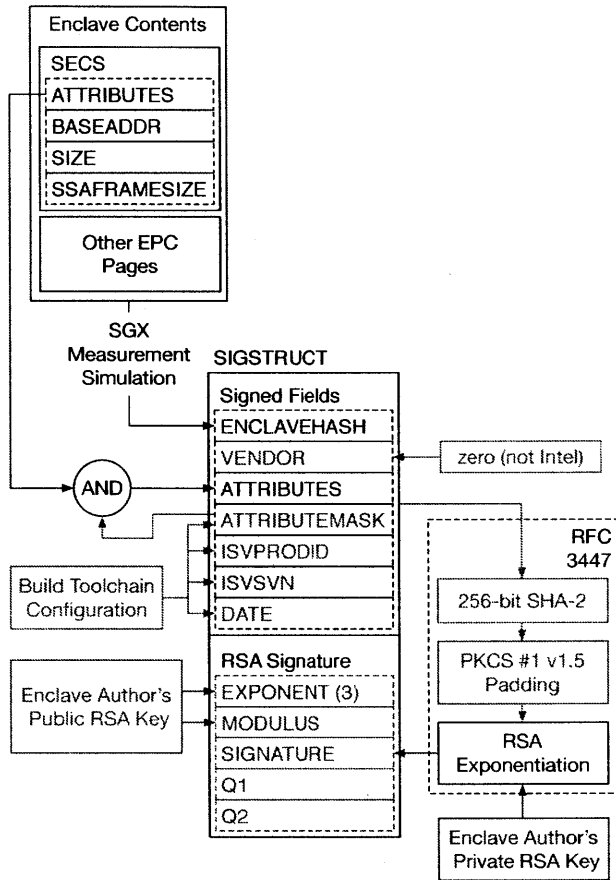


Figure C-20: An enclave's Signature Structure (SIGSTRUCT) is intended to be generated by an enclave building toolchain that has access to the enclave author's private RSA key.

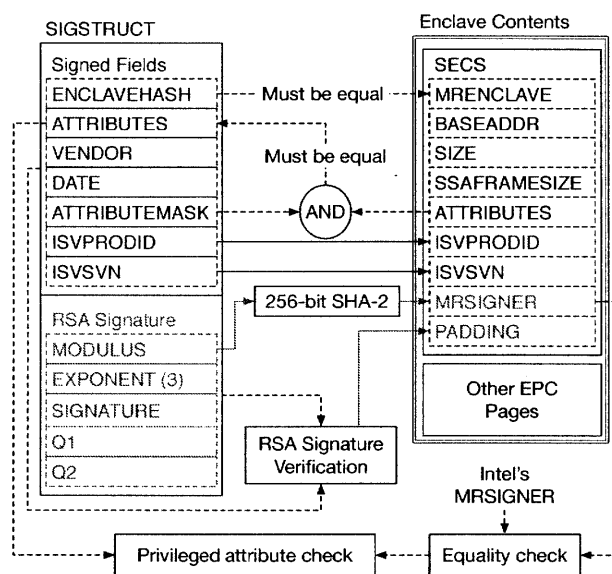


Figure C-21: EINIT verifies the RSA signature in the enclave's certificate. If the certificate is valid, the information in it is used to populate the SECS fields that make up the enclave's certificate-based identity.

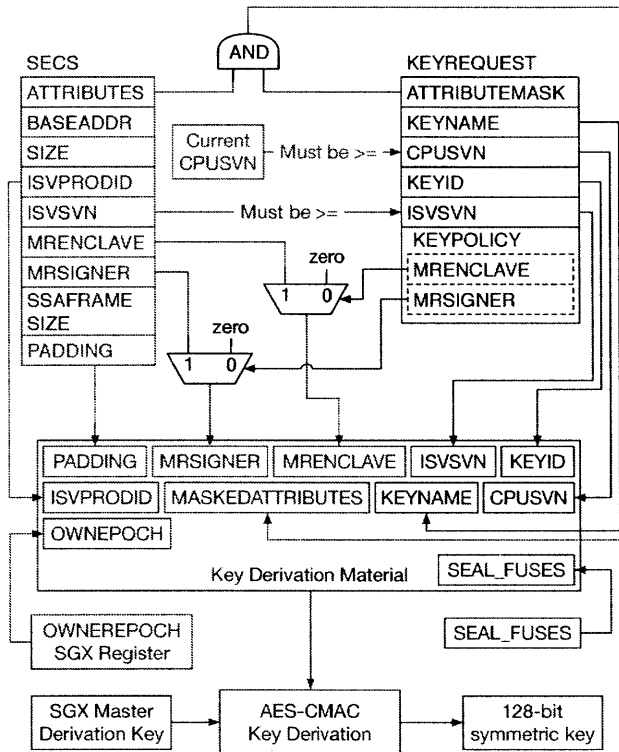


Figure C-22: EGETKEY implements a key derivation service that is primarily used by SGX's secret migration feature. The key derivation material is drawn from the SECS of the calling enclave, the information in a Key Request structure, and secure storage inside the CPU's hardware.

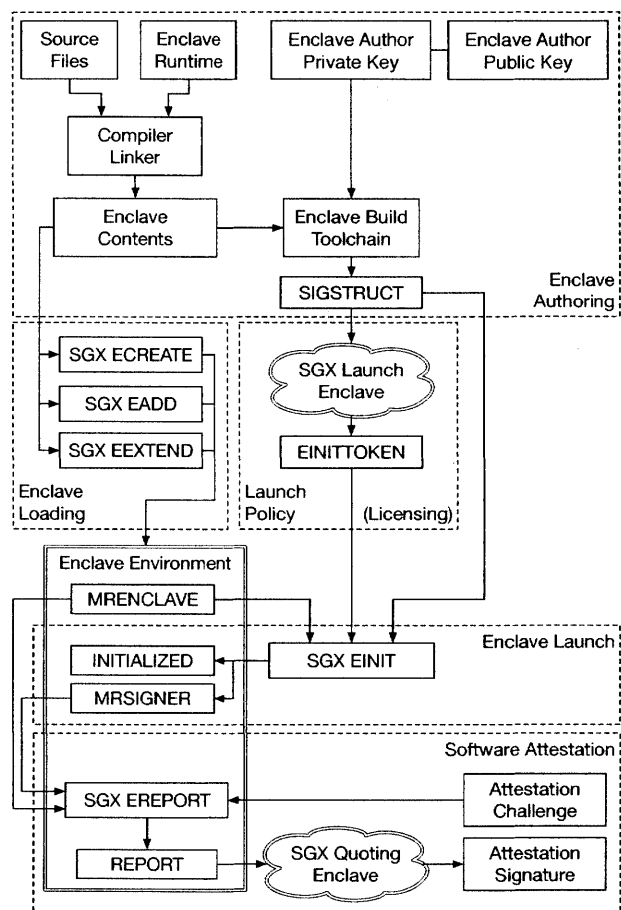


Figure C-23: Setting up an SGX enclave and undergoing the software attestation process involves the SGX instructions EINIT and EREPORT, and two special enclaves authored by Intel, the SGX Launch Enclave and the SGX Quoting Enclave.

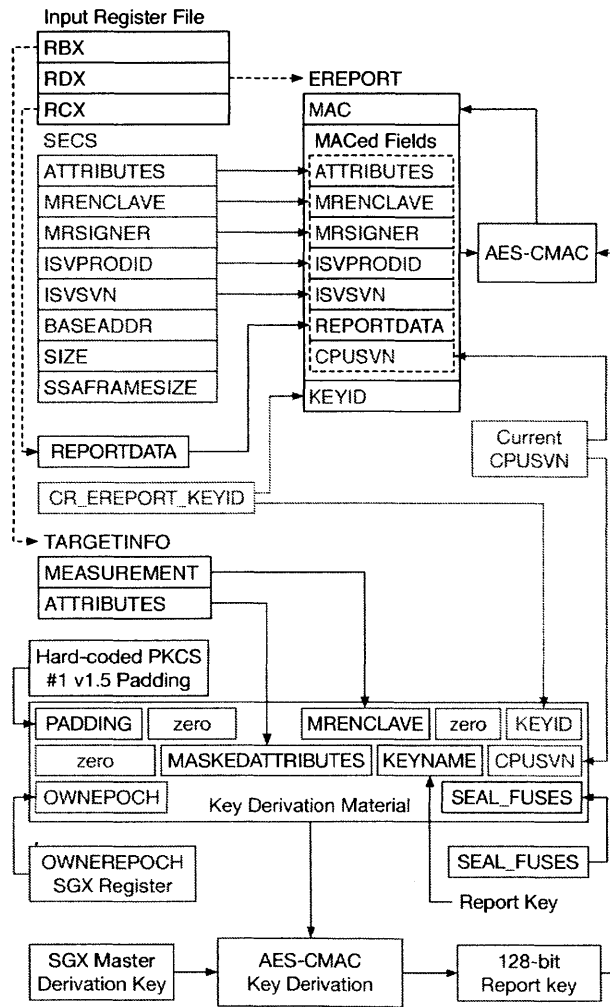


Figure C-24: EREPORT data flow

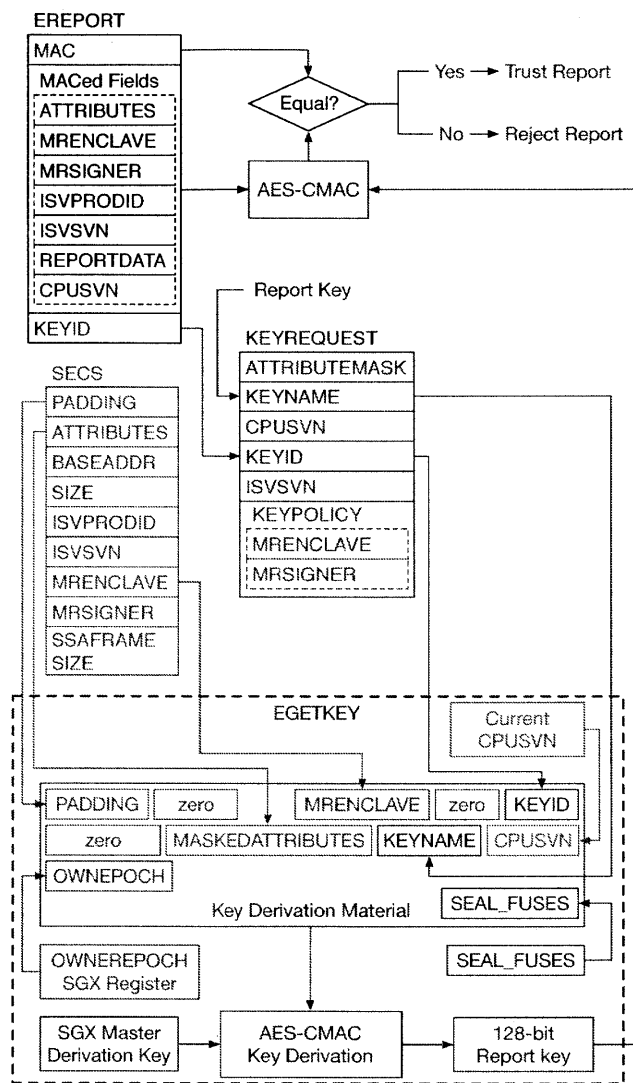


Figure C-25: The authenticity of the REPORT structure created by EREPORT can and should be verified by the report's target enclave. The target's code uses EGETKEY to obtain the key used for the MAC tag embedded in the REPORT structure, and then verifies the tag.

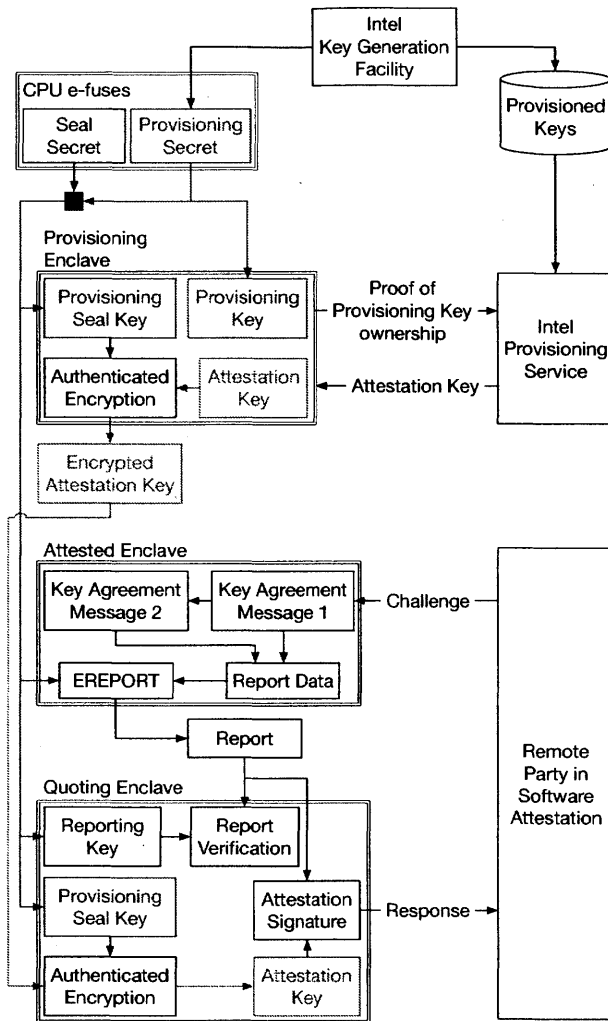


Figure C-26: SGX's software attestation is based on two secrets stored in e-fuses inside the processor's die, and on a key received from Intel's provisioning service.

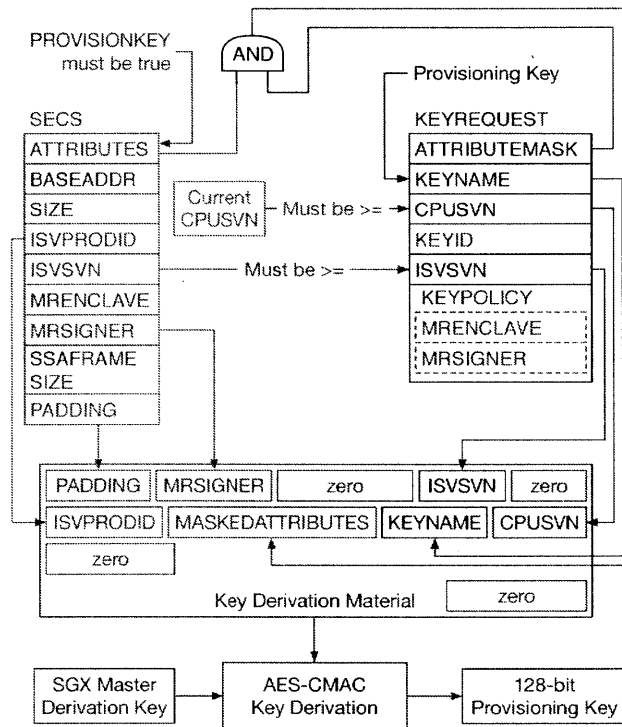


Figure C-27: When EGETKEY is asked to derive a Provisioning key, it does not use the Seal Secret or OWNEREPOCH. The Provisioning key does, however, depend on MRSIGNER and on the SVN of the SGX implementation.

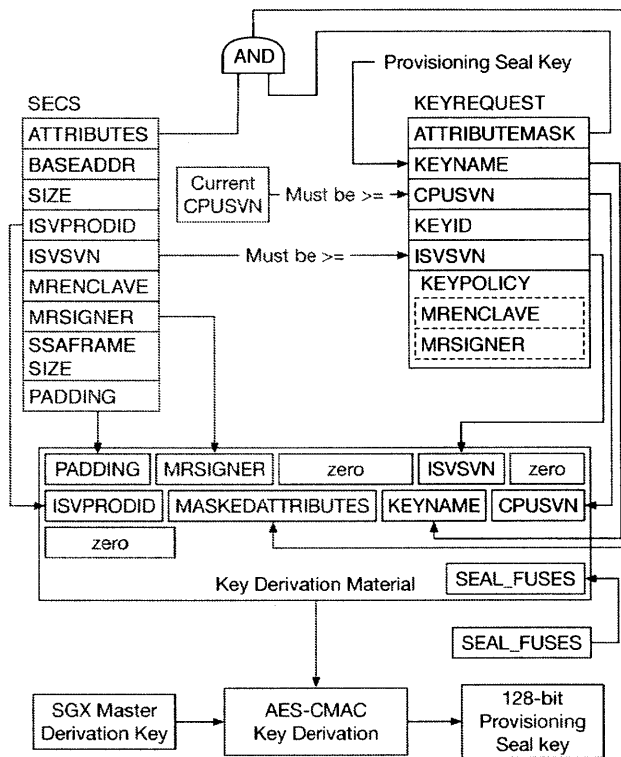


Figure C-28: The derivation material used to produce Provisioning Seal keys does not include the OWNEREPOCH value, so the keys survive computer ownership changes.

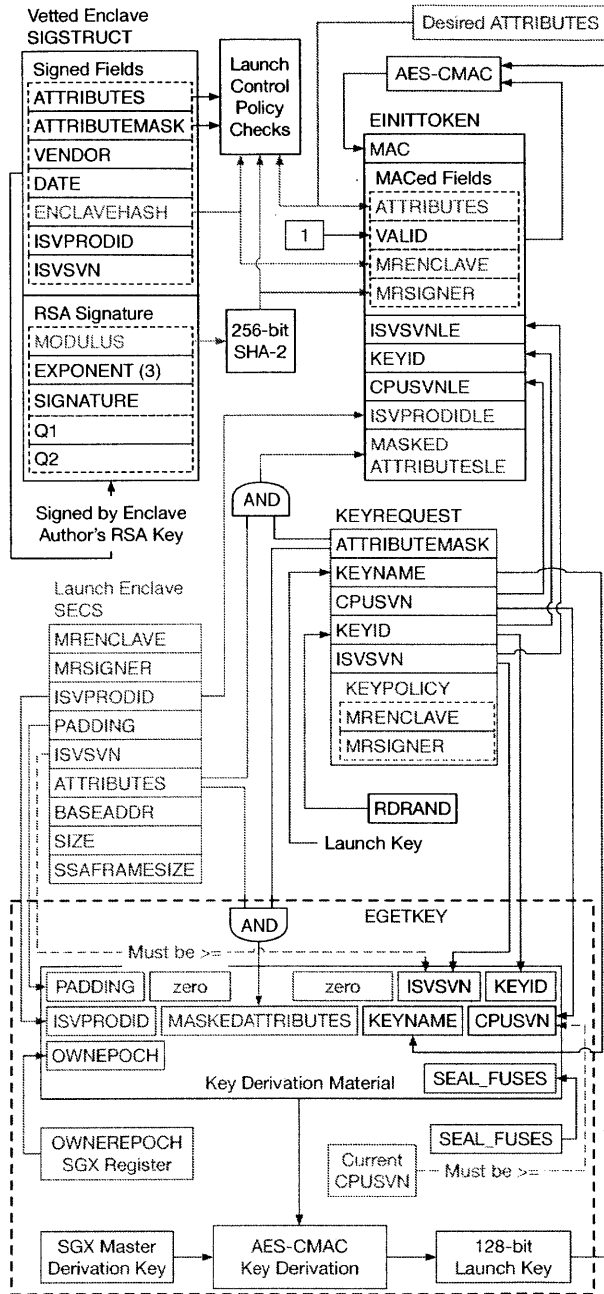


Figure C-29: The SGX Launch Enclave computes the EINITTOKEN.

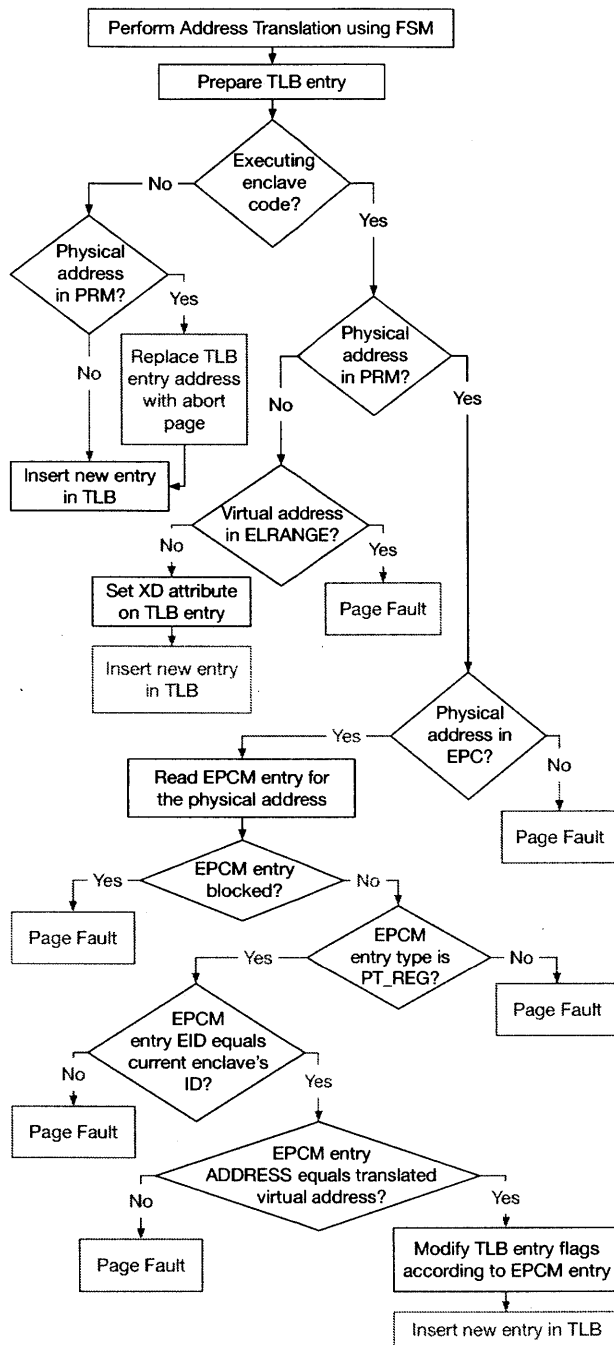


Figure C-30: SGX adds a few security checks to the PMH. The checks ensure that all the TLB entries created by the address translation unit meet SGX's memory access restrictions.

ECREATE(*SECS*)

▷ Initialize the SECS state used for tracking.

- 1 *SECS* . *tracking* ← FALSE
- 2 *SECS* . *done-tracking* ← FALSE
- 3 *SECS* . *active-threads* ← 0
- 4 *SECS* . *tracked-threads* ← 0
- 5 *SECS* . *lp-mask* ← 0

Figure C-31: The algorithm used to initialize the SECS fields used by the TLB flush tracking method presented in this section.

ETRACK(*SECS*)

▷ Abort if tracking is already active.

- 1 **if** *SECS* . *tracking* = TRUE
- 2 **then return** SGX-PREV-TRK-INCMPL

▷ Activate TLB flush tracking.

- 3 *SECS* . *tracking* ← TRUE
- 4 *SECS* . *done-tracking* ← FALSE
- 5 *SECS* . *tracked-threads* ←
 ATOMIC-READ(*SECS* . *active-threads*)
- 6 **for** *i* ← 0 **to** MAX-LP-ID
- 7 **do** ATOMIC-CLEAR(*SECS* . *lp-mask*[*i*])

Figure C-32: The algorithm used by ETRACK to activate TLB flush tracking.

ENCLAVE-EXIT(*SECS*)

▷ Track an enclave exit.

- 1 ATOMIC-DECREMENT(*SECS* . *active-threads*)
- 2 **if** ATOMIC-TEST-AND-SET(
 SECS . *lp-mask*[LP-ID])
- 3 **then** ATOMIC-DECREMENT(
 SECS . *tracked-threads*)
- 4 **if** *SECS* . *tracked-threads* = 0
- 5 **then** *SECS* . *done-tracking* ← TRUE

Figure C-33: The algorithm that updates the TLB flush tracking state when an LP exits an enclave via EEXIT or AEX.

```

ENCLAVE-ENTER(SECS)
  ▷ Track an enclave entry.
  1 ATOMIC-INCREMENT(SECS. active-threads)
  2 ATOMIC-SET(SECS. lp-mask[LP-ID])

```

Figure C-34: The algorithm that updates the TLB flush tracking state when an LP enters an enclave via EENTER or ERESUME.

```

EWB-VERIFY(virtual-addr)
  1 physical-addr ← TRANSLATE(virtual-addr)
  2 epcm-slot ← EPCM-SLOT(physical-addr)
  3 if EPCM[slot]. BLOCKED = FALSE
  4   then return SGX-NOT-BLOCKED
  5 SECS ← EPCM-ADDR(
      EPCM[slot]. ENCLAVESECS)
  ▷ Verify that the EPC page can be evicted.
  6 if SECS. tracking = FALSE
  7   then return SGX-NOT-TRACKED
  8 if SECS. done-tracking = FALSE
  9   then return SGX-NOT-TRACKED

```

Figure C-35: The algorithm that ensures that all LPs running an enclave’s code when ETRACK was executed have exited enclave mode at least once.

```

EBLOCK(virtual-addr)
  1 physical-addr ← TRANSLATE(virtual-addr)
  2 epcm-slot ← EPCM-SLOT(physical-addr)
  3 if EPCM[slot]. BLOCKED = TRUE
  4   then return SGX-BLKSTATE
  5 if SECS. tracking = TRUE
  6   then if SECS. done-tracking = FALSE
  7     then return SGX-ENTRYEPOCH-LOCKED
  8     SECS. tracking ← FALSE
  9 EPCM[slot]. BLOCKED ← TRUE

```

Figure C-36: The algorithm that marks the end of a TLB flushing cycle when EBLOCK is executed.

1. Compute $u \leftarrow s \times s$ and $v \leftarrow q_1 \times m$
2. If $u < v$, abort. q_1 must be incorrect.
3. Compute $w \leftarrow u - v$
4. If $w \geq m$, abort. q_1 must be incorrect.
5. Compute $x \leftarrow w \times s$ and $y \leftarrow q_2 \times m$
6. If $x < y$, abort. q_2 must be incorrect.
7. Compute $z \leftarrow x - y$.
8. If $z \geq m$, abort. q_2 must be incorrect.
9. Output z .

Figure C-37: An RSA signature verification algorithm specialized for the case where the public exponent is 3. s is the RSA signature and m is the RSA key modulus. The algorithm uses two additional inputs, q_1 and q_2 .

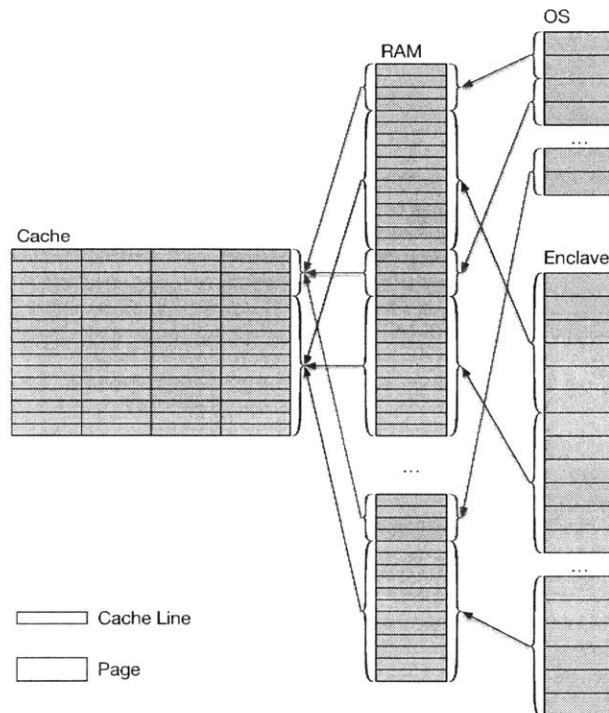


Figure C-38: A malicious OS can partition a cache between the software running inside an enclave and its own malicious code. Both the OS and the enclave software have cache sets dedicated to them. When allocating DRAM to itself and to the enclave software, the malicious OS is careful to only use DRAM regions that map to the appropriate cache sets. On a system with an Intel CPU, the OS can partition the L2 cache by manipulating the page tables in a way that is completely oblivious to the enclave's software.

Bibliography

- [1] *FIPS 140-2 Consolidated Validation Certificate No. 0003*. 2011.
- [2] *IBM 4765 Cryptographic Coprocessor Security Module - Security Policy*. Dec 2012.
- [3] Sha1 deprecation policy. <http://blogs.technet.com/b/pki/archive/2013/11/12/sha1-deprecation-policy.aspx>, 2013. [Online; accessed 4-May-2015].
- [4] 7-zip lzma benchmark: Intel haswell. <http://www.7-cpu.com/cpu/Haswell.html>, 2014. [Online; accessed 10-February-2015].
- [5] Bios freedom status. <https://puri.sm/posts/bios-freedom-status/>, Nov 2014. [Online; accessed 2-Dec-2015].
- [6] Gradually sunseting sha-1. <http://googleonlinesecurity.blogspot.com/2014/09/gradually-sunseting-sha-1.html>, 2014. [Online; accessed 4-May-2015].
- [7] Ipc2 hardware specification. <http://fit-pc.com/download/intense-pc2/documents/ipc2-hw-specification.pdf>, Sep 2014. [Online; accessed 2-Dec-2015].
- [8] Linux kernel: Cve security vulnerabilities, versions and detailed reports. http://www.cvedetails.com/product/47/Linux-Linux-Kernel.html?vendor_id=33, 2014. [Online; accessed 27-April-2015].
- [9] Nist's policy on hash functions. <http://csrc.nist.gov/groups/ST/hash/policy.html>, 2014. [Online; accessed 4-May-2015].
- [10] Xen: Cve security vulnerabilities, versions and detailed reports. http://www.cvedetails.com/product/23463/XEN-XEN.html?vendor_id=6276, 2014. [Online; accessed 27-April-2015].
- [11] Spec cpu 2006. Technical report, Standard Performance Evaluation Corporation, May 2015.
- [12] Xen project software overview. http://wiki.xen.org/wiki/Xen_Project_Software_Overview, 2015. [Online; accessed 27-April-2015].

- [13] Seth Abraham. Time to revisit rep;movs - comment. <https://software.intel.com/en-us/forums/topic/275765>, Aug 2006. [Online; accessed 23-January-2015].
- [14] Tiago Alves and Don Felton. Trustzone: Integrated hardware and software security. *Information Quarterly*, 3(4):18–24, 2004.
- [15] Ittai Anati, Shay Gueron, Simon P Johnson, and Vincent R Scarlata. Innovative technology for cpu based attestation and sealing. In *Proceedings of the 2nd International Workshop on Hardware and Architectural Support for Security and Privacy, HASP*, volume 13, 2013.
- [16] Ross Anderson. *Security engineering: A guide to building dependable distributed systems*. Wiley, 2001.
- [17] Sebastian Anthony. Who actually develops linux? the answer might surprise you. <http://www.extremetech.com/computing/175919-who-actually-develops-linux>, 2014. [Online; accessed 27-April-2015].
- [18] Gorka Irazoqui Apecechea, Mehmet Sinan Inci, Thomas Eisenbarth, and Berk Sunar. Fine grain cross-vm attacks on xen and vmware are possible! *Cryptology ePrint Archive*, Report 2014/248, 2014. <http://eprint.iacr.org/>.
- [19] ARM Limited. *AMBA® AXI Protocol*, Mar 2004. Reference no. IHI 0022B, IHI 0024B, AR500-DA-10004.
- [20] ARM Limited. *ARM Security Technology Building a Secure System using TrustZone® Technology*, Apr 2009. Reference no. PRD29-GENC-009492C.
- [21] Sebastian Banescu. Cache timing attacks. 2011. [Online; accessed 26-January-2014].
- [22] Elaine Barker, William Barker, William Burr, William Polk, and Miles Smid. Recommendation for key management part 1: General (revision 3). *Federal Information Processing Standards (FIPS) Special Publications (SP)*, 800-57, Jul 2012.
- [23] Elaine Barker, William Barker, William Burr, William Polk, and Miles Smid. Secure hash standard (shs). *Federal Information Processing Standards (FIPS) Publications (PUBS)*, 180-4, Aug 2015.
- [24] Friedrich Beck. *Integrated Circuit Failure Analysis: a Guide to Preparation Techniques*. John Wiley & Sons, 1998.
- [25] Daniel Bleichenbacher. Chosen ciphertext attacks against protocols based on the rsa encryption standard pkcs# 1. In *Advances in Cryptology CRYPTO'98*, pages 1–12. Springer, 1998.
- [26] D.D. Boggs and S.D. Rodgers. Microprocessor with novel instruction for signaling event occurrence and for providing event handling information in response thereto, 1997. US Patent 5,625,788.

- [27] Joseph Bonneau and Ilya Mironov. Cache-collision timing attacks against aes. In *Cryptographic Hardware and Embedded Systems-CHES 2006*, pages 201–215. Springer, 2006.
- [28] Ernie Brickell and Jiangtao Li. Enhanced privacy id from bilinear pairing. *IACR Cryptology ePrint Archive*, 2009.
- [29] Billy Bob Brumley and Nicola Tuveri. Remote timing attacks are still practical. In *Computer Security-ESORICS 2011*, pages 355–371. Springer, 2011.
- [30] David Brumley and Dan Boneh. Remote timing attacks are practical. *Computer Networks*, 48(5):701–716, 2005.
- [31] John Butterworth, Corey Kallenberg, Xenon Kovah, and Amy Herzog. Bios chronomancy: Fixing the core root of trust for measurement. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & Communications Security*, pages 25–36. ACM, 2013.
- [32] J Lawrence Carter and Mark N Wegman. Universal classes of hash functions. In *Proceedings of the 9th annual ACM Symposium on Theory of Computing*, pages 106–112. ACM, 1977.
- [33] David Champagne and Ruby B Lee. Scalable architectural support for trusted software. In *High Performance Computer Architecture (HPCA), 2010 IEEE 16th International Symposium on*, pages 1–12. IEEE, 2010.
- [34] Daming D Chen and Gail-Joon Ahn. Security analysis of x86 processor microcode. 2014. [Online; accessed 7-January-2015].
- [35] Haogang Chen, Yandong Mao, Xi Wang, Dong Zhou, Nikolai Zeldovich, and M Frans Kaashoek. Linux kernel vulnerabilities: State-of-the-art defenses and open problems. In *Proceedings of the Second Asia-Pacific Workshop on Systems*, page 5. ACM, 2011.
- [36] Lily Chen. Recommendation for key derivation using pseudorandom functions. *Federal Information Processing Standards (FIPS) Special Publications (SP)*, 800-108, Oct 2009.
- [37] Coreboot. Developer manual, Sep 2014. [Online; accessed 4-March-2015].
- [38] M.P. Cornaby and B. Chaffin. Microinstruction pointer stack including speculative pointers for out-of-order execution, 2007. US Patent 7,231,511.
- [39] Intel Corporation. *Intel® Xeon® Processor E5 v3 Family Uncore Performance Monitoring Reference Manual*, Sep 2014. Reference no. 331051-001.
- [40] Victor Costan, Ilia Lebedev, and Srinivas Devadas. Sanctum: Minimal hardware extensions for strong software isolation. *Cryptology ePrint Archive*, Report 2015/564, 2015. <http://eprint.iacr.org/>.

- [41] J. Daemen and V. Rijmen. Aes proposal: Rijndael, aes algorithm submission, Sep 1999.
- [42] S.M. Datta and M.J. Kumar. Technique for providing secure firmware, 2013. US Patent 8,429,418.
- [43] S.M. Datta, V.J. Zimmer, and M.A. Rothman. System and method for trusted early boot flow, 2010. US Patent 7,752,428.
- [44] Shaun Davenport. Sgx: the good, the bad and the downright ugly. *Virus Bulletin*, 2014.
- [45] Pete Dice. Booting an intel architecture system, part i: Early initialization. *Dr. Dobb's*, Dec 2011. [Online; accessed 2-Dec-2015].
- [46] Whitfield Diffie and Martin E Hellman. New directions in cryptography. *Information Theory, IEEE Transactions on*, 22(6):644–654, 1976.
- [47] Loïc Duflot, Daniel Etiemble, and Olivier Grumelard. Using cpu system management mode to circumvent operating system security functions. *CanSecWest/core06*, 2006.
- [48] Alan Dunn, Owen Hofmann, Brent Waters, and Emmett Witchel. Cloaking malware with the trusted platform module. In *USENIX Security Symposium*, 2011.
- [49] Morris Dworkin. Recommendation for block cipher modes of operation: Methods and techniques. *Federal Information Processing Standards (FIPS) Special Publications (SP)*, 800-38A, Dec 2001.
- [50] Morris Dworkin. Recommendation for block cipher modes of operation: The cmac mode for authentication. *Federal Information Processing Standards (FIPS) Special Publications (SP)*, 800-38B, May 2005.
- [51] Morris Dworkin. Recommendation for block cipher modes of operation: Galois/-counter mode (gcm) and gmac. *Federal Information Processing Standards (FIPS) Special Publications (SP)*, 800-38D, Nov 2007.
- [52] D. Eastlake and P. Jones. RFC 3174: US Secure Hash Algorithm 1 (SHA1). *Internet RFCs*, 2001.
- [53] Shawn Embleton, Sherri Sparks, and Cliff C Zou. Smm rootkit: a new breed of os independent malware. *Security and Communication Networks*, 2010.
- [54] Dmitry Evtvushkin, Jesse Elwell, Meltem Ozsoy, Dmitry Ponomarev, Nael Abu Ghazaleh, and Ryan Riley. Iso-x: A flexible architecture for hardware-managed isolated execution. In *Microarchitecture (MICRO), 2014 47th annual IEEE/ACM International Symposium on*, pages 190–202. IEEE, 2014.
- [55] Niels Ferguson, Bruce Schneier, and Tadayoshi Kohno. *Cryptography Engineering: Design Principles and Practical Applications*. John Wiley & Sons, 2011.

- [56] Christopher W Fletcher, Marten van Dijk, and Srinivas Devadas. A secure processor architecture for encrypted computation on untrusted programs. In *Proceedings of the Seventh ACM Workshop on Scalable Trusted Computing*, pages 3–8. ACM, 2012.
- [57] Agner Fog. Instruction tables - lists of instruction latencies, throughputs and micro-operation breakdowns for intel, amd and via cpus. Dec 2014. [Online; accessed 23-January-2015].
- [58] Andrew Furtak, Yuriy Bulygin, Oleksandr Bazhaniuk, John Loucaides, Alexander Matrosov, and Mikhail Gorobets. Bios and secure boot attacks uncovered. *The 10th ekoparty Security Conference*, 2014. [Online; accessed 22-October-2015].
- [59] William Futral and James Greene. *Intel® Trusted Execution Technology for Server Platforms*. Apress Open, 2013.
- [60] Blaise Gassend, Dwaine Clarke, Marten Van Dijk, and Srinivas Devadas. Silicon physical random functions. In *Proceedings of the 9th ACM Conference on Computer and Communications Security*, pages 148–160. ACM, 2002.
- [61] Blaise Gassend, G Edward Suh, Dwaine Clarke, Marten Van Dijk, and Srinivas Devadas. Caches and hash trees for efficient memory integrity verification. In *Proceedings of the 9th International Symposium on High-Performance Computer Architecture*, pages 295–306. IEEE, 2003.
- [62] Daniel Genkin, Lev Pachmanov, Itamar Pipman, and Eran Tromer. Stealing keys from pcs using a radio: Cheap electromagnetic attacks on windowed exponentiation. Cryptology ePrint Archive, Report 2015/170, 2015.
- [63] Daniel Genkin, Itamar Pipman, and Eran Tromer. Get your hands off my laptop: Physical side-channel key-extraction attacks on pcs. Cryptology ePrint Archive, Report 2014/626, 2014.
- [64] Daniel Genkin, Adi Shamir, and Eran Tromer. Rsa key extraction via low-bandwidth acoustic cryptanalysis. Cryptology ePrint Archive, Report 2013/857, 2013.
- [65] Craig Gentry. *A fully homomorphic encryption scheme*. PhD thesis, Stanford University, 2009.
- [66] R.T. George, J.W. Brandt, K.S. Venkatraman, and S.P. Kim. Dynamically partitioning pipeline resources, 2009. US Patent 7,552,255.
- [67] A. Glew, G. Hinton, and H. Akkary. Method and apparatus for performing page table walks in a microprocessor capable of processing speculative instructions, 1997. US Patent 5,680,565.
- [68] A.F. Glew, H. Akkary, R.P. Colwell, G.J. Hinton, D.B. Papworth, and M.A. Fetterman. Method and apparatus for implementing a non-blocking translation lookaside buffer, 1996. US Patent 5,564,111.

- [69] Oded Goldreich. Towards a theory of software protection and simulation by oblivious rams. In *Proceedings of the 19th annual ACM symposium on Theory of Computing*, pages 182–194. ACM, 1987.
- [70] J.R. Goodman and H.H.J. Hum. Mesif: A two-hop cache coherency protocol for point-to-point interconnects. 2009.
- [71] K.C. Gotze, G.M. Iovino, and J. Li. Secure provisioning of secret keys during integrated circuit manufacturing, 2014. US Patent App. 13/631,512.
- [72] K.C. Gotze, J. Li, and G.M. Iovino. Fuse attestation to secure the provisioning of secret keys during integrated circuit manufacturing, 2014. US Patent 8,885,819.
- [73] Joe Grand. Advanced hardware hacking techniques, Jul 2004.
- [74] David Grawrock. *Dynamics of a Trusted Platform: A building block approach*. Intel Press, 2009.
- [75] Trusted Computing Group. Tpm main specification. http://www.trustedcomputinggroup.org/resources/tpm_main_specification, 2003.
- [76] Daniel Gruss, Clémentine Maurice, and Stefan Mangard. Rowhammer.js: A remote software-induced fault attack in javascript. *CoRR*, abs/1507.06955, 2015.
- [77] Shay Gueron. Quick verification of rsa signatures. In *8th International Conference on Information Technology: New Generations (ITNG)*, pages 382–386. IEEE, 2011.
- [78] Ben Hawkes. Security analysis of x86 processor microcode. 2012. [Online; accessed 7-January-2015].
- [79] John L Hennessy and David A Patterson. *Computer Architecture - a Quantitative Approach (5 ed.)*. Morgan Kaufmann, 2012.
- [80] Christoph Herbst, Elisabeth Oswald, and Stefan Mangard. An aes smart card implementation resistant to power analysis attacks. In *Applied cryptography and Network security*, pages 239–252. Springer, 2006.
- [81] G. Hildesheim, I. Anati, H. Shafi, S. Raikin, G. Gerzon, U.R. Savagaonkar, C.V. Rozas, F.X. McKeen, M.A. Goldsmith, and D. Prashant. Apparatus and method for page walk extension for enhanced security checks, 2014. US Patent App. 13/730,563.
- [82] Matthew Hoekstra, Reshma Lal, Pradeep Pappachan, Vinay Phegade, and Juan Del Cuvillo. Using innovative instructions to create trustworthy software solutions. In *Proceedings of the 2nd International Workshop on Hardware and Architectural Support for Security and Privacy, HASP*, volume 13, 2013.
- [83] Gael Hofemeier. Intel manageability firmware recovery agent. Mar 2013. [Online; accessed 2-Dec-2015].

- [84] George Hotz. Ps3 glitch hack. 2010. [Online; accessed 7-January-2015].
- [85] Andrew Huang. *Hacking the Xbox: an Introduction to Reverse Engineering*. No Starch Press, 2003.
- [86] C.J. Hughes, Y.K. Chen, M. Bomb, J.W. Brandt, M.J. Buxton, M.J. Charney, S. Chen-nupaty, J. Corbal, M.G. Dixon, M.B. Girkar, et al. Gathering and scattering multiple data elements, 2013. US Patent 8,447,962.
- [87] IEEE Computer Society. *IEEE Standard for Ethernet*, Dec 2012. IEEE Std. 802.3-2012.
- [88] Mehmet Sinan Inci, Berk Gulmezoglu, Gorka Irazoqui, Thomas Eisenbarth, and Berk Sunar. Seriously, get off my cloud! cross-vm rsa key recovery in a public cloud. Cryptology ePrint Archive, Report 2015/898, 2015.
- [89] Intel Corporation. *Intel® Processor Serial Number*, Mar 1999. Order no. 245125-001.
- [90] Intel Corporation. *Intel® architecture Platform Basics*, Sep 2010. Reference no. 324377.
- [91] Intel Corporation. *Intel® Core 2 Duo and Intel® Core 2 Solo Processor for Intel® Centrino® Duo Processor Technology Intel® Celeron® Processor 500 Series - Specification Update*, Dec 2010. Reference no. 314079-026.
- [92] Intel Corporation. *Intel® Trusted Execution Technology (Intel® TXT) LAB Handout*, 2010. [Online; accessed 2-July-2015].
- [93] Intel Corporation. *Intel® Xeon® Processor 7500 Series Uncore Programming Guide*, Mar 2010. Reference no. 323535-001.
- [94] Intel Corporation. *An Introduction to the Intel® QuickPath Interconnect*, Mar 2010. Reference no. 323535-001.
- [95] Intel Corporation. *Minimal Intel® Architecture Boot LoaderBare Bones Functionality Required for Booting an Intel® Architecture Platform*, Jan 2010. Reference no. 323246.
- [96] Intel Corporation. *Intel® 7 Series Family - Intel® Management Engine Firmware 8.1 - 1.5MB Firmware Bring Up Guide*, Jul 2012. Revision 8.1.0.1248 - PV Release.
- [97] Intel Corporation. *Intel® Xeon® Processor E5-2600 Product Family Uncore Performance Monitoring Guide*, Mar 2012. Reference no. 327043-001.
- [98] Intel Corporation. *Software Guard Extensions Programming Reference*, 2013. Reference no. 329298-001US.
- [99] Intel Corporation. *Intel® 64 and IA-32 Architectures Optimization Reference Manual*, Sep 2014. Reference no. 248966-030.

- [100] Intel Corporation. *Intel® Xeon® Processor 7500 Series Datasheet - Volume Two*, Mar 2014. Reference no. 329595-002.
- [101] Intel Corporation. *Intel® Xeon® Processor E7 v2 2800/4800/8800 Product Family Datasheet - Volume Two*, Mar 2014. Reference no. 329595-002.
- [102] Intel Corporation. *Software Guard Extensions Programming Reference*, 2014. Reference no. 329298-002US.
- [103] Intel Corporation. *Intel® 100 Series Chipset Family Platform Controller Hub (PCH) Datasheet - Volume One*, Aug 2015. Reference no. 332690-001EN.
- [104] Intel Corporation. *Intel® 64 and IA-32 Architectures Software Developer's Manual*, Sep 2015. Reference no. 325462-056US.
- [105] Intel Corporation. *Intel® C610 Series Chipset and Intel® X99 Chipset Platform Controller Hub (PCH) Datasheet*, Oct 2015. Reference no. 330788-003.
- [106] Intel Corporation. *Intel® Software Guard Extensions (Intel® SGX)*, Jun 2015. Reference no. 332680-002.
- [107] Intel Corporation. *Intel® Xeon® Processor 5500 Series - Specification Update*, 2 2015. Reference no. 321324-018US.
- [108] Intel Corporation. *Intel® Xeon® Processor E5-1600, E5-2400, and E5-2600 v3 Product Family Datasheet - Volume Two*, Jan 2015. Reference no. 330784-002.
- [109] Intel Corporation. *Intel® Xeon® Processor E5 Product Family - Specification Update*, Jan 2015. Reference no. 326150-018.
- [110] Intel Corporation. *Mobile 4th Generation Intel® Core® Processor Family I/O Datasheet*, Feb 2015. Reference no. 329003-003.
- [111] Bruce Jacob and Trevor Mudge. Virtual memory: Issues of implementation. *Computer*, 31(6):33–43, 1998.
- [112] Simon P Johnson, Uday R Savagaonkar, Vincent R Scarlata, Francis X McKeen, and Carlos V Rozas. Technique for supporting multiple secure enclaves, Dec 2010. US Patent 8,972,746.
- [113] Jakob Jonsson and Burt Kaliski. RFC 3447: Public-Key Cryptography Standards (PKCS) #1: RSA Cryptography Specifications Version 2.1. *Internet RFCs*, Feb 2003.
- [114] Burt Kaliski. RFC 2313: PKCS #1: RSA Encryption Version 1.5. *Internet RFCs*, Mar 1998.
- [115] Burt Kaliski and Jessica Staddon. RFC 2437: PKCS #1: RSA Encryption Version 2.0. *Internet RFCs*, Oct 1998.

- [116] Corey Kallenberg, Xeno Kovah, John Butterworth, and Sam Cornwell. Extreme privilege escalation on windows 8/uefi systems, 2014.
- [117] Emilia Käsper and Peter Schwabe. Faster and timing-attack resistant aes-gcm. In *Cryptographic Hardware and Embedded Systems-CHES 2009*, pages 1–17. Springer, 2009.
- [118] Jonathan Katz and Yehuda Lindell. *Introduction to modern cryptography*. CRC Press, 2014.
- [119] Richard E Kessler and Mark D Hill. Page placement algorithms for large real-indexed caches. *ACM Transactions on Computer Systems (TOCS)*, 10(4):338–359, 1992.
- [120] Taesoo Kim and Nickolai Zeldovich. Practical and effective sandboxing for non-root users. In *USENIX Annual Technical Conference*, pages 139–144, 2013.
- [121] Yoongu Kim, Ross Daly, Jeremie Kim, Chris Fallin, Ji Hye Lee, Donghyuk Lee, Chris Wilkerson, Konrad Lai, and Onur Mutlu. Flipping bits in memory without accessing them: An experimental study of dram disturbance errors. In *Proceeding of the 41st annual International Symposium on Computer Architecture*, pages 361–372. IEEE Press, 2014.
- [122] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduwe, Kai Engelhardt, Rafal Kolanski, Michael Norrish, et al. sel4: Formal verification of an os kernel. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*, pages 207–220. ACM, 2009.
- [123] L.A. Knauth and P.J. Irelan. Apparatus and method for providing eventing ip and source data address in a statistical sampling infrastructure, 2014. US Patent App. 13/976,613.
- [124] N. Koblitz. Elliptic curve cryptosystems. *Mathematics of Computation*, 48(177):203–209, 1987.
- [125] Paul Kocher, Joshua Jaffe, and Benjamin Jun. Differential power analysis. In *Advances in Cryptology (CRYPTO)*, pages 388–397. Springer, 1999.
- [126] Paul C Kocher. Timing attacks on implementations of diffie-hellman, rsa, dss, and other systems. In *Advances in Cryptology CRYPTO96*, pages 104–113. Springer, 1996.
- [127] Hugo Krawczyk, Ran Canetti, and Mihir Bellare. Hmac: Keyed-hashing for message authentication. 1997.
- [128] Markus G Kuhn. Electromagnetic eavesdropping risks of flat-panel displays. In *Privacy Enhancing Technologies*, pages 88–107. Springer, 2005.
- [129] Tsvika Kurts, Guillermo Savransky, Jason Ratner, Eilon Hazan, Daniel Skaba, Sharon Elmosnino, and Geeyarpuram N Santhanakrishnan. Generic debug external connection (gdx) for high integration integrated circuits, 2011. US Patent 8,074,131.

- [130] Yunsup Lee, Andrew Waterman, Rimas Avizienis, Henry Cook, Chen Sun, Vladimir Stojanovic, and Krste Asanovic. A 45nm 1.3 ghz 16.7 double-precision gflops/w risc-v processor with vector accelerators. In *European Solid State Circuits Conference (ESSCIRC), ESSCIRC 2014-40th*, pages 199–202. IEEE, 2014.
- [131] David Levinthal. Performance analysis guide for intel® core i7 processor and intel® xeon 5500 processors. https://software.intel.com/sites/products/collateral/hpc/vtune/performance_analysis_guide.pdf, 2010. [Online; accessed 26-January-2015].
- [132] David Lie, Chandramohan Thekkath, Mark Mitchell, Patrick Lincoln, Dan Boneh, John Mitchell, and Mark Horowitz. Architectural support for copy and tamper resistant software. *ACM SIGPLAN Notices*, 35(11):168–177, 2000.
- [133] Jiang Lin, Qingda Lu, Xiaoning Ding, Zhao Zhang, Xiaodong Zhang, and P Sadayappan. Gaining insights into multicore cache partitioning: Bridging the gap between simulation and real systems. In *14th International IEEE Symposium on High Performance Computer Architecture (HPCA)*, pages 367–378. IEEE, 2008.
- [134] Barbara Liskov and Stephen Zilles. Programming with abstract data types. In *ACM Sigplan Notices*, volume 9, pages 50–59. ACM, 1974.
- [135] Fangfei Liu, Yuval Yarom, Qian Ge, Gernot Heiser, and Ruby B Lee. Last-level cache side-channel attacks are practical. In *Security and Privacy (SP), 2015 IEEE Symposium on*, pages 143–158. IEEE, 2015.
- [136] Martin Maas, Eric Love, Emil Stefanov, Mohit Tiwari, Elaine Shi, Krste Asanovic, John Kubiatowicz, and Dawn Song. Phantom: Practical oblivious computation in a secure processor. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 311–324. ACM, 2013.
- [137] R. Maes, P. Tuyls, and I. Verbauwhede. Low-Overhead Implementation of a Soft Decision Helper Data Algorithm for SRAM PUFs. In *Cryptographic Hardware and Embedded Systems (CHES)*, pages 332–347, 2009.
- [138] James Manger. A chosen ciphertext attack on rsa optimal asymmetric encryption padding (oaep) as standardized in pkcs# 1 v2.0. In *Advances in Cryptology CRYPTO 2001*, pages 230–238. Springer, 2001.
- [139] Clmentine Maurice, Nicolas Le Scouarnec, Christoph Neumann, Olivier Heen, and Aurlien Francillon. Reverse engineering intel last-level cache complex addressing using performance counters. In *Proceedings of the 18th International Symposium on Research in Attacks, Intrusions and Defenses (RAID)*, 2015.
- [140] Jonathan M McCune, Yanlin Li, Ning Qu, Zongwei Zhou, Anupam Datta, Virgil Gligor, and Adrian Perrig. Trustvisor: Efficient tcb reduction and attestation. In *Security and Privacy (SP), 2010 IEEE Symposium on*, pages 143–158. IEEE, 2010.

- [141] David McGrew and John Viega. The galois/counter mode of operation (gcm). 2004. [Online; accessed 28-December-2015].
- [142] Francis X McKeen, Carlos V Rozas, Uday R Savagaonkar, Simon P Johnson, Vincent Scarlata, Michael A Goldsmith, Ernie Brickell, Jiang Tao Li, Howard C Herbert, Prashant Dewan, et al. Method and apparatus to provide secure application execution, Dec 2009. US Patent 9,087,200.
- [143] Frank McKeen, Ilya Alexandrovich, Alex Berenzon, Carlos V Rozas, Hisham Shafi, Vedvyas Shanbhogue, and Uday R Savagaonkar. Innovative instructions and software model for isolated execution. *HASP*, 13:10, 2013.
- [144] Michael Naehrig, Kristin Lauter, and Vinod Vaikuntanathan. Can homomorphic encryption be practical? In *Proceedings of the 3rd ACM workshop on Cloud computing security workshop*, pages 113–124. ACM, 2011.
- [145] National Institute of Standards and Technology (NIST). The advanced encryption standard (aes). *Federal Information Processing Standards (FIPS) Publications (PUBS)*, 197, Nov 2001.
- [146] National Institute of Standards and Technology (NIST). The digital signature standard (dss). *Federal Information Processing Standards (FIPS) Processing Standards Publications (PUBS)*, 186-4, Jul 2013.
- [147] National Security Agency (NSA) Central Security Service (CSS). Cryptography today on suite b phase-out. https://www.nsa.gov/ia/programs/suiteb_cryptography/, Aug 2015. [Online; accessed 28-December-2015].
- [148] M.S. Natu, S. Datta, J. Wiedemeier, J.R. Vash, S. Kottapalli, S.P. Bobholz, and A. Baum. Supporting advanced ras features in a secured computing system, 2012. US Patent 8,301,907.
- [149] Yossef Oren, Vasileios P Kemerlis, Simha Sethumadhavan, and Angelos D Keromytis. The spy in the sandbox – practical cache attacks in javascript. *arXiv preprint arXiv:1502.07373*, 2015.
- [150] Dag Arne Osvik, Adi Shamir, and Eran Tromer. Cache attacks and countermeasures: the case of aes. In *Topics in Cryptology–CT-RSA 2006*, pages 1–20. Springer, 2006.
- [151] Scott Owens, Susmit Sarkar, and Peter Sewell. A better x86 memory model: x86-tso (extended version). *University of Cambridge, Computer Laboratory, Technical Report*, (UCAM-CL-TR-745), 2009.
- [152] Emmanuel Owusu, Jun Han, Sauvik Das, Adrian Perrig, and Joy Zhang. Accessory: password inference using accelerometers on smartphones. In *Proceedings of the Twelfth Workshop on Mobile Computing Systems & Applications*, page 9. ACM, 2012.

- [153] D.B. Papworth, G.J. Hinton, M.A. Fetterman, R.P. Colwell, and A.F. Glew. Exception handling in a processor that performs speculative out-of-order instruction execution, 1999. US Patent 5,987,600.
- [154] David A Patterson and John L Hennessy. *Computer Organization and Design: the hardware/software interface*. Morgan Kaufmann, 2013.
- [155] P. Pessl, D. Gruss, C. Maurice, M. Schwarz, and S. Mangard. Reverse engineering intel dram addressing and exploitation. *ArXiv e-prints*, Nov 2015.
- [156] Stefan M Petters and Georg Farber. Making worst case execution time analysis for hard real-time tasks on state of the art processors feasible. In *Sixth International Conference on Real-Time Computing Systems and Applications*, pages 442–449. IEEE, 1999.
- [157] S.A. Qureshi and M.O. Nicholes. System and method for using a firmware interface table to dynamically load an acpi ssdt, 2006. US Patent 6,990,576.
- [158] S. Raikin, O. Hamama, R.S. Chappell, C.B. Rust, H.S. Luu, L.A. Ong, and G. Hildesheim. Apparatus and method for a multiple page size translation lookaside buffer (tlb), 2014. US Patent App. 13/730,411.
- [159] S. Raikin and R. Valentine. Gather cache architecture, 2014. US Patent 8,688,962.
- [160] Stefan Reinauer. x86 intel: Add firmware interface table support. <http://review.coreboot.org/#/c/2642/>, 2013. [Online; accessed 2-July-2015].
- [161] Thomas Ristenpart, Eran Tromer, Hovav Shacham, and Stefan Savage. Hey, you, get off of my cloud: Exploring information leakage in third-party compute clouds. In *Proceedings of the 16th ACM Conference on Computer and Communications Security*, pages 199–212. ACM, 2009.
- [162] RL Rivest, A. Shamir, and L. Adleman. A method for obtaining digital signatures and public-key cryptosystems. *Communications of the ACM*, 21(2):120–126, 1978.
- [163] S.D. Rodgers, K.K. Tiruvallur, M.W. Rhodehamel, K.G. Konigsfeld, A.F. Glew, H. Akkary, M.A. Karnik, and J.A. Brayton. Method and apparatus for performing operations based upon the addresses of microinstructions, 1997. US Patent 5,636,374.
- [164] S.D. Rodgers, R. Vidwans, J. Huang, M.A. Fetterman, and K. Huck. Method and apparatus for generating event handler vectors based on both operating mode and event type, 1999. US Patent 5,889,982.
- [165] M. Rosenblum and T. Garfinkel. Virtual machine monitors: current technology and future trends. *Computer*, 38(5):39–47, May 2005.
- [166] Xiaoyu Ruan. *Platform Embedded Security Technology Revealed*. Apress, 2014.
- [167] Joanna Rutkowska. Thoughts on intel’s upcoming software guard extensions (part 2). *Invisible Things Lab*, 2013.

- [168] Joanna Rutkowska. Intel x86 considered harmful. Oct 2015. [Online; accessed 2-Nov-2015].
- [169] Joanna Rutkowska and Rafał Wojtczuk. Preventing and detecting xen hypervisor subversions. *Blackhat Briefings USA*, 2008.
- [170] Jerome H Saltzer and M Frans Kaashoek. *Principles of Computer System Design: An Introduction*. Morgan Kaufmann, 2009.
- [171] Mark Seaborn and Thomas Dullien. Exploiting the dram rowhammer bug to gain kernel privileges. <http://googleprojectzero.blogspot.com/2015/03/exploiting-dram-rowhammer-bug-to-gain.html>, Mar 2015. [Online; accessed 9-March-2015].
- [172] V. Shanbhogue, J.W. Brandt, and J. Wiedemeier. Protecting information processing system secrets from debug attacks, 2015. US Patent 8,955,144.
- [173] V. Shanbhogue and S.J. Robinson. Enabling virtualization of a processor resource, 2014. US Patent 8,806,104.
- [174] Stephen Shankland. Itanium: A cautionary tale. Dec 2005. [Online; accessed 11-February-2015].
- [175] Alan Jay Smith. Cache memories. *ACM Computing Surveys (CSUR)*, 14(3):473–530, 1982.
- [176] Sean W Smith, Ron Perez, Steve Weingart, and Vernon Austel. Validating a high-performance, programmable secure coprocessor. In *22nd National Information Systems Security Conference*. IBM Thomas J. Watson Research Division, 1999.
- [177] Sean W Smith and Steve Weingart. Building a high-performance, programmable secure coprocessor. *Computer Networks*, 31(8):831–860, 1999.
- [178] Emil Stefanov, Marten Van Dijk, Elaine Shi, Christopher Fletcher, Ling Ren, Xiangyao Yu, and Srinivas Devadas. Path oram: An extremely simple oblivious ram protocol. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 299–310. ACM, 2013.
- [179] Marc Stevens, Pierre Karpman, and Thomas Peyrin. Free-start collision on full sha-1. Cryptology ePrint Archive, Report 2015/967, 2015.
- [180] G Edward Suh, Dwaine Clarke, Blaise Gassend, Marten Van Dijk, and Srinivas Devadas. Aegis: architecture for tamper-evident and tamper-resistant processing. In *Proceedings of the 17th annual international conference on Supercomputing*, pages 160–171. ACM, 2003.
- [181] G Edward Suh and Srinivas Devadas. Physical unclonable functions for device authentication and secret key generation. In *Proceedings of the 44th annual Design Automation Conference*, pages 9–14. ACM, 2007.

- [182] G. Edward Suh, Charles W. O'Donnell, Ishan Sachdev, and Srinivas Devadas. Design and Implementation of the AEGIS Single-Chip Secure Processor Using Physical Random Functions. In *Proceedings of the 32nd ISCA'05*. ACM, June 2005.
- [183] George Taylor, Peter Davies, and Michael Farmwald. The tlb slice - a low-cost high-speed address translation mechanism. *SIGARCH Computer Architecture News*, 18(2SI):355–363, 1990.
- [184] Alexander Tereshkin and Rafal Wojtczuk. Introducing ring-3 rootkits. Master's thesis, 2009.
- [185] Kris Tiri, Moonmoon Akmal, and Ingrid Verbauwhede. A dynamic and differential cmos logic with signal independent power consumption to withstand differential power analysis on smart cards. In *Proceedings of the 28th European Solid-State Circuits Conference (ESSCIRC)*, pages 403–406. IEEE, 2002.
- [186] UEFI Forum. *Unified Extensible Firmware Interface Specification, Version 2.5*, 2015. [Online; accessed 1-Jul-2015].
- [187] Rich Uhlig, Gil Neiger, Dion Rodgers, Amy L Santoni, Fernando CM Martins, Andrew V Anderson, Steven M Bennett, Alain Kagi, Felix H Leung, and Larry Smith. Intel virtualization technology. *Computer*, 38(5):48–56, 2005.
- [188] Wim Van Eck. Electromagnetic radiation from video display units: an eavesdropping risk? *Computers & Security*, 4(4):269–286, 1985.
- [189] Amit Vasudevan, Jonathan M McCune, Ning Qu, Leendert Van Doorn, and Adrian Perrig. Requirements for an integrity-protected hypervisor on the x86 hardware virtualized architecture. In *Trust and Trustworthy Computing*, pages 141–165. Springer, 2010.
- [190] Sathish Venkataramani. *Advanced Board Bring Up - Power Sequencing Guide for Embedded Intel Architecture*. Intel Corporation, Apr 2011. Reference no. 325268.
- [191] Vassilios Ververis. Security evaluation of intel's active management technology. 2010.
- [192] Andrew Waterman, Yunsup Lee, Rimas Avizienis, David A. Patterson, and Krste Asanovic. The risc-v instruction set manual volume ii: Privileged architecture version 1.7. Technical Report UCB/EECS-2015-49, EECS Department, University of California, Berkeley, May 2015.
- [193] Andrew Waterman, Yunsup Lee, and et al. Celio, Christopher. Risc-v proxy kernel and boot loader. Technical report, EECS Department, University of California, Berkeley, May 2015.
- [194] Andrew Waterman, Yunsup Lee, David A. Patterson, and Krste Asanovic. The risc-v instruction set manual, volume i: User-level isa, version 2.0. Technical Report UCB/EECS-2014-54, EECS Department, University of California, Berkeley, May 2014.

- [195] Filip Wecherowski. A real smm rootkit: Reversing and hooking bios smi handlers. *Phrack Magazine*, 13(66), 2009.
- [196] Mark N Wegman and J Lawrence Carter. New hash functions and their use in authentication and set equality. *Journal of Computer and System Sciences*, 22(3):265–279, 1981.
- [197] Rafal Wojtczuk and Joanna Rutkowska. Attacking intel trusted execution technology. *Black Hat DC*, 2009.
- [198] Rafal Wojtczuk and Joanna Rutkowska. Attacking smm memory via intel cpu cache poisoning. *Invisible Things Lab*, 2009.
- [199] Rafal Wojtczuk and Joanna Rutkowska. Attacking intel txt via sinit code execution hijacking, 2011.
- [200] Rafal Wojtczuk, Joanna Rutkowska, and Alexander Tereshkin. Another way to circumvent intel® trusted execution technology. *Invisible Things Lab*, 2009.
- [201] Rafal Wojtczuk and Alexander Tereshkin. Attacking intel® bios. *Invisible Things Lab*, 2010.
- [202] Y. Wu and M. Breternitz. Genetic algorithm for microcode compression, 2008. US Patent 7,451,121.
- [203] Y. Wu, S. Kim, M. Breternitz, and H. Hum. Compressing and accessing a microcode rom, 2012. US Patent 8,099,587.
- [204] Yuanzhong Xu, Weidong Cui, and Marcus Peinado. Controlled-channel attacks: Deterministic side channels for untrusted operating systems. In *Proceedings of the 36th IEEE Symposium on Security and Privacy (Oakland)*. IEEE Institute of Electrical and Electronics Engineers, May 2015.
- [205] Yuval Yarom and Katrina E Falkner. Flush+ reload: a high resolution, low noise, l3 cache side-channel attack. *IACR Cryptology ePrint Archive*, 2013:448, 2013.
- [206] Yuval Yarom, Qian Ge, Fangfei Liu, Ruby B. Lee, and Gernot Heiser. Mapping the intel last-level cache. *Cryptology ePrint Archive*, Report 2015/905, 2015.
- [207] Bennet Yee. *Using secure coprocessors*. PhD thesis, Carnegie Mellon University, 1994.
- [208] Bennet Yee, David Sehr, Gregory Dardyk, J Bradley Chen, Robert Muth, Tavis Ormandy, Shiki Okasaka, Neha Narula, and Nicholas Fullagar. Native client: A sandbox for portable, untrusted x86 native code. In *Security and Privacy, 2009 30th IEEE Symposium on*, pages 79–93. IEEE, 2009.

- [209] Marcelo Yuffe, Ernest Knoll, Moty Mehalel, Joseph Shor, and Tsvika Kurts. A fully integrated multi-cpu, gpu and memory controller 32nm processor. In *Solid-State Circuits Conference Digest of Technical Papers (ISSCC), 2011 IEEE International*, pages 264–266. IEEE, 2011.
- [210] Xiantao Zhang and Yaozu Dong. Optimizing xen vmm based on intel® virtualization technology. In *Internet Computing in Science and Engineering, 2008. ICICSE'08. International Conference on*, pages 367–374. IEEE, 2008.
- [211] Li Zhuang, Feng Zhou, and J Doug Tygar. Keyboard acoustic emanations revisited. *ACM Transactions on Information and System Security (TISSEC)*, 13(1):3, 2009.
- [212] V.J. Zimmer and S.H. Robinson. Methods and systems for microcode patching, 2012. US Patent 8,296,528.
- [213] V.J. Zimmer and J. Yao. Method and apparatus for sequential hypervisor invocation, 2012. US Patent 8,321,931.