



Advancing storage &
information technology

NVM Programming Model (NPM)

Version 1.1

Abstract: This SNIA document defines recommended behavior for software supporting Non-Volatile Memory (NVM).

Publication of this Working Draft for review and comment has been approved by the NVM Programming TWG. This draft represents a “best effort” attempt by the NVM Programming TWG to reach preliminary consensus, and it may be updated, replaced, or made obsolete at any time. This document should not be used as reference material or cited as other than a “work in progress.” Suggestion for revision should be directed to <http://www.snia.org/feedback/>.

Working Draft

November 7, 2014

2 USAGE

3 The SNIA hereby grants permission for individuals to use this document for personal use only,
4 and for corporations and other business entities to use this document for internal use only
5 (including internal copying, distribution, and display) provided that:

- 6 1. Any text, diagram, chart, table or definition reproduced shall be reproduced in its
7 entirety with no alteration, and,
8
- 9 2. Any document, printed or electronic, in which material from this document (or any
10 portion hereof) is reproduced shall acknowledge the SNIA copyright on that material,
11 and shall credit the SNIA for granting permission for its reuse.
12

13 Other than as explicitly provided above, you may not make any commercial use of this
14 document, sell any or this entire document, or distribute this document to third parties. All
15 rights not explicitly granted are expressly reserved to SNIA.

16 Permission to use this document for purposes other than those enumerated above may be
17 requested by e-mailing tcmd@snia.org. Please include the identity of the requesting individual
18 and/or company and a brief description of the purpose, nature, and scope of the requested
19 use.

20 All code fragments, scripts, data tables, and sample code in this SNIA document are made
21 available under the following license:

22 BSD 3-Clause Software License

23 Copyright (c) 2014, the Storage Networking Industry Association.

24
25 Redistribution and use in source and binary forms, with or without modification, are permitted provided that the
26 following conditions are met:

27
28 * Redistributions of source code must retain the above copyright notice, this list of conditions and the following
29 disclaimer.

30
31 * Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following
32 disclaimer in the documentation and/or other materials provided with the distribution.

33
34 * Neither the name of The Storage Networking Industry Association (SNIA) nor the names of its contributors may be
35 used to endorse or promote products derived from this software without specific prior written permission.

36
37 THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY
38 EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF
39 MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL
40 THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL,
41 SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO,
42 PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS
43 INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT,
44 STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE
45 USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.
46
47

48 **DISCLAIMER**

49 The information contained in this publication is subject to change without notice. The SNIA
50 makes no warranty of any kind with regard to this specification, including, but not limited to, the
51 implied warranties of merchantability and fitness for a particular purpose. The SNIA shall not
52 be liable for errors contained herein or for incidental or consequential damages in connection
53 with the furnishing, performance, or use of this specification.

54 Suggestions for revisions should be directed to <http://www.snia.org/feedback/>.

55 Copyright © 2014 SNIA. All rights reserved. All other trademarks or registered trademarks are
56 the property of their respective owners.

57 **Revision History**

58 **Update 1 Revision 1**

59 **Date**

60 November 7, 2014

61

62 **Changes Incorporated**

63

- 64 • USAGE: added BSD 3-clause license (new SNIA template)
- 65 • 1 Scope: last sentence: clarified expectation that sharing data was consistent with native
- 66 hardware (as well as OS) behavior
- 67 • 3.4 Conventions: remove text on numeric conventions because they are not used
- 68 • 4.1 bullet #1: changed hard-coded section number to xref
- 69 • 4.3.3 remove bullet saying NVM.PM.VOLUME addresses errors
- 70 • 6.8 last paragraph: changed *Programming* to have 2 “m”s
- 71 • 6,10 Aligned operations on fundamental data types::moved here from former Consistency
- 72 Annex (Ballot-Proposal-00008)
- 73 • 7.2.7 NVM.BLOCK.SCAR last paragraph: added “r” to “scarred” in “A block stays scarred
- 74 until it is updated by a write operation.”
- 75 • 8.4.3.2.3, Failure Scenarios, paragraph 1: removed word “fundamental”
- 76 • 9.1 remove bullet saying NVM.PM.VOLUME addresses errors
- 77 • 10.1 NVM.PM.FLE Overview: changes “bold red line” to “red wavy line”
- 78 • 10.1.1 Applications and PM Consistency:: Content moved from former Consistency Annex
- 79 (Ballot-Proposal-00008)
- 80 • 10.2.4 NVM.PM.FILE.SYNC paragraph 3: deleted “An annex to this specification is
- 81 proposed to address this.” This refers to former Consistency Annex content which is now in
- 82 10.1.1 (Ballot-Proposal-00008)
- 83 • 10.2.23 NVM.PM.FILE.MAP paragraph 3: add OPTIMIZED_FLUSH and
- 84 OPTIMIZED_FLUSH_AND_SYNC to the list of sync actions (Ballot-Proposal-00006)
- 85 • 10.2.5 NVM.PM.FILE.OPTIMIZED_FLUSH and 10.2.7
- 86 NVM.PM.FILE.OPTIMIZED_FLUSH_AND_SYNC: added “Requires: NVM.PM.FILE.MAP”
- 87 (Ballot-Proposal-00006)
- 88 • 10.3.3 NVM.PM.FILE.INTERRUPTED_STORE_ATOMIcity first sentence: add reference
- 89 to new 6.10 discussion of aligned operations on fundamental data types (Ballot-Proposal-
- 90 00008)
- 91 • Annex B Consistency deleted (Ballot-Proposal-00008)

92

94	FOREWORD	9
95	1 SCOPE	10
96	2 REFERENCES	11
97	3 DEFINITIONS, ABBREVIATIONS, AND CONVENTIONS	12
98	3.1 DEFINITIONS	12
99	3.2 KEYWORDS.....	13
100	3.3 ABBREVIATIONS	13
101	3.4 CONVENTIONS	14
102	4 OVERVIEW OF THE NVM PROGRAMMING MODEL (INFORMATIVE)	15
103	4.1 HOW TO READ AND USE THIS SPECIFICATION.....	15
104	4.2 NVM DEVICE MODELS.....	15
105	4.3 NVM PROGRAMMING MODES.....	17
106	4.4 INTRODUCTION TO ACTIONS, ATTRIBUTES, AND USE CASES.....	19
107	5 COMPLIANCE TO THE PROGRAMMING MODEL	21
108	5.1 OVERVIEW.....	21
109	5.2 DOCUMENTATION OF MAPPING TO APIS	21
110	5.3 COMPATIBILITY WITH UNSPECIFIED NATIVE ACTIONS	21
111	5.4 MAPPING TO NATIVE INTERFACES	21
112	6 COMMON PROGRAMMING MODEL BEHAVIOR	22
113	6.1 OVERVIEW.....	22
114	6.2 CONFORMANCE TO MULTIPLE FILE MODES	22
115	6.3 DEVICE STATE AT SYSTEM STARTUP.....	22
116	6.4 SECURE ERASE	22
117	6.5 ALLOCATION OF SPACE	22
118	6.6 INTERACTION WITH I/O DEVICES	22
119	6.7 NVM STATE AFTER A MEDIA OR CONNECTION FAILURE	23

120	6.8	ERROR HANDLING FOR PERSISTENT MEMORY.....	23
121	6.9	PERSISTENCE DOMAIN.....	23
122	6.10	ALIGNED OPERATIONS ON FUNDAMENTAL DATA TYPES.....	23
123	6.11	COMMON ACTIONS.....	24
124	6.12	COMMON ATTRIBUTES.....	25
125	6.13	USE CASES.....	25
126	7	NVM.BLOCK MODE.....	27
127	7.1	OVERVIEW.....	27
128	7.2	ACTIONS.....	29
129	7.3	ATTRIBUTES.....	32
130	7.4	USE CASES.....	36
131	8	NVM.FILE MODE.....	40
132	8.1	OVERVIEW.....	40
133	8.2	ACTIONS.....	40
134	8.3	ATTRIBUTES.....	42
135	8.4	USE CASES.....	43
136	9	NVM.PM.VOLUME MODE.....	50
137	9.1	OVERVIEW.....	50
138	9.2	ACTIONS.....	50
139	9.3	ATTRIBUTES.....	53
140	9.4	USE CASES.....	55
141	10	NVM.PM.FILE.....	58
142	10.1	OVERVIEW.....	58
143	10.2	ACTIONS.....	61
144	10.3	ATTRIBUTES.....	66
145	10.4	USE CASES.....	68

146	ANNEX A (INFORMATIVE) PM POINTERS	78
147	ANNEX B (INFORMATIVE) PM ERROR HANDLING.....	79
148	ANNEX C (INFORMATIVE) DEFERRED BEHAVIOR	83
149	D.1 REMOTE SHARING OF NVM.....	83
150	D.2 MAP_CACHED OPTION FOR NVM.PM.FILE.MAP.....	83
151	D.3 NVM.PM.FILE.DURABLE.STORE.....	83
152	D.4 ENHANCED NVM.PM.FILE.WRITE	83
153	D.5 MANAGEMENT-ONLY BEHAVIOR	83
154	D.6 ACCESS HINTS	83
155	D.7 MULTI-DEVICE ATOMIC MULTI-WRITE ACTION.....	83
156	D.8 NVM.BLOCK.DISCARD_IF_YOU_MUST ACTION	83
157	D.9 ATOMIC WRITE ACTION WITH ISOLATION.....	85
158	D.10 ATOMIC SYNC/FLUSH ACTION FOR PM.....	85
159	D.11 HARDWARE-ASSISTED VERIFY	85
160		

161

Table of Figures

162	Figure 1 Block NVM example	16
163	Figure 2 PM example.....	16
164	Figure 3 Block volume using PM HW	16
165	Figure 4 NVM.BLOCK and NVM.FILE mode examples.....	17
166	Figure 5 NVM.PM.VOLUME and NVM.PM.FILE mode examples	18
167	Figure 6 NVM.BLOCK mode example	27
168	Figure 7 SSC in a storage stack	36
169	Figure 8 SSC software cache application	37
170	Figure 9 SSC with caching assistance.....	37
171	Figure 10 NVM.FILE mode example	40
172	Figure 11 NVM.PM.VOLUME mode example.....	50
173	Figure 12 Zero range offset example.....	54
174	Figure 13 Non-zero range offset example	54
175	Figure 14 NVM.PM.FILE mode example	58
176	Figure 15 Linux Machine Check error flow with proposed new interface	81
177		

178 FOREWORD

179 The SNIA NVM Programming Technical Working Group was formed to address the ongoing
180 proliferation of new non-volatile memory (NVM) functionality and new NVM technologies. An
181 extensible NVM Programming Model is necessary to enable an industry wide community of
182 NVM producers and consumers to move forward together through a number of significant
183 storage and memory system architecture changes.

184 This SNIA specification defines recommended behavior between various user space and
185 operating system (OS) kernel components supporting NVM. This specification does not
186 describe a specific API. Instead, the intent is to enable common NVM behavior to be exposed
187 by multiple operating system specific interfaces.

188 After establishing context, the specification describes several operational modes of NVM
189 access. Each mode is described in terms of use cases, actions and attributes that inform user
190 and kernel space components of functionality that is provided by a given compliant
191 implementation.

192 Acknowledgements

193 The SNIA NVM Programming Technical Working Group, which developed and reviewed this
194 standard, would like to recognize the significant contributions made by the following members:

195	<i>Organization Represented</i>	<i>Name of Representative</i>
196	EMC	Bob Beauchamp
197	Hewlett Packard	Hans Boehm
198	NetApp	Steve Byan
199	Hewlett Packard	Joe Foster
200	Fusion-io	Walt Hubis
201	Red Hat	Jeff Moyer
202	Fusion-io	Ned Plasson
203	Rougs, LLC	Tony Roug
204	Intel Corporation	Andy Rudoff
205	Microsoft	Spencer Shepler
206	Fusion-io	Nisha Talagata
207	Hewlett Packard	Doug Voigt
208	Intel Corporation	Paul von Behren

209 **1 Scope**

210 This specification is focused on the points in system software where NVM is exposed either as
211 a hardware abstraction within an operating system kernel (e.g., a volume) or as a data
212 abstraction (e.g., a file) to user space applications. The technology that motivates this
213 specification includes flash memory packaged as solid state disks and PCI cards as well as
214 other solid state non-volatile devices, including those which can be accessed as memory.

215 It is not the intent to exhaustively describe or in any way deprecate existing modes of NVM
216 access. The goal of the specification is to augment the existing common storage access
217 models (e.g., volume and file access) to add new NVM access modes. Therefore this
218 specification describes the discovery and use of capabilities of NVM media, connections to the
219 NVM, and the system containing the NVM that are emerging in the industry as vendor specific
220 implementations. These include:

- 221 • supported access modes,
- 222 • visibility in memory address space,
- 223 • atomicity and durability,
- 224 • recognizing, reporting, and recovering from errors and failures,
- 225 • data granularity, and
- 226 • capacity reclamation.

227 This revision of the specification focuses on NVM behaviors that enable user and kernel space
228 software to locate, access, and recover data. It does not describe behaviors that are specific to
229 administrative or diagnostic tasks for NVM. There are several reasons for intentionally leaving
230 administrative behavior out of scope.

- 231 • For new types of storage programming models, the access must be defined and agreed on
232 before the administration can be defined. Storage management behavior is typically
233 defined in terms of how it enables and diagnoses the storage programming model.
- 234 • Administrative tasks often require human intervention and are bound to the syntax for the
235 administration. This document does not define syntax. It focuses only on the semantics of
236 the programming model.
- 237 • Defining diagnostic behaviors (e.g., wear-leveling) as vendor-agnostic is challenging across
238 all vendor implementations. A common recommended behavior may not allow an approach
239 optimal for certain hardware.

240
241 This revision of the specification does not address sharing data across computing nodes. This
242 revision of the specification assumes that sharing data between processes and threads follows
243 the native OS and hardware behavior.

244 **2 References**

245 The following referenced documents are indispensable for the application of this document.

246 For references available from ANSI, contact ANSI Customer Service Department at (212) 642-
247 49004980 (phone), (212) 302-1286 (fax) or via the World Wide Web at <http://www.ansi.org>.

- SPC-3 ISO/IEC 14776-453, SCSI Primary Commands – 3 [ANSI INCITS 408-2005]
Approved standard, available from ANSI.
- SBC-2 ISO/IEC 14776-322, SCSI Block Commands - 2 [T10/BSR INCITS 514]
Approved standard, available from ANSI.
- ACS-2 ANSI INCITS 482-2012, Information technology - ATA/ATAPI Command Set -2
Approved standard, available from ANSI.
- NVMe 1.1 NVM Express Revision 1.1,
Approved standard, available from <http://nvmexpress.org>
- SPC-4 ISO/IEC 14776-454, SCSI Primary Commands - 4 (SPC-4) (T10/1731-D)
Under development, available from <http://www.t10.org>.
- SBC-4 ISO/IEC 14776-324, SCSI Block Commands - 4 (SBC-4) [BSR INCITS 506]
Under development, available from <http://www.t10.org>.
- T10 13-064r0 T10 proposal 13-064r0, Rob Elliot, Ashish Batwara, SBC-4 SPC-5 Atomic writes
Proposal, available from <http://www.t10.org>.
- ACS-2 r7 Information technology - ATA/ATAPI Command Set – 2 r7 (ACS-2)
Under development, available from <http://www.t13.org>.
- Intel SPG Intel Corporation, *Intel 64 and IA-32 Architectures Software Developer's Manual Combined Volumes 3A, 3B, and 3C: System Programming Guide*, Parts 1 and 2, available from <http://download.intel.com/products/processor/manual/325384.pdf>

248

249 **3 Definitions, abbreviations, and conventions**

250 For the purposes of this document, the following definitions and abbreviations apply.

251 **3.1 Definitions**

252 3.1.1 **durable**

253 committed to a persistence domain (see 3.1.7)

254 3.1.2 **load and store operations**

255 commands to move data between CPU registers and memory

256 3.1.3 **memory-mapped file**

257 segment of virtual memory which has been assigned a direct byte-for-byte correlation with
258 some portion of a file

259 3.1.4 **non-volatile memory**

260 any type of memory-based, persistent media; including flash memory packaged as solid state
261 disks, PCI cards, and other solid state non-volatile devices

262 3.1.5 **NVM block capable driver**

263 driver supporting the native operating system interfaces for a block device

264 3.1.6 **NVM volume**

265 subset of one or more NVM devices, treated by software as a single logical entity

266 See 4.2 NVM device models

267 3.1.7 **persistence domain**

268 location for data that is guaranteed to preserve the data contents across a restart of the device
269 containing the data

270 See 6.9 Persistence domain

271 3.1.8 **persistent memory**

272 storage technology with performance characteristics suitable for a load and store programming
273 model

274 3.1.9 **programming model**

275 set of software interfaces that are used collectively to provide an abstraction for hardware with
276 similar capabilities

277 **3.2 Keywords**

278 In the remainder of the specification, the following keywords are used to indicate text related to
279 compliance:

280 **3.2.1 mandatory**

281 a keyword indicating an item that is required to conform to the behavior defined in this
282 standard

283 **3.2.2 may**

284 a keyword that indicates flexibility of choice with no implied preference; “may” is equivalent to
285 “may or may not”

286 **3.2.3 may not**

287 keywords that indicate flexibility of choice with no implied preference; “may not” is equivalent to
288 “may or may not”

289 **3.2.4 need not**

290 keywords indicating a feature that is not required to be implemented; “need not” is equivalent
291 to “is not required to”

292 **3.2.5 optional**

293 a keyword that describes features that are not required to be implemented by this standard;
294 however, if any optional feature defined in this standard is implemented, then it shall be
295 implemented as defined in this standard

296 **3.2.6 shall**

297 a keyword indicating a mandatory requirement; designers are required to implement all such
298 mandatory requirements to ensure interoperability with other products that conform to this
299 standard

300 **3.2.7 should**

301 a keyword indicating flexibility of choice with a strongly preferred alternative

302 **3.3 Abbreviations**

303 ACID Atomicity, Consistency, Isolation, Durability

304 NVM Non-Volatile Memory

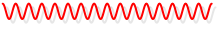
305 PM Persistent Memory

306 SSD Solid State Disk

307 **3.4 Conventions**

308 **Representation of modes in figures**

309 Modes are represented by red, wavy lines in figures, as shown below:

310 

311 The wavy lines have labels identifying the mode name (which in turn, identifies a clause of the
312 specification).

313 **4 Overview of the NVM Programming Model (informative)**

314 **4.1 How to read and use this specification**

315 Documentation for I/O programming typically consists of a set of OS-specific Application
316 Program Interfaces (APIs). API documentation describes the syntax and behavior of the API.
317 This specification intentionally takes a different approach and describes the behavior of NVM
318 programming interfaces, but allows the syntax to integrate with similar operating system
319 interfaces. A recommended approach for using this specification is:

- 320 1. Determine which mode applies (read 4.3 NVM programming modes).
- 321 2. Refer to the mode section to learn about the functionality provided by the mode and
322 how it relates to native operating system APIs; the use cases provide examples. The mode
323 specific section refers to other specification sections that may be of interest to the developer.
- 324 3. Determine which mode actions and attributes relate to software objectives.
- 325 4. Locate the vendor/OS mapping document (see 5.2) to determine which APIs map to the
326 actions and attributes.

327 For an example, a developer wants to update an existing application to utilize persistent
328 memory hardware. The application is designed to bypass caches to assure key content is
329 durable across power failures; the developer wants to learn about the persistent memory
330 programming model. For this example:

- 331 1. The NVM programming modes section identifies NVM.PM.FILE mode (see 10
332 NVM.PM.FILE) as the starting point for application use of persistent memory.
- 333 2. The NVM.PM.FILE mode text describes the general approach for accessing PM (similar
334 to native memory-mapped files) and the role of PM aware file system.
- 335 3. The NVM.PM.FILE mode identifies the NVM.PM.FILE.MAP and NVM.PM.FILE.SYNC
336 actions and attributes that allow an application to discover support for optional features.
- 337 4. The operating system vendor's mapping document describes the mapping between
338 NVM.PM.FILE.MAP/SYNC and API calls, and also provides information about supported PM-
339 aware file systems.

340 **4.2 NVM device models**

341 **4.2.1 Overview**

342 This section describes device models for NVM to help readers understand how key terms in
343 the programming model relate to other software and hardware. The models presented here
344 generally apply across operating systems, file systems, and hardware; but there are
345 differences across implementations. This specification strives to discuss the model generically,
346 but mentions key exceptions.

347 One of the challenges discussing the software view of NVM is that the same terms are often
348 used to mean different things. For example, between commonly used management
349 applications, programming interfaces, and operating system documentation, *volume* may refer
350 to a variety of things. Within this specification, *NVM volume* has a specific meaning.

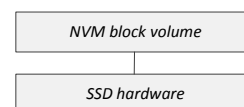
351 *An NVM volume* is a subset of one or more NVM devices, treated by software as a single
352 logical entity. For the purposes of this specification, a volume is a container of storage. A
353 volume may be block capable and may be persistent memory capable. The consumer of a
354 volume sees its content as a set of contiguous addresses, but the unit of access for a volume
355 differs across different modes and device types. Logical addressability and physical allocation
356 may be different.

357 In the examples in this section, “NVM block device” refers to NVM hardware that emulates a
358 disk and is accessed in software by reading or writing ranges of blocks. “PM device” refers to
359 NVM hardware that may be accessed via load and store operations.

360 4.2.2 Block NVM example

361 Consider a single drive form factor SSD where the entire SSD
362 capacity is dedicated to a file system. In this case, a single NVM block
363 volume maps to a single hardware device. A file system (not depicted) is
364 mounted on the NVM block volume.

Figure 1 Block NVM example

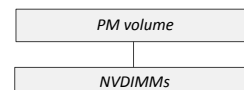


365 The same model may apply to NVM block hardware other than an SSD (including flash on
366 PCIe cards).

367 4.2.3 Persistent memory example

368 This example depicts a NVDIMM and PM volume. A PM-aware file system
369 (not depicted) would be mounted on the PM volume.

Figure 2 PM example

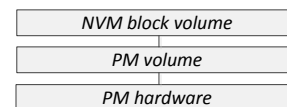


370 The same model may apply to PM hardware other than an NVDIMM (including SSDs, PCIe
371 cards, etc.).

372 4.2.4 NVM block volume using PM hardware

Figure 3 Block volume using PM HW

373 In this example, the persistent memory implementation includes a driver
374 that uses a range of persistent memory (a PM volume) and makes it
375 appear to be a block NVM device in the legacy block stack. This
376 emulated block device could be aggregated or de-aggregated like legacy
377 block devices. In this example, the emulated block device is mapped 1-1 to an NVM block
378 volume and non-PM file system.



379 Note that there are other models for connecting a non-PM file system to PM hardware.

380 **4.3 NVM programming modes**

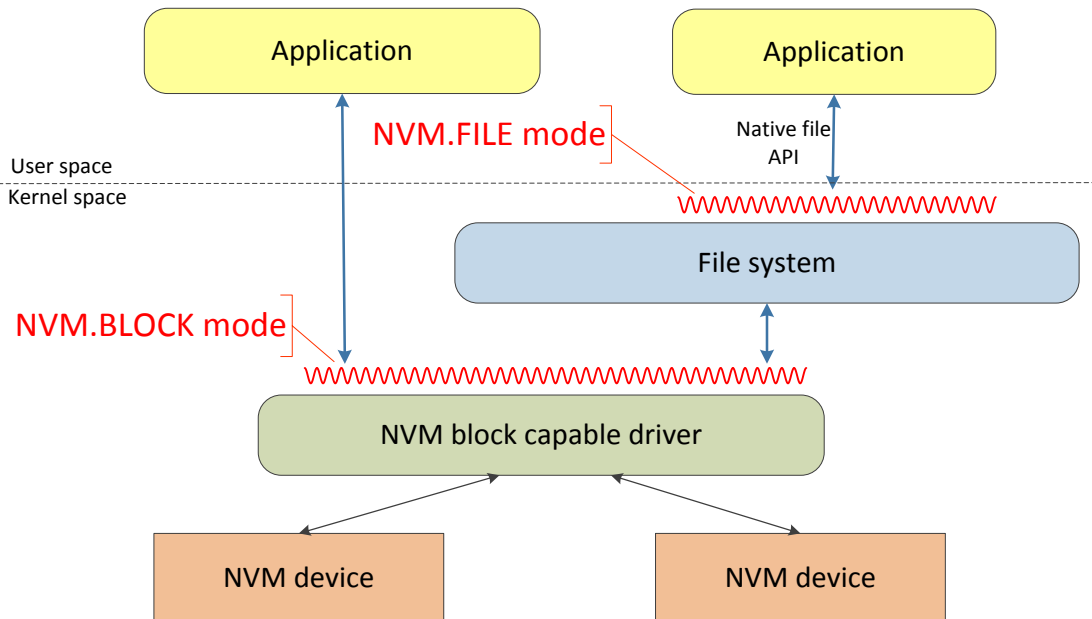
381 **4.3.1 NVM.BLOCK mode overview**

382 NVM.BLOCK and NVM.FILE modes are used when NVM devices provide block storage
383 behavior to software (in other words, emulation of hard disks). The NVM may be exposed as a
384 single or as multiple NVM volumes. Each NVM volume supporting these modes provides a
385 range of logically-contiguous blocks. NVM.BLOCK mode is used by operating system
386 components (for example, file systems) and by applications that are aware of block storage
387 characteristics and the block addresses of application data.

388 This specification does not document existing block storage software behavior; the
389 NVM.BLOCK mode describes NVM extensions including:

- 390 • Discovery and use of atomic write and discard features
391 • The discovery of granularities (length or alignment characteristics)
392 • Discovery and use of ability for applications or operating system components to mark
393 blocks as unreadable
394
395

Figure 4 NVM.BLOCK and NVM.FILE mode examples



396

397 **4.3.2 NVM.FILE mode overview**

398 NVM.FILE mode is used by applications that are not aware of details of block storage
399 hardware or addresses. Existing applications written using native file I/O behavior should work
400 unmodified with NVM.FILE mode; adding support in the application for NVM extensions may
401 optimize the application.

402 An application using NVM.FILE mode may or may not be using memory-mapped file I/O
403 behavior.

404 The NVM.FILE mode describes NVM extensions including:

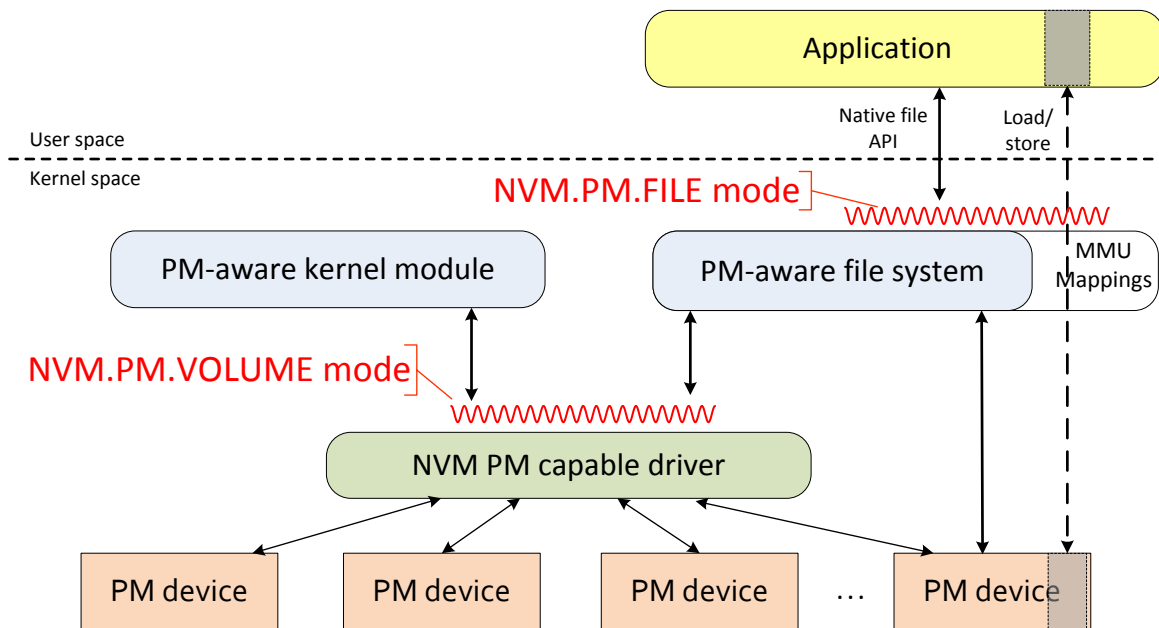
- 405 • Discovery and use of atomic write features
- 406 • The discovery of granularities (length or alignment characteristics)

407 4.3.3 NVM.PM.VOLUME mode overview

408 NVM.PM.VOLUME mode describes the behavior for operating system components (such as
409 file systems) accessing persistent memory. NVM.PM.VOLUME mode provides a software
410 abstraction for Persistent Memory hardware and profiles functionality for operating system
411 components including:

- 412 • the list of physical address ranges associated with each PM volume

413 **Figure 5 NVM.PM.VOLUME and NVM.PM.FILE mode examples**
414



415

416 4.3.4 NVM.PM.FILE mode overview

417 NVM.PM.FILE mode describes the behavior for applications accessing persistent memory.
418 The commands implementing NVM.PM.FILE mode are similar to those using NVM.FILE mode,
419 but NVM.PM.FILE mode may not involve I/O to the page cache. NVM.PM.FILE mode
420 documents behavior including:

- 421 • mapping PM files (or subsets of files) to virtual memory addresses
- 422 • syncing portions of PM files to the persistence domain

423 **4.4 Introduction to actions, attributes, and use cases**

424 **4.4.1 Overview**

425 This specification uses four types of elements to describe NVM behavior. Use cases are the
426 highest order description. They describe complete scenarios that accomplish a goal. Actions
427 are more specific in that they describe an operation that represents or interacts with NVM.
428 Attributes comprise information about NVM. Property Group Lists describe groups of related
429 properties that may be considered attributes of a data structure or class; but the specification
430 allows flexibility in the implementation.

431 **4.4.2 Use cases**

432 In general, a use case states a goal or trigger and a result. It captures the intent of an
433 application and describes how actions are used to accomplish that intent. Use cases illustrate
434 the use of actions and help to validate action definitions. Use cases also describe system
435 behaviors that are not represented as actions. Each use case includes the following
436 information:

- 437 • a purpose and context including actors involved in the use case;
- 438 • triggers and preconditions indicating when a use case applies;
- 439 • inputs, outputs, events and actions that occur during the use case;
- 440 • references to related materials or concepts including other use cases that use or extend the
441 use case.

442 **4.4.3 Actions**

443 Actions are defined using the following naming convention:

444 <context>.<mode>.<verb>

445 The actions in this specification all have a context of “NVM”. The mode refers to one of the
446 NVM models documented herein (or “COMMON” for actions used in multiple modes). The verb
447 states what the action does. Examples of actions include “NVM.COMMON.GET_ATTRIBUTE”
448 and “NVM.FILE.ATOMIC_WRITE”. In some cases native actions that are not explicitly
449 specified by the programming model are referenced to illustrate usage.

450 The description of each action includes:

- 451 • parameters and results of the action
- 452 • details of the action’s behavior
- 453 • compatibility of the action with pre-existing APIs in the industry

454 A number of actions involve options that can be specified each time the action is used. The
455 options are given names that begin with the name of the action and end with a descriptive term
456 that is unique for the action. Examples include NVM.PM.FILE.MAP_COPY_ON_WRITE and
457 NVM.PM.FILE.MAP_SHARED.

458 A number of actions are optional. For each of these, there is an attribute that indicates whether
459 the action is supported by the implementation in question. By convention these attributes end
460 with the term “CAPABLE” such as NVM.BLOCK.ATOMIC_WRITE_CAPABLE. Supported
461 options are also enumerated by attributes that end in “CAPABLE”.

462 4.4.4 **Attributes**

463 Attributes describe properties or capabilities of a system. This includes indications of which
464 actions can be performed in that system and variations on the internal behavior of specific
465 actions. For example attributes describe which NVM modes are supported in a system, and
466 the types of atomicity guarantees available.

467 In this programming model, attributes are not arbitrary key value pairs that applications can
468 store for unspecified purposes. Instead the NVM attributes are intended to provide a common
469 way to discover and configure certain aspects of systems based on agreed upon
470 interpretations of names and values. While this can be viewed as a key value abstraction it
471 does not require systems to implement a key value repository. Instead, NVM attributes are
472 mapped to a system’s native means of describing and configuring those aspects associated
473 with said attributes. Although this specification calls out a set of attributes, the intent is to allow
474 attributes to be extended in vendor unique ways through a process that enables those
475 extensions to become attributes and/or attribute values in a subsequent version of the
476 specification or in a vendor’s mapping document.

477 4.4.5 **Property group lists**

478 A **property group** is set of property values used together in lists; typically **property group**
479 **lists** are inputs or outputs to actions. The implementation may choose to implement a property
480 group as a new data structure or class, use properties in existing data structures or classes, or
481 other mechanisms as long as the caller can determine which collection of values represent the
482 members of each list element.

483 **5 Compliance to the programming model**

484 **5.1 Overview**

485 Since a programming model is intentionally abstract, proof of compliance is somewhat indirect.
486 The intent is that a compliant implementation, when properly configured, can be used in such a
487 way as to exhibit the behaviors described by the programming model without unnecessarily
488 impacting other aspects of the implementation.

489 Compliance of an implementation shall be interpreted as follows.

490 **5.2 Documentation of mapping to APIs**

491 In order to be considered compliant with this programming model, implementations must
492 provide documentation of the mapping of attributes and actions in the programming model to
493 their counterparts in the implementation.

494 **5.3 Compatibility with unspecified native actions**

495 Actions and attributes of the native block and file access methods that correspond to the
496 modes described herein shall continue to function as defined in those native methods. This
497 specification does not address unmodified native actions except in passing to illustrate their
498 usage.

499 **5.4 Mapping to native interfaces**

500 Implementations are expected to provide the behaviors specified herein by mapping them as
501 closely as possible to native interfaces. An implementation is not required to have a one-to-one
502 mapping between actions (or attributes) and APIs – for example, an implementation may have
503 an API that implements multiple actions.

504 NVM Programming Model action descriptions do not enumerate all possible results of each
505 action. Only those that modify programming model specific behavior are listed. The results that
506 are referenced herein shall be discernible from the set of possible results returned by the
507 native action in a manner that is documented with action mapping.

508 Attributes with names ending in `_CAPABLE` are used to inform a caller whether an optional
509 action or attribute is supported by the implementations. The mandatory requirement for
510 `_CAPABLE` attributes can be met by the mapping document describing the implementation's
511 default behavior for reporting unsupported features. For example: the mapping document
512 could state that if a flag with a name based on the attribute is undefined, then the
513 action/attribute is not supported.

514 **6 Common programming model behavior**

515 **6.1 Overview**

516 This section describes behavior that is common to multiple modes and also behavior that is
517 independent from the modes.

518 **6.2 Conformance to multiple file modes**

519 A single computer system may include implementations of both NVM.FILE and NVM.PM.FILE
520 modes. A given file system may be accessed using either or both modes provided that the
521 implementations are intended by their vendor(s) to interoperate. Each implementation shall
522 specify its own mapping to the NVM Programming Model.

523 A single file system implementation may include both NVM.FILE and NVM.PM.FILE modes.
524 The mapping of the implementation to the NVM Programming Model must describe how the
525 actions and attributes of different modes are distinguished from one another.

526 Implementation specific errors may result from attempts to use NVM.PM.FILE actions on files
527 that were created in NVM.FILE mode or vice versa. The mapping of each implementation to
528 the NVM Programming Model shall specify any limitations related multi-mode access.

529 **6.3 Device state at system startup**

530 Prior to use, a file system is associated with one or more volumes and/or NVM devices.

531 The NVM devices shall be in a state appropriate for use with file systems. For example, if
532 transparent RAID is part of the solution, components implementing RAID shall be active so the
533 file system sees a unified virtual device rather than individual RAID components.

534 **6.4 Secure erase**

535 Secure erase of a volume or device is an administrative act with no defined programming
536 model action.

537 **6.5 Allocation of space**

538 Following native operating system behavior, this programming model does not define specific
539 actions for allocating space. Most allocation behavior is hidden from the user of the file, volume
540 or device.

541 **6.6 Interaction with I/O devices**

542 Interaction between Persistent Memory and I/O devices (for example, DMA) shall be
543 consistent with native operating system interactions between devices and volatile memory.

544 **6.7 NVM State after a media or connection failure**

545 There is no action defined to determine the state of NVM for circumstances such as a media or
546 connection failure. Vendors may provide techniques such as redundancy algorithms to
547 address this, but the behavior is outside the scope of the programming model.

548 **6.8 Error handling for persistent memory**

549 The handling of errors in memory-mapped file implementations varies across operating
550 systems. Existing implementations support memory error reporting however there is not
551 sufficient similarity for a uniform approach to persistent memory error handling behavior.
552 Additional work is required to define an error handling approach. The following factors are to
553 be taken into account when dealing with errors.

- 554 • The application is in the best position to perform recovery as it may have access to
555 additional sources of data necessary to rewrite a bad memory address.
- 556 • Notification of a given memory error occurrence may need to be delivered to both kernel
557 and user space consumers (e.g., file system and application)
- 558 • Various hardware platforms have different capabilities to detect and report memory errors
- 559 • Attributes and possibly actions related to error handling behavior are needed in the NVM
560 Programming model

561 A proposal for persistent memory error handling appears as an appendix; see Annex B.

562 **6.9 Persistence domain**

563 NVM PM hardware supports the concept of a persistence domain. Once data has reached a
564 persistence domain, it may be recoverable during a process that results from a system restart.
565 Recoverability depends on whether the pattern of failures affecting the system during the
566 restart can be tolerated by the design and configuration of the persistence domain.

567 Multiple persistence domains may exist within the same system. It is an administrative act to
568 align persistence domains with volumes and/or file systems. This must be done in such a way
569 that NVM Programming Model behavior is assured from the point of view of each compliant
570 volume or file system.

571 **6.10 Aligned operations on fundamental data types**

572 Data alignment means putting the data at a memory offset equal to some multiple of the word
573 size, which increases the system's performance due to the way the CPU handles memory
574 (from Wikipedia "Data structure alignment"). Data types are *fundamental* when they are native
575 to programming languages or libraries.

576 Aligned operations on data types are usually exactly the same operations that under normal
577 operation become visible to other threads/data producers atomically. They are already well-
578 defined for most settings:

- 579 • Instruction Set Architectures already define them.
580 ○ E.g., for x86, MOV instructions with naturally aligned operands of at most 64 bits
581 qualify.
- 582 • They're generated by known high-level language constructs, e.g.:
583 ○ C++11 lock-free atomic<T>, C11 _Atomic(T), Java & C# volatile, OpenMP atomic
584 directives.

585 For optimal performance, fundamental data types fit within CPU cache lines.

586 **6.11 Common actions**

587 **6.11.1 NVM.COMMON.GET_ATTRIBUTE**

588 Requirement: mandatory

589 Get the value of one or more attributes. Implementations conforming to the specification shall
590 provide the get attribute behavior, but multiple programmatic approaches may be used.

591 **Inputs:**

- 592 • reference to appropriate instance (for example, reference to an NVM volume)
593 • attribute name

594 **Outputs:**

- 595 • value of attribute

596 The vendor's mapping document shall describe the possible errors reported for all applicable
597 programmatic approaches.

598 **6.11.2 NVM.COMMON.SET_ATTRIBUTE**

599 Requirement: optional

600 Note: at this time, no settable attributes are defined in this specification, but they may be
601 added in a future revision.

602 Set the value of one attribute. Implementations conforming to the specification shall provide
603 the set attribute behavior, but multiple programmatic approaches may be used.

604 **Inputs:**

- 605 • reference to appropriate instance
606 • attribute name
607 • value to be assigned to the attribute

608 The vendor's mapping document shall describe the possible errors reported for all applicable
609 programmatic approaches.

610 **6.12 Common attributes**

611 **6.12.1 NVM.COMMON.SUPPORTED_MODES**

612 Requirement: mandatory

613 SUPPORTED_MODES returns a list of the modes supported by the NVM implementation.

614 Possible values: NVM.BLOCK, NVM.FILE, NVM.PM.FILE, NVM.PM.VOLUME

615 NVM.COMMON.SET_ATTRIBUTE is not supported for
616 NVM.COMMON.SUPPORTED_MODES.

617 **6.12.2 NVM.COMMON.FILE_MODE**

618 Requirement: mandatory if NVM.FILE or NVM.PM.FILE is supported

619 Returns the supported file modes (NVM.FILE and/or NVM.PM.FILE) provided by a file system.

620 Target: a file path

621 Output value: a list of values: “NVM.FILE” and/or “NVM.PM.FILE”

622 See 6.2 Conformance to multiple file modes.

623 **6.13 Use cases**

624 **6.13.1 Application determines which mode is used to access a file system**

625 **Purpose/triggers:**

626 An application needs to determine whether the underlying file system conforms to NVM.FILE
627 mode, NVM.PM.FILE mode, or both.

628 **Scope/context:**

629 Some actions and attributes are defined differently in NVM.FILE and NVM.PM.FILE;
630 applications may need to be designed to handle these modes differently. This use case
631 describes steps in an application’s initialization logic to determine the mode(s) supported by
632 the implementation and set a variable indicating the preferred mode the application will use in
633 subsequent actions. This application prefers to use NVM.PM.FILE behavior if both modes are
634 supported.

635 **Success scenario:**

- 636 1) Invoke NVM.COMMON.GET_ATTRIBUTE (NVM.COMMON.FILE_MODE) targeting a
637 file path; the value returned provides information on which modes may be used to
638 access the data.
639 2) If the response includes “NVM.FILE”, then the actions and attributes described for the
640 NVM.FILE mode are supported. Set the preferred mode for this file system to
641 NVM.FILE.

642 3) If the response includes “NVM.PM.FILE”, then the actions and attributes described for
643 the NVM.PM.FILE mode are supported. Set the preferred mode for this file system to
644 NVM.PM.FILE.

645 **Outputs:**

646 **Postconditions:**

647 A variable representing the preferred mode for the file system has been initialized.

648 See also:

649 6.2 Conformance to multiple file modes

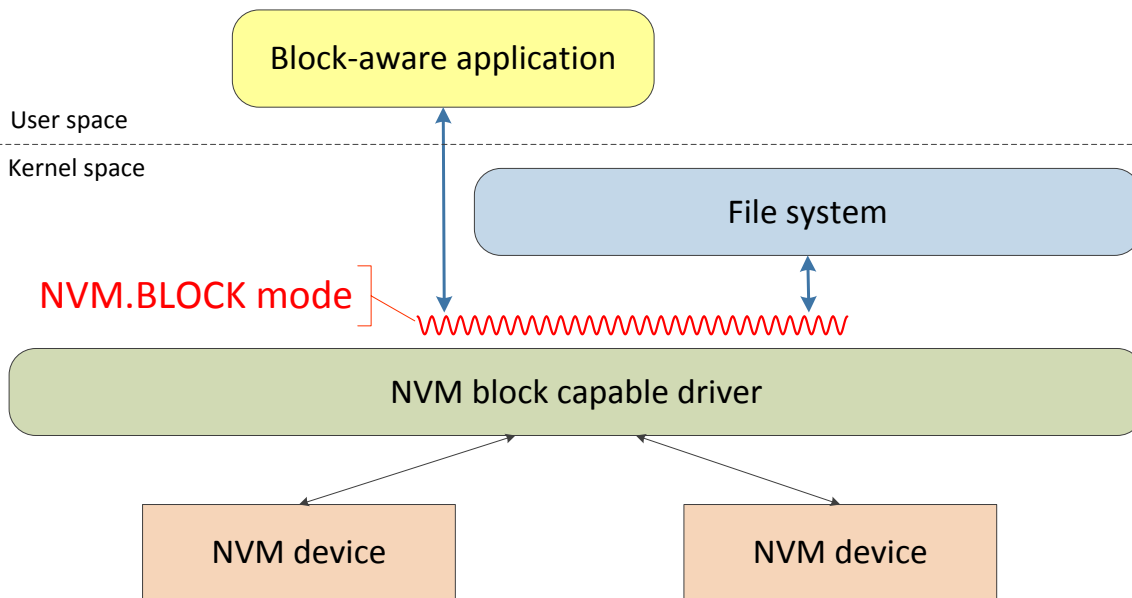
650 6.12.2 NVM.COMMON.FILE_MODE

651 **7 NVM.BLOCK mode**

652 **7.1 Overview**

653 NVM.BLOCK mode provides programming interfaces for NVM implementations behaving as
654 block devices. The programming interfaces include the native operating system behavior for
655 sending I/O commands to a block driver and adds NVM extensions. To support this mode, the
656 NVM devices are supported by an NVM block capable driver that provides the command
657 interface to the NVM. This specification does not document the native operating system block
658 programming capability; it is limited to the NVM extensions.

659 **Figure 6 NVM.BLOCK mode example**



660

661 Support for NVM.BLOCK mode requires that the NVM implementation support all behavior not
662 covered in this section consistently with the native operating system behavior for native block
663 devices.

664 The NVM extensions supported by this mode include:

- 665
- 666 • Discovery and use of atomic write and discard features
 - 667 • The discovery of granularities (length or alignment characteristics)
 - 668 • Discovery and use of per-block metadata used for verifying integrity
 - 669 • Discovery and use of ability for applications or operating system components to mark
670 blocks as unreadable

671 **7.1.1 Discovery and use of atomic write features**

672 Atomic Write support provides applications with the capability to assure that all the data for an
673 operation is written to the persistence domain or, if a failure occurs, it appears that no
674 operation took place. Applications may use atomic write operations to assure consistent

675 behavior during a failure condition or to assure consistency between multiple processes
676 accessing data simultaneously.

677

678 7.1.2 The discovery of granularities

679 Attributes are introduced to allow applications to discover granularities associated with NVM
680 devices.

681

682 7.1.3 Discovery and use of capability to mark blocks as unreadable

683 An action (NVM.BLOCK.SCAR) is defined allowing an application to mark blocks as
684 unreadable.

685

686 7.1.4 NVM.BLOCK consumers: operating system and applications

687 NVM.BLOCK behavior covers two types of software: NVM-aware operating system
688 components and block-optimized applications.

689 7.1.4.1 NVM.BLOCK operating system components

690 NVM-aware operating system components use block storage and have been enhanced to take
691 advantage of NVM features. Examples include file systems, logical volume managers,
692 software RAID, and hibernation logic.

693 7.1.4.2 Block-optimized applications

694 Block-optimized applications use a hybrid behavior utilizing files and file I/O operations, but
695 construct file I/O commands in order to cause drivers to issue desired block commands.
696 Operating systems and file systems typically provide mechanisms to enable block-optimized
697 application. The techniques are system specific, but may include:

- 698 • A mechanism for a block-optimized application to request that the file system move data
699 directly between the device and application memory, bypassing the buffering typically
700 provided by the file system.
- 701 • The operating system or file system may require the application to align requests on block
702 boundaries.

703 The file system and operating system may allow block-optimized applications to use memory-
704 mapped files.

705 7.1.4.3 Mapping documentation

706 NVM.BLOCK operating system components may use I/O commands restricted to kernel space
707 to send I/O commands to drivers. NVM.BLOCK applications may use a constrained set of file
708 I/O operations to send commands to drivers. As applicable, the implementation shall provide
709 documentation mapping actions and/or attributes for all supported techniques for NVM.BLOCK
710 behavior.

711 The implementation shall document the steps to utilize supported capabilities for block-
712 optimized applications and the constraints (e.g., block alignment) compared to NVM.FILE
713 behavior.

714 **7.2 Actions**

715 **7.2.1 Actions that apply across multiple modes**

716 The following actions apply to NVM.BLOCK mode as well as other modes.

717 NVM.COMMON.GET_ATTRIBUTE (see 6.11.1)

718 NVM.COMMON.SET_ATTRIBUTE (see 6.11.2)

719 **7.2.2 NVM.BLOCK.ATOMIC_WRITE**

720 Requirement: mandatory if ATOMIC_WRITE_CAPABLE (see 7.3.1) is true

721 Block-optimized applications or operating system components may use ATOMIC_WRITE to
722 assure consistent behavior during a power failure condition. This specification does not specify
723 the order in which this action occurs relative to other I/O operations, including other
724 ATOMIC_WRITE or ATOMIC_MULTIWRITE actions. This specification does not specify when
725 the data written becomes visible to other threads.

726 **Inputs:**

- 727 • the starting memory address
- 728 • a reference to the block device
- 729 • the starting block address
- 730 • the length

731 The interpretation of addresses and lengths (block or byte, alignment) should be consistent
732 with native write actions. Implementations shall provide documentation on the requirements for
733 specifying the starting addresses, block device, and length.

734 **Return values:**

- 735 • Success shall be returned if all blocks are updated in the persistence domain
- 736 • an error shall be reported if the length exceeds ATOMIC_WRITE_MAX_DATA_LENGTH
737 (see 7.3.3)
- 738 • an error shall be reported if the starting address is not evenly divisible by
739 ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY (see 7.3.4)
- 740 • an error shall be reported if the length is not evenly divisible by
741 ATOMIC_WRITE_LENGTH_GRANULARITY (see 7.3.5)
- 742 • If anything does or will prevent all of the blocks from being updated in the persistence
743 domain before completion of the operation, an error shall be reported and all the logical
744 blocks affected by the operation shall contain the data that was present before the device
745 server started processing the write operation (i.e., the old data, as if the atomic write
746 operation had no effect). If the NVM and processor are both impacted by a power failure,
747 no error will be returned since the execution context is lost.
- 748 • the different errors described above shall be discernible by the consumer and shall be
749 discernible from media errors

750 **Relevant attributes:**

751 ATOMIC_WRITE_MAX_DATA_LENGTH (see 7.3.3)

752 ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY (see 7.3.4)
753 ATOMIC_WRITE_LENGTH_GRANULARITY (see 7.3.5)
754 ATOMIC_WRITE_CAPABLE (see 7.3.1)

755 7.2.3 NVM.BLOCK.ATOMIC_MULTIWRITE

756 Requirement: mandatory if ATOMIC_MULTIWRITE_CAPABLE (see 7.3.6) is true

757 Block-optimized applications or operating system components may use
758 ATOMIC_MULTIWRITE to assure consistent behavior during a power failure condition. This
759 action allows a caller to write non-adjacent extents atomically. The caller of
760 ATOMIC_MULTIWRITE provides a Property Group List (see 4.4.5) where the properties
761 describe the memory and block extents (see Inputs below); all of the extents are written as a
762 single atomic operation. This specification does not specify the order in which this action
763 occurs relative to other I/O operations, including other ATOMIC_WRITE or
764 ATOMIC_MULTIWRITE actions. This specification does not specify when the data written
765 becomes visible to other threads.

766 **Inputs:**

767 A Property Group List (see 4.4.5) where the properties are:

- 768 • memory address starting address
- 769 • length of data to write (in bytes)
- 770 • a reference to the device being written to
- 771 • the starting LBA on the device

772 Each property group represents an I/O. The interpretation of addresses and lengths (block or
773 byte, alignment) should be consistent with native write actions. Implementations shall provide
774 documentation on the requirements for specifying the ranges.

775 **Return values:**

- 776 • Success shall be returned if all block ranges are updated in the persistence domain
- 777 • an error shall be reported if the block ranges overlap
- 778 • an error shall be reported if the total size of memory input ranges exceeds
779 ATOMIC_MULTIWRITE_MAX_DATA_LENGTH (see 7.3.8)
- 780 • an error shall be reported if the starting address in any input memory range is not evenly
781 divisible by ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY (see 7.3.9)
- 782 • an error shall be reported if the length in any input range is not evenly divisible by
783 ATOMIC_MULTIWRITE_LENGTH_GRANULARITY (see 7.3.10)
- 784 • If anything does or will prevent all of the writes from being applied to the persistence
785 domain before completion of the operation, an error shall be reported and all the logical
786 blocks affected by the operation shall contain the data that was present before the device
787 server started processing the write operation (i.e., the old data, as if the atomic write
788 operation had no effect). If the NVM and processor are both impacted by a power failure,
789 no error will be returned since the execution context is lost.
- 790 • the different errors described above shall be discernible by the consumer

- 791 **Relevant attributes:**
792 ATOMIC_MULTIWRITE_MAX_IOS (see 7.3.7)
793 ATOMIC_MULTIWRITE_MAX_DATA_LENGTH (see 7.3.8)
794 ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY (see 7.3.9)
795 ATOMIC_MULTIWRITE_LENGTH_GRANULARITY (see 7.3.10)
796 ATOMIC_MULTIWRITE_CAPABLE (see 7.3.6)

797 7.2.4 **NVM.BLOCK.DISCARD_IF_YOU_CAN**

798 Requirement: mandatory if DISCARD_IF_YOU_CAN_CAPABLE (see 7.3.17) is true

799 This action notifies the NVM device that some or all of the blocks which constitute a volume
800 are no longer needed by the application. This action is a hint to the device.

801 Although the application has logically discarded the data, it may later read this range. Since
802 the device is not required to physically discard the data, its response is undefined: it may
803 return successful response status along with unknown data (e.g., the old data, a default
804 “undefined” data, or random data), or it may return an unsuccessful response status with an
805 error.

806
807 Inputs: a range of blocks (starting LBA and length in logical blocks)

808 Status: Success indicates the request is accepted but not necessarily acted upon.

809 7.2.5 **NVM.BLOCK.DISCARD_IMMEDIATELY**

810 Requirement: mandatory if DISCARD_IMMEDIATELY_CAPABLE (see 7.3.18) is true

811 Requires that the data block be unmapped (see NVM.BLOCK.EXISTS 7.2.6) before the next
812 READ or WRITE reference even if garbage collection of the block has not occurred yet,

813 DISCARD_IMMEDIATELY commands cannot be acknowledged by the NVM device until the
814 DISCARD_IMMEDIATELY has been durably written to media in a way such that upon
815 recovery from a power-fail event, the block is guaranteed to remain discarded.

816 Inputs: a range of blocks (starting LBA and length in logical blocks)

817 The values returned by subsequent read operations are specified by the
818 DISCARD_IMMEDIATELY_RETURNS (see 7.3.19) attribute.

819 Status: Success indicates the request is completed.

820 See also EXISTS (7.2.6), DISCARD_IMMEDIATELY_RETURNS (7.3.19),
821 DISCARD_IMMEDIATELY_CAPABLE (7.3.18).

822 7.2.6 **NVM.BLOCK.EXISTS**

823 Requirement: mandatory if EXISTS_CAPABLE (see 7.3.12) is true

824 An NVM device may allocate storage through a thin provisioning mechanism or one of the
825 discard actions. As a result, a block can exist in one of three states:

- 826 • **Mapped**: the block has had data written to it
- 827 • **Unmapped**: the block has not been written, and there is no memory allocated
- 828 • **Allocated**: the block has not been written, but has memory allocated to it

829 The EXISTS action allows the NVM user to determine if a block has been allocated.

830 Inputs: an LBA

831 Output: the state (mapped, unmapped, or allocated) for the input block

832 Result: the status of the action

833 7.2.7 **NVM.BLOCK.SCAR**

834 Requirement: mandatory if SCAR_CAPABLE (see 7.3.13) is true

835 This action allows an application to request that subsequent reads from any of the blocks in
836 the address range will cause an error. This action uses an implementation-dependent means
837 to insure that all future reads to any given block from the scarred range will cause an error until
838 new data is stored to any given block in the range. A block stays scarred until it is updated by a
839 write operation.

840 Inputs: reference to a block volume, starting offset, length

841 Outputs: status

842 Relevant attributes:

843 NVM.BLOCK.SCAR_CAPABLE (7.3.13) – Indicates that the SCAR action is supported.

844 7.3 **Attributes**

845 7.3.1 **Attributes that apply across multiple modes**

846 The following attributes apply to NVM.BLOCK mode as well as other modes.

847 NVM.COMMON.SUPPORTED_MODES (see 6.12.1)

848 7.3.2 **NVM.BLOCK.ATOMIC_WRITE_CAPABLE**

849 Requirement: mandatory

850 This attribute indicates that the implementation is capable of the
851 NVM.BLOCK.ATOMIC_WRITE action.

852 7.3.3 **NVM.BLOCK.ATOMIC_WRITE_MAX_DATA_LENGTH**

853 Requirement: mandatory if ATOMIC_WRITE_CAPABLE (see 7.3.1) is true.

854 ATOMIC_WRITE_MAX_DATA_LENGTH is the maximum length of data that can be
855 transferred by an ATOMIC_WRITE action.

856 7.3.4 **NVM.BLOCK.ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY**

857 Requirement: mandatory if ATOMIC_WRITE_CAPABLE (see 7.3.1) is true.

858 ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY is the granularity of the starting
859 memory address for an ATOMIC_WRITE action. Address inputs to ATOMIC_WRITE shall be
860 evenly divisible by ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY.

861 7.3.5 **NVM.BLOCK.ATOMIC_WRITE_LENGTH_GRANULARITY**

862 Requirement: mandatory if ATOMIC_WRITE_CAPABLE (see 7.3.1) is true.

863 ATOMIC_WRITE_LENGTH_GRANULARITY is the granularity of the length of data transferred
864 by an ATOMIC_WRITE action. Length inputs to ATOMIC_WRITE shall be evenly divisible by
865 ATOMIC_WRITE_LENGTH_GRANULARITY.

866 7.3.6 **NVM.BLOCK.ATOMIC_MULTIWRITE_CAPABLE**

867 Requirement: mandatory

868 ATOMIC_MULTIWRITE_CAPABLE indicates that the implementation is capable of the
869 NVM.BLOCK.ATOMIC_MULTIWRITE action.

870 7.3.7 **NVM.BLOCK.ATOMIC_MULTIWRITE_MAX_IOS**

871 Requirement: mandatory if ATOMIC_MULTIWRITE_CAPABLE (see 7.3.6) is true

872 ATOMIC_MULTIWRITE_MAX_IOS is the maximum length of the number of IOs (i.e., the size
873 of the Property Group List) that can be transferred by an ATOMIC_MULTIWRITE action.

874 7.3.8 **NVM.BLOCK.ATOMIC_MULTIWRITE_MAX_DATA_LENGTH**

875 Requirement: mandatory if ATOMIC_MULTIWRITE_CAPABLE (see 7.3.6) is true

876 ATOMIC_MULTIWRITE_MAX_DATA_LENGTH is the maximum length of data that can be
877 transferred by an ATOMIC_MULTIWRITE action.

878 7.3.9 **NVM.BLOCK.ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY**

879 Requirement: mandatory if ATOMIC_MULTIWRITE_CAPABLE (see 7.3.6) is true

880 ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY is the granularity of the
881 starting address of ATOMIC_MULTIWRITE inputs. Address inputs to ATOMIC_MULTIWRITE
882 shall be evenly divisible by ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY.

883 7.3.10 **NVM.BLOCK.ATOMIC_MULTIWRITE_LENGTH_GRANULARITY**

884 Requirement: mandatory if ATOMIC_MULTIWRITE_CAPABLE (see 7.3.6) is true

885 ATOMIC_MULTIWRITE_LENGTH_GRANULARITY is the granularity of the length of
886 ATOMIC_MULTIWRITE inputs. Length inputs to ATOMIC_MULTIWRITE shall be evenly
887 divisible by ATOMIC_MULTIWRITE_LENGTH_GRANULARITY.

888 7.3.11 **NVM.BLOCK.WRITE_ATOMICITY_UNIT**

889 Requirement: mandatory

890 If a write is submitted of this size or less, the caller is guaranteed that if power is lost before the
891 data is completely written, then the NVM device shall ensure that all the logical blocks affected
892 by the operation contain the data that was present before the device server started processing
893 the write operation (i.e., the old data, as if the atomic write operation had no effect).

894 If the NVM device can't assure that at least one LOGICAL_BLOCKSIZE (see 7.3.14) extent
895 can be written atomically, WRITE_ATOMICITY_UNIT shall be set to zero.

896 The unit is NVM.BLOCK.LOGICAL_BLOCKSIZE (see 7.3.14).

897 7.3.12 **NVM.BLOCK.EXISTS_CAPABLE**

898 Requirement: mandatory

899 This attribute indicates that the implementation is capable of the NVM.BLOCK.EXISTS action.

900 7.3.13 **NVM.BLOCK.SCAR_CAPABLE**

901 Requirement: mandatory

902 This attribute indicates that the implementation is capable of the NVM.BLOCK.SCAR (see
903 7.2.7) action.

904 7.3.14 **NVM.BLOCK.LOGICAL_BLOCK_SIZE**

905 Requirement: mandatory

906 LOGICAL_BLOCK_SIZE is the smallest unit of data (in bytes) that may be logically read or
907 written from the NVM volume.

908 7.3.15 **NVM.BLOCK.PERFORMANCE_BLOCK_SIZE**

909 Requirement: mandatory

910 PERFORMANCE_BLOCK_SIZE is the recommended granule (in bytes) the caller should use
911 in I/O requests for optimal performance; starting addresses and lengths should be multiples of
912 this attribute. For example, this attribute may help minimizing device-implemented
913 read/modify/write behavior.

914 7.3.16 **NVM.BLOCK.ALLOCATION_BLOCK_SIZE**

915 Requirement: mandatory

916 ALLOCATION_BLOCK_SIZE is the recommended granule (in bytes) for allocation and
917 alignment of data. Allocations smaller than this attribute (even if they are multiples of
918 LOGICAL_BLOCK_SIZE) may work, but may not yield optimal lifespan.

919 7.3.17 **NVM.BLOCK.DISCARD_IF_YOU_CAN_CAPABLE**

920 Requirement: mandatory

921 DISCARD_IF_YOU_CAN_CAPABLE shall be set to true if the implementation supports
922 DISCARD_IF_YOU_CAN.

923 7.3.18 **NVM.BLOCK.DISCARD_IMMEDIATELY_CAPABLE**

924 Requirement: mandatory

925 Returns true if the implementation supports DISCARD_IMMEDIATELY.

926 7.3.19 **NVM.BLOCK.DISCARD_IMMEDIATELY_RETURNS**

927 Requirement: mandatory if DISCARD_IMMEDIATELY_CAPABLE (see 7.3.18) is true

928 The value returned from read operations to blocks specified by a DISCARD_IMMEDIATELY
929 action with no subsequent write operations. The possible values are:

- 930 • A value that is returned to each read of an unmapped block (see NVM.BLOCK.EXISTS
931 7.2.6) until the next write action
- 932 • Unspecified

933 7.3.20 **NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE**

934 Requirement: mandatory

935 FUNDAMENTAL_BLOCK_SIZE is the number of bytes that may become unavailable due to
936 an error on an NVM device.

937 A zero value means that the device is unable to provide a guarantee on the number of
938 adjacent bytes impacted by an error.

939 This attribute is relevant when the device does not support write atomicity.

940 If FUNDAMENTAL_BLOCK_SIZE is smaller than LOGICAL_BLOCK_SIZE (see 7.3.14), an
941 application may organize data in terms of FUNDAMENTAL_BLOCK_SIZE to avoid certain torn
942 write behavior. If FUNDAMENTAL_BLOCK_SIZE is larger than LOGICAL_BLOCK_SIZE, an
943 application may organize data in terms of FUNDAMENTAL_BLOCK_SIZE to assure two key
944 data items do not occupy an extent that is vulnerable to errors.

945 **7.4 Use cases**

946 **7.4.1 Flash as cache use case**

947 **Purpose/triggers:**

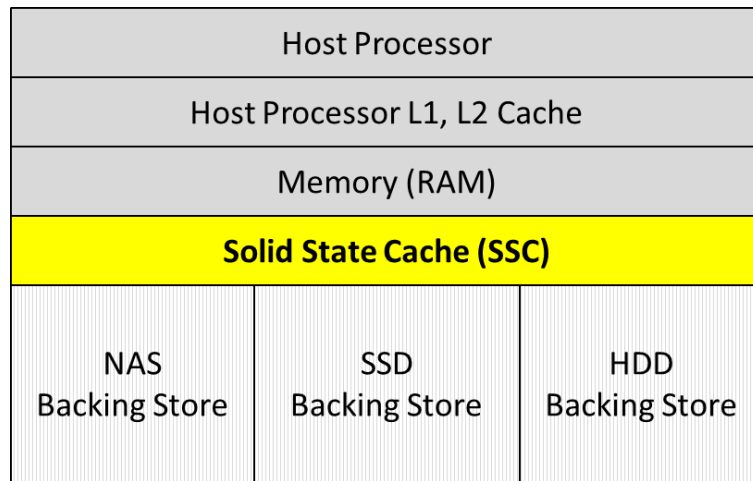
948 Use Flash based NVM as a data cache.

949 **Scope/context:**

950 Flash memory's fast random I/O performance and non-volatile characteristic make it a good
951 candidate as a Solid State Cache device (SSC). This use case is described in Figure 7 SSC in
952 a storage stack.

953

Figure 7 SSC in a storage stack



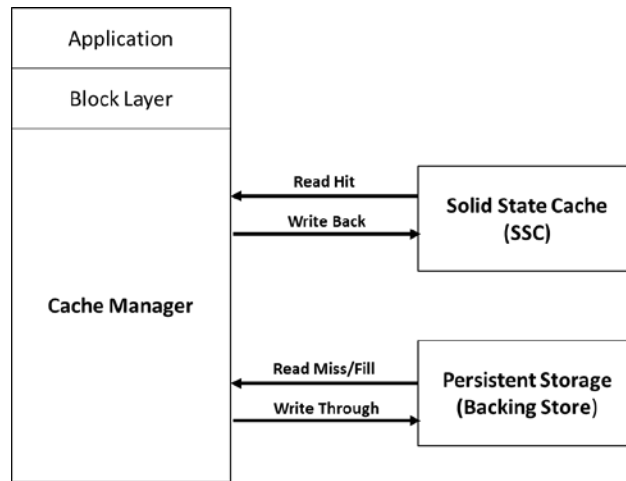
954

955

956 A possible software application is shown in Figure 8 SSC software cache application. In this
957 case, the cache manager employs the Solid State Cache to improve caching performance and
958 to maintain persistence and cache coherency across power fail.

959

Figure 8 SSC software cache application

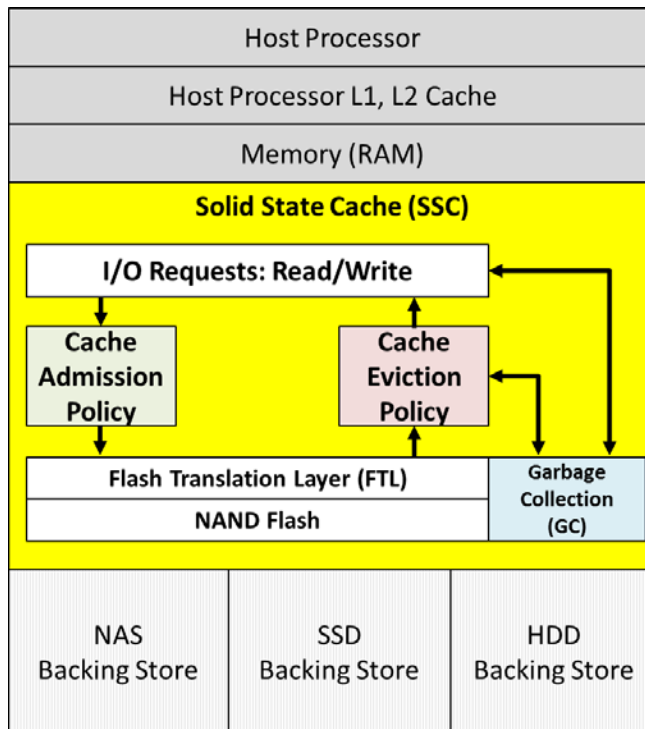


960

961 It is also possible to use an enhanced SSC to perform some of the functions that the cache
962 manager must normally contend with as shown in Figure 9 SSC with caching assistance.

963

Figure 9 SSC with caching assistance



964

965 In this use case, the Solid State Cache (SSC) provides a sparse address space that may be
966 much larger than the amount of physical NVM memory and manages the cache through its
967 own admission and eviction policies. The backing store is used to persist the data when the
968 cache becomes full. As a result, the block state for each block of virtual storage in the cache
969 must be maintained by the SSC. The SSC must also present a consistent cache interface that
970 can persist the cached data across a power fail and never returns stale data.

971 In either of these cases, two important extensions to existing storage commands must be
972 present:

973 **Eviction:** An explicit eviction mechanism is required to invalidate cached data in the
974 SSC to allow the cache manager to precisely control the contents of the SSC. This
975 means that the SSC must insure that the eviction is durable before completing the
976 request. This mechanism is generally referred to as a persistent trim. This is the
977 NVM.BLOCK.DISCARD_IMMEDIATELY functionality.

978 **Exists:** The EXISTS action allows the cache manager to determine the state of a block,
979 or of a range of blocks, in the SSC. This action is used to test for the presence of data in
980 the cache, or to determine which blocks in the SSC are dirty and need to be flushed to
981 backing storage. This is the NVM.BLOCK.EXISTS functionality.

982 The most efficient mechanism for a cache manager would be to simply read the requested
983 data from the SSC which would the return either the data or an error indicated that the
984 requested data was not in the cache. This approach is problematic, since most storage drivers
985 and software require reads to be successful and complete by returning data - not an error.
986 Device that return errors for normal read operations are usually put into an offline state by the
987 system drivers. Further, the data that a read returns must be consistent from one read
988 operation to the next, provided that no intervening writes occur. As a result, a two stage
989 process must be used by the cache manager. The cache manager first issues an EXISTS
990 action to determine if the requested data is present in the cache. Based on the result, the
991 cache manager decides whether to read the data from the SSC or from the backing storage.

992 **Success scenario:**

993 The requested data is successfully read from or written to the SSC.

994 **See also:**

- 995 • 7.2.5 NVM.BLOCK.DISCARD_IMMEDIATELY
- 996 • 7.2.6 NVM.BLOCK.EXISTS
- 997 • Ptrim() + Exists(): Exposing New FTL Primitives to Applications, David Nellans, Michael
998 Zappe, Jens Axboe, David Flynn, 2011 Non-Volatile Memory Workshop. See:
999 <http://david.nellans.org/files/NVMW-2011.pdf>
- 1000 • FlashTier: a Lightweight, Consistent, and Durable Storage Cache, Mohit Saxena,
1001 Michael M. Swift and Yiying Zhang, University of Wisconsin-Madison. See:
1002 http://pages.cs.wisc.edu/~swift/papers/eurosys12_flashtier.pdf
- 1003 HEC: Improving Endurance of High Performance Flash-based Cache Devices, Jingpei
1004 Yang, Ned Plasson, Greg Gillis, Nisha Talagala, Swaminathan Sundararaman, Robert
1005 Wood, Fusion-io, Inc., SYSTOR '13, June 30 - July 02 2013, Haifa, Israel
- 1006 • Unioning of the Buffer Cache and Journaling Layers with Non-volatile Memory, Eunji
1007 Lee, Hyokyung Bahn, and Sam H. Noh. See:
1008 https://www.usenix.org/system/files/conference/fast13/fast13-final114_0.pdf

1009 7.4.2 SCAR use case

1010 **Purpose/triggers:**

1011 Demonstrate the use of the SCAR action

1012 **Scope/context:**

1013 This generic use case for SCAR involves two processes.

- 1014 • The “detect block errors process” detects errors in certain NVM blocks, and uses SCAR to
1015 communicate to other processes that the contents of these blocks cannot be reliably read,
1016 but can be safely re-written.
- 1017 • The “recover process” sees the error reported as the result of SCAR. If this process can
1018 regenerate the contents of the block, the application can continue with no error.

1019 For this use case, the “detect block errors process” is a RAID component doing a background
1020 scan of NVM blocks. In this case, the NVM is not in a redundant RAID configuration so block
1021 READ errors can’t be transparently recovered. The “recover process” is a cache component
1022 using the NVM as a cache for RAID volumes. Upon receipt of the SCAR error on a read, this
1023 component evaluates whether the block contents also reside on the cached volume; if so, it
1024 can copy the corresponding volume block to the NVM. This write to NVM will clear the SCAR
1025 error condition.

1026 **Preconditions:**

1027 The “detect block errors process” detected errors in certain NVM blocks, and used SCAR to
1028 mark these blocks.

1029 **Success scenario:**

- 1030 1. The cache manager intercepts a read request from an application
- 1031 2. The read request to the NVM cache returns a status indicating the requested blocks
1032 have been marked by a SCAR action
- 1033 3. The cache manager uses an implementation-specific technique and determines the
1034 blocks marked by a SCAR are also available on the cached volume
- 1035 4. The cache manager copies the blocks from the cached volume to the NVM
- 1036 5. The cache manager returns the requested block to the application with a status
1037 indicating the read succeeded

1038 **Postconditions:**

1039 The blocks previously marked with a SCAR action have been repaired.

1040 **Failure Scenario:**

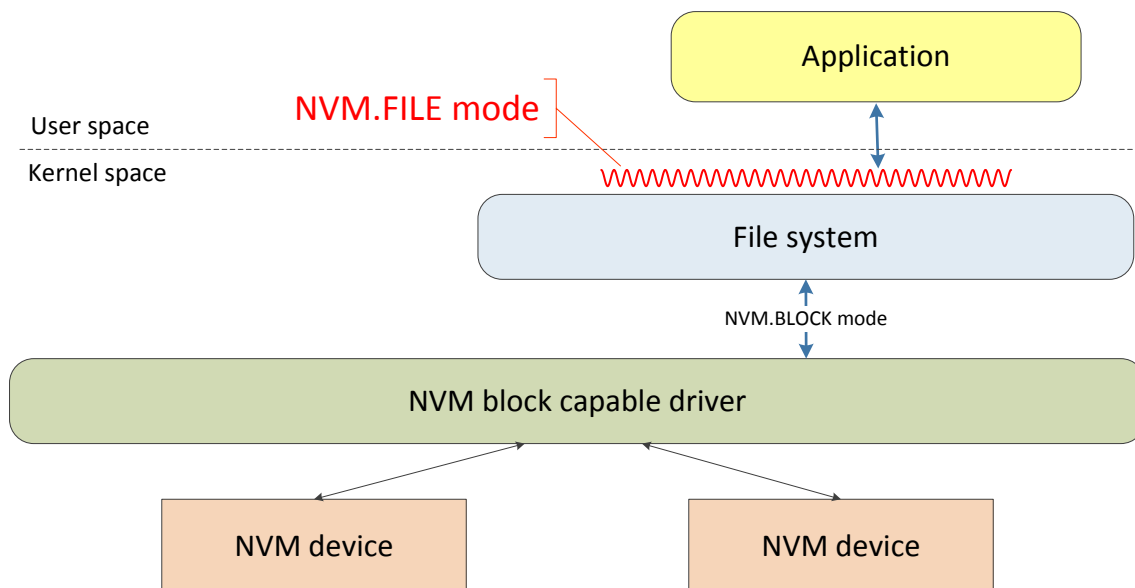
- 1041 1. In Success Scenario step 3 or 4, the cache manager discovers the corresponding
1042 blocks on the volume are invalid or cannot be read.
- 1043 2. The cache manager returns a status to the application indicating the blocks cannot be
1044 read.

1045 **8 NVM.FILE mode**

1046 **8.1 Overview**

1047 NVM.FILE mode addresses NVM-specification extensions to native file I/O behavior (the
1048 approach to I/O used by most applications). Support for NVM.FILE mode requires that the
1049 NVM solution ought to support all behavior not covered in this section consistently with the
1050 native operating system behavior for native block devices.

1051 **Figure 10 NVM.FILE mode example**



1052
1053 **8.1.1 Discovery and use of atomic write features**

1054 Atomic Write features in NVM.FILE mode are available to block-optimized applications (see
1055 7.1.4.2 Block-optimized applications).

1056 **8.1.2 The discovery of granularities**

1057 The NVM.FILE mode exposes the same granularity attributes as NVM.BLOCK.

1058 **8.1.3 Relationship between native file APIs and NVM.BLOCK.DISCARD**

1059 NVM.FILE mode does not define specific action that cause TRIM/DISCARD behavior. File
1060 systems may invoke NVM.BLOCK DISCARD actions when native operating system APIs
1061 (such as POSIX truncate or Windows SetEndOfFile).

1062 **8.2 Actions**

1063 **8.2.1 Actions that apply across multiple modes**

1064 The following actions apply to NVM.FILE mode as well as other modes.

1065 NVM.COMMON.GET_ATTRIBUTE (see 6.11.1)

1066 NVM.COMMON.SET_ATTRIBUTE (see 6.11.2)

1067 **8.2.2 NVM.FILE.ATOMIC_WRITE**

1068 Requirement: mandatory if ATOMIC_WRITE_CAPABLE (see 8.3.2) is true

1069 Block-optimized applications may use ATOMIC_WRITE to assure consistent behavior during a
1070 failure condition. This specification does not specify the order in which this action occurs
1071 relative to other I/O operations, including other ATOMIC_WRITE and ATOMIC_MULTIWRITE
1072 actions. This specification does not specify when the data written becomes visible to other
1073 threads.

1074 The inputs, outputs, and error conditions are similar to those for
1075 NVM.BLOCK.ATOMIC_WRITE, but typically the application provides file names and file
1076 relative block addresses rather than device name and LBA.

1077 Relevant attributes:

1078 ATOMIC_WRITE_MAX_DATA_LENGTH
1079 ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY
1080 ATOMIC_WRITE_LENGTH_GRANULARITY
1081 ATOMIC_WRITE_CAPABLE

1082 **8.2.3 NVM.FILE.ATOMIC_MULTIWRITE**

1083 Requirement: mandatory if ATOMIC_MULTIWRITE_CAPABLE (see 8.3.6) is true

1084 Block-optimized applications may use ATOMIC_MULTIWRITE to assure consistent behavior
1085 during a failure condition. This action allows a caller to write non-adjacent extents atomically.
1086 The caller of ATOMIC_MULTIWRITE provides properties defining memory and block extents;
1087 all of the extents are written as a single atomic operation. This specification does not specify
1088 the order in which this action occurs relative to other I/O operations, including other
1089 ATOMIC_WRITE and ATOMIC_MULTIWRITE actions. This specification does not specify
1090 when the data written becomes visible to other threads.

1091 The inputs, outputs, and error conditions are similar to those for
1092 NVM.BLOCK.ATOMIC_MULTIWRITE, but typically the application provides file names and file
1093 relative block addresses rather than device name and LBA.

1094 Relevant attributes:

1095 ATOMIC_MULTIWRITE_MAX_IOS
1096 ATOMIC_MULTIWRITE_MAX_DATA_LENGTH
1097 ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY
1098 ATOMIC_MULTIWRITE_LENGTH_GRANULARITY
1099 ATOMIC_MULTIWRITE_CAPABLE

1100 **8.3 Attributes**

1101 Some attributes share behavior with their NVM.BLOCK counterparts. NVM.FILE attributes are
1102 provided because the actual values may change due to the implementation of the file system.

1103 **8.3.1 Attributes that apply across multiple modes**

1104 The following attributes apply to NVM.FILE mode as well as other modes.

1105 NVM.COMMON.SUPPORTED_MODES (see 6.12.1)

1106 NVM.COMMON.FILE_MODE (see 6.12.2)

1107 **8.3.2 NVM.FILE.ATOMIC_WRITE_CAPABLE**

1108 Requirement: mandatory

1109 This attribute indicates that the implementation is capable of the
1110 NVM.BLOCK.ATOMIC_WRITE action.

1111 **8.3.3 NVM.FILE.ATOMIC_WRITE_MAX_DATA_LENGTH**

1112 Requirement: mandatory

1113 ATOMIC_WRITE_MAX_DATA_LENGTH is the maximum length of data that can be
1114 transferred by an ATOMIC_WRITE action.

1115 **8.3.4 NVM.FILE.ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY**

1116 Requirement: mandatory

1117 ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY is the granularity of the starting
1118 memory address for an ATOMIC_WRITE action. Address inputs to ATOMIC_WRITE shall be
1119 evenly divisible by ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY.

1120 **8.3.5 NVM.FILE.ATOMIC_WRITE_LENGTH_GRANULARITY**

1121 Requirement: mandatory

1122 ATOMIC_WRITE_LENGTH_GRANULARITY is the granularity of the length of data transferred
1123 by an ATOMIC_WRITE action. Length inputs to ATOMIC_WRITE shall be evenly divisible by
1124 ATOMIC_WRITE_STARTING_ADDRESS_GRANULARITY.

1125 **8.3.6 NVM.FILE.ATOMIC_MULTIWRITE_CAPABLE**

1126 Requirement: mandatory

1127 This attribute indicates that the implementation is capable of the
1128 NVM.FILE.ATOMIC_MULTIWRITE action.

1129 **8.3.7 NVM.FILE.ATOMIC_MULTIWRITE_MAX_IOS**

1130 Requirement: mandatory

1131 ATOMIC_MULTIWRITE_MAX_IOS is the maximum length of the number of IOs (i.e., the size
1132 of the Property Group List) that can be transferred by an ATOMIC_MULTIWRITE action.

1133 8.3.8 **NVM.FILE.ATOMIC_MULTIWRITE_MAX_DATA_LENGTH**

1134 Requirement: mandatory

1135 ATOMIC_MULTIWRITE_MAX_DATA_LENGTH is the maximum length of data that can be
1136 transferred by an ATOMIC_MULTIWRITE action.

1137 8.3.9 **NVM.FILE.ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY**

1138 Requirement: mandatory

1139 ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY is the granularity of the
1140 starting address of ATOMIC_MULTIWRITE inputs. Address inputs to ATOMIC_MULTIWRITE
1141 shall be evenly divisible by ATOMIC_MULTIWRITE_STARTING_ADDRESS_GRANULARITY.

1142 8.3.10 **NVM.FILE.ATOMIC_MULTIWRITE_LENGTH_GRANULARITY**

1143 Requirement: mandatory

1144 ATOMIC_MULTIWRITE_LENGTH_GRANULARITY is the granularity of the length of
1145 ATOMIC_MULTIWRITE inputs. Length inputs to ATOMIC_MULTIWRITE shall be evenly
1146 divisible by ATOMIC_MULTIWRITE_LENGTH_GRANULARITY.

1147 8.3.11 **NVM.FILE.WRITE_ATOMICITY_UNIT**

1148 See 7.3.11 NVM.BLOCK.WRITE_ATOMICITY_UNIT

1149 8.3.12 **NVM.FILE.LOGICAL_BLOCK_SIZE**

1150 See 7.3.14 NVM.BLOCK.LOGICAL_BLOCK_SIZE

1151 8.3.13 **NVM.FILE.PERFORMANCE_BLOCK_SIZE**

1152 See 7.3.15 NVM.BLOCK.PERFORMANCE_BLOCK_SIZE

1153 8.3.14 **NVM.FILE.LOGICAL_ALLOCATION_SIZE**

1154 See 7.3.16 NVM.BLOCK.ALLOCATION_BLOCK_SIZE

1155 8.3.15 **NVM.FILE.FUNDAMENTAL_BLOCK_SIZE**

1156 See 7.3.20 NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE

1157 **8.4 Use cases**

1158 8.4.1 **Block-optimized application updates record**

1159 Update a record in a file without using a memory-mapped file

1160 **Purpose/triggers:**
1161 An application using block NVM updates an existing record. The application requests that the
1162 file system bypass cache; the application conforms to native API requirements when
1163 bypassing cache – this may mean that read and write actions must use multiples of a page
1164 cache size. For simplicity, this application uses fixed size records. The record size is defined
1165 by application data considerations, not disk or page block sizes. The application factors in the
1166 PERFORMANCE_BLOCK_SIZE granularity to avoid device-side inefficiencies such as
1167 read/modify/write.

1168 **Scope/context:**
1169 Block NVM context; this shows basic behavior.

1170 **Preconditions:**
1171 - The administrator created a file and provided its name to the application; this name is
1172 accessible to the application – perhaps in a configuration file
1173 - The application has populated the contents of this file
1174 - The file is not in use at the start of this use case (no sharing considerations)

1175 **Inputs:**
1176 The content of the record, the location (relative to the file) where the record resides

1177 **Success scenario:**
1178 1) The application uses the native OPEN action, passing in the file name and specifying
1179 appropriate options to bypass the file system cache
1180 2) The application acquires the device’s optimal I/O granule size by using the
1181 GET_ATTRIBUTE action for the PERFORMANCE_BLOCK_SIZE.
1182 3) The application allocates sufficient memory to contain all of the blocks occupied by the
1183 record to be updated.
1184 a. The application determines the offset within the starting block of the record and uses
1185 the length of the block to determine the number of partial blocks.
1186 b. The application allocates sufficient memory for the record plus enough additional
1187 memory to accommodate any partial blocks.
1188 c. If necessary, the memory size is increased to assure that the starting address and
1189 length read and write actions are multiples of PERFORMANCE_BLOCK_SIZE.
1190 4) The application uses the native READ action to read the record by specifying the starting
1191 disk address and the length (the same length as the allocated memory buffer). The
1192 application also provides the allocated memory address; this is where the read action will
1193 put the record.
1194 5) The application updates the record in the memory buffer per the inputs
1195 6) The application uses the native write action to write the updated block(s) to the same disk
1196 location they were read from.
1197 7) The application uses the native file SYNC action to assure the updated blocks are written to
1198 the persistence domain

1199 8) The application uses the native CLOSE action to clean up.

1200 **Failure Scenario 1:**

1201 The native read action reports a hardware error. If the unreadable block corresponds to blocks
1202 being updated, the application may attempt recovery (write/read/verify), or preventative
1203 maintenance (scar the unreadable blocks). If the unreadable blocks are needed for a
1204 read/modify/write update and the application lacks an alternate source; the application may
1205 inform the user that an unrecoverable hardware error has occurred.

1206 **Failure Scenario 2:**

1207 The native write action reports a hardware error. The application may be able to recover by
1208 rewriting the block. If the rewrite fails, the application may be able to scar the bad block and
1209 write to a different location.

1210 **Postconditions:**

1211 The record is updated.

1212 **8.4.2 Atomic write use case**

1213 **Purpose/triggers:**

1214 Used by a block-optimized application (see Block-optimized applications) striving for durability
1215 of on-disk data

1216 **Scope/context:**

1217 Assure a record is written to disk in a way that torn writes can be detected and rolled back (if
1218 necessary). If the device supports atomic writes, they will be used. If not, a double write buffer
1219 is used.

1220 **Preconditions:**

1221 The application has taken steps (based on NVM.BLOCK attributes) to assure the record being
1222 written has an optimal memory starting address, starting disk LBA and length.

1223 **Success scenario:**

- 1224 • Use GET_ATTRIBUTE to determine whether the device is ATOMIC_WRITE_CAPABLE
- 1225 (or ATOMIC_MULTIWRITE_CAPABLE)
- 1226 • Is so, use the appropriate atomic write action to write the record to NVM
- 1227 • If the device does not support atomic write, then
 - 1228 ○ Write the page to the double write buffer
 - 1229 ○ Wait for the write to complete
 - 1230 ○ Write the page to the final destination
- 1231 • At application startup, if the device does not support atomic write
 - 1232 • Scan the double write buffer and for each valid page in the buffer check if the page
 - 1233 in the data file is valid too.

1234 **Postconditions:**
1235 After application startup recovery steps, there are no inconsistent records on disk after a failure
1236 caused the application (and possibly system) to restart.

1237 8.4.3 Block and File Transaction Logging

1238 **Purpose/Triggers:**

1239 An application developer wishes to implement a transaction log that maintains data integrity
1240 through system crashes, system resets, and power failures. The underlying storage is block-
1241 granular, although it may be accessed via a file system that simulates byte-granular access to
1242 files.

1243 **Scope/Context:**

1244 NVM.BLOCK or NVM.FILE (all the NVM.BLOCK attributes mentioned in the use case are also
1245 defined for NVM.FILE mode).

1246 For notational convenience, this use case will use the term “file” to apply to either a file in the
1247 conventional sense which is accessed through the NVM.FILE interface, or a specific subset of
1248 blocks residing on a block device which are accessed through the NVM.BLOCK interface.

1249 **Inputs:**

- 1250 • A set of changes to the persistent state to be applied as a single transaction.
- 1251 • The data and log files.

1252 **Outputs:**

- 1253 • An indication of transaction commit or abort

1254 **Postconditions:**

- 1255 • If an abort indication was returned, the data was not committed and the previous
1256 contents have not been modified.
- 1257 • If a commit indication was returned, the data has been entirely committed.
- 1258 • After a system crash, reset, or power failure followed by system restart and execution of
1259 the application transaction recovery process, the data has either been entirely
1260 committed or the previous contents have not been modified.

1261 **Success Scenario:**

1262 The application transaction logic uses a log file in combination with its data file to atomically
1263 update the persistent state of the application. The log may implement a before-image log or a
1264 write-ahead log. The application transaction logic should configure itself to handle torn or
1265 interrupted writes to the log or data files.

1266 **8.4.3.1 NVM.BLOCK.WRITE_ATOMICITY_UNIT >= 1**

1267 If the NVM.BLOCK.WRITE_ATOMICITY_UNIT is one or greater, then writes of a single logical
1268 block cannot be torn or interrupted.

1269 In this case, if the log or data record size is less than or equal to the
1270 NVM.BLOCK.LOGICAL_BLOCK_SIZE, the application need not handle torn or interrupted
1271 writes to the log or data files.

1272 If the log or data record size is greater than the NVM.BLOCK.LOGICAL_BLOCK_SIZE, the
1273 application should be prepared to detect a torn write of the record and either discard or recover
1274 such a torn record during the recovery process. One common way of detecting such a torn
1275 write is for the application to compute hash of the record and record the hash in the record.
1276 Upon reading the record, the application re-computes the hash and compares it with the
1277 recorded hash; if they do not match, the record has been torn. Another method is for the
1278 application to insert the transaction identifier within each logical block. Upon reading the
1279 record, the application compares the transaction identifiers in each logical block; if they do not
1280 match, the record has been torn. Another method is for the application to use the
1281 NVM.BLOCK.ATOMIC_WRITE action to perform the writes of the record.

1282 **8.4.3.2 NVM.BLOCK.WRITE_ATOMICITY_UNIT = 0**

1283 If the NVM.BLOCK.WRITE_ATOMICITY_UNIT is zero, then writes of a single logical block can
1284 be torn or interrupted and the application should handle torn or interrupted writes to the log or
1285 data files.

1286 In this case, if a logical block were to contain data from more than one log or data record, a
1287 torn or interrupted write could corrupt a previously-written record. To prevent propagating an
1288 error beyond the record currently being written, the application aligns the log or data records
1289 with the NVM.BLOCK.LOGICAL_BLOCK_SIZE and pads the record size to be an integral
1290 multiple of NVM.BLOCK.LOGICAL_BLOCK_SIZE. This prevents more than one record from
1291 residing in the same logical block and therefore a torn or interrupted write may only corrupt the
1292 record being written.

1293 *8.4.3.2.1 NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE >=* 1294 *NVM.BLOCK.LOGICAL_BLOCK_SIZE*

1295 If the NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE is greater than or equal to the
1296 NVM.BLOCK.LOGICAL_BLOCK_SIZE, the application should be prepared to handle an
1297 interrupted write. An interrupted write results when the write of a single
1298 NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE unit is interrupted by a system crash, system
1299 reset, or power failure. As a result of an interrupted write, the NVM device may return an error
1300 when any of the logical blocks comprising the NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE
1301 unit are read. (See also SQLite.org, *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>.)
1302 This presents a danger to the integrity of previously written records that, while residing in
1303 differing logical blocks, share the same fundamental block. An interrupted write may prevent
1304 the reading of those previously written records in addition to preventing the read of the record
1305 in the process of being written.

1306 One common way of protecting previously written records from damage due to an interrupted
1307 write is to align the log or data records with the NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE
1308 and pad the record size to be an integral multiple of
1309 NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE. This prevents more than one record from

1310 residing in the same fundamental block. The application should be prepared to discard or
1311 recover the record if the NVM device returns an error when subsequently reading the record
1312 during the recovery process.

1313 8.4.3.2.2 *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE <* 1314 *NVM.BLOCK.LOGICAL_BLOCK_SIZE*

1315 If the *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE* is less than the
1316 *NVM.BLOCK.LOGICAL_BLOCK_SIZE*, the application should be prepared to handle both
1317 interrupted writes and torn writes within a logical block.

1318 As a result of an interrupted write, the NVM device may return an error when the logical block
1319 containing the *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE* unit which was being written at
1320 the time of the system crash, system reset, or power failure is subsequently read. The
1321 application should be prepared to discard or recover the record in the logical block if the NVM
1322 device returns an error when subsequently reading the logical block during the recovery
1323 process.

1324 A torn write results when an integral number of *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE*
1325 units are written to the NVM device but the entire *NVM.BLOCK.LOGICAL_BLOCK_SIZE* has
1326 not been written. In this case, the NVM device may not return an error when the logical block is
1327 read. The application should therefore be prepared to detect a torn write of a logical block and
1328 either discard or recover such a torn record during the recovery process. One common way of
1329 detecting such a torn write is for the application to compute a hash of the record and record the
1330 hash in the record. Upon reading the record, the application re-computes the hash and
1331 compares it with the recorded hash; if they do not match, a logical block within the record has
1332 been torn. Another method is for the application to insert the transaction identifier within each
1333 *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE* unit. Upon reading the record, the application
1334 compares the transaction identifiers in each *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE*
1335 unit; if they do not match, the logical block has been torn.

1336 8.4.3.2.3 *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE = 0*

1337 If the *NVM.BLOCK.FUNDAMENTAL_BLOCK_SIZE* is zero, the application lacks sufficient
1338 information to handle torn or interrupted writes to the log or data files.

1339 **Failure Scenarios:**

1340 Consider the recovery of an error resulting from an interrupted write on a device where the
1341 *NVM.BLOCK.WRITE_ATOMICITY_UNIT* is zero. This error may be persistent and may be
1342 returned whenever the affected block is read. To repair this error, the application should be
1343 prepared to overwrite such a block.

1344 One common way of ensuring that the application will overwrite a block is by assigning it to the
1345 set of internal free space managed by the application, which is never read and is available to
1346 be allocated and overwritten at some point in the future. For example, the block may be part of
1347 a circular log. If the block is marked as free, the transaction log logic will eventually allocate
1348 and overwrite that block as records are written to the log.

1349 Another common way is to record either a before-image or after-image of a data block in a log.
1350 During recovery after a system crash, system reset, or power failure, the application replays
1351 the records in the log and overwrites the data block with either the before-image contents or
1352 the after-image contents.

1353 **See also:**

- 1354 • SQLite.org, *Atomic Commit in SQLite*, <http://www.sqlite.org/atomiccommit.html>
- 1355 • SQLite.org, *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>
- 1356 • SQLite.org, *Write-Ahead Logging*, <http://www.sqlite.org/wal.html>

1357 **9 NVM.PM.VOLUME mode**

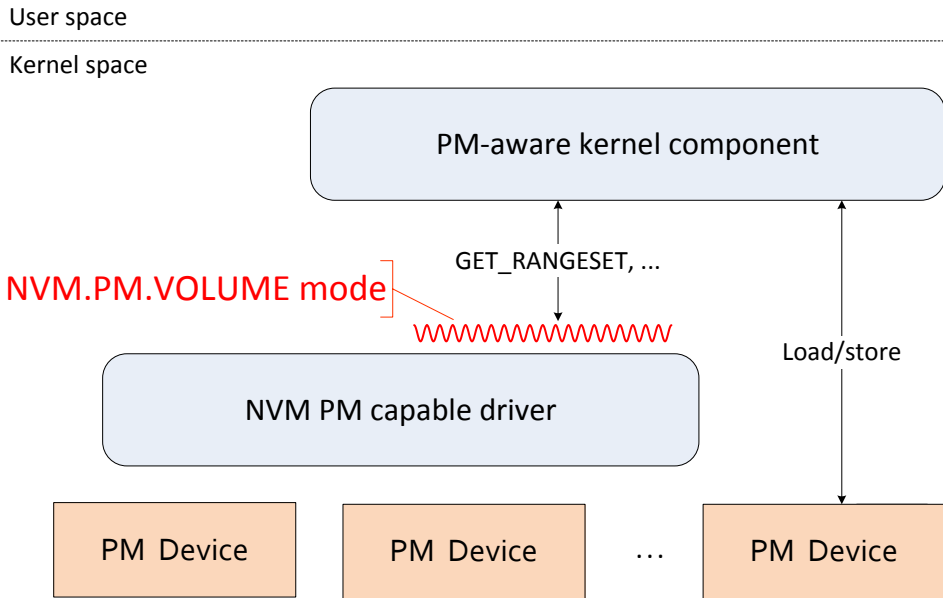
1358 **9.1 Overview**

1359 NVM.PM.VOLUME mode describes the behavior to be consumed by operating system
1360 abstractions such as file systems or pseudo-block devices that build their functionality by
1361 directly accessing persistent memory. NVM.PM.VOLUME mode provides a software
1362 abstraction (a PM volume) for persistent memory hardware and profiles functionality for
1363 operating system components including:

- 1364 • list of physical address ranges associated with each PM volume

1366 The PM volume provides memory mapped capability in a fashion that traditional CPU load and
1367 store operations are possible. This PM volume may be provided via the memory channel of the
1368 CPU or via a PCIe memory mapping or other methods. Note that there should not be a
1369 requirement for an operating system context switch for access to the PM volume.

1370 **Figure 11 NVM.PM.VOLUME mode example**



1371

1372 **9.2 Actions**

1373 **9.2.1 Actions that apply across multiple modes**

1374 The following actions apply to NVM.PM.VOLUME mode as well as other modes.

- 1375 NVM.COMMON.GET_ATTRIBUTE (see 6.11.1)
- 1376 NVM.COMMON.SET_ATTRIBUTE (see 6.11.2)

1377 **9.2.2 NVM.PM.VOLUME.GET_RANGESET**

1378 Requirement: mandatory

1379 The purpose of this action is to return a set of processor physical address ranges (and relate
1380 properties) representing all of the content for the identified volume.

1381 When interpreting the set of physical addresses as a contiguous, logical address range; the
1382 data underlying that logical address range will always be the same and in the same sequence
1383 across PM volume instantiations.

1384 Due to physical memory reconfiguration, the number and sizes of ranges may change in
1385 successive get ranges calls, however the total number of bytes in the sum of the ranges does
1386 not change, and the order of the bytes spanning all of the ranges does not change. The space
1387 defined by the list of ranges can always be addressed relative to a single base which
1388 represents the beginning of the first range.

1389 Input: a reference to the PM volume

1390 Returns a Property Group List (see 4.4.5) where the properties are:

- 1391 • starting physical address (byte address)
- 1392 • length (in bytes)
- 1393 • connection type – see below
- 1394 • sync type – see below

1395 For this revision of the specification, the following values (in text) are valid for connection type:

- 1396 • *“memory”*: for persistent memory attached to a system memory channel
- 1397 • *“PCIe”*: for persistent memory attached to a PCIe extension bus

1398 For this revision of the specification, the following values (in text) are valid for sync type:

- 1399 • *“none”*: no device-specific sync behavior is available – implies no entry to
1400 NVM.PM.VOLUME implementation is required for flushing
- 1401 • *“VIRTUAL_ADDRESS_SYNC”*: the caller needs to use VIRTUAL_ADDRESS_SYNC (see
1402 9.2.3) to assure sync is durable
- 1403 • *“PHYSICAL_ADDRESS_SYNC”*: the caller needs to use PHYSICAL_ADDRESS_SYNC
1404 (see 9.2.4) to assure sync is durable

1405 **9.2.3 NVM.PM.VOLUME.VIRTUAL_ADDRESS_SYNC**

1406 Requirement: optional

1407 The purpose of this action is to invoke device-specific actions to synchronize persistent
1408 memory content to assure durability and enable recovery by forcing data to reach the
1409 persistence domain. VIRTUAL_ADDRESS_SYNC is used by a caller that knows the
1410 addresses in the input range are virtual memory addresses.

1411 Input: virtual address and length (range)

1412 See also: PHYSICAL_ADDRESS_SYNC

- 1413 **9.2.4 NVM.PM.VOLUME.PHYSICAL_ADDRESS_SYNC**
- 1414 Requirement: optional
- 1415 The purpose of this action is to synchronize persistent memory content to assure durability and
1416 enable recovery by forcing data to reach the persistence domain. This action is used by a
1417 caller that knows the addresses in the input range are physical memory addresses.
- 1418 See also: VIRTUAL_ADDRESS_SYNC
- 1419 Input: physical address and length (range)
- 1420 **9.2.5 NVM.PM.VOLUME.DISCARD_IF_YOU_CAN**
- 1421 Requirement: mandatory if DISCARD_IF_YOU_CAN_CAPABLE (see 9.3.6) is true
- 1422 This action notifies the NVM device that the input range (volume offset and length) are no
1423 longer needed by the caller. This action may not result in any action by the device, depending
1424 on the implementation and the internal state of the device. This action is meant to allow the
1425 underlying device to optimize the data stored within the range. For example, the device can
1426 use this information in support of functionality like thin provisioning or wear-leveling.
- 1427 Inputs: a range of addresses (starting address and length in bytes). The address shall be a
1428 logical memory address offset from the beginning of the volume.
- 1429 Status: Success indicates the request is accepted but not necessarily acted upon.
- 1430 **9.2.6 NVM.PM.VOLUME.DISCARD_IMMEDIATELY**
- 1431 Requirement: mandatory if DISCARD_IMMEDIATELY_CAPABLE (see 9.3.7) is true
- 1432 This action notifies the NVM device that the input range (volume offset and length) are no
1433 longer needed by the caller. Similar to DISCARD_IF_YOU_CAN, but the implementation is
1434 required to unmap the range before the next READ or WRITE action, even if garbage
1435 collection of the range has not occurred yet.
- 1436 Inputs: a range of addresses (starting address and length in bytes). The address shall be a
1437 logical memory address offset from the beginning of the volume.
- 1438 The values returned by subsequent read operations are specified by the
1439 DISCARD_IMMEDIATELY_RETURNS (see 9.3.8) attribute.
- 1440 Status: Success indicates the request is completed.
- 1441 **9.2.7 NVM.PM.VOLUME.EXISTS**
- 1442 Requirement: mandatory if EXISTS_CAPABLE (see 9.3.9) is true
- 1443 A PM device may allocate storage through a thin provisioning mechanism or one of the discard
1444 actions. As a result, memory can exist in one of three states:

- 1445 • **Mapped:** the range has had data written to it
 - 1446 • **Unmapped:** the range has not been written, and there is no memory allocated
 - 1447 • **Allocated:** the range has not been written, but has memory allocated to it
- 1448 The EXISTS action allows the NVM user to determine if a range of bytes has been allocated.
- 1449 Inputs: a range of bytes (starting byte address and length in bytes)
- 1450 Output: a Property Group List (see 4.4.5) where the properties are the starting address, length
 1451 and state. State is a string equal to “mapped”, “unmapped”, or “allocated”.
- 1452 Result: the status of the action

1453 **9.3 Attributes**

1454 **9.3.1 Attributes that apply across multiple modes**

1455 The following attributes apply to NVM.PM.VOLUME mode as well as other modes.
 1456 NVM.COMMON.SUPPORTED_MODES (see 6.12.1)

1458 **9.3.2 NVM.PM.VOLUME.VOLUME_SIZE**

1459 Requirement: mandatory

1460 VOLUME_SIZE is the volume size in units of bytes. This shall be the same as the sum of the
 1461 lengths of the ranges returned by the GET_RANGES action.

1462 **9.3.3 NVM.PM.VOLUME.INTERRUPTED_STORE_ATOMICITY**

1463 Requirement: mandatory

1464 INTERRUPTED_STORE_ATOMICITY indicates whether the device supports power fail
 1465 atomicity of store actions.

1466 A value of true indicates that after a store interrupted by reset, power loss or system crash;
 1467 upon restart the contents of persistent memory reflect either the state before the store or the
 1468 state after the completed store. A value of false indicates that after a store interrupted by reset,
 1469 power loss or system crash, upon restart the contents of memory may be such that
 1470 subsequent loads may create exceptions depending on the value of the
 1471 FUNDAMENTAL_ERROR_RANGE attribute (see 9.3.4).

1472 **9.3.4 NVM.PM.VOLUME.FUNDAMENTAL_ERROR_RANGE**

1473 Requirement: mandatory

1474 FUNDAMENTAL_ERROR_RANGE is the number of bytes that may become unavailable due
 1475 to an error on an NVM device.

1476 This attribute is relevant when the device does not support write atomicity.

1477 A zero value means that the device is unable to provide a guarantee on the number of
 1478 adjacent bytes impacted by an error.

1479 A caller may organize data in terms of FUNDAMENTAL_ERROR_RANGE to avoid certain torn
 1480 write behavior.

1481 **9.3.5 NVM.PM.VOLUME.FUNDAMENTAL_ERROR_RANGE_OFFSET**

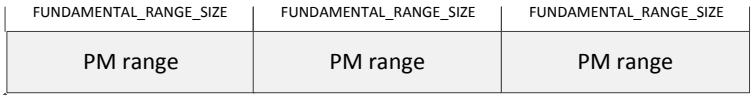
1482 Requirement: mandatory

1483 The number of bytes offset from the beginning of a volume range (as returned by
 1484 GET_RANGESET) before FUNDAMENTAL_ERROR_RANGE_SIZE intervals apply.

1485 A fundamental error range is not required to start at a byte address evenly divisible by
 1486 FUNDAMENTAL_ERROR_RANGE. FUNDAMENTAL_ERROR_RANGE_OFFSET shall be set
 1487 to the difference between the starting byte address of a fundamental error range rounded
 1488 down to a multiple of FUNDAMENTAL_ERROR_RANGE.

1489 Figure 12 Zero range offset example depicts an implementation where fundamental error
 1490 ranges start at bye address zero; the implementation shall return zero for
 1491 FUNDAMENTAL_ERROR_RANGE_OFFSET.

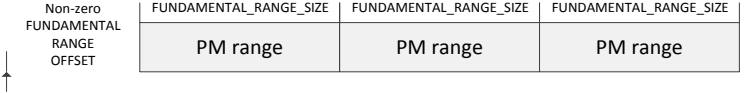
1492 **Figure 12 Zero range offset example**



1493

1494 Figure 13 Non-zero range offset example depicts an implementation where fundamental error
 1495 ranges start at a non-zero offset; the implementation shall return the difference between the
 1496 starting byte address of a fundamental error range rounded down to a multiple of
 1497 FUNDAMENTAL_ERROR_RANGE.

1498 **Figure 13 Non-zero range offset example**



1499

1500 **9.3.6 NVM.PM.VOLUME.DISCARD_IF_YOU_CAN_CAPABLE**

1501 Requirement: mandatory

1502 Returns true if the implementation supports DISCARD_IF_YOU_CAN.

1503 **9.3.7 NVM.PM.VOLUME.DISCARD_IMMEDIATELY_CAPABLE**

1504 Requirement: mandatory

1505 Returns true if the implementation supports DISCARD_IMMEDIATELY.

1506 9.3.8 **NVM.PM.VOLUME.DISCARD_IMMEDIATELY_RETURNS**

1507 Requirement: mandatory if DISCARD_IMMEDIATELY_CAPABLE (see 9.3.7) is true

1508 The value returned from read operations to bytes specified by a DISCARD_IMMEDIATELY
1509 action with no subsequent write operations. The possible values are:

- 1510 • A value that is returned to each load of bytes in an unmapped range until the next store
1511 action
- 1512 • Unspecified

1513 9.3.9 **NVM.PM.VOLUME.EXISTS_CAPABLE**

1514 Requirement: mandatory

1515 This attribute indicates that the implementation is capable of the NVM.PM.VOLUME.EXISTS
1516 action.

1517 **9.4 Use cases**

1518 9.4.1 **Initialization steps for a PM-aware file system**

1519 **Purpose/triggers:**

1520 Steps taken by a file system when a PM-aware volume is attached to a PM volume.

1521 **Scope/context:**

1522 NVM.PM.VOLUME mode

1523 **Preconditions:**

- 1524 • The administrator has defined a PM volume
- 1525 • The administrator has completed one-time steps to create a file system on the PM
1526 volume

1527 **Inputs:**

- 1528 • A reference to a PM volume
- 1529 • The name of a PM file system

1530 **Success scenario:**

- 1531 1. The file system issues a GET_RANGESET action to retrieve information about the
1532 ranges comprised by the PM volume.
- 1533 2. The file system uses the range information from GET_RANGESET to determine
1534 physical address range(s) and offset(s) of the file system's primary metadata (for
1535 example, the primary superblock), then loads appropriate metadata to determine no
1536 additional validity checking is needed.
- 1537 3. The file system sets a flag in the metadata indicating the file system is mounted by
1538 storing the updated status to the appropriate location

- 1539 a. If the range containing this location requires VIRTUAL_ADDRESS_SYNC or
1540 PHYSICAL_ADDRESS_SYNC is needed (based on GET_RANGESET's sync
1541 mode property), the file system invokes the appropriate SYNC action

1542 **Postconditions:**

1543 The file system is usable by applications.

1544 **9.4.2 Driver emulates a block device using PM media**

1545 **Purpose/triggers:**

1546 The steps supporting an application write action from a driver that emulates a block device
1547 using PM as media.

1548 **Scope/context:**

1549 NVM.PM.VOLUME mode

1550 **Preconditions:**

1551 PM layer FUNDAMENTAL_SIZE reported to driver is cache line size.

1552 **Inputs:**

1553 The application provides:

- 1554 • the starting address of the memory (could be volatile) memory containing the data to
1555 write
- 1556 • the length of the memory range to be written,
- 1557 • an OS-specific reference to a block device (the virtual device backed by the PM
1558 volume),
- 1559 • the starting LBA within that block device

1560 **Success scenario:**

- 1561 1. The driver registers with the OS-specific component to be notified of errors on the PM
1562 volume. PM error handling is outside the scope of this specification, but may be similar to
1563 what is described in (and above) Figure 15 Linux Machine Check error flow with proposed
1564 new interface.
- 1565 2. Using information from a GET_RANGESET response, the driver splits the write operating
1566 into separate pieces if the target PM addresses (corresponding to application target LBAs)
1567 are in different ranges with different “sync type” values. For each of these pieces:
 - 1568 a. Using information from a GET_RANGESET response, the driver determines the PM
1569 memory address corresponding to the input starting LBA, and performs a memory
1570 copy operation from the callers input memory to the PM
 - 1571 b. The driver then performs a platform-specific flush operation
 - 1572 c. Using information from a GET_RANGESET response, the driver invokes the
1573 PHYSICAL_ADDRESS_SYNC or VIRTUAL_ADDRESS_SYNC action as needed

1574 3. No PM errors are reported by the PM error component, the driver reports that the write
1575 action succeeded.

1576 **Alternative Scenario 1:**

1577 In step 3 in the Success Scenario, the PM error component reports a PM error. The driver
1578 verifies that this error impacts the PM range being written and returns an error to the caller.

1579 **Postconditions:**

1580 The target PM range (i.e., the block device LBA range) is updated.

1581 **See also:**

1582 4.2.4 NVM block volume using PM hardware

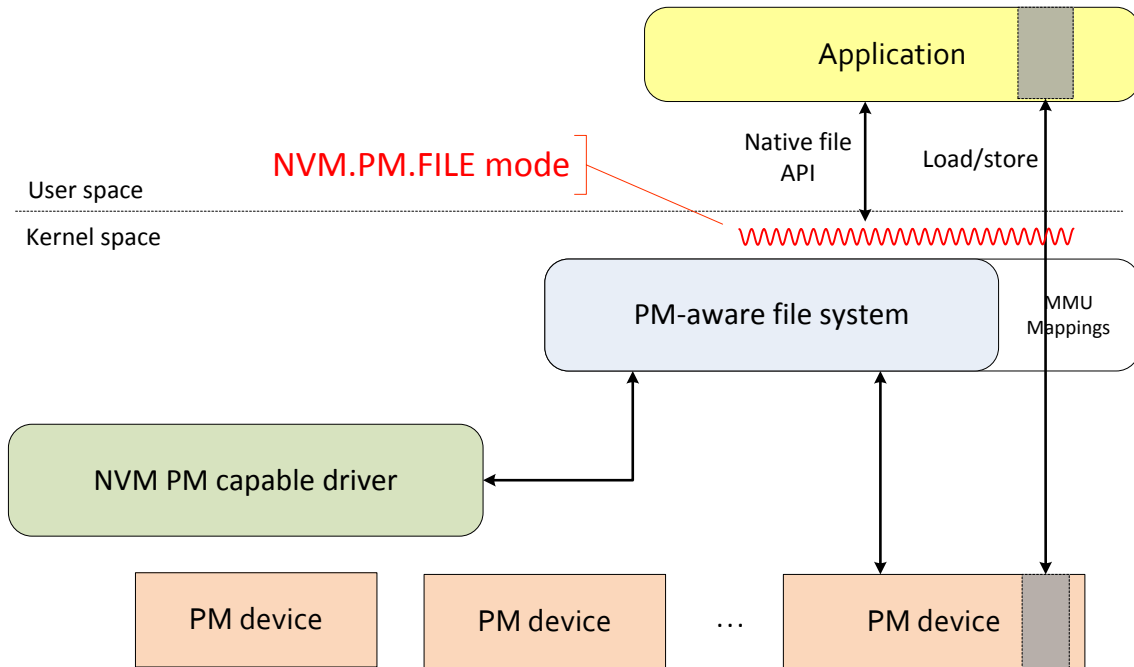
1583 **10 NVM.PM.FILE**

1584 **10.1 Overview**

1585 The NVM.PM.FILE mode access provides a means for user space applications to directly
1586 access NVM as memory. Most of the standard actions in this mode are intended to be
1587 implemented as APIs exported by existing file systems. An NVM.PM.FILE implementation
1588 behaves similarly to preexisting file system implementations, with minor exceptions. This
1589 section defines extensions to the file system implementation to accommodate persistent
1590 memory mapping and to assure interoperability with NVM.PM.FILE mode applications.

1591 Figure 14 NVM.PM.FILE mode example shows the context surrounding the point in a system
1592 (the red, wavy line) where the NVM.PM.FILE mode programming model is exposed by a PM-
1593 aware file system. A user space application consumes the programming model as is typical for
1594 current file systems. This example is not intended to preclude the possibility of a user space
1595 PM-aware file system implementation. It does, however presume that direct load/store access
1596 from user space occurs within a memory-mapped file context. The PM-aware file system
1597 interacts with an NVM PM capable driver to achieve any non-direct-access actions needed to
1598 discover or configure NVM. The PM-aware file system may access NVM devices for purposes
1599 such as file allocation, free space or other metadata management. The PM-aware file system
1600 manages additional metadata that describes the mapping of NVM device memory locations
1601 directly into user space.

1602 **Figure 14 NVM.PM.FILE mode example**



1603

1604 Once memory mapping occurs, the behavior of the NVM.PM.FILE mode diverges from
1605 NVM.FILE mode because accesses to mapped memory are in the persistence domain as soon
1606 as they reach memory. This is represented in Figure 14 NVM.PM.FILE mode example by the
1607 arrow that passes through the “MMU Mappings” extension of the file system. As a result of

1608 persistent memory mapping, primitive ACID properties arise from CPU and memory
1609 infrastructure behavior as opposed to disk drive or traditional SSD behavior. Note that writes
1610 may still be retained within processor resident caches or memory controller buffers before they
1611 reach a persistence domain. As with NVM.FILE.SYNC, the possibility remains that memory
1612 mapped writes may become persistent before a corresponding NVM.PM.FILE.SYNC action.

1613 The following actions have behaviors specific to the NVM.PM.FILE mode:

1614 NVM.PM.FILE.MAP – Add a subset of a PM file to application's address space for
1615 load/store access.

1616 NVM.PM.FILE.SYNC – Synchronize persistent memory content to assure durability and
1617 enable recovery by forcing data to reach the persistence domain.

1618 10.1.1 Applications and PM Consistency

1619 Applications (either directly or using services of a library) rely on CPU and kernel tools to
1620 achieve consistency of data in PM. These tools cause PM to exhibit certain data consistency
1621 properties enabling applications to operate correctly:

- 1622 • PM is usable as volatile (not just persistent) memory
- 1623 • Data residing in PM is consistent and durable even after a failure

1624 Consistency is defined relative to the application's objectives and design. For example, an
1625 application can utilize a write-ahead log (see SQLite.org, Write-Ahead Logging,
1626 <http://www.sqlite.org/wal.html>); when the application starts, recovery logic uses the write-ahead
1627 log to determine whether store operations completed and modifies data to achieve
1628 consistency. Similarly, durability objectives vary with applications. For database software,
1629 durability typically means that once a transaction has been committed it will remain so, even in
1630 the event of unexpected restarts. Other applications use a checkpoint mechanism other than
1631 transactions to define durable data states.

1632 When persistence behavior is ignored, memory-mapped PM is expected to operate like volatile
1633 memory. Compiled code without durability expectations is expected to continue to run
1634 correctly.

1635 This includes the following:

- 1636 • Accessible through load, store, and atomic read/modify/write instructions
- 1637 • Subject to existing processor cache coherency and “uncacheable” models
1638 (uncacheable models do not require a cache flush instruction to assure data is
1639 written to memory)
- 1640 • Load, store, and atomic read/modify/write instructions retain their current semantics
 - 1641 ○ Even when accessed from multiple threads
 - 1642 ○ Even if locks or lock-protected data live in PM
- 1643 • Able to use existing code (e.g., sort function) on PM data
- 1644 • Applies for all data producers: CPU and, where relevant, I/O
- 1645 • “Execute In Place” capability

- 1646
- Supports pointers to PM data structures

1647 At the implementation level, the behavior for fence instructions in libraries and thread visibility
1648 behavior is the same for data in PM as for data in volatile memory.

1649 Two properties assure data is consistent and durable even after failures:

- 1650
- Atomicity: some stores can't be partly visible even after a failure
 - Strict write ordering
- 1651

1652 EXAMPLE - This is a pseudo C language example of atomicity and strict ordering. In this
1653 example, `msync` implements `NVM.PM.FILE.SYNC`:

```
1654 // a, a_end in PM
1655 a[0] = foo();
1656 msync(&a[0], ...);
1657 a_end = 0;
1658 msync(&a_end, ...);
1659 . . .
1660 n = a_end + 1;
1661 a[n] = foo();
1662 msync(&a[n], ...);
1663 a_end = n;
1664 msync(&a_end, ...);
```

1665 For correctness of this example, the following assertions apply:

- 1666
- `a[0 .. a_end]` always contains valid data, even after a failure in this code.
 - `a_end` is written atomically to PM, so that the second store to `a_end` occurs no earlier than the store to `a[n]`.
- 1667
- 1668

1669 To achieve failure atomicity, aligned stores of fundamental data types (see 6.10) reach PM
1670 atomically. After a failure (allowed by the failure model), each such store is fully reflected in the
1671 resulting PM state or not at all.

1672 At least two facilities are useful to achieve strict ordering:

- 1673
- `msync`: Wait for all writes in a range to complete
 - optimization using an intra-cache-line ordering guarantee.
- 1674

1675 To elaborate on these, `msync(address_range)` ensures that if any effects from code
1676 following the call are visible, then so are all stores to `address_range` (from any thread) which
1677 precede the call to `msync` .

1678 With intra-cache-line ordering, thread-ordered stores to a single cache line become visible in
1679 PM in the order in which they are issued. The term “thread-ordered” refers to certain stores
1680 that are already known in today’s implementations to reach coherent cache in order, such as:

- 1681
- x86 MOV
 - some C11, C++11 atomic stores
- 1682

1683 • Java & C# volatile stores.

1684

1685 The CPU core and compiler do not reorder these. Within a single cache line, this order is
1686 normally preserved when the lines are evicted to PM. This last point is a critical consideration
1687 as the preservation of thread-ordered stores during eviction to PM is sometimes not
1688 guaranteed.

1689 **10.2 Actions**

1690 The following actions are mandatory for compliance with the NVM Programming Model
1691 NVM.PM.FILE mode.

1692 **10.2.1 Actions that apply across multiple modes**

1693 The following actions apply to NVM.PM.FILE mode as well as other modes.

1694 NVM.COMMON.GET_ATTRIBUTE (see 6.11.1)

1695 NVM.COMMON.SET_ATTRIBUTE (see 6.11.2)

1696 **10.2.2 Native file system actions**

1697 Native actions shall apply with unmodified syntax and semantics provided that they are
1698 compatible with programming model specific actions. This is intended to support traditional file
1699 operations allowing many applications to use PM without modification. This specifically
1700 includes mandatory implementation of the native synchronization of mapped files. As always,
1701 specific implementations may choose whether or not to implement optional native operations.

1702 **10.2.3 NVM.PM.FILE.MAP**

1703 Requirement: mandatory

1704 The mandatory form of this action shall have the same syntax found in a pre-existing file
1705 system, preferably the operating system's native file map call. The specified subset of a PM file
1706 is added to application's address space for load/store access. The semantics of this action are
1707 unlike the native MAP action because NVM.PM.FILE.MAP causes direct load/store access.
1708 For example, the role of the page cache might be reduced or eliminated. This reduces or
1709 eliminates the consumption of volatile memory as a staging area for non-volatile data. In
1710 addition, by avoiding demand paging, direct access can enable greater uniformity of access
1711 time across volatile and non-volatile data.

1712 PM mapped file operation may not provide the access time and modify time behavior typical of
1713 native file systems.

1714 PM mapped file operation may not provide the normal semantics for the native file
1715 synchronization actions (e.g., POSIX fsync and fdatasync and Win32 FlushFileBuffers). If a file
1716 is mapped at the time when the native file synchronization action is invoked, the normal
1717 semantics apply. However if the file had been mapped, data had been written to the file
1718 through the map, the data had not been synchronized by use of the NVM.PM.FILE.SYNC
1719 action, the NVM.PM.FILE.OPTIMIZED_FLUSH action, the

1720 NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY action, or the native mapped file sync
1721 action, and the mapping had been removed prior to the execution of the native file
1722 synchronization action, the action is not required to synchronize the data written to the map.

1723 Requires NVM.PM.FILE.OPEN

1724 Inputs: align with native operating system's map

1725 Outputs: align with native operating system's map

1726 Relevant Options:

1727 All of the native file system options should apply.

1728 NVM.PM.FILE.MAP_SHARED (Mandatory) – This existing native option shall be
1729 supported by the NVM.PM.FILE.MAP action. This option indicates that user space
1730 processes other than the writer can see any changes to mapped memory immediately.

1731 NVM.PM.FILE.MAP_COPY_ON_WRITE (Optional)– This existing native option
1732 indicates that any write after mapping will cause a copy on write to volatile memory, or
1733 PM that is discarded during any type of restart. The copy is only visible to the writer.
1734 The copy is not folded back into PM during the sync command.

1735 Relevant Attributes:

1736 NVM.PM.FILE.MAP_COPY_ON_WRITE_CAPABLE (see 10.3.2) - Native operating
1737 system map commands make a distinction between MAP_SHARED and
1738 MAP_COPY_ON_WRITE. Both are supported with native semantics under the NVM
1739 Programming Model. This attribute indicates whether the MAP_COPY_ON_WRITE
1740 mapping mode is supported. All NVM.PM.FILE.MAP implementations shall support the
1741 MAP_SHARED option.

1742 Error handling for mapped ranges of persistent memory is unlike I/O, in that there is no
1743 acknowledgement to a load or store instruction. Instead processors equipped to detect
1744 memory access failures respond with machine checks. These can be routed to user threads as
1745 asynchronous events. With memory-mapped PM, asynchronous events are the primary means
1746 of discovering the failure of a load to return good data. Please refer to
1747 NVM.PM.FILE.GET_ERROR_INFO (section 10.2.6) for more information on error handling
1748 behavior.

1749 Depending on memory configuration, CPU memory write pipelines may effectively preclude
1750 application level error handling during memory accesses that result from store instructions. For
1751 example, errors detected during the process of flushing the CPU's write pipeline are more
1752 likely to be associated with that pipeline than the NVM itself. Errors that arise within the CPU's
1753 write pipeline generally do not enable application level recovery at the point of the error. As a
1754 result application processes may be forced to restart when these errors occur (see PM Error
1755 Handling Annex B). Such errors should appear in CPU event logs, leading to an administrative
1756 response that is outside the scope of this specification.

1757 Applications needing timely assurance that recently stored data is recoverable should use the
1758 NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY action to verify data from NVM after it is
1759 flushed (see 10.2.7). Errors during verify are handled in the manner described in this annex.

1760 10.2.4 **NVM.PM.FILE.SYNC**

1761 Requirement: mandatory

1762 The purpose of this action is to synchronize persistent memory content to assure durability and
1763 enable recovery by forcing data to reach the persistence domain.

1764 The native file system sync action may be supported by implementations that also support
1765 NVM.PM.FILE.SYNC. The intent is that the semantics of NVM.PM.FILE.SYNC match native
1766 sync operation on memory-mapped files however because persistent memory is involved,
1767 NVM.PM.FILE implementations need not flush full pages. Note that writes may still be subject
1768 to functionality that may mask whether stored data has reached the persistence domain (such
1769 as caching or buffering within processors or memory controllers). NVM.PM.FILE.SYNC is
1770 responsible for insuring that data within the processor or memory buffers reaches the
1771 persistence domain.

1772 A number of boundary conditions can arise regarding interoperability of PM and non-PM
1773 implementation components. The following limitations apply:

- 1774 • The behavior of an NVM.PM.FILE.SYNC action applied to a range in a file that was not
1775 mapped using NVM.PM.FILE.MAP is unspecified.
- 1776 • The behavior of NVM.PM.FILE.SYNC on non-persistent memory is unspecified.

1777 In both the PM and non-PM modes, updates to ranges mapped as shared can and may
1778 become persistent in any order before a sync requires them all to become persistent. The sync
1779 action applied to a shared mapping does not guarantee write atomicity. The byte range
1780 referenced by the sync parameter may have reached a persistence domain prior to the sync
1781 command. The sync action guarantees only that the range referenced by the sync action will
1782 reach the persistence domain before the successful completion of the sync action. Any
1783 atomicity that is achieved is not caused by the sync action itself.

1784 Requires: NVM.PM.FILE.MAP

1785 Inputs: Align with native operating system's sync with the exception that alignment restrictions
1786 are relaxed.

1787 Outputs: Align with native operating system's sync with the addition that it shall return an error
1788 code.

1789 Users of the NVM.PM.FILE.SYNC action should be aware that for files that are mapped as
1790 shared, there is no requirement to buffer data on the way to the persistence domain. Although
1791 data may traverse a processor's write pipeline and other buffers within memory controllers
1792 these are more transient than the disk I/O buffering that is common in NVM.FILE
1793 implementations.

1794 **10.2.5 Error handling related to this action is expected to be derived from ongoing work**
1795 **that begins with Annex B (Informative) PM error**
1796 **handling.NVM.PM.FILE.OPTIMIZED_FLUSH**

1797 Requirement: mandatory if NVM.PM.OPTIMIZED_FLUSH_CAPABLE is set.

1798 The purpose of this action is to synchronize multiple ranges of persistent memory content to
1799 assure durability and enable recovery by forcing data to reach the persistence domain. This
1800 action has the same effect as NVM.PM.FILE.SYNC however it is intended to allow additional
1801 implementation optimization by excluding options supported by sync and by allowing multiple
1802 byte ranges to be synchronized during a single action. Page oriented alignment constraints
1803 imposed by the native definition are lifted. Because of this, implementations might be able to
1804 use underlying persistent memory more optimally than they could with the native sync. In
1805 addition some implementations may enable this action to avoid context switches into kernel
1806 space. With the exception of these differences all of the content of the NVM.PM.FILE.SYNC
1807 action description also applies to NVM.PM.FILE.OPTIMIZED_FLUSH.

1808 Requires: NVM.PM.FILE.MAP

1809 Inputs: Identical to NVM.PM.FILE.SYNC except that an array of byte ranges is specified and
1810 options are precluded. A reference to the array and the size of the array are input instead of a
1811 single address and length. Each element of the array contains an address and length of a
1812 range of bytes to be synchronized.

1813 Outputs: Align with native OS's sync with the addition that it shall return an error code.

1814 Relevant attributes: NVM.PM.FILE.OPTIMIZED_FLUSH_CAPABLE – Indicates whether this
1815 action is supported by the NVM.PM.FILE implementation (see 10.3.5).

1816 NVM.PM.FILE.OPTIMIZED_FLUSH provides no guarantee of atomicity within or across the
1817 synchronized byte ranges. Neither does it provide any guarantee of the order in which the
1818 bytes within the ranges of the action reach a persistence domain.

1819 In the event of failure the progress of the action is indeterminate. Various byte ranges may or
1820 may not have reached a persistence domain. There is no indication as to which byte ranges
1821 may have been synchronized.

1822 **10.2.6 NVM.PM.FILE.GET_ERROR_EVENT_INFO**

1823 Requirement: mandatory if NVM.PM.ERROR_EVENT_CAPABLE is set.

1824 The purpose of this action is to provide a sufficient description of an error event to enable
1825 recovery decisions to be made by an application. This action is intended to originate during an
1826 application event handler in response to a persistent memory error. In some implementations
1827 this action may map to the delivery of event description information to the application at the
1828 start of the event handler rather than a call made by the event handler. The error information
1829 returned is specific to the memory error that caused the event.

1830 Inputs: It is assumed that implementations can extract the information output by this action
1831 from the event being handled.

1832 Outputs:

1833 1 – An indication of whether or not execution of the application can be resumed from the point
1834 of interruption. If execution cannot be resumed then the process running the application should
1835 be restarted for full recovery.

1836 2 – An indication of error type enabling the application to determine whether an address is
1837 provided and the direction of data flow (load/verify vs. store) when the error was detected.

1838 3 – The memory mapped address and length of the byte range where data loss was detected
1839 by the event.

1840 Relevant attributes:

1841 NVM.PM.FILE.ERROR_EVENT_CAPABLE – Indicates whether load error event handling and
1842 this action are supported by the NVM.PM.FILE implementation (see 10.3.6).

1843 This action is used to obtain information about an error that caused a machine check involving
1844 memory mapped persistent memory. This is necessary because with persistent memory there
1845 is no opportunity to provide error information as part of a function call or I/O. The intent is to
1846 allow sophisticated error handling and recovery to occur before the application sees the event
1847 by allowing the NVM.PM.FILE implementation to handle it first. It is expected that after
1848 NVM.PM.FILE has completed whatever recovery is possible, the application error handler will
1849 be called and use the error information described here to stage subsequent recovery actions,
1850 some of which may occur after the application’s process is restarted.

1851 In some implementations the same event handler may be used for many or all memory errors.
1852 Therefore this action may arise from memory accesses unrelated to NVM. It is the application
1853 event handler’s responsibility to determine whether the memory range indicated is relevant for
1854 recovery. If the memory range is irrelevant then the event should be ignored other than as a
1855 potential trigger for a restart.

1856 In some systems, errors related to memory stores may not provide recovery information to the
1857 application unless and until load instructions attempt to access the memory locations involved.
1858 This can be accomplished using the NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY action
1859 (section 10.2.7).

1860 For more information on the circumstances which may surround this action please refer to PM
1861 Error Handling Annex B.

1862 **10.2.7 NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY**

1863 Requirement: mandatory if NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY_CAPABLE is
1864 set.

1865 The purpose of this action is to synchronize multiple ranges of persistent memory content to
1866 assure durability and enable recovery by forcing data to reach the persistence domain.
1867 Furthermore, this action verifies that data was written correctly by verifying it. The intent is to
1868 supply a mechanism whereby the application can receive data integrity assurance on writes to
1869 memory-mapped PM prior to completion of this action. This is the PM equivalent to the POSIX
1870 definition of synchronized I/O which clarifies that the intent of synchronized I/O data integrity
1871 completion is "so that an application can ensure that the data being manipulated is physically
1872 present on secondary mass storage devices".

1873 Except for the additional verification of flushed data, this action has the same effect as
1874 NVM.PM.FILE.OPTIMIZED_FLUSH.

1875 Requires: NVM.PM.FILE.MAP

1876 Inputs: Identical to NVM.PM.FILE.OPTIMIZED_FLUSH.

1877 Outputs: Align with native OS's sync with the addition that it shall return an error code. The
1878 error code indicates whether or not all data in the indicated range set is readable.

1879 Relevant attributes:

1880 NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY_CAPABLE – Indicates whether this action
1881 is supported by the NVM.PM.FILE implementation (see 10.3.7).

1882 OPTIMIZED_FLUSH_AND_VERIFY shall assure that data has been verified to be readable.
1883 Any errors discovered during verification should be logged for administrative attention.
1884 Verification shall occur across all data ranges specified in the action regardless of when they
1885 were actually flushed. Verification shall complete prior to completion of the action.

1886 In the event of failure the progress of the action is indeterminate.

1887

1888 **10.3 Attributes**

1889 **10.3.1 Attributes that apply across multiple modes**

1890 The following attributes apply to NVM.PM.FILE mode as well as other modes.

1891 NVM.COMMON.SUPPORTED_MODES (see 6.12.1)

1892 NVM.COMMON.FILE_MODE (see 6.12.2)

1893 **10.3.2 NVM.PM.FILE.MAP_COPY_ON_WRITE_CAPABLE**

1894 Requirement: mandatory

1895 This attribute indicates that MAP_COPY_ON_WRITE option is supported by the
1896 NVM.PM.FILE.MAP action.

1897 10.3.3 **NVM.PM.FILE.INTERRUPTED_STORE_ATOMICITY**

1898 Requirement: mandatory

1899 INTERRUPTED_STORE_ATOMICITY indicates whether the volume supports power fail
1900 atomicity of aligned store operations on fundamental data types. To achieve failure atomicity,
1901 aligned operations on fundamental data types reach NVM atomically. Formally “aligned
1902 operations on fundamental data types” is implementation defined. See 6.10.

1903 A value of true indicates that after an aligned store of a fundamental data type is interrupted by
1904 reset, power loss or system crash; upon restart the contents of persistent memory reflect either
1905 the state before the store or the state after the completed store. A value of false indicates that
1906 after a store interrupted by reset, power loss or system crash, upon restart the contents of
1907 memory may be such that subsequent loads may create exceptions. A value of false also
1908 indicates that after a store interrupted by reset, power loss or system crash; upon restart the
1909 contents of persistent memory may not reflect either the state before the store or the state after
1910 the completed store.

1911 The value of this attribute is true only if it’s true for all ranges in the file system.

1912 10.3.4 **NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE**

1913 Requirement: mandatory

1914 FUNDAMENTAL_ERROR_RANGE is the number of bytes that may become unavailable due
1915 to an error on an NVM device.

1916 An application may organize data in terms of FUNDAMENTAL_ERROR_RANGE to assure
1917 two key data items are not likely to be affected by a single error.

1918 Unlike NVM.PM.VOLUME (see 9), NVM.PM.FILE does not associate an offset with the
1919 FUNDAMENTAL_ERROR_RANGE because the file system is expected to handle any volume
1920 mode offset transparently to the application. The value of this attribute is the maximum of the
1921 values for all ranges in the file system.

1922 10.3.5 **NVM.PM.FILE.OPTIMIZED_FLUSH_CAPABLE**

1923 Requirement: mandatory

1924 This attribute indicates that the OPTIMIZED_FLUSH action is supported by the NVM.PM.FILE
1925 implementation.

1926 10.3.6 **NVM.PM.FILE.ERROR_EVENT_CAPABLE**

1927 Requirement: mandatory

1928 This attribute indicates that the NVM.PM.FILE implementation is capable of handling error
1929 events in such a way that, in the event of data loss, those events are subsequently delivered to
1930 applications. If error event handling is supported then NVM.PM.FILE.GET_ERROR_INFO
1931 action shall also be supported.

1932 **10.3.7 NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY_CAPABLE**

1933 Requirement: mandatory

1934 This attribute indicates that the OPTIMIZED_FLUSH_AND_VERIFY action is supported by the
1935 NVM.PM.FILE implementation.

1936 **10.4 Use cases**

1937 **10.4.1 Update PM File Record**

1938 Update a record in a PM file.

1939 **Purpose/triggers:**

1940 An application using persistent memory updates an existing record. For simplicity, this
1941 application uses fixed size records. The record size is defined by application data
1942 considerations.

1943 **Scope/context:**

1944 Persistent memory context; this use case shows basic behavior.

1945 **Preconditions:**

- 1946 • The administrator created a PM file and provided its name to the application; this name is
- 1947 accessible to the application – perhaps in a configuration file
- 1948 • The application has populated the PM file contents
- 1949 • The PM file is not in use at the start of this use case (no sharing considerations)

1950 **Inputs:**

1951 The content of the record, the location (relative to the file) where the record resides

1952 **Success scenario:**

- 1953 1) The application uses the native OPEN action, passing in the file name
- 1954 2) The application uses the NVM.PM.FILE.MAP action, passing in the file descriptor returned
1955 by the native OPEN. Since the records are not necessarily page aligned, the application
1956 maps the entire file.
- 1957 3) The application registers for memory hardware exceptions
- 1958 4) The application stores the new record content to the address returned by
1959 NVM.PM.FILE.MAP offset by the record’s location
- 1960 5) The application uses NVM.PM.FILE.SYNC to flush the updated record to the persistence
1961 domain
 - 1962 a. The application may simply sync the entire file

- 1963 b. Alternatively, the application may limit the range to be sync'd
1964 6) The application uses the native UNMAP and CLOSE actions to clean up.

1965 **Failure Scenario:**

1966 While reading PM content (accessing via a load operation), a memory hardware exception is
1967 reported. The application's event handler is called with information about the error as
1968 described in NVM.PM.FILE.GET_ERROR_INFO. Based on the information provided, the
1969 application records the error for subsequent recovery and determines whether to restart or
1970 continue execution.

1971 **Postconditions:**

1972 The record is updated.

1973 10.4.2 **Direct load access**

1974 **Purpose/triggers:**

1975 An application developer wishes to retrieve data from a persistent memory-mapped file using
1976 direct memory load instruction access with error handling for uncorrectable errors.

1977 **Scope/context:**

1978 NVM.PM.FILE

1979 **Inputs:**

- 1980 • Virtual address of the data.

1981 **Outputs:**

- 1982 • Data from persistent memory if successful
1983 • Error code if an error was detected within the accessed memory range.

1984 **Preconditions:**

- 1985 • The persistent memory file must be mapped into a region of virtual memory.
1986 • The virtual address must be within the mapped region of the file.

1987 **Postconditions:**

- 1988 • If an error was returned, the data may be unreadable. Future load accesses may
1989 continue to return an error until the data is overwritten to clear the error condition
1990 • If no error was returned, there is no postcondition.

1991 **Success and Failure Scenarios:**

1992 Consider the following fragment of example source code, which is simplified from the code for
1993 the function that reads SQLite's transaction journal:

```
1994     retCode = pread(journalFD, magic, 8, off);  
1995     if (retCode != SQLITE_OK) return retCode;  
1996
```

```

1997     if (memcmp(magic, journalMagic, 8) != 0)
1998         return SQLITE_DONE;

```

1999 This example code reads an eight-byte magic number from the journal header into an eight-
 2000 byte buffer named *magic* using a standard file *read* call. If an error is returned from the *read*
 2001 system call, the function exits with an error return code indicating that an I/O error occurred. If
 2002 no error occurs, it then compares the contents of the *magic* buffer against the expected magic
 2003 number constant named *journalMagic*. If the contents of the buffer do not match the expected
 2004 magic number, the function exits with an error return code.

2005 An equivalent version of the function using direct memory load instruction access to a mapped
 2006 file is:

```

2007     volatile siginfo_t errContext;
2008     ...
2009     int retCode = SQLITE_OK;
2010
2011     TRY
2012     {
2013         if (memcmp(journalMmapAddr + off, journalMagic, 8) != 0)
2014             retCode = SQLITE_DONE;
2015     }
2016     CATCH(BUS_MCEERR_AR)
2017     {
2018         if ((errContext.si_code == BUS_MCEERR_AR) &&
2019             (errContext.si_addr >= journalMmapAddr) &&
2020             (errContext.si_addr < (journalMmapAddr + journalMmapSize))) {
2021             retCode = SQLITE_IOERR;
2022         } else {
2023             signal(errContext.si_signo, SIG_DFL);
2024             raise(errContext.si_signo);
2025         }
2026     }
2027     ENDTRY;
2028
2029     if (retCode != SQLITE_OK) return retCode;

```

2030 The mapped file example compares the magic number in the header of the journal file against
 2031 the expected magic number using the *memcmp* function by passing a pointer containing the
 2032 address of the magic number in the mapped region of the file. If the contents of the magic
 2033 number member of the file header do not match the expected magic number, the function exits
 2034 with an error return code.

2035 The application-provided TRY/CATCH/ENDTRY macros implement a form of exception
 2036 handling using POSIX *sigsetjmp* and *siglongjmp* C library functions. The TRY macro initializes
 2037 a *sigjmp_buf* by calling *sigsetjmp*. When a SIGBUS signal is raised, the signal handler calls
 2038 *siglongjmp* using the *sigjmp_buf* set by the *sigsetjmp* call in the TRY macro. Execution then
 2039 continues in the CATCH clause. (Caution: the code in the TRY block should not call library
 2040 functions as they are not likely to be exception-safe.) Code for the Windows platform would be
 2041 similar except that it would use the standard Structured Exception Handling *try-except*
 2042 statement catching the EXCEPTION_IN_PAGE_ERROR exception rather than application-
 2043 provided TRY/CATCH/ENDTRY macros.

2044 If an error occurs during the read of the magic number data from the mapped file, a SIGBUS
2045 signal will be raised resulting in the transfer of control to the CATCH clause. The address of
2046 the error is compared against the range of the memory-mapped file. In this example the error
2047 address is assumed to be in the process's logical address space. If the error address is within
2048 the range of the memory-mapped file, the function returns an error code indication that an I/O
2049 error occurred. If the error address is outside the range of the memory-mapped file, the error is
2050 assumed to be for some other memory region such as the program text, stack, or heap, and
2051 the signal or exception is re-raised. This is likely to result in a fatal error for the program.

2052 **See also:**

- 2053
 - Microsoft Corporation, Reading and Writing From a File View (Windows), available from
2054 <http://msdn.microsoft.com/en-us/library/windows/desktop/aa366801.aspx>

2055 **10.4.3 Direct store access**

2056 **Purpose/triggers:**

2057 An application developer wishes to place data in a persistent memory-mapped file using direct
2058 memory store instruction access.

2059 **Scope/context:**

2060 NVM.PM.FILE

2061 **Inputs:**

- 2062
 - Virtual address of the data.
 - The data to store.

2064 **Outputs:**

- 2065
 - Error code if an error occurred.

2066 **Preconditions:**

- 2067
 - The persistent memory file must be mapped into a region of virtual memory.
 - The virtual address must be within the mapped region of the file.

2069 **Postconditions:**

- 2070
 - If an error was returned, the state of the data recorded in the persistence domain is
2071 indeterminate.
 - If no error was returned, the specified data is either recorded in the persistence domain
2072 or an undiagnosed error may have occurred.

2074 **Success and Failure Scenarios:**

2075 Consider the following fragment of example source code, which is simplified from the code for
2076 the function that writes to SQLite's transaction journal:

```
2077     ret = pwrite(journalFD, dbPgData, dbPgSize, off);  
2078     if (ret != SQLITE_OK) return ret;  
2079     ret = write32bits(journalFD, off + dbPgSize, cksum);
```

```
2080     if (ret != SQLITE_OK) return ret;
2081     ret = fdatsync(journalFD);
2082     if (ret != SQLITE_OK) return ret;
```

2083 This example code writes a page of data from the database cache to the journal using a
2084 standard file *write* call. If an error is returned from the *write* system call, the function exits with
2085 an error return code indicating that an I/O error occurred. If no error occurs, the function then
2086 appends the checksum of the data, again using a standard file *write* call. If an error is returned
2087 from the *write* system call, the function exits with an error return code indicating that an I/O
2088 error occurred. If no error occurs, the function then invokes the *fdatsync* system call to flush
2089 the written data from the file system buffer cache to the persistence domain. If an error is
2090 returned from the *fdatsync* system call, the function exits with an error return code indicating
2091 that an I/O error occurred. If no error occurs, the written data has been recorded in the
2092 persistence domain.

2093 An equivalent version of the function using direct memory store instruction access to a
2094 memory-mapped file is:

```
2095     memcpy(journalMmapAddr + off, dbPgData, dbPgSize);
2096     PM_track_dirty_mem(dirtyLines, journalMmapAddr + off, dbPgSize);
2097
2098     store32bits(journalMmapAddr + off + dbPgSize, cksum);
2099     PM_track_dirty_mem(dirtyLines, journalMmapAddr + off + dbPgSize, 4);
2100
2101     ret = PM_optimized_flush(dirtyLines, dirtyLinesCount);
2102
2103     if (ret == SQLITE_OK) dirtyLinesCount = 0;
2104
2105     return ret;
```

2106 The memory-mapped file example writes a page of data from the database cache to the
2107 journal using the *memcpy* function by passing a pointer containing the address of the page
2108 data field in the mapped region of the file. It then appends the checksum using direct stores to
2109 the address of the checksum field in the mapped region of the file.

2110 The code calls the application-provided *PM_track_dirty_mem* function to record the virtual
2111 address and size of the memory regions that it has modified. The *PM_track_dirty_mem*
2112 function constructs a list of these modified regions in the *dirtyLines* array.

2113 The function then calls the *PM_optimized_flush* function to flush the written data to the
2114 persistence domain. If an error is returned from the *PM_optimized_flush* call, the function exits
2115 with an error return code indicating that an I/O error occurred. If no error occurs, the written
2116 data is either recorded in the persistence domain or an undiagnosed error may have occurred.
2117 Note that this postcondition is weaker than the guarantee offered by the *fdatsync* system call
2118 in the original example.

2119 See also:

- 2120 • Microsoft Corporation, Reading and Writing From a File View (Windows), available from
2121 <http://msdn.microsoft.com/en-us/library/windows/desktop/aa366801.aspx>

2122 10.4.4 Direct store access with synchronized I/O data integrity completion

2123 **Purpose/triggers:**

2124 An application developer wishes to place data in a persistent memory-mapped file using direct
2125 memory store instruction access with synchronized I/O data integrity completion.

2126 **Scope/context:**

2127 NVM.PM.FILE

2128 **Inputs:**

- 2129 • Virtual address of the data.
- 2130 • The data to store.

2131 **Outputs:**

- 2132 • Error code if an error occurred.

2133 **Preconditions:**

- 2134 • The persistent memory file must be mapped into a region of virtual memory.
- 2135 • The virtual address must be within the mapped region of the file.

2136 **Postconditions:**

- 2137 • If an error was returned, the state of the data recorded in the persistence domain is
2138 indeterminate.
- 2139 • If no error was returned, the specified data is recorded in the persistence domain.

2140 **Success and Failure Scenarios:**

2141 Consider the following fragment of example source code, which is simplified from the code for
2142 the function that writes to SQLite's transaction journal:

```
2143     ret = pwrite(journalFD, dbPgData, dbPgSize, off);  
2144     if (ret != SQLITE_OK) return ret;  
2145     ret = write32bits(journalFD, off + dbPgSize, cksum);  
2146     if (ret != SQLITE_OK) return ret;  
2147  
2148     ret = fdatasync(journalFD);  
2149     if (ret != SQLITE_OK) return ret;
```

2150 This example code writes a page of data from the database cache to the journal using a
2151 standard file *write* call. If an error is returned from the *write* system call, the function exits with
2152 an error return code indicating that an I/O error occurred. If no error occurs, the function then
2153 appends the checksum of the data, again using a standard file *write* call. If an error is returned
2154 from the *write* system call, the function exits with an error return code indicating that an I/O
2155 error occurred. If no error occurs, the function then invokes the *fdatasync* system call to flush
2156 the written data from the file system buffer cache to the persistence domain. If an error is
2157 returned from the *fdatasync* system call, the function exits with an error return code indicating

2158 that an I/O error occurred. If no error occurs, the written data has been recorded in the
2159 persistence domain.

2160 An equivalent version of the function using direct memory store instruction access to a
2161 memory-mapped file is:

```
2162 memcpy(journalMmapAddr + off, dbPgData, dbPgSize);  
2163 PM_track_dirty_mem(dirtyLines, journalMmapAddr + off, dbPgSize);  
2164  
2165 store32bits(journalMmapAddr + off + dbPgSize, cksum);  
2166 PM_track_dirty_mem(dirtyLines, journalMmapAddr + off + dbPgSize, 4);  
2167  
2168 ret = PM_optimized_flush_and_verify(dirtyLines, dirtyLinesCount);  
2169  
2170 if (ret == SQLITE_OK) dirtyLinesCount = 0;  
2171  
2172 return ret;
```

2173 The memory-mapped file example writes a page of data from the database cache to the
2174 journal using the *memcpy* function by passing a pointer containing the address of the page
2175 data field in the mapped region of the file. It then appends the checksum using direct stores to
2176 the address of the checksum field in the mapped region of the file.

2177 The code calls the application-provided *PM_track_dirty_mem* function to record the virtual
2178 address and size of the memory regions that it has modified. The *PM_track_dirty_mem*
2179 function constructs a list of these modified regions in the *dirtyLines* array.

2180 The function then calls the *PM_optimized_flush_and_verify* function to flush the written data to
2181 the persistence domain. If an error is returned from the *PM_optimized_flush_and_verify* call,
2182 the function exits with an error return code indicating that an I/O error occurred. If no error
2183 occurs, the written data has been recorded in the persistence domain. Note that this
2184 postcondition is equivalent to the guarantee offered by the *fdatsync* system call in the original
2185 example.

2186 See also:

- 2187 • Microsoft Corp, FlushFileBuffers function (Windows),
2188 <http://msdn.microsoft.com/en-us/library/windows/desktop/aa364439.aspx>
- 2189 • Oracle Corp, Synchronized I/O section in the Programming Interfaces Guide, available
2190 from
2191 <http://docs.oracle.com/cd/E19683-01/816-5042/chap7rt-57/index.html>
- 2192 • The Open Group, “The Open Group Base Specification Issue 6”, section 3.373
2193 “Synchronized Input and Output”, available from
2194 http://pubs.opengroup.org/onlinepubs/007904975/basedefs/xbd_chap03.html#tag_03_3
2195 [73](#)

2196 10.4.5 Persistent Memory Transaction Logging

2197 **Purpose/Triggers:**

2198 An application developer wishes to implement a transaction log that maintains data integrity
2199 through system crashes, system resets, and power failures. The underlying storage is byte-
2200 granular persistent memory.

2201 **Scope/Context:**

2202 NVM.PM.VOLUME and NVM.PM.FILE

2203 For notational convenience, this use case will use the term “file” to apply to either a file in the
2204 conventional sense which is accessed through the NVM.PM.FILE interface, or a specific
2205 subset of memory ranges residing on an NVM device which are accessed through the
2206 NVM.BLOCK interface.

2207 **Inputs:**

- 2208 • A set of changes to the persistent state to be applied as a single transaction.
- 2209 • The data and log files.

2210 **Outputs:**

- 2211 • An indication of transaction commit or abort.

2212 **Postconditions:**

- 2213 • If an abort indication was returned, the data was not committed and the previous
2214 contents have not been modified.
- 2215 • If a commit indication was returned, the data has been entirely committed.
- 2216 • After a system crash, reset, or power failure followed by system restart and execution of
2217 the application transaction recovery process, the data has either been entirely
2218 committed or the previous contents have not been modified.

2219 **Success Scenario:**

2220 The application transaction logic uses a log file in combination with its data file to atomically
2221 update the persistent state of the application. The log may implement a before-image log or a
2222 write-ahead log. The application transaction logic should configure itself to handle torn or
2223 interrupted writes to the log or data files.

2224 Since persistent memory may be byte-granular, torn writes may occur at any point during a
2225 series of stores. The application should be prepared to detect a torn write of the record and
2226 either discard or recover such a torn record during the recovery process. One common way of
2227 detecting such a torn write is for the application to compute a hash of the record and record the
2228 hash in the record. Upon reading the record, the application re-computes the hash and
2229 compares it with the recorded hash; if they do not match, the record has been torn.

2230 **10.4.5.1 NVM.PM.FILE.INTERRUPTED_STORE_ATOMICITY is true**

2231 If the NVM.PM.FILE.INTERRUPTED_STORE_ATOMICITY is true, then writes which are
2232 interrupted by a system crash, system reset, or power failure occur atomically. In other words,
2233 upon restart the contents of persistent memory reflect either the state before the store or the
2234 state after the completed store.

2235 In this case, the application need not handle interrupted writes to the log or data files.

2236 **10.4.5.2 NVM.PM.FILE.INTERRUPTED_STORE_ATOMICITY is false**

2237 NVM.PM.FILE.INTERRUPTED_STORE_ATOMICITY is false, then writes which are
2238 interrupted by a system crash, system reset, or power failure do not occur atomically. In other
2239 words, upon restart the contents of persistent memory may be such that subsequent loads
2240 may create exceptions depending on the value of the FUNDAMENTAL_ERROR_RANGE
2241 attribute.

2242 In this case, the application should be prepared to handle an interrupted write to the log or data
2243 files.

2244 *10.4.5.2.1 NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE > 0*

2245 If the NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE is greater than zero, the application
2246 should align the log or data records with the
2247 NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE and pad the record size to be an integral
2248 multiple of NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE. This prevents more than one
2249 record from residing in the same fundamental error range. The application should be prepared
2250 to discard or recover the record if a load returns an exception when subsequently reading the
2251 record during the recovery process. (See also SQLite.org, *Powersafe Overwrite*,
2252 <http://www.sqlite.org/psow.html>.)

2253 *10.4.5.2.2 NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE = 0*

2254 If the NVM.PM.FILE.FUNDAMENTAL_ERROR_RANGE is zero, the application lacks sufficient
2255 information to handle interrupted writes to the log or data files.

2256 **Failure Scenarios:**

2257 Consider the recovery of an error resulting from an interrupted write on a persistent memory
2258 volume or file system where the NVM.PM.FILE.INTERRUPTED_STORE_ATOMICITY is false.
2259 This error may be persistent and may be returned whenever the affected fundamental error
2260 range is read. To repair this error, the application should be prepared to overwrite such a
2261 range.

2262 One common way of ensuring that the application will overwrite a range is by assigning it to
2263 the set of internal free space managed by the application, which is never read and is available
2264 to be allocated and overwritten at some point in the future. For example, the range may be part
2265 of a circular log. If the range is marked as free, the transaction log logic will eventually allocate
2266 and overwrite that range as records are written to the log.

2267 Another common way is to record either a before-image or after-image of a data range in a log.
2268 During recovery after a system crash, system reset, or power failure, the application replays
2269 the records in the log and overwrites the data range with either the before-image contents or
2270 the after-image contents.

2271 **See also:**

- 2272 • SQLite.org, *Atomic Commit in SQLite*, <http://www.sqlite.org/atomiccommit.html>
- 2273 • SQLite.org, *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>
- 2274 • SQLite.org, *Write-Ahead Logging*, <http://www.sqlite.org/wal.html>

2275 **Annex A (Informative) PM pointers**

2276 Pointers are data types that hold virtual addresses of data in memory. When applications use
2277 pointers with volatile memory, the value of the pointer must be re-assigned each time the
2278 program is run (a consequence of the memory being volatile). When applications map a file (or
2279 a portion of a file) residing in persistent memory to virtual addresses, it may or may not be
2280 assigned the same virtual address. If not, then pointers to values in that mapped memory will
2281 not reference the same data. There are several possible solutions to this problem:

- 2282 1) Relative pointers
- 2283 2) Regions are mapped at fixed addresses
- 2284 3) Pointers are relocated when region is remapped

2285 All three approaches are problematic, and involve different challenges that have not been fully
2286 addressed.

2287 None, except perhaps the third one, handles C++ vtable pointers inside persistent memory, or
2288 pointers to string constants, where the string physically resides in the executable, and not the
2289 memory-mapped file. Both of those issues are common.

2290 Option (1) implies that no existing pointer-containing library data structures can be stored in
2291 PM, since pointer representations change. Option (2) requires careful management of virtual
2292 addresses to ensure that memory-mapped files that may need to be accessed simultaneously
2293 are not assigned to the same address. It may also limit address space layout randomization.
2294 Option (3) presents challenges in, for example, a C language environment in which pointers
2295 may not be unambiguously identifiable, and where they may serve as hash table indices or the
2296 like. Pointer relocation would invalidate such hash tables. It may be significantly easier in the
2297 context of a Java-like language.

Annex B (Informative) PM error handling

Persistent memory error handling for NVM.PM.FILE.MAP ranges is unique because unlike I/O, there is no acknowledgement to a load or store instruction. Instead processors equipped to detect memory access failures respond with machine checks. In some cases these can be routed to user threads as asynchronous events.

This annex only describes the handling of errors resulting from load instructions that access memory. As will be described later in this annex, no application level recovery is enabled at the point of a store error. These errors should appear in CPU event logs, leading to an administrative response that is outside the scope of this annex.

Applications needing timely assurance that recently stored data is recoverable should use the NVM.PM.FILE.OPTIMIZED_FLUSH_AND_VERIFY (see 10.2.7) action to read data back from NVM after it is flushed. Errors during verify are handled in the manner described in this annex.

There are several scenarios that can arise in the handling of machine checks related to persistent memory errors while reading data from memory locations such as can occur during “load” instructions. Concepts are introduced here in an attempt to advance the state of the art in persistent memory error handling. The goal is to provide error reporting and recovery capability to applications that is equivalent to the current practice for I/O.

We need several definitions to assist in reasoning about asynchronous events.

- Machine check: an interrupt. In this case interrupts that result from memory errors are of specific interest.
- Precise machine check – an interrupt that allows an application to resume at the interrupted instruction
- Error containment – this is an indication of how well the system can determine the extent of an error. This enables a range of memory affected by an error that caused an interrupt to be returned to the application.
- Real time error recovery – This refers to scenarios in which the application can continue execution after an error as opposed to being restarted.
- Asynchronous event handler – This refers to code provided by an application that runs in response to an asynchronous event, in this case an event that indicates a memory error. An application’s event handler uses information about the error to determine whether execution can safely continue from within the application or whether a partial or full restart of the application is required to recover from the error.

The ability to handle persistent memory errors depends on the capability of the processor and memory system. It is useful to categorize error handling capability into three levels:

- No memory error detection – the lowest end systems have little or no memory error detection or correction capability such as ECC, CRC or parity.
- Non-precise or uncontained memory error detection – these systems detect memory errors but they do not provide information about the location of the error and/or fail to offer enough information to resume execution from the interrupted instruction.

2337 - Precise, contained memory error detection – these systems detect memory errors and
2338 report their locations in real time. These systems are also able to contain many errors more
2339 effectively. This increases the range of errors that allowing applications to continue
2340 execution rather than resetting the application or the whole system. This capability is
2341 common when using higher RAS processors.

2342 Only the last category of systems can, with appropriate operating system software
2343 enhancement, meet the error reporting goal stated above. The other two categories of systems
2344 risk scenarios where persistent memory errors are forced to repeatedly reset threads or
2345 processors, rendering them incapable of doing work. Unrecovered persistent memory errors
2346 are more problematic than volatile memory errors because they are less likely to disappear
2347 during a processor reset or application restart.

2348 Systems with precise memory error detection capability can experience a range of scenarios
2349 depending on the nature of the error. These can be categorized into three types.

- 2350 - Platform can't capture error
 - 2351 • Perhaps application or operating system dies
 - 2352 • Perhaps hardware product include diagnostic utilities
- 2353 - Platform can capture error, considered fatal
 - 2354 • Operating system crashes
 - 2355 • Address info potentially stored by operating system or hardware/firmware
 - 2356 • Application could use info on restart
- 2357 - Platform can capture error & deliver to application
 - 2358 • Reported to application using asynchronous "event"
 - 2359 • Example: SIGBUS on UNIX w/address info

2360 If the platform can't capture the error then no real time recovery is possible. The system may
2361 function intermittently or not at all until diagnostics can expose the problem. The same thing
2362 happens whether the platform lacks memory error detection capability or the platform has the
2363 capability but was unable to use it due to a low probability error scenario.

2364 If the platform can capture the error but it is fatal then real time recovery is not possible,
2365 however then the system may make information about the error available after system or
2366 application restart. For this scenario, actions are proposed below to obtain error descriptions.

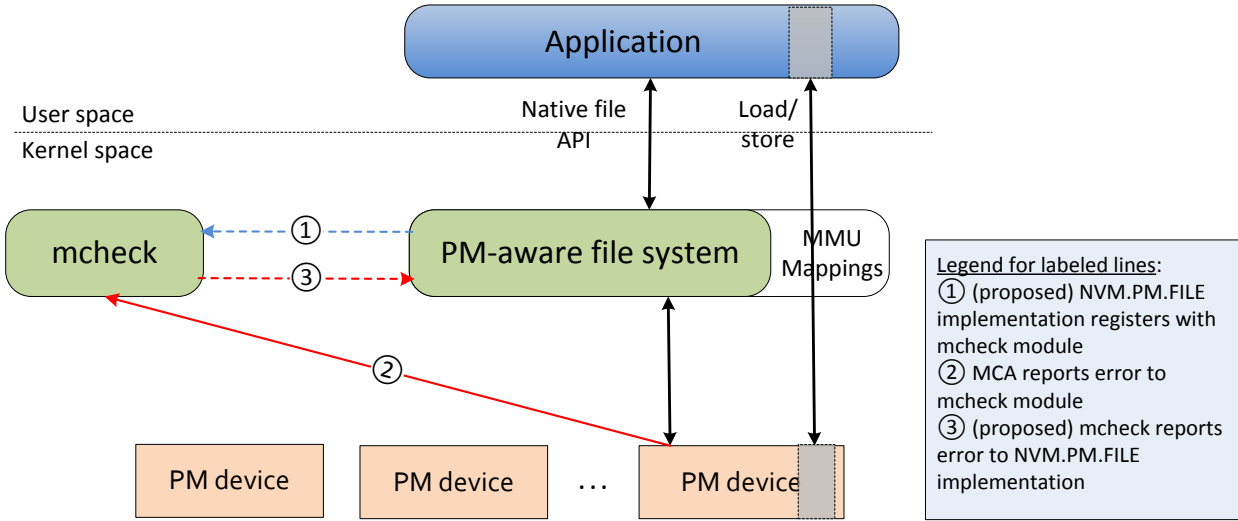
2367 If the platform can deliver the error to the application then real time recovery may be possible.
2368 An action is proposed below to represent the means that the application uses to obtain error
2369 information immediately after the failure.

2370 As stated at the beginning of this annex, only errors during load are addressed by this annex.
2371 As with other storage media, little or no error checking occurs during store instructions (aka
2372 writes). In addition, memory write pipelines within CPU's effectively preclude error handling
2373 during memory accesses that result from store instructions. For example, errors detected
2374 during the process of flushing the CPU's write pipeline are more likely to be associated with
2375 that pipeline than the NVM itself. Errors that arise within the CPU's write pipeline are generally
2376 not contained so no application level recovery is enabled at the point of the error.

2377 Continuing to analyze the real time error delivery scenario, the handling of errors on load
 2378 instructions is sufficient in today's high RAS systems to avoid the consumption of erroneous
 2379 data by the application. Several enhancements are required to meet the goal of I/O-like
 2380 application recoverability.

2381 Using Linux running on the Intel architecture as an example, memory errors are reported using
 2382 Intel's Machine Check Architecture (MCA). When the operating system enables this feature,
 2383 the error flow on an uncorrectable error is shown by the solid red arrow (labeled ②) in Figure
 2384 15 Linux Machine Check error flow with proposed new interface, which depicts the mcheck
 2385 component getting notified when the bad location in PM is accessed.

2386 **Figure 15 Linux Machine Check error flow with proposed new interface**



2387 As mentioned above, sending the application a SIGBUS (a type of asynchronous event) allows
 2388 the application to decide what to do. However, in this case, remember that the NVM.PM.FILE
 2389 manages the PM and that the location being accessed is part of a file on that file system. So
 2390 even if the application gets a signal preventing it from using corrupted data, a method for
 2391 recovering from this situation must be provided. A system administrator may try to back up rest
 2392 of the data in the file system before replacing the faulty PM, but with the error mechanism
 2393 we've described so far, the backup application would be sent a SIGBUS every time it touched
 2394 the bad location. What is needed in this case is a way for the NVM.PM.FILE implementation to
 2395 be notified of the error so it can isolate the affected PM locations and then continue to provide
 2396 access to the rest of the PM file system. The dashed arrows in Figure 15 show the necessary
 2397 modification to the machine check code in Linux. On start-up, the NVM.PM.FILE
 2398 implementation registers with the machine code to show it has responsibility for certain ranges
 2399 of PM. Later, when the error occurs, NVM.PM.FILE gets called back by the mcheck
 2400 component and has a chance to handle the error.
 2401

2402 This suggested machine check flow change enables the file system to participate in recovery
 2403 while not eliminating the ability to signal the error to the application. The application view of
 2404 errors not corrected by the file system depends on whether the error handling was precise and
 2405 contained. Imprecise error handling precludes resumption of the application, in which case the

2406 one recovery method available besides restart is a non-local go-to. This resumes execution at
2407 an application error handling routine which, depending on the design of the application, may be
2408 able to recover from the error without resuming from the point in the code that was interrupted.

2409 Taking all of this into account, the proposed application view of persistent memory errors is as
2410 described by the NVM.PM.FILE.MAP action (section 10.2.3) and the
2411 NVM.PM.FILE.GET_ERROR_INFO action (section 10.2.6).

2412 The following actions have been proposed to provide the application with the means necessary
2413 to obtain error information after a fatal error.

- 2414 • PM.FILE.ERROR_CHECK(file, offset, length): Discover if range has any outstanding
2415 errors. Returns a list of errors referenced by file and offset.
- 2416 • PM.FILE.ERROR_CLEAR(file, offset, length): Reset error state (and data) for a range: may
2417 not succeed

2418 The following attributes have been proposed to enable application to discover the error
2419 reporting capabilities of the implementation.

- 2420 • NVM.PM.FILE.ERROR_CHECK_CAPABLE - System supports asking if range is in error
2421 state

2422 **Annex C (Informative) Deferred behavior**

2423 This annex lists some behaviors that are being considered for future specifications.

2424 **D.1 Remote sharing of NVM**

2425 This version of the specification talks about the relationship between DMA and persistent
2426 memory (see 6.6 Interaction with I/O devices) which should enable a network device to access
2427 NVM devices. But no comprehensive approach to remote share of NVM is addressed in this
2428 version of the specification.

2429 **D.2 MAP_CACHED OPTION FOR NVM.PM.FILE.MAP**

2430 This would enable memory mapped ranges to be either cached or uncached by the CPU.

2431 **D.3 NVM.PM.FILE.DURABLE.STORE**

2432 This might imply that through this action things become durable and visible at the same time,
2433 or not visible until it is durable. Is there a special case for atomic write that, by the time the
2434 operation completes, it is both visible and durable? The prospective use case is an opportunity
2435 for someone with a hardware implementation that does not require separation of store and
2436 sync. This is not envisioned as the same as a file system write. It still implies a size of the
2437 store. The use case for NVM.FILE.DURABLE.STORE is to force access to the persistence
2438 domain.

2439 **D.4 Enhanced NVM.PM.FILE.WRITE**

2440 Add an NVM.PM.FILE.WRITE action where the only content describes error handling.

2441 **D.5 Management-only behavior**

2442 Several management-only behaviors have been discussed, but deferred to a future revision;
2443 including:

- 2444 • Secure Erase
- 2445 • Behavior enabling management application to discover PM devices (and behavior to fill
2446 gaps in the discovery of block NVM attributes)
- 2447 • Attribute exposing flash erase block size for management of disk partitions

2448 **D.6 Access hints**

2449 Allow applications to suggest how data is placed on storage

2450 **D.7 Multi-device atomic multi-write action**

2451 Perform an atomic write to multiple extents in different devices.

2452 **D.8 NVM.BLOCK.DISCARD_IF_YOU_MUST action**

2453 The text below was partially developed, before being deferred to a future revision.

2454 10.4.6 **NVM.BLOCK.DISCARD_IF_YOU_MUST**

2455 Proposed new name MARK_DISCARDABLE

2456 Purpose - discard blocks to prevent write amplification

2457 This action notifies the NVM device that some or all of the blocks which constitute a volume
2458 are no longer needed by the application, but the NVM device should defer changes to the
2459 blocks as long as possible. This action is a hint to the device.

2460 If the data has been retained, a subsequent read shall return “success” along with the data.
2461 Otherwise, it shall return an error indicating the data does not exist (and the data buffer area
2462 for that block is undefined).

2463 Inputs: a range of blocks (starting LBA and length in logical blocks)

2464 Status: Success indicates the request is accepted but not necessarily acted upon.

2465 Existing implementations of TRIM may work this way.

2466 10.4.7 **DISCARD_IF_YOU_MUST use case**

2467 **Purpose/triggers:**

2468 An NVM device may allocate blocks of storage from a common pool of storage. The device
2469 may also allocate storage through a thin provisioning mechanism. In each of these cases, it is
2470 useful to provide a mechanism which allows an application or NVM user to notify the NVM
2471 storage system that some or all of the blocks which constitute the volume are no longer
2472 needed by the application. This allows the NVM device to return the memory allocated for the
2473 unused blocks to the free memory pool and make the unused blocks available for other
2474 consumers to use.

2475 DISCARD_IF_YOU_MUST operation informs the NVM device that that the specified blocks
2476 are no longer required. DISCARD_IF_YOU_MUST instructs the NVM device to release
2477 previously allocated blocks to the NVM device’s free memory pool. The NVM device releases
2478 the used memory to the free storage pool based on the specific implementation of that device.
2479 If the device cannot release the specified blocks, the DISCARD_IF_YOU_MUST operation
2480 returns an error.

2481 **Scope/context:**

2482 This use case describes the capabilities of an NVM device that the NVM consumer can
2483 determine.

2484 **Inputs:**

2485 The range to be freed.

2486 **Success scenario:**

2487 The operation succeeds unless an invalid region is specified or the NVM device is unable to
2488 free the specified region.

2489 **Outputs:**
2490 The completion status.

2491 **Postconditions:**
2492 The specified region is erased and released to the free storage pool.

2493 **See also:**
2494 DISCARD_IF_YOU_CAN
2495 EXISTS

2496 **D.9 Atomic write action with Isolation**

2497 Offer alternatives to ATOMIC_WRITE and ATOMIC_MULTIWRITE that also include isolation
2498 with respect to other atomic write actions. Issues to consider include whether order is required,
2499 whether isolation applies across multiple paths, and how isolation applies to file mapped I/O.

2500 **D.10 Atomic Sync/Flush action for PM**

2501 The goal is a mechanism analogous to atomic writes for persistent memory. Since stored
2502 memory may be implicitly flushed by a file system, defining this mechanism may be more
2503 complex than simply defining an action.

2504 **D.11 Hardware-assisted verify**

2505 Future PM device implementations may provide a capability to perform the verify step of
2506 OPTIMIZED_FLUSH_AND_VERIFY without requiring an explicit load instruction. This
2507 capability may require the addition of actions and attributes in NVM.PM.VOLUME mode; this
2508 change is deferred until we have examples of this type of device.