

# Teaching an Old Dog New Tricks: Verifiable FHE Using Commodity Hardware

Jules Drean  
MIT

Srinivas Devadas  
MIT

Fisher Jepsen  
MIT

Aamer Jaleel  
NVIDIA

Edward Suh  
NVIDIA

Gururaj Saileshwar  
NVIDIA

## ABSTRACT

We present Argos, a simple approach for adding verifiability to fully homomorphic encryption (FHE) schemes using trusted hardware. Traditional approaches to verifiable FHE require expensive cryptographic proofs, which incur an overhead of up to seven orders of magnitude *on top of FHE*, making them impractical.

With Argos, we show that trusted hardware can be securely used to provide verifiability for FHE computations, with minimal overhead relative to the baseline FHE computation. An important contribution of Argos is showing that the major security pitfall associated with trusted hardware, *microarchitectural* side channels, can be completely mitigated by excluding any secrets from the CPU and the memory hierarchy. This is made possible by focusing on building a platform that only enforces program and data *integrity* and *not* confidentiality (which is sufficient for verifiable FHE, since all data remain encrypted at all times). All secrets related to the attestation mechanism are kept in a separate coprocessor (e.g., a TPM)—inaccessible to any software-based attacker. Relying on a discrete TPM typically incurs significant performance overhead, which is why (insecure) software-based TPMs are used in practice. As a second contribution, we show that for FHE applications, the attestation protocol can be adapted to only incur a fixed cost.

Argos requires no dedicated hardware extensions and is supported on commodity processors from 2008 onward. Our prototype implementation introduces 6% overhead to the FHE evaluation, and 8% for more complex protocols. In particular, we show that Argos can be used for real-world applications of FHE, such as private information retrieval (PIR) and private set intersection (PSI), where providing verifiability is imperative. By demonstrating how to combine cryptography with trusted hardware, Argos paves the way for widespread deployment of FHE-based protocols beyond the semi-honest setting, without the overhead of cryptographic proofs.

## 1 INTRODUCTION

Fully homomorphic encryption (FHE) [28, 33, 46, 55] makes it possible to evaluate any logical circuit directly on encrypted data. In the client-server setting, it can be used to build other powerful cryptographic primitives such as private set intersection (PSI) [2, 35, 64, 80, 81], private information retrieval (PIR) [27, 29, 35, 69] or multi-party computation (MPC) [5, 41]. FHE also has many practical use cases including private contact discovery [84], private smart contracts [105], or private inference [87]. Reducing performance overhead (currently 3 to 7 orders of magnitudes over non-private computation) has been the main focus of the last decade of research, but other limitations of FHE are now becoming more relevant.

At the forefront of these issues, the fact that FHE-schemes are not secure when considering a malicious evaluator, or that existing FHE schemes lack notions of integrity, are significant limitations for real-world deployment [6, 24, 114]. First, if a malicious server is able to supply malformed ciphertext for a client to decrypt, it can mount key recovery attacks [31, 61] and break all security (if the attacker recovers the private key, it can now decrypt the client’s private requests). Second, if FHE enables delegation of computation on confidential data to an untrusted party, existing constructions cannot help the user *verify* that the correct function was evaluated. In other words, FHE schemes always assumes honest-but-curious attackers. In the case of smart contracts, for instance, a malicious evaluator can completely change the functionality of the contract, with the option to even reveal confidential inputs and violate privacy. A straightforward solution is to add a layer of verifiable computation on top of the FHE scheme. Unfortunately, existing solutions suffer from major limitations. Cryptographic proofs incur an overhead of 4 to 6 orders of magnitude on top of FHE [114]. Another solution is to replicate the evaluation across several non-colluding parties, but distributed-trust setups are often impractical in the real world [77].

Trusted execution environments (TEEs) or secure enclaves [32, 39, 94, 104] have also been considered as a potential solution [71, 87, 114]. These technologies implement remote attestation, providing processors the capability to attest the code they are running, even in the presence of a privileged software attacker. TEEs offer better performance than cryptographic proofs, but come with weaker security guarantees: an attacker does not need to break a cryptographic assumption to forge a proof. Instead, the attacker needs to hack the platform or steal the attestation key. Still, these platforms initially seemed to offer a reasonable compromise. Unfortunately, a long series of published attacks (the SoK of Li et al. [75] covers 43 of them) quickly revealed that they were extensively vulnerable to microarchitectural side channels and transient execution attacks [26, 57, 102, 110, 112, 120]. These powerful attacks can be mounted in software and exploit microarchitectural structures and speculative mechanisms present in modern processors to extract secrets from trusted environments, not only violating confidentiality of the protected programs but also extracting keys used for attestation. New microarchitectural side channels are discovered every year [82, 116] and the current consensus is that it will take many years of research to build software and hardware defense mechanisms that can efficiently eliminate these threats from modern processors. On top of this existential threat, the availability of TEEs is also limited. Modern platforms require dedicated hardware extensions and are currently only available on server-grade processors, limiting adoption.

**This Work.** We present Argos, the first enclave platform specifically designed to build maliciously-secure and verifiable FHE. As observed in previous work [71, 87, 108, 114], for use cases where all application data is encrypted, such as in FHE, the only secrets exposed on the server are those used for attestation. We go further and show how this is a unique opportunity to rethink TEE architecture and design *integrity-only* enclaves that focus on data and code integrity *without* offering any confidentiality guarantee. Lessening the requirements for enclaves gives us the opportunity to 1) completely eliminate microarchitectural side channels by secluding all secret key material in a physically separated coprocessor (TPM), 2) use a hypervisor-based TEE architecture that offers better performance and hardware compatibility, and 3) build a simplified attestation scheme that makes it possible to easily prove the security of our construction and amortize the cost of relying on a slow coprocessor for key operations.

**Eliminating Microarchitectural Side Channels.** An important observation we make is that side-channel attacks are “read-only” gadgets and can be eliminated simply by excluding any secrets from the CPU and the memory hierarchy (see fig. 1). Software-mounted side channels require an attacker program to share hardware resources with the victim programs in order to extract secrets. If all secrets are secluded on a separate chip, such as a physical trusted platform module (TPM) or a secure coprocessor [3, 37, 67], no attacker program can ever share microarchitectural resources with the cryptographic algorithm and extract secrets. Leveraging this insight, we design Argos to store all attestation secrets inside a discrete TPM and delegate all sensitive operations outside of the CPU. This simple principle makes Argos secure *by construction* against side-channels. The only (side) channel left between the secluded cryptographic keys and a potential attacker is completion time [83] which is addressed by using constant time cryptography in the TPM (see our threat model in section 2).

**Rehabilitating Hypervisor-based TEEs.** Argos relies on a hypervisor-based TEE design. A security monitor runs with hypervisor privilege and provides code integrity and isolation to enclave programs. Such TEEs were introduced by TrustVisor [78] and are still widely used by cloud vendors to implement confidential computing [4]. These architectures are usually considered less secure than dedicated hardware platforms such as Intel SGX or AMD-SEV as they do not encrypt main memory and offer no security against Coldboot [62], one of the most practical physical attacks (see table 7). In our context, this is not an issue as no secrets are ever exposed to main memory. As a result, despite our strong threat model, Argos can adapt this architecture to bring compatibility with most commodity hardware. In terms of performance, existing platforms typically rely on virtual TPMs [99]—one per execution environment—that are endorsed by the root TPM but executed on CPU [8, 10]. While this approach avoids the performance bottleneck of a single physical TPM [78], it makes them vulnerable to microarchitectural side-channels. In contrast, Argos prefers a secure-by-design approach, centered on the discrete TPM.

**Performance Optimization With a Discrete Chip.** Relying on a single discrete TPM can have a significant impact on performance [79]. Nevertheless, we can still virtualize most TPM resources in the security monitor as measurement (e.g., hashing) does not require any secret manipulation. Our use case also excludes the

use of the TPM “sealed” secret storage. That means we only use the TPM as a “signing oracle.” This can still be expensive, especially when considering interactive cryptographic protocols that require repeated and attested message exchanges.

**Simple and Efficient Attestation Scheme.** We show that for our FHE applications, we can build a simple attestation scheme that only ever requires one TPM signature. Unlike in other enclave platforms, the remote user does not need to see any attestation proof before sending its sensitive data, as it will always stay encrypted. This also applies to intermediate protocol messages that do not need to be attested but simply added to a transcript whose hash is extended in the security monitor. At the end of the FHE computation (or more precisely, before any FHE ciphertext needs to be decrypted), the transcript’s hash will be signed using the TPM. We show how our simplified transcript-based attestation can be model as a proof system, and show it is sufficient to achieve malicious security needed for verifiable FHE. We also show how Argos can be extended to support batched verifiable FHE. As a result, the overhead of using a physical TPM becomes a fixed cost, and performance becomes similar to that of a virtual TPM, while enforcing security-with-abort in the presence of a malicious attacker [56].

**Extending Argos to FHE-Based Applications.** Circuit-level security sometimes differs from application-level security. In many FHE-based applications, a malicious server can provide corrupted inputs and gain some information from how the client behaves following decryption. This can have devastating downstream security implications. We show how Argos can easily be extended to support complex protocols in the malicious setting, such as authenticated PIR or authenticated PSI with almost no overhead compared to the semi-honest schemes.

**Implementation And Evaluation.** We implement and evaluate Argos on real hardware. Our prototype runs on Intel x86 platforms, but Argos is easily adaptable to other vendors and architectures (e.g., AMD X86, ARM or RISC-V) and compatible with most commodity machines past 2008. To guarantee the TCB integrity at boot-time, we use commodity hardware root-of-trust technologies. Our prototype uses a TPM, but our architecture is adaptable to other hardware roots of trust and secure coprocessors such as the Apple Secure Enclave [3] or Open Titan [67]. Our security monitor is an open source fork of Tyche [23], modified to support our attestation scheme and remove side-channel protections, but our approach is also compatible with other micro-hypervisors like seL4 [70]. We implement a custom minimal runtime for FHE and the SEAL library [28], emulating system calls and providing randomness through a hardware random number generator (i.e., RDRAND instructions) that is not under OS control. For more complex applications that require a broader class of system calls, Gramine [109] can be used for better compatibility, at the price of a larger TCB and some performance overhead (5x slower startup time, for example). Argos has a limited attack surface with a minimal trusted code base (TCB) of less than 18KLOC, plus the target FHE application ( $\approx$  50KLOC). Our evaluation shows that Argos is 80 times faster than previous work leveraging Intel SGX for FHE integrity [114] with a minimal average performance overhead of 3% for FHE evaluation. We show that Argos can be used to implement more complex protocols such as attested PIR and attested PSI with performance overheads under 8%

and without the significant offline communication costs incurred by cryptographic solutions.

To summarize, our main contributions are:

- Argos, the first *integrity-only* enclave platform designed to build maliciously-secure verifiable FHE, which:
  - can be used to build fully malicious and authenticated PSI and PIR schemes;
  - is secure by construction against microarchitectural side channels and transient execution attacks;
  - requires no specialized hardware and is compatible with commodity processors from 2008 onward; and
  - only incurs 3% average performance overhead for FHE computation, less than 8% performance overhead for complex protocols, and virtually no communication overhead;

## 2 THREAT MODEL

Our goal is to ensure integrity of the data and computation running inside our enclave environment. Most importantly, we do not protect confidentiality of our programs nor do we protect them against denial of service. We assume a strong adversary collocated on the server that has compromised the majority of the software stack including the OS, and can mount any microarchitectural side-channel and transient execution attacks. Because Argos is secure against all software-mounted side channels that underlie known transient execution attacks, we will refer to all these attacks as “microarchitectural side channels” throughout the rest of the paper. Our trusted computing base (TCB) consists only of a security monitor (18KLOC) and our application, both assumed to be bug-free. We assume that the hardware is *functionally* correct and that the TPM’s cryptography is properly implemented with constant-time programming. Rowhammer and fault injection attacks are not as practical and are considered out of scope. We also protect against the majority of physical attacks (cold boot, BIOS tampering, physical side channels) on the main processor and DRAM, but for the rest of this paper, we will focus on software-mounted attacks and consider them out of scope. A detailed discussion on the remaining attack surface can be found in section 9.7.

## 3 BACKGROUND & MOTIVATION

### 3.1 FHE Schemes And Client-Server Setup

We consider a client-server setup where a server evaluates a logical circuit on some client ciphertext obtained using an FHE scheme. An FHE scheme is usually defined as a tuple of algorithms  $(\mathcal{E}.\text{Gen}, \mathcal{E}.\text{Enc}, \mathcal{E}.\text{Eval}, \mathcal{E}.\text{Dec})$ . The client generates a key pair using  $\mathcal{E}.\text{Gen}$  and encrypts some input  $x$  using  $\mathcal{E}.\text{Enc}$ . The server then uses  $\mathcal{E}.\text{Eval}$  to evaluate a circuit  $f$  on the ciphertext  $c_x$  and sends the resulting ciphertext  $c_y$  back to the client to decrypt using  $\mathcal{E}.\text{Dec}$ . Concretely:

- $\mathcal{E}.\text{Gen}(1^\lambda) \rightarrow (\text{pk}, \text{sk})$
- $\mathcal{E}.\text{Enc}_{\text{key}}(x) \rightarrow c_x$
- $\mathcal{E}.\text{Eval}_{\text{pk}}(c_x, f) \rightarrow c_y$  where  $y = f(x)$
- $\mathcal{E}.\text{Dec}_{\text{sk}}(c_y) \rightarrow y$

FHE schemes must enforce the following properties:

**Correctness.** A scheme is *correct* if any honest computation will decrypt to the expected result (i.e.  $\text{Dec}_{\text{sk}}(c_y) = f(x)$  with high probability).

**Security.** All FHE schemes are secure against chosen plaintext attacks, which means that an attacker with the public key cannot distinguish between the encryption of two different messages of its choice.

### 3.2 Semi-Honest vs. Malicious Security

In a honest-but-curious or semi-honest setting, an attacker is allowed to observe intermediate messages and passively infer information in order to break security but cannot deviate from the agreed-upon protocol. In a malicious setting, the attacker is much more powerful and allowed to deviate arbitrarily from the protocol, e.g., to lie to the victim or temper with messages. The semi-honest setting is often used in theory but has limited viability in a deployment context where assuming a malicious attacker is more realistic.

### 3.3 Why Does FHE Need Verifiability?

Most state-of-the-art FHE schemes are *insecure* to use in real-world settings [28, 33, 46, 55]. Because these schemes are not CCA-secure [13], if an attacker is able to supply malformed ciphertexts and observe decryption outcomes (for example, by observing client behavior), it can mount key recovery attacks and completely compromise privacy [25, 31, 50, 61]. Beyond the “circuit-level” security, the absence of integrity for FHE schemes also has implications for FHE-based applications. A malicious server can arbitrarily modify the circuit they evaluate on the encrypted data. Attacks on correctness are an obvious issue, but these can also translate to attacks on privacy at the application level. Let us take the example of private contact discovery. Alice wants to discover who among her contacts is using a private messaging service [84]. Alice encrypts her address book, which contains all of her contacts’ private information, and sends it to the server. The server now has access to Alice’s encrypted set of contacts. If malicious, it could simply return the entire set to Alice. This could mislead Alice into thinking that one of her friends (Bob) uses the messaging service. Consequently, Alice contacts Bob via the service, leaking his phone number. This would violate Bob’s privacy and defeat the use of private contact discovery.

### 3.4 Why do we need a new TEE platform?

Trusted execution environments (TEEs) provide hardware-based isolation for confidential programs while minimizing the trusted code base (see fig. 1). The existing TEE landscape spans from platforms that protect small “enclave” programs [7, 39, 94] to those that implement confidential computing for entire virtual machines [4, 32, 76, 104]. However, the discovery of numerous microarchitectural side-channel and transient execution attacks has severely compromised their confidentiality guarantees. Side channel attacks exploit shared microarchitectural structures such as memory caches [19, 57], translation look-aside buffers (TLB) [58], branch predictors [1, 49], and DRAM controllers [115] to enable information leakage between security domains.

TEEs also provide integrity guarantees through remote attestation, allowing a remote client to verify the authenticity of the hardware and the initial state of the execution environment. However, even these attestation mechanisms have proven vulnerable

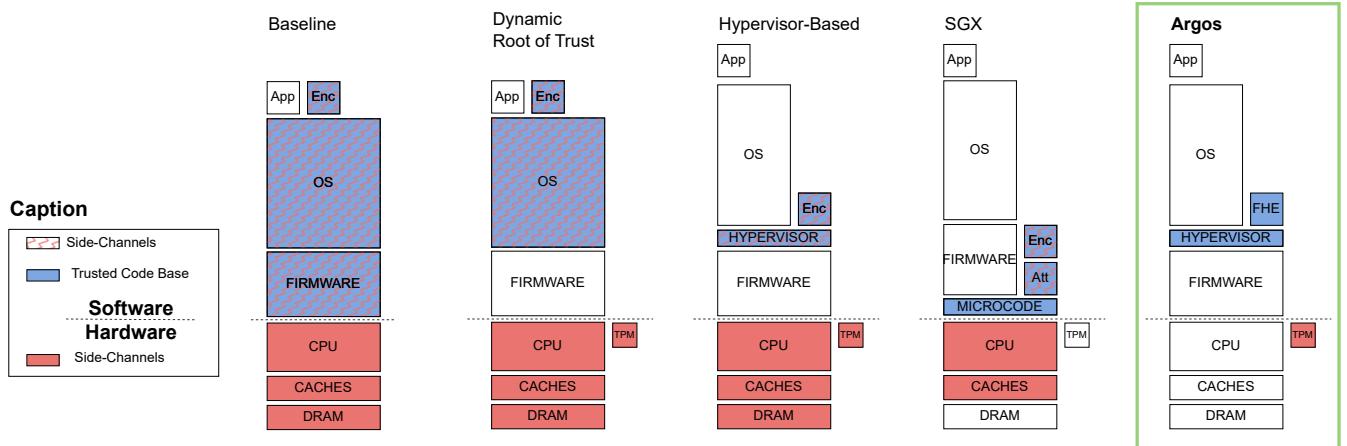


Figure 1: Evolution of the attack surface on TEE platforms. Enc: enclave program, Att: attestation enclave.

to side-channel attacks that can extract secret keys used by platforms to sign attestation reports [26, 102, 112], leading the cryptographic community to lose faith in TEEs. On the other hand, discrete trusted platform modules (TPMs) and associated dynamic root of trust (DRoT) have long served as dedicated hardware for platform integrity and secure storage against software-based attackers. However, their discrete nature also creates a severe performance bottleneck when attesting multiple security domains (Flicker [79] incurs 2-3 order of magnitude overheads). Although TPM virtualization [8, 10, 78, 99] can address these performance issues, in turn they expose sensitive key material to microarchitectural side channels (see section 10 for more detailed comparison of Argos and other platforms).

To rebuild trust in hardware security primitives, we must develop new systems that are both inherently secure against microarchitectural side channels and capable of outperforming functionally equivalent cryptographic solutions. We address this challenge by focusing specifically on FHE applications and building a TEE platform that *only enforces integrity* of the enclave program, which is sufficient for verifiable FHE since data remains end-to-end encrypted throughout the computation.

## 4 INSIGHTS & DESIGN PRINCIPLES

### 4.1 FHE Applications Do Not Expose Secrets

In FHE applications, no sensitive data is ever manipulated in clear. This has several implications that can help simplify the design of a TEE platform and the attestation protocol.

**No Secrets Exposed.** In standard TEE platform, a central feature is to enforce the confidentiality of the execution environment. In FHE, all data is encrypted. As a result, the content of our enclave program (i.e., the state of  $E.\text{Eval}$ ) is public. This means our application is not vulnerable to side channels, and we do not need to harden the execution environment against the usual threats.

**No Long-Term Secret Storage.** In usual TEE platforms, an important feature is the ability for an enclave program to *seal* secrets and recover them later. This property is usually tied to the identity (or the measurement) of the enclave program, and a secret can only

be unsealed by a program that matches the correct measurement. This is, for instance, what is used in Bitlocker to only ever release the disk encryption key to a correctly booted system. TPMs and TEE platforms such as SGX all offer a seal operation. FHE evaluation does not require long-term secret storage, which significantly reduces the features required by the TEE.

**No Need To Establish Trust.** In usual remote attestation protocols, the remote client first needs to verify the TEE attestation to establish trust. Only once it trusts the TEE to have been correctly setup will it send its encrypted private input. This is because the inputs will then be decrypted inside of the TEE. In our case, inputs are never decrypted, which means that the remote client does not need to establish trust before sending the encrypted input. This makes it possible to amortize the cost of attestation by only signing the final transcript of program inputs and outputs (see section 5).

### 4.2 Eliminating Microarchitectural Side Channels Using A Physical TPM

Because our FHE-applications do not expose secrets, the only secret ever present on the server is the private signing key used in the attestation scheme. This is a much smaller attack surface than in usual TEE systems. We can take advantage of that opportunity to completely eliminate microarchitectural side channels by placing the signing key in a microarchitecturally isolated physical TPM. Because no secrets are left on the CPU or in main memory, this effectively makes our platform completely secure against microarchitectural side channels.

**Side Channels Are Read-Only.** One important element to keep in mind is that side channels are “read-only” gadgets and can only extract secrets from a victim program. Specifically, a side channel cannot modify the state of a victim’s program memory or change its control flow. These attacks are called fault injections and have not been as successful at attacking TEE platforms (see section 9.7).

**Side Channels Require Shared Resources.** To mount a microarchitectural side channel, an attacker needs to trigger a transmitter that will access the secret and modulate a channel to transmit information to a receiver under the attacker’s control. That means

the transmitter and receiver programs need to share some microarchitectural state in order for the transmission to be possible.

**Physical TPMs Are Microarchitecturally Isolated.** TPM are commonly implemented on a discrete chip or as firmware on a secure co-processor (we cannot secure software TPM running on the CPU). That means they are communicating with the main CPU using a hardware bus. They do not share cache nor any processor resources with the CPU. They use their own private resources, like private memory and registers, that are only accessible to the TPM. As a result, if the secret key is kept in a physical TPM, it is microarchitecturally isolated from the attacker and there is no possibility for the secret to be extracted through a microarchitectural side channel. Intrusive physical attacks might still be possible but are outside of our threat model (see section 9.7).

### 4.3 System Overview

Our system is composed of several hardware and software elements (see fig. 1). The security monitor is a small ( $\approx 12\text{KLOC}$ ) trusted piece of code running at the hypervisor privilege level. It is in charge of enforcing isolation between the different security domains and integrity of the enclave programs. Isolation is enforced by hardware mechanisms (see section 8.2) and integrity is enforced through our attestation scheme (section 5). The security monitor does not manipulate any secret key material and delegates all cryptographic operations that require a key to the physical TPM. Integrity of the security monitor is enforced by the dynamic root-of-trust mechanism (section 8.1). Running on top of the security monitor, a large untrusted operating system is in charge of all resource allocations. It contains all the complex logic, for instance, for memory allocation or scheduling. It also provides many services for (enclave) applications such as networking, or IO.

**Running an Enclave Program.** To setup an trusted execution environment and securely run an enclave program, the untrusted OS interacts with the security monitor through vmcalls. On behalf of the OS, the security monitor will reclaim some resources (memory and cores) and allocate them to the enclave. Thanks to the isolation mechanisms under the control of the security monitor, the OS is now unable to modify the state of the enclave memory. The security monitor will also enforce some invariants (setting the page table correctly and zeroing regions where no content is loaded), ensure the correct binary is loaded, and attest the enclave initial state (see section 5.4). Once loaded, the enclave can be started. The remote client sends inputs over the network. They are received by the OS and placed in shared memory so that the enclave can access them. Inputs are systematically copied to enclave's private memory before being used. Any output is also copied to shared memory and sent back to the client by the OS.

## 5 THE ARGOS ATTESTATION SCHEME

### 5.1 Existing Attestation Protocols

Attestation schemes are generally interactive protocols that involve several rounds of communication between a client and a server. First, before any sensitive data is sent, trust needs to be established between the execution environment running server side and the remote client. The client will send a challenge and receive an attestation report, endorsed by the manufacturer and attesting that a

genuine piece of hardware has been correctly set up to instantiate a trusted execution environment. Once trust has been established, the enclave program and the remote client can perform key exchange and establish a trusted communication channel. They can then start communicating and running the expected application. Any message coming out of the enclave needs to be encrypted and attested using the private key generated inside of the enclave to prevent a malicious attacker from observing or tempering with the content of the messages. Similarly, every message from the remote client needs to be decrypted and its integrity verified.

### 5.2 Overhead of Using a Discrete TPM

One of Argos's design principles dictates that no secret key material should ever be present on the main CPU. This means that each enclave would need to delegate all key operations to a single physical TPM. Even if only considering integrity protection (signature of messages and attestation protocol), this implies that the TPM becomes a serious bottleneck for the platform. Performing a signature using a discrete TPM takes approximately 196 ms, and the numbers are similar for integrated TPMs. This can potentially be a significant overhead for any application that requires several rounds of communication. However, Argos's focus on FHE applications allows us to minimize much of this overhead.

### 5.3 Attested Transcript for FHE Applications

Trust does not need to be established before sending encrypted data. Because intermediate results are also encrypted, privacy cannot be compromised until a ciphertext is decrypted. That means integrity only needs to be checked *just before* decrypting an output (see section 6 for our security analysis). In other words, each message of the FHE application does not need to be individually attested. Instead, the integrity-protected security monitor can keep a transcript of all relevant data (i.e. all application inputs, intermediate results, and final output), and only perform a single TPM signature over that transcript. This simplifies the attestation scheme and drastically reduces the cost of using a physical TPM as remote attestation is now a simple fixed cost.

### 5.4 Argos Attestation Scheme

**Measured Boot.** Upon boot, the hardware DRoT is used to derive a measurement of the Argos security monitor (see section 8.1). This measurement is securely stored in the TPM registers (PCRs).

**Enclave Setup.** When a new enclave is initialized, the security monitor will allocate some of its private memory to the enclave transcript. During enclave set up, the security monitor measures the enclave binary and stores it in the transcript. In our specific case, the binary contains the FHE circuit, small public inputs (e.g., FHE parameters), alongside the logic to evaluate the circuit (the SEAL library in our case). It also contains the rest of the application logic. All the values of the relevant architectural registers that represent the initial state of the execution environment are also added to this measurement. For instance, the enclave initial program counter, page tables, the stack pointer and the addresses of shared memory use to communicate with the OS. The measurement of the binary and initial state unambiguously represent a function that takes inputs from shared memory and outputs results to it.

**Enclave Lifetime.** During the lifetime of the enclave, the transcript is appended with hashes of the data received and sent by the enclave. This includes all inputs and outputs.

**Attesting the Transcript.** Once the final output has been generated and the transcript extended with it, the application performs a special vmcall to the security monitor to request the attested transcript. The security monitor then leverages the TPM to obtain a signature over the transcript. The application will then be able to transmit the attested transcript to the remote client, along with the final output.

**Client Verification.** Upon receiving the result(s) and the attested transcript, the client first verifies the attestation. The verification of the attestation goes as follows: First, the client verifies that the key used to sign the transcript belongs to a genuine TPM (see appendix A for more details). It then checks that the measurement for the security monitor contained in the transcript signature matches a reference measurement that is known to be correct (see section 9.7 on how such a reference measurement can be obtained). Finally, it checks that the binary references in the transcript corresponds to the correct FHE application and that it was evaluated on the input it provided. If all checks are satisfied, the verification is successful. In case of a verification failure, the server is behaving maliciously and the client aborts.

**Output Decryption.** If and only if verification succeeds, the client decrypts the FHE output.

### 5.5 Remote Attestation as a Proof System

Our simplified attestation scheme is now functionally equivalent to a proof system. We can formalize it as a tuple  $(\Pi.\text{Gen}, \Pi.\text{Prove}, \Pi.\text{Verify})$ . The server runs  $\Pi.\text{Gen}$  to generate a key pair. It can then use the private key to run  $\Pi.\text{Prove}$  and generate an attested transcript  $(\pi_y)$  over some client input  $x$ , an application  $g$  (formally a function, in practice a binary and an initial state for the execution environment), and the resulting output  $y$ . Given the public key, the client can then run  $\Pi.\text{Verify}$  to verify that the attested transcript  $\pi_y$  is correct, that is,  $y = g(x)$ . In summary, we have:

- $\Pi.\text{Gen}(1^\lambda) \rightarrow (\text{pk}, \text{sk})$
- $\Pi.\text{Prove}_{\text{sk}}(x, g) \rightarrow (y, \pi_y)$  where  $y = g(x)$
- $\Pi.\text{Verify}_{\text{pk}}(y, x, \pi_y) \rightarrow \text{true/false}$  where the client accept the output is true

Additionally, our scheme satisfies these two properties:

**Completeness.** The system is *complete* if  $\Pi.\text{Verify}$  will always accept an honestly computed result i.e., a correctly formed transcript where  $y = g(x, w)$ .

**Soundness.** The system is *sound* if an adversary cannot make  $\Pi.\text{Verify}$  accept an incorrect answer i.e., a transcript where  $y \neq g(x, w)$ .

**Modeling Several Rounds of Communication.** For applications that require several rounds of communication, we only generate one attested transcript at the end (see section 5.3). That means  $x$  is the concatenation of all client inputs/messages. Our specific implementation simply shares intermediate results (the server messages) with the client. It also caches some internal state for  $\Pi.\text{Prove}$  (such as the measurement of  $g$ ) so it can efficiently generate the final

attested transcript once provided with all client inputs.

**Security Analysis.** Completeness of our attestation scheme reduces to the completeness of the underlying signature scheme used in the TPM. Soundness is more complex to establish, as it requires informally laying down all of our system security. First, as we explained earlier, an attacker (as described in our threat model) is not able to extract the private key from the TPM. Second, the security of our DRoT guarantees that the security monitor was correctly loaded and its integrity cannot be compromised. Third, our assumption of the TCB tells us that the security monitor is correctly implemented and bug-free. As a result, attested transcripts will only be generated for enclaves that are correctly set up, loaded, and executed. Our assumption on the underlying hardware isolation mechanisms (execution modes and nested page tables) guarantees that the integrity of the enclave code and data is maintained over the course of its lifetime. Finally, the functional correctness of the hardware guarantees that when an enclave is run that implements functionality  $f$  on input  $x$ , the resulting output  $y = f(x)$ . With all of this in place, the soundness of our attestation scheme reduces to the soundness of the TPM signature scheme.

## 6 CIRCUIT-LEVEL VERIFIABLE FHE

We first show how, using Argos, we can take a semi-honest FHE scheme and build a circuit-level verifiable FHE scheme that achieves malicious security. In section 7, we show how we can extend Argos to satisfy application-level security.

In our client-server setup, all state-of-the-art FHE schemes [28, 33, 46] are *insecure* when considering a malicious server. Let us take a simple example where a client wants to outsource some sensitive computation to a server using FHE. Here, the server might be able to infer information regarding the decryption of a ciphertext by observing the client reaction. For instance, in FHE, if the noise in the ciphertext overflows a given threshold, client-side decryption might fail, which could lead to the client emitting a new request to the server. This is not a problem if assuming an honest-but-curious server. However, if the server is malicious, it can tweak the response to the client in arbitrary ways and observe its reaction. This is equivalent to providing our attacker with access to a (limited) decryption oracle that outputs “success” or “failure” for ciphertexts under the attacker’s control. This is already enough to mount key-recovery attacks [31, 61] and break all security for the FHE scheme (if the attacker recovers the private key, it can now decrypt the client’s private requests). In fact, this is much more problematic than a simple “correctness” issue where the client would just get a wrong output. In this simple outsourcing scenario, enforcing integrity at the *FHE-circuit-level* is enough, as it prevents the attacker from providing malformed ciphertexts to be decoded. When the client receives the response from the Argos server, it will first verify the attested transcript to ensure that the correct circuit was evaluated on the correct inputs. If verification fails, the client aborts, does not decrypt the ciphertext, and detects malicious behavior from the server. If and only if verification succeeds, the client may decrypt the ciphertext with the guarantee that it was correctly formed.

## 6.1 Definition

Inspired by the formalism of Viand et. al. [114], we define circuit-level verifiable FHE (vFHE) as follow. A vFHE scheme is defined as a tuple of algorithms ( $\text{Gen}$ ,  $\text{Enc}$ ,  $\text{Eval}$ ,  $\text{Verify}$ ,  $\text{Dec}$ ). The client is equipped with a key pair  $(\text{pk}_c, \text{sk}_c)$  and the server with  $(\text{pk}_s, \text{sk}_s)$ . The client uses  $\text{Enc}$  to encrypt some input  $x$  and obtains a ciphertext  $c_x$ . Using  $\text{Eval}$ , the server evaluates a circuit  $f$  over  $c_x$  and generates an encrypted output  $c_y$ , along with an attested transcript  $\pi_y$ . When the client receives the output, it will first  $\text{verify}$  the transcript is correct ( $c_y = f(c_x)$ ) and if and only if it is, use  $\text{Dec}$  to decrypt  $c_y$  and recover  $y$ . In summary:

- $\text{Gen}(1^\lambda) \rightarrow ((\text{pk}_c, \text{sk}_c), (\text{pk}_s, \text{sk}_s))$
- $\text{Enc}_{\text{pk}_c}(x) \rightarrow c_x$
- $\text{Eval}_{\text{pk}_c, \text{sk}_s}(c_x, f) \rightarrow (c_y, \pi_y)$  where  $y = f(x)$
- $\text{Verify}_{\text{pk}_s}(c_y, c_x, \pi_y) \rightarrow \text{true/false}$  where the client accept if true, aborts otherwise.
- $\text{Dec}_{\text{sk}_c}(c_y) \rightarrow y$

Correctness is defined identically to a semi-honest FHE scheme (see section 3.1). Security differs as our attacker is now malicious and can send arbitrary  $c_y$  values. We also add two properties compared to the vanilla FHE-scheme.

**Completeness.** A vFHE scheme is *complete* if  $\text{Verify}$  will always accept an honestly computed result.

**Soundness.** Finally, a scheme is *sound* if an adversary cannot make  $\text{Verify}$  accept an incorrect answer.

## 6.2 Putting Everything Together

Let's consider an FHE scheme  $(\mathcal{E}.\text{Gen}, \mathcal{E}.\text{Enc}, \mathcal{E}.\text{Eval}, \mathcal{E}.\text{Dec})$  and our remote attestation scheme  $(\Pi.\text{Gen}, \Pi.\text{Prove}, \Pi.\text{Verify})$ . Argos' construction is as follow:

- $\text{Gen}(1^\lambda) \rightarrow ((\text{pk}_\mathcal{E}, \text{sk}_\mathcal{E}), (\text{pk}_\Pi, \text{sk}_\Pi))$   
with  $(\text{pk}_\mathcal{E}, \text{sk}_\mathcal{E}) = \mathcal{E}.\text{Gen}(1^\lambda)$  and  $(\text{pk}_\Pi, \text{sk}_\Pi) = \Pi.\text{Gen}(1^\lambda)$
- $\text{Enc}(x) \rightarrow c_x$  with  $c_x = \mathcal{E}.\text{Enc}(x)$
- $\text{Eval}(c_x, f) \rightarrow (c_y, \pi_y)$  with  $(c_y, \pi_y) = \Pi.\text{Prove}(c_x, \mathcal{E}.\text{Eval}(\cdot, f))$
- $\text{Verify}(c_y, c_x, \pi_y) \rightarrow b$  with  $b = \Pi.\text{Verify}(c_y, c_x, \pi_y)$ .
- $\text{Dec}(c_y) \rightarrow y$  with  $y = \mathcal{E}.\text{Dec}(c_y)$

The key here is to use the attestation mechanism to attest the correct execution of  $\mathcal{E}.\text{Eval}(\cdot, f)$  represented by the binary and initial state of the execution environment.

**Security Proof.** Correctness of our construction reduces to the correctness of the initial FHE scheme. Similarly, completeness and soundness directly derive from the completeness and soundness of the attestation scheme. Security is more subtle. Remember that our initial FHE scheme is not secure if provided with malformed ciphertexts. That means we can only rely on the security of the underlying FHE scheme if our attacker cannot gain control over the ciphertexts. This is enforced by the verification step. If the client runs  $\text{Dec}$  on an output received from the server, that means the output and the transcript had to successfully pass  $\text{Verify}$ . If this is true and the attacker was able to supply a ciphertext of its choice, i.e.,  $c_y \neq \mathcal{E}.\text{Eval}(c_x, f)$ , that means the attacker was able to break the soundness of the underlying attestation system, which is a contradiction. As a result, our construction is secure,

and we have successfully built a circuit-level verifiable FHE scheme.

**Batching FHE Evaluations.** Running  $\Pi.\text{Prove}$  is expensive (it requires a TPM sign operation). Because evaluations of different FHE circuits are independent, we can batch them for performance optimization and only produce one attested transcript for all these evaluations.  $\text{Enc}$ ,  $\text{Verify}$ , and  $\text{Dec}$  are the same, and security still holds as long as  $\text{Dec}$  is never run if verification was not successful.

## 7 EXTENDING ARGOS TO FHE-BASED APPLICATIONS

Circuit-level security sometimes differs from application-level security. For instance, attacks on correctness can escalate to attacks on privacy (section 3.3). Additionally, the security of our FHE schemes breaks if a malicious server is able to supply malformed ciphertexts for decryption section 3.3. This is potentially dangerous when, in many FHE-applications, the server itself can provide public or private inputs and act as a party in the computation. We explore two such applications, 1) private information retrieval, and 2) private set intersection. For each application, we discuss the challenges of deploying FHE-based schemes a real-world setting and show how Argos can be extended to enforce malicious security.

### 7.1 Malicious Private Information Retrieval

Private information retrieval (PIR) allows a client to retrieve an element from a public database hosted on a remote server, without the server learning which element was requested by the client. PIR has many applications such as public-key directory [34], certificate transparency logs [64, 103], private web search [63] or private analytics [59]. PIR can be instantiated in a multi-server or single-server setting using a range of cryptographic primitives, but in this paper we focus on the state of the art for single-server PIR which are FHE-based schemes [2, 64, 80, 81].

Most schemes assume a honest-but-curious server; however, this is not a realistic setup for real-world deployment. Furthermore, circuit-level integrity alone is not enough to enforce malicious security. A malicious server can always modify the values in the database, which could break not only the correctness (send the wrong answer for a given query), but most importantly security (i.e., privacy of client input) by mounting *selective failure attacks* [72]. We show how Argos can be extended to secure semi-honest FHE-based PIR scheme in a malicious setting.

**Defending Against Selective Failure Attacks.** In these attacks, a malicious server chooses an entry of the database and replaces the entry with a value that it knows will cause a decryption failure. This is possible because in FHE, using an input that is not in the correct plaintext space (for instance, an integer that is too large) can lead to a malformed ciphertext which will not decrypt properly. Because our attacker might be able to observe the client reaction to decrypting the output (e.g., the client will send the request again), if decryption fails, the attacker can infer with high probability that the client queried the malformed entry.

Circuit integrity is not enough here. Argos also needs to preprocess the public database and check that the server's public inputs belong to the plaintext space. If a check fails, Argos should abort the server execution and not serve any client request. Formally,

this is equivalent to allowing server inputs  $w$  for  $f$ , and in our construction, wrapping  $\mathcal{E}.\text{Eval}$  in a function that evaluates  $\mathcal{E}.\text{Eval}$  if and only if  $w$  is well formed and aborts otherwise.

**Authenticated PIR.** So far, our malicious PIR scheme is only able to guarantee that the server performed the protocol using a well-formed database but not a specific one. Adding authenticity to a PIR scheme is a desirable property and might be essential for some real-world deployments. For instance, in the case of a transparency log or a public-key directory, it might be essential to guarantee the correct database was queried. We take inspiration from previous work and solve this problem by having the server initially commit to the database. This assumes a reference commitment is available for the client to fetch out-of-band.<sup>1</sup> Argos can easily be extended to support this functionality by adding a fresh commitment to the database (i.e. a hash or the root of a Merkle tree) to the attestation transcript. When receiving the response to its query, a client should first verify the attestation transcript, including that the database commitment matches the reference one. If the verification fails, it should abort. If (and only if) it succeeds, it should then decrypt the response and proceed. In our formalism, this is equivalent to extending Verify to take a reference commitment  $H(w)$  as an input, and tweaking soundness to ensure an adversary cannot make Verify accept a transcript  $\pi_y$  that does not contain the correct commitment.

## 7.2 Malicious Private Set Intersection

Private set intersection (PSI) allows two distrustful parties with some private sets to compute the intersection of these sets without learning more information than the intersection itself. PSI has many applications, such as private contact discovery [84] or compromised credential checking [90].

Most efficient PSI protocols do not achieve (full) malicious security [27, 29, 35, 69]. Some constructions do [43, 88, 93, 98, 100, 101, 106] but at the price of high offline communication cost (see section 9.5). Using a state-of-the-art semi-honest PSI [35] scheme based on FHE, we show how Argos can be extended to make the scheme fully secure against a malicious server. Here, the main difference from the PIR setup is that the server inputs are not public, but private. State-of-the-art FHE schemes are not *circuit private*. That means the client might be able to learn information regarding the server’s private input by looking at the noise contained in the output. As it is standard in FHE-based PSI schemes, this requires to mask the server’s input using oblivious pseudorandom functions (OPRFs) to provide privacy for the server even in the presence of malicious clients. Moreover, security of the FHE scheme is enough to also provide privacy for the client, even in the presence of a “malicious” server. This is because the ideal PSI functionality considered in these works assumes a limited attacker with no observability of the client reactions [27]. However, this might not hold in a real-world deployment, as we explained before.

**Verifiable PSI.** Once deployed in a real-world setting, the attacker might be able to observe some of the client’s reactions. Attacks on correctness might have serious security implications, sometimes translating to attacks on privacy (see example for private contact

discovery in section 3.3). Fortunately, circuit-level integrity—as provided by Argos—can guarantee *verifiable* PSI, which means that verification of the attestation transcript succeeds if and only if the server correctly executed the PSI circuit on a private set. Nevertheless, for full malicious security, we still need to ensure that the server evaluated the over the *correct* set.

**Authenticated PSI.** If the server’s input set is not authenticated, a malicious server can arbitrarily modify the dataset without detection. This can have serious security implications in deployment. For instance, in the case of compromised credential checking, Alice reaches out to Mozilla’s server to ensure that her passwords were not compromised and are safe to use. If Alice indeed uses some compromised credential and the server is malicious, it could evaluate the PSI circuit on her private set, but using a malicious input set, crafted server side, ensuring no intersection is found. This might lead Alice to use her compromised credentials.

Similarly to Authenticated PIR, Argos can be extended to include a commitment to the dataset (e.g. Merkle tree or hash-based commitment scheme) in the attestation transcript. Here too, a public commitment needs to be made available out-of-band to the client, through a transparency log, for instance. Upon receiving a response to its request, the client first verifies the transcript and aborts if it does not match the correct commitment. Formally, the extensions we made for PIR are almost sufficient, we just need to tweak the wrapper function over  $\mathcal{E}.\text{Eval}$  to match our PSI functionality. Instead of only checking that  $w$  are well formed, the wrapping function should compute the corresponding OPRF values for  $w$ . The server’s privacy is guaranteed by the OPRF, the soundness of the attestation scheme, and the fact that commitment schemes are hiding i.e., they do not leak information about the server’s set.

## 8 IMPLEMENTATION DETAILS

### 8.1 Security Monitor

Our security monitor is based on the X86 version of Tyche [23], a Rust micro-hypervisor that exposes hardware resources and isolation primitives to better protect trust domains. Note that other “small” hypervisors such as XEN [12] are at least one order of magnitude larger or closed source and micro-kernels such as SeL4 [70]. We implement our attestation scheme on top of Tyche, including support for the hardware TPM and the interface to manage attested transcripts. We are also able to simplify some aspects of the platform, for instance, by removing side channel protections from Tyche (e.g., cache and CPU flush). One of the implementation challenges was to adapt Tyche and the Linux driver to support large (i.e., several GBs) enclaves with our custom runtime and manage memory allocation accordingly. To enforce the integrity of the security monitor, we use the TPM and a dynamic root of trust ((e.g., Intel TXT or AMD DRMT [53, 89])) to implement standard measure boot (see primer in appendix A). At the end of this process, two TPM registers (PCRs) are set with hash values that uniquely identify the security monitor binary. Later, when the security monitor needs to attest an application transcript, it will load the transcript hash in a third PCR and request a TPM Quote (a signature over select PCRs).

<sup>1</sup>This assumption is made by Authenticated PIR [34] but recent works [42, 45] have also shown how using a limited number of “validation queries”, the client can verify with high probability that the committed database is correct.

## 8.2 Memory Isolation Primitives

We use hardware extensions for virtualization to enforce isolation between security domains [111]. Extended or nested page tables [15], as implemented by Intel VT-x or AMD-V, make it possible for a guest domain to have an entire address space available and create a level of indirection for memory management that is exclusively under the control of the hypervisor. That means no execution domain can access other execution domains' address space nor modify its content. Other virtualization technologies such as AMD-Vi and Intel VT-d make it possible to protect domains from arbitrary DMA accesses. We use these technologies in Argos to enforce memory isolation for our enclave programs.

## 8.3 Custom Runtime

The execution environments provided by Argos are similar to bare-metal environments with no OS support. That means no syscalls, no program loader, no memory allocation, nor access to shared libraries. A runtime is required to support programs and modern libraries like SEAL that leverage many services from the OS. Gramine [109] is such a runtime and already supported by Tyche. It was conceived to securely run off-the-shelf applications on Intel SGX. As a result, it implements all Linux syscalls and also allow extensive support for shared library and memory management. However, because of its versatility, it is rather large (20KLOC) and adds some non-negligible overhead [96]. In order to reduce the TCB and improve performance, we implement our own custom runtime. It is minimal (870LOC) and consists of a simple memory allocator [92] and a handful of syscall handlers required by our application. For instance, we implement our own open and read syscalls to supply randomness from RDRAND when an application attempts to read `/dev/urandom`. We write our runtime to interface with MUSL [86] by interjecting on system calls and providing our own handlers. That means all applications and libraries (including SEAL [28]) need to be compiled statically and linked using our MUSL library, which requires some engineering efforts to port applications but makes it possible to obtain high performance *and* a very light runtime.

## 9 EVALUATION

### 9.1 Testbench

We prototype Argos on a **local Dell Optiplex** machine from 2017 equipped with an Intel Core i7-7700 processor clocked at 3.6GHz, a discrete TPM 2.0, and 8GB of RAM. We run our security monitor (adapted from Tyche [23] see section 8.1) directly on the hardware and evaluate our benchmarks using our custom SEAL runtime or Gramine [109]. We run Ubuntu 22.04 with the Linux kernel v6.2 and SEAL 4.1 with the BFV scheme. All benchmarks are all single-threaded, compiled using clang++ and using O3 optimization level. We also compare Argos with two other TEE platforms that have the following specs. SGXv2 running on an **Azure server** equipped with an Intel Xeon Platinum 8370C CPU with a frequency of 2.80GHz. AWS Nitro Enclaves [7] running on an **AWS server** equipped with an Intel Xeon Platinum 8259CL CPU with a frequency of 2.50GHz. The SGXv2 setup uses Gramine as a runtime while Nitro Enclave uses a docker container environment that runs its own Linux kernel. Note that under our threat model, both SGXv2 and AWS Nitro

Component	LOC
BIOS	1.5M
Linux	28M
Security Monitor	18K
Runtime (Custom / Gramine)	870 / 20K
SEAL Library [28]	20K
Application	1K–20K

**Table 1: Software Components and TCB Breakdown.**  
Red columns are excluded from the TCB.

Platform	Time
Software	42
vTPM	136
dTPM	195752

**Table 2: Signing operation ( $\mu$ s).**

	Setup	Attest	Term
SGX2	953	7	274
Nitro	3759	1	363
Argos+G	466	196	5
Argos	85	196	10

**Table 3: Fix-costs operations on a 1GB enclave (ms).**

Enclaves are *insecure* and are only used as comparison points for performance. See table 7 for a more complete comparison. All numbers are obtained by averaging values over 10 runs with minimal variance.

### 9.2 TCB Evaluation

Detailed count for lines of code (LOCs) can be found in 1. Our TCB is small with a minimum of 40KLOC for our PIR application and a total of 60KLOC for our PSI application. Note that we did not account for the microcode used in the DRoT or the TPM as it is not public information, and we consider these elements as part of the hardware. This should be compared to other TEEs and system projects that try to keep a small footprint. For instance, Nitro Enclaves include an entire Linux kernel which blows up their TCB in the range of MLOC, several orders of magnitude bigger than us. On the other hand, SGX would not include the security monitor (18KLOC), but relies on a significant amount of microcode. As a point of comparison, XEN [12] (a popular mini-hypervisor), an Coreboot [16] (a minimal BIOS) each total 200KLOC.

### 9.3 Microbenchmarks

**TPM Sign Operation.** We evaluate the cost of signing a 32-byte message (i.e., a hash) using a physical TPM. We also compare the incurred overhead with two insecure baselines: 1) software signature in the application and 2) software signature performed by the security monitor (vTPM). Results can be found in table 2. The vTPM signature requires a vmcall and a context switch to the security monitor, which incurs a small overhead (about 100us) compared to the application-level signature. Operations with the physical TPMs are comparatively slow: we measure 196ms to perform a signature (which in our case boils down to loading the hash value in a PCR and requesting a TPM quote). This is 3 orders of magnitude slower than a vTPM. Some modern chipset-integrated hardware TPMs are faster than discrete TPMs and are still microarchitecturally isolated from the CPU. However, they are only 5x faster than discrete TPM at best. These results hint at the importance of Argos' optimized

attestation scheme, which only requires at most one signature per FHE evaluation and even less when using batching.

**Enclave Fixed Costs.** We evaluate the different fixed costs of different operations over the lifetime of an enclave. To obtain measurements, we instrument a dummy enclave for which we allocate 1GB of memory. Setup time is measured between the time a command line program requests the operating system to start an enclave and when execution reaches the main function of the enclave application. This means it includes the allocation of enclave resources, loading of the binary, page table setup, binary measurement, and runtime initialization. We can see that Nitro Enclaves have the highest setup time mostly due to the fact that they contain an entire linux kernel that needs to be booted. SGXv2 also presents some large startup time. This is mostly due to all the special instructions required to start an enclave and the time it takes Gramine to initialize. As indicated in Gramine’s documentation, initialization can initiate more than 200 OCALL which all incur expensive context switches, with architectural flush (CPU caches and state) that cost 8,000 – 12,000 cycles each [96]. This overhead can also be observed when running Argos with Gramine. In comparison, our custom runtime is initialized with minimal overhead, which might be essential for latency-sensitive applications.

Attestation time is the time it takes for an enclave to request an attestation quote. For SGXv2, support for remote attestation service by Intel is deprecated [38] so we use the analogue Microsoft Azure Attestation[95] scheme. Once provisioned, signing of the attestation quote is performed locally by the quoting enclave, hence the high performance similar to using a vTPM. For Nitro Enclaves, the attestation is performed by the hypervisor using a vTPM [8] which also explains the high performance. Due to Argos’s reliance on a physical TPM, we have a comparably slow attestation time.

Termination time is measured between when the enclave main function exits and when the control returns to the command-line program that launched the enclave. It includes all the time required to deallocate, clear resources, and return them to the OS. For SGX, that also implies a lot of cleanup at the hardware level. For Nitro Enclaves, it requires to wait for Linux to shutdown, which also adds delay. In Argos, since all sensitive data is encrypted, memory does not need to be sanitized before being reallocated to the OS.

**Transcript Size.** Our transcript is quite small and composed of 2 PCR values used by the DRoT, the initial value of some key registers for the execution environment, the hash of our binary, the hash of messages exchanged by the application, the hash of some server input if needed. It also contains the signature and the TPM public key along with the TPM certificates endorsed by the manufacturer. Concretely, we measure a small transcript size of only 1,407 bits, negligible in front of the size of an FHE ciphertext ( $\approx 500\text{KB}$ ).

## 9.4 FHE Evaluation

We evaluate the overhead of circuit-level vFHE in Argos compared to semi-honest FHE. Results can be seen in table 4. We use the same benchmarks than previous work on vFHE [114] and compile and evaluate the different circuits using the SEAL library [28]. The results from Viand. et al. highlight the impractical performance

Platform	Tiny	Small	Medium
Baseline [114]	2ms	11ms	14ms
Bulletproofs [22]	7569s	278s	8697s
Aurora [14]	1554s	3750s	5028s
Groth16 [60]	196s	473s	634s
Rinocchio [54]	320ms	305s	443s
SGXv1 [114]	154ms	1100ms	1260ms
Baseline Azure	283us	1727us	3170us
SGXv2 Azure	290us (+3%)	1840us (+7%)	3638us (+15%)
Baseline AWS	324us	1889us	3456us
Nitro Enclave	317us (-2%)	1827us (-3%)	3450us (+0%)
Baseline Local	351us	2341us	4376us
Argos+G	392us (+12%)	2702us (+16%)	5202us (+19%)
Argos	352us (+0%)	2480us (+6%)	4447us (+2%)

**Table 4: Time for FHE evaluation for three circuits. Greyed values are reported from [114]. We indicate the overhead compared to the baseline on the same machine.**

	Per Enclave			Per Query		Per Batch
	Setup	Pre-Proc.	Hash	Proc.	Hash	Attest.
Baseline	3	2835	N/A	1593	N/A	N/A
Argos+G	1559	3192	465	1812	1	196
Argos	158	2821	465	1610	1	196

**Table 5: Breakdown of end-to-end execution times for authenticated private information retrieval (ms).**

#C	Enclave			Query		Batch
	Setup	Pre-Proc.	Hash	Proc.	Hash	Attest.
Baseline	1	4	40351	N/A	3231	N/A
Argos+G	1	1576	44530	132	3732	187
Argos	1	179	37054	105	3354	148
Baseline	3K	4	40348	N/A	3276	N/A
Argos+G	3K	1530	44467	131	3740	187
Argos	3K	175	36933	104	3403	148

**Table 6: Breakdown of end-to-end execution times for private set intersection (ms). #C: size of the client set.**

overheads when using cryptographic proofs to build vFHE: between 6 and 7 orders of magnitude, with the exception of Pinocchio for the Tiny benchmark. Their baseline is also relatively slow, hinting that they might have disable hardware acceleration. Their implementation of FHE-in-TEE also shows poor performance probably due to the limitation of SGXv1 (limited 128MB EPC memory and unreliable timers). For a fair comparison, we run the benchmarks of SGXv2 and Nitro Enclaves. We compiled the benchmarks on one machine and used the same binary across or different setups. Here, the performance looks much better with single-digit overheads and even small speedups for Nitro Enclaves. Indeed, once setup, running a program in an enclave is equivalent to running it with bare-metal performance. We tried to keep the runtime as similar as possible between the different configurations, but small differences might explain these results. For Argos, we evaluate two different configurations. One using the Gramine runtime (Argos+G) and one using our custom runtime. Our custom runtime performs better than Gramine. This might be for a couple reasons including that the

MUSL library tends to be faster than Glibc, our custom runtime simply ignores interrupts sent by the OS scheduler (the enclave is not de-scheduled), and our custom syscall handlers handle less corner cases and are more efficient (for instance, to allocate memory). In conclusion, Argos shows performance comparable to commercial TEE while offering better security. Our custom runtime makes it possible to reduce the TCB and improve performance at the cost of some engineering effort to port applications.

## 9.5 Authenticated Private Information Retrieval

We evaluate the overhead of authenticated PIR using Argos compared to semi-honest unauthenticated PSI. We use the SealPIR library from Angle et. al. [2]. We use the BFV scheme with polynomials of degree 4096. The plaintext modulus is set to 20. We evaluate performance for a database of size  $2^{20} \approx 1M$  where each element measures 128 bytes. We allocate 4GB of static memory to our enclave. We instrument the code to hash the database, queries, and responses and to attest the final transcript once all queries have been processed. Results can be found in table 5. We compare two configurations for Argos. One using the Gramine runtime (Argos+G) and one using our custom runtime. For Argos+G, the server setup –which is a fixed cost– shows an overhead of +83% mostly due to overhead setting up the enclave (+55% alone) and hashing the database (+16%). For query processing, we measure an overhead of 13% for fetching one element. Here, hashing the response before sending it to the client takes a negligible amount of time. Taking into account the signing of the final transcript, the overhead goes up to 26% for retrieving one element. Batching queries here would amortize costs down to a minimum of 13% for numerous queries. For Argos using our custom runtime, overhead goes down to 21% for setup (mostly due to hashing of the batabase) and 1% for query processing. This is due to our more efficient handling of system calls. The only cost left is attestation that can be amortized by batching requests (see section 6.2). Once again, our custom runtime offers better performance than Gramine. For all setups, Argos’s transcript incurs negligible ( $>1\%$ ) communication overhead.

Recent work on single-server authenticated PIR using cryptographic constructions by Colombo. et al., while providing better security guarantees, show a performance overhead of 30-100x [34]. Other recent work, VeriSimplePIR [42], show more reasonable online overheads (13%-20%) but incur significant offline communication in the same order of magnitude than the database itself (2GiB communication for a 4GiB dataset). This also makes updating the database extremely difficult, while our scheme only requires updating a Merkle tree and communicating the commitment.

## 9.6 Authenticated Private Set Intersection

We evaluate the overhead of authenticated PSI using Argos compared to semi-honest non-attested PSI. We use the APSI library compatible with SEAL that implements the state-of-the-art FHE-based scheme from Cong et. al. [35]. We use the BFV scheme with parameters 8192 for the polynomial degree modulus and the coefficient modulus. The plaintext modulus is set to 65537. We evaluated the performance of the unlabeled scheme with a server set size of  $2^{20} \approx 1M$ . We allocate 4GB of static memory to our enclave. We instrument the code to hash the database, queries, and responses and

to attest the final transcript once all queries have been processed. Results can be found in table 6. For Argos+Gramine, the server setup phase, we measure an overhead of 15% mainly due to preprocessing operations, but with our custom runtime, we see a *speedup* of 8% due to our fast handling of memory allocation, and other syscalls. Hashing the database is negligible here. For query time, we see +16% and +8% overhead for Gramine and our custom runtime, respectively. For all setups, Argos incurs negligible communication ( $>1\%$ ). In comparison, the state-of-the-art for malicious PSI using a cryptographic construction [98] shows a 19% performance overhead and a 67% communication overhead compared to their own semi-honest construction for a dataset of  $2^{20}$  elements. However, they cannot guarantee that the server uses the same dataset across clients and require to compute and transfer a large commitment for each client.

**Attested Transcript.** Our remote attestation scheme makes it possible to amortize the cost of attestation over messages and queries. In the PSI scheme we use, each query requires two rounds of communication (one round for OPRF, and one for the query itself), plus one initial attestation, which would bring the overhead for one query to +27%. Instead, we only need to pay the cost of one attestation operation per batch. The PSI scheme we use already supports batching but only up to a certain client set size (this is a limitation of FHE-based PSI schemes). For our given server set size, we empirically measure a non-negligible failure rate for client set size above 3000 elements. A client who would want to compute the intersection of a bigger set would have to do so using several queries. That means that the maximum overhead of 15% for a client set size of 1 can be amortized for larger client set sizes. For instance, a client set size of 30000, computing the complete intersection would take 10 queries and overhead would be reduced back to 8%. This phenomenon would be heightened for more efficient PSI schemes as the cost of TPM signing is constant.

## 9.7 Analysis of the Remaining Attack Surface

Argos primarily defends against microarchitectural side channels, which constitute the vast majority of TEE attacks so far (43 published according to [75]). While Argos is vulnerable to fault injection [68, 85, 97] and Rowhammer attacks (less studied on TEEs), these threats are not fatal – software can be hardened against fault injections [20], and no Rowhammer attacks have demonstrated the capabilities needed to compromise Argos (gaining hypervisor privilege on a hardware VM). Regarding physical attacks, Argos protects against ColdBoot [62] and physical side channels on the CPU [73], but remains vulnerable to physical fault injection attacks [21, 30]. The TPM only exposes a minimal attack surface for physical side channels, and modern TPMs are now directly integrated into the SoC, significantly raising the bar for these attacks [66]. Finally, our TCB might still contain bugs, and attacks that exploit implementation mistakes have been demonstrated on the Intel DRoT and the associated secure coprocessor [36, 48, 118, 119].

## 10 RELATED WORK

**TPMs and Hypervisor-based TEEs.** Trusted platform modules (TPMs) combined with dynamic root-of-trust technologies have long served as dedicated hardware to enforce platform integrity.

TEE Platform	Security						Usability			Performance		
	TCB	$\mu$ -arch SC	Cold Boot	Phys. SC CPU	Fault Inject.	Phys. TPM	Availability	Dedic. HW	Implementation	Setup	Attest.	Comp.
ZK Proofs	Null	P	P	P	P	P	SW	SW	Open Source	---	---	---
Nitro Enclaves	Large	V	V	V	V	V	AWS	Yes	Closed Source	---	++	++
Arm TrustZone	Small	V	V	V	V	V	High	Yes	HW Closed Source	+	+	-
Intel SGX V1	Small	V	P	V	V	V	Deprecated	Yes	HW Closed-Source	--	+	--
Intel SGX V2	Small	V	P	V	V	V	Cloud	Yes	HW Closed-Source	--	+	+
AMD SEV	Large	V	P	V	V	V	Cloud	Yes	HW Closed-Source	--	+	++
Intel TDX	Large	V	P	V	V	V	Coming	Yes	HW Closed-Source	--	+	++
ARM CCA	Large	V	P	V	V	V	Coming	Yes	HW Closed-Source	--	+	++
TrustVisor[78]	Small	V	V	V	V	V	Deprecated	No	Open Source	-	++	++
Flicker [79]	Tiny	P	P	P	V	V	Deprecated	No	Open Source	--	---	--
Argos	Small	P	P	P	V	V	High	No	Open Source	+	-	++

**Table 7: Comparison of TEE Platforms across Security, Usability and Performance.**  
SC: Side Channel, P: Protected, V: Vulnerable

Flicker [79]) proposes to use these two technologies combined to isolate and attest small pieces of application logic, or PALs while reducing the software TCB to a handful of lines of code. No keys ever leave the TPM in Flicker, which gives it the same resistance to microarchitectural side channels as Argos (even if not discussed in the paper at the time). However, their complete reliance on hardware mechanisms creates significant overheads between 2-3 orders of magnitude. Virtualization technologies (making their appearance in processors around 2008) aim at isolating different trusted domains from each other in the form of virtual machines (VMs). They create a new privilege execution mode dedicated to the hypervisor and provide hardware primitives such as extended (or nested) page tables to enforce memory isolation between different security domains. These technologies are quickly identified as an opportunity to virtualize TPM [91] and provide confidentiality and integrity for trusted domains with minimal performance overhead. TrustVisor [78] is the first work to introduce this hypervisor-based architecture which is eventually adopted by cloud providers such as AWS for their Nitro Enclaves [7] or Apple for their private cloud [4]. Virtual TPMs are also the standard for trusted computing in modern clouds [8, 10, 99]. As a result, all these architectures are insecure under our threat model as they all expose sensitive key material to microarchitectural side channels.

**Hardware TEE Platforms.** To the best of our knowledge, all commercial hardware TEE platforms are vulnerable to microarchitectural side channels. This includes, but is not limited to, Intel SGX [39], Intel TDX [32], AND SEV-SNP [104], ARM TrustZone [94] or ARM CCA [76]. Beyond their vulnerability to side channels, commercial platforms also suffer from availability issues as they require dedicated hardware and are not available on consumer-grade processors. For example, Intel announced that SGX would be deprecated on non-server-grade processors [65] and also deprecated its remote attestation service [38]. An exhaustive comparison of existing platforms with Argos can be found in table 7. Academic platform [9, 11, 18, 40, 44, 51, 51, 52, 74, 113, 117] are limited in term

of how much they can customize the hardware to defend against microarchitectural side channels. As a result, the vast majority simply consider side channels out of scope, with the exception of MI6 [17] which is built on top of an open-source RISC-V processor.

**Combining Cryptography and TEEs.** Previous work has identified that protecting remote attestation mechanisms alone against microarchitectural side channels was easier than enforcing confidentiality of entire enclave programs [108]. Some even suggest to combine trusted hardware with cryptography to build hybrid mechanisms, but all show important performance overheads, use (insecure) SGX and are vulnerable to microarchitectural side channels. CryptoFlow [71] looks at secure inference and proposes to build maliciously-secure MPC from semi-honest MPC schemes using SGX. They show a 3x performance overhead over the semi-honest scheme. Chex-Mix [87] leverages SGX with FHE to build secure inference, but considers microarchitectural side channels out of scope and even relies on SGX confidentiality for the privacy of the model weights. They show a 142% performance overhead compared to semi-honest FHE and consider key recovery attacks out of scope. Viand et. al. [114] introduces many useful formalism we get inspiration from in Argos. Their paper primarily focuses on combining FHE with cryptographic proofs but they also consider a FHE-in-TEE approach and design a special protocol to optimize its performance. Nevertheless, they use SGX, do not adapt their formalism to the FHE-in-TEE approach, and their results show an 80x performance overhead when compared to Argos. In contrast, Argos presents a new TEE design for malicious and verifiable FHE with a dedicated remote attestation scheme. Argos is secure *by design* against microarchitectural side channels and shows minimal performance overhead compared to semi-honest baselines.

## 11 CONCLUSION

We present Argos, the first *integrity-only* enclave platform designed to build maliciously-secure verifiable FHE. Argos is secure *by design* against microarchitectural side channels and transient execution attacks and can be used to build fully malicious and authenticated

PSI and PIR schemes. It requires no specialized hardware and is compatible with commodity processors from 2008 onward. Argos only incurs 3% average performance overhead for FHE computation, less than 8% performance overhead for complex protocols, and virtually no communication overhead. By demonstrating how to combine cryptography with trusted hardware, Argos paves the way for widespread deployment of FHE-based protocols beyond the semi-honest setting, without the overhead of cryptographic proofs.

## ACKNOWLEDGMENTS

The authors thank Charly Castes, Adrien Ghosn, Neelu S. Kalani, and the authors of Tyche [23] for providing early access to their code and for their invaluable assistance in setting up Tyche on our hardware and adapting the platform to our needs. We thank Sacha Servan-Schreiber for helpful discussions on fully homomorphic encryption and useful feedback. The authors used Claude to revise the text in most of this paper to correct for typos, grammatical errors, and awkward phrasing.

## REFERENCES

- [1] Onur Aci̇mez, Çetin Kaya Koç, and Jean-Pierre Seifert. 2006. Predicting secret keys via branch prediction. In *Topics in Cryptology—CT-RSA 2007: The Cryptographers’ Track at the RSA Conference 2007, San Francisco, CA, USA, February 5–9, 2007. Proceedings*. Springer, 225–242.
- [2] Sebastian Angel, Hao Chen, Kim Laine, and Srinath Setty. 2018. PIR with compressed queries and amortized query processing. In *2018 IEEE symposium on security and privacy (SP)*. IEEE, 962–979.
- [3] Apple. 2022. Apple Platform Security. <https://help.apple.com/pdf/security/en-US/apple-platform-security-guide.pdf>. Accessed on 28.04.2023.
- [4] Apple. 2024. Private Cloud Compute: A new frontier for AI privacy in the cloud. <https://security.apple.com/blog/private-cloud-compute/>.
- [5] Gilad Asharov, Abhishek Jain, Adriana López-Alt, Eran Tromer, Vinod Vaikuntanathan, and Daniel Wichs. 2012. Multiparty computation with low communication, computation and interaction via threshold FHE. In *Advances in Cryptology—EUROCRYPT 2012: 31st Annual International Conference on the Theory and Applications of Cryptographic Techniques, Cambridge, UK, April 15–19, 2012. Proceedings 31*. Springer, 483–501.
- [6] Shahla Atapoor, Karim Bagheri, Hilder VL Pereira, and Jannik Spiessens. 2024. Verifiable FHE via Lattice-based SNARKs. *Cryptography ePrint Archive* (2024).
- [7] AWS. 2024. AWS Nitro Enclaves Documentation. <https://docs.aws.amazon.com/enclaves/latest/user/nitro-enclave.html>.
- [8] AWS. 2024. AWS Nitro TPM Documentation. <https://aws.amazon.com/blogs/compute/deep-dive-into-nitrotpm-and-uefi-secure-boot-support-in-amazon-ec2/>.
- [9] Ahmed M Azab, Peng Ning, and Xiaolan Zhang. 2011. Sice: a hardware-level strongly isolated computing environment for x86 multi-core platforms. In *Proceedings of the 18th ACM conference on Computer and communications security*. 375–388.
- [10] Azure. 2024. Virtual TPMs in Azure confidential VMs. <https://learn.microsoft.com/en-us/azure/confidential-computing/virtual-tpms-in-azure-confidential-vm>.
- [11] Raad Bahmani, Ferdinand Brasser, Ghada Dessouky, Patrick Jauernig, Matthias Klimmek, Ahmad-Reza Sadeghi, and Emmanuel Stafp. 2021. CURE: A Security Architecture with Customizable and Resilient Enclaves.. In *USENIX Security Symposium*. 1073–1090.
- [12] Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, and Andrew Warfield. 2003. Xen and the art of virtualization. *ACM SIGOPS operating systems review* 37, 5 (2003), 164–177.
- [13] Mihir Bellare, Anand Desai, David Pointcheval, and Phillip Rogaway. 1998. Relations among notions of security for public-key encryption schemes. In *Advances in Cryptology—CRYPTO’98: 18th Annual International Cryptology Conference Santa Barbara, California, USA August 23–27, 1998 Proceedings 18*. Springer, 26–45.
- [14] Eli Ben-Sasson, Alessandro Chiesa, Michael Riabzev, Nicholas Spooner, Madars Virza, and Nicholas P Ward. 2019. Aurora: Transparent succinct arguments for RICS. In *Advances in Cryptology—EUROCRYPT 2019: 38th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Darmstadt, Germany, May 19–23, 2019, Proceedings, Part I 38*. Springer, 103–128.
- [15] Ravi Bhargava, Benjamin Serebrin, Francesco Spadini, and Srilatha Manne. 2008. Accelerating two-dimensional page walks for virtualized systems. In *Proceedings of the 13th international conference on Architectural support for programming languages and operating systems*. 26–35.
- [16] Anton Borisov. 2009. Coreboot at your service! *Linux Journal* 2009, 186 (2009), 1.
- [17] Thomas Bourgeat, Ilia Lebedev, Andrew Wright, Sizhuo Zhang, and Srinivas Devadas. 2019. Mi6: Secure enclaves in a speculative out-of-order processor. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*. 42–56.
- [18] Ferdinand Brasser, David Gens, Patrick Jauernig, Ahmad-Reza Sadeghi, and Emmanuel Stafp. 2019. SANCTUARY: ARMing TrustZone with User-space Enclaves.. In *NDSS*.
- [19] Ferdinand Brasser, Urs Müller, Alexandra Dmitrienko, Kari Kostiainen, Srdjan Capkun, and Ahmad-Reza Sadeghi. 2017. Software grand exposure: {SGX} cache attacks are practical. In *11th USENIX workshop on offensive technologies (WOOT 17)*.
- [20] Jakub Breier and Xiaolu Hou. 2022. How practical are fault injection attacks, really? *IEEE Access* 10 (2022), 113122–113130.
- [21] Robert Buhren, Hans-Niklas Jacob, Thilo Krachenfels, and Jean-Pierre Seifert. 2021. One glitch to rule them all: Fault injection attacks against amd’s secure encrypted virtualization. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. 2875–2889.
- [22] Benedikt Bünz, Jonathan Bootle, Dan Boneh, Andrew Poelstra, Pieter Wuille, and Greg Maxwell. 2018. Bulletproofs: Short proofs for confidential transactions and more. In *2018 IEEE symposium on security and privacy (SP)*. IEEE, 315–334.
- [23] Charly Castes, Adrien Ghosn, Neelu S Kalani, Yuchen Qian, Marios Kogias, Mathias Payer, and Edouard Bugnion. 2023. Creating Trust by Abolishing Hierarchies. In *Proceedings of the 19th Workshop on Hot Topics in Operating Systems*. 231–238.
- [24] Melissa Chase, Hao Chen, Jintai Ding, Shafi Goldwasser, Sergey Gorbunov, Jeffrey Hoffstein, Kristin Lauter, Satya Lokam, Dustin Moody, Travis Morrison, et al. 2017. Security of homomorphic encryption. *HomomorphicEncryption.org, Redmond WA, Tech. Rep* (2017).
- [25] Bhuvnesh Chaturvedi, Anirban Chakraborty, Ayantika Chatterjee, and Debdip Mukhopadhyay. 2022. A practical full key recovery attack on tfhe and fheb by inducing decryption errors. *Cryptography ePrint Archive* (2022).
- [26] Guoxing Chen, Sanchuan Chen, Yuan Xiao, Yinqian Zhang, Zhiqiang Lin, and Ten H Lai. 2019. Sgxpectre: Stealing intel secrets from sgx enclaves via speculative execution. In *2019 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 142–157.
- [27] Hao Chen, Zhicong Huang, Kim Laine, and Peter Rindal. 2018. Labeled PSI from fully homomorphic encryption with malicious security. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. 1223–1237.
- [28] Hao Chen, Kim Laine, and Rachel Player. 2017. Simple encrypted arithmetic library—SEAL v2. 1. In *Financial Cryptography and Data Security: FC 2017 International Workshops, WAHC, BITCOIN, VOTING, WTSC, and TA, Siemaa, Malta, April 7, 2017, Revised Selected Papers 21*. Springer, 3–18.
- [29] Hao Chen, Kim Laine, and Peter Rindal. 2017. Fast private set intersection from homomorphic encryption. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. 1243–1255.
- [30] Zitai Chen, Georgios Vasilakis, Kit Murdock, Edward Dean, David Oswald, and Flavio D Garcia. 2021. {VoltPillager}: Hardware-based fault injection attacks against Intel {SGX} Enclaves using the {SVID} voltage scaling interface. In *30th USENIX Security Symposium (USENIX Security 21)*. 699–716.
- [31] Massimo Chenal and Qiang Tang. 2015. On key recovery attacks against existing somewhat homomorphic encryption schemes. In *Progress in Cryptology—LATINCRYPT 2014: Third International Conference on Cryptology and Information Security in Latin America Florianópolis, Brazil, September 17–19, 2014 Revised Selected Papers 3*. Springer, 239–258.
- [32] Pau-Chen Cheng, Wojciech Ozga, Enriquillo Valdez, Salman Ahmed, Zhongshu Gu, Hani Jamjoom, Hubertus Franke, and James Bottomley. 2024. Intel tdx demystified: A top-down approach. *Comput. Surveys* 56, 9 (2024), 1–33.
- [33] Ilaria Chillotti, Nicolas Gama, Mariya Georgieva, and Malika Izabachène. 2020. TFHE: fast fully homomorphic encryption over the torus. *Journal of Cryptology* 33, 1 (2020), 34–91.
- [34] Simone Colombo, Kirill Nikitin, Henry Corrigan-Gibbs, David J Wu, and Bryan Ford. 2023. Authenticated private information retrieval. In *32nd USENIX security symposium (USENIX Security 23)*. 3835–3851.
- [35] Kelong Cong, Radames Cruz Moreno, Mariana Botelho da Gama, Wei Dai, Ilia Iliashenko, Kim Laine, and Michael Rosenberg. 2021. Labeled PSI from homomorphic encryption with reduced computation and communication. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. 1135–1150.
- [36] Intel Corporation. 2019. Intel CSME, Intel SPS, Intel TXE, Intel ‘DAL, and Intel AMT 2019.1 QSR advisory. <https://www.intel.com/content/www/us/en/security-center/advisory/intel-sa-00213.html>.

- [37] Intel Corporation. 2022. Intel® Converged Security and Management Engine (Intel® CSME) Security. <https://www.intel.com/content/dam/www/public/us/en/security-advisory/documents/intel-csme-security-white-paper.pdf>.
- [38] Intel Corporation. 2024. Intel® Software Guard Extensions Attestation Service Using Intel® Enhanced Privacy Identification End-of-Life Timeline. <https://www.intel.com/content/www/us/en/developer/articles/technical/software-security-guidance/resources/sgx-ias-using-epid-eol-timeline.html>.
- [39] Victor Costan and Srinivas Devadas. 2016. Intel SGX explained. *Cryptology ePrint Archive* (2016).
- [40] Victor Costan, Ilya Lebedev, and Srinivas Devadas. 2016. Sanctum: Minimal hardware extensions for strong software isolation. In *25th USENIX Security Symposium (USENIX Security 16)*. 857–874.
- [41] Morten Dahl, Daniel Demmler, Sarah El Kazdadi, Arthur Meyre, Jean-Baptiste Orfila, Dragos Rotaru, Nigel P Smart, Samuel Tap, and Michael Walter. 2023. Noah’s Ark: Efficient Threshold-FHE Using Noise Flooding. In *Proceedings of the 11th Workshop on Encrypted Computing & Applied Homomorphic Cryptography*. 35–46.
- [42] Leo de Castro and Keewoo Lee. 2024. VeriSimplePIR: verifiability in simplePIR at no online cost for honest servers. *Cryptology ePrint Archive* (2024).
- [43] Emiliano De Cristofaro, Jihye Kim, and Gene Tsudik. 2010. Linear-complexity private set intersection protocols secure in malicious model. In *International Conference on the Theory and Application of Cryptology and Information Security*. Springer, 213–231.
- [44] Liang Deng, Qingkai Zeng, Weiguang Wang, and Yao Liu. 2014. EqualVisor: Providing memory protection in an untrusted commodity hypervisor. In *2014 IEEE 13th International Conference on Trust, Security and Privacy in Computing and Communications*. IEEE, 300–309.
- [45] Marian Dietz and Stefano Tessaro. 2024. Fully malicious authenticated PIR. In *Annual International Cryptology Conference*. Springer, 113–147.
- [46] Léo Ducas and Daniele Micciancio. 2015. FHEW: bootstrapping homomorphic encryption in less than a second. In *Annual international conference on the theory and applications of cryptographic techniques*. Springer, 617–640.
- [47] Alexander Ermolov. 2016. Safeguarding rootkits: Intel BootGuard. *Remote access (01 Sep 2021): https://2016.zeronights.ru/wp-content/uploads/2017/03/Intel-BootGuard.pdf* (2016).
- [48] Mark Ermolov and Maxim Goryachy. 2017. How to hack a turned-off computer, or running unsigned code in intel management engine. *Black Hat Europe* (2017).
- [49] Dmitry Etyushkin, Ryan Riley, Nael CSE Abu-Ghazaleh, ECE, and Dmitry Ponomarev. 2018. Branchscope: A new side-channel attack on directional branch predictor. *ACM SIGPLAN Notices* 53, 2 (2018), 693–707.
- [50] Prastudy Fauzi, Martha Norberg Hovd, and Håvard Raddum. 2022. On the IND-CCA1 security of FHE schemes. *Cryptography* 6, 1 (2022), 13.
- [51] Erhu Feng, Xu Lu, Dong Du, Bicheng Yang, Xueqiang Jiang, Yubin Xia, Bin Yu Zang, and Haibo Chen. 2021. Scalable Memory Protection in the {PENGLAI} Enclave. In *15th {USENIX} Symposium on Operating Systems Design and Implementation ({OSDI} 21)*. 275–294.
- [52] Andrew Ferraiuolo, Andrew Baumann, Chris Hawblitzel, and Bryan Parno. 2017. Komodo: Using verification to disentangle secure-enclave hardware from software. In *Proceedings of the 26th Symposium on Operating Systems Principles*. 287–305.
- [53] William Futral and James Greene. 2013. *Intel Trusted Execution Technology for Server Platforms: A Guide to More Secure Datacenters*. Springer Nature.
- [54] Chaya Ganesh, Anca Nitulescu, and Eduardo Soria-Vazquez. 2023. Rinocchio: SNARKs for ring arithmetic. *Journal of Cryptology* 36, 4 (2023), 41.
- [55] Craig Gentry. 2009. Fully homomorphic encryption using ideal lattices. In *Proceedings of the forty-first annual ACM symposium on Theory of computing*. 169–178.
- [56] Oded Goldreich, Silvio Micali, and Avi Wigderson. 2019. How to play any mental game, or a completeness theorem for protocols with honest majority. In *Providing Sound Foundations for Cryptography: On the Work of Shafi Goldwasser and Silvio Micali*. 307–328.
- [57] Johannes Götzfried, Moritz Eckert, Sebastian Schinzel, and Tilo Müller. 2017. Cache attacks on Intel SGX. In *Proceedings of the 10th European Workshop on Systems Security*. 1–6.
- [58] Ben Gras, Kaveh Razavi, Herbert Bos, Cristiano Giuffrida, et al. 2018. Translation Leak-aside Buffer: Defeating Cache Side-channel Protections with TLB Attacks. In *USENIX Security Symposium*, Vol. 216.
- [59] Matthew Green, Watson Ladd, and Ian Miers. 2016. A protocol for privately reporting ad impressions at scale. In *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*. 1591–1601.
- [60] Jens Groth. 2016. On the size of pairing-based non-interactive arguments. In *Advances in Cryptology—EUROCRYPT 2016: 35th Annual International Conference on the Theory and Applications of Cryptographic Techniques, Vienna, Austria, May 8–12, 2016, Proceedings, Part II* 35. Springer, 305–326.
- [61] Qian Guo, Denis Nabokov, Elias Suvanto, and Thomas Johansson. 2024. Key Recovery Attacks on Approximate Homomorphic Encryption with Non-Worst-Case Noise Flooding Countermeasures. In *Usenix Security*.
- [62] J Alex Halderman, Seth D Schoen, Nadia Heninger, William Clarkson, William Paul, Joseph A Calandrino, Ariel J Feldman, Jacob Appelbaum, and Edward W Felten. 2009. Let’s remember: cold-boot attacks on encryption keys. *Commun. ACM* 52, 5 (2009), 91–98.
- [63] Alexandra Henzinger, Emma Dauterman, Henry Corrigan-Gibbs, and Nickolai Zeldovich. 2023. Private web search with Tiptoe. In *Proceedings of the 29th symposium on operating systems principles*. 396–416.
- [64] Alexandra Henzinger, Matthew M Hong, Henry Corrigan-Gibbs, Sarah Meiklejohn, and Vinod Vaikuntanathan. 2023. One Server for the Price of Two: Simple and Fast {Single-Server} Private Information Retrieval. In *32nd USENIX Security Symposium (USENIX Security 23)*. 3889–3905.
- [65] Intel. 2024. Core Processors. Deprecated Technologies. <https://www.intel.com/content/www/us/en/design/ipla/software-development-platforms/client/platforms/alder-lake-desktop/12th-generation-intel-core-processors-datasheet-volume-1-of-2/004/deprecated-technologies/>.
- [66] Hans Niklas Jacob, Christian Werling, Robert Buhren, and Jean-Pierre Seifert. 2023. faulTPM: Exposing AMD fTPMs’ Deepest Secrets. In *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 1128–1142.
- [67] Scott Johnson, Dominic Rizzo, Partha Sarathy Ranganathan, Jon McCune, and Richard Ho. 2018. Titan: enabling a transparent silicon root of trust for cloud. In *Hot Chips: A Symposium on High Performance Chips*, Vol. 194. 10.
- [68] Zijo Kenjar, Tommaso Frassetto, David Gens, Michael Franz, and Ahmad-Reza Sadeghi. 2020. {VOLTpwn}: Attacking x86 processor integrity from software. In *29th USENIX Security Symposium (USENIX Security 20)*. 1445–1461.
- [69] Ágnes Kiss, Jian Liu, Thomas Schneider, N Asokan, and Benny Pinkas. 2017. Private set intersection for unequal set sizes with mobile applications. In *Privacy Enhancing Technologies Symposium*. De Gruyter, 177–197.
- [70] Gerwin Klein, Kevin Elphinstone, Gernot Heiser, June Andronick, David Cock, Philip Derrin, Dhammika Elkaduve, Kai Engelhardt, Rafal Kolanski, Michael Norrish, et al. 2009. seL4: Formal verification of an OS kernel. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*. 207–220.
- [71] Nishant Kumar, Mayank Rathee, Nishanth Chandran, Divya Gupta, Aseem Rastogi, and Rahul Sharma. 2020. Cryptflow: Secure tensorflow inference. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 336–353.
- [72] Eyal Kushilevitz and Rafail Ostrovsky. 1997. Replication is not needed: Single database, computationally-private information retrieval. In *Proceedings 38th annual symposium on foundations of computer science*. IEEE, 364–373.
- [73] Dayeol Lee, Dongha Jung, Ian T Fang, Chia-Che Tsai, and Raluca Ada Popa. 2020. An {Off-Chip} attack on hardware enclaves via the memory bus. In *29th USENIX Security Symposium (USENIX Security 20)*.
- [74] Dayeol Lee, David Kohlbrenner, Shweta Shinde, Krste Asanović, and Dawn Song. 2020. Keystone: An open framework for architecting trusted execution environments. In *Proceedings of the Fifteenth European Conference on Computer Systems*. 1–16.
- [75] Mengyuan Li, Yuheng Yang, Guoxing Chen, Mengjia Yan, and Yingqian Zhang. 2024. SoK: Understanding Design Choices and Pitfalls of Trusted Execution Environments. In *Proceedings of the 19th ACM Asia Conference on Computer and Communications Security*. 1600–1616.
- [76] Xupeng Li, Xuheng Li, Christoffer Dall, Ronghui Gu, Jason Nieh, Yousuf Sait, and Gareth Stockwell. 2022. Design and verification of the arm confidential compute architecture. In *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. 465–484.
- [77] Yehuda Lindell, David Cook, Tim Geoghegan, Sarah Gran, Rolfe Schmidt, Ehren Kret, Darya Kaviani, and Raluca Ada Popa. 2023. The deployment dilemma: Merits & challenges of deploying MPC. <https://mpc.cs.berkeley.edu/blog/deployment-dilemma.html>.
- [78] Jonathan M McCune, Yanlin Li, Ning Qu, Zongwei Zhou, Anupam Datta, Virgil Gligor, and Adrian Perrig. 2010. TrustVisor: Efficient TCB reduction and attestation. In *2010 IEEE Symposium on Security and Privacy*. IEEE, 143–158.
- [79] Jonathan M McCune, Bryan J Parno, Adrian Perrig, Michael K Reiter, and Hiroshi Isozaki. 2008. Flicker: An execution infrastructure for TCB minimization. In *Proceedings of the 3rd ACM SIGOPS/EuroSys European Conference on Computer Systems 2008*. 315–328.
- [80] Samir Jordan Menon and David J Wu. 2022. Spiral: Fast, high-rate single-server PIR via FHE composition. In *2022 IEEE Symposium on Security and Privacy (SP)*. IEEE, 930–947.
- [81] Samir Jordan Menon and David J Wu. 2024. YPIR: High-Throughput Single-Server PIR with Silent Preprocessing. *Cryptology ePrint Archive* (2024).
- [82] Daniel Moghimi. 2023. Downfall: Exploiting speculative data gathering. In *32nd USENIX Security Symposium (USENIX Security 23)*. 7179–7193.
- [83] Daniel Moghimi, Berk Sunar, Thomas Eisenbarth, and Nadia Heninger. 2020. TPM-FAIL: TPM meets Timing and Lattice Attacks. In *29th USENIX Security Symposium (USENIX Security 20)*. 2057–2073.
- [84] moxie0. 2017. Technology preview: Private contact discovery for Signal. <https://signal.org/blog/private-contact-discovery/>. Accessed on 28.04.2023.
- [85] Kit Murdock, David Oswald, Flavio D Garcia, Jo Van Bulck, Daniel Gruss, and Frank Piessens. 2020. Plundervolt: Software-based fault injection attacks

- against Intel SGX. In *2020 IEEE Symposium on Security and Privacy (SP)*. IEEE, 1466–1482.
- [86] musl. 2024. Musl Libc. <https://musl.libc.org>.
- [87] Deepika Natarajan, Andrew Loveless, Wei Dai, and Ronald Dreslinski. 2021. Chex-mix: Combining homomorphic encryption with trusted execution environments for two-party oblivious inference in the cloud. *Cryptography ePrint Archive* (2021).
- [88] Ofri Nevo, Ni Trieu, and Avishay Yanai. 2021. Simple, fast malicious multiparty private set intersection. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. 1151–1165.
- [89] Cong Nie. 2007. Dynamic root of trust in trusted computing. In *TKK T1105290 Seminar on Network Security*. Citeseer.
- [90] Bijet Pal, Mazharul Islam, Marina Sanusi Bohuk, Nick Sullivan, Luke Valenta, Tara Whalen, Christopher Wood, Thomas Ristenpart, and Rahul Chatterjee. 2022. Might I get pwned?: A second generation compromised credential checking service. In *31st USENIX Security Symposium (USENIX Security 22)*. 1831–1848.
- [91] Ronald Perez, Reiner Sailer, Leendert van Doorn, et al. 2006. vTPM: virtualizing the trusted platform module. In *Proc. 15th Conf. on USENIX Security Symposium*. 305–320.
- [92] James L Peterson and Theodore A Norman. 1977. Buddy systems. *Commun. ACM* 20, 6 (1977), 421–431.
- [93] Benny Pinkas, Mike Rosulek, Ni Trieu, and Avishay Yanai. 2020. PSI from PaXoS: fast, malicious private set intersection. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 739–767.
- [94] Sandro Pinto and Nuno Santos. 2019. Demystifying arm trustzone: A comprehensive survey. *ACM Computing Surveys (CSUR)* 51, 6 (2019), 1–36.
- [95] Gramine Project. 2024. Gramine Documentation. Microsoft Azure Attestation (MAA) Integration. [https://github.com/gramineproject/contrib/tree/master/Integrations/azure/ra\\_tls\\_maa](https://github.com/gramineproject/contrib/tree/master/Integrations/azure/ra_tls_maa).
- [96] Gramine Project. 2024. Gramine Documentation. Performance tuning and analysis. <https://gramine.readthedocs.io/en/stable/performance.html>.
- [97] Pengfei Qiu, Dongsheng Wang, Yongqiang Lyu, and Gang Qu. 2019. VoltJockey: Breaking SGX by software-controlled voltage-induced hardware faults. In *2019 Asian Hardware Oriented Security and Trust Symposium (AsianHOST)*. IEEE, 1–6.
- [98] Srinivasan Raghuraman and Peter Rindal. 2022. Blazing fast PSI from improved OKVS and subfield VOLE. In *Proceedings of the 2022 ACM SIGSAC Conference on Computer and Communications Security*. 2505–2517.
- [99] Himanshu Raj, Stefan Saroiu, Alec Wolman, Ronald Aigner, Jeremiah Cox, Paul England, Chris Fennner, Kinshuman Kinshumann, Jork Loeser, Dennis Mattoon, et al. 2016. F TPM: A Software-Only Implementation of a TPM Chip. In *USENIX Security Symposium*, Vol. 16. 841–856.
- [100] Peter Rindal and Philipp Schoppmann. 2021. VOLE-PSI: fast OPRF and circuit-PSI from vector-OLE. In *Annual International Conference on the Theory and Applications of Cryptographic Techniques*. Springer, 901–930.
- [101] Mike Rosulek and Ni Trieu. 2021. Compact and malicious private set intersection for small sets. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*. 1166–1181.
- [102] Keegan Ryan. 2019. Hardware-backed heist: Extracting ECDSA keys from qualcomm’s trustzone. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*. 181–194.
- [103] Mark D Ryan. 2013. Enhanced certificate transparency and end-to-end encrypted mail. *Cryptography ePrint Archive* (2013).
- [104] AMD SEV-SNP. 2020. Strengthening VM isolation with integrity protection and more. *White Paper, January* (2020).
- [105] Nigel P Smart. 2023. Practical and Efficient FHE-based MPC. In *IMA International Conference on Cryptography and Coding*. Springer, 263–283.
- [106] Yunqing Sun, Jonathan Katz, Mariana Raykova, Philipp Schoppmann, and Xiao Wang. 2024. Actively Secure Private Set Intersection in the Client-Server Setting. *Cryptography ePrint Archive*, Paper 2024/570. <https://eprint.iacr.org/2024/570>.
- [107] Cheng Tan, Lijun Zhang, and Liang Bao. 2020. A Deep Exploration of BitLocker Encryption and Security Analysis. In *2020 IEEE 20th International Conference on Communication Technology (ICCT)*. IEEE, 1070–1074.
- [108] Florian Tramer, Fan Zhang, Huang Lin, Jean-Pierre Hubaux, Ari Juels, and Elaine Shi. 2017. Sealed-glass proofs: Using transparent enclaves to prove and sell knowledge. In *2017 IEEE European Symposium on Security and Privacy (EuroS&P)*. IEEE, 19–34.
- [109] Chia-Che Tsai, Donald E Porter, and Mona Vij. 2017. {Graphene-SGX}: A practical library {OS} for unmodified applications on {SGX}. In *2017 USENIX Annual Technical Conference (USENIX ATC 17)*. 645–658.
- [110] Jo Van Bulck, Marina Minkin, Ofir Weisse, Daniel Genkin, Baris Kasikci, Frank Piessens, Mark Silberstein, Thomas F Wenisch, Yuval Yarom, and Raoul Strackx. 2018. Foreshadow: Extracting the keys to the intel {SGX} kingdom with transient out-of-order execution. In *27th { USENIX } Security Symposium ({ USENIX } Security 18)*. 991–1008.
- [111] Leendert Van Doorn. 2006. Hardware virtualization trends. In *ACM/Usenix International Conference On Virtual Execution Environments: Proceedings of the 2nd international conference on Virtual execution environments*, Vol. 14. 45–45.
- [112] Stephan van Schaik, Andrew Kwong, Daniel Genkin, and Yuval Yarom. 2020. SGAXe: How SGX fails in practice.
- [113] Thomas Van Strydonck, Job Noorman, Jennifer Jackson, Leonardo Alves Dias, Robin Vanderstraeten, David Oswald, Frank Piessens, and Dominique Devriese. 2023. CHERI-TrEE: Flexible enclaves on capability machines. In *2023 IEEE 8th European Symposium on Security and Privacy (EuroS&P)*. IEEE, 1143–1159.
- [114] Alexander Viand, Christian Knabenhaus, and Anwar Hithnawi. 2023. Verifiable fully homomorphic encryption. *arXiv preprint arXiv:2301.07041* (2023).
- [115] Yao Wang, Andrew Ferraiuolo, and G Edward Suh. 2014. Timing channel protection for a shared memory controller. In *2014 IEEE 20th International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 225–236.
- [116] Yingchen Wang, Riccardo Paccagnella, Elizabeth Tang He, Hovav Shacham, Christopher W Fletcher, and David Kohlbrenner. 2022. Hertzbleed: Turning power {Side-Channel} attacks into remote timing attacks on x86. In *31st USENIX Security Symposium (USENIX Security 22)*. 679–697.
- [117] Samuel Weiser, Mario Werner, Ferdinand Brasser, Maja Malenko, Stefan Mangard, and Ahmad-Reza Sadeghi. 2019. Timber-v: Tag-isolated memory bringing fine-grained enclaves to risc-v. In *NDSS*.
- [118] Rafal Wojtczuk and Joanna Rutkowska. 2009. Attacking intel trusted execution technology. *Black Hat DC* 2009 (2009), 1–6.
- [119] Rafal Wojtczuk, Joanna Rutkowska, and Alexander Tereshkin. 2009. Another way to circumvent Intel trusted execution technology. *Invisible Things Lab* (2009), 1–8.
- [120] Tianhong Xu, Aidong Adam Ding, and Yunsi Fei. 2024. TrustZoneTunnel: A Cross-World Pattern History Table-Based Microarchitectural Side-Channel Attack. In *2024 IEEE International Symposium on Hardware Oriented Security and Trust (HOST)*. IEEE, 01–11.
- [121] Vincent Zimmer and Michael Krau. 2016. Establishing the root of trust. *UEFI.org document dated August* (2016).

## A PRIMER ON MEASURED BOOT

Measure boot is obtained using a chain of trust that ties the measurement of the operating system or hypervisor to a component that is inherently trusted by the system (the root). Starting from the root, every element of the chain measures (hashes) and signs the binary of the next element, ensuring its integrity. This process is sometimes referred to as secure boot, as, historically, the chain of trust was composed of the entire boot chain. The root of trust (RoT) is the base case of this recursive security argument, often a hardware mechanism (e.g., read only memory) that enforces integrity for the first instructions executed at boot and which integrity is supported by a certificate signed by the manufacturer [47, 121].

When elements of the boot chain are measured (for instance, before any memory can be trusted), these measurements need to be stored in a secure location where it cannot be tampered with, for instance by a compromised BIOS. To serve that functionality, a small chip or coprocessor called the trusted platform module (TPM) will be in charge of securely storing hash-based measurements in special registers (PCRs). The interface exposed by the TPM is quite limited. For example, measurements can only be hash-extended (never reset). TPMs also have other capabilities, like long-term secrets storage referred to a *secret sealing* (e.g., a key for disk encryption [107]) that can only be accessed if the PCRs are in a specific state, i.e., the correct system has been booted. It also has attestation capabilities, which will use a private key tied to the TPM to sign or *attest* the values of the PCRs and the state of the system. Nevertheless, the TPM is merely a subordinate to the RoT and the software in the chain of trust, as it needs to be instructed and fed inputs to extend its internal measurements.

Any code measured through that process, up to the final domain, is considered part of the trusted code base (TCB). Code in the TCB is critical for the security of the system (it enforces the integrity of the next step in the chain of trust), and as a result, code in the TCB is expected to be bug-free. This assumption is practically impossible

to enforce for a TCB that contains millions of lines of code (e.g., the Linux kernel with 28MLOC), hence the incentive to exclude as much code as possible from the TCB. A dynamic root of trust (DRoT) makes it possible to remove elements of the boot chain from the chain of trust, hence reducing the TCB. DRoT is implemented as a hardware extension (e.g. Intel TXT or AMD DRMT [53, 89]) that can instruct and control the TPM independently of any firmware or software. However, DRoT and TPM were originally conceived to attest to one trusted domain (the OS). That means 1) there is no mechanism to isolate more than one domain at a time 2) TPM hardware resources (PCRs, key storage, bandwidth) are sized for one domain and extremely constrained.