1 Storage Networking Industry Association

2 Object-Based Storage Devices

3

# Object Store Security Document          T10/03-279r0

5

**Revision: 8**
**Last Revised: 8/12/2003**

8

## Abstract

10

This document presents the requirements, motivation and a proposal for the security protocol for object store. This protocol is based upon the original Network Attached Storage Device (NASD) work [8] as well as other work on secure object stores 9.

14

## Related Documents

- The OSD White Paper offers an introduction to OSD and its applications.
- The OSD Requirements Document discusses requirements of the OSD applications discussed in the white paper.
- The T10 SCSI Draft Standard for OSD implements the OSD framework for the SCSI architecture model.

21

# Table of Contents

2

# 0   Revision History

## 0.1  Revision 1

**Authors:** Dalit Naor, Michael Factor, Julian Satran (IBM), Don Beaver, Erik Riedel (Seagate), and David Nagle (Panasas).

## 0.2  Revision 2

Authors: Dalit Naor and Michael Factor

## 0.3  Revision 3

Authors: Erik Riedel

Changes: Added description of Levels 2 and 3. Added sequence diagrams and detailed message arguments. Changed "client" to "host" throughout for consistency. Use "OSD" instead of "object store". Reorganized items in the introduction. Added more white space for better readability.

## 0.4  Revision 4

Author: Dalit Naor.

Changes: Incorporated presentation changes and open issues from comments submitted by May 6 from David Nagle, Erik Riedel, Dalit Naor, Michael Factor.

## 0.5  Revision 5

Author: Michael Factor

Changes: Incorporate changes to open issues as discussed in the SNIA Symposium in Boston. Main changes include, added section on protocol between OSD and Security Manager, added description of error codes, cleaned up presentation of sections 2 and 3

## 0.6  Revision 6

Author: Michael Factor

Changes: Fixed formatting.  Cleanup of Chapter 1 (provided an introduction, deleted objectives section as being repetitive with requirements, deleted some discussion of levels 4-7 of security, editorial corrections).  Moved key management to Chapter 8 and added Chapter 8.  Added a separate chapter on Nonces (Chapter 4) in place of the section on nonces which was in the discussion of level 2 security.

## 0.7 Revision 7

Authors: Michael Factor, Dalit Naor

Minor changes from prior revision: Integrate editorial comments, in particular cleanup usage of capability and credential. Add section 2.5 to describe credentials for creating objects without specifying an object ID.

## 0.8 Revision 6

Author: Michael Factor

Additional minor changes: clarify object version number in credential (and renamed to object version tag), (re)add creation time to credential, clarify which keys protect credentials for commands that are not scoped to a partition. Also remove descriptions of how do we know the appropriate security level; this will not be addressed in the first version of the standard

# 1  Introduction

118

119  Object storage is a new storage paradigm (in particular for network accessible storage) in which
120  the abstraction of an array of blocks is replaced with an abstraction of a collection of objects.
121  In object storage, a client accesses data by specifying the identity of an object along with an
122  offset in the object, and the implementation of the storage is responsible for mapping the offset
123  to the actual location on the physical storage.  From a security perspective, the main change
124  between object storage and today's block storage paradigm is that every command is
125  accompanied by a cryptographically secure capability.  Thus, object storage provides the means
126  of having secure, fine-grained access to storage.

127  This document presents the requirements, motivation and a description of the object store
128  security protocol.    The goals of this document are multifold.  First, it is intended to specify the
129  behavior of the high-level protocol in sufficient detail to allow a direct mapping to a standard
130  specification in a particular transport (*e.g.,* in SCSI).  It is also intended to explain the protocol
131  in a way that it can be shared with security experts, outside of the OSD community, to allow an
132  independent review of its correctness.  Finally, it is intended as a general background material to
133  explain OSD security.

134  One major goal for OSD security is to work well both on top of a secure network infrastructure
135  and in environments where there is no such infrastructure.  This requirement has led us to define
136  multiple of levels of security which reflect the assumptions on the underlying infrastructure and
137  the protection required.

138  This document is organized as follows.   This chapter describes the basic security model, and
139  the requirements we imposed upon ourselves.  The next chapter describes the structure of the
140  capabilities/credentials and the basic message flow; this structure and flow is common for all
141  levels of security.  Chapter 3 describes the details of the security level which ensures integrity of
142  the security mechanism; this level is ideally suited for use on top of a secure network
143  infrastructure, but it also can be used in environments where there is no concern of network-
144  type attacks.   Chapters 5 and 6 describe two different levels of security intended for use on
145  insecure networks; they differ in whether or not they secure the data.  Chapters 7 and 8
146  describe security aspects that are not on the main data path.

## 1.1.1  Basic Security Model

147

148  The object store security model is a credential-based access control system composed of three
149  active entities: the object store, a security manager, and a client/host. Each entity plays a
150  different role.

151  As a credential-based access control system, all requests to the object store must be
152  accompanied with a valid capability that allows the host to perform the requested operation. A
153  *credential* is a cryptographically secured capability and a *capability* is a set of rights the holder
154  has on an object (or set of objects).

155  The role of the security manager is to generate credentials for authorized hosts at the request of
156  the host. The protocol between the host and the security manager is not defined as part of the

157 OSD protocol; however, the structure of the credential returned from the security manager to
158 the host is defined.  In addition, the protocol between the OSD and the security manager is
159 specified.

160 The role of the OSD is to validate a capability presented by a host:

161 1.  The requested operation is permitted by the capability based on *a)* the type of operation
162     (*e.g.,* read, write) and *b)* a logical match of the specified attributes
163 2.  The capability has not been tampered with, *i.e.,* it was generated by the security manager
164     and was rightfully obtained by the host that presents it (either directly or via delegation).

165 The object store can validate that a host rightfully obtained a capability since a credential
166 contains both the capability and a secret part (CAP_Key – see section 2), which the host uses
167 to sign its messages to the object store.  Without this secret part, which should be transferred on
168 an encrypted channel from the security manager to the host, the host cannot generate validly
169 signed messages.  Note this protocol does allow delegation of a credential if a host transfers
170 both the secret part of the credential as well as the public capability arguments.

171 The role of the host is to follow the protocol. While the host is not trusted to follow the
172 protocol, the protocol is structured in such a way that it is in the host's self-interest to follow the
173 protocol. In other words, if the host does not follow the protocol, it will not receive service from
174 the OSD. The figure below shows this basic flow.



175

176                          **Figure 1.    Basic System Structure**

177 We specify seven levels of security, of which only the first three are within the scope of the
178 current proposal:

179             Level 1 – Integrity of capability
180             Level 2 – Integrity of arguments
181             Level 3 – Integrity of data in transit
182             Level 4 – Privacy of arguments
183             Level 5 – Privacy of data in transit

7

186 Levels 2-4 correspond to the security levels defined in the original NASD work [8]. Level 1 is
187 best suited for the case where the network between the OSD and the host is secured; it can be
188 used as another layer on top of the network security 9.

189 With Level 1, only *access security* is handled within the OSD specification, and *network*
190 *security* is handled by an external, network-specific means (*e.g.,* IPSec or FCS).

191 In order to implement Level 3 efficiently, the authentication hashes for user data must be carried
192 by the underlying transports. The structure and interpretation of these hashes will be specified in
193 this document, but an efficient mapping to a particular network transport layer (*e.g.,* FC or
194 TCP/IP) is left to external specifications.[1]

## 1.1.2  Trust Assumptions

196 Trust assumptions describe how each element of the system trusts the other elements of the
197 system. The OSD is a trusted component. This means that once a host authenticates that it is
198 talking to a specific OSD, it trusts the OSD to:

199 1. provide integrity  for the data while stored
200 2. follow the protocol
201 3. not be controlled by an adversary

202 The host can authenticate that it is talking to the intended OSD, *i.e.,* the one for which the
203 security manager has granted it credentials, either via the use of an externally provided
204 authenticated channel or as part of each command using mechanisms defined in this protocol.

205 The security manager is also a trusted component. After it is authenticated,[2] it is trusted to:

206 1. safely store long-lived keys
207 2. compute access controls correctly according to the policy it implements[3]
208 3. follow the protocol
209 4. not be controlled by an adversary

210 The trust assumption on the host is that a user trusts their own operating system to protect them
211 from malicious clients on the same machine, *e.g.,* protect its CAP_Key.  We do not trust the
212 host to correctly follow the protocol; however, the host will not receive service if it does not
213 follow the protocol.

## 1.1.3  Security Flow and Channel Requirements

215 As mentioned above, when a host wishes to access an object (or set of objects), it makes a
216 request to the security manager for a credential allowing the intended operation.[4]  In this

---

[1] The difficulty is in the ordering of the hashes with respect to the data while in transit and during
verification at the device. This is discussed further in Section 6.1.
[2] Authentication of the Security Manager by the host is out-of-scope of this protocol.
[3] The definition of this policy is outside of the scope of this proposal.

217 request, the host must specify the OSD and partition (see section 2.3.2) on which it wishes to
218 perform the operation; the identity of the object(s) it wishes to access; and the operation(s) it
219 wishes to perform. The security manager upon receiving this request may need to authenticate
220 the host making the request.[5]  After authenticating the host, the security manager applies its
221 policy to determine whether the client is authorized to perform the requested operation(s) on the
222 indicated object(s). If not, the security manager will fail the request for the credential.
223 Otherwise, the security manager will generate a credential including the requested capability; this
224 credential is cryptographically secured by a secret shared between the security manager and the
225 OSD. The credential is then sent from the security manager to the host over a channel which is
226 encrypted and authenticated. Other than specifying the structure of the credential returned from
227 the security manager to the client, the protocol between the client and security manager is not
228 defined by the OSD protocol.

229 The host must present a capability on each operation it executes against the OSD.  When the
230 OSD receives the capability, it verifies that it has not been modified, using the secret it shares
231 with the security manager.[6] If the credential has not been modified (and is properly held by the
232 requesting client), the OSD will permit the operation based upon the rights encoded in the
233 capability.

234 When using Level 1 protection, we assume an existing network infrastructure that provides
235 secure channels (*e.g.,* IPsec) between the OSD and the host. More precisely, if we are running
236 over a secure channel, we require both parties of the communication to know that they are
237 communicating with the parties that originally established the channel (an authenticated but
238 anonymous channel). We do not imply a requirement for privacy, *i.e.,* we assume it is still
239 possible for a malicious party to eavesdrop on the channel.

240 We also assume there is a channel for communication between the OSD and the security
241 manager. Looking at the bandwidth and latency requirements of the various channels, the
242 channel between the security manager and the object store has the least stringent requirements.
243 This channel is used only for a periodic key exchange[7] and other administrative security
244 operations (see chapter 7). We believe that the performance of this channel is not an issue.

245 The channel between the security manager and host has medium network requirements, since it
246 is used for a message exchange for each unique credential required by the client. In some
247 configurations this could become a performance issue, since it is expected this channel be
248 encrypted.

249 The channel between the host and the OSD has the most stringent bandwidth requirements as
250 every request to the OSD flows on this channel. Because of the heavy traffic on this channel, it
251 is not reasonable to assume that by default this channel is encrypted.

---

[4] The host may request a broader set of rights than what is required for the operation it currently wishes to perform.

[5] It is conceivable that an authentication is not required, *e.g*., an object with world-wide read permission.

[6] When caching of credentials is possible, some verification steps can be omitted.

[7] The protocol does not specify this period, but we believe tens of minutes or longer would be reasonable.

252 **1.1.4  Layered Approach to Protocol Definition**

253 We take a layered approach to defining the protocol for object store security.  This allows an
254 implementation to provide only the desired level(s) of (internal) security and to surface the
255 various layers in a consistent manner.[8]  We also want to ensure a consistent message exchange
256 between the elements of the system, regardless of what level of security is supported.

257 An object store implementation defines the levels of security it supports.  By enriching the
258 information included in a message a higher level of security can be internally provided, as
259 opposed to leveraging an external network security mechanism.

260 In taking this approach, we want to provide flexibility in choosing how to secure the transport,
261 either internal or external, while allowing an installation to pay only for the level of security
262 needed. This should enable a simplified solution in certain glass-house environments (where no
263 network attacks are expected). It also should enable leveraging existing infrastructure for
264 network security and privacy while avoiding the cost of duplicate mechanisms. At the same
265 time, we must define a mechanism, which an object store can optionally implement to provide
266 network security as part of protocol for use where no secure transport exists.

267 ## *1.2  Levels of Security*

268 We consider several different levels of security that an OSD could provide.  In the first version
269 of the protocol we only directly provide the first three of these security levels.   Privacy can be
270 provided through external mechanisms, *e.g.,* running the protocol on an encrypted channel.
271 The levels are incremental and support all the protections of the level below them.

272 The particular level of security to be used for accessing a set of objects will be defined using a
273 mechanism not specified by the initial version of the standard.

274 **1.2.1  No Security**

275 In the no-security level, the same message structure will be used.  However, when an object
276 store is running with no security, the host must place zeros in the message related fields and the
277 object store must not examine these fields.

278 This is not considered a security level and its support is optional.

279 **1.2.2  Level 1 – Integrity of Capability (Access Control Security)**

280 Access Control Security is the common component to all the levels.

281 Integrity of capabilities by itself is most useful when the channel between the object store and
282 client is externally secured.  In this case, *e.g.,* where we have an authenticated IPSec channel,
283 we still need a mechanism that prevents a host from forging or otherwise modifying a credential
284 and/or replaying a credential over a different authenticated channel.  In addition, we need to
285 verify that the host rightfully possesses the credential it is presenting. Without a secure network,

---

[8] An implementation need not be layered

10

286 using only integrity of capability leaves an installation susceptible to certain network attacks,
287 *e.g.,* man-in-the-middle, replay, *etc.*

288 Support for this security level is optional. However, its functionality must be supported in
289 conjunction with all other security levels.

### 1.2.3 Level 2 – Integrity of Command and Arguments

291 Integrity of command and arguments is most useful when the channel between the object store
292 and the client is not externally secured and where providing integrity (hashes) for both
293 commands and data would be too expensive.

294 With integrity of arguments, malicious hosts cannot replay command parameters, even when
295 running on unsecured networks, but they can use network attacks on the data portion of the
296 messages exchanged between the client and the OSD.

297 Integrity of arguments prevents a malicious host from accessing a portion of object which was
298 not accessed by some client with a valid credential for the object, or changing a read operation
299 into a write, but it does not prevent a malicious host from modifying the data read from or
300 written to the object.

301 Support for this security level is optional. If supported, it must be supported in conjunction with
302 the functionality of integrity of capability.

### 1.2.4 Level 3 – Integrity of Data (Access Control and Internal End-To-End Security)

305 We assume that integrity of the data includes integrity of the command, *i.e.,* there is no point in
306 protecting the data if the command parameters describing to which object the datum belongs is
307 not also protected. This level of security provides security similar to integrity of capability when
308 the channel between the object store and the client is authenticated. The exact comparison
309 between the two depends on the level of network security that is provided by the external
310 security mechanism. The difference is that this level of the security internally secures the network
311 as an integral part of the object store protocol, thereby defining an end-to-end solution at the
312 storage layer as opposed to building upon pre-existing mechanisms for secure network
313 channels.

314 Support for this security level is optional. If supported, it must be supported in conjunction with
315 the functionality of the two prior levels.

### 1.2.5 Privacy

317 Providing privacy, *i.e.,* encryption, to the command and data, either in flight or at rest is beyond
318 the scope of the current proposal. This includes Levels 4 and 5. Note, there is nothing in this
319 proposal that precludes building upon external mechanisms for encryption.

### 1.2.6  Summary of Security Levels

All of the security levels are summarized in the table below. The table shows each level on its own as well as each layer when combined with a network security mechanism (such as IPSec) providing integrity and (separately) encryption as well as integrity.

| | | w/o a secure network | w/ a secure network (integrity) | W/ a secure network (encryption) |
|---|---|---|---|---|
| None | No Security | No security | Network-level integrity | Network-level privacy |
| Level 1 | Access Security | End-to-end verification of credentials | + Protection from network attacks | + Protection from network snooping |
| Level 2 | + Command Integrity | Protection from mistakes | + Protection from network attacks (some duplicated work) | + Protection from network snooping |
| Level 3 | + Data Integrity | End-to-end verification of requests | Duplicated work | + Protection from network snooping |
| Level 4 | + Command Privacy | Protection from traffic analysis on commands | Duplicated work | + Protection from snooping of data |
| Level 5 | + Data Privacy | End-to-end protection from snooping of data | Duplicated work | Duplicated work |
| Level 6 | + Data Integrity at Rest | Protection from modification of data on physical attack | Duplicated work | Duplicated work |
| Level 7 | + Data Privacy at Rest | Protection from leaking of data on physical attack | Duplicated work | Duplicated work |

## *1.3  Requirements Summary*

We have defined a set of requirements for the OSD security model; these requirements attempt to address a range of target platforms for implementing OSD.

On the one hand we believe it is important to enable efficient implementations of the object storage interface in storage controllers; such storage controllers are relatively resource rich, and it is reasonable to envision them containing support for standard network security, *e.g.,* hardware support for IPSec. We wish to be able to use an existing network security infrastructure (when practical) to take advantage of the development and design effort, as well as the administrative and support tools developed for such an infrastructure, *i.e.,* we do not want to (needlessly) reinvent the wheel.

On the other hand there is a requirement to enable efficient implementation in low-end storage devices. These devices are resource poor and the developers of these devices do not want to add additional hardware without a clear justification. These devices will not always support standard network security and in such environments it is necessary to provide end-to-end security against attacks without depending on an external mechanism to secure the network.

We have defined the following set of requirements that must be met by the OSD security model. We distinguish in defining these requirements between access control security (security which is

342 directly tied to the semantics of object storage) and network security (security which is related
343 primarily to network protocols and could be handled separately from the semantics of OSD).
344 The requirements we define are:

345 • Must prevent against attacks on individual objects. Such attacks include both intentional and
346   inadvertent access to an object in a way not authorized by the security manager.  In
347   particular, we must address malicious hosts forging or modifying a credential, a host stealing
348   a credential from the channel between the object store and client,[9] *etc*.

349 • Must enable protection against attacks on the network such as man-in-the-middle (*e.g.,* a
350   computer posing as an object store), replay, *etc.*

351 • Must provide a stand-alone solution that works in the event there is no existing network
352   security infrastructure or for whatever reasons the implementer desires not to use an
353   externally secured network.

354 • Must provide a solution that can use an existing standard network security infrastructure.

355 • Must not duplicate the cost of security, where it can be avoided. *e.g.,* if the host is running
356   over a secure network with Level 1, it should not incur a higher overhead than a host
357   running over a non-secure network with Level 3.

358 • Must allow low cost implementation of the critical path.

359 • Must be simple. In particular, we should use the same structures and same message flow
360   across all the protocol levels.

361 • Should allow efficient implementation on existing network transports.

## 1.4  Limitations in the Proposed Version of Object Store Protocol

364 The version of the protocol defined in the following sections of this document is a first step
365 towards OSD security. As such, it has the following limitations:

366 • It does not internally support privacy on the channel between the object store and the client

367 • It does not support privacy for the data at rest

368 • It requires a communication channel between the object store and the security manager.
369   This channel must be capable of carrying authenticated and encrypted messages.

370 • Ability to define capabilities that apply to multiple objects where the object to which a
371   capability applies is defined by a predicate on the object's attributes.  Note this is not the
372   same as commands which apply to multiple objects.

373 • Ability to define a capability which applies to only a portion of an object or to only certain
374   object attributes.

375 • It does not provide a means of determining from the object store what security level should
376   be used.

---

[9] As stated above, the assumption is that the channel between the host and the OSD is not encrypted, and
thus it is possible for a malicious host to eavesdrop on this channel.

# 2 Structure of Credentials and Basic Message Flow

## 2.1 Introduction

To enforce legitimate use of capabilities, the client receives from the security manager (over a secure channel) both the *capability (CAP_Args) and* some associated secret information, a *capability key (CAP_Key)*. Together the capability and capability key are the credential. The client sends a capability to the object store as part of each request. The client uses the capability key to compute a *validation tag*, which it appends to each request. The structure of this validation tag depends upon whether an existing network security infrastructure is being used, or whether the network security is provided internally by the protocol. Among other semantics depending upon the security level, the validation tag ensures the capability has not been modified. Using the protocol appropriate for the security level, the object store validates the validation tag and checks whether the operation requested by the command is indeed permissible. Note that the object store does not need to authenticate the client or to have a notion of "client identity".

## 2.2 Cryptographic Building Blocks

The cryptographic primitive that is used throughout this protocol is a keyed message authentication code. The protocol uses an HMAC-SHA1 [9][10] whose output is 160 bits long. When applicable, the final output of 160 bits is truncated into 96 bits. The HMAC-SHA1 key is 160 bits long. The means by which the HMAC-SHA1 key is generated is not specified by the protocol. Later versions of this protocol may allow an object store to specify alternate cryptographic primitives (see section 8.8).

## 2.3 Key Management Overview

The credential is based on a secret key that is shared between the object store and the security manager. For each object store $s_j$, $K_{secret\ key\_j}$ is an authentication key shared between $s_j$ and the security manager. For clarity, when concentrating on a specific object store we omit the index $j$ where no ambiguity arises. In particular, $K_{secret\ key}$ is a 160-bits long SHA1 key. More accurately, there is a hierarchy of keys shared between the object store and the security manager.

This protocol exchanges a secret key between the object store and the security manager:

1. The security manager sends a secret key to the object store along with the key's version number.
2. The object store updates its key, removes any cached credentials established with the previous key, and acknowledges receipt of key.

In chapter 8, we elaborate on the hierarchy of keys and the protocol for exchanging keys.

In a later version of the protocol we may define a mechanism for piggybacking the exchanges of keys over the client-object store channel without requiring a separate channel for the

413 communication between the object store and security manager.  As we describe below, since a
414 channel between the object store and security manager is needed for other reasons, we take
415 advantage of this channel for the key exchange.

### 2.3.1  Maintaining two valid keys $K_{secret\ key}$ simultaneously

417 A key refresh event between the object store and the security manager invalidates all credentials
418 at once.  This results in heavy communication traffic between all clients and the security
419 manager; moreover, all new credentials must be explicitly validated (via MAC calculation)
420 before being cached. This phenomenon may cause undesired performance degradation after the
421 key refresh. To mediate this effect, we allow an object store to declare the last two (or more
422 generally $n$) refreshed versions of $K_{secret\ key}$ as valid, instead of just the latest one. As a result,
423 the process of validating a credential requires a *key_version* field in the credential to enable the
424 object store to know which key to use in validating the credential.

425 The number of key versions used is configured between the OSD and the security manager.
426 The OSD implementation can specify the maximum number of key versions it supports; one is a
427 legal value.  The maximum number of key versions supported by the protocol is 16.

### 2.3.2  Partitions

429 An object store is divided into multiple partitions, each of which carries its own keys for security
430 purposes.  Instead of having a separate secret key for each object store $s_j$, there is a distinct
431 secret key for each two-tuple of object store $s_j$ and partition $p_k$.

432 All commands other than key exchange commands (see chapter 8) come with credentials which
433 are protected by the key associated with a specific partition.  For most commands, e.g., those
434 that operate on a specific object, the partition used is the partition containing the object being
435 operated upon.   For those commands which operate at the level of an entire object store, *e.g.,*
436 the commands for formatting the object store or creating/removing partitions, we use the keys
437 associated with partition zero.  Since the commands that operate at the level of the object store
438 and not at the level of individual objects are by their nature very powerful, we want to limit the
439 use of the keys associated with credentials used to execute these commands.  We thus define
440 that partition zero should not contain user objects; in addition, to solve the problem of
441 bootstrapping, an object store must always contain a partition zero (e.g., to allow formatting the
442 object store).  We note that a realization of the object store standard can define an identity
443 between the root object and partition zero.

## 2.4  Capability Argument and Capability Key

445 Define:

446 • *Type* of the credential (4 bits), which must currently all, be zero.  This is intended to allow
447     future extension to different types of credentials.

448 • *MAC Function* is a four bit field indicating the cryptographic primitive used to construct the
449     credential.  In the initial version of the protocol, the value of this field must be zero and the
450     HMAC SHA-1 must be used (see section 2.2)

- *Partion ID* is the identity of the partition for which this capability is being generated. Note we do not include the object store ID in the capability under the assumption that it is passed on all commands as part of the addressing.

- *Capability Nonce* to be an *l*-bits nonce (*l*=128) chosen uniquely by the security manager for each credential. The *nonce* may be a counter. We do not specify the means of generating this nonce, leaving the mechanism up to the implementer of the security manager. The role of this nonce is twofold: *1)* to ensure that every credential generated by the security manager is unique; this prevents a host from masquerading as an OSD to another host, which would be possible if both hosts received the same exact credentials and *2)* to serve as an audit field for allowing management applications to track the client which received a capability.

  This nonce has the following structure

  - *Audit tag* is a 32 bit value which the security manager uses in an implementation defined way to associate a credential with the client to which it granted the credential. The correctness of the system will not be dependent upon the value the security manager places in this *audit tag*. However, the overall performance and usability of the system can be improved if this field is used as a audit tag. This field can be used for purposes of auditing and report generation. It can also be used by an object store to better manage nonces in level 2 and level 3 of the protocol.

  - *Random bits* a 96 bit value which must be unique across all credentials with the same audit tag and values for the other fields.

- *Rights string* specifies the rights and object(s) to which they apply. At this point we propose the following structure for the rights string:[10]

  - *Type* – the implementation of the rights string; this is four bits with the following values

    - 0 – a specific object and set of operations is specified
    - 1-15 – reserved

  - *Operations* – a bitmap with one bit per OSD command; this bitmap should contain additional reserved bits for potential extension, without requiring a change in the size of credentials.

  - If the type == 0, then the following additional field is defined

    - *Object* – the local ID of the object to which this command applies.[11]

- *Object Version Tag* – a *k*-bits value (*k*=32) that is maintained as an attribute for each object. It is used to invalidate credentials, which have been issued earlier for the same object. If the security manager wishes to invalidate all credentials it had previously generated for an object, it modifies the value of the attribute associated with the object (see section 7.1); the new value should never have been previously used in a credential for this

---

[10] The size of the rights string is the sum of the sizes of its component fields with any necessary padding.
[11] The space required to encode the local ID will be used for pattern matching on attributes for future types of credentials to be defined.

488      object ID.  To allow resumed access to the object, the security manager should use this
489      new value in future credentials it generates for this object .

490    •   *Creation Time* – the time the object was created provided as an attribute by the object
491      store.  If object IDs are reused, then two creates for an object in a given partition which use
492      the same ID must have different values for the create time.  Note, it is clearly acceptable for
493      this value to be unique for every object created in an object store.  The size and resolution
494      of this value will be as defined for the creation time attribute of objects.

495    •   *Key_version* – a four bit index indicating the key version of $K_{secret\ key}$. The key version is set
496      at every key refresh between the object store and the security manager.  See also section
497      2.3.1 and chapter 8.

498    •   *Expiry Time* – a 48 bit field giving the time the credential expires in milliseconds since
499      January 1, 1970.  The security manager should generate this time.  By using an expiry time
500      we allow the security manager to give different lifetimes to different credentials.  We assume
501      a weakly synchronized clock between the security manager and the object store. No
502      assumptions are made on the client's clocks.  The OSD should not accept a capability with
503      an expiry time in the past.

504 The credential *C* that the security manager issues for a client is comprised of two components, a
505 "public token" *CAP_Args* and a "secret extra information" *CAP_Key*.

506     *CAP_Args* º *[rights string, Key_version, Nonce, Object Version Tag, creation time,*
507       *expiry time, partition ID, object store ID]*
508     *CAP_Key* º *MAC_K $_{secret\ key}$ (CAP_Args)*

509 *CAP_Key* is the 160-bits long output of the HMAC-SHA1 computation on CAP_Args along
510 with the implicit parameters of the object store ID and partition ID. Note that CAP_Key cannot
511 be truncated (to 96 bits) as it is used later in the protocol as a key to another MAC
512 computation.  It is the host's responsibility to keep CAP_Key secret; if CAP_Key is
513 compromised, than it is possible for an adversary to issue requests using the capability if it
514 determines CAP_Args, which are passed on the wire between the OSD and host in the clear.

515 We note that not all of the fields in the CAP_Args need to be passed explicitly on the wire.  In
516 particular, since the object store knows the creation time and desired version tag for each
517 object, it is not necessary to pass these values.  Instead, given the object ID, the object store
518 can determine which object version tag and creation time to use in calculating the CAP_Key.  If
519 the host had a credential created using different values for these fields, a MAC calculation
520 would fail and the command would be rejected.

521 In a similar vein, the partition ID and object store ID do not need to be passed as part of the
522 capability for each command.  This is because these fields are part of addressability and will

523 need to be passed as part of the basic command (even if running with no security). In other
524 words, there is no need to pass partition ID and object store ID twice.

525 The precise treatment of the object version tag, creation time, partition ID and object store ID
526 will be defined by each realization of a concrete object store standard, however, we
527 recommend that they not be passed as part of a capability on each command.

528 Since the credential includes information which is stored as attributes for the objects (namely the
529 creation time and version tag), we may have a problem of bootstrapping, in particularly if the
530 security manager does not have this information in its memory. How does the security manager
531 generate a credential to read these attributes if it does not know these attributes? In addition, in
532 certain usage scenarios, *e.g.,* all object IDs assigned by an external cataloging entity, the use of
533 the creation time may require additional message exchanges and provide no benefit.

534 To address this, if a credential generated by the security manager uses zero for the version
535 tag/creation time, then when calculating the CAP_Args the object store should not take into
536 account the actual value of the respective attribute associated with the object but rather will use
537 zero (of the appropriate number of bits). When used with the version tag, this essentially
538 creates a credential which cannot be invalidated (other than by a key exchange which invalidates
539 all credentials for the partition generated with the same working key). Note that if a realization
540 of this work as a concrete standard does not pass the complete values of version tag or creation
541 time with each command (see above), it must pass an indication of whether or not these fields
542 should be treated as zero.

543 To delegate a credential *C* to another host, a host must transfer both the CAP_Args and
544 CAP_Key. While beyond the scope of this protocol, to ensure security, such delegation should
545 be done over an encrypted channel.

## 2.5  Anonymous Object Creation

547 To support creating an object where the OSD provides the object ID, the security manager
548 should generate a capability in which the object ID embedded in the rights string is zero and the
549 only right specified in the operations bitmap is object creation. The OSD must not allow such a
550 capability to be used more than once. To minimize the memory requirements the OSD must
551 dedicate to ensuring that such capabilities are used at most once, it is strongly recommended
552 that the security manager construct such capabilities with expiry times very close to the current
553 time.

## 2.6  Message Flow

555 Prior to sending an object store command to a target, the client must request the credential from
556 the security manager and in return the security manager sends back both the public part of the
557 credential, *CAP_Args*, as well as the private part, *CAP_Key*. *CAP_Key* should be sent to the
558 client over an authenticated and encrypted, channel to maintain its secrecy. To establish this
559 channel (and also to let the security manager identify the client), the client and the security
560 manager should authenticate each other in a preliminary step. The implementation of this channel
561 and its protocol are not part of the object store protocol.

562　When the client executes the actual object store command, the object store should validate:

563　　1.　The integrity of the public credential *CAP_Args*
564　　2.　That the public credential *CAP_Args* is used by a client that legitimately received it.
565　　3.　The integrity of the command itself (command and data), as required by the security
566　　　level.

567　For that, the client sends, along with the command, the public credential *CAP_Args* along with
568　a MAC-based validation tag, which is computed using *CAP_Key*.  Since *CAP_Key* can be
569　computed from *CAP_Args* and the secret shared between the security manager and the object
570　store, the validation tag is also computable by the object store.

571　The structure of the validation tag and its usage depends on the security level being used.
572　Section 3 describes the validation tag if we assume an external mechanism for the integrity of
573　data and command, namely an authenticated channel such as an IPSec authenticated channel. In
574　Sections 5and 6 no such external mechanism is assumed and therefore the validation tag as well
575　as its validation at the object store is more elaborate.

## 2.7  Credential Invalidation

577　As described above, the object store protocol provides two means for invalidating a credential.
578　By the use of object version tag in each credential, we can invalidate all of the outstanding
579　credentials for an object.  By a key exchange between the security manager and the object
580　store we can invalidate all credentials a security manager had generated for a particular object
581　store partition.  Note, by explicit decision, we have decided not to support an efficient means of
582　externally invalidating all of the credentials given to a particular host by the security manager (but
583　see section 4.4).

## 2.8  Security Related Error Status

585　The following error responses related to security can be returned by the OSD to the host.
586　Some of these responses are limited to specific security levels as indicated:

587　• *NOT_SUPPORTED_CREDENTIAL_TYPE* – the type of the credential is not supported
588　　by the object store.

589　• *CAPABILITY_MISMATCH* – the requested operation is not allowed by the rights string

590　• *INVALID_MAC* – the message authentication code (MAC) included in the request is not
591　　consistent with the credential included in the message; in other words, the *CAP_Key*
592　　calculated based upon the credential cannot be used to compute the same MAC as the one
593　　included in the message.  In the event that the MAC is invalid either this error or
594　　INVALID_KEY must be returned (regardless of other errors detected).

595　• *INVALID_VERSION* – the request includes a credential with a object version tag which is
596　　no longer being accepted

597　• *INVALID_KEY* – the key as indicated by key_version in the credential is no longer valid;
598　　the host must retrieve a new credential from the security manager prior to retrying the

599      operation. An OSD implementation does not need to be able to distinguish this situation

600      from the situation reported by *INVALID_MAC;* in which case it should report

601      *INVALID_MAC.*

602 • *EXPIRED_CREDENTIAL* – based upon the expiry time, the credential has expired.

603 • *INVALID_NONCE* – the nonce does not contain a valid timestamp; a recommended time

604      stamp will be returned with this code. This may only be returned for level 2 or level 3.

605 • *NONCE_NOT_UNIQUE* – a message with this per request nonce has previously been

606      seen by this object store. This may only be returned for level 2 or level 3.

607 • *CAPABILITY_BLOCKED* **–** The capability was blocked, *e.g.,* based on the capability

608      audit tag. Note that the reason for the "blockade" is not given. This may only be returned

609      for level 2 or level 3.

610 In addition to these error responses which are specific to security, the following additional

611 errors, which we are not specifically related to security, can be returned:

612 • *INSUFFICIENT_RESOURCES* – a temporary condition exists which does not allow

613      processing this request due to a lack of resources

614 • *INVALID_MESSAGE_STRUCTURE* – the structure of the message is not syntactically

615      valid.

# 3  Level 1 – Integrity of Capabilities

This security level is useful in two scenarios: *1)* where no network attacks are expected to take place (such as a 'glass house' scenario) and *2)* where an authenticated channel between the host and the OSD is assumed. The mechanism for establishing this channel is beyond the scope of the OSD protocol.

## *3.1  Level 1 Security with Authenticated Channel*

For Level 1 security with an authenticated channel, the channel provides integrity of messages as well as an anti-replay mechanism (for the given channel).  The OSD-specific protocol prevents copying messages from one channel to another by tying the message to the channel via a *validation tag*; this tag is computed as $MAC_{CAP\_key}(ChannelID)$, where *ChannelID* identifies the communication channel.  Given that the OSD knows the channel on which a request was received, the OSD can validate that the $MAC_{CAP\_key}(ChannelID)$ included in a message is for the *ChannelID* associated with the channel on which the message was received.   The same validation tag can be used with all requests based upon a given credential.

We do not need this validation tag on OSD responses since *1)* the authenticated channel ensures the host any responses it receives are received from the intended OSD and *2)* we trust the OSDs to not copy messages between channels (see section 1.1.2).

The *ChannelID* is  a name for the channel that is unique to this channel between the client and the object store and is known to both ends.  The size of the *ChannelID* is transport dependent.  The lifetime of the *ChannelID* is no greater than the lifetime of the channel;[14] the lifetime of the *ChannelID* is independent of the lifetime of the key $K_{secret\ key.}$ See section 3.3 for a more precise definition of the assumptions on the channel and *ChannelID*.

Most likely, a value that can be used as a *ChannelID* already exists and was created at the time the channel was established. Otherwise, it requires another message exchange between the client and target, where:

1.  Client requests 'open security window' with the object store.
2.  Object store responds with a randomly chosen *m*-bit channel name *ChannelID*.[15]

At this point we do not architect such a flow to explicitly have the object store provide the *ChannelID* but rather we assume the channel provides the *ChannelID*.

Below is the protocol flow of messages along with a table that explains the messages and their corresponding arguments. Note that the *OpenWindow* message is *not* needed for every *ReqCap* message. Furthermore, it may not be needed at all if a *ChannelID* is already exchanged.

---

[14] Note it would be permissible to change the *ChannelID* for an existing channel; this would invalidate cached credentials.

[15] Analogously, a 'close security window' clears knowledge of the session at the object store.

**Sec Mgr**   **Host**   Open Window   **ObS**

Response: Channel ID

Req. Cap: Obj ID, ObS ID

Response: Cap Args, Cap Key

Read: Command, Cap Args, Channel ID, $MAC_{CAP\_Key}(0, ChannelID)$

Response: Status, $MAC_{CAP\_Key}(0, ChannelID)$

Read: ... $MAC_{CAP\_Key}(0, ChannelID)$

Req. Cap

Response   Read

Response

649
650

651   **Figure 2.   Flow of Messages for Level 1 Security.   The establishment of the authenticated channel**
652   **is not shown.   The *OpenWindow* exchange will not be required for most channels.**

653

| Message | Argument | | Explanation |
|---|---|---|---|
| | | | |
| **ReqCap** | ObjID | | Object identifier |
| | Partition ID | | Partition ID |
| | ObSID | | Object Store identifier |
| | | | |
| **ReqReturn** | Version | | |
| | Rights | **CAP-Args** | |
| | Expiration | | |
| | Partition ID | | |
| | Creation | | |
| | Capability Nonce | | |
| | CAP_Key | | $MAC_{secret\_key}(CAP\_Arguments)$ |
| | | | |
| **OpenWindow** | | | |
| | | | |
| **WindowReturn** | ChannelID | | |
| | | | |
| **ReadData** | ObjID | | |
| | Partition ID | | Partition ID |
| | ObSID | | Object Store identifier |
| | CAP-Args | | |

22

| | Offset | |
|---|---|---|
| | Length | |
| | Nonce | Ignored (all zeros) |
| | ReqMac | $MAC_{CAP\_Key}(ChannelID)$ |
| | | |
| **ReadReturn** | Status | |
| | RetMac[16] | $MAC_{CAP\_Key}(ChannelID)$ |

654

## 3.2 Level 1 - security without network security

656 As noted above, this level of security is also useful when no network attacks are expected to
657 take place. In the event no secure network infrastructure is used, level 1 security protects the
658 integrity of the capability. The protocol is identical to the one described above. However, since
659 the *ChannelID* is not in practice tied to a channel and there is no true means for tying a message
660 to the channel. An OSD implementation cannot, however, ignore the value of the validation tag
661 if level 1 is being used without a secure network infrastructure since the validation tag is also
662 used to validate the capability has not been modified. In this case, zero should be used as the
663 value of the *ChannelID*.

## 3.3 Assumptions on Network Infrastructure for End-to-End Security

666 We place the following requirements on the channel and *ChannelID* if we want to ensure and
667 end-to-end security solution using level 1 of OSD security:

668 • Within the lifetime of a key, $K_{secret\ key,}$ all channels established with a given object store from
669 any host must receive unique channel IDs.

670 • There must be a means for the host and OSD to associate a received message with the
671 *ChannelID* for the channel on which the message was received.

672 • Assuming that the channel provides the value of the *ChannelID*, this value must be non-
673 forgeable.

674 • The channel must be authenticated (although it may be anonymous) in the sense that it must
675 ensure both parties can be guaranteed all messages in a session come from the same party.

676 The channel must ensure message integrity, *i.e.,* non-modification of message contents by the
677 network.

## 3.4 Client-Object Store Message and Flow

679    *1.* Client sends a command to the object store, along with the public token *CAP_Args*
680       (defined above) and a 96-bits long validation tag $V = MAC_{CAP\_key}(ChannelID)$. *V* is
681       computed using HMAC-SHA1 on the *ChannelID,* truncated to 96 bits.

---

[16] As discussed above, this MAC is not necessary – is it only used for symmetry with level 2 and level 3
security

682    2. Verification at object store:

683        1. The validation tag $V$ equals $MAC_{CAP\_key}$ *(ChannelID)*, where *CAP_Key* is obtained
684            as $MAC_{K\_secret\;key}$ *(CAP_Args)*.
685        *2.* The rights string in the *CAP_Args* allows the requested operation
686        3. The *key_version* is current.
687        4. The capability's *Version Tag* is either zero or equal to the version tag attribute of the
688            object. [17]
689        5. The capability's *Creation Time* is either zero or equal to the creation time attribute of
690            the object.

691        If any of the checks fails, the request is denied. If checks (a) and (c) pass, the object
692        store may cache the token *CAP-Args* associated with channel *ChannelID*.

693  An object store implementation may cache the validation calculations. In particular, if
694  *CAP_Args* has ever been presented to the object store on this channel within the lifetime of
695  current *ChannelID*, the request may be granted without re-validation (*i.e., without redoing*
696  *step 1)*. The authenticated channel assures that another client is not replaying *CAP_Args* on
697  this channel, rather it is currently presented by the same entity that presented it in the past, and
698  hence a re-validation is not necessary.



699

700               **Figure 3.    Flow for Level 1 Security**

---

[17] This check can be implicit in checking the validation tag as the object version tag is part of the CAP_Args and if the version tag is incorrect, the object store will not be using the same CAP_Key as the host. This comment applies as well to the creation time. It also applies to the other security levels. It is not repeated.

701 ### *3.5 Performance Considerations*

702 For level 1 security, we have the following performance considerations:

703 • A client does not need to request a new credential on every command; rather the client can
704 reuse the *CAP_Args* and *CAP_Key* on multiple commands for the same object(s).

705 • The client does not need to recalculate the ReqMAC on each command; rather, this needs
706 to be calculated only once per credential.

707 • The object store does not need to recalculate X and Y on each exchange with a client.
708 Rather since we assume a secure channel, these values need to only be calculated the first
709 time object store sees a given capability.

# 4 Per Request Nonces for Level 2 and Level 3

710

711 Level 2 and Level 3 of the security protocol use Nonces included in each request to prevent
712 replay. The requirements for correctness of a nonce-based approach to preventing replay are
713 as follows:

714 - The object store must not accept the same nonce more than once.

715 - The object store must not accept a nonce that was rejected in the past.[18]

716 In addition, it is acceptable for an implementation to reject valid requests with unseen nonces if
717 necessary to ensure that the two basic requirements are met.

718 There are three main means of generating nonces:

719 - Random
720 - Session based
721 - Time based

722 We believe it is fairly easy to argue that the time-based protocol has better performance and
723 space requirements than either a session-based or random generation protocol if all entities in
724 the system are well-behaved. On the other hand, the time-based protocol can have extremely
725 large memory requirements or require frequent changes of the secret keys if enough clients in the
726 system are not well-behaved. In addition, assumptions on strong clock synchronization
727 between the clients and object store are problematic both from a practical and security
728 perspective.

729 The approach to nonces we define is a time-based approach modified to have only weak
730 dependencies upon client clocks and augments to minimize the impact of poorly behaved
731 entities.

## 4.1 Background

732

733 We define a system running level 2 or level 3 security as *well-behaved* if at any given time $t$, the
734 total number of *far-in-the-future messages* (messages with a nonce whose time is greater than
735 $t + d$) which have been sent to an object store from all clients, is less than $k$, for
736 implementation-defined values $d$ and $k$. In other words, a system is well-behaved if the number
737 of far-in-the-future messages, which an object store has received, is bounded. Similarly a
738 client running level 2 or level 3 security is defined as *well-behaved* if it does not send any
739 nonces for a time greater than $t + d$ where $t$ is the current time of the target object store. An
740 *ill-behaved client* and *ill-behaved system* have the obvious definitions. Malicious intent is not
741 required for a client to be *ill-behaved*. Also note that if there is malicious intent, the
742 maliciousness is not necessarily directly from the ill-behaved client; for instance, a malicious
743 time-server can causes clients to be ill-behaved.

744 In the worst case, with the time-based nonces, an object store implementation must ensure that
745 it has sufficient memory[19] to remember the nonce from each message it could receive in the

---

[18] This holds regardless of the reason the message was rejected

746 period between working key exchanges. This is to prevent messages from being replayed: the
747 implementation must also ensure that any message which has ever been rejected will never
748 become valid in the future. In other words, given an object store which can handle $n$ messages a
749 second and a key exchange every $e$ seconds, the object store needs to be able to remember $ne$
750 nonces.[20]   This, admittedly unlikely worst case, would occur if every message received was for
751 the end of the period in which the key was valid.   Note an alternative would be to allocate a
752 fixed amount of memory, much less than for $ne$ nonces and if this memory fills up, for the object
753 store to force a key exchange.   This leaves open a denial of service (DOS) attack in which all
754 existing capabilities  are invalidated.  The goal of our modifications to a pure-time-based nonce
755 protocol is to reduce the easy of this DOS attack

756 One way to mitigate the amount of memory required to handle ill-behaved systems[21] is to design
757 the messages in such a way that the object store would be able to reduce its memory
758 requirements by organizing the nonces into groups. If the far-in-the-future messages are limited
759 to a subset of the groups of nonces, the implementation can decide to reject nonces belonging to
760 the problematic groups, while continuing to accept other nonces.  Clearly, the efficiency of such
761 an approach depends on the accuracy of the grouping.  We leverage the audit tag field of the
762 nonce[22] in the *CAP_Args* for this purpose; see section 2.4.

## 4.2  Requirements

764 In addition to the general requirements listed above, we place the following requirements on the
765 protocol:

766 • The changes to allow better behavior in ill-behaved systems should incur no additional cost
767 in the case of a well-behaved system.

768 • An implementation that chooses so must be able to bound the amount of memory required
769 for correct behavior independent of the frequency in which the key is exchanged between
770 the object store and security manager (*i.e.,* safety with bounded memory), while sill ensuring
771 liveness for well-behaved clients in many scenarios where there are ill-behaved clients.

772 • We must allow freedom to the object store implementer to trade-off between
773 implementation complexity and overall system behavior in the event that clients are not well-
774 behaved.

## 4.3  Structure of the Per Command Nonce

776 When working with the time-based nonces, on each request, the host generates a nonce by
777 combining a 48-bit time representing the number of milliseconds since January 1, 1970.  The
778 nonce also includes a 48-bit random number.

---

[19] Clearly various compression techniques could be used; for example see [11].

[20] This is the number of nonces that must be remembered; the memory that is required is implementation dependent and may need to take into account compression techniques.

[21] Although there are still scenarios in which correct behavior entails either remembering all nonces or forcing a key exchange.

[22] Not to be confused with the per command nonce described in this section

## 779  *4.4  Use of Nonce for Anti Replay*

780   Define the *current interval* for an object store whose clock currently is at time $t$ as the period
781   of time beginning with $t-d_1$ and ending with $t+d_2$, where $d_1$ and $d_2$ are values determined by the
782   object store implementation.   The current interval defines the time-based nonces the object
783   store expects to receive from well-behaved clients.   The object store can accept any valid
784   request received in this time range.  To prevent replay, the object store must have sufficient
785   resources to remember all nonces seen in this range.  Messages received with nonces less than
786   $t-d_1$ do not need to be remembered.[23]   Define as *far-in-the-future* a nonce for a time greater
787   than $t+d_2$.  By definition, such nonces will only be sent by ill-behaved clients.

788   When the object store receives a request in level 2 or level 3 with a nonce in the current interval,
789   the object store must remember the nonce in a *current interval nonce list*.[24]  While the only
790   requirement from the protocol is that anti-replay be provided, the space allocated to the current
791   interval nonce list should be sufficient[25] to hold the number of nonces that can be received by the
792   object store during the time of the current interval, *i.e.,* a function of the size of a nonce (12
793   bytes) times the number of messages the object store can receive in time $d_1 + d_2$.

794   Note, the object store must remember the nonce even if the message fails verification of the
795   MAC.  This is required to prevent the following, replay-like attack.  Assume an adversary
796   hijacks a request to the object store, modifies the command portion of the request and forwards
797   the request to the object store.  The object store will send an *INVALID_MAC* error response
798   to the client.  The client may then decide to regenerate the request with a new nonce and MAC.
799   Assuming this request is executed, the adversary can now replay the original request.  We
800   should point out that a client should be suspicious of an *INVALID_MAC* response which does
801   not itself contain a valid MAC.

802   If the nonce in a request is for a time that is older than the current interval, the object store
803   rejects the request without further processing with an *INVALID_NONCE* error message.   The
804   *INVALID_NONCE* response includes the current time of the object store, allowing the client to
805   try again with a nonce that will fall in the current interval.  The object store does not need to
806   remember the nonce.  A well-behaved client will (logically) reset its clock to be that of the
807   object store for future messages it sends.

808   Finally, if the nonce in a request is a *far-in-the-future* nonce, the object store must remember
809   the nonce in the *far-in-the-future nonce list*.[26]  The object store implementation may reject the
810   command with an *INVALID_NONCE* status or it may decide to process the request as
811   described for messages received with a nonce in the current interval, as long as the nonce
812   uniqueness is guaranteed.[27]  If an *INVALID_NONCE* response is returned, as above, it will

---

[23] Note we assume that if the clock of an object store is set backwards, a key exchange with the security
manager will also take place.

[24] Note, the reference to a current internal nonce list is for explanatory purposes only; an implementation
may choose any mechanism to remember previously seen nonce as long as the basic requirements are met.

[25] After any compression techniques

[26] Note, the reference to a far-in-the-future nonce list is for explanatory purposes only; an implementation
may choose any mechanism to remember previously seen nonce as long as the basic requirements are met.

[27] But this does not enable the client to be informed that it should update its clock

813 include the current time of the object store and a well-behaved client will logically reset its clock
814 to be that of the object store.

815 We define the size of the *far-in-the-future nonce list* to be large enough to hold some number,
816 $k$, of nonces,[28] where $k$ is implementation dependent, not specified by the protocol, and may
817 vary at different times for a given implementation.   Clearly, a nonce can be removed from the
818 far-in-the-future nonce list when the nonce represents a time prior to the start of the current
819 interval.  If an implementation ensures the basic requirements, a nonce can be removed from the
820 far-in-the-future nonce list at other times.

821 To verify that a nonce has not previously been seen, the object store must look in both the
822 current and far-in-the-future nonce lists.[29]

823 If the object store receives more than $k$ far in the future nonces, *i.e.,* the object store has run out
824 of resources to remember far-in-the-future nonce, the object store implementation has several
825 options, as long as it guarantees the basic requirements of not accepting the same nonce more
826 than once and not accepting a nonce that was previously rejected.

827 One option, the "big hammer" option, is for the object store to refuse to accept any more
828 messages using the same working key which was used for the capabilities in the messages with
829 the far in the future nonces.   In this case, the object store may return an indication of
830 *INVALID_KEY* when it receives requests with this working key.  It is implementation
831 dependent as to how the security manager is notified that the working key needs updating.
832 Options include (but are neither limited to, nor required to include) having the security manager
833 poll the object store and having the client pass on an indication to the security manager.

834 The drawback of the "big hammer" option is that it invalidates all capabilities  whose
835 corresponding credential was created with the given working key.  In other words, all clients
836 which have capabilities for the given object store partition created with the same version of the
837 working key are impacted.

838 To mitigate the likelihood an implementation needs to resort to the big hammer, the
839 implementation can organize the far-in-the-future nonce list based upon the architected audit tag
840 that the security manager places in the credential.[30]  One option an implementation can choose is
841 to partition this nonce list based upon the audit tag.  For instance, if the object store receives
842 more than $c$ far-in-the-future nonces with a given audit tag created by the same working key,
843 the object store can refuse to receive additional requests with the given audit tag until the oldest
844 request in the far-in-the-future nonce list for this audit tag is older than the start of the current
845 interval.   If the object store is refusing to receive requests with a given audit tag or capability, it
846 should return *CAPABILITY_BLOCKED*.  For this to work, the object store must always
847 remember the $c$ newest far-in-the-future nonces received with a given audit tag.  In this case, the

---

[28] Again, this may be after compression

[29] The description of separate current and far-in-the-future nonces lists is for explanatory reasons only; an implementation that ensures the basic requirements need not have separate lists.

[30] The implementation may arrange the *far in the future set* in any manner, e.g., it according to the nonce hash value. However using audit tags is a reasonable choice as they identify the "source" of the attack.

848 object store only needs to "drop the hammer" if more than $k/c$ clients are not well behaved.
849 Other implementations are clearly possibly as long as they meet the base requirements.

850 We require that $c$ be a value that is visible to a client. Clients may send a batch of requests
851 without waiting for a response. In this case, a client needs to be able to determine how many
852 outstanding requests it can send to an object store without risking having the object store decide
853 it is ill-behaved and thus refusing to accept requests from it.

## 4.5  Host Protocol

855 To prevent replay of responses, hosts must maintain nonce lists in the same way the object store
856 supports nonce lists

## 4.6  Use of Time

858 The only requirement for the time used to determine nonce timestamps is that it be
859 monotonically increasing, although weakly synchronized clocks between the OSD and hosts will
860 avoid additional messages. This time must never go backwards without a key exchange. In
861 order to catch the time up to an external "real time", the OSD may choose to accelerate or
862 decelerate the passage of time until it has caught up or the "real time" has caught up. Any OSD
863 that is unsure of the time, or concerned about a time-based attack, may choose to expand the
864 size of its nonce lists as it sees fit. This may slow performance, but does not affect security.

## 4.7  Additional Attributes on Partition Object

866 To allow implementing a complete solution, an object store implementing level 2 or level 3
867 security, must define the following attributes on a partition object:

868 • *NUM_REQS_BEFORE_BAD* – the minimum number of requests which are far-in-the-
869 future which a client may send, prior to the object store determining that the client is ill-
870 behaved.[31] This guarantee will only hold if there are not too many clients sending
871 NUM_REQS_BEFORE_BAD at the same time. Note that one and zero are legal values.

872 • *WORKING_KEY_FROZEN(i)* – an array of $n=16$ Boolean attributes, where the $i$'th
873 attribute is true if an object store needs to "drop the hammer" and refuse any credentials
874 created with the $i$'th version of the working key (as indicated in the key_version) field of the
875 credential. An OSD sets bit $i$ when it, of its own initiative, invalidates working key $i$ and an
876 OSD unsets bit $i$ when it receives and accepts a key management command that defines a
877 new value for working key $i$.

878 • *OLDEST_VALID_NONCE* – the minimum number of milliseconds older than the object
879 store's current time a nonce that is received will be considered valid; this attribute maps to
880 the value $d_1$ defined above. Zero is a legal value implying the absence of information.

---

[31] This is the "maximum" number of requests that a client trying to be well-behaved can issue without
receiving a response from any, and be confident that the OSD will not invalidate the associated working key
in the case that its nonces are in fact far-in-the-future relative to the OSD clock.

881 • *NEWEST_VALID_NONCE* – the minimum number of milliseconds newer than the object
882    store's current time a nonce that is received will be considered valid; note an object store
883    implementation may decide to treat as valid nonces that are even newer than this.  This
884    attribute maps to the value $d_2$ defined above. Zero is a legal value implying the absence of
885    information.

# 5  Level 2 – Integrity of Arguments

887 This security level does not make any assumption about the security of the underlying network
888 and internally provides end-to-end protection for the arguments at the level of the OSD
889 protocol.

890 The Host makes a request for a capability to the Manager and the manager returns the
891 credential composed of *CAP_Args* as well as the *CAP_Key*.

892 The Host then presents the command, including the *CAP_Args* and the *Cmd_Args*, along with
893 the *ReqMac* to the OSD. The *ReqMac* is a MAC using the *CAP_Key* of the *Cmd_Args* and a
894 *nonce* constructed as described in the prior chapter. This ensures that the *Cmd_Args* are not
895 modified in transit. The *Nonce* ensures that the command is not being replayed from some point
896 in the past.

897 The OSD then verifies that:

898   1. the *Nonce* is fresh, *i.e.,* it has not been seen before
899   2. the *Cmd_Args* are compatible with the *CAP_Args* (i.e., the rights string permits the
900       operation)
901   3. the Version Tag and Creation Time are valid
902   4. *CapY* matches *ReqMac* as sent by the host

903 Where *CapY* is calculated using

904   *CapX*  = a MAC computed using the *secret_key* on the *CAP_Args* (this is the *CAP_Key)*
905   *CapY*  = a MAC using CapX on the Cmd_Args and the Nonce

906 If any of these conditions cannot be verified, the request is rejected and no further command
907 processing is performed other than processing related to the nonce as described above.
908 Nonce related failures are handled as described in the prior section.  Other failures are reported
909 with a *Status* as described in Section 2.8.

910 In all cases:

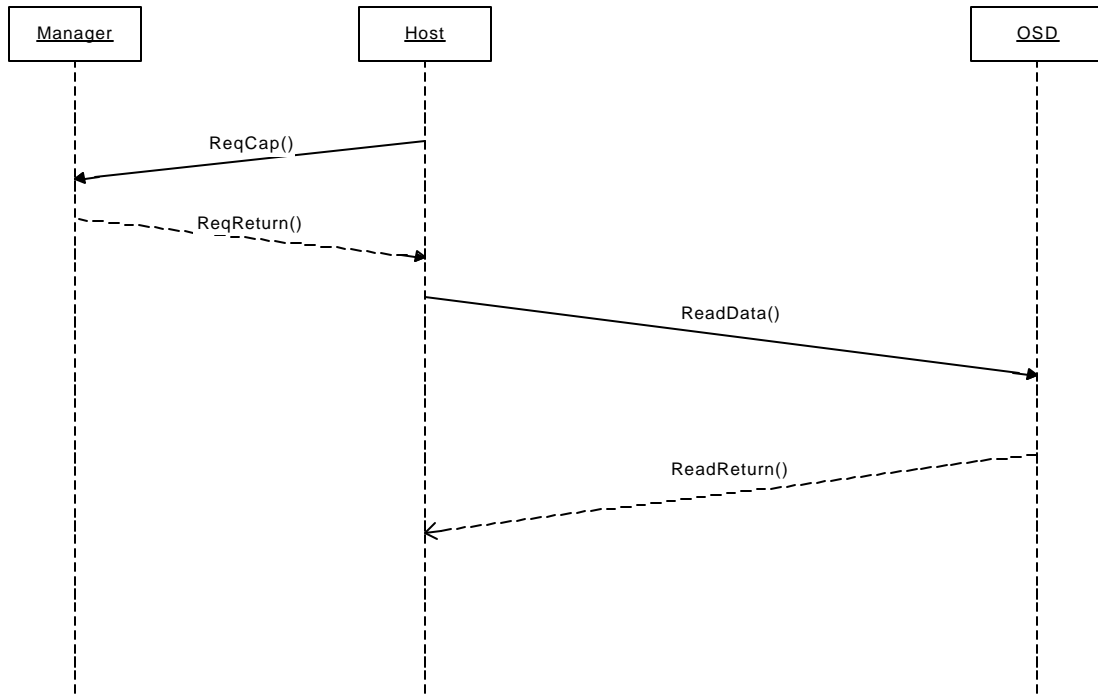911   1. A *RetMac* is computed using *CapX* on the *Status* and the *Nonce* (from the original
912       request) to allow the host to verify the response

913 Note we can safely apply this MAC to all messages, including with a status of *INVALID_MAC*
914 without becoming susceptible to a black box attack due to the properties of HMAC we are
915 using.[32]  See section 2.2.

---

[32] The Message Authentication Code (MAC) has the Computation - Resistance property [1], namely, given
text -MAC pairs $(x_i, h_k(x_i))$, it is computationally infeasible to compute any other text -MAC pair $(x_j, h_k(x_j))$ for any new input $x_j$ [3] .

Manager    Host    OSD

ReqCap()

ReqReturn()

ReadData()

ReadReturn()

916
917
918

| Message | Arguments | Explanation |
|---|---|---|
|  |  |  |
| **ReqCap** | ObjID | Object identifier |
|  | Partition ID | Partition identifier |
|  | ObSID | Object Store identifier |
|  |  |  |
| **ReqReturn** | Version |  |
|  | Rights       **CAP-Args** |  |
|  | Expiration |  |
|  | Partition ID |  |
|  | Creation |  |
|  | Capability Nonce | Capability nonce |
|  | CAP_Key | $MAC_{secret\_key}(CAP\_Arguments)$ |
|  |  |  |
| **ReadData** | ObjID       **CmdArguments** |  |
|  | Partition |  |
|  | Offset |  |
|  | Length |  |
|  | Nonce | Per command nonce |
|  | CapArgumentsCAP_Args |  |
|  | ReqMac | $MAC_{CapKeyCAP\_Key}(ObjID, Offset, Length, Nonce)$ |
|  |  |  |
| **ReadReturn** | Status | return code from the request, success or failure |
|  | RetMac | $MAC_{CapKeyCAP\_Key}(Status, Nonce)$ |

919

33

## 5.1 Performance Considerations

920

921 For level 2 security, we have the following performance considerations:

922 • A client does not need to request a new credential on every command; rather the client can
923 reuse the *CAP_Args* and *CAP_Key* on multiple commands for the same object(s).

924 • The object store does not need to recalculate *CapX* on each exchange with a client.

# 6   Level 3 – Integrity of Arguments and Data

This security level does not make any assumption about the security of the underlying network and provides end-to-end protection at the level of the OSD protocol. In addition to the protection of Level 2, this level also includes integrity checking of the data portion of the command.

The Host makes a request for a capability to the Security Manager and the host returns a credential, namely the *CAP_Args* as well as the *CAP_Key*.   As in level 2, the Host then presents the command, including the *CAP_Args* and the *Cmd_Args*, along with the *ReqMac* to the OSD. The *ReqMac* is a MAC using the *CAP_Key* of the *Cmd_Args* and the *Nonce*.

In addition, the *DataMac* is a MAC using the *CAP_Key* of the *Data* and the *Nonce*. On a WRITE command, the *DataMac* is computed by the Host, on a READ it is calculated by the OSD.

The OSD then verifies that:

1. the *Nonce* is fresh, it has not been seen before
2. the *Cmd_Args* are compatible with the *CAP_Args* (i.e., the rights string permits the operation)
3. the Version Tag and Creation Time are valid
4. *CapY* matches *ReqMac* as sent by the host

Where *CapY* is calculated using

*CapX*  = a MAC computed using the *secret_key* on the *CAP_Args* (this is the *CAP_Key)*
*CapY*  = a MAC using CapX on the Cmd_Args and the Nonce

In addition for Writes the OSD verifies

5.   (WRITE) *DataZ* matches *DataMac* as sent by the Host

Where

*DataZ*  =  a MAC computed using *CAP_Key* on the *Data* and *Nonce*

If any of these conditions cannot be verified, the request is rejected and no further command processing is performed other than processing related to the nonce as described above. Nonce related failures are handled as described in the section 4.  Other failures are reported with a *Status* as described in Section 2.8.
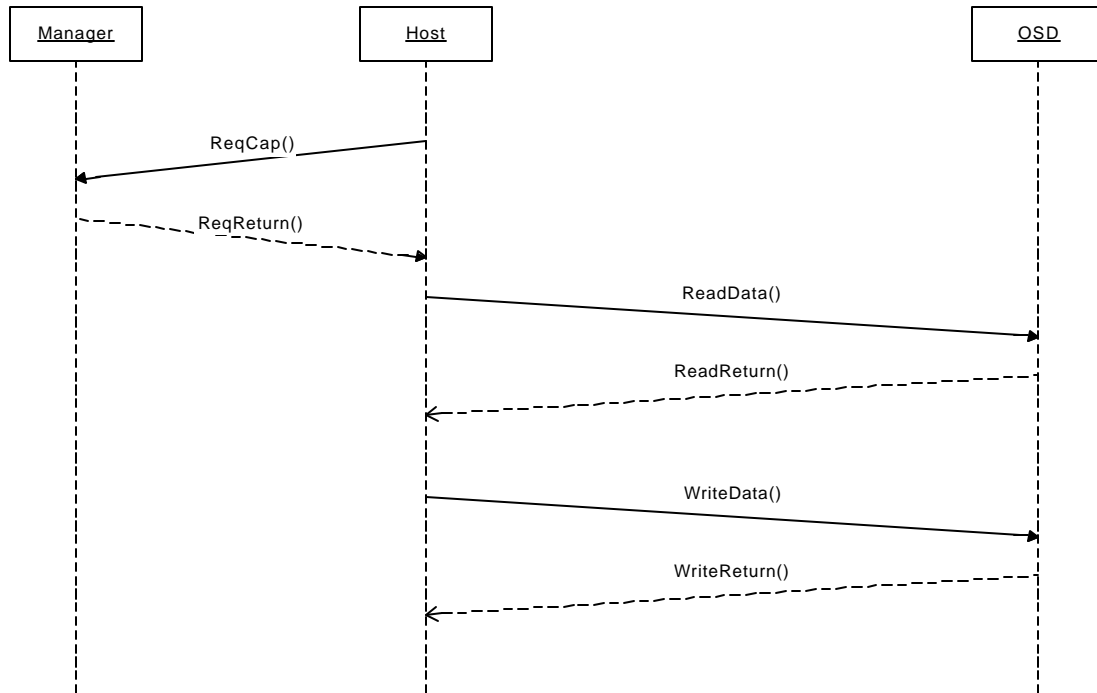
In all cases:

1. a *RetMac* is computed using *CapX* on the *Status* and the *Nonce*

In addition for successful read commands, the OSD returns

2. a *DataMac* is computed using *CAP_Key* on the *Data* and *Nonce*

958    to allow the host to verify the response.



959
960

| Message | Arguments | Explanation |
|---|---|---|
| | | |
| **ReqCap** | ObjID | Object identifier |
| | Partition ID | Partition Id |
| | ObSID | Object Store identifier |
| | | |
| **ReqReturn** | Version | |
| | Rights | |
| | Expiration | |
| | Partition ID | |
| | Creation | |
| | Capability Nonce | Capability nonce |
| | CAP_Key | $MAC_{secret\_key}(CAP\_Args)$ |
| | | |
| **ReadData** | ObjID | |
| | Partition | |
| | ObSID | |
| | Offset | |
| | Length | |
| | Nonce | Time-based nonce |
| | CAP_Args | |
| | ReqMac | $MAC_{CAP\_Key}(Cmd\_Args, Nonce)$ |
| | | |
| **ReadReturn** | Status | return code from the request, success or failure |
| | DataMac | $MAC_{CAP\_Key}(Data, Nonce)$ |
| | RetMac | $MAC_{CAP\_Key}(Status, Nonce)$ |
| | | |
| **WriteData** | Cmd_Args | |

In the ReqReturn rows, Version, Rights, Expiration, Partition ID are braced as **CAP_Args**.
In the ReadData rows, ObjID, Partition, ObSID, Offset, Length are braced as **Cmd_Args**.

36

| | Nonce | Time-based nonce |
|---|---|---|
| | CAP_Args | |
| | ReqMac | $MAC_{CAP\_Key}$(Cmd_Args, Nonce) |
| | DataMac | $MAC_{CAP\_Key}$(Data, Nonce) |
| | | |
| **WriteReturn** | Status | |
| | RetMac | $MAC_{CAP\_Key}$(Status, Nonce) |

961

## 6.1 Implementation Efficiency

963 The efficient computation of the *DataMac* is straightforward in the case of READ. As data is
964 read from the media, the MAC is computed and it is sent as part of the status message at the
965 end of the command.

966 The case of WRITE is more difficult. If the *DataMac* is sent in the same message as the
967 command, then the Host must make two passes over the data – one to compute the MAC and
968 a second to send the data. In order to avoid this, there must be an additional message as shown
969 in the following.

970



971
972

| **ReadVerify** | DataMac | $MAC_{CAP\_Key}$(Data, Nonce) |
|---|---|---|
| | | |
| **WriteVerify** | DataMac | $MAC_{CAP\_Key}$(Data, Nonce) |

973

974 Implementation of this additional message must be supported by the underlying transport in
975 order to achieve the necessary efficiency.

# 7 Security Manger – OSD protocol

977 While the precise behavior and policies applied by the security manager are not defined by this
978 protocol, the interactions between the security manager and the OSD are defined.

979 The OSD treats commands from the security manager in the same way it processes commands
980 received from a host. In other words, these commands must contain a valid capability
981 authorizing the operation. A security manager must use the appropriate level of security as
982 specified for the partition with which it is interacting.

## 7.1 Invalidation of capabilities for a Specific Object

984 The security manager can invalidate all previously issued capabilities for a given object by
985 informing the OSD that it should only accept capabilities for the object with a given object
986 version tag. The parameters that must be provided in this command include:

987 • *Object* – the identity of the object to which this command applies. This should include the
988 partition ID and local ID

989 • *Object Version Tag* - the value below which no capability will be accepted for this object.

990 This function will be realized as a set attribute on the indicated object. In addition to allowing
991 set attributes, the capabilitiy provided for this function must include administrative rights.

## 7.2 Clocks and Expiry Time

993 The OSD must reject any capabilities that have expired. Since the time placed in the capability
994 comes from the security manager's clock, for the OSD to be able to properly interpret the
995 expiry time in the capability, we require some degree of synchronization between the clocks of
996 the OSD and Security manager.

997 The protocol for synchronizing the clocks is not specified as part of the object store protocol.
998 The expectation is that a standard clock synchronization protocol will be used; we also believe it
999 makes sense to allow multiple such protocols to be implemented. The specification of the
1000 protocol is beyond the scope of this document.

1001 We do, however, assume that this protocol will be implemented in a secure manner, *i.e.,* we do
1002 not want an adversary to be able to change the time for the OSD or Security Manager. Such
1003 an action could constitute an attack, which increased the effective lifetime of legitimately issued
1004 capabilities. Depending upon the implementation, it could also extent the time during which a
1005 secret key is used.

# 8 Key Management

The credential is based on a secret key that is shared between the object store and the security manager. In order to prevent an adversary from obtaining too many credentials generated with the same key, keys must be refreshed regularly. Thus, a key management scheme is required.

## 8.1 Requirements

- The security manager should be able to replace the object store keys in a secure manner even if the channel it has with the object store is not secure.

- The security manager (or a higher level authority) should be able to divide the drive into multiple partitions. Each partition should carry its own keys for security purposes. Thus, a credential generated for one partition cannot be valid for another.

- A key refresh should invalidate all the credentials generated by that key.

- The key refresh scheme should not necessarily lead to a surge in the communication caused by clients requesting a new valid credential.

- The security manager has a source for random bits.

- The object store is not required to have a source for generating random bits.

- The drive manufacturer cannot assume to know the identity of the drive purchaser.

- The drive manufacturer should not have control over the drive once it is initialized. *i.e.*, the manufacturer should not be able to know the secret keys that are used to generate credentials.

- A drive crash should not necessarily invalidate valid credentials.

- Provisioning a new drive should not require mechanical actions to configure the security mechanism.

## 8.2 Key Hierarchy

We suggest using the key hierarchy proposed by Gobioff in [7]. The key hierarchy is comprised of 4 layers as described below:

- **Master key** – held by the disk owner. Used to initialize the drive and to create the drive key. This key does not change unless the drive owner is changed. As the top most key in the hierarchy it should be used as little as possible in order to reduce its exposure, and it would be preferable if this key could be immutable as long as the drive does not change owners.

- **Drive key** – held by the disk owner, used to divide the drive into multiple partitions and to create the partition keys. This key is used very rarely and is changed only if either it is (suspected to be) compromised, or the drive owner changes, or a (rare) key refresh operation is carried in order to increase security.

1040       • **Partition keys** – held by the (partition's) security manager. Used solely to create
1041        the working keys. The partition keys are changed infrequently, but in a regular
1042        manner to increase security.

1043       • **Working keys** – held by the (partition's) security manager. Used to generate the
1044        cap-keys. The working keys are refreshed frequently (e.g., on an hourly or daily
1045        basis) in order to limit the number of credentials that are generated by the same key.

### 1046   8.2.1   Master Key

1047 The *master key* is the topmost key in the hierarchy. It allows unrestricted access to the drive.
1048 Its loss is considered a catastrophic event. Due to the importance of the master key, it is desired
1049 to limit its use as much as possible. Thus, the only use of the master key is to initialize the drive
1050 and to set the drive key. This master key does not change unless the drive owner is changed,
1051 *e.g.,* the drive is sold. We denote the Master key by **Km**.

### 1052   8.2.2   Drive Key

1053 The *drive key* provides an unrestricted access to the drive, very much like the master key,
1054 except that it cannot be used either to initialize the drive or to set another master key or a new
1055 drive key. Once the drive key is set it can be used to divide the drive into partitions and to set
1056 the partitions' keys. The drive key can be changed in case it was compromised, or as part of a
1057 scheduled update operation in order to maintain security. We denote the drive key by **Kd**.

### 1058   8.2.3   Partition Key

1059 An object store can be divided into multiple partitions, formerly known as *security classes*,
1060 which carry their own keys for security purposes.  From the perspective of the security
1061 manager, it will have a distinct secret key for each two-tuple of object store $s_j$ and partition $c_k$.
1062 We denote the key of partition j by **Kp$_j$**.

### 1063   8.2.4   Working Key

1064 The working keys are used to generate the capability keys for a particular partition; hence they
1065 should be refreshed very frequently, *e.g.,* on an hourly basis. However, since a key refresh
1066 event between the object store and the security manager invalidates all credentials generated by
1067 that key at once, a simplistic scheme which keeps only a single working key for each partition
1068 would result in an undesired performance degradation as all the clients would be required to
1069 communicate with the security manager in order to get new credentials; moreover, all new
1070 credentials must be explicitly validated (via MAC calculation) before being cached by the object
1071 store. To mitigate the undesired effects of a key refresh, the following optimization, as suggested
1072 in [8], can be used: an object store may declare the last two (or more generally *n*) refreshed
1073 versions of the *working key* as valid, instead of just the latest one. As a result, the process of
1074 validating a capability requires a *key_version* field to be incorporated in the capability indicating
1075 which key should be used in the validation process.[33]

1076

---

[33] For more details on this mechanism, see the object store security document

1077 The number of active key versions used is configured between the OSD and the security
1078 manager. When setting a new working key, the security manager tags the key with a version
1079 number (between 0 and 15); the object store uses this tag to determine which key to use in
1080 validating a command. The OSD implementation can specify the maximum number of key
1081 versions it supports; one is a legal value. The maximum number of key versions supported by
1082 the protocol is 16.

1083 We denote the working key of partition $j$ with version $i$ by $\mathbf{Kw_{j,i}}$.

## 8.3 Key Exchange Protocol

1085 We present a protocol for key exchange that applies well-known techniques for key updates[34]
1086 and does not use encryption.

1087 The protocol has the following characteristics:

1088     •   Except for the topmost key, keys of one level can be replaced only by using a
1089        higher-level key. We describe how the master key is set in the Drive Initialization
1090        section.

1091     •   The compromise of a key at a given level does not reveal information on keys in
1092        higher levels, or on other keys (if multiple key versions exist) at the same level.

1093     •   The exchange of a key at a given level invalidates all keys at lower levels (e.g., a new
1094        partition key invalidates all working keys).

1095 We propose that the drive use a pseudo random number generator to generate the keys using a
1096 random string (a seed) which is sent to it by the drive owner / security manager. Note that the
1097 security manager and the drive must use the same generation procedure.

1098 A cryptographic pseudo random number generator may be constructed either from a good
1099 MAC function, *e.g.,* SHA1, or a block cipher function, *e.g.,* AES. The specific cryptographic
1100 pseudorandom number generator we propose is one that utilizes the cryptographic hash SHA-
1101 1, as defined in FIPS 186, Section 3.3. Upon selecting the seed $s$, it basically applies the MAC
1102 function to the values $s$ and $s+1$ using a shared (secret) key.

1103 Again, *TimeNonce* refers to the 12-bytes nonce structure defines in the OSD protocol (a 32-
1104 bits timestamp followed by 64 random bits).

1105 We require that at each level, there will be two keys rather then one. The first key is used for
1106 message authentication and the second for key generation. For example, instead of having one
1107 master key, Km, we have two keys, a keyed MAC key, denoted $\mathbf{K_{m\_A}}$, used for message
1108 authentication and a second key for the pseudo random number generator, denoted $\mathbf{K_{m\_G}}$, used
1109 for key generation. The same scheme holds for every level. As before, we defer the discussion
1110 on how to set the master keys to the *Drive Initialization* section.

---

[34] See for instance section 12.3.1 of [3], Remark 12.19 (pp. 498-490), states that the confidentiality of the key
update is not necessary, and that it may be avoided by employing instead a key derivation from a
pseudorandom permutation.

1111 Note that the protocol does not describe how random seeds are generated. It is the
1112 responsibility of the security manager to create them as random as possible.

### 8.3.1  Setting the Drive Key

1114 In order to set the drive key, a **SetKey** message is sent (as described below), protected by the
1115 master key.  This command will include a *Seed, which* is a random string of length 160 bits
1116 computed by the drive owner; LSB (least significant bit) of the seed must be zero.

1117 The new drive authentication key and generation key are computed by applying the generator
1118 function on the seed to obtain two distinct pseudo random numbers as follows:

1119 $\quad$ $K_{d\_G} = G_{Km\_G}(Seed\ or\ 0x01)$
1120 $\quad$ $K_{d\_A} = G_{Km\_G}(Seed)$

### 8.3.2  Setting a Partition Key

1122 In order to set the keys of a specific partition, a **SetKey** message is sent (as described below),
1123 protected by the drive key.  The command will include a seed as defined above as well as a
1124 *Partition Number, which* is the number of the partition for which the key is to be set.

1125 The new partition authentication key and generation key are computed by:

1126 $\quad$ $K_{p,partition\ number\_G} = G_{Kd\_G}(Seed\ or\ 0x01)$
1127 $\quad$ $K_{p,partition\ number\_A} = G_{Kd\_G}(Seed\ )$

1128 Note, setting a partition key invalidates all working keys for the partition and thus all capability
1129 keys for the partition.

### 8.3.3  Setting a Working Key

1131 In order to set the working keys of a specific partition, a **SetKey** is sent (as described below),
1132 protected by the partition key, *e.g.,* for partition j, the security manager uses Kp,j.  The
1133 command will include a seed and partition number as defined above, as well as a *Version*
1134 *Number,* which is the version number of the key to be set.

1135 The new working authentication key and generation key are computed by:

1136 $\quad$ $K_{w,j,version\ number\_G} = G_{K\ p,j\ \_G}(Seed\ or\ 0x01)$
1137 $\quad$ $K_{w,j,\ version\ number\_A} = G_{K\ p,j\ \_G}(Seed)$

## *8.4  Using the standard protocol to Set Keys*

1139 Instead of defining a set of specific protocol messages to be used for key management, we can
1140 use a single new SetKey command along with the basic OSD security mechanisms.  We assume
1141 that we have objects (or pseudo objects) with known identifiers representing the object store as
1142 a whole as well as each partition.  The partition and working keys are set by invoking SetKey
1143 on the object for the partition and the drive key by invoking SetKey on the object for the object
1144 store as a whole.

1145 The parameters of the command are:

1146 • One of the following, DriveKey , PatritionKey, or WorkingKey depending upon the
1147 key being set

1148 • an 8-byte string composed of a 1-byte KeyVersion[35] followed by 7 bytes that
1149 uniquely identifies the key (a counter will do). In particular, the key identifier
1150 indicates the Partition number.  This information can be used for auditing and other
1151 reporting purposes.

1152 • the information that is needed to infer the next key, i.e., Value is set to be the Seed
1153 that is used to generate the two corresponding keys (message authentication key and
1154 key generation).[36]

1155 The command is sent using the OSD security protocol as appropriate for the level of security
1156 being used by the object store.   For messages sent to set the key for the drive, the object
1157 representing the drive must be queried to determine the appropriate security level.  The CAP-
1158 Args right-string must contain an indication that keys can be set.   Note that the *CAP_Key* that
1159 corresponds to the credential issued on this command is computed using $K_{higher\_A}$. Specifically,
1160 $CAP\_key = MAC\_K_{higher} (CAP\text{-}Args)$.

## 8.5  Drive Initialization

1162 The protocol gives full power over the drive to the possessor of the master key. Thus, using and
1163 setting the master key should be done in the most secure environment possible.  To allow setting
1164 the master key after the drive is obtained from a vendor, we assume that the drive comes from
1165 the manufacturer with an initial master key built-in. This master key is also provided *in a secure*
1166 *manner* (*e.g.,* a floppy, a separate email message) to the owner.  Before the drive is used for
1167 storing the client data, the drive must be initialized. The initialization is done by replacing the
1168 initial master key with a new one, generated by the security manager / drive owner.  Note that

1169 • The manufacturer cannot access the drive if initialization was done properly since the
1170 new Master Key is known only to the owner.

1171 • If the drive has been initialized elsewhere (mistakenly or maliciously) this will be
1172 detected by the owner as the initial Master Key that was provided to the owner will
1173 no longer work.

1174 The following command will be used to set the master key. The message is authenticated using
1175 the previous master key denote by   Km_A_old

1176 **SetMasterKey** msg $M_{Km\_A\_previous}$ (**SetMasterKey**, msg)
1177 Where msg = Seed, TimeNonce

1178 • Seed is a random string of length 160 bits computed by the drive owner; the LSB
1179 (least significant bit) of the seed is zero.

---

[35] In the range 0-15.
[36] There is an assumption for Level 2 security that the attribute value is part of the command parameters and thus protected by the per command MAC.

1180   The effect of this command is to set the master key as follows:

1181      $K_{m\_G\_new} = G_{Km\_G\_previous}$ (*Seed or 0x01*)
1182      $K_{m\_A\_new} = G_{Km\_G\_previous}$ (*Seed*)

1183   Note, if one is concerned that an entity may listen on the wire as well as steal the master key
1184   provided by the object store manufacturer, there is nothing that prevents sending these
1185   commands via a direction connection and not over a network.

1186   We point out that this differs from the suggestion in [8] is that the drive comes in an
1187   *uninitialized state,* where it has no partitions and no valid keys. Here, before the drive is
1188   placed in the general network, the owner initializes it using a secure network, e.g., a cable
1189   directly attached from the owner laptop to the drive.

## 8.6  Storing Long Lived Keys

1191   The drive keys are considered highly secret information. It is important to protect them from
1192   being leaked to an adversary.  In order to protect the drive the keys should be stored in a
1193   tamper resistant[37] nonvolatile manner and maybe even protected by tamper resistant software
1194   shield. Note that only the master key must be remembered in a tamper resistant manner.  The
1195   seeds that were used to create all other keys can be saved in a nonvolatile memory and used to
1196   recompute the keys in case of a drive crash.

1197   Note, the object store should not remember the messages sent to set the master key in a
1198   manner that could be externally accessible.[38]

## 8.7  Secure Computation

1200   In order to conform with FIPS 140-1 [5] level 4, storing keys, computing the credential keys
1201   and the key exchange protocol should be done in a secure coprocessor.

## 8.8  Parameterizing Cryptographic Primitives

1203   We would like to provide the flexibility of having an object store support multiple
1204   implementations of the cryptographic primitives, i.e., MAC functions.  To do this, a root object
1205   will support an attribute which provides the cryptographic primitives an object store prefers; this
1206   will be provided as an ordered and numbered list of primitives, where number zero is the highest
1207   preference.  We will allow an object store to support up to sixteen primitives.  Note all objects
1208   stores must support an HMAC SHA-1.

---

[37] See *Security Engineering - A guide to building dependable distributed systems*, by Ross Anderson, John Wiley & Sons, Inc. pp.277-304.
[38] The actual requirement for correctness may be slightly weaker than this, but this seems to be sufficient, if not completely necessary.

1209 When the user gets the initial key for the object store, the key will also specify which
1210 cryptographic primitives to use with the initial key exchange; the number of this combination will
1211 also be specified.

1212 The CAP_Args includes a four bit field indicating the cryptographic primitive used to construct
1213 the credential. The security manager will place in this field the number of the cryptographic
1214 primitives used in constructing thecredential. The security manager will need to take into
1215 account the clients capabilities when it gives a credential to the client. The client will need to use
1216 the cryptographic primitive upon which it agreed with the security manager. The intent of this
1217 approach was to allow a smooth upgrade of a system, in which some clients may not support a
1218 newer cryptographic primitive.

1219 In the first version of the standard we will only support a single MAC function. Later versions
1220 of the standard will need to address the security issues that arise in using multiple MAC
1221 functions with a single key.

1222

# 9 References

*[1]* Azagury, R. Canetti, M. Factor, S. Halevi, E. Henis, D. Naor, N. Rinetzky, O. Rodeh, J. Satran, "A Two Layered Approach for Securing an Object Store Network," *First IEEE International Security In Storage Workshop*, Greenbelt, MD, Dec 2002

*[2]* A. J. Menezes, P. C. Van Oorschot, and S. A. Vanstone, *Handbook of Applied Cryptography,* CRC Press 1996. pp. 325.

*[3]* A. J. Menezes, P. C. Van Oorschot, S. A. Vanstone*, Handbook of Applied Cryptography*, by, CRC Press 1996. Section 12.3.1

*[4]* AES Advanced encryption standard

*[5]* FIPS 140-1 security standard

*[6]* FIPS Publication 186, Section 3.3

*[7]* H. Gobioff, *et al.,* Security for network attached storage devices, Technical report, CMU-CS-97-185.ps

*[8]* H. Gobioff, *Security for a High Performance Commodity Storage Subsystem,* PhD thesis, Carnegie Mellon University, 1999.

*[9]* H. Krawczyk, M. Bellare, R. Canetti, "HMAC: Keyed-Hashing for Message Authentication", RFC 2104, http://www.ietf.org/rfc/rfc2104.txt

*[10]* M. Bellare, R. Canetti, H. Krawczyk, "The HMAC Construction", *Cryptobytes* Vol. 2, No. 1, Spring 1996.

*[11]* M.K. Aguilera, M. Ji, M. Lillibridge, J. MacCormick, E. Oertli, D. Andersen, M. Burrows, T. Mann, C.A. Thekkath, *"Block-Level Security for Network-Attached Disks," 2nd Usenix Conference on File and Storage Technology*, San Francisco, CA, March 2003.
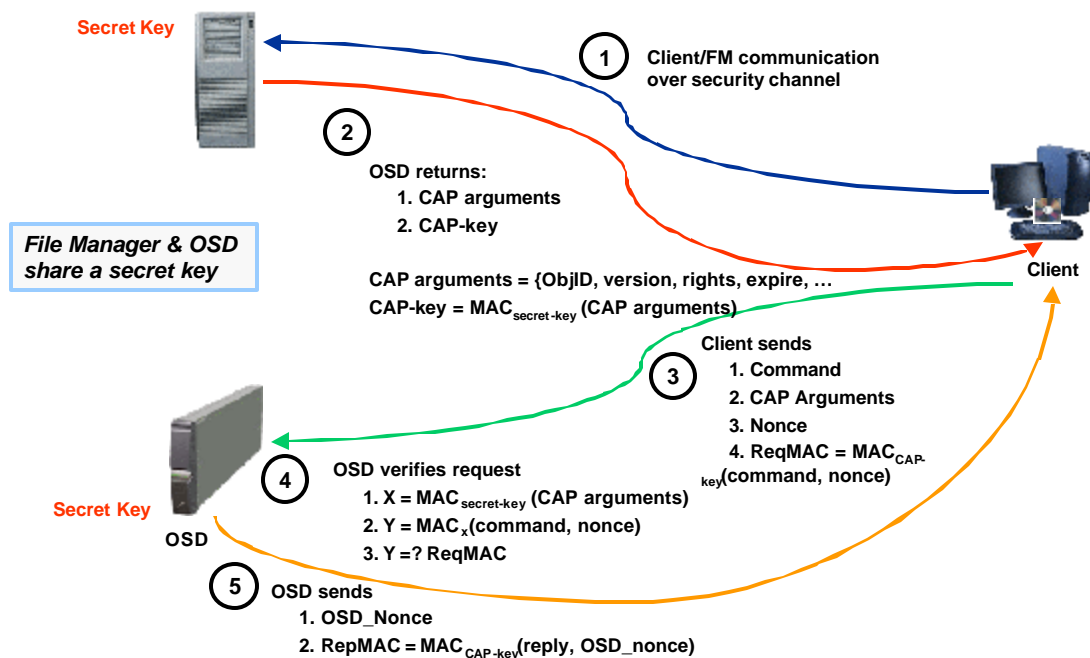
# 10 Appendix

## 10.1 Comparison to Original Approach

We now summarize the original NASD protocol and describe the differences between the original protocol and the current proposal.

### 10.1.1 Original NASD Proposal

As we stated, we want to enable providing only integrity of capabilities (if desired). In the original NASD work 9, which is the starting point for this work, integrity of capabilities is intertwined with network security. To access an object, a host receives a credential composed of a capability and a CAP-key from a Security Manager; the CAP-key is derived from the capability using a secret shared between the object store and the security manager. On each request, the CAP-key is used to authenticate the request: the CAP-key is used as key for a MAC on the nonce and the command/data. The nonce provides anti-replay, *i.e.,* provides a function of network security. The MAC on the command/data provides integrity of command/data. The use of the CAP-key for computing the MAC implicitly provides integrity of capability; if the capability had been modified then the object store would fail in its attempt to validate the MAC of the command/data. The CAP-key is also used by the object store to authenticate its reply to the client.



**Secret Key**

**(1) Client/FM communication over security channel**

**(2) OSD returns:**
1. CAP arguments
2. CAP-key

*File Manager & OSD share a secret key*

CAP arguments = {ObjID, version, rights, expire, …
CAP-key = $MAC_{secret\text{-}key}$ (CAP arguments)

**(3) Client sends**
1. Command
2. CAP Arguments
3. Nonce
4. ReqMAC = $MAC_{CAP\text{-}key}$(command, nonce)

**(4) OSD verifies request**
1. X = $MAC_{secret\text{-}key}$ (CAP arguments)
2. Y = $MAC_{x}$(command, nonce)
3. Y =? ReqMAC

**Secret Key**

**OSD**

**(5) OSD sends**
1. OSD_Nonce
2. RepMAC = $MAC_{CAP\text{-}key}$(reply, OSD_nonce)

**Client**

1267 This approach requires both the host and the object store to calculate a new MAC for each
1268 command. However, if we have a secure or trusted network, a direct application of original
1269 NASD protocol involves redundant computation. In particular, if we were running on top of an
1270 IPSec authenticated channel we would have:
1271
1272 • Two mechanisms for anti-replay
1273 • Two mechanisms for integrity of data
1274
1275 This leads to our challenge: Define integrity of capabilities  and integrity of command/data such
1276 that integrity of capabilities uses a subset of the cryptographic structure.
1277 In addition to this major challenge, there are some additional minor issues with the original
1278 definition of the protocol. These issues led to additional changes from the original NASD
1279 protocol in the version of the object store security protocol presented in the following sections.

## 1280 10.1.2  Ability to Use Either Channel ID or Command Unique Nonce

1281 By replacing the command unique nonce with a channel ID, we are able to extend the original
1282 NASD protocol into a protocol that supports running on an externally secured channel without
1283 incurring unnecessary overhead.   Since the channel ID does not change on each command, it is
1284 not necessary to recalculate a MAC that involves this channel ID on each command.
1285 However, since the channel ID is tied to the channel and the channel is authenticated, receipt of
1286 a MAC based upon this channel ID enables the object store to be certain that the capability it is
1287 receiving was legitimately obtained by the client.

## 1288 10.1.3  Unique Value Added to *CAP_Args*

1289 To avoid scenarios in which the same *CAP_Args* and *CAP_key* is given by the security
1290 manager to different clients requesting the same rights to the same (set of) object(s), we add a
1291 unique value to each *CAP_Args*.  This change closes a potential security hole in the version of
1292 the protocol using internal security.   Without this change a client could masquerade as an object
1293 store for another client, if both clients get the same authorization for a given object.