

Guardat: A foundation for policy-protected data

Anjo Vahldiek Eslam Elnikety Aastha Mehta Deepak Garg Peter Druschel Ansley Post (MPI-SWS)
Rodrigo Rodrigues (CITI FCT/UNL) Johannes Gehrke (Microsoft/Cornell)

Technical Report: MPI-SWS-2014-002

January 31, 2014

Abstract

We present Guardat, an architecture that enforces rich data access policies at the storage layer. Users, application developers and system administrators can provide per-file policies to Guardat. Guardat enforces these policies and provides attestations about the state of stored files. With Guardat, the data integrity, confidentiality and access accounting rules for a collection of files can be stated as a single declarative policy. Policy enforcement relies only on the integrity of the Guardat controller and any external policy dependencies; it does not depend on correct software, configuration and operator actions in other parts of a system. Guardat allows developers, system administrators and third-party hosting platform providers to enforce concise, system-wide data protection policies based on a small trusted computing base (TCB), and to demonstrate their compliance to any party that trusts the Guardat layer. We present a design and prototype implementation of Guardat, show experimentally that the space and time overhead of making policy checks is low, and discuss applications and policies.

1 Introduction

As the volume and value of digitally stored assets keep increasing, so do the risks to the integrity and confidentiality of said data. Computer and storage systems are increasing in complexity, exposing data to risks from software bugs, security vulnerabilities and human error. In addition, data is increasingly stored on third-party platforms, introducing additional risks like unauthorized data use by a third party.

In today's systems, data confidentiality and integrity depend on the absence of design errors, bugs, malware and operator mistakes in most components of a system. Moreover, the applicable policies for a collection of data files may be implicit in the code, and their specification and enforcement spread over different subsystems, increasing the risk of misconfigurations. For data stored on third-party platforms, data confidentiality and integrity, as well as proper accounting of data use, additionally depend on the reliability of the third-party provider.

To address these challenges, we present *Guardat*, an architecture that includes a policy interpreter, crypto and enforcement logic at the storage layer. With Guardat, users, developers and administrators can state the in-

tegrity, confidentiality, and accounting rules for a collection of data files using a concise, declarative policy language. Applications communicate with Guardat through secure channels, tunneling through untrusted system layers like storage servers or hosting platforms. Applications send policies, commands and evidence of policy compliance (e.g., proof of authentication) to Guardat and request attestations of stored data and their policies from Guardat. Guardat enforces the policies while relying only on its own interpreter, enforcement logic and any explicit policy dependencies, thus minimizing the size and attack surface of the computing base relied upon for enforcement.

A Guardat policy specifies the conditions under which a file may be read, updated, or have its policy changed. These conditions, written in a declarative language, may depend on client authentication, the initial and final states of the file (size and content) in an update transaction, or signed statements by external trusted components (certifying, for instance, the current wall-clock time). Guardat stores the policy as part of its own metadata and ensures that each access to the file complies with the policy.

Following are some example Guardat policies that mitigate important threats: System binaries can be protected from viruses through a policy that permits updates only when signed by a trusted party; system log corruption and tampering can be avoided through a Guardat-enforced append-only policy; accidental deletion or corruption of backup data can be prevented by a policy that prevents modification for a specific period of time; confidentiality of a user's private data can be enforced by allowing reads only in a session authenticated by the user's public key; and, accesses to a data file may require that a corresponding record is added to an append-only log file, enforcing mandatory access logging.

While these policies can be implemented in higher software layers, with Guardat the policy applicable to a collection of files can be specified using a concise, declarative language, and enforced by a small trusted computing base (TCB) with a small attack surface. Guardat complements existing techniques for ensuring the reliability of data processing systems, including software testing, verification, security auditing, sealed data and trusted computing. While no technique can provide comprehensive protection, Guardat provides a safety net that protects a system's persistent data from a wide range of threats. Moreover, Guardat can demonstrate compliance

with client and provider policies, as well as legal requirements, to any party that trusts Guardat.

The contributions of this work include (1) the Guardat architecture and API; (2) a declarative policy language, which balances expressiveness and efficient evaluation; and, (3) an experimental evaluation of a Guardat prototype, which shows that Guardat policies can be enforced with low overhead. Section 2 describes the design of Guardat and its policy language. Section 3 presents example policies and the guarantees they provide. In Section 4, we present results from experiments with a prototype implementation in the iSCSI IET SAN server. We cover related work in Section 5 and conclude in Section 6.

2 Guardat design

The Guardat design is guided by three principles. First, Guardat enforces policies entirely in the *storage layer* to minimize the TCB and its attack surface. Second, we keep policy specifications concise and separate from code by expressing policies in a domain-specific *declarative policy language*. Third, in the interest of a small TCB, the Guardat policy language provides only a minimal set of primitives sufficient to *check* a rich set of policies, but we rely on untrusted code to specify *how to satisfy* a policy. We point out instances of such design economy throughout this section, as we describe the Guardat API and policy language.

Design overview Guardat’s program logic, called the Guardat controller or GDC, executes just above or inside the storage layer and enforces the file’s policy on every read and write access to data. The GDC exports an *extended block-device API* that allows users, applications and system administrators to (a) create, delete, read and update files, (b) cryptographically authenticate and establish secure sessions to tunnel commands and data through other untrusted software and hardware, (c) associate policies with files, (d) provide credentials and other information to satisfy these policies during subsequent access, and (e) obtain Guardat attestation on stored files and their policies.

Data stored in Guardat is organized into files. For each policy-protected file, Guardat maintains its own metadata, consisting of an ordered list of extents, a unique numeric identifier, a textual name string (typically used to store the file’s pathname), and a reference to a policy in effect for the file. The set of numeric identifiers form a flat namespace, while the set of names can encode a hierarchy maintained by an untrusted filesystem. Each file can have its own policy, but typically, a collection of files share the same policy.

The policy of a file consists of four rules, one for each

of the permissions **read**, **update**, **destroy** and **setpolicy**. Each rule specifies conditions on the context and environment under which the respective permission holds. Abstractly, the **read** rule represents the file’s confidentiality policy; the **update** rule encodes the file’s integrity policy; the **destroy** rule governs when the file’s identifier (name) can be recycled; and the **setpolicy** rule describes when the policy can be changed. API calls that read or update a file or its metadata check conditions of the corresponding policy rules.

Besides a device for storing data, the GDC requires a small amount of fast, persistent memory like Flash for storing policies and other metadata. Flash memory is widely available now; hybrid disks even combine a HDD and Flash in a single enclosure [41]. To authenticate itself as a legitimate Guardat device, sign attestations and encrypt data, the GDC includes a manufacturer-provided unique key pair and certificate.

Implementation, threat model and scope Guardat can be implemented in different ways depending on the deployment. For example, the GDC can be implemented (a) in a SAN server for use in a data center, or (b) integrated with the microcontroller of a hybrid disk for use in an individual machine. A possible third implementation of the GDC is a (c) trustlet within a virtual machine monitor or operating system, isolated using trusted hardware features like Intel SGX [23] or ARM TrustZone [7].

In each implementation, the GDC, metadata and data must be protected from unauthorized access and undetected tampering. In implementation (a), which is the basis of our Guardat prototype described in Section 4, the SAN server includes the GDC, data and metadata storage devices, and must be physically protected, e.g., in a machine room where access is restricted to trusted staff. In this scenario, the Guardat policies are enforced despite any bugs, misconfigurations, or security incidents outside the SAN server, and regardless of actions by employees without access to the machine room.

In implementation (b), the GDC is implemented as part of a microcontroller embedded in a hybrid disk. Here, the metadata and data are encrypted and authenticated to protect them from unauthorized access and undetected tampering. The microcontroller implements the GDC and stores its private key in an embedded TPM. In this scenario, the Guardat policies are enforced as long as the microcontroller has not been physically tampered with. While we have not attempted this implementation, we believe it is feasible with a high-end microcontroller that has on-chip hardware support for secure hashing and cryptography, as well as a TPM. Implementation (c) has similar security properties, except that the GDC executes on the main CPU and trust is derived from this CPU’s trusted isolation capabilities.

	Deployment objective	Guardat implementation	Trust assumption	Who trusts?	How trust is discharged
A.	User wishes to protect her data from bugs, errors and opportunistic employees at a reputable Cloud provider	Cloud storage servers	Only trusted staff has physical access to servers	User	Provider restricts physical access to servers
B.	Data center wishes protection from bugs, misconfigurations, disgruntled employees	Storage servers	ditto	Data center	Center restricts physical access to servers
C.	Service provider wishes to protect proprietary content cached on user's machine	Microcontroller in user's disk	User cannot compromise the controller	Provider	User is unable to tamper with controller chip
D.	User wants to protect data on her machine from bugs, viruses and mistakes	Microcontroller in user's disk	None needed	–	–

Table 1: Guardat deployment scenarios and trust assumptions

Table 1 lists examples of deployment scenarios, their threat models and trust assumptions. For instance, if a reputable Cloud provider deploys Guardat-enabled SAN servers to protect user data from bugs and misconfigurations in its infrastructure and from opportunistic access by employees, then the user must trust the Cloud provider to prevent physical access to the SAN by all but trusted employees (line A in Table 1). Similarly, if a digital content provider locally caches copyrighted content in a Guardat-enabled disk on a user's machine, then it must trust that the user is unable to tamper with the controller chip, much in same way trusted computing applications trust that users are unable to tamper with a TPM.

We also make standard assumptions about policies: correct policies must be installed when data is first stored, and external dependencies of policies (e.g., time servers, client's authentication keys) must be trustworthy. Under these assumptions, Guardat defends against threats to confidentiality and integrity of stored data or to data in transit to storage. This includes threats due to bugs and vulnerabilities in intermediate software layers including operating systems, filesystems, storage services built on top of Guardat, and networks, and threats due to human negligence and opportunistic malice. Guardat is not concerned with data availability. To mask the effects of a hardware or media failure, loss, or destruction of a Guardat device, data must be replicated on multiple devices with independent failure modes.

2.1 Guardat API

The Guardat API extends the standard block device API, which is maintained for backward compatibility, with means to establish sessions, provide policies, provide information to satisfy policies, and obtain attestations.

Session API: A user or an application (generically called *client*) interacts with Guardat in a *session*. (To access files whose policies do not require authentication, it is possible to communicate with Guardat outside of a ses-

sion. Such communication is conceptually treated as part of a default, untrusted session.)

A session is established with a handshake protocol in which the client and Guardat authenticate each other using their private keys. As part of the protocol, new, session-specific keys are created. These keys are used to encrypt and/or authenticate (through message authentication codes) all subsequent communication in the session. This protects in-transit data and commands from snooping and modification in intermediate layers. Moreover, the public key of the client (which acts as an identifier for the client) becomes available during every policy evaluation in the session; hence, Guardat can enforce policies that restrict access to specific users. At the end of the handshake, Guardat returns a unique session identifier that links later API calls to the session. The handshake protocol and its API are described in the appendix (Section 7).

Batch API:

The call `openBatch(sessionId, objName)` starts a new batch on the file named `objName`. Generating the `objName` is up to the (untrusted) higher layers, e.g., the filesystem. If `objName` does not exist, a new empty file is created and given this name (this is the only way to create a file in Guardat). The call returns a batch id that links later calls to the batch and the session. Subsequently, the call `readBatched(batchId, off, len, buf)` reads `len` bytes of the file starting at logical offset `off` in the file and returns the result to the buffer `buf`. The **read** rule of the file's policy is evaluated before writing to `buf`; if it denies access, the call fails. This enforces data confidentiality. Note that we allow byte-level addressing on files.

An file is updated by reusing content from its *current* version and adding fresh content to create a *new* version. The call `reuse(batchId, off, len, off')` takes content in the logical range `[off, off+len-1]` from the current version and inserts it at offset `off'` in the new version (insertion is purely a metadata operation). The call `fresh(batchId, off, len, buf, off')` writes `len` bytes from buffer `buf` to the

extent starting at byte number `off` on disk and adds the resulting extent to the new version at logical offset `off'`. Before writing the extent, Guardat checks that it is not occupied by any file (including the file being modified). The new version of the file may be given a new policy with the call `setPolicy(batchId, new_policy)`.

The updates in a batch are committed with the call `endBatch(batchId)`. Guardat evaluates the **update** rule of the file's policy before committing the new version. This enforces data integrity. The **update** rule has access to the current and new content of the file, as well as relevant metadata, e.g., the offsets and lengths of reads and writes in the batch. Additionally, if the policy has been updated, Guardat evaluates the **setpolicy** rule of the file's policy; this protects the policy from unauthorized changes.

The call `destroy(objName)` deletes `objName` and all its metadata. The file's **destroy** policy authorizes the call. The call also requires that `objName` be *empty*. This design, following our goal of concise policy representation, ensures that the integrity policy of an file is represented entirely in the **update** permission; **destroy** only controls removal of the file's name from Guardat metadata.

Content caches: Two Guardat caches buffer file content for use in policy evaluation. There is a per-session cache of two types of records: `(obj,off,len,content)`, where `content` is the sequence of `len` bytes stored starting at offset `off` within the file `obj`; and `(obj,((off1,len1),..., (offn,lenn)),hash)`, where `hash` is the 32-byte SHA-256 hash over the bytes stored at the specified ordered list of `(off,len)` extents within the file. The session cache reflects the current committed content of the referenced files. In addition, there is a separate per-batch cache of records `(off,len,content)` and `((off1,len1),..., (offn,lenn)),hash` which represent the new uncommitted content of the file manipulated within the batch. Records are inserted into the session cache as a side-effect of the `readBatched()` call, while records are inserted into the batch cache as a side-effect of the `fresh()` call. Flags to these calls indicate whether content, hashes or neither should be inserted (these flags are described in Appendix 7). When a batch commits, any records in the batch cache are moved into the session cache, and any existing session cache records they supersede are evicted. When a batch aborts, the records in the batch cache are discarded.

All relevant file content must exist in the caches before policy evaluation, else access is denied. We rely on the untrusted client for this: The client must set appropriate cache flags in `readBatched()` and `fresh()` calls. This is in line with our principle of economy in design: denying access for lack of content in the cache is *safe*, whereas searching for that information on disk is inefficient and can easily lead to DoS attacks. Policy-aware wrapper libraries could alleviate the need for adding cache flags

in every application's code.

Certificate API: The call `setCertificate(certificate)` forwards a third-party certificate to Guardat for use in subsequent policy evaluations, whereas `getNonce()` returns a fresh nonce, which can be embedded in a subsequent certificate. Third-party certificates are described further in Section 2.2.

The call `attest(objName, nonce)` returns a Guardat-signed certificate that attests the existence of `objName`, its extents and its policy. Optionally, the certificate also includes a hash of any of the file's data in the session cache. The attestation certificate embeds a client-provided nonce, which is useful for preventing replay attacks in protocols built over Guardat. The **read** policy rule authorizes this call.

Pickle API: The pickle APIs allows untrusted platform software to securely manage the replication and migration of policy-protected files among Guardat devices, without access to their cleartext contents. A file copy succeeds only if the file's policy allows it, and if the integrity of the file's contents, name and policy are maintained during the transfer. The pickle operation invoked at a source device encrypts a file and its policy for a specific target Guardat device, while the unpickle operation installs the file at the target device. An attestation from the target device can then be used to prove to the source device that the file resides on the target device. A file's policy controls if, when and where a file can be migrated or replicated.

Backwards compatibility: For compatibility with existing systems, Guardat supports the standard block-device API calls `read(blk,cnt,buf)` and `write(blk,cnt,buf)`. The `read()` call reads `cnt` blocks sequentially from disk starting at block `blk` and returns the data in buffer `buf`. The `write()` call is dual. In executing these calls, Guardat uses its metadata to find all files that intersect the extent being read or written. It evaluates the respective **read** or **update** policy rule of all these files, and fails with an error if any evaluation denies access. Disk blocks not associated with any file can be accessed without restriction through the `read()` and `write()` calls. Hence, Guardat may be configured to selectively protect only a part of a storage disk. Also, Guardat can interoperate with existing, unmodified file systems using an application library. More details can be found in the appendix (Section 7).

2.2 Guardat policy language

Guardat file access policies are specified in an expressive and simple declarative language. Each file's policy contains four *rules*, one for each of the permissions **read**, **update**, **destroy** and **setpolicy**. Each rule specifies the *conditions* under which the respective permission holds.

A rule has the form `(perm :- conds)` and means that

permission “perm” is granted if the conditions “conds” are satisfied. The conditions “conds” consist of *atomic facts* connected with conjunction (“and”, written \wedge) and disjunction (“or”, written \vee). Operationally, policy rules are clauses of constrained Datalog, with all atomic facts in conditions treated as external [25]. Datalog is a standard foundation for writing access policies [8, 13, 36], known for its clarity, high-level of abstraction and ease of implementation.

Each atomic fact contains a predicate that relates file names, content, public keys, extent lists, etc. to each other. The expressiveness of the Guardat policy language stems from the wide range of available predicates. *Universal predicates* are available in all policy rules. The predicate `session_is(K)` checks that the ongoing session is authenticated with the public key K and `file_name_is(O)` means that the file being accessed has name O . The predicate `(obj, off, len) says R` provides access to file data. It means that a record in the session data cache states that file `obj` has content R at offset `off`. Similarly, `(obj, ((off1, len1), ... (offn, lenn))) hasHash H` provides access to hashes of file data. Through these predicates, a file’s policy may test the content of another file. We find this useful in representing many policies, including mandatory access logging (Section 3).

Additionally, *contextual predicates* provide information specific to a policy rule. In **read**, this information includes the length of the read and its logical and physical offsets. As a result of such fine-grained information, confidentiality policies may be specified at the granularity of bytes. In **update**, contextual predicates provide information about the current and new extents of the file, the current and new file sizes, and access to the data cache through the predicates `(off, len) willsay R` (the new content at offset `off` will be R) and `((off1, len1), ... (offn, lenn)) willHaveHash H` (the new bytes stored in the list of `(off, len)` pairs will have hash H). This facilitates rich integrity policies that correlate old and new file content as well as content across two different parts of an file. Again, we find this handy for many policies, including mandatory access logging. All contextual predicates are described in the appendix (Section 8).

Finally, Guardat policies may contain arbitrary *uninterpreted predicates* that are established through signed third-party certificates. These include time-server certificates to establish wall clock time. When a third-party certificate is provided to Guardat through the `setCertificate()` API call, Guardat checks its signature using standard certificate chain verification [12] and stores its content and its signer’s public key in a certificate cache. This cache is available during policy evaluation, through two types of uninterpreted predicates:

- Public key binding, `key_is(k, a)`, which states that

public key k has attribute a . For example, a may be “TimeServer”, suggesting that k is a time server’s public key. The corresponding certificate must be signed by a certifying authority (CA) or its delegatee.

- Signed relation, k signs $r(t_1, \dots, t_n)$ at t : There is a certificate verified by public key k and received at time t , which contains the relation $r(t_1, \dots, t_n)$.

The policy designer and certificate issuers must agree on the meaning of the attribute a in the first point and of the relation r in second point. Guardat treats both a and r as bitstrings. Section 3 illustrates this further. To prevent certificate replay attacks, each certificate must include a Guardat-generated nonce, obtained through the API call `getNonce()`.

To enforce time-sensitive policies, Guardat relies on time-server certificates and an internal *timing counter*. When a time-server certificate is received, its cache entry is stamped with the value of the timing counter. Later, clock time can be estimated by adding the difference of the then-value of the timing counter and the value of this stamp to the time mentioned in the time-server certificate. This timing counter need not be very precise because it can be reset periodically to prevent a large drift. Whenever the timing counter is reset, all time-server certificates must be evicted from the cache.

Following our principle of economy in design, Guardat does not include logic to contact third-parties to obtain relevant policy certificates. Instead, populating the cache with relevant certificates before access is the responsibility of the untrusted Guardat client. If required certificates are missing, access is denied. (When a certificate issuer is offline, access to files that rely on certificates from that issuer may be denied, but access to other files remains unaffected.)

3 Policy examples

We illustrate Guardat’s capabilities by presenting several example policies. For brevity, we introduce the following convention to omit default policy rules: If the rule for the **read** or **update** permissions is omitted, then the permission is always allowed and if the rule for the **setpolicy** or **destroy** permission is omitted, then that permission is never allowed.

Protected executables For a binary file, it is desirable to defend against accidental or malicious overwriting or rollback to a prior version. A representative Guardat policy to accomplish this is shown below. The policy states that the new content of the binary after any update must be signed by the software vendor (called “Vendor”) as being version 10 or later. Moreover, any changes to the

policy must be certified with the administrator’s key, k_{ad} .

```

update :- file_name_is( $O$ )  $\wedge$  new_length_is( $L$ )  $\wedge$ 
  ( $O, L$ ) willHaveHash  $Nh$   $\wedge$  key_is( $K$ , “Vendor”)  $\wedge$ 
   $K$  signs ok_hash( $O, N, Nh$ )  $\wedge$  ( $N \geq 10$ )
setpolicy :- file_name_is( $O$ )  $\wedge$ 
  new_pol_hash_is( $Nph$ )  $\wedge$ 
   $k_{\text{ad}}$  signs good_policy( $O, Nph$ )

```

The first rule allows an update to the file only if there is a public key K belonging to “Vendor” (condition $\text{key_is}(K, \text{“Vendor”})$), which signs that the file’s new content hash, Nh , is the N th version of the binary (condition $K \text{ signs ok_hash}(O, N, Nh)$) and $N \geq 10$. The uninterpreted predicates $\text{key_is}(K, \text{“Vendor”})$ and $K \text{ signs ok_hash}(O, N, Nhash)$ are verified from client-provided certificates signed by a certifying authority and the vendor, respectively. Because the vendor’s certificate contains a single hash over the entire file content, a batch update is needed to satisfy this policy.

The second rule allows a change to the binary’s policy only if the hash of the new policy, called Nph , has been certified by the administrator (condition $k_{\text{ad}} \text{ signs good_policy}(O, Nph)$).

Properties: As long as the integrity of the vendor’s and admin’s keys is maintained, files protected by the policy cannot be overwritten except with content signed by the vendor and version ≥ 10 , even if the entire system is compromised (write integrity). Moreover, a client operating system can ensure it executes only trusted executables despite a compromised storage service (read integrity) as follows: before executing a binary, it obtains an attestation certificate for the file from Guardat, verifies the policy and file name (full path name) in the certificate, and compares the hash of the data delivered by the storage service with that manifest in the certificate.

Append-only logs The following policy specifies an append-only file that may be extended by anyone but modified in-place (e.g., rotated) only by an administrator identified by the public key k_{ad} . The policy would prevent accidental or malicious record deletion from system log files.

```

update :- session_is( $k_{\text{ad}}$ )  $\vee$ 
  (old_length_is( $Lo$ )  $\wedge$  new_length_is( $Ln$ )  $\wedge$  ( $Ln \geq Lo$ )  $\wedge$ 
  updated_locations_are( $M$ )  $\wedge$  disjoint( $M, [0, Lo]$ ))

```

The policy allows an update if either the session is authenticated by the administrator (condition $\text{session_is}(k_{\text{ad}})$) or the file’s new length Ln exceeds its current length Lo and the first Lo bytes of the file are not modified.

Properties: As long as the integrity of the admin’s key is maintained, append-only writes are ensured for any

file with the policy, even if the entire remaining system is compromised.

Protected backup Backup files can be protected from accidental or malicious modification for a fixed period of time using the following policy.

```

update :- key_is( $K$ , “TimeServer”)  $\wedge$ 
   $K$  signs time( $T$ ) at  $T_i$   $\wedge$ 
  time_is( $T_j$ )  $\wedge$  ( $T + T_j - T_i > \text{endT}$ )

```

The policy allows modification to the file only if the current time exceeds a pre-determined time endT . To enforce such policies, Guardat relies on periodic certificates from time servers and a short-range internal timing counter. In detail, the policy says that there should be a key K belonging to a time server (condition $\text{key_is}(K, \text{“TimeServer”})$), which issued a certificate that the time was T when the Guardat internal clock had value T_i (condition $K \text{ signs time}(T) \text{ at } T_i$), the current internal clock value is T_j (condition $\text{time_is}(T_j)$) and the current time (calculated as $T + T_j - T_i$) exceeds the backup end time endT .

Properties: As long as the integrity of the time server and its signing key is maintained, a file with this policy cannot be modified before the designated time, even if the system, the admin’s and the file owner’s private keys are compromised.

Mandatory access logging Legislation and organizational policies often mandate that all read and write access to sensitive information like medical records be logged to a separate file. Although application-level solutions to enforce such mandatory access logging (MAL) exist, enforcing the policy in Guardat is desirable because it would result in a smaller trusted computing base.

For this exposition, let P be the sensitive file which must be protected by MAL and let L be its log file. We assume that the log file is append-only, through the policy described earlier. The MAL requirement is three-fold: 1) (Completeness) For every read on P , an entry in L should describe *who* read and from *where* in P . For every write, a similar entry must exist in L and it must additionally contain a hash of the content written. 2) (Causality) Given two write entries in L , the order in which they were applied to P should be evident and, similarly for a read and a write entry. 3) (Precision) Call a write entry in L *dangling* if it does not correspond to an actual write on P . Then, either dangling entries should not be allowed in L or they should be detectable.¹

¹Dangling read entries are usually not a problem, because it is in the client’s interest to establish that it did *not* read certain data and, hence, not create dangling read entries. We also describe later how read entries can be made precise.

We start with an obvious strawman policy for P , which is complete, but does not provide causality and precision. We refine the design later. We define two kinds of entries for L : $\text{may_read}(K, S)$, which indicates that the client with public key K has potentially read the set S of (off, len) ranges from P ; and $\text{change}(K, S, H)$, which states that content with hash H has been written to the ranges in S . To force logging of reads, we require in the **read** rule of P 's policy that if the range R is read by client K , then an entry $\text{may_read}(K, S)$ with $R \subseteq S$ exist in L . Similarly, write logging could be forced through P 's **update** rule.

This strawman policy for P can be expressed in the Guardat policy language because the set R of locations read or updated is available through contextual predicates in the policy language, the client K is available through the predicate $\text{is_session}(K)$ and L 's content is available through the session cache (predicate says). The policy can also be easily satisfied by the client: Prior to reading or writing, the client could append an appropriate entry to L and have it cached for P 's subsequent policy evaluation. Even though this policy satisfies the MAL requirement of completeness, it does not satisfy causality and precision. Nothing in L 's policy prevents the client from creating entries that are never used and such entries cannot be distinguished from others (this violates precision). Moreover, nothing in P 's policy prevents use of L 's entries out-of-order, which violates causality.

To obtain causality and precision, we refine this strawman design. We embed a counter in each entry in L and enforce through L 's policy that the counter increase by 1 at each successive change entry and remain the same at each may_read entry. We enforce through P 's policy that the value of the counter in the last change entry that has already been applied to P be written at a designated locus in P . Further, the entry used to justify a read must have a counter number that matches the current counter in P . We describe below how we enforce these requirements. Assuming that they have been enforced, both causality and precision are satisfied. Causality holds because the policies just described force that change entries apply to P in increasing order of their counter numbers, and that a read corresponding to a may_read is used after all change entries with smaller or equal counter numbers have been applied. Precision holds because a change entry is dangling if and only if its counter number is higher than the counter in P .

The log's entries are revised to include counter numbers. They take the forms $\text{may_read}(N, K, S)$ and $\text{change}(N, K, S, H)$, where N denotes a counter. We reserve a fixed locus in P for a counter, called C . The log is initialized with a dummy entry with $N = 0$ and P is initialized with $C = 0$. We describe relevant policies of L and P in words, omitting symbolic representations for

clarity. (We have formally represented these policies in our prototype implementation; experimental results are presented in Section 4.)

L 's update policy: Only appends are allowed and only entries of the two designated forms may be added. If the added entry has the form $\text{may_read}(N, \dots)$, then N must be copied from the previous entry and if the added entry has the form $\text{change}(N, \dots)$, then N must be one more than the previous entry's counter. (These requirements can be represented in the Guardat policy language because the previous entry and the new entry are accessible through the session and batch caches, respectively, during evaluation of the **update** rule.)

P 's read policy: L must contain a may_read entry with the same counter number as C and range set larger than the actual range read. (L 's relevant entry and C are accessible through the session cache during P 's policy evaluation. In particular, C can be referenced because Guardat supports byte-level addressing on files and the locus of C is fixed in advance. The client is responsible for specifying which entry of L in the session cache satisfies the policy.)

P 's update policy: L must contain an entry describing the update precisely. The counter in the entry must be one more than C . The update must also increment C by 1. (When evaluating P 's policy, L 's relevant entry and the old value of C are accessible through the session cache. The new value of C is accessible through the batch cache.)

Properties: The policies enforce all MAL requirements as long as the Guardat controller is not compromised.

3.1 Expressiveness and overhead

The example policies described in this section illustrate the expressiveness of Guardat's policy language and the effort required to write policies. To the best of our knowledge, with the exception of [15, 48], no prior work can express any of these policies. The append-only log policy with administrator rotation and the MAL policy cannot be expressed in any prior work.

Similar to Datalog, the worst case complexity for evaluating a Guardat policy is $m(D^n)$ where m is the size of the policy, D is the size of the domain of terms (bounded by the size of Guardat's cache) and n is the number of variables in the policy. As we will show in Section 4, however, policy evaluation has low overhead for reasonable policies. To prevent DoS attacks, Guardat denies access if a policy evaluation exceeds a certain time limit.

Policies can use external verifiers to "outsource" decisions that cannot be expressed in the policy language or exceed the computational resources of the Guardat policy interpreter. However, external verifiers come at the cost of adding external dependencies to the TCB. An ex-

ample of a policy that can't be expressed in our language (without an external verifier) is an integrity policy that requires that the content of a file is well-formed XML.

4 Experimental evaluation

In this section, we present results of an experimental evaluation of a Guardat prototype.

4.1 Prototype Implementation

Our prototype is a modified iSCSI Enterprise Target (IET) SAN server. IET implements the server-side iSCSI protocol, which provides SCSI block storage access via Ethernet. IET is in production use and available for many Linux distributions, e.g., SUSE, RHEL and Debian.

IET consists of a kernel module, which implements block accesses, and a user-level daemon process, which implements iSCSI management functions. To implement Guardat, we extended the kernel module and added a second user-level daemon process, which implements the metadata structures, Guardat API and policy evaluation. The kernel module performs upcalls to determine if iSCSI block accesses should be allowed. The server has access to an SSD for storing Guardat metadata and one or more magnetic disks for payload data.

The Guardat daemon maintains two B-tree index structures: a block-to-file index to find the file and policy associated with a given disk location (block id), and a name-to-file index to retrieve the file information (set of extents, policy, etc.) given a file name. For performance, the Guardat daemon maintains a write-through DRAM cache of B-tree nodes and policies, backed by the SSD. Updates to these data structures are persisted on the SSD during a batch commit.

When the kernel module receives a disk access request, it passes the access type (read/write) and location (disk offset, length) to the multi-threaded Guardat daemon, which consults the block-to-file index. If the disk location is not associated with any file, the access is granted. Otherwise, the daemon retrieves and evaluates relevant policies, and returns the result to the kernel module. For read requests, the disk read is performed while checking the permission. This may result in some wasted work if the read is denied, but results in lower latency if it is not. During a write request, the disk write must be deferred until the daemon grants the permission.

Our prototype's attack surface consists of the IET management interface, the block-device interface, the Guardat interface extensions as well as the policy language. Despite the relatively large IET codebase, which includes a minimally configured Linux kernel, the resulting attack surface is likely to be significantly smaller than that of the systems and applications built on top of Guar-

dat in most cases. Our Guardat implementation adds less than 20,000 LOC to the existing IET codebase, plus the OpenSSL and glib libraries it relies on.

4.2 Experimental setup

In our setup, the Guardat enhanced IET SAN server (based on version 1.4.20.3-9.6.1) [47] runs on a separate physical server connected to the client via 10Gbit Ethernet links. The client software runs on OpenSuse Linux 12.1 (kernel version 3.1.10-1.16, x86-64). The Linux iSCSI client connects to the IET server, and appears to the Linux filesystems as a locally connected SCSI block device.

The IET server and the Linux client each run on a Dell Precision T1600 workstation with an Intel Xeon 3.1Ghz quad core CPU (AES and AVX instruction set) and 8GB main memory. The server has a 500GB disk drive with the server OS installation, and two disks that are used for Guardat. Data is stored on a Seagate Barracuda 2TB 7200 rpm hard drive with a 64MB cache [40], and the Guardat metadata is stored on a OCZ Deneva 2 C SLC 60GB (raw 64GB) SSD [30]. Only 4GB of the SSD is actually used for Guardat metadata; the remaining capacity is available for general use by clients. Reducing the size of the SSD available for general use by 4GB is a small cost for the added security.

The openssl crypto library [31], Intel AES encryption library [21], and Intel's fast SHA256 implementation [22] are used for Guardat cryptographic operations.

4.3 Microbenchmarks

We performed a series of microbenchmarks to quantify the overheads incurred by the Guardat prototype in terms of storage space, read/write latency and throughput, policy-evaluation latency and Flash memory wear.

4.3.1 Space requirements for metadata

First, we quantify the metadata storage requirements. Because the metadata size depends on the structure of the payload data, we analyzed the metadata space requirements for 70,825 filesystem snapshots collected by Agrawal et al. [1]. The snapshots were taken from Windows systems within Microsoft corporation between 2000 and 2004, and contain between 30k and 90k files each with an average file size between 108KB and 189KB. For evaluation purposes, we give each file in each snapshot an integrity policy that disallows modification prior to a given date.

As a point of reference, the ratio of solid state to magnetic disk capacity in commercially available hybrid disks is at least 0.8% [41] at the time of this writing. At this ratio, the required metadata can be accommodated

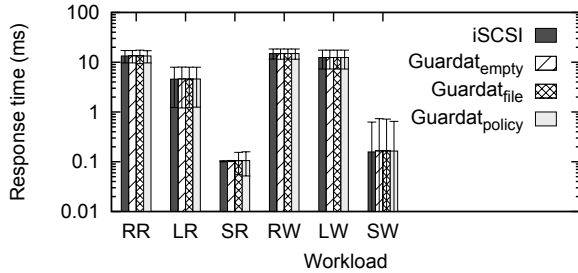


Figure 1: End-to-end I/O latency for synthetic workloads

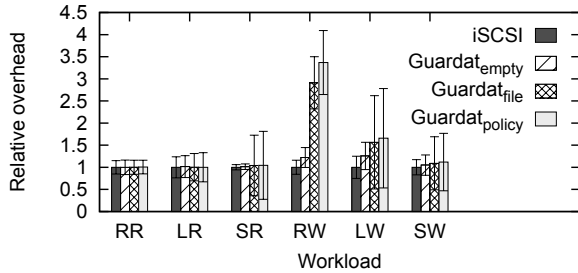


Figure 2: Normalized I/O latency for synthetic workloads excluding disk access latency

in the solid state memory for 99.99% of the snapshots. Newer combinations of Flash/disk devices like Apple’s Fusion Drive (128GB Flash/1TB HDD) achieve much higher Flash to disk capacity ratios, and would easily accommodate the metadata for all snapshots.

4.3.2 Read/write latency

Next, we examine the read/write latency of the Guardat prototype under synthetic workloads. For this experiment, we fill the 2TB disk with 3.8 million files, each spanning a single 512KB extent, and compare the read/write latency of the Guardat prototype with the original IET under three different configurations:

iSCSI: The plain IET iSCSI implementation.

Guardat_{empty}: No files are protected by a policy. The overheads are limited to the cost of communication between the kernel module and the Guardat daemon, and the (negative) check for a policy.

Guardat_{file}: An “allow all” policy is associated with each file. Each access to a disk block requires the Guardat daemon to lookup the metadata associated with the file and interpret the null policy.

Guardat_{policy}: Each file is protected by a policy selected at random from a set of 40,000 different policies, each of which allows access after a past date. The additional overhead includes fetching and interpreting the different policies.

Each configuration is exercised with three different access patterns (**Sequential**: blocks accessed in order of increasing block id, **Local**: each accessed block chosen

randomly within 40,000 block ids of the previous block, **Random**: each accessed block chosen randomly on the entire disk), and two access types (**Read** and **Write**). Each access reads or writes a 512B disk block.

For the different configurations, access patterns and types, Figure 1 show the absolute end-to-end access latency. Five runs were performed with each configuration, for a total of 100,000 accesses. Each run was started at a randomly chosen location on the disk. The bars show the average of the measured latencies; error bars indicate the standard deviation.

The Guardat latency overheads for local and random read/write accesses are negligible (below 1%), because they are dominated by the access latency of the magnetic disk. The Guardat overheads are more noticeable during the much faster sequential accesses (2.9% for **SR** and 4% for **SW** in the Guardat_{policy} configuration). During read accesses, Guardat can partially hide the policy check latency by issuing the disk read in parallel with the policy check, and squashing the read if the check is negative. Writes, on the other hand, cannot be scheduled safely until the policy check completes, thus the higher overhead. There is room for further optimization, for instance, by caching the results of previous policy evaluations in the kernel module while they remain valid, thus avoiding upcalls in the common case.

To zoom in on the delays introduced by Guardat, Figure 2 shows the measured end-to-end latencies *excluding the disk access latency*, and normalized to the latency of the unmodified IET server without the disk access latency. We see that the overheads relative to the original IET are small, except during random writes and, to a lesser extent, local writes. As noted above, the policy checks cannot be overlapped with the disk access during writes. Moreover, during random and local writes, the Guardat metadata lookups tend to miss the cache more often, leading to average latencies of up to 930 μ s compared to the IET latency of 278 μ s in the Guardat_{policy} configuration on random writes (**RW**). However, this overhead has little impact on the end-to-end latency, which is dominated by disk access latency, as shown in Figure 1.

4.3.3 Read/write throughput

Next, we examine the read/write throughput of the Guardat prototype, using the same configurations used in the latency experiment. However, instead of issuing 512B accesses sequentially, the test client issues four 64KB requests concurrently, which suffices to achieve maximal throughput in the experiment.

Figure 3 shows the absolute throughput of the various configurations and workloads. The results shown are the averages of 5 runs, each starting at a block id picked ran-

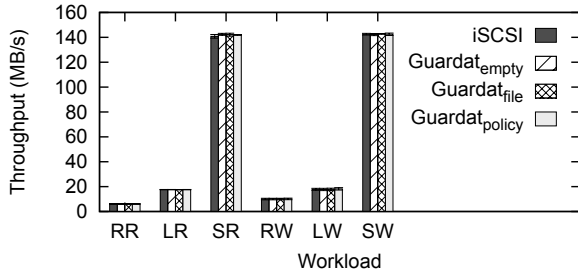


Figure 3: I/O throughput for synthetic workloads

Policy size	Domain size (number of entries)				
	1	2	4	8	16
1	2.2	3.4	5.8	10.7	20.4
2	4.6	10.4	28.9	95.1	345.8
3	7.0	24.0	121.2	770.5	5,518.1
4	9.4	50.9	485.3	6,156.4	88,319.3
5	11.9	104.9	1,951.3	49,234.7	1,411,800.8

Table 2: Evaluation latency in μs for varying policy size (number of predicates/variables) and domain size (number of cache entries per predicate/variable)

domly within the disk and comprising 20,000 accesses. Error bars indicate the standard deviation of all 100,000 accesses. The Guardat overheads for all access patterns are negligible (below 0.6%), because the Guardat policy checks are largely overlapped by disk accesses.

4.3.4 Policy evaluation latency

As discussed in section 3.1, worst case policy evaluation time is $m(D^n)$, where m is the size of the policy, D is the size of the Guardat cache and n is the number of variables in the policy. In Table 2, we show the measured worst case policy evaluation time for different D (in columns) and different m and n (in rows; $m = n$ in all experiments). The results correspond to the formula $m(D^n)$.

Although the table indicates (correctly) that policy evaluation could be a substantial bottleneck for some policies, in practice, we do not observe this bottleneck. The average policy evaluation latency of our most complex policy, MAL (Section 4.5) is only 27.7 μs , even though the policy has $m = 4$, $n = 4$ and $D = 40$. This is because of a careful implementation of the policy interpreter to consider more recent cache entries first. Our other example policies evaluate even faster; the average evaluation time of Guardat_policy (Section 4.3.2) is only 3.7 μs .

4.3.5 Flash memory wear

Because Flash memory can endure only a limited number of erase/program cycles, it must not wear out before the expected lifetime of the magnetic disk. To be con-

servative, we assume that the Flash must last at least 10 years. The lifetime is influenced by the size of the metadata, the rate of metadata updates (more updates require more Flash writes), and the Flash capacity. (A smaller capacity causes the Flash log to wrap around faster and leads to higher utilization, which in turn reduces cleaning efficiency and requires even more Flash writes.)

Under the configuration of Guardat_policy used above in 4.3.2, we keep track of how much wear the Flash experiences while presented with a series of metadata updates, i.e., adding and removing extents to a content file picked at random. Using only 4GB of Flash memory with a nominal lifetime of 100,000 erase/program cycles, we can accommodate up to 19.5M updates per day (225 per second). This is an extraordinarily high update rate and can accommodate even the most write-intensive applications.

4.4 Use case: Web server

Next, we study the capabilities and performance of the Guardat prototype as part of a Web server. The server holds a 220GB static snapshot of English language Wikipedia articles from 2008 [52] and images from 2005 [51], containing 15 million files with an average file size of 15KB and maximum file size of $\sim 500\text{KB}$.

We use three different machines connected by 10 Gbit Ethernet links to run the IET storage server, the Apache Web server (version 2.2.23)[4] and Apache HTTPSync-Client for Java (version 4.0-beta-1)[5]. The Web server either fetches the contents from the mounted iSCSI device or from Linux’s buffer cache. Apache does not use a dedicated disk cache.

The HTTP client asynchronously requests HTML pages from the Web server, using a workload based on the actual access counts of Wikipedia pages during one hour on April 1, 2012 [53]. Because our snapshot is much older and had fewer articles at the time, we ignore accesses to non-existing pages. In total, about 350,000 different pages were accessed in the trace, of which 250,000 are part of the 2008 snapshot. Since we do not have access to time stamps, we distributed the individual accesses evenly within an hour, and replayed the first 100,000 page requests.

We use the following Guardat policies: **Static content pages:** Allow updates signed by a fixed owner. We use 40,000 different owner identities and randomly assign them to the content files. **System binaries:** Allow updates signed by a special vendor signature only.

The workload and policies used are particularly challenging for Guardat, due to the large number of small protected files and a disk intensive workload.

Figure 4 shows the average throughput of three runs as a function of the number of concurrent HTTP accesses.

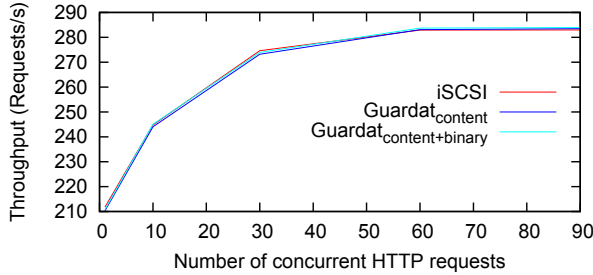


Figure 4: Throughput of 100,000 page loads, as a function of the number of concurrent accesses. Average of five runs, variations were below 0.6%

Each run loads 100,000 Wikipedia pages. Results are shown for four different configurations:

iSCSI: The plain IET iSCSI implementation.

Guardat_content: Guardat protecting content.

Guardat_content+binary: Guardat protecting content and Apache binaries (e.g. apachectl, httpd, libapr).

There is little difference (less than 1%) between the throughput achieved by the unmodified IET server and the Guardat_content and Guardat_content+binary configurations, confirming that the Guardat overheads are mostly hidden by disk access latencies.

To summarize, the use of Guardat in the Web server use case has negligible performance overhead. In terms of functionality, Guardat protects content files from tampering and corruption by unauthorized parties, and prevents intruders from modifying executables to install Trojan horses.

4.5 Mandatory access logging (MAL)

In the final experiment, we perform accesses to a file with a mandatory access logging policy. The MAL policy requires adding the proper log entry to a separate log file in order for an access to be allowed by Guardat. We use a 64MB file with or without the MAL policy in place. The primary file and log file reside on different disks attached to the same Guardat IET server. The version counter embedded in the primary file is stored in a block of available Flash memory. The client connects to the Guardat device and accesses the file as follows:

no_log: the file is accessed without any logging.

log: the file is accessed and the accesses are logged without policy enforcement.

Guardat_MAL: the policy-protected file is accessed and the accesses logged, the policy is enforced by Guardat.

Figures 5 and 6 show the access latency for sequential 4KB reads and writes, respectively, of the same location within the file. We vary the number of accesses per log entry. The bars show the average latency of 100,000 accesses and the error bars show the 98th percentile. In the

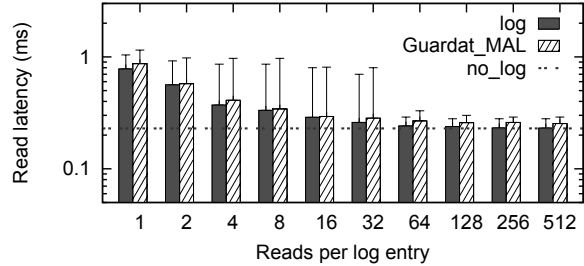


Figure 5: Read latency with MAL, voluntary and no logging

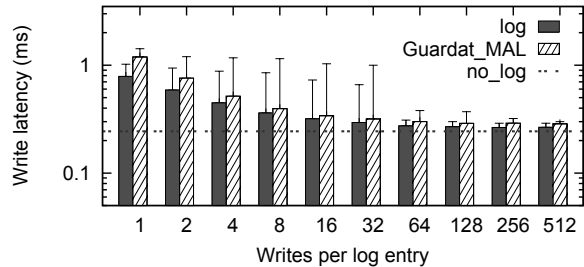


Figure 6: Write latency with MAL, voluntary and no logging

case of a single read/write, voluntary logging increases the latency from 0.20/0.16ms to 0.78/0.79ms, and policy enforcement increases the latency to 0.87/1.19ms. The higher cost of enforced logged writes reflects the need to update the version number. The synchronous log write contributes a significant part of the overhead; policy enforcement increases the read/write latency by 11.5% and 50.6%, respectively, for individually logged operations. As shown, the cost of MAL can be amortized by logging several operations with a single log entry.

4.6 Future Work

While clients connect directly to the Guardat prototype server to issue Guardat API, we're in the process of making the Btrfs filesystem Guardat-aware. Given a Guardat-aware filesystem, clients can access the Guardat API using I/O controls. As a result clients benefit from advanced filesystem features such as prefetching or buffering for better performance including Guardat API calls. Furthermore a Guardat-aware file system simplifies the access to Guardat API and the general programmability of systems using Guardat.

5 Related work

Storage work group specification. Although developed independently, the Guardat architecture bears some resemblance to a set of specifications for storage devices

standardized by the storage work group of the trusted computing group (TCG) [45]. Similar to Guardat, the TCG standard prescribes session-based communication with storage devices and access control on all calls to them. This industry interest supports the case for Guardat’s architecture. Unlike our work, however, the TCG standard does not describe a concrete design, implementation, or policy language, leaving these to device vendors; nor does it include certification of stored data by the storage device. Implementations exist for a subset of the TCG specification [44], providing full-disk encryption to preserve confidentiality of data upon device theft, loss or end of life. They do not include secure sessions, universal access checks or integrity policies, all of which Guardat does.

Guardat vs. trusted computing. At a high level, trusted computing (TC) relies on a trusted platform module (TPM) attached to a computer’s motherboard to provide a hardware root-of-trust [33], while Guardat relies on a controller (GDC) attached to a storage device, enclosure or server. While TC provides remote attestation of the hardware/software executing on a computer, Guardat attests the state of stored files, and enforces an application-defined policy for read and write accesses to files. TC provides sealed storage, where disk data is encrypted with a key stored within the TPM and released only when the computer runs a specific, trusted software configuration. Guardat instead enforces a declarative policy on all data accesses from untrusted client code. Compared to TC, Guardat can reduce the size of the TCB and its attack surface. Depending on the policy and implementation, the TCB can be as small as the Guardat controller. Lastly, TC can complement Guardat: A Guardat policy for access to a file can require that trusted software, verified via TC remote attestation, execute on the client computer. Conversely, TC can be used to protect the GDC.

Related trusted computing proposals. Building on TC, semantic attestation [18] enforces properties of a computation by a runtime verification substrate within a virtual machine monitor. Guardat provides a limited form of semantic attestation that enforces a data access policy, and does not require machine virtualization.

Excalibur [38] extends sealed storage with a primitive that binds cryptographically sealed data to a policy, such that a node can decrypt the data only if a trusted authority states that it obeys the policy (e.g., “this node is in Europe” or “this node is running Xen”). Guardat can be used to implement a similar capability, possibly with the help of a trusted authority as required by the Excalibur policy. However, Guardat can enforce many rich policies directly, without requiring an external trusted authority to map high-level policy descriptions to the nodes that meet those requirements.

Pasture [24] is a messaging and logging library that

enables data to be stored on an untrusted client machine while ensuring that a user cannot access the data without logging that access. Furthermore, users can delete unaccessed data in a way that provably prevents future access. The protocol relies on a TPM on the client machine. In Section 3, we describe a Guardat policy that can enforce the more general property of mandatory access logging for a collection of files.

Protecting data integrity and confidentiality. Butler et al. [9, 10, 11] describe storage devices that control access to storage segments contingent on the presence of a hardware token, or on successful remote attestation of the host computer. In Guardat, these forms of access control can be expressed as policies. Moreover, Guardat supports richer and per-file policies that can additionally depend on client authentication, wall-time clock or file contents, and Guardat supports certification of file states.

Commercially available self-encrypting disks [39] encrypt data to ensure its confidentiality when the device is lost or stolen. Guardat includes this capability as well, and additionally enforces rich data access policies. Web storage services like Amazon S3 [3] provide access control to a client’s data based on user identities, groups and roles, encryption for secure data storage and transit, and access logging. Guardat can enforce these (and many other) policies and provides file attestations. Because it operates at the storage layer, it does not require trust in the Cloud provider’s remaining platform.

In capability-based network-attached storage (NAS) [17, 2, 14], individual access requests include a capability, i.e., a tamper-proof description of client access rights. This capability is created out of band by a policy manager, a trusted component that serves all storage devices in a data center. A Guardat device, on the other hand, can interpret and enforce rich policies without relying on an external policy manager; thus, Guardat can operate in an otherwise untrusted or offline environment (unless a policy specifically delegates to an external verifier). Guardat can enforce state-based policies and certify the state of files and their policies, which capability-based NAS cannot.

Type-safe disks (TSD) [42] track the filesystem’s relationship among disk blocks using an extended block interface. Thus, a TSD can enforce basic filesystem integrity invariants, such as preventing access to unlinked blocks. A security extension called ACCESS adds read and write capabilities to selected disk blocks, thus enabling access control for entire files and directories. Guardat additionally supports cryptography and secure channels, which provides stronger security with respect to compromised hosts, buggy filesystems and operator mistakes. Also, Guardat’s policy language can support rich policies beyond filesystem metadata integrity.

Storage systems such as Self Securing Storage

(S4) [46] and NetApp's SnapVault [19] RAID storage server retain shadow copies of overwritten data or disable writes for a given period of time to address the specific problem of accidental or malicious corruption of data. Guardat can enforce these and much richer integrity constraints (Section 3), as well as confidentiality and accountability.

jVPFS [50, 49] is a stacked, microkernel-based filesystem that combines a small, isolated trusted component with a conventional untrusted filesystem. jVPFS uses encryption, hash trees and logging to ensure data confidentiality and integrity. Guardat can provide similar functionality at the storage layer, and supports a much richer set of policies.

The PCFS [15] and PFS [48] filesystems enforce integrity and confidentiality policies expressed in rich policy languages similar to Guardat's. However, unlike Guardat, PCFS and PFS trust the entire storage stack below the filesystem, cannot enforce integrity policies that depend on the content or size of files, do not certify the state of stored files, and can be bypassed by booting into a different configuration. PFS policies are expressed in a formal logic similar in expressiveness to the Guardat policy language. PCFS uses a formal logic that is more expressive than the Guardat policy language, but much more expensive (in terms of time and space) to implement. The Guardat policy language deliberately avoids policy features like recursive predicates that increase complexity but are rarely used in practice.

Protecting data availability. Storage systems like RAID [34], snapshotting filesystems [20, 32, 28] and some backup utilities [6, 29] use redundancy to ensure data availability. Guardat addresses the orthogonal problem of ensuring integrity, confidentiality and access accounting in the face of human error, adversarial threats and software bugs (e.g., a bug in a backup application that overwrites backed up data [16]). In practice, Guardat must be combined with redundant storage to ensure the availability of data in case of a media failure, loss, destruction or failure of a Guardat device.

Extended storage functionality. Commercial hybrid disks [41] package a magnetic disk drive with a modest amount of NAND Flash memory, used as a non-volatile write-back cache to increase performance. Guardat uses a comparable amount of Flash memory to store its policy metadata but, in addition, protects data. Object-based storage devices replace the traditional block-based with an object-based interface [27]. These systems offer capability-based security for whole objects, which we already compared to. Part of the Guardat API is also object-based, and could therefore be integrated with an emerging object-based storage standard. Several storage subsystems like active disks [37], semantically smart disks [43] and differentiated storage services [26] in-

clude program logic to improve performance. Guardat addresses the orthogonal concerns of data confidentiality, integrity and access accounting.

Pennington et al. [35] describe an intrusion detection system (IDS) at the storage layer, which raises an alarm when an access matches a per-file or global rule. Guardat instead is able to *enforce* per-file security policies, and these policies can be richer than the rules of an IDS system. However, intrusion detection rules could be specified as Guardat policies that allow offending accesses but log an alarm record.

6 Conclusion

To the best of our knowledge, Guardat is the first system that enforces rich, per-file confidentiality, integrity and access accounting policies at the storage layer, and attests the state of files. Policies are expressed in a concise declarative language and can be predicated on a wide range of conditions, including client authentication, remote attestation, physical authorization tokens, trusted wall-clock time, and the state (content) of files, even at sub-file granularity. Enforcing policies at the storage layer reduces the attack surface and, in many cases, the size of the TCB relied upon for enforcement. Existing techniques, on the other hand, either rely on a larger TCB, spread the specification and enforcement of policies affecting a given data file over many different components and layers of a system, or support a smaller set of policies. The Guardat design is rich enough to enable powerful policies, primitives and applications, yet is amenable to an efficient implementation, as demonstrated by an experimental evaluation.

References

- [1] AGRAWAL, N., BOLOSKEY, W. J., DOUCEUR, J. R., AND LORCH, J. R. A five-year study of file-system metadata. *Trans. Storage* 3, 3 (2007).
- [2] AGUILERA, M. K., JI, M., LILLIBRIDGE, M., MACCORMICK, J., OERTLI, E., ANDERSEN, D. G., BURROWS, M., MANN, T., AND THEKKATH, C. Block-level security for network-attached disks. In *Proc. of the 2nd USENIX FAST* (2003).
- [3] Amazon simple storage service (S3). <http://aws.amazon.com/s3/>.
- [4] APACHE SOFTWARE FOUNDATION. Apache http server. <http://httpd.apache.org/>, 2012.
- [5] APACHE SOFTWARE FOUNDATION. Apache httpasyncclient. <http://hc.apache.org/httpcomponents-asyncclient-dev/index.html>, 2012.
- [6] APPLE INC. Time Machine. <http://www.apple.com/osx/what-is/>.
- [7] ARM. ARM Security Technology. http://infocenter.arm.com/help/topic/com.arm.doc/prd29-genc-009492c/PRD29-GENC-009492C_trustzone_security_whitepaper.pdf, ARM Technical White Paper, 2009.
- [8] BECKER, M. Y., FOURNET, C., AND GORDON, A. D. Design

- and semantics of a decentralized authorization language. In *20th IEEE Computer Security Foundations Symposium (2007)*.
- [9] BUTLER, K., MCLAUGHLIN, S., MOYER, T., AND MCDANIEL, P. New security architectures based on emerging disk functionality. IEEE Computer Society.
 - [10] BUTLER, K. R. B., MCLAUGHLIN, S. E., AND MCDANIEL, P. D. Rootkit-resistant disks. In *Proc. of the ACM CCS (2008)*.
 - [11] BUTLER, K. R. B., MCLAUGHLIN, S. E., AND MCDANIEL, P. D. Kells: a protection framework for portable data. In *Proc. of the ACSAC (2010)*.
 - [12] COOPER, D., SANTESSON, S., FARRELL, S., BOEYEN, S., HOUSLEY, R., AND POLK, W. Internet X.509 public key infrastructure certificate and certificate revocation list (CRL) profile. RFC 5280. <http://www.ietf.org/rfc/rfc5280.txt>, 2008.
 - [13] DETREVILLE, J. Binder, a logic-based security language. In *Proc. of the S&P (2002)*.
 - [14] FACTOR, M., NAOR, D., ROM, E., SATRAN, J., AND TAL, S. Capability based secure access control to networked storage devices. In *Proc. of the 24th IEEE MSST (2007)*.
 - [15] GARG, D., AND PFENNING, F. A proof-carrying file system. In *Proc. of the 31st IEEE S&P (2010)*.
 - [16] GARRET, R. A Time Machine time bomb. <http://blog.rongarret.info/2009/09/time-machine-time-bomb.html>.
 - [17] GIBSON, G. A., NAGLE, D. F., AMIRI, K., BUTLER, J., CHANG, F. W., GOBIOFF, H., HARDIN, C., RIEDEL, E., ROCHBERG, D., AND ZELENKA, J. A cost-effective, high-bandwidth storage architecture. In *Proc. of the 8th ACM ASPLOS (1998)*.
 - [18] HALDAR, V., CHANDRA, D., AND FRANZ, M. Semantic remote attestation: A virtual machine directed approach to trusted computing. In *Proc. of the 3rd USENIX Virtual Machine Research And Technology Symposium (2004)*.
 - [19] HAYAKAWA, M. WORM Storage on Magnetic Disks Using SnapLock Compliance and SnapLock Enterprise. Tech. Rep. TR-3263, Network Appliance, 2007.
 - [20] HITZ, D., LAU, J., AND MALCOLM, M. File system design for an NFS file server appliance. In *Proc. of the USENIX Winter Technical Conference (1994)*.
 - [21] INTEL CORP. AESNI library. <http://software.intel.com/en-us/articles/download-the-intel-aesni-sample-library/>, 2011.
 - [22] INTEL CORP. Fast SHA256. <http://download.intel.com/embedded/processor/whitpaper/327457.pdf>, 2012.
 - [23] INTEL CORP. Software Guard Extension Programming Reference. <http://software.intel.com/sites/default/files/329298-001.pdf>, 2012.
 - [24] KOTLA, R., RODEHEFFER, T., ROY, I., STUEDI, P., AND WESTER, B. Pasture: Secure offline data access using commodity trusted hardware. In *Proc. of the 10th USENIX OSDI (2012)*.
 - [25] LI, N., AND MITCHELL, J. C. Datalog with constraints: A foundation for trust management languages. In *Proc. of the 5th Symposium on Practical Aspects of Declarative Languages (2003)*.
 - [26] MESNIER, M., CHEN, F., LUO, T., AND AKERS, J. B. Differentiated storage services. In *Proc. of the 23rd ACM SOS (2011)*.
 - [27] MESNIER, M., GANGER, G., AND RIEDEL, E. Object-based storage. *Communications Magazine* 41, 8 (2003).
 - [28] MICROSOFT CORP. What is volume shadow copy service? [http://technet.microsoft.com/en-us/library/cc757854\(ws.10\).aspx](http://technet.microsoft.com/en-us/library/cc757854(ws.10).aspx).
 - [29] MICROSOFT CORP. Windows Backup and Restore. <http://www.microsoft.com/athome/setup/backupdata.aspx#fbid=17X90d97a1I>.
 - [30] OCZ TECHNOLOGY INC. Deneva 2 data sheet. <http://www.oczenterprise.com/ssd-products/deneva-2-c-sata-6g-2.5-slc.html>, 2011.
 - [31] OPENSLL CRYPTOGRAPHIC LIBRARY. <http://www.openssl.org/docs/crypto/crypto.html>, 2012.
 - [32] ORACLE CORPORATION. Solaris ZFS. <http://www.oracle.com/us/products/servers-storage/storage/storage-software/031857.htm>.
 - [33] PARNO, B., MCCUNE, J. M., AND PERRIG, A. Bootstrapping trust in commodity computers. In *Proc. of the 31st IEEE S&P (2010)*.
 - [34] PATTERSON, D. A., GIBSON, G., AND KATZ, R. H. A case for redundant arrays of inexpensive disks (RAID). In *Proc. of ACM SIGMOD (1988)*.
 - [35] PENNINGTON, A. G., GRIFFIN, J. L., BUCY, J. S., STRUNK, J. D., AND GANGER, G. R. Storage-based intrusion detection. *ACM Trans. Inf. Syst. Secur.* 13, 4 (Dec. 2010).
 - [36] PIMLOTT, A., AND KISELYOV, O. Soutei, a logic-based trust-management system. In *Proc. of the 8th International Symposium on Functional and Logic Programming (FLOPS) (2006)*.
 - [37] RIEDEL, E., FALOUTSOS, C., GIBSON, G., AND NAGLE, D. Active disks for large-scale data processing. *IEEE Computer* 34, 6 (2001).
 - [38] SANTOS, N., RODRIGUES, R., GUMMADI, K. P., AND SAROUI, S. Policy-sealed data: A new abstraction for building trusted cloud services. In *Proc. of the 21st USENIX Security Symposium (2012)*.
 - [39] SEAGATE TECHNOLOGY LLC. Self-encrypting hard disk drives in the data center. Tech. Rep. TP583, 2007.
 - [40] SEAGATE TECHNOLOGY LLC. Barracuda Data Sheet. <http://www.seagate.com/files/staticfiles/docs/pdf/datasheet/disc/barracuda-xt-ds1696.3-1102us.pdf>, 2011.
 - [41] SEAGATE TECHNOLOGY LLC. Momentus XT Data Sheet. <http://www.seagate.com/files/staticfiles/docs/pdf/datasheet/disc/momentus-xt-data-sheet-ds1704-4-1209-us.pdf>, 2012.
 - [42] SIVATHANU, G., SUNDARARAMAN, S., AND ZADOK, E. Type-safe disks. In *Proc. of the 7th USENIX OSDI (2006)*.
 - [43] SIVATHANU, M., PRABHAKARAN, V., POPOVICI, F. I., DENEHY, T. E., ARPACI-DUSSEAU, A. C., AND ARPACI-DUSSEAU, R. H. Semantically-smart disk systems. In *Proc. of the 2nd USENIX FAST (2003)*.
 - [44] STORAGE WORK GROUP OF THE TRUSTED COMPUTING GROUP. Self-encrypting drives take off for strong data protection. http://www.trustedcomputinggroup.org/community/2010/03/selfencrypting_drives_take_off_for_strong_data_protection, 2011.
 - [45] STORAGE WORK GROUP OF THE TRUSTED COMPUTING GROUP. TCG storage architecture core specification. http://www.trustedcomputinggroup.org/resources/tcg_storage_architecture_core_specification, 2011.
 - [46] STRUNK, J. D., GOODSON, G. R., SCHEINHOLTZ, M. L., SOULES, C. A. N., AND GANGER, G. R. Self-securing storage: Protecting data in compromised systems. In *Proc. of the 4th USENIX OSDI (2000)*.
 - [47] THE ISCSI ENTERPRISE TARGET PROJECT. <http://iscsitarget.sourceforge.net/>, 2011.
 - [48] WALSH, K., AND SCHNEIDER, F. B. Costs of security in the PFS file system. Tech. rep., Computing and Information Science, Cornell University, 2012.
 - [49] WEINHOLD, C., AND HÄRTIG, H. VPFS: Building a virtual private file system with a small trusted computing base. In *Proc. of the 3rd ACM SIGOPS EuroSys (2008)*.
 - [50] WEINHOLD, C., AND HÄRTIG, H. jVPFS: Adding robustness to a secure stacked file system with untrusted local storage components. In *Proc. of the USENIX ATC (2011)*.
 - [51] WIKIMEDIA FOUNDATION. Image Dump. <http://archive.org/details/wikimedia-image-dump-2005-11>, 2005.
 - [52] WIKIMEDIA FOUNDATION. Static HTML dump. <http://dumps.wikimedia.org/>, 2008.
 - [53] WIKIMEDIA FOUNDATION. Page view statistics April 2012.

7 Details of Guardat API calls

This section summarizes the Guardat API calls and describes the session handshake protocol.

Session API We describe the handshake protocol that authenticates the client and Guardat to each other and establishes session keys. We use some abbreviations. C and G denote the client and Guardat respectively. For a principal A , K_A denotes its public key and K_A^{-1} denotes its corresponding private key. Guardat's public key K_G is assumed to be known to everyone (through a manufacturer-provided certificate). The client and Guardat use two freshly chosen random nonces, denoted N_C and N_G respectively. $\text{sig}(K, M)$ denotes a message M and its signature with private key K . $\text{enc}(K, M)$ denotes the encryption of message M with public key K . (M_1, M_2) denotes concatenation of messages M_1 and M_2 . $A \rightarrow B : M$ means that A sends message M to B . Our handshake protocol is:

$$\begin{aligned} C \rightarrow G &: K_C, \text{sig}(K_C^{-1}, \text{enc}(K_G, N_C)) \\ G \rightarrow C &: \text{sig}(K_G^{-1}, (0, \text{enc}(K_C, (N_C, N_G)))) \\ C \rightarrow G &: \text{sig}(K_C^{-1}, (1, \text{enc}(K_G, (N_C, N_G)))) \end{aligned}$$

At the end of the protocol, both C and G have seen their nonces signed by the other party; this authenticates them to each other. The sequence numbers 0 and 1 in the second and third messages distinguish the two messages from each other, preventing a man in the middle from causing confusion through replays. The session keys are derived by the client and Guardat using key derivation functions on the concatenation of the two nonces, (N_C, N_G) .

This protocol is based on the station-to-station protocol and provides direct bidirectional authentication and perfect forward secrecy. The handshake protocol is implemented using two API calls, named `handshake1` and `handshake2`. The first call, `handshake1`, corresponds to the first message in the protocol and its return value corresponds to the second message. The second call, `handshake2`, corresponds to the third message and its return value is just a confirmation that the session has been established. (The return value is irrelevant to security of the protocol, but it is needed to tell the client to proceed.)

- `handshake1(message)`: The message should be of the form of the first message in the protocol. If correct, the return value is the second message of the protocol and a session id to link the second call.
- `handshake2(sessionId, message)`: The message should be of the form of the third message in the

protocol. The return value is either success or failure. If the value is success, then `sessionId` is used as the identifier for the rest of the session.

The API call `endSession(sessionId)` ends a session.

Object API To allow flexible content hashing during `readBatched()` and `fresh()` calls, Guardat provides hash computation buffers to which the client can selectively choose to add data during these calls. Once the client has added all the data it needs to a buffer, it finalizes the buffer, which moves the hash of the content in the buffer to either the session cache or the batch cache (depending on whether the hash buffer has current or new object content). In the implementation, each non-finalized hash buffer is an open hash computation and newly appended content is hashed immediately. Hash buffers are accessible through the following API.

- `initHash(batchId, currOrNew)`: Initialize a new hash buffer for the object associated with the open batch `batchId`. The Boolean flag `currOrNew` indicates whether the buffer will hold a hash of the object's current version or its updated version. This choice is enforced in calls that add to the buffer. Returns a unique identifier for the hash buffer, `hashId`.
- `closeHash(hashId)`: Finalize the hash buffer `hashId` and move the hash in it to the session cache if the buffer has current object content or to the batch cache if the buffer has new object content.

API calls to start and end batches and to read and update objects were described in Section 2.1. We summarize them below with details of caching flags that we omitted from Section 2.1.

- `openBatch(sessionId, objName)`: Start batch on object with name `objName` in session identified by `sessionId`. If `objName` does not exist, create it. Returns a new `batchId` on success.
- `readBatched(batchId, off, len, buf, cacheFlags, cacheIntervals, hashId)`: Read `len` bytes from logical offset `off` of the object associated with `batchId` and return the result in `buf`. The `cacheFlags` indicate whether or not the read content should be added to the session cache and whether or not it should be added to a hash buffer. If either is the case, then `cacheIntervals` specifies which logical ranges of the read content need to be added. If content is to be added to a hash buffer, `hashId` identifies the buffer. This call evaluates the **read** policy rule.
- `reuse(batchId, off, len, off')`: Take content in the logical range `[off, off+len-1]` from the current version of the object associated with `batchId` and insert it at offset `off'` in the new version.

- `fresh(batchId, off, len, buf, off', cacheFlags, cacheIntervals, hashId)`: Write `len` bytes from buffer `buf` to the extent starting at byte offset `off` on disk and add the resulting extent to the new version at offset `off'`. The arguments `cacheFlags`, `cacheIntervals` and `hashId` have roles similar to those in `readBatched()`, but the batch cache, and not the session cache, is affected.
- `setPolicy(batchId, new_policy)`: Set the policy of the new version to `new_policy`.
- `endBatch(batchId)`: Close the batch identified by `batchId`. Evaluates the **update** policy rule and, if the policy is being modified, also the **setpolicy** rule.
- `destroy(objName)`: Destroy the metadata of the empty object named `objName`. Evaluates the **destroy** policy rule.

Certificate API The certificate API generates nonces, attests content and forwards third-party certificates to Guardat for policy evaluation.

- `getNonce(sessionId)`: Return a new nonce that is to be embedded in a third-party certificate later.
- `setCertificate(sessionId, certificate)`: Provide certificate for use in later policy evaluation.
- `attest(sessionId, objName, attestFlags, hashId, nonce)`: Attest `objName`'s metadata and optionally a hash of selected content. The argument `attestFlags` specifies which of the following need to be attested: the list of physical extents, policy, policy hash and content hash. If content hash is to be attested, `hashId` points to an entry in the session cache which has that hash. Guardat returns a single certificate containing all requested vectors and the client-provided nonce. The **read** policy rule of `objName` is evaluated.

Backwards compatibility Standard block-device API calls `read(blk,cnt,buf)` (read `cnt` disk blocks starting at block `blk` into buffer `buf`) and `write(blk,cnt,buf)` can be used to read and write disk extents. The **read** or **update** policy rule of all extents that overlap the accessed extents is evaluated.

8 Details of the Guardat policy language

We summarize in this section predicates available in the Guardat policy language. In addition to these predicates, policies may use any uninterpreted predicates established through third-party certificates.

The following universal predicates are available in all policy rules.

- `object_name_is(O)`: The object being accessed has name `O`.
- `session_is(K)`: The current session has been authenticated with public key `K`.
- `time_is(T)`: The internal timing counter is currently `T`.
- `guardat_key_is(K)`: The public key of this Guardat installation is `K`.
- Arithmetic, string comparison: $t_1 == t_2$, $t_1 < t_2$, $t_1 \leq t_2$.
- `is_subset(R1, R2)`: Range set `R1` is a subset of range set `R2`.
- `disjoint(R1, R2)`: Range sets `R1` and `R2` are disjoint.
- `(obj, off, len) says C`: The content at offset `off` in object `obj` is `C`. (Based on the session cache.)
- `(obj, ((off1, len1), ..., (offn, lenn))) hasHash H`: The bytes stored in the given list of `(off,len)` pairs in object `obj` have hash `H`. (Based on the session cache.)

Specific contextual predicates are available in each policy rule. We list these below, categorized by the policy rules. The **destroy** rule admits no contextual predicates.

Read rule The following predicates are available in the **read** policy rule.

- `isAttest()`: True if the **read** rule is being evaluated in an `attest()` call, and false otherwise.
- `access_locations_are(R)`: The set of logical `(off,len)` pairs read from the object is `R`.
- `access_physical_extents_are(E)`: The set of physical extents read is `E`.
- `access_length_is(L)`: The number of bytes read is `L`.

Update/setpolicy rule The rules for **update** and **setpolicy** evaluate in the same call, `endBatch()`, and, hence, they admit the same contextual predicates with only one exception that is shown later. When the **update** rule evaluates in the call `write()`, some of these predicates always evaluate to false. These predicates are marked with `*`.

- `isCreate*`: True if and only if the batch executed on an object that did not already exist.

- `current_length_is(L)`: The length of the current version of the object is L bytes.
- `new_length_is*(L)`: The length of the new version of the object is L bytes.
- `current_physical_extents_are(E)`: The current version of the object spans the set E of physical extents.
- `new_physical_extents_are*(E)`: The new version of the object spans the set E of physical extents.
- `updated_locations_are(M)`: The set of logical (off, len) pairs updated during the batch is M .
- `read_locations_are*(R)`: The set of logical (off, len) pairs read during the batch is R .
- `current_pol_hash_is(H)`: The current policy has hash H .
- `new_pol_hash_is*(H)`: The new policy has hash H .
- `(off, len) willsay*C`: The new content at offset off is C . (Based on the batch cache.)
- `((off1, len1), ..., (offn, lenn)) willHaveHash* H`: The new bytes stored in the given list of (off, len) pairs have hash H . (Based on the batch cache.)

The following predicate is available only in the **update** policy rule, not the **setpolicy** rule, and can be used to distinguish evaluation of **update** in the `endBatch()` call from that in the `write()` call.

- `isWrite()`: True if and only if evaluation is part of the `write()` call.

9 Compatibility with existing filesystems

We sketch how our Guardat prototype can be used with an existing, unmodified filesystem, which is not aware of Guardat and issues only ordinary block reads and writes. In this compatibility mode, applications are linked with a library, which implements the standard POSIX filesystem interface, and provides additional operations for applications to authenticate, set a policy for a file, provide certificates required by policies, and request attestations. The library interacts with the Guardat userspace daemon directly and makes API calls to associate block read and write operations issued by the filesystem with an object,

client session and batch. We note that the library is untrusted and does not require extra privileges. In particular, the library only executes Guardat calls on behalf of applications that the applications are allowed to execute themselves.

To determine if a block read operation is allowed, the Guardat daemon maps the requested block number to the associated object (if one exists) using the block-to-object B-tree. To further be able to map the read operation to an authenticated session, we impose the limitation that only a single batch or session can be open for a given object at any given time in compatibility mode.

Write operations may refer to an extent not currently associated with any object. (When a file is extended, the filesystem allocates new blocks.) Therefore, prior to writing new data to a file, the application library provides the Guardat daemon with a vector of hashes of aligned blocks containing the new data. This vector enables the daemon to associate subsequent writes issued by the filesystem with the correct object, offset, session and batch. In order to avoid ambiguity, two blocks with the same hash value may not be outstanding at the same time. The daemon enforces this condition by temporarily refusing to accept a block hash vector that contains an element that is already present among the current set of outstanding vectors.

When the kernel module receives a write command, it computes the hash of blocks to be written, and sends the hash to the daemon along with the request. The daemon matches write requests with the list of hashes provided by the compatibility library. Computing hashes in the kernel avoids sending data from the kernel to the userspace daemon.

The compatibility mode has limitations. As mentioned above, only a single session and batch may be active for any given object, which can lead to some loss of performance in workloads with concurrent accesses to the same file. Also, because the filesystem is unaware of Guardat, any attempt by the filesystem to relocate a file with an associated integrity policy may fail. As a result, defragmentation of an unmodified filesystem requires a modified defragmentation utility. Object data encrypted with a session key must be communicated between library and the Guardat daemon without going through the iSCSI driver, to avoid polluting the filesystem's buffer cache with session-encrypted data. These limitations can be lifted by modifying a filesystem to use the extended Guardat API, which is a subject of ongoing work.