



# NVM Programming Model

Version 1.0.0 Revision 10

*Abstract: This SNIA Working Draft 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.” Suggestions for revision should be directed to <http://www.snia.org/feedback/>

***Working Draft***

September 30, 2013

2

## Revision history

### 3 **Revision 0**

#### 4 **Date**

5 March 12, 2013

#### 6 **Changes incorporated**

7 Reworked as first draft for a single Programming Model specification (based on 4  
8 modes)

### 9 **Revision 1**

#### 10 **Date**

11 March 26, 2013

#### 12 **Changes incorporated**

- 13 - Primarily consists of changes discussed at March face-to-face meeting
- 14 - Added temporary annex with templates for fixmes, actions, ... easy to copy-n-paste
- 15 where needed
- 16 - Added bibliography
- 17 - express third state of atomicity outcome as an attribute (tentative)
- 18 - Changed text talking about saving to media to say persistence domain
- 19 - Added TRIM and friends to NVM.BLOCK, included moving Walt's use cases from
- 20 NVM.PM.FILE
- 21 - Added Atomic Write to NVM.BLOCK
- 22 - Added SCAR to NVM.BLOCK (relocated use case here)
- 23 - Added a minimal placeholder for access hints for NVM.FILE
- 24 - Added TBD for per-block metadata for NVM.BLOCK
- 25 - Added discovery of granularities to NVM.BLOCK and NVM.FILE
- 26 - Added atomic write to NVM/FILE
- 27 - Extensive changes in NVM.PM.FILE per meeting discussion

28

### 29 **Revision 2**

#### 30 **Date**

31 April 16, 2013

#### 32 **Changes incorporated**

- 33 - Addressed comments from rev 1
- 34 - Added Compliance clause, relocated compliance related text from 1.1 and 1.2
- 35 - Relocated non-Scope content (1.3 to 1.13) to Common Programming Model
- 36 Behavior clause
- 37 - Added Block-optimized applications
- 38 - Added Deferred Behavior annex

- 39 - Removed Editorial Help Annex (moving to separate document)

40

### 41 **Revision 3**

#### 42 **Date**

43 May 8, 2013

#### 44 **Changes incorporated**

- 45 - Removed definition for page cache and virtual memory – commonly used external
- 46 definitions match our usage
- 47 - Remove Block Programming Overview text and examples
- 48 - Reworked NVM modes overview as a short description of key behavior from each
- 49 mode
- 50 - Reworked NVM.BLOCK and NVM.FILE mode overviews with sub-sections for types
- 51 of behavior covered (atomicity, granularities, ...)
- 52 - Remove text related to per-block metadata; created annex
- 53 - NVM.BLOCK: Removed access hints TDB and added access hints for deferred
- 54 behaviors annex
- 55 - Another pass at clarifying block-optimized application behavior
- 56 - Rework of Atomic Write content to address numerous reviewer comments
- 57 - Simplified description of file system at system startup (6.2)
- 58 - Relaxed “shall” clause and clarified relationship to other file implementations (10.1)
- 59 - Added quick reference on new NVM.PM.FILE actions. (10.1)
- 60 - Clarified reference to CPU cache (10.1)
- 61 - Changed COPY map option to COPY\_ON\_WRITE (10.2.3)
- 62 - Made MAP\_SHARED NVM.PM.FILE.MAP option mandatory and deleted
- 63 MAP\_SHARED\_CAPABLE attribute. (10.3)
- 64 - Moved UNCACHED NVM.PM.FILE.MAP option and
- 65 NVM.PM.FILE.DURABLE.WRITE to future work. (10.2)
- 66 - Clarified wording of error reporting situation (6.7)
- 67 - Declared certain permutations of interoperability between PM.FILE and other modes
- 68 to be unspecified (10.2.4)
- 69 - Changed ERASE option on NVM.PM.FILE.SET\_END\_OF\_FILE to ZERO (10.2.5)
- 70 - Simplified description of memory mapping to reference performance characteristics
- 71 suitable for memory programming models. (1)
- 72 - Removed Mapping to Native APIs Annex – this was intended to be short-lived

### 73 **Revision 4**

#### 74 **Date**

75 May 21, 2013

#### 76 **Changes incorporated**

- 77 - Updates to NVM.PM.BLOCK mode per discussions at May meeting
- 78 ○ Add ATOMIC\_WRITE action

- 79 ○ ATOMIC\_MULTIWRITE action – add attributes describing implementations limits
- 80 on number of ranges, sizes
- 81 ○ ISOLATED\_ATOMIC\_MULTIWRITE moved to Deferred Behavior annex
- 82 ○ BLOCK.DISCARD\_IF\_YOU\_MUST moved to Deferred Behavior annex
- 83 ○ Added DISCARD\_IMMEDIATELY\_RETURNS attribute: values are “zeros” or
- 84 “unspecified”, reference from DISCARD\_IMMEDIATELY
- 85 ○ ATOMICITY\_ERROR\_BEHAVIOR – delete this attribute (clean up “see also”
- 86 references)
- 87 ○ ECC\_BLOCK\_SIZE description added
- 88 - Added mandatory/optional info to everything in NVM.BLOCK and NVM.FILE
- 89 - Updates to NVM.PM.FILE mode per discussions at May meeting
- 90 - Updated NVM Pointer Annex
- 91 - Removed Per-block Metadata Annex
- 92 - Temporarily Removed Persistent Memory Error Handling Annex (updated version
- 93 being reviewed as separate document)

## 94 **Revision 5**

### 95 **Date**

96 June 12, 2013

### 97 **Changes incorporated**

- 98 - Updated cover and footer as Working Draft
- 99 - Updated ACS-2 Reference under Development to clarify this is rev 7
- 100 - Removed TBDs to add cross-reference to Error Handling Annex (which is not
- 101 included in this revision)
- 102 - Spelled out some abbreviations
- 103 - Increased size of diagrams that were fuzzy in PDFs
- 104 - Corrected heading level for several use cases

## 105 **Revision 6**

### 106 **Date**

107 - June 26, 2013

### 108 **Changes incorporated**

- 109 - Simplified Scope section 1 based on comments received on Rev 3.
- 110 - Added clarifications in NVM.PM.FILE overview section 10.1 regarding relationship to
- 111 the functionality of existing access methods and specificity to direct access via
- 112 memory mapping.
- 113 - Section 10.3.3 - Added PM file mapping exception to native file mapping behavior
- 114 emulation related to unmapping before synchronization.
- 115 - Section 10.3.4 – Added text to PM.SYNC reinforcing the notion that multiple
- 116 implementations may co-reside in a system based on implementation specific
- 117 mappings of the sync action.

- 118 - Removed section 10.4.2 on the ERASE CAPABLE attribute because the action that
- 119 used the attribute was removed earlier.
- 120 - Added revised PM Error Handling appendix D and corrected references there-to
- 121 - Changed NVM.ATTRIBUTE.GET to NVM.COMMON.SET\_ATTRIBUTE (similar for
- 122 GET) for consistency with other action/attribute names
- 123 - Moved all mandatory/optional designations for attributes/actions to the first line after
- 124 the heading
- 125 - Change “required” to “mandatory” in contexts where it’s used as a keyword
- 126 - Clean up text about common attributes/actions in all four modes
- 127 - Fix cross references throughout NVM.BLOCK and NVM.FILE
- 128 - Make text related to EXISTS states more consistent in several places in spec
- 129 - Changed “system” to “Implementation” in multiple places the text was inadvertently
- 130 limiting implementation options
- 131 - Added second version of device model, incorporating comments from TWG.
- 132 Examples in the device model are not complete

## 133 **Revision 7**

### 134 **Date**

- 135 - July 17, 2013

### 136 **Changes incorporated**

- 137 - Reformatted References section; merged in bibliography
- 138 - Corrected typos in scope section 1
- 139 - Added ECC block size and offset attributes to NVM.PM.FILE section
- 140 - Added section 6.2 describing interoperability between NVM.FILE and NVM.PM.FILE
- 141 modes
- 142 - Introduced *property group lists* and reworked references to parallel arrays of action
- 143 inputs/outputs to use property group lists
- 144 - Added NVM.BLOCK use cases demonstrating EXISTS/DISCARD actions and
- 145 SCAR
- 146 - Updated ECC\_BLOCK\_SIZE definitions to factor in power protection
- 147 - Filled in device models

## 148 **Revision 8**

### 149 **Date**

- 150 - July 28, 2013

### 151 **Changes incorporated**

- 152 - Added PERFORMANCE\_BLOCK\_SIZE attribute usage to BLOCK “update a record”
- 153 use case
- 154 - Filled in NVM.FILE atomic write use case
- 155 - Filled in some details to deferred behavior annex
- 156 - Changed ECC\_BLOCK\_SIZE to FUNDAMENTAL\_BLOCK\_SIZE and updated
- 157 related text to clarify that this may apply to other errors that had a “blast radius”

- 158 - Simplified working of “conformance with multiple file modes” section 6.2
- 159 - Added NVM.PM.FILE.OPTIMIZED\_FLUSH (section 10.2.5) and
- 160 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY (section 10.2.7) and related
- 161 attributes (section 10.3)
- 162 - Added NVM.PM.FILE.GET\_ERROR\_INFO (section 10.2.6) and related error
- 163 handling descriptions in NVM.PM.FILE.MAP (section 10.2.3) and the PM Error
- 164 Handling Annex.
- 165 - Removed text from NVM.PM.FILE.SYNC (section 10.2.4) that implied unintended
- 166 deviation from the native sync.
- 167 - Corrected major formatting error that caused displacement of section 10.2.4 in rev 7.

## 168 **Revision 9**

### 169 **Date**

- 170 - September 5, 2013

### 171 **Changes incorporated**

- 172 - Removed NVM Device definition and related text in NVM Device Models after
- 173 discussion at July face-to-face meeting
- 174 - Add a note to SET\_ATTRIBUTE action description saying it’s not used at this time
- 175 - Remove paragraph from NVM.PM.VOLUME saying management behavior is out-of-
- 176 scope (which is already stated in Scope clause).
- 177 - Scope clause – clarify that sharing NVM is not in scope.
- 178 - Incorporated use cases for NVM.PM.FILE, flash-as-cache, and NVM.PM.VOLUME
- 179 - Incorporated changes related to rev 8 ballot comments
- 180 - Generalized references to native file actions in NVM.PM.FILE
- 181 - Removed NVM.PM.FILE overview reference to device model
- 182 - Removed forward reference to “contained errors” in NVM.PM.FILE.MAP
- 183 - Generalized reference to “errors” in NVM.PM.FILE.ERROR\_EVENT\_CAPABLE
- 184 - Added file to address for error\_check and error\_clear actions in error handling
- 185 annex.
- 186 - Changed “Reasoning about Consistency” to “consistency”
- 187 - Removed reference to persistent media in
- 188 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY
- 189 - Added reference to persistence domain to
- 190 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY
- 191 - Changed references to media within NVM.PM use case sections 10.4.3 and 10.4.4
- 192 to refer to persistence domain
- 193 - Added an additional indication of error type to the outputs from
- 194 NVM.PM.FILE.GET\_ERROR\_INFO
- 195 - Removed load specific qualifiers from NVM.PM.FILE.GET\_ERROR\_INFO, enabling
- 196 (but not mandating) more general use
- 197 - Redid all images to address poor quality in generated PDF

198 **Revision 10**

199 **Date**

200 - September 30, 2013

201 **Changes incorporated**

202 - Added Acknowledgements

203 - Changes several occurrences of “this revision of the specification” to “this  
204 specification” in preparation for publication

205 - In 6.9 Persistence Domain, rewrote long “Once data...” sentence as two sentences

206 - clarified user of “mapped” (may refer to memory mapped files or non-discarded  
207 blocks/ranges)

208 - added “Atomic Sync/Flush action for PM” in deferred behavior

209 - removed text with requirements for file system mount from Device state at system  
210 startup

211 - removed gradients from several figures to enhance readability when printed

212 - fixed Word cross references that did not identify the target by section name

213 - fixed spelling of INTERRUPT

214 - removed extraneous space in NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE

215 - added “layout” to “address space layout randomization” in Annex A

216 - added transactional logging use cases

217 - In Scope, reworked last statement into separate sentences about remote and  
218 thread/process sharing.

219 - Section 10.2.7 – changed “read back” to “verified and eliminated text regarding  
220 hardware scope.

221 - Section 10.2.7 – changed reporting method of verify failures to an error code

222 - Section 10.3.3 – added attribute NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMIcity

223 - Section 10.3.4 – Added indication that fundamental error offset is not needed

224 - Section 10.4.2 – Clarified that error addresses are in the processor’s logical address  
225 space

226 - Added cross-references to section 10.

227 Suggestions for changes or modifications to this document should be submitted at  
228 <http://www.snia.org/feedback>. The SNIA hereby grants permission for individuals to use  
229 this document for personal use only, and for corporations and other business entities to  
230 use this document for internal use only (including internal copying, distribution, and  
231 display) provided that:

232 1. Any text, diagram, chart, table or definition reproduced must be reproduced in its  
233 entirety with no alteration, and,

234 2. Any document, printed or electronic, in which material from this document (or any  
235 portion hereof) is reproduced must acknowledge the SNIA copyright on that  
236 material, and must credit the SNIA for granting permission for its reuse.

237 Other than as explicitly provided above, you may not make any commercial use of this  
238 document, sell any or this entire document, or distribute this document to third parties.  
239 All rights not explicitly granted are expressly reserved to SNIA. Permission to use this  
240 document for purposes other than those enumerated above may be requested by e-  
241 mailing [tcmd@snia.org](mailto:tcmd@snia.org). Please include the identity of the requesting individual and/or  
242 company and a brief description of the purpose, nature, and scope of the requested  
243 use.

244 Copyright © 2013 Storage Networking Industry Association



246	<b>FOREWORD</b> .....	<b>13</b>
247	<b>1 SCOPE</b> .....	<b>15</b>
248	<b>2 REFERENCES</b> .....	<b>16</b>
249	<b>3 DEFINITIONS, ABBREVIATIONS, AND CONVENTIONS</b> .....	<b>17</b>
250	3.1 DEFINITIONS.....	17
251	3.2 KEYWORDS .....	18
252	3.3 ABBREVIATIONS .....	18
253	3.4 CONVENTIONS .....	19
254	<b>4 OVERVIEW OF THE NVM PROGRAMMING MODEL (INFORMATIVE)</b> .....	<b>20</b>
255	4.1 HOW TO READ AND USE THIS SPECIFICATION .....	20
256	4.2 NVM DEVICE MODELS .....	20
257	4.3 NVM PROGRAMMING MODES .....	22
258	4.4 INTRODUCTION TO ACTIONS, ATTRIBUTES, AND USE CASES .....	24
259	<b>5 COMPLIANCE TO THE PROGRAMMING MODEL</b> .....	<b>26</b>
260	5.1 OVERVIEW .....	26
261	5.2 DOCUMENTATION OF MAPPING TO APIS .....	26
262	5.3 COMPATIBILITY WITH UNSPECIFIED NATIVE ACTIONS .....	26
263	5.4 MAPPING TO NATIVE INTERFACES .....	26
264	<b>6 COMMON PROGRAMMING MODEL BEHAVIOR</b> .....	<b>27</b>
265	6.1 OVERVIEW .....	27
266	6.2 CONFORMANCE TO MULTIPLE FILE MODES .....	27
267	6.3 DEVICE STATE AT SYSTEM STARTUP .....	27
268	6.4 SECURE ERASE.....	27
269	6.5 ALLOCATION OF SPACE.....	27
270	6.6 INTERACTION WITH I/O DEVICES .....	28
271	6.7 NVM STATE AFTER A MEDIA OR CONNECTION FAILURE .....	28

272	6.8	ERROR HANDLING FOR PERSISTENT MEMORY .....	28
273	6.9	PERSISTENCE DOMAIN.....	28
274	6.10	COMMON ACTIONS .....	29
275	6.11	COMMON ATTRIBUTES .....	29
276	6.12	USE CASES .....	30
277	<b>7</b>	<b>NVM.BLOCK MODE.....</b>	<b>32</b>
278	7.1	OVERVIEW .....	32
279	7.2	ACTIONS .....	34
280	7.3	ATTRIBUTES .....	38
281	7.4	USE CASES .....	41
282	<b>8</b>	<b>NVM.FILE MODE.....</b>	<b>46</b>
283	8.1	OVERVIEW .....	46
284	8.2	ACTIONS .....	46
285	8.3	ATTRIBUTES .....	48
286	8.4	USE CASES .....	50
287	<b>9</b>	<b>NVM.PM.VOLUME MODE.....</b>	<b>56</b>
288	9.1	OVERVIEW .....	56
289	9.2	ACTIONS .....	56
290	9.3	ATTRIBUTES .....	59
291	9.4	USE CASES .....	61
292	<b>10</b>	<b>NVM.PM.FILE .....</b>	<b>64</b>
293	10.1	OVERVIEW .....	64
294	10.2	ACTIONS .....	65
295	10.3	ATTRIBUTES .....	71
296	10.4	USE CASES .....	72
297	<b>ANNEX A</b>	<b>(INFORMATIVE) NVM POINTERS .....</b>	<b>82</b>

298	<b>ANNEX B (INFORMATIVE) CONSISTENCY .....</b>	<b>83</b>
299	<b>ANNEX C (INFORMATIVE) PM ERROR HANDLING .....</b>	<b>87</b>
300	<b>ANNEX D (INFORMATIVE) DEFERRED BEHAVIOR.....</b>	<b>91</b>
301	D.1 REMOTE SHARING OF NVM .....	91
302	D.2 MAP_CACHED OPTION FOR NVM.PM.FILE.MAP .....	91
303	D.3 NVM.PM.FILE.DURABLE.STORE .....	91
304	D.4 ENHANCED NVM.PM.FILE.WRITE.....	91
305	D.5 MANAGEMENT-ONLY BEHAVIOR.....	91
306	D.6 ACCESS HINTS.....	91
307	D.7 MULTI-DEVICE ATOMIC MULTI-WRITE ACTION.....	91
308	D.8 NVM.BLOCK.DISCARD_IF_YOU_MUST ACTION.....	91
309	D.9 ATOMIC WRITE ACTION WITH ISOLATION .....	93
310	D.10 ATOMIC SYNC/FLUSH ACTION FOR PM .....	93
311	D.11 HARDWARE-ASSISTED VERIFY .....	93
312		

313	Figure 1 Block NVM example.....	21
314	Figure 2 PM example .....	21
315	Figure 3 Block volume using PM HW .....	21
316	Figure 4 NVM.BLOCK and NVM.FILE mode examples .....	22
317	Figure 5 NVM.PM.VOLUME and NVM.PM.FILE mode examples.....	23
318	Figure 6 NVM.BLOCK mode example.....	32
319	Figure 7 SSC in a storage stack.....	42
320	Figure 8 SSC software cache application.....	42
321	Figure 9 SSC with caching assistance .....	43
322	Figure 10 NVM.FILE mode example .....	46
323	Figure 11 NVM.PM.VOLUME mode example .....	56
324	Figure 12 Zero range offset example .....	60
325	Figure 13 Non-zero range offset example .....	60
326	Figure 14 NVM.PM.FILE mode example.....	64
327	Figure 15 Consistency overview .....	83
328	Figure 16 Linux Machine Check error flow with proposed new interface.....	89
329		

330 **FOREWORD**

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

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

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

344 **Acknowledgements**

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

348	<i>Organization Represented</i>	<i>Name of Representative</i>
349	EMC	Bob Beauchamp
350	Hewlett Packard	Hans Boehm
351	NetApp	Steve Byan
352	Hewlett Packard	Joe Foster
353	Fusion-io	Walt Hubis
354	Red Hat	Jeff Moyer
355	Fusion-io	Ned Plasson
356	Rougs, LLC	Tony Roug
357	Intel Corporation	Andy Rudoff
358	Microsoft	Spencer Shepler
359	Fusion-io	Nisha Talagata
360	Hewlett Packard	Doug Voigt
361	Intel Corporation	Paul von Behren

362  
363 **SNIA Web Site**

364 Current SNIA practice is to make updates and other information available through their  
365 web site at <http://www.snia.org>

366 **SNIA Address**

367 Requests for interpretation, suggestions for improvement and addenda, or defect  
368 reports are welcome. They should be sent via the SNIA Feedback Portal at  
369 <http://www.snia.org/feedback/> or by mail to the Storage Networking Industry  
370 Association, 425 Market Street, Suite 1020, San Francisco, CA 94105, U.S.A.

## 371 **1 Scope**

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

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

- 384 • supported access modes,
- 385 • visibility in memory address space,
- 386 • atomicity and durability,
- 387 • recognizing, reporting, and recovering from errors and failures,
- 388 • data granularity, and
- 389 • capacity reclamation.

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

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

404  
405 This revision of the specification does not address sharing data across computing  
406 nodes. This revision of the specification assumes that sharing data between processes  
407 and threads follows the native OS behavior.

## 408 **2 References**

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

411 For references available from ANSI, contact ANSI Customer Service Department at  
412 (212) 642-4900/4980 (phone), (212) 302-1286 (fax) or via the World Wide Web at  
413 <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 SO/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>

414



## 415 **3 Definitions, abbreviations, and conventions**

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

### 417 **3.1 Definitions**

#### 418 3.1.1 **durable**

419 committed to a persistence domain

#### 420 3.1.2 **load and store operations**

421 commands to move data between CPU registers and memory

#### 422 3.1.3 **memory-mapped file**

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

#### 425 3.1.4 **non-volatile memory**

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

#### 428 3.1.5 **NVM block capable driver**

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

#### 430 3.1.6 **NVM volume**

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

432 See 4.2 NVM device models

#### 433 3.1.7 **persistence domain**

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

436 See 6.9 Persistence domain

#### 437 3.1.8 **persistent memory**

438 storage technology with performance characteristics suitable for a load and store  
439 programming model

#### 440 3.1.9 **programming model**

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

443 **3.2 Keywords**

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

446 **3.2.1 mandatory**

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

449 **3.2.2 may**

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

452 **3.2.3 may not**

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

455 **3.2.4 need not**

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

458 **3.2.5 optional**

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

462 **3.2.6 shall**

463 a keyword indicating a mandatory requirement; designers are required to implement all  
464 such mandatory requirements to ensure interoperability with other products that  
465 conform to this standard

466 **3.2.7 should**

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

468 **3.3 Abbreviations**

469 ACID Atomicity, Consistency, Isolation, Durability

470 NVM Non-Volatile Memory

471 PM Persistent Memory

472 SSD Solid State Disk

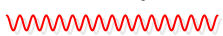
473 **3.4 Conventions**

474 The nomenclature used for binary power multiplier values in this standard is based on  
475 IEC 60027:2000, *Letter symbols to be used in electrical technology - Part 2:*  
476 *Telecommunications and electronics:*

- 477 • one kibibit is 1 Kib is 1,024 bits
- 478 • one mebibyte is 1 MiB is 1,048 576 bytes
- 479 • one gibibyte is 1 GiB is 1,073,741,824 bytes

480 **Representation of modes in figures**

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

482 

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

## 485 **4 Overview of the NVM Programming Model (informative)**

### 486 **4.1 How to read and use this specification**

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

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

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

- 504 1. The NVM programming modes section identifies section 10 NVM.PM.FILE mode  
505 as the starting point for application use of persistent memory.
- 506 2. The NVM.PM.FILE mode text describes the general approach for accessing PM  
507 (similar to native memory-mapped files) and the role of PM aware file system.
- 508 3. The NVM.PM.FILE mode identifies the NVM.PM.FILE.MAP and  
509 NVM.PM.FILE.SYNC actions and attributes that allow an application to discover support  
510 for optional features.
- 511 4. The operating system vendor's mapping document describes the mapping  
512 between NVM.PM.FILE.MAP/SYNC and API calls, and also provides information about  
513 supported PM-aware file systems.

### 514 **4.2 NVM device models**

#### 515 **4.2.1 Overview**

516 This section describes device models for NVM to help readers understand how key  
517 terms in the programming model relate to other software and hardware. The models  
518 presented here generally apply across operating systems, file systems, and hardware;

519 but there are differences across implementations. This specification strives to discuss  
520 the model generically, but mentions key exceptions.

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

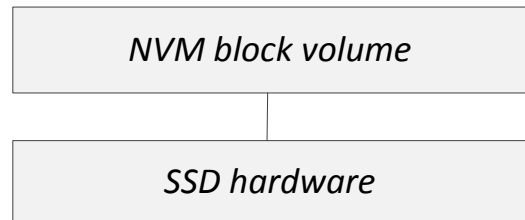
526 *An NVM volume* is a subset of one or more NVM devices, treated by software as a  
527 single logical entity. For the purposes of this specification, a volume is a container of  
528 storage. A volume may be block capable and may be persistent memory capable. The  
529 consumer of a volume sees its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

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

#### 535 4.2.2 Block NVM example

536 Consider a single drive form factor SSD where  
537 the entire SSD capacity is dedicated to a file  
538 system. In this case, a single NVM block volume  
539 maps to a single hardware device. A file system  
540 (not depicted) is mounted on the NVM block  
541 volume.

Figure 1 Block NVM example

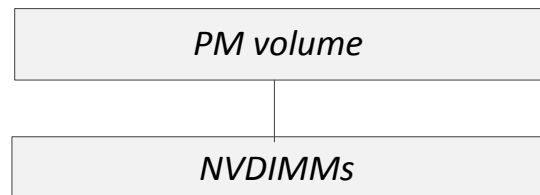


542 The same model may apply to NVM block hardware other than an SSD (including flash  
543 on PCIe cards).

#### 544 4.2.3 Persistent memory example

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

Figure 2 PM example

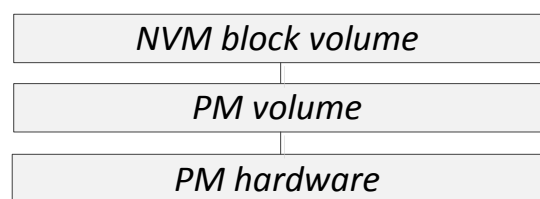


548 The same model may apply to PM hardware  
549 other than an NVDIMM (including SSDs, PCIe  
550 cards, etc.).

#### 551 4.2.4 NVM block volume using PM hardware

552 In this example, the persistent memory  
553 implementation includes a driver that uses a  
554 range of persistent memory (a PM volume) and  
555 makes it appear to be a block NVM device in  
556 the legacy block stack. This emulated block  
557 device could be aggregated or de-aggregated

Figure 3 Block volume using PM HW



558 like legacy block devices. In this example, the emulated block device is mapped 1-1 to  
559 an NVM block volume and non-PM file system.

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

### 561 4.3 NVM programming modes

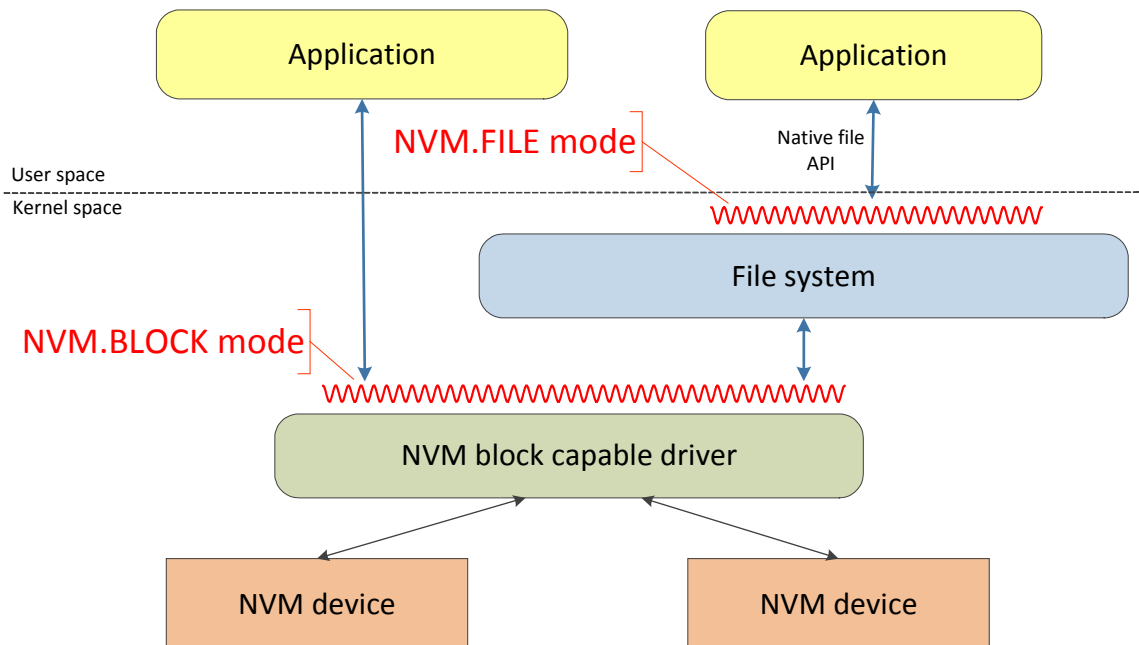
#### 562 4.3.1 NVM.BLOCK mode overview

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

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

- 571 • Discovery and use of atomic write and discard features
- 572 • The discovery of granularities (length or alignment characteristics)
- 573 • Discovery and use of ability for applications or operating system components to  
574 mark blocks as unreadable

575  
576 **Figure 4 NVM.BLOCK and NVM.FILE mode examples**



#### 577 578 4.3.2 NVM.FILE mode overview

579 NVM.FILE mode is used by applications that are not aware of details of block storage  
580 hardware or addresses. Existing applications written using native file I/O behavior

581 should work unmodified with NVM.FILE mode; adding support in the application for  
582 NVM extensions may optimize the application.

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

585 The NVM.FILE mode describes NVM extensions including:

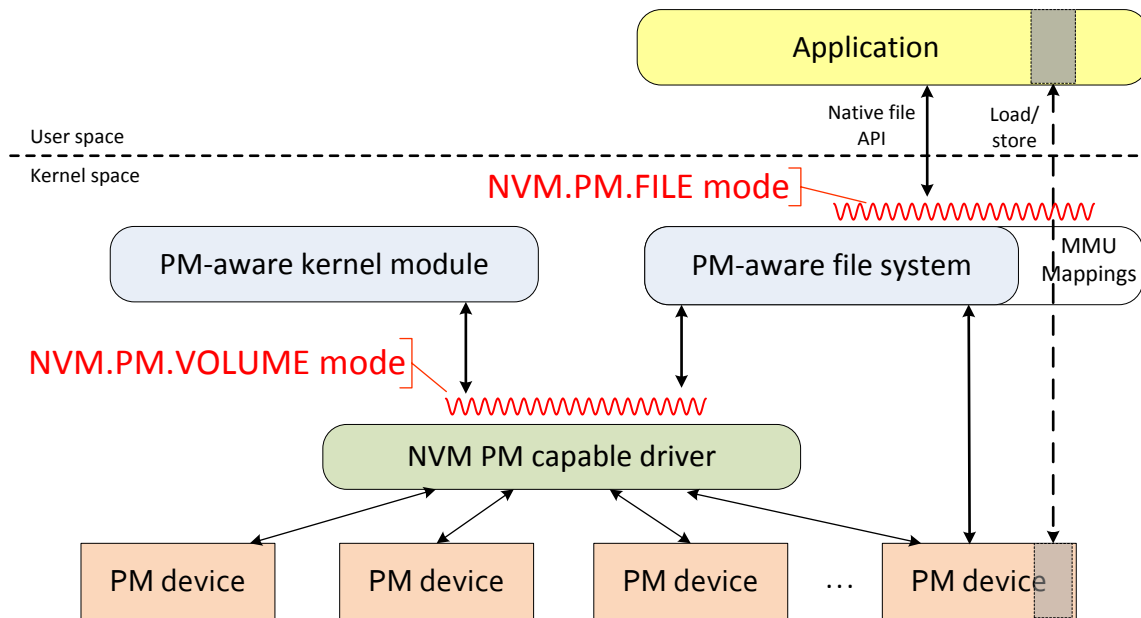
- 586 • Discovery and use of atomic write features
- 587 • The discovery of granularities (length or alignment characteristics)

### 588 4.3.3 NVM.PM.VOLUME mode overview

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

- 593 • the list of physical address ranges associated with each PM volume
- 594 • the capability to determine whether PM errors have been reported

595  
596 **Figure 5 NVM.PM.VOLUME and NVM.PM.FILE mode examples**



597

### 598 4.3.4 NVM.PM.FILE mode overview

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

- 603 • mapping PM files (or subsets of files) to virtual memory addresses

- 604 • syncing portions of PM files to the persistence domain

## 605 **4.4 Introduction to actions, attributes, and use cases**

### 606 4.4.1 **Overview**

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

### 613 4.4.2 **Use cases**

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

- 619 • a purpose and context including actors involved in the use case;
- 620 • triggers and preconditions indicating when a use case applies;
- 621 • inputs, outputs, events and actions that occur during the use case;
- 622 • references to related materials or concepts including other use cases that use or  
623 extend the use case.

### 624 4.4.3 **Actions**

625 Actions are defined using the following naming convention:

626 <context>.<mode>.<verb>

627 The actions in this specification all have a context of “NVM”. The mode refers to one of  
628 the NVM models documented herein (or “COMMON” for actions used in multiple  
629 modes). The verb states what the action does. Examples of actions include  
630 “NVM.COMMON.GET\_ATTRIBUTE” and “NVM.FILE.ATOMIC\_WRITE”. In some cases  
631 native actions that are not explicitly specified by the programming model are referenced  
632 to illustrate usage.

633 The description of each action includes:

- 634 • parameters and results of the action
- 635 • details of the action’s behavior
- 636 • compatibility of the action with pre-existing APIs in the industry

637 A number of actions involve options that can be specified each time the action is used.  
638 The options are given names that begin with the name of the action and end with a



639 descriptive term that is unique for the action. Examples include  
640 NVM.PM.FILE.MAP\_COPY\_ON\_WRITE and NVM.PM.FILE.MAP\_SHARED.

641 A number of actions are optional. For each of these, there is an attribute that indicates  
642 whether the action is supported by the implementation in question. By convention these  
643 attributes end with the term “CAPABLE” such as  
644 NVM.BLOCK.ATOMIC\_WRITE\_CAPABLE. Supported options are also enumerated by  
645 attributes that end in “CAPABLE”.

#### 646 4.4.4 **Attributes**

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

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

#### 661 4.4.5 **Property group lists**

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

## 667 **5 Compliance to the programming model**

### 668 **5.1 Overview**

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

673 Compliance of an implementation shall be interpreted as follows.

### 674 **5.2 Documentation of mapping to APIs**

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

### 678 **5.3 Compatibility with unspecified native actions**

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

### 683 **5.4 Mapping to native interfaces**

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

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

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

## 699 **6 Common programming model behavior**

### 700 **6.1 Overview**

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

### 703 **6.2 Conformance to multiple file modes**

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

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

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

### 716 **6.3 Device state at system startup**

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

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

### 722 **6.4 Secure erase**

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

### 725 **6.5 Allocation of space**

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

## 729 **6.6 Interaction with I/O devices**

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

## 733 **6.7 NVM State after a media or connection failure**

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

## 738 **6.8 Error handling for persistent memory**

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

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

752 A proposal for persistent memory error handling appears as an appendix; see Annex C.

## 753 **6.9 Persistence domain**

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

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

763 **6.10 Common actions**

764 6.10.1 **NVM.COMMON.GET\_ATTRIBUTE**

765 Requirement: mandatory

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

769 **Inputs:**

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

772 **Outputs:**

- 773 • value of attribute

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

776 6.10.2 **NVM.COMMON.SET\_ATTRIBUTE**

777 Requirement: optional

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

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

782 **Inputs:**

- 783 • reference to appropriate instance
- 784 • attribute name
- 785 • value to be assigned to the attribute

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

788 **6.11 Common attributes**

789 6.11.1 **NVM.COMMON.SUPPORTED\_MODES**

790 Requirement: mandatory

791 SUPPORTED\_MODES returns a list of the modes supported by the NVM  
792 implementation.

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

794 NVM.COMMON.SET\_ATTRIBUTE is not supported for  
795 NVM.COMMON.SUPPORTED\_MODES.

## 796 6.11.2 **NVM.COMMON.FILE\_MODE**

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

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

800 Target: a file path

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

802 See 6.2 Conformance to multiple file modes.

## 803 **6.12 Use cases**

### 804 6.12.1 **Application determines which mode is used to access a file system**

#### 805 **Purpose/triggers:**

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

#### 808 **Scope/context:**

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

#### 815 **Preconditions:**

816 None

#### 817 **Inputs:**

818 None

#### 819 **Success scenario:**

- 820 1) Invoke NVM.COMMON.GET\_ATTRIBUTE (NVM.COMMON.FILE\_MODE)  
821 targeting a file path; the value returned provides information on which modes  
822 may be used to access the data.
- 823 2) If the response includes “NVM.FILE”, then the actions and attributes described  
824 for the NVM.FILE mode are supported. Set the preferred mode for this file  
825 system to NVM.FILE.

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

829 **Outputs:**

830 **Postconditions:**

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

832 See also:

833 6.2 Conformance to multiple file modes

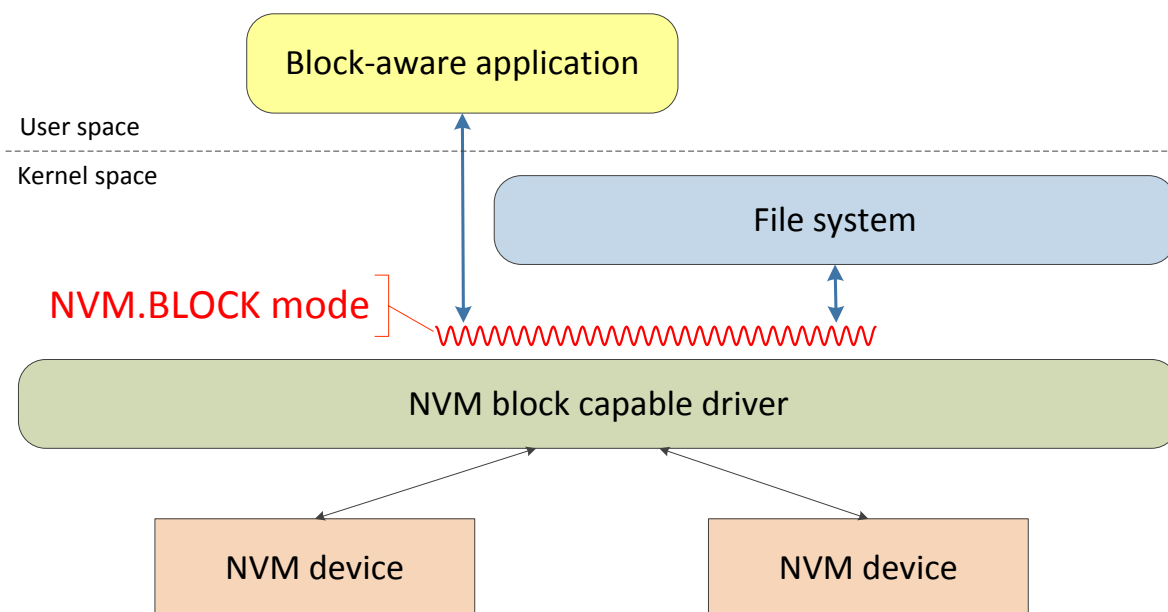
834 6.11.2 NVM.COMMON.FILE\_MODE

## 835 7 NVM.BLOCK mode

### 836 7.1 Overview

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

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



845

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

849 The NVM extensions supported by this mode include:

- 850 • Discovery and use of atomic write and discard features
- 851 • The discovery of granularities (length or alignment characteristics)
- 852 • Discovery and use of per-block metadata used for verifying integrity
- 853 • Discovery and use of ability for applications or operating system components to  
854 mark blocks as unreadable

855

#### 856 7.1.1 Discovery and use of atomic write features

857 Atomic Write support provides applications with the capability to assure that all the data  
858 for an operation is written to the persistence domain or, if a failure occurs, it appears



859 that no operation took place. Applications may use atomic write operations to assure  
860 consistent behavior during a failure condition or to assure consistency between multiple  
861 processes accessing data simultaneously.  
862

### 863 **7.1.2 The discovery of granularities**

864 Attributes are introduced to allow applications to discover granularities associated with  
865 NVM devices.  
866

### 867 **7.1.3 Discovery and use of capability to mark blocks as unreadable**

868 An action (NVM.BLOCK.SCAR) is defined allowing an application to mark blocks as  
869 unreadable.  
870

### 871 **7.1.4 NVM.BLOCK consumers: operating system and applications**

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

#### 874 **7.1.4.1 NVM.BLOCK operating system components**

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

#### 878 **7.1.4.2 Block-optimized applications**

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

- 883 • A mechanism for a block-optimized application to request that the file system move  
884 data directly between the device and application memory, bypassing the buffering  
885 typically provided by the file system.
- 886 • The operating system or file system may require the application to align requests on  
887 block boundaries.

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

#### 890 **7.1.4.3 Mapping documentation**

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

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

## 899 **7.2 Actions**

### 900 **7.2.1 Actions that apply across multiple modes**

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

902 NVM.COMMON.GET\_ATTRIBUTE (see 6.10.1)

903 NVM.COMMON.SET\_ATTRIBUTE (see 6.10.2)

### 904 **7.2.2 NVM.BLOCK.ATOMIC\_WRITE**

905 Requirement: mandatory if ATOMIC\_WRITE\_CAPABLE (see 7.3.1) is true

906 Block-optimized applications or operating system components may use  
907 ATOMIC\_WRITE to assure consistent behavior during a power failure condition. This  
908 specification does not specify the order in which this action occurs relative to other I/O  
909 operations, including other ATOMIC\_WRITE or ATOMIC\_MULTIWRITE actions. This  
910 specification does not specify when the data written becomes visible to other threads.

#### 911 **Inputs:**

- 912 • the starting memory address
- 913 • a reference to the block device
- 914 • the starting block address
- 915 • the length

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

919

920 **Output:** none

#### 921 **Return values:**

- 922 • Success shall be returned if all blocks are updated in the persistence domain
- 923 • an error shall be reported if the length exceeds  
924 ATOMIC\_WRITE\_MAX\_DATA\_LENGTH (see 7.3.3)
- 925 • an error shall be reported if the starting address is not evenly divisible by  
926 ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY (see 7.3.4)
- 927 • an error shall be reported if the length is not evenly divisible by  
928 ATOMIC\_WRITE\_LENGTH\_GRANULARITY (see 7.3.5)
- 929 • If anything does or will prevent all of the blocks from being updated in the  
930 persistence domain before completion of the operation, an error shall be reported  
931 and all the logical blocks affected by the operation shall contain the data that was  
932 present before the device server started processing the write operation (i.e., the old  
933 data, as if the atomic write operation had no effect). If the NVM and processor are

934 both impacted by a power failure, no error will be returned since the execution  
935 context is lost.  
936 • the different errors described above shall be discernible by the consumer and shall  
937 be discernible from media errors

938 **Relevant attributes:**

939 ATOMIC\_WRITE\_MAX\_DATA\_LENGTH (see 7.3.3)  
940 ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY (see 7.3.4)  
941 ATOMIC\_WRITE\_LENGTH\_GRANULARITY (see 7.3.5)  
942 ATOMIC\_WRITE\_CAPABLE (see 7.3.1)

943 **7.2.3 NVM.BLOCK.ATOMIC\_MULTIWRITE**

944 Requirement: mandatory if ATOMIC\_MULTIWRITE\_CAPABLE (see 7.3.6) is true

945 Block-optimized applications or operating system components may use  
946 ATOMIC\_MULTIWRITE to assure consistent behavior during a power failure condition.  
947 This action allows a caller to write non-adjacent extents atomically. The caller of  
948 ATOMIC\_MULTIWRITE provides a Property Group List (see 4.4.5) where the properties  
949 describe the memory and block extents (see Inputs below); all of the extents are written  
950 as a single atomic operation. This specification does not specify the order in which this  
951 action occurs relative to other I/O operations, including other ATOMIC\_WRITE or  
952 ATOMIC\_MULTIWRITE actions. This specification does not specify when the data  
953 written becomes visible to other threads.

954 **Inputs:**

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

- 956 • memory address starting address
- 957 • length of data to write (in bytes)
- 958 • a reference to the device being written to
- 959 • the starting LBA on the device

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

964 **Output:** none

965 **Return values:**

- 966 • Success shall be returned if all block ranges are updated in the persistence domain
- 967 • an error shall be reported if the block ranges overlap
- 968 • an error shall be reported if the total size of memory input ranges exceeds  
969 ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH (see 7.3.8)
- 970 • an error shall be reported if the starting address in any input memory range is not  
971 evenly divisible by  
972 ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY (see 7.3.9)

- 973 • an error shall be reported if the length in any input range is not evenly divisible by  
974 ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY (see 7.3.10)
- 975 • If anything does or will prevent all of the writes from being applied to the persistence  
976 domain before completion of the operation, an error shall be reported and all the  
977 logical blocks affected by the operation shall contain the data that was present  
978 before the device server started processing the write operation (i.e., the old data, as  
979 if the atomic write operation had no effect). If the NVM and processor are both  
980 impacted by a power failure, no error will be returned since the execution context is  
981 lost.
- 982 • the different errors described above shall be discernible by the consumer

983 **Relevant attributes:**

- 984 ATOMIC\_MULTIWRITE\_MAX\_IOS (see 7.3.7)
- 985 ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH (see 7.3.8)
- 986 ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY (see 7.3.9)
- 987 ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY (see 7.3.10)
- 988 ATOMIC\_MULTIWRITE\_CAPABLE (see 7.3.6)

989 **7.2.4 NVM.BLOCK.DISCARD\_IF\_YOU\_CAN**

990 Requirement: mandatory if DISCARD\_IF\_YOU\_CAN\_CAPABLE (see 7.3.17) is true

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

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

998  
999 Inputs: a range of blocks (starting LBA and length in logical blocks)

1000 Status: Success indicates the request is accepted but not necessarily acted upon.

1001 **7.2.5 NVM.BLOCK.DISCARD\_IMMEDIATELY**

1002 Requirement: mandatory if DISCARD\_IMMEDIATELY\_CAPABLE (see 7.3.18) is true

1003 Requires that the data block be unmapped (see NVM.BLOCK.EXISTS 7.2.6) before the  
1004 next READ or WRITE reference even if garbage collection of the block has not occurred  
1005 yet,

1006 DISCARD\_IMMEDIATELY commands cannot be acknowledged by the NVM device  
1007 until the DISCARD\_IMMEDIATELY has been durably written to media in a way such  
1008 that upon recovery from a power-fail event, the block is guaranteed to remain discarded.

1009 Inputs: a range of blocks (starting LBA and length in logical blocks)

1010 The values returned by subsequent read operations are specified by the  
1011 DISCARD\_IMMEDIATELY\_RETURNS (see 7.3.19) attribute.

1012 Status: Success indicates the request is completed.

1013 See also EXISTS (7.2.6), DISCARD\_IMMEDIATELY\_RETURNS (7.3.19),  
1014 DISCARD\_IMMEDIATELY\_CAPABLE (7.3.18).

1015 **7.2.6 NVM.BLOCK.EXISTS**

1016 Requirement: mandatory if EXISTS\_CAPABLE (see 9.3.9) is true

1017 An NVM device may allocate storage through a thin provisioning mechanism or one of  
1018 the discard actions. As a result, a block can exist in one of three states:

1019 • **Mapped**: the block has had data written to it  
1020 • **Unmapped**: the block has not been written, and there is no memory allocated  
1021 • **Allocated**: the block has not been written, but has memory allocated to it

1022 The EXISTS action allows the NVM user to determine if a block has been allocated.

1023 Inputs: an LBA

1024 Output: the state (mapped, unmapped, or allocated) for the input block

1025 Result: the status of the action

1026 **7.2.7 NVM.BLOCK.SCAR**

1027 Requirement: mandatory if SCAR\_CAPABLE (see 7.3.13) is true

1028 This action allows an application to request that subsequent reads from any of the  
1029 blocks in the address range will cause an error. This action uses an implementation-  
1030 dependent means to insure that all future reads to any given block from the scarred  
1031 range will cause an error until new data is stored to any given block in the range. A  
1032 block stays scared until it is updated by a write operation.

1033 Inputs: reference to a block volume, starting offset, length

1034 Outputs: status

1035 Relevant attributes:

1036 NVM.BLOCK.SCAR\_CAPABLE (7.3.13) – Indicates that the SCAR action is  
1037 supported.

1038 **7.3 Attributes**

1039 **7.3.1 Attributes that apply across multiple modes**

1040 The following attributes apply to NVM.BLOCK mode as well as other modes.

1041 NVM.COMMON.SUPPORTED\_MODES (see 6.11.1)

1042 **7.3.2 NVM.BLOCK.ATOMIC\_WRITE\_CAPABLE**

1043 Requirement: mandatory

1044 This attribute indicates that the implementation is capable of the  
1045 NVM.BLOCK.ATOMIC\_WRITE action.

1046 **7.3.3 NVM.BLOCK.ATOMIC\_WRITE\_MAX\_DATA\_LENGTH**

1047 Requirement: mandatory if ATOMIC\_WRITE\_CAPABLE (see 7.3.1) is true.

1048 ATOMIC\_WRITE\_MAX\_DATA\_LENGTH is the maximum length of data that can be  
1049 transferred by an ATOMIC\_WRITE action.

1050 **7.3.4 NVM.BLOCK.ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY**

1051 Requirement: mandatory if ATOMIC\_WRITE\_CAPABLE (see 7.3.1) is true.

1052 ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY is the granularity of the  
1053 starting memory address for an ATOMIC\_WRITE action. Address inputs to  
1054 ATOMIC\_WRITE shall be evenly divisible by  
1055 ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY.

1056 **7.3.5 NVM.BLOCK.ATOMIC\_WRITE\_LENGTH\_GRANULARITY**

1057 Requirement: mandatory if ATOMIC\_WRITE\_CAPABLE (see 7.3.1) is true.

1058 ATOMIC\_WRITE\_LENGTH\_GRANULARITY is the granularity of the length of data  
1059 transferred by an ATOMIC\_WRITE action. Length inputs to ATOMIC\_WRITE shall be  
1060 evenly divisible by ATOMIC\_WRITE\_LENGTH\_GRANULARITY.

1061 **7.3.6 NVM.BLOCK.ATOMIC\_MULTIWRITE\_CAPABLE**

1062 Requirement: mandatory

1063 ATOMIC\_MULTIWRITE\_CAPABLE indicates that the implementation is capable of the  
1064 NVM.BLOCK.ATOMIC\_MULTIWRITE action.

1065 **7.3.7 NVM.BLOCK.ATOMIC\_MULTIWRITE\_MAX\_IOS**

1066 Requirement: mandatory if ATOMIC\_MULTIWRITE\_CAPABLE (see 7.3.6) is true

1067 ATOMIC\_MULTIWRITE\_MAX\_IOS is the maximum length of the number of IOs (i.e.,  
1068 the size of the Property Group List) that can be transferred by an  
1069 ATOMIC\_MULTIWRITE action.

1070 **7.3.8 NVM.BLOCK.ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH**

1071 Requirement: mandatory if ATOMIC\_MULTIWRITE\_CAPABLE (see 7.3.6) is true

1072 ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH is the maximum length of data that can  
1073 be transferred by an ATOMIC\_MULTIWRITE action.

1074 **7.3.9 NVM.BLOCK.ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARIT**  
1075 **Y**

1076 Requirement: mandatory if ATOMIC\_MULTIWRITE\_CAPABLE (see 7.3.6) is true

1077 ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY is the granularity of  
1078 the starting address of ATOMIC\_MULTIWRITE inputs. Address inputs to  
1079 ATOMIC\_MULTIWRITE shall be evenly divisible by  
1080 ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY.

1081 **7.3.10 NVM.BLOCK.ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY**

1082 Requirement: mandatory if ATOMIC\_MULTIWRITE\_CAPABLE (see 7.3.6) is true

1083 ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY is the granularity of the length of  
1084 ATOMIC\_MULTIWRITE inputs. Length inputs to ATOMIC\_MULTIWRITE shall be  
1085 evenly divisible by ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY.

1086 **7.3.11 NVM.BLOCK.WRITE\_ATOMICITY\_UNIT**

1087 Requirement: mandatory

1088 If a write is submitted of this size or less, the caller is guaranteed that if power is lost  
1089 before the data is completely written, then the NVM device shall ensure that all the  
1090 logical blocks affected by the operation contain the data that was present before the  
1091 device server started processing the write operation (i.e., the old data, as if the atomic  
1092 write operation had no effect).

1093 If the NVM device can't assure that at least one LOGICAL\_BLOCKSIZE (see 7.3.14)  
1094 extent can be written atomically, WRITE\_ATOMICITY\_UNIT shall be set to zero.

1095 The unit is NVM.BLOCK.LOGICAL\_BLOCKSIZE (see 7.3.14).

1096 **7.3.12 NVM.BLOCK.EXISTS\_CAPABLE**

1097 Requirement: mandatory

1098 This attribute indicates that the implementation is capable of the NVM.BLOCK.EXISTS  
1099 action.

1100 **7.3.13 NVM.BLOCK.SCAR\_CAPABLE**

1101 Requirement: mandatory

1102 This attribute indicates that the implementation is capable of the NVM.BLOCK.SCAR  
1103 (see 7.2.7) action.

1104 **7.3.14 NVM.BLOCK.LOGICAL\_BLOCK\_SIZE**

1105 Requirement: mandatory

1106 LOGICAL\_BLOCK\_SIZE is the smallest unit of data (in bytes) that may be logically read  
1107 or written from the NVM volume.

1108 **7.3.15 NVM.BLOCK.PERFORMANCE\_BLOCK\_SIZE**

1109 Requirement: mandatory

1110 PERFORMANCE\_BLOCK\_SIZE is the recommended granule (in bytes) the caller  
1111 should use in I/O requests for optimal performance; starting addresses and lengths  
1112 should be multiples of this attribute. For example, this attribute may help minimizing  
1113 device-implemented read/modify/write behavior.

1114 **7.3.16 NVM.BLOCK.ALLOCATION\_BLOCK\_SIZE**

1115 Requirement: mandatory

1116 ALLOCATION\_BLOCK\_SIZE is the recommended granule (in bytes) for allocation and  
1117 alignment of data. Allocations smaller than this attribute (even if they are multiples of  
1118 LOGICAL\_BLOCK\_SIZE) may work, but may not yield optimal lifespan.

1119 **7.3.17 NVM.BLOCK.DISCARD\_IF\_YOU\_CAN\_CAPABLE**

1120 Requirement: mandatory

1121 DISCARD\_IF\_YOU\_CAN\_CAPABLE shall be set to true if the implementation supports  
1122 DISCARD\_IF\_YOU\_CAN.

1123 **7.3.18 NVM.BLOCK.DISCARD\_IMMEDIATELY\_CAPABLE**

1124 Requirement: mandatory

1125 Returns true if the implementation supports DISCARD\_IMMEDIATELY.

1126 **7.3.19 NVM.BLOCK.DISCARD\_IMMEDIATELY\_RETURNS**

1127 Requirement: mandatory if DISCARD\_IMMEDIATELY\_CAPABLE (see 7.3.18) is true

1128 The value returned from read operations to blocks specified by a  
1129 DISCARD\_IMMEDIATELY action with no subsequent write operations. The possible  
1130 values are:

1131 • A value that is returned to each read of an unmapped block (see  
1132 NVM.BLOCK.EXISTS 7.2.6) until the next write action



1133 • Unspecified

### 1134 7.3.20 **NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE**

1135 Requirement: mandatory

1136 FUNDAMENTAL\_BLOCK\_SIZE is the number of bytes that may become unavailable  
1137 due to an error on an NVM device.

1138 A zero value means that the device is unable to provide a guarantee on the number of  
1139 adjacent bytes impacted by an error.

1140 This attribute is relevant when the device does not support write atomicity.

1141 If FUNDAMENTAL\_BLOCK\_SIZE is smaller than LOGICAL\_BLOCK\_SIZE (see  
1142 7.3.14), an application may organize data in terms of FUNDAMENTAL\_BLOCK\_SIZE to  
1143 avoid certain torn write behavior. If FUNDAMENTAL\_BLOCK\_SIZE is larger than  
1144 LOGICAL\_BLOCK\_SIZE, an application may organize data in terms of  
1145 FUNDAMENTAL\_BLOCK\_SIZE to assure two key data items do not occupy an extent  
1146 that is vulnerable to errors.

## 1147 **7.4 Use cases**

### 1148 7.4.1 **Flash as cache use case**

#### 1149 **Purpose/triggers:**

1150 Use Flash based NVM as a data cache.

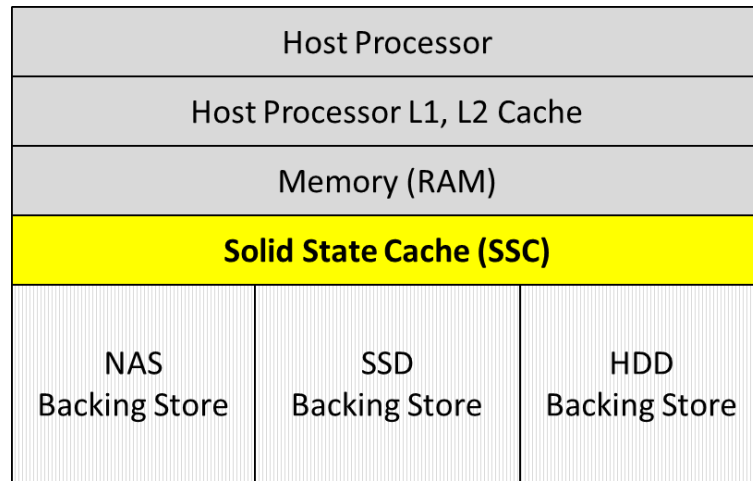
#### 1151 **Scope/context:**

1152 Flash memory's fast random I/O performance and non-volatile characteristic make it a  
1153 good candidate as a Solid State Cache device (SSC). This use case is described in  
1154 Figure 7 SSC in a storage stack.

1155

1156

Figure 7 SSC in a storage stack



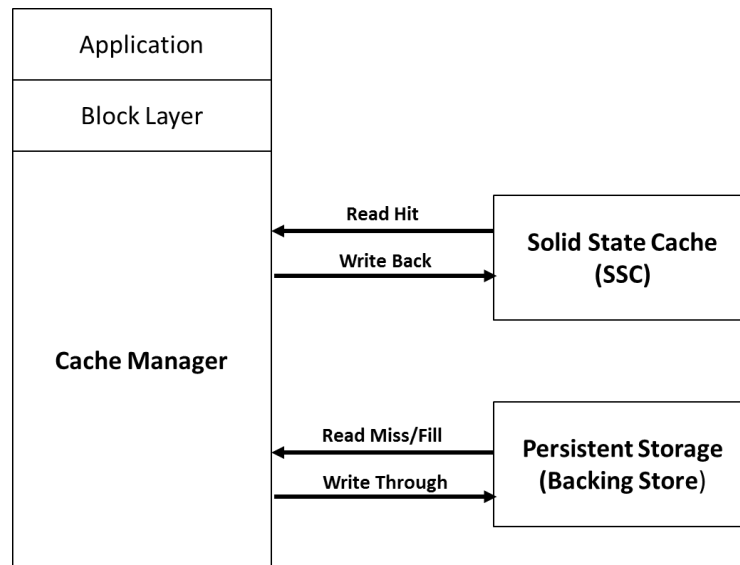
1157

1158

1159 A possible software application is shown in Figure 8 SSC software cache application. In  
1160 this case, the cache manager employs the Solid State Cache to improve caching  
1161 performance and to maintain persistence and cache coherency across power fail.

1162

Figure 8 SSC software cache application

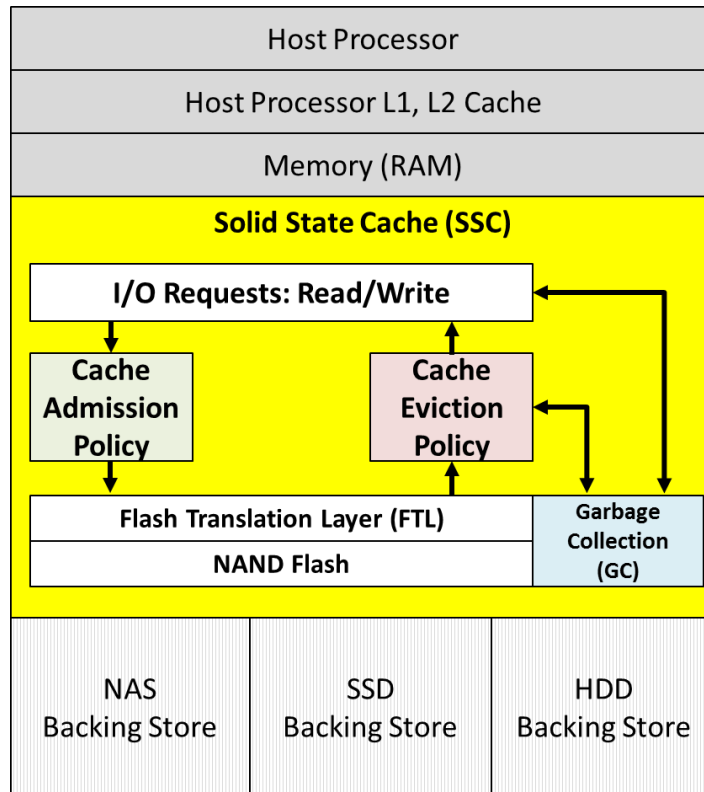


1163

1164 It is also possible to use an enhanced SSC to perform some of the functions that the  
1165 cache manager must normally contend with as shown in Figure 9 SSC with caching  
1166 assistance.

1167

Figure 9 SSC with caching assistance



1168

1169 In this use case, the Solid State Cache (SSC) provides a sparse address space that  
 1170 may be much larger than the amount of physical NVM memory and manages the cache  
 1171 through its own admission and eviction policies. The backing store is used to persist the  
 1172 data when the cache becomes full. As a result, the block state for each block of virtual  
 1173 storage in the cache must be maintained by the SSC. The SSC must also present a  
 1174 consistent cache interface that can persist the cached data across a power fail and  
 1175 never returns stale data.

1176 In either of these cases, two important extensions to existing storage commands must  
 1177 be present:

1178 **Eviction:** An explicit eviction mechanism is required to invalidate cached data in  
 1179 the SSC to allow the cache manager to precisely control the contents of the SSC.  
 1180 This means that the SSC must insure that the eviction is durable before  
 1181 completing the request. This mechanism is generally referred to as a persistent  
 1182 trim. This is the NVM.BLOCK.DISCARD\_IMMEDIATELY functionality.

1183 **Exists:** The *exists* operation allows the cache manager to determine the state of  
 1184 a block, or of a range of blocks, in the SSC. This operation is used to test for the  
 1185 presence of data in the cache, or to determine which blocks in the SSC are dirty  
 1186 and need to be flushed to backing storage. This is the NVM.BLOCK.EXISTS  
 1187 functionality.

1188 The most efficient mechanism for a cache manager would be to simply read the  
1189 requested data from the SSC which would the return either the data or an error  
1190 indicated that the requested data was not in the cache. This approach is problematic,  
1191 since most storage drivers and software require reads to be successful and complete by  
1192 returning data - not an error. Device that return errors for normal read operations are  
1193 usually put into an offline state by the system drivers. Further, the data that a read  
1194 returns must be consistent from one read operation to the next, provided that no  
1195 intervening writes occur. As a result, a two stage process must be used by the cache  
1196 manager. The cache manager first issues an *exists* command to determine if the  
1197 requested data is present in the cache. Based on the result, the cache manager decides  
1198 whether to read the data from the SSC or from the backing storage.

1199 **Preconditions:**

1200 N/A

1201 **Success scenario:**

1202 The requested data is successfully read from or written to the SSC.

1203 **See also:**

- 1204 • NVM.BLOCK.DISCARD\_IMMEDIATELY
- 1205 • NVM.BLOCK.EXISTS
- 1206 • Ptrim() + Exists(): Exposing New FTL Primitives to Applications, David Nellans,  
1207 Michael Zappe, Jens Axboe, David Flynn, 2011 Non-Volatile Memory Workshop.  
1208 See: <http://david.nellans.org/files/NVMW-2011.pdf>
- 1209 • FlashTier: a Lightweight, Consistent, and Durable Storage Cache, Mohit Saxena,  
1210 Michael M. Swift and Yiying Zhang, University of Wisconsin-Madison. See:  
1211 [http://pages.cs.wisc.edu/~swift/papers/eurosys12\\_flashtier.pdf](http://pages.cs.wisc.edu/~swift/papers/eurosys12_flashtier.pdf)  
1212 HEC: Improving Endurance of High Performance Flash-based Cache Devices,  
1213 Jingpei Yang, Ned Plasson, Greg Gillis, Nisha Talagala, Swaminathan  
1214 Sundararaman, Robert Wood, Fusion-io, Inc., SYSTOR '13, June 30 - July 02  
1215 2013, Haifa, Israel
- 1216 • Unioning of the Buffer Cache and Journaling Layers with Non-volatile Memory,  
1217 Eunji Lee, Hyokyung Bahn, and Sam H. Noh. See:  
1218 [https://www.usenix.org/system/files/conference/fast13/fast13-final114\\_0.pdf](https://www.usenix.org/system/files/conference/fast13/fast13-final114_0.pdf)

1219 **7.4.2 SCAR use case**

1220 **Purpose/triggers:**

1221 Demonstrate the use of the SCAR action

1222 **Scope/context:**

1223 This generic use case for SCAR involves two processes.

- 1224 • The “detect block errors process” detects errors in certain NVM blocks, and uses  
1225 SCAR to communicate to other processes that the contents of these blocks cannot  
1226 be reliably read, but can be safely re-written.

1227 • The “recover process” sees the error reported as the result of SCAR. If this process  
1228 can regenerate the contents of the block, the application can continue with no error.

1229 For this use case, the “detect block errors process” is a RAID component doing a  
1230 background scan of NVM blocks. In this case, the NVM is not in a redundant RAID  
1231 configuration so block READ errors can’t be transparently recovered. The “recover  
1232 process” is a cache component using the NVM as a cache for RAID volumes. Upon  
1233 receipt of the SCAR error on a read, this component evaluates whether the block  
1234 contents also reside on the cached volume; if so, it can copy the corresponding volume  
1235 block to the NVM. This write to NVM will clear the SCAR error condition.

1236 **Preconditions:**

1237 The “detect block errors process” detected errors in certain NVM blocks, and used  
1238 SCAR to mark these blocks.

1239 **Inputs:**

1240 None

1241 **Success scenario:**

- 1242 1. The cache manager intercepts a read request from an application
- 1243 2. The read request to the NVM cache returns a status indicating the requested  
1244 blocks have been marked by a SCAR action
- 1245 3. The cache manager uses an implementation-specific technique and determines  
1246 the blocks marked by a SCAR are also available on the cached volume
- 1247 4. The cache manager copies the blocks from the cached volume to the NVM
- 1248 5. The cache manager returns the requested block to the application with a status  
1249 indicating the read succeeded

1250 **Postconditions:**

1251 The blocks previously marked with a SCAR action have been repaired.

1252 **Failure Scenario:**

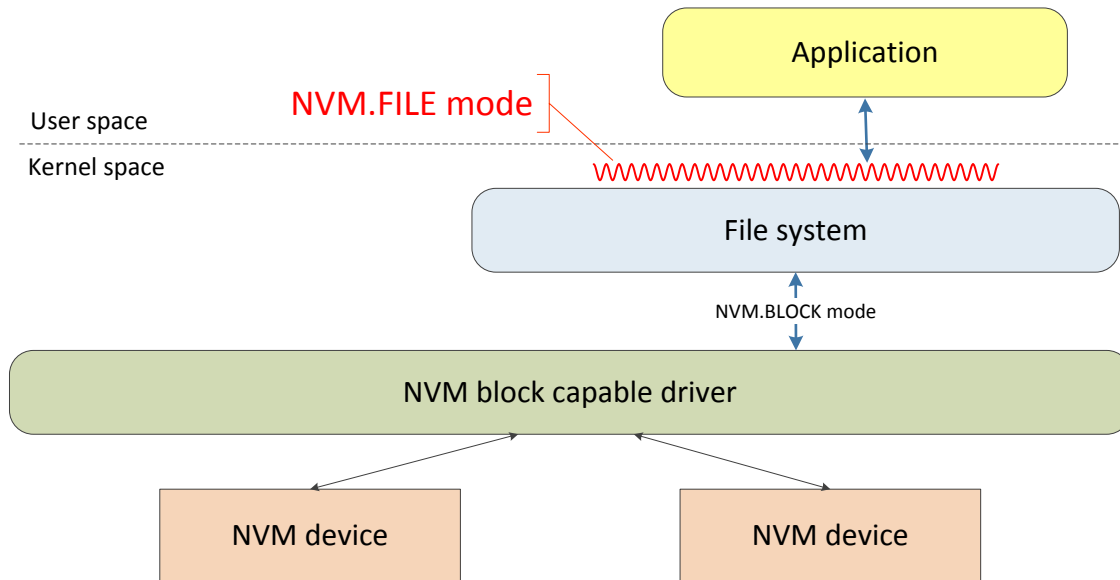
- 1253 1. In Success Scenario step 3 or 4, the cache manager discovers the  
1254 corresponding blocks on the volume are invalid or cannot be read.
- 1255 2. The cache manager returns a status to the application indicating the blocks  
1256 cannot be read.

## 1257 8 NVM.FILE mode

### 1258 8.1 Overview

1259 NVM.FILE mode addresses NVM-specification extensions to native file I/O behavior  
1260 (the approach to I/O used by most applications). Support for NVM.FILE mode requires  
1261 that the NVM solution ought to support all behavior not covered in this section  
1262 consistently with the native operating system behavior for native block devices.

1263 **Figure 10 NVM.FILE mode example**



1264

#### 1265 8.1.1 Discovery and use of atomic write features

1266 Atomic Write features in NVM.FILE mode are available to block-optimized applications  
1267 (see 7.1.4.2 Block-optimized applications).

#### 1268 8.1.2 The discovery of granularities

1269 The NVM.FILE mode exposes the same granularity attributes as NVM.BLOCK.

#### 1270 8.1.3 Relationship between native file APIs and NVM.BLOCK.DISCARD

1271 NVM.FILE mode does not define specific action that cause TRIM/DISCARD behavior.  
1272 File systems may invoke NVM.BLOCK DISCARD actions when native operating system  
1273 APIs (such as POSIX truncate or Windows SetEndOfFile).

### 1274 8.2 Actions

#### 1275 8.2.1 Actions that apply across multiple modes

1276 The following actions apply to NVM.FILE mode as well as other modes.

1277 NVM.COMMON.GET\_ATTRIBUTE (see 6.10.1)

1278 NVM.COMMON.SET\_ATTRIBUTE (see 6.10.2)

1279 **8.2.2 NVM.FILE.ATOMIC\_WRITE**

1280 Requirement: mandatory if ATOMIC\_WRITE\_CAPABLE (see 8.3.2) is true

1281 Block-optimized applications may use ATOMIC\_WRITE to assure consistent behavior  
1282 during a failure condition. This specification does not specify the order in which this  
1283 action occurs relative to other I/O operations, including other ATOMIC\_WRITE and  
1284 ATOMIC\_MULTIWRITE actions. This specification does not specify when the data  
1285 written becomes visible to other threads.

1286 The inputs, outputs, and error conditions are similar to those for  
1287 NVM.BLOCK.ATOMIC\_WRITE, but typically the application provides file names and file  
1288 relative block addresses rather than device name and LBA.

1289 Relevant attributes:

1290        ATOMIC\_WRITE\_MAX\_DATA\_LENGTH  
1291        ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY  
1292        ATOMIC\_WRITE\_LENGTH\_GRANULARITY  
1293        ATOMIC\_WRITE\_CAPABLE

1294 **8.2.3 NVM.FILE.ATOMIC\_MULTIWRITE**

1295 Requirement: mandatory if ATOMIC\_MULTIWRITE\_CAPABLE (see 8.3.6) is true

1296 Block-optimized applications may use ATOMIC\_MULTIWRITE to assure consistent  
1297 behavior during a failure condition. This action allows a caller to write non-adjacent  
1298 extents atomically. The caller of ATOMIC\_MULTIWRITE provides properties defining  
1299 memory and block extents; all of the extents are written as a single atomic operation.  
1300 This specification does not specify the order in which this action occurs relative to other  
1301 I/O operations, including other ATOMIC\_WRITE and ATOMIC\_MULTIWRITE actions.  
1302 This specification does not specify when the data written becomes visible to other  
1303 threads.

1304 The inputs, outputs, and error conditions are similar to those for  
1305 NVM.BLOCK.ATOMIC\_MULTIWRITE, but typically the application provides file names  
1306 and file relative block addresses rather than device name and LBA.

1307 Relevant attributes:

1308        ATOMIC\_MULTIWRITE\_MAX\_IOS  
1309        ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH  
1310        ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY  
1311        ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY  
1312        ATOMIC\_MULTIWRITE\_CAPABLE

1313 **8.3 Attributes**

1314 Some attributes share behavior with their NVM.BLOCK counterparts. NVM.FILE  
1315 attributes are provided because the actual values may change due to the  
1316 implementation of the file system.

1317 **8.3.1 Attributes that apply across multiple modes**

1318 The following attributes apply to NVM.FILE mode as well as other modes.

1319 NVM.COMMON.SUPPORTED\_MODES (see 6.11.1)

1320 NVM.COMMON.FILE\_MODE (see 6.11.2)

1321 **8.3.2 NVM.FILE.ATOMIC\_WRITE\_CAPABLE**

1322 Requirement: mandatory

1323 This attribute indicates that the implementation is capable of the  
1324 NVM.BLOCK.ATOMIC\_WRITE action.

1325 **8.3.3 NVM.FILE.ATOMIC\_WRITE\_MAX\_DATA\_LENGTH**

1326 Requirement: mandatory

1327 ATOMIC\_WRITE\_MAX\_DATA\_LENGTH is the maximum length of data that can be  
1328 transferred by an ATOMIC\_WRITE action.

1329 **8.3.4 NVM.FILE.ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY**

1330 Requirement: mandatory

1331 ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY is the granularity of the  
1332 starting memory address for an ATOMIC\_WRITE action. Address inputs to  
1333 ATOMIC\_WRITE shall be evenly divisible by  
1334 ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY.

1335 **8.3.5 NVM.FILE.ATOMIC\_WRITE\_LENGTH\_GRANULARITY**

1336 Requirement: mandatory

1337 ATOMIC\_WRITE\_LENGTH\_GRANULARITY is the granularity of the length of data  
1338 transferred by an ATOMIC\_WRITE action. Length inputs to ATOMIC\_WRITE shall be  
1339 evenly divisible by ATOMIC\_WRITE\_STARTING\_ADDRESS\_GRANULARITY.

1340 **8.3.6 NVM.FILE.ATOMIC\_MULTIWRITE\_CAPABLE**

1341 Requirement: mandatory

1342 This attribute indicates that the implementation is capable of the  
1343 NVM.FILE.ATOMIC\_MULTIWRITE action.



- 1344 **8.3.7 NVM.FILE.ATOMIC\_MULTIWRITE\_MAX\_IOS**
- 1345 Requirement: mandatory
- 1346 ATOMIC\_MULTIWRITE\_MAX\_IOS is the maximum length of the number of IOs (i.e.,  
1347 the size of the Property Group List) that can be transferred by an  
1348 ATOMIC\_MULTIWRITE action.
- 1349 **8.3.8 NVM.FILE.ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH**
- 1350 Requirement: mandatory
- 1351 ATOMIC\_MULTIWRITE\_MAX\_DATA\_LENGTH is the maximum length of data that can  
1352 be transferred by an ATOMIC\_MULTIWRITE action.
- 1353 **8.3.9 NVM.FILE.ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY**
- 1354 Requirement: mandatory
- 1355 ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY is the granularity of  
1356 the starting address of ATOMIC\_MULTIWRITE inputs. Address inputs to  
1357 ATOMIC\_MULTIWRITE shall be evenly divisible by  
1358 ATOMIC\_MULTIWRITE\_STARTING\_ADDRESS\_GRANULARITY.
- 1359 **8.3.10 NVM.FILE.ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY**
- 1360 Requirement: mandatory
- 1361 ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY is the granularity of the length of  
1362 ATOMIC\_MULTIWRITE inputs. Length inputs to ATOMIC\_MULTIWRITE shall be  
1363 evenly divisible by ATOMIC\_MULTIWRITE\_LENGTH\_GRANULARITY.
- 1364 **8.3.11 NVM.FILE.WRITE\_ATOMICITY\_UNIT**
- 1365 See 7.3.11 NVM.BLOCK.WRITE\_ATOMICITY\_UNIT
- 1366 **8.3.12 NVM.FILE.LOGICAL\_BLOCK\_SIZE**
- 1367 See 7.3.14 NVM.BLOCK.LOGICAL\_BLOCK\_SIZE
- 1368 **8.3.13 NVM.FILE.PERFORMANCE\_BLOCK\_SIZE**
- 1369 See 7.3.15 NVM.BLOCK.PERFORMANCE\_BLOCK\_SIZE
- 1370 **8.3.14 NVM.FILE.LOGICAL\_ALLOCATION\_SIZE**
- 1371 See 7.3.16 NVM.BLOCK.ALLOCATION\_BLOCK\_SIZE
- 1372 **8.3.15 NVM.FILE.FUNDAMENTAL\_BLOCK\_SIZE**
- 1373 See 7.3.20 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE

1374 **8.4 Use cases**

1375 **8.4.1 Block-optimized application updates record**

1376 Update a record in a file without using a memory-mapped file

1377 **Purpose/triggers:**

1378 An application using block NVM updates an existing record. The application requests  
1379 that the file system bypass cache; the application conforms to native API requirements  
1380 when bypassing cache – this may mean that read and write actions must use multiples  
1381 of a page cache size. For simplicity, this application uses fixed size records. The record  
1382 size is defined by application data considerations, not disk or page block sizes. The  
1383 application factors in the PERFORMANCE\_BLOCK\_SIZE granularity to avoid device-  
1384 side inefficiencies such as read/modify/write.

1385 **Scope/context:**

1386 Block NVM context; this shows basic behavior.

1387 **Preconditions:**

- 1388 - The administrator created a file and provided its name to the application; this name is
- 1389 accessible to the application – perhaps in a configuration file
- 1390 - The application has populated the contents of this file
- 1391 - The file is not in use at the start of this use case (no sharing considerations)

1392 **Inputs:**

1393 The content of the record, the location (relative to the file) where the record resides

1394 **Success scenario:**

- 1395 1) The application uses the native OPEN action, passing in the file name and  
1396 specifying appropriate options to bypass the file system cache
- 1397 2) The application acquires the device's optimal I/O granule size by using the  
1398 GET\_ATTRIBUTE action for the PERFORMANCE\_BLOCK\_SIZE.
- 1399 3) The application allocates sufficient memory to contain all of the blocks occupied by  
1400 the record to be updated.
  - 1401 a. The application determines the offset within the starting block of the record  
1402 and uses the length of the block to determine the number of partial blocks.
  - 1403 b. The application allocates sufficient memory for the record plus enough  
1404 additional memory to accommodate any partial blocks.
  - 1405 c. If necessary, the memory size is increased to assure that the starting address  
1406 and length read and write actions are multiples of  
1407 PERFORMANCE\_BLOCK\_SIZE.
- 1408 4) The application uses the native READ action to read the record by specifying the  
1409 starting disk address and the length (the same length as the allocated memory)

- 1410 buffer). The application also provides the allocated memory address; this is where  
1411 the read action will put the record.
- 1412 5) The application updates the record in the memory buffer per the inputs
  - 1413 6) The application uses the native write action to write the updated block(s) to the same  
1414 disk location they were read from.
  - 1415 7) The application uses the native file SYNC action to assure the updated blocks are  
1416 written to the persistence domain
  - 1417 8) The application uses the native CLOSE action to clean up.

1418 **Failure Scenario 1:**

1419 The native read action reports a hardware error. If the unreadable block corresponds to  
1420 blocks being updated, the application may attempt recovery (write/read/verify), or  
1421 preventative maintenance (scar the unreadable blocks). If the unreadable blocks are  
1422 needed for a read/modify/write update and the application lacks an alternate source; the  
1423 application may inform the user that an unrecoverable hardware error has occurred.

1424 **Failure Scenario 2:**

1425 The native write action reports a hardware error. The application may be able to recover  
1426 by rewriting the block. If the rewrite fails, the application may be able to scar the bad  
1427 block and write to a different location.

1428 **Outputs:**

1429 None

1430 **Postconditions:**

1431 The record is updated.

1432 **8.4.2 Atomic write use case**

1433 **Purpose/triggers:**

1434 Used by a block-optimized application (see Block-optimized applications) striving for  
1435 durability of on-disk data

1436 **Scope/context:**

1437 Assure a record is written to disk in a way that torn writes can be detected and rolled  
1438 back (if necessary). If the device supports atomic writes, they will be used. If not, a  
1439 double write buffer is used.

1440 **Preconditions:**

1441 The application has taken steps (based on NVM.BLOCK attributes) to assure the record  
1442 being written has an optimal memory starting address, starting disk LBA and length.

1443 **Inputs:**

1444 None

- 1445 **Success scenario:**
- 1446 • Use GET\_ATTRIBUTE to determine whether the device is
  - 1447 ATOMIC\_WRITE\_CAPABLE (or ATOMIC\_MULTIWRITE\_CAPABLE)
  - 1448 • Is so, use the appropriate atomic write action to write the record to NVM
  - 1449 • If the device does not support atomic write, then
  - 1450 ○ Write the page to the double write buffer
  - 1451 ○ Wait for the write to complete
  - 1452 ○ Write the page to the final destination
  - 1453 • At application startup, if the device does not support atomic write
  - 1454 • Scan the double write buffer and for each valid page in the buffer check if the
  - 1455 page in the data file is valid too.

1456 **Outputs:**

1457 None

1458 **Postconditions:**

1459 After application startup recovery steps, there are no inconsistent records on disk after a  
1460 failure caused the application (and possibly system) to restart.

1461 **8.4.3 Block and File Transaction Logging**

1462 **Purpose/Triggers:**

1463 An application developer wishes to implement a transaction log that maintains data  
1464 integrity through system crashes, system resets, and power failures. The underlying  
1465 storage is block-granular, although it may be accessed via a file system that simulates  
1466 byte-granular access to files.

1467 **Scope/Context:**

1468 NVM.BLOCK or NVM.FILE (all the NVM.BLOCK attributes mentioned in the use case  
1469 are also defined for NVM.FILE mode).

1470 For notational convenience, this use case will use the term “file” to apply to either a file  
1471 in the conventional sense which is accessed through the NVM.FILE interface, or a  
1472 specific subset of blocks residing on a block device which are accessed through the  
1473 NVM.BLOCK interface.

1474 **Inputs:**

- 1475 • A set of changes to the persistent state to be applied as a single transaction.
- 1476 • The data and log files.

1477 **Outputs:**

- 1478 • An indication of transaction commit or abort.

1479 **Postconditions:**  
1480     • If an abort indication was returned, the data was not committed and the previous  
1481     contents have not been modified.  
1482     • If a commit indication was returned, the data has been entirely committed.  
1483     • After a system crash, reset, or power failure followed by system restart and  
1484     execution of the application transaction recovery process, the data has either  
1485     been entirely committed or the previous contents have not been modified.

1486 **Success Scenario:**

1487 The application transaction logic uses a log file in combination with its data file to  
1488 atomically update the persistent state of the application. The log may implement a  
1489 before-image log or a write-ahead log. The application transaction logic should  
1490 configure itself to handle torn or interrupted writes to the log or data files.

1491 **8.4.3.1 NVM.BLOCK.WRITE\_ATOMICITY\_UNIT >= 1**

1492 If the NVM.BLOCK.WRITE\_ATOMICITY\_UNIT is one or greater, then writes of a single  
1493 logical block cannot be torn or interrupted.

1494 In this case, if the log or data record size is less than or equal to the  
1495 NVM.BLOCK.LOGICAL\_BLOCK\_SIZE, the application need not handle torn or  
1496 interrupted writes to the log or data files.

1497 If the log or data record size is greater than the NVM.BLOCK.LOGICAL\_BLOCK\_SIZE,  
1498 the application should be prepared to detect a torn write of the record and either discard  
1499 or recover such a torn record during the recovery process. One common way of  
1500 detecting such a torn write is for the application to compute hash of the record and  
1501 record the hash in the record. Upon reading the record, the application re-computes the  
1502 hash and compares it with the recorded hash; if they do not match, the record has been  
1503 torn. Another method is for the application to insert the transaction identifier within each  
1504 logical block. Upon reading the record, the application compares the transaction  
1505 identifiers in each logical block; if they do not match, the record has been torn. Another  
1506 method is for the application to use the NVM.BLOCK.ATOMIC\_WRITE action to  
1507 perform the writes of the record.

1508 **8.4.3.2 NVM.BLOCK.WRITE\_ATOMICITY\_UNIT = 0**

1509 If the NVM.BLOCK.WRITE\_ATOMICITY\_UNIT is zero, then writes of a single logical  
1510 block can be torn or interrupted and the application should handle torn or interrupted  
1511 writes to the log or data files.

1512 In this case, if a logical block were to contain data from more than one log or data  
1513 record, a torn or interrupted write could corrupt a previously-written record. To prevent  
1514 propagating an error beyond the record currently being written, the application aligns  
1515 the log or data records with the NVM.BLOCK. LOGICAL\_BLOCK\_SIZE and pads the  
1516 record size to be an integral multiple of NVM.BLOCK. LOGICAL\_BLOCK\_SIZE. This

1517 prevents more than one record from residing in the same logical block and therefore a  
1518 torn or interrupted write may only corrupt the record being written.

1519 *8.4.3.2.1 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE >=*  
1520 *NVM.BLOCK.LOGICAL\_BLOCK\_SIZE*

1521 If the NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE is greater than or equal to the  
1522 NVM.BLOCK.LOGICAL\_BLOCK\_SIZE, the application should be prepared to handle an  
1523 interrupted write. An interrupted write results when the write of a single  
1524 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE unit is interrupted by a system crash,  
1525 system reset, or power failure. As a result of an interrupted write, the NVM device may  
1526 return an error when any of the logical blocks comprising the  
1527 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE unit are read. (See also SQLite.org,  
1528 *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>.) This presents a danger to the  
1529 integrity of previously written records that, while residing in differing logical blocks, share  
1530 the same fundamental block. An interrupted write may prevent the reading of those  
1531 previously written records in addition to preventing the read of the record in the process  
1532 of being written.

1533 One common way of protecting previously written records from damage due to an  
1534 interrupted write is to align the log or data records with the  
1535 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE and pad the record size to be an integral  
1536 multiple of NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE. This prevents more than one  
1537 record from residing in the same fundamental block. The application should be prepared  
1538 to discard or recover the record if the NVM device returns an error when subsequently  
1539 reading the record during the recovery process.

1540 *8.4.3.2.2 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE <*  
1541 *NVM.BLOCK.LOGICAL\_BLOCK\_SIZE*

1542 If the NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE is less than the  
1543 NVM.BLOCK.LOGICAL\_BLOCK\_SIZE, the application should be prepared to handle  
1544 both interrupted writes and torn writes within a logical block.

1545 As a result of an interrupted write, the NVM device may return an error when the logical  
1546 block containing the NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE unit which was  
1547 being written at the time of the system crash, system reset, or power failure is  
1548 subsequently read. The application should be prepared to discard or recover the record  
1549 in the logical block if the NVM device returns an error when subsequently reading the  
1550 logical block during the recovery process.

1551 A torn write results when an integral number of  
1552 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE units are written to the NVM device but  
1553 the entire NVM.BLOCK.LOGICAL\_BLOCK\_SIZE has not been written. In this case, the  
1554 NVM device may not return an error when the logical block is read. The application  
1555 should therefore be prepared to detect a torn write of a logical block and either discard  
1556 or recover such a torn record during the recovery process. One common way of  
1557 detecting such a torn write is for the application to compute a hash of the record and

1558 record the hash in the record. Upon reading the record, the application re-computes the  
1559 hash and compares it with the recorded hash; if they do not match, a logical block within  
1560 the record has been torn. Another method is for the application to insert the transaction  
1561 identifier within each NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE unit. Upon reading  
1562 the record, the application compares the transaction identifiers in each  
1563 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE unit; if they do not match, the logical  
1564 block has been torn.

#### 1565 8.4.3.2.3 NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE = 0

1566 If the NVM.BLOCK.FUNDAMENTAL\_BLOCK\_SIZE is zero, the application lacks  
1567 sufficient information to handle torn or interrupted writes to the log or data files.

#### 1568 **Failure Scenarios:**

1569 Consider the recovery of an error resulting from an interrupted write on a device where  
1570 the NVM.BLOCK.WRITE\_ATOMICITY\_UNIT is zero. This error may be persistent and  
1571 may be returned whenever the affected fundamental block is read. To repair this error,  
1572 the application should be prepared to overwrite such a block.

1573 One common way of ensuring that the application will overwrite a block is by assigning  
1574 it to the set of internal free space managed by the application, which is never read and  
1575 is available to be allocated and overwritten at some point in the future. For example, the  
1576 block may be part of a circular log. If the block is marked as free, the transaction log  
1577 logic will eventually allocate and overwrite that block as records are written to the log.

1578 Another common way is to record either a before-image or after-image of a data block  
1579 in a log. During recovery after a system crash, system reset, or power failure, the  
1580 application replays the records in the log and overwrites the data block with either the  
1581 before-image contents or the after-image contents.

#### 1582 **See also:**

- 1583 • SQLite.org, *Atomic Commit in SQLite*, <http://www.sqlite.org/atomiccommit.html>
- 1584 • SQLite.org, *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>
- 1585 • SQLite.org, *Write-Ahead Logging*, <http://www.sqlite.org/wal.html>

1586 **9 NVM.PM.VOLUME mode**

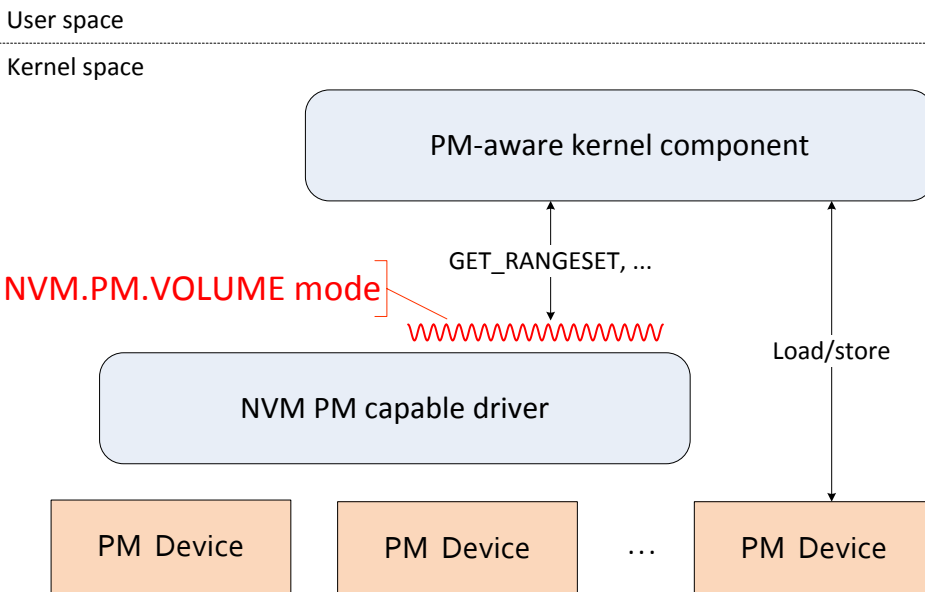
1587 **9.1 Overview**

1588 NVM.PM.VOLUME mode describes the behavior to be consumed by operating system  
1589 abstractions such as file systems or pseudo-block devices that build their functionality  
1590 by directly accessing persistent memory. NVM.PM.VOLUME mode provides a software  
1591 abstraction (a PM volume) for persistent memory hardware and profiles functionality for  
1592 operating system components including:

- 1593 • list of physical address ranges associated with each PM volume,  
1594 • capability to determine whether PM errors have been reported

1595 The PM volume provides memory mapped capability in a fashion that traditional CPU  
1596 load and store operations are possible. This PM volume may be provided via the  
1597 memory channel of the CPU or via a PCIe memory mapping or other methods. Note  
1598 that there should not be a requirement for an operating system context switch for  
1599 access to the PM volume.

1600 **Figure 11 NVM.PM.VOLUME mode example**



1601

1602 **9.2 Actions**

1603 **9.2.1 Actions that apply across multiple modes**

1604 The following actions apply to NVM.PM.VOLUME mode as well as other modes.

1605 NVM.COMMON.GET\_ATTRIBUTE (see 6.10.1)

1606 NVM.COMMON.SET\_ATTRIBUTE (see 6.10.2)



1607 **9.2.2 NVM.PM.VOLUME.GET\_RANGESET**

1608 Requirement: mandatory

1609 The purpose of this action is to return a set of processor physical address ranges (and  
1610 relate properties) representing all of the content for the identified volume.

1611 When interpreting the set of physical addresses as a contiguous, logical address range;  
1612 the data underlying that logical address range will always be the same and in the same  
1613 sequence across PM volume instantiations.

1614 Due to physical memory reconfiguration, the number and sizes of ranges may change in  
1615 successive get ranges calls, however the total number of bytes in the sum of the ranges  
1616 does not change, and the order of the bytes spanning all of the ranges does not  
1617 change. The space defined by the list of ranges can always be addressed relative to a  
1618 single base which represents the beginning of the first range.

1619 Input: a reference to the PM volume

1620 Returns a Property Group List (see 4.4.5) where the properties are:

- 1621 • starting physical address (byte address)
- 1622 • length (in bytes)
- 1623 • connection type – see below
- 1624 • sync type – see below

1625 For this revision of the specification, the following values (in text) are valid for  
1626 connection type:

- 1627 • “*memory*”: for persistent memory attached to a system memory channel
- 1628 • “*PCIe*”: for persistent memory attached to a PCIe extension bus

1629 For this revision of the specification, the following values (in text) are valid for sync type:

- 1630 • “none”: no device-specific sync behavior is available – implies no entry to  
1631 NVM.PM.VOLUME implementation is required for flushing
- 1632 • “VIRTUAL\_ADDRESS\_SYNC”: the caller needs to use VIRTUAL\_ADDRESS\_SYNC  
1633 (see 9.2.3) to assure sync is durable
- 1634 • “PHYSICAL\_ADDRESS\_SYNC”: the caller needs to use  
1635 PHYSICAL\_ADDRESS\_SYNC (see 9.2.4) to assure sync is durable

1636 **9.2.3 NVM.PM.VOLUME.VIRTUAL\_ADDRESS\_SYNC**

1637 Requirement: optional

1638 The purpose of this action is to invoke device-specific actions to synchronize persistent  
1639 memory content to assure durability and enable recovery by forcing data to reach the  
1640 persistence domain. VIRTUAL\_ADDRESS\_SYNC is used by a caller that knows the  
1641 addresses in the input range are virtual memory addresses.

1642 Input: virtual address and length (range)

1643 See also: PHYSICAL\_ADDRESS\_SYNC

#### 1644 9.2.4 **NVM.PM.VOLUME.PHYSICAL\_ADDRESS\_SYNC**

1645 Requirement: optional

1646 The purpose of this action is to synchronize persistent memory content to assure  
1647 durability and enable recovery by forcing data to reach the persistence domain. This  
1648 action is used by a caller that knows the addresses in the input range are physical  
1649 memory addresses.

1650 See also: VIRTUAL\_ADDRESS\_SYNC

1651 Input: physical address and length (range)

#### 1652 9.2.5 **NVM.PM.VOLUME.DISCARD\_IF\_YOU\_CAN**

1653 Requirement: mandatory if DISCARD\_IF\_YOU\_CAN\_CAPABLE (see 9.3.6) is true

1654 This action notifies the NVM device that the input range (volume offset and length) are  
1655 no longer needed by the caller. This action may not result in any action by the device,  
1656 depending on the implementation and the internal state of the device. This action is  
1657 meant to allow the underlying device to optimize the data stored within the range. For  
1658 example, the device can use this information in support of functionality like thin  
1659 provisioning or wear-leveling.

#### 1660 9.2.6 **NVM.PM.VOLUME.DISCARD\_IMMEDIATELY**

1661 Requirement: mandatory if DISCARD\_IMMEDIATELY\_CAPABLE (see 9.3.7) is true

1662 This action notifies the NVM device that the input range (volume offset and length) are  
1663 no longer needed by the caller. Similar to DISCARD\_IF\_YOU\_CAN, but the  
1664 implementation is required to unmap the range before the next READ or WRITE action,  
1665 even if garbage collection of the range has not occurred yet.

#### 1666 9.2.7 **NVM.PM.VOLUME.EXISTS**

1667 Requirement: mandatory if EXISTS\_CAPABLE (see 7.3.12) is true

1668 A PM device may allocate storage through a thin provisioning mechanism or one of the  
1669 discard actions. As a result, memory can exist in one of three states:

- 1670 • **Mapped**: the range has had data written to it
- 1671 • **Unmapped**: the range has not been written, and there is no memory allocated
- 1672 • **Allocated**: the range has not been written, but has memory allocated to it

1673 The EXISTS action allows the NVM user to determine if a range of bytes has been  
1674 allocated.

1675 Inputs: a range of bytes (starting byte address and length in bytes)  
1676 Output: a Property Group List (see 4.4.5) where the properties are the starting address,  
1677 length and state. State is a string equal to “mapped”, “unmapped”, or “allocated”.  
1678 Result: the status of the action

## 1679 **9.3 Attributes**

### 1680 **9.3.1 Attributes that apply across multiple modes**

1681 The following attributes apply to NVM.PM.VOLUME mode as well as other modes.  
1682 NVM.COMMON.SUPPORTED\_MODES (see 6.11.1)

1683

### 1684 **9.3.2 NVM.PM.VOLUME.VOLUME\_SIZE**

1685 Requirement: mandatory

1686 VOLUME\_SIZE is the volume size in units of bytes. This shall be the same as the sum  
1687 of the lengths of the ranges returned by the GET\_RANGES action.

### 1688 **9.3.3 NVM.PM.VOLUME.INTERRUPTED\_STORE\_ATOMICALITY**

1689 Requirement: mandatory

1690 INTERRUPTED\_STORE\_ATOMICALITY indicates whether the device supports power fail  
1691 atomicity of store actions.

1692 A value of true indicates that after a store interrupted by reset, power loss or system  
1693 crash; upon restart the contents of persistent memory reflect either the state before the  
1694 store or the state after the completed store. A value of false indicates that after a store  
1695 interrupted by reset, power loss or system crash, upon restart the contents of memory  
1696 may be such that subsequent loads may create exceptions depending on the value of  
1697 the FUNDAMENTAL\_ERROR\_RANGE attribute (see 9.3.4).

### 1698 **9.3.4 NVM.PM.VOLUME.FUNDAMENTAL\_ERROR\_RANGE**

1699 Requirement: mandatory

1700 FUNDAMENTAL\_ERROR\_RANGE is the number of bytes that may become  
1701 unavailable due to an error on an NVM device.

1702 This attribute is relevant when the device does not support write atomicity.

1703 A zero value means that the device is unable to provide a guarantee on the number of  
1704 adjacent bytes impacted by an error.

1705 A caller may organize data in terms of FUNDAMENTAL\_ERROR\_RANGE to avoid  
1706 certain torn write behavior.

1707 **9.3.5 NVM.PM.VOLUME.FUNDAMENTAL\_ERROR\_RANGE\_OFFSET**

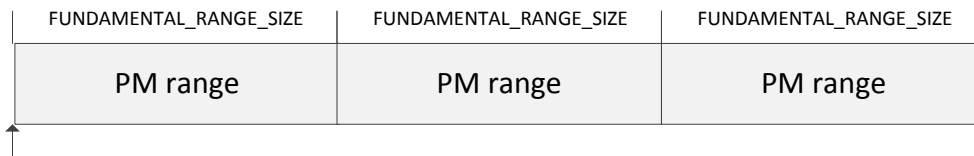
1708 Requirement: mandatory

1709 The number of bytes offset from the beginning of a volume range (as returned by  
1710 GET\_RANGESET) before FUNDAMENTAL\_ERROR\_RANGE\_SIZE intervals apply.

1711 A fundamental error range is not required to start at a byte address evenly divisible by  
1712 FUNDAMENTAL\_ERROR\_RANGE. FUNDAMENTAL\_ERROR\_RANGE\_OFFSET shall  
1713 be set to the difference between the starting byte address of a fundamental error range  
1714 rounded down to a multiple of FUNDAMENTAL\_ERROR\_RANGE.

1715 Figure 12 Zero range offset example depicts an implementation where fundamental  
1716 error ranges start at byte address zero; the implementation shall return zero for  
1717 FUNDAMENTAL\_ERROR\_RANGE\_OFFSET.

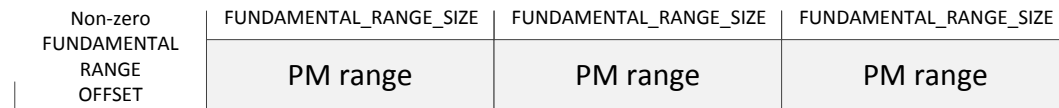
1718 **Figure 12 Zero range offset example**



1719 Byte Address Zero

1720 Figure 13 Non-zero range offset example depicts an implementation where fundamental  
1721 error ranges start at a non-zero offset; the implementation shall return the difference  
1722 between the starting byte address of a fundamental error range rounded down to a  
1723 multiple of FUNDAMENTAL\_ERROR\_RANGE.

1724 **Figure 13 Non-zero range offset example**



1725 Byte Address Zero

1726 **9.3.6 NVM.PM.VOLUME.DISCARD\_IF\_YOU\_CAN\_CAPABLE**

1727 Requirement: mandatory

1728 Returns true if the implementation supports DISCARD\_IF\_YOU\_CAN.

1729 **9.3.7 NVM.PM.VOLUME.DISCARD\_IMMEDIATELY\_CAPABLE**

1730 Requirement: mandatory

1731 Returns true if the implementation supports DISCARD\_IMMEDIATELY.

1732 **9.3.8 NVM.PM.VOLUME.DISCARD\_IMMEDIATELY\_RETURNS**

1733 Requirement: mandatory if DISCARD\_IMMEDIATELY\_CAPABLE (see 9.3.7) is true

1734 The value returned from read operations to bytes specified by a  
1735 DISCARD\_IMMEDIATELY action with no subsequent write operations. The possible  
1736 values are:

- 1737 • A value that is returned to each load of bytes in an unmapped range until the next  
1738 store action
- 1739 • Unspecified

1740 **9.3.9 NVM.PM.VOLUME.EXISTS\_CAPABLE**

1741 Requirement: mandatory

1742 This attribute indicates that the implementation is capable of the  
1743 NVM.PM.VOLUME.EXISTS action.

1744 **9.4 Use cases**

1745 **9.4.1 Initialization steps for a PM-aware file system**

1746 **Purpose/triggers:**

1747 Steps taken by a file system when a PM-aware volume is attached to a PM volume.

1748 **Scope/context:**

1749 NVM.PM.VOLUME mode

1750 **Preconditions:**

- 1751 • The administrator has defined a PM volume
- 1752 • The administrator has completed one-time steps to create a file system on the  
1753 PM volume

1754 **Inputs:**

- 1755 • A reference to a PM volume
- 1756 • The name of a PM file system

1757 **Success scenario:**

- 1758 1. The file system issues a GET\_RANGESET action to retrieve information about  
1759 the ranges comprised by the PM volume.
- 1760 2. The file system uses the range information from GET\_RANGESET to determine  
1761 physical address range(s) and offset(s) of the file system's primary metadata (for  
1762 example, the primary superblock), then loads appropriate metadata to determine  
1763 no additional validity checking is needed.

- 1764 3. The file system sets a flag in the metadata indicating the file system is mounted  
1765 by storing the updated status to the appropriate location  
1766 a. If the range containing this location requires VIRTUAL\_ADDRESS\_SYNC  
1767 or PHYSICAL\_ADDRESS\_SYNC is needed (based on  
1768 GET\_RANGESET's sync mode property), the file system invokes the  
1769 appropriate SYNC action

1770 **Postconditions:**

1771 The file system is usable by applications.

1772 9.4.2 **Driver emulates a block device using PM media**

1773 **Purpose/triggers:**

1774 The steps supporting an application write action from a driver that emulates a block  
1775 device using PM as media.

1776 **Scope/context:**

1777 NVM.PM.VOLUME mode

1778 **Preconditions:**

1779 PM layer FUNDAMENTAL\_SIZE reported to driver is cache line size.

1780 **Inputs:**

1781 The application provides:

- 1782 • the starting address of the memory (could be volatile) memory containing the  
1783 data to write  
1784 • the length of the memory range to be written,  
1785 • an OS-specific reference to a block device (the virtual device backed by the PM  
1786 volume),  
1787 • the starting LBA within that block device

1788 **Success scenario:**

- 1789 1. The driver registers with the OS-specific component to be notified of errors on the  
1790 PM volume. PM error handling is outside the scope of this specification, but may be  
1791 similar to what is described in (and above) Figure 16 Linux Machine Check error flow  
1792 with proposed new interface.  
1793 2. Using information from a GET\_RANGESET response, the driver splits the write  
1794 operating into separate pieces if the target PM addresses (corresponding to  
1795 application target LBAs) are in different ranges with different "sync type" values. For  
1796 each of these pieces:

- 1797           a. Using information from a GET\_RANGESET response, the driver determines  
1798           the PM memory address corresponding to the input starting LBA, and  
1799           performs a memory copy operation from the callers input memory to the PM  
1800           b. The driver then performs a platform-specific flush operation  
1801           c. Using information from a GET\_RANGESET response, the driver invokes the  
1802           PHYSICAL\_ADDRESS\_SYNC or VIRTUAL\_ADDRESS\_SYNC action as  
1803           needed  
1804        3. No PM errors are reported by the PM error component, the driver reports that the  
1805        write action succeeded.

1806        **Alternative Scenario 1:**

1807        In step 3 in the Success Scenario, the PM error component reports a PM error. The  
1808        driver verifies that this error impacts the PM range being written and returns an error to  
1809        the caller.

1810        **Postconditions:**

1811        The target PM range (i.e., the block device LBA range) is updated.

1812        **See also:**

1813        4.2.4 NVM block volume using PM hardware

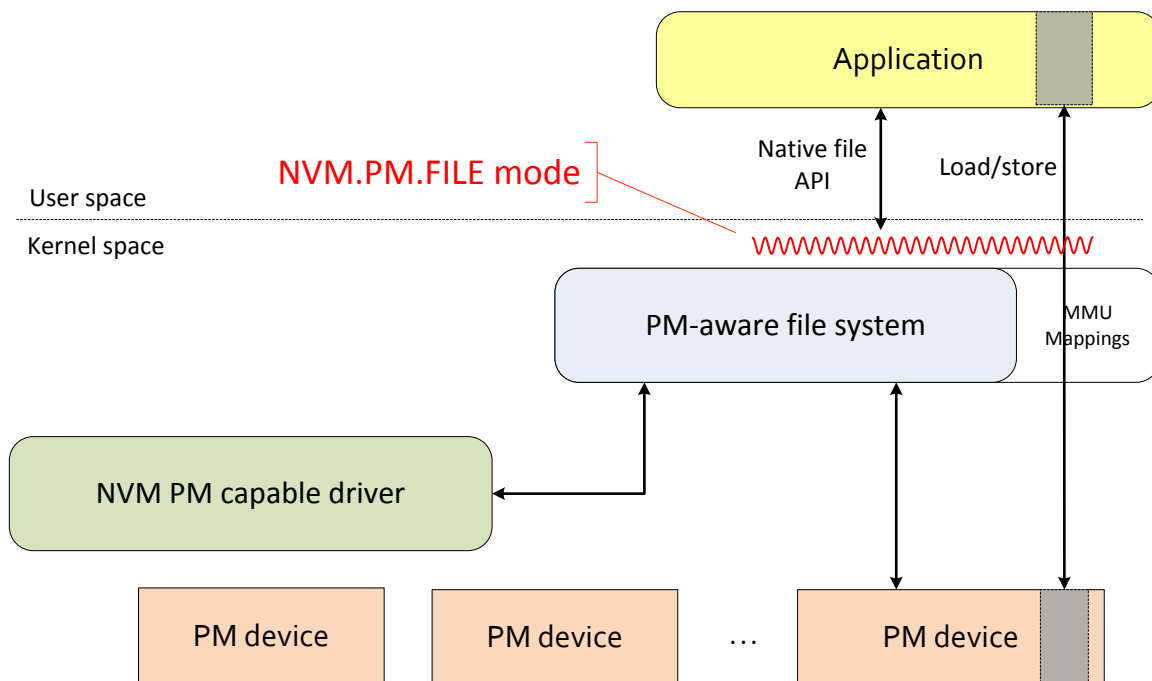
1814 **10 NVM.PM.FILE**

1815 **10.1 Overview**

1816 The NVM.PM.FILE mode access provides a means for user space applications to  
1817 directly access NVM as memory. Most of the standard actions in this mode are intended  
1818 to be implemented as APIs exported by existing file systems. An NVM.PM.FILE  
1819 implementation behaves similarly to preexisting file system implementations, with minor  
1820 exceptions. This section defines extensions to the file system implementation to  
1821 accommodate persistent memory mapping and to assure interoperability with  
1822 NVM.PM.FILE mode applications.

1823 Figure 14 NVM.PM.FILE mode example shows the context surrounding the point in a  
1824 system (the bold red line) where the NVM.PM.FILE mode programming model is  
1825 exposed by a PM-aware file system. A user space application consumes the  
1826 programming model as is typical for current file systems. This example is not intended  
1827 to preclude the possibility of a user space PM-aware file system implementation. It  
1828 does, however presume that direct load/store access from user space occurs within a  
1829 memory-mapped file context. The PM-aware file system interacts with an NVM PM  
1830 capable driver to achieve any non-direct-access actions needed to discover or configure  
1831 NVM. The PM-aware file system may access NVM devices for purposes such as file  
1832 allocation, free space or other metadata management. The PM-aware file system  
1833 manages additional metadata that describes the mapping of NVM device memory  
1834 locations directly into user space.

1835 **Figure 14 NVM.PM.FILE mode example**



1836



1837 Once memory mapping occurs, the behavior of the NVM.PM.FILE mode diverges from  
1838 NVM.FILE mode because accesses to mapped memory are in the persistence domain  
1839 as soon as they reach memory. This is represented in Figure 14 NVM.PM.FILE mode  
1840 example by the arrow that passes through the “MMU Mappings” extension of the file  
1841 system. As a result of persistent memory mapping, primitive ACID properties arise from  
1842 CPU and memory infrastructure behavior as opposed to disk drive or traditional SSD  
1843 behavior. Note that writes may still be retained within processor resident caches or  
1844 memory controller buffers before they reach a persistence domain. As with  
1845 NMV.FILE.SYNC, the possibility remains that memory mapped writes may become  
1846 persistent before a corresponding NVM.PM.FILE.SYNC action.

1847 The following actions have behaviors specific to the NVM.PM.FILE mode:

1848       NVM.PM.FILE.MAP – Add a subset of a PM file to application's address space for  
1849       load/store access.

1850       NVM.PM.FILE.SYNC – Synchronize persistent memory content to assure  
1851       durability and enable recovery by forcing data to reach the persistence domain.

## 1852 **10.2 Actions**

1853 The following actions are mandatory for compliance with the NVM Programming Model  
1854 NVM.PM.FILE mode.

### 1855 **10.2.1 Actions that apply across multiple modes**

1856 The following actions apply to NVM.PM.FILE mode as well as other modes.

1857       NVM.COMMON.GET\_ATTRIBUTE (see 6.10.1)

1858       NVM.COMMON.SET\_ATTRIBUTE (see 6.10.2)

### 1859 **10.2.2 Native file system actions**

1860 Native actions shall apply with unmodified syntax and semantics provided that they are  
1861 compatible with programming model specific actions. This is intended to support  
1862 traditional file operations allowing many applications to use PM without modification.  
1863 This specifically includes mandatory implementation of the native synchronization of  
1864 mapped files. As always, specific implementations may choose whether or not to  
1865 implement optional native operations.

### 1866 **10.2.3 NVM.PM.FILE.MAP**

1867 Requirement: mandatory

1868 The mandatory form of this action shall have the same syntax found in a pre-existing file  
1869 system, preferably the operating system's native file map call. The specified subset of a  
1870 PM file is added to application's address space for load/store access. The semantics of  
1871 this action are unlike the native MAP action because NVM.PM.FILE.MAP causes direct  
1872 load/store access. For example, the role of the page cache might be reduced or  
1873 eliminated. This reduces or eliminates the consumption of volatile memory as a staging

1874 area for non-volatile data. In addition, by avoiding demand paging, direct access can  
1875 enable greater uniformity of access time across volatile and non-volatile data.

1876 PM mapped file operation may not provide the access time and modify time behavior  
1877 typical of native file systems.

1878 PM mapped file operation may not provide the normal semantics for the native file  
1879 synchronization actions (e.g., POSIX fsync and fdatasync and Win32 FlushFileBuffers).  
1880 If a file is mapped at the time when the native file synchronization action is invoked, the  
1881 normal semantics apply. However if the file had been mapped, data had been written to  
1882 the file through the map, the data had not been synchronized by use of the  
1883 NVM.PM.FILE.SYNC action or the native mapped file sync action, and the mapping had  
1884 been removed prior to the execution of the native file synchronization action, the action  
1885 is not required to synchronize the data written to the map.

1886 Requires NVM.PM.FILE.OPEN

1887 Inputs: align with native operating system's map

1888 Outputs: align with native operating system's map

1889 Relevant Options:

1890 All of the native file system options should apply.

1891 NVM.PM.FILE.MAP\_SHARED (Mandatory) – This existing native option shall be  
1892 supported by the NVM.PM.FILE.MAP action. This option indicates that user  
1893 space processes other than the writer can see any changes to mapped memory  
1894 immediately.

1895 NVM.PM.FILE.MAP\_COPY\_ON\_WRITE (Optional)– This existing native option  
1896 indicates that any write after mapping will cause a copy on write to volatile  
1897 memory, or PM that is discarded during any type of restart. The copy is only  
1898 visible to the writer. The copy is not folded back into PM during the sync  
1899 command.

1900 Relevant Attributes:

1901 NVM.PM.FILE.MAP\_COPY\_ON\_WRITE\_CAPABLE (see 10.3.2) - Native  
1902 operating system map commands make a distinction between MAP\_SHARED  
1903 and MAP\_COPY\_ON\_WRITE. Both are supported with native semantics under  
1904 the NVM Programming Model. This attribute indicates whether the  
1905 MAP\_COPY\_ON\_WRITE mapping mode is supported. All NVM.PM.FILE.MAP  
1906 implementations shall support the MAP\_SHARED option.

1907 Error handling for mapped ranges of persistent memory is unlike I/O, in that there is no  
1908 acknowledgement to a load or store instruction. Instead processors equipped to detect  
1909 memory access failures respond with machine checks. These can be routed to user

1910 threads as asynchronous events. With memory-mapped PM, asynchronous events are  
1911 the primary means of discovering the failure of a load to return good data. Please refer  
1912 to NVM.PM.FILE.GET\_ERROR\_INFO (section 10.2.6) for more information on error  
1913 handling behavior.

1914 Depending on memory configuration, CPU memory write pipelines may effectively  
1915 preclude application level error handling during memory accesses that result from store  
1916 instructions. For example, errors detected during the process of flushing the CPU's write  
1917 pipeline are more likely to be associated with that pipeline than the NVM itself. Errors  
1918 that arise within the CPU's write pipeline generally do not enable application level  
1919 recovery at the point of the error. As a result application processes may be forced to  
1920 restart when these errors occur (see PM Error Handling Annex C). Such errors should  
1921 appear in CPU event logs, leading to an administrative response that is outside the  
1922 scope of this specification.

1923 Applications needing timely assurance that recently stored data is recoverable should  
1924 use the NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY action to verify data from  
1925 NVM after it is flushed (see 10.2.7). Errors during verify are handled in the manner  
1926 described in this annex.

#### 1927 10.2.4 **NVM.PM.FILE.SYNC**

1928 Requirement: mandatory

1929 The purpose of this action is to synchronize persistent memory content to assure  
1930 durability and enable recovery by forcing data to reach the persistence domain.

1931 The native file system sync action may be supported by implementations that also  
1932 support NVM.PM.FILE.SYNC. The intent is that the semantics of NVM.PM.FILE.SYNC  
1933 match native sync operation on memory-mapped files however because persistent  
1934 memory is involved, NVM.PM.FILE implementations need not flush full pages. Note that  
1935 writes may still be subject to functionality that may mask whether stored data has  
1936 reached the persistence domain (such as caching or buffering within processors or  
1937 memory controllers). NVM.PM.FILE.SYNC is responsible for insuring that data within  
1938 the processor or memory buffers reaches the persistence domain.

1939 A number of boundary conditions can arise regarding interoperability of PM and non-PM  
1940 implementation components. An annex to this specification is being proposed to  
1941 address this. The following limitations apply:

- 1942 • The behavior of an NVM.PM.FILE.SYNC action applied to a range in a file that was  
1943 not mapped using NVM.PM.FILE.MAP is unspecified.
- 1944 • The behavior of NVM.PM.FILE.SYNC on non-persistent memory is unspecified.

1945 In both the PM and non-PM modes, updates to ranges mapped as shared can and may  
1946 become persistent in any order before a sync requires them all to become persistent.

1947 The sync action applied to a shared mapping does not guarantee write atomicity. The  
1948 byte range referenced by the sync parameter may have reached a persistence domain

1949 prior to the sync command. The sync action guarantees only that the range referenced  
1950 by the sync action will reach the persistence domain before the successful completion  
1951 of the sync action. Any atomicity that is achieved is not caused by the sync action itself.

1952 Requires: NVM.PM.FILE.MAP

1953 Inputs: Align with native operating system's sync with the exception that alignment  
1954 restrictions are relaxed.

1955 Outputs: Align with native operating system's sync with the addition that it shall return  
1956 an error code.

1957 Users of the NVM.PM.FILE.SYNC action should be aware that for files that are mapped  
1958 as shared, there is no requirement to buffer data on the way to the persistence domain.  
1959 Although data may traverse a processor's write pipeline and other buffers within  
1960 memory controllers these are more transient than the disk I/O buffering that is common  
1961 in NVM.FILE implementations.

1962 Error handling related to this action is expected to be derived from ongoing work that  
1963 begins with Annex C (Informative) PM error handling.

#### 1964 10.2.5 **NVM.PM.FILE.OPTIMIZED\_FLUSH**

1965 Requirement: mandatory if NVM.PM.OPTIMIZED\_FLUSH\_CAPABLE is set.

1966 The purpose of this action is to synchronize multiple ranges of persistent memory  
1967 content to assure durability and enable recovery by forcing data to reach the  
1968 persistence domain. This action has the same effect as NVM.PM.FILE.SYNC however it  
1969 is intended to allow additional implementation optimization by excluding options  
1970 supported by sync and by allowing multiple byte ranges to be synchronized during a  
1971 single action. Page oriented alignment constraints imposed by the native definition are  
1972 lifted. Because of this, implementations might be able to use underlying persistent  
1973 memory more optimally than they could with the native sync. In addition some  
1974 implementations may enable this action to avoid context switches into kernel space.  
1975 With the exception of these differences all of the content of the NVM.PM.FILE.SYNC  
1976 action description also applies to NVM.PM.FILE.OPTIMIZED\_FLUSH.

1977 Inputs: Identical to NVM.PM.FILE.SYNC except that an array of byte ranges is specified  
1978 and options are precluded. A reference to the array and the size of the array are input  
1979 instead of a single address and length. Each element of the array contains an address  
1980 and length of a range of bytes to be synchronized.

1981 Outputs: Align with native OS's sync with the addition that it shall return an error code.

1982 Relevant attributes: NVM.PM.FILE.OPTIMIZED\_FLUSH\_CAPABLE – Indicates whether  
1983 this action is supported by the NVM.PM.FILE implementation (see 10.3.5).

1984 NVM.PM.FILE.OPTIMIZED\_FLUSH provides no guarantee of atomicity within or across  
1985 the synchronized byte ranges. Neither does it provide any guarantee of the order in  
1986 which the bytes within the ranges of the action reach a persistence domain.

1987 In the event of failure the progress of the action is indeterminate. Various byte ranges  
1988 may or may not have reached a persistence domain. There is no indication as to which  
1989 byte ranges may have been synchronized.

#### 1990 10.2.6 **NVM.PM.FILE.GET\_ERROR\_EVENT\_INFO**

1991 Requirement: mandatory if NVM.PM.ERROR\_EVENT\_CAPABLE is set.

1992 The purpose of this action is to provide a sufficient description of an error event to  
1993 enable recovery decisions to be made by an application. This action is intended to  
1994 originate during an application event handler in response to a persistent memory error.  
1995 In some implementations this action may map to the delivery of event description  
1996 information to the application at the start of the event handler rather than a call made by  
1997 the event handler. The error information returned is specific to the memory error that  
1998 caused the event.

1999 Inputs: It is assumed that implementations can extract the information output by this  
2000 action from the event being handled.

2001 Outputs:

2002 1 – An indication of whether or not execution of the application can be resumed from the  
2003 point of interruption. If execution cannot be resumed then the process running the  
2004 application should be restarted for full recovery.

2005 2 – An indication of error type enabling the application to determine whether an address  
2006 is provided and the direction of data flow (load/verify vs. store) when the error was  
2007 detected.

2008 3 – The memory mapped address and length of the byte range where data loss was  
2009 detected by the event.

2010 Relevant attributes:

2011 NVM.PM.FILE.ERROR\_EVENT\_CAPABLE – Indicates whether load error event  
2012 handling and this action are supported by the NVM.PM.FILE implementation (see  
2013 10.3.6).

2014 This action is used to obtain information about an error that caused a machine check  
2015 involving memory mapped persistent memory. This is necessary because with  
2016 persistent memory there is no opportunity to provide error information as part of a  
2017 function call or I/O. The intent is to allow sophisticated error handling and recovery to  
2018 occur before the application sees the event by allowing the NVM.PM.FILE  
2019 implementation to handle it first. It is expected that after NVM.PM.FILE has completed

2020 whatever recovery is possible, the application error handler will be called and use the  
2021 error information described here to stage subsequent recovery actions, some of which  
2022 may occur after the application's process is restarted.

2023 In some implementations the same event handler may be used for many or all memory  
2024 errors. Therefore this action may arise from memory accesses unrelated to NVM. It is  
2025 the application event handler's responsibility to determine whether the memory range  
2026 indicated is relevant for recovery. If the memory range is irrelevant then the event  
2027 should be ignored other than as a potential trigger for a restart.

2028 In some systems, errors related to memory stores may not provide recovery information  
2029 to the application unless and until load instructions attempt to access the memory  
2030 locations involved. This can be accomplished using the  
2031 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY action (section 10.2.7).

2032 For more information on the circumstances which may surround this action please refer  
2033 to PM Error Handling Annex C.

#### 2034 10.2.7 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY

2035 Requirement: mandatory if  
2036 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY\_CAPABLE is set.

2037 The purpose of this action is to synchronize multiple ranges of persistent memory  
2038 content to assure durability and enable recovery by forcing data to reach the  
2039 persistence domain. Furthermore, this action verifies that data was written correctly by  
2040 verifying it. The intent is to supply a mechanism whereby the application can receive  
2041 data integrity assurance on writes to memory-mapped PM prior to completion of this  
2042 action. This is the PM equivalent to the POSIX definition of synchronized I/O which  
2043 clarifies that the intent of synchronized I/O data integrity completion is "so that an  
2044 application can ensure that the data being manipulated is physically present on  
2045 secondary mass storage devices".

2046 Except for the additional verification of flushed data, this action has the same effect as  
2047 NVM.PM.FILE.OPTIMIZED\_FLUSH.

2048 Inputs: Identical to NVM.PM.FILE.OPTIMIZED\_FLUSH.

2049 Outputs: Align with native OS's sync with the addition that it shall return an error code.  
2050 The error code indicates whether or not all data in the indicated range set is readable.

2051 Relevant attributes:

2052 NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY\_CAPABLE – Indicates whether this  
2053 action is supported by the NVM.PM.FILE implementation (see 10.3.7).

2054 OPTIMIZED\_FLUSH\_AND\_VERIFY shall assure that data has been verified to be  
2055 readable. Any errors discovered during verification should be logged for administrative

2056 attention. Verification shall occur across all data ranges specified in the action  
2057 regardless of when they were actually flushed. Verification shall complete prior to  
2058 completion of the action.

2059 In the event of failure the progress of the action is indeterminate. .

## 2060 **10.3 Attributes**

### 2061 **10.3.1 Attributes that apply across multiple modes**

2062 The following attributes apply to NVM.PM.FILE mode as well as other modes.

2063 NVM.COMMON.SUPPORTED\_MODES (see 6.11.1)

2064 NVM.COMMON.FILE\_MODE (see 6.11.2)

### 2065 **10.3.2 NVM.PM.FILE.MAP\_COPY\_ON\_WRITE\_CAPABLE**

2066 Requirement: mandatory

2067 This attribute indicates that MAP\_COPY\_ON\_WRITE option is supported by the  
2068 NVM.PM.FILE.MAP action.

### 2069 **10.3.3 NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMICITY**

2070 Requirement: mandatory

2071 INTERRUPTED\_STORE\_ATOMICITY indicates whether the volume supports power  
2072 fail atomicity of aligned store operations on fundamental data types. To achieve failure  
2073 atomicity, aligned operations on fundamental data types reach NVM atomically.  
2074 Formally “aligned operations on fundamental data types” is implementation defined. See  
2075 Annex B(Informative) Consistency.

2076 A value of true indicates that after an aligned store of a fundamental data type is  
2077 interrupted by reset, power loss or system crash; upon restart the contents of persistent  
2078 memory reflect either the state before the store or the state after the completed store. A  
2079 value of false indicates that after a store interrupted by reset, power loss or system  
2080 crash, upon restart the contents of memory may be such that subsequent loads may  
2081 create exceptions. A value of false also indicates that after a store interrupted by reset,  
2082 power loss or system crash; upon restart the contents of persistent memory may not  
2083 reflect either the state before the store or the state after the completed store.

2084 The value of this attribute is true only if it’s true for all ranges in the file system.

### 2085 **10.3.4 NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE**

2086 Requirement: mandatory

2087 FUNDAMENTAL\_ERROR\_RANGE is the number of bytes that may become  
2088 unavailable due to an error on an NVM device.

2089 An application may organize data in terms of `FUNDAMENTAL_ERROR_RANGE` to  
2090 assure two key data items are not likely to be affected by a single error.

2091 Unlike `NVM.PM.VOLUME`, `NVM.PM.FILE` does not associate an offset with the  
2092 `FUNDAMENTAL_ERROR_RANGE` (see section 9.3.5). because the file system is  
2093 expected to handle any volume mode offset transparently to the application. The value  
2094 of this attribute is the maximum of the values for all ranges in the file system.

### 2095 10.3.5 **NVM.PM.FILE.OPTIMIZED\_FLUSH\_CAPABLE**

2096 Requirement: mandatory

2097 This attribute indicates that the `OPTIMIZED_FLUSH` action is supported by the  
2098 `NVM.PM.FILE` implementation.

### 2099 10.3.6 **NVM.PM.FILE.ERROR\_EVENT\_CAPABLE**

2100 Requirement: mandatory

2101 This attribute indicates that the `NVM.PM.FILE` implementation is capable of handling  
2102 error events in such a way that, in the event of data loss, those events are subsequently  
2103 delivered to applications. If error event handling is supported then  
2104 `NVM.PM.FILE.GET_ERROR_INFO` action shall also be supported.

### 2105 10.3.7 **NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY\_CAPABLE**

2106 Requirement: mandatory

2107 This attribute indicates that the `OPTIMIZED_FLUSH_AND_VERIFY` action is supported  
2108 by the `NVM.PM.FILE` implementation.

## 2109 **10.4 Use cases**

### 2110 10.4.1 **Update PM File Record**

2111 Update a record in a PM file.

#### 2112 **Purpose/triggers:**

2113 An application using persistent memory updates an existing record. For simplicity, this  
2114 application uses fixed size records. The record size is defined by application data  
2115 considerations.

#### 2116 **Scope/context:**

2117 Persistent memory context; this use case shows basic behavior.

#### 2118 **Preconditions:**

- 2119
- 2120 • The administrator created a PM file and provided its name to the application; this  
2121 name is accessible to the application – perhaps in a configuration file
  - The application has populated the PM file contents



2122 • The PM file is not in use at the start of this use case (no sharing considerations)

2123 **Inputs:**

2124 The content of the record, the location (relative to the file) where the record resides

2125 **Success scenario:**

- 2126 1) The application uses the native OPEN action, passing in the file name
- 2127 2) The application uses the NVM.PM.FILE.MAP action, passing in the file descriptor
- 2128 returned by the native OPEN. Since the records are not necessarily page aligned,
- 2129 the application maps the entire file.
- 2130 3) The application registers for memory hardware exceptions
- 2131 4) The application stores the new record content to the address returned by
- 2132 NVM.PM.FILE.MAP offset by the record's location
- 2133 5) The application uses NVM.PM.FILE.SYNC to flush the updated record to the
- 2134 persistence domain
- 2135 a. The application may simply sync the entire file
- 2136 b. Alternatively, the application may limit the range to be sync'd
- 2137 6) The application uses the native UNMAP and CLOSE actions to clean up.

2138 **Failure Scenario:**

2139 While reading PM content (accessing via a load operation), a memory hardware

2140 exception is reported. The application's event handler is called with information about

2141 the error as described in NVM.PM.FILE.GET\_ERROR\_INFO. Based on the information

2142 provided, the application records the error for subsequent recovery and determines

2143 whether to restart or continue execution.

2144 **Outputs:**

2145 None

2146 **Postconditions:**

2147 The record is updated.

2148 **10.4.2 Direct load access**

2149 **Purpose/triggers:**

2150 An application developer wishes to retrieve data from a persistent memory-mapped file

2151 using direct memory load instruction access with error handling for uncorrectable errors.

2152 **Scope/context:**

2153 NVM.PM.FILE

2154 **Inputs:**

- 2155 • Virtual address of the data.

- 2156 **Outputs:**
- 2157 • Data from persistent memory if successful
  - 2158 • Error code if an error was detected within the accessed memory range.
- 2159 **Preconditions:**
- 2160 • The persistent memory file must be mapped into a region of virtual memory.
  - 2161 • The virtual address must be within the mapped region of the file.
- 2162 **Postconditions:**
- 2163 • If an error was returned, the data may be unreadable. Future load accesses may
  - 2164 continue to return an error until the data is overwritten to clear the error condition
  - 2165 • If no error was returned, there is no postcondition.

2166 **Success and Failure Scenarios:**

2167 Consider the following fragment of example source code, which is simplified from the  
 2168 code for the function that reads SQLite’s transaction journal:

```

2169     retCode = pread(journalFD, magic, 8, off);
2170     if (retCode != SQLITE_OK) return retCode;
2171
2172     if (memcmp(magic, journalMagic, 8) != 0)
2173         return SQLITE_DONE;
  
```

2174 This example code reads an eight-byte magic number from the journal header into an  
 2175 eight-byte buffer named *magic* using a standard file *read* call. If an error is returned from  
 2176 the *read* system call, the function exits with an error return code indicating that an I/O  
 2177 error occurred. If no error occurs, it then compares the contents of the *magic* buffer  
 2178 against the expected magic number constant named *journalMagic*. If the contents of the  
 2179 buffer do not match the expected magic number, the function exits with an error return  
 2180 code.

2181 An equivalent version of the function using direct memory load instruction access to a  
 2182 mapped file is:

```

2183     volatile siginfo_t errContext;
2184     ...
2185     int retCode = SQLITE_OK;
2186
2187     TRY
2188     {
2189         if (memcmp(journalMmapAddr + off, journalMagic, 8) != 0)
2190             retCode = SQLITE_DONE;
2191     }
2192     CATCH(BUS_MCEERR_AR)
2193     {
2194         if ((errContext.si_code == BUS_MCEERR_AR) &&
2195             (errContext.si_addr >= journalMmapAddr) &&
2196             (errContext.si_addr < (journalMmapAddr + journalMmapSize))){
2197             retCode = SQLITE_IOERR;
2198         } else {
  
```

```

2199         signal(errContext.si_signo, SIG_DFL);
2200         raise(errContext.si_signo);
2201     }
2202 }
2203 ENDTRY;
2204
2205     if (retCode != SQLITE_OK) return retCode;

```

2206 The mapped file example compares the magic number in the header of the journal file  
 2207 against the expected magic number using the *memcmp* function by passing a pointer  
 2208 containing the address of the magic number in the mapped region of the file. If the  
 2209 contents of the magic number member of the file header do not match the expected  
 2210 magic number, the function exits with an error return code.

2211 The application-provided TRY/CATCH/ENDTRY macros implement a form of exception  
 2212 handling using POSIX *sigsetjmp* and *siglongjmp* C library functions. The TRY macro  
 2213 initializes a *sigjmp\_buf* by calling *sigsetjmp*. When a SIGBUS signal is raised, the signal  
 2214 handler calls *siglongjmp* using the *sigjmp\_buf* set by the *sigsetjmp* call in the TRY  
 2215 macro. Execution then continues in the CATCH clause. (Caution: the code in the TRY  
 2216 block should not call library functions as they are not likely to be exception-safe.) Code  
 2217 for the Windows platform would be similar except that it would use the standard  
 2218 Structured Exception Handling *try-except* statement catching the  
 2219 EXCEPTION\_IN\_PAGE\_ERROR exception rather than application-provided  
 2220 TRY/CATCH/ENDTRY macros.

2221 If an error occurs during the read of the magic number data from the mapped file, a  
 2222 SIGBUS signal will be raised resulting in the transfer of control to the CATCH clause.  
 2223 The address of the error is compared against the range of the memory-mapped file. In  
 2224 this example the error address is assumed to be in the process's logical address space.  
 2225 If the error address is within the range of the memory-mapped file, the function returns  
 2226 an error code indication that an I/O error occurred. If the error address is outside the  
 2227 range of the memory-mapped file, the error is assumed to be for some other memory  
 2228 region such as the program text, stack, or heap, and the signal or exception is re-raised.  
 2229 This is likely to result in a fatal error for the program.

2230 **See also:**

- 2231 • Microsoft Corporation, Reading and Writing From a File View (Windows),  
 2232 available from  
 2233 <http://msdn.microsoft.com/en-us/library/windows/desktop/aa366801.aspx>

2234 **10.4.3 Direct store access**

2235 **Purpose/triggers:**

2236 An application developer wishes to place data in a persistent memory-mapped file using  
 2237 direct memory store instruction access.

2238 **Scope/context:**

2239 NVM.PM.FILE

2240 **Inputs:**

- 2241     • Virtual address of the data.  
2242     • The data to store.

2243 **Outputs:**

- 2244     • Error code if an error occurred.

2245 **Preconditions:**

- 2246     • The persistent memory file must be mapped into a region of virtual memory.  
2247     • The virtual address must be within the mapped region of the file.

2248 **Postconditions:**

- 2249     • If an error was returned, the state of the data recorded in the persistence domain  
2250       is indeterminate.  
2251     • If no error was returned, the specified data is either recorded in the persistence  
2252       domain or an undiagnosed error may have occurred.

2253 **Success and Failure Scenarios:**

2254 Consider the following fragment of example source code, which is simplified from the  
2255 code for the function that writes to SQLite's transaction journal:

```
2256     ret = pwrite(journalFD, dbPgData, dbPgSize, off);  
2257     if (ret != SQLITE_OK) return ret;  
2258     ret = write32bits(journalFD, off + dbPgSize, cksum);  
2259     if (ret != SQLITE_OK) return ret;  
2260     ret = fdatasync(journalFD);  
2261     if (ret != SQLITE_OK) return ret;
```

2262 This example code writes a page of data from the database cache to the journal using a  
2263 standard file *write* call. If an error is returned from the *write* system call, the function  
2264 exits with an error return code indicating that an I/O error occurred. If no error occurs,  
2265 the function then appends the checksum of the data, again using a standard file *write*  
2266 call. If an error is returned from the *write* system call, the function exits with an error  
2267 return code indicating that an I/O error occurred. If no error occurs, the function then  
2268 invokes the *fdatasync* system call to flush the written data from the file system buffer  
2269 cache to the persistence domain. If an error is returned from the *fdatasync* system call,  
2270 the function exits with an error return code indicating that an I/O error occurred. If no  
2271 error occurs, the written data has been recorded in the persistence domain.

2272 An equivalent version of the function using direct memory store instruction access to a  
2273 memory-mapped file is:

```
2274     memcpy(journalMmapAddr + off, dbPgData, dbPgSize);  
2275     PM_track_dirty_mem(dirtyLines, journalMmapAddr + off, dbPgSize);  
2276  
2277     store32bits(journalMmapAddr + off + dbPgSize, cksum);  
2278     PM_track_dirty_mem(dirtyLines, journalMmapAddr + off + dbPgSize, 4);  
2279
```

```
2280     ret = PM_optimized_flush(dirtyLines, dirtyLinesCount);
2281
2282     if (ret == SQLITE_OK) dirtyLinesCount = 0;
2283
2284     return ret;
```

2285 The memory-mapped file example writes a page of data from the database cache to the  
2286 journal using the *memcpy* function by passing a pointer containing the address of the  
2287 page data field in the mapped region of the file. It then appends the checksum using  
2288 direct stores to the address of the checksum field in the mapped region of the file.

2289 The code calls the application-provided *PM\_track\_dirty\_mem* function to record the  
2290 virtual address and size of the memory regions that it has modified. The  
2291 *PM\_track\_dirty\_mem* function constructs a list of these modified regions in the  
2292 *dirtyLines* array.

2293 The function then calls the *PM\_optimized\_flush* function to flush the written data to the  
2294 persistence domain. If an error is returned from the *PM\_optimized\_flush* call, the  
2295 function exits with an error return code indicating that an I/O error occurred. If no error  
2296 occurs, the written data is either recorded in the persistence domain or an undiagnosed  
2297 error may have occurred. Note that this postcondition is weaker than the guarantee  
2298 offered by the *fdatsync* system call in the original example.

2299 **See also:**

- 2300
  - Microsoft Corporation, Reading and Writing From a File View (Windows),  
2301 available from  
2302 <http://msdn.microsoft.com/en-us/library/windows/desktop/aa366801.aspx>

2303 **10.4.4 Direct store access with synchronized I/O data integrity completion**

2304 **Purpose/triggers:**

2305 An application developer wishes to place data in a persistent memory-mapped file using  
2306 direct memory store instruction access with synchronized I/O data integrity completion.

2307 **Scope/context:**

2308 NVM.PM.FILE

2309 **Inputs:**

- 2310
  - Virtual address of the data.
  - The data to store.

2312 **Outputs:**

- 2313
  - Error code if an error occurred.

2314 **Preconditions:**

- 2315
  - 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.

- 2317 **Postconditions:**
- 2318 • If an error was returned, the state of the data recorded in the persistence domain
  - 2319 is indeterminate.
  - 2320 • If no error was returned, the specified data is recorded in the persistence domain.

2321 **Success and Failure Scenarios:**

2322 Consider the following fragment of example source code, which is simplified from the

2323 code for the function that writes to SQLite's transaction journal:

```
2324     ret = pwrite(journalFD, dbPgData, dbPgSize, off);
2325     if (ret != SQLITE_OK) return ret;
2326     ret = write32bits(journalFD, off + dbPgSize, cksum);
2327     if (ret != SQLITE_OK) return ret;
2328
2329     ret = fdatsync(journalFD);
2330     if (ret != SQLITE_OK) return ret;
```

2331 This example code writes a page of data from the database cache to the journal using a

2332 standard file *write* call. If an error is returned from the *write* system call, the function

2333 exits with an error return code indicating that an I/O error occurred. If no error occurs,

2334 the function then appends the checksum of the data, again using a standard file *write*

2335 call. If an error is returned from the *write* system call, the function exits with an error

2336 return code indicating that an I/O error occurred. If no error occurs, the function then

2337 invokes the *fdatsync* system call to flush the written data from the file system buffer

2338 cache to the persistence domain. If an error is returned from the *fdatsync* system call,

2339 the function exits with an error return code indicating that an I/O error occurred. If no

2340 error occurs, the written data has been recorded in the persistence domain.

2341 An equivalent version of the function using direct memory store instruction access to a

2342 memory-mapped file is:

```
2343     memcpy(journalMmapAddr + off, dbPgData, dbPgSize);
2344     PM_track_dirty_mem(dirtyLines, journalMmapAddr + off, dbPgSize);
2345
2346     store32bits(journalMmapAddr + off + dbPgSize, cksum);
2347     PM_track_dirty_mem(dirtyLines, journalMmapAddr + off + dbPgSize, 4);
2348
2349     ret = PM_optimized_flush_and_verify(dirtyLines, dirtyLinesCount);
2350
2351     if (ret == SQLITE_OK) dirtyLinesCount = 0;
2352
2353     return ret;
```

2354 The memory-mapped file example writes a page of data from the database cache to the

2355 journal using the *memcpy* function by passing a pointer containing the address of the

2356 page data field in the mapped region of the file. It then appends the checksum using

2357 direct stores to the address of the checksum field in the mapped region of the file.

2358 The code calls the application-provided *PM\_track\_dirty\_mem* function to record the

2359 virtual address and size of the memory regions that it has modified. The

2360 *PM\_track\_dirty\_mem* function constructs a list of these modified regions in the  
2361 *dirtyLines* array.

2362 The function then calls the *PM\_optimized\_flush\_and\_verify* function to flush the written  
2363 data to the persistence domain. If an error is returned from the  
2364 *PM\_optimized\_flush\_and\_verify* call, the function exits with an error return code  
2365 indicating that an I/O error occurred. If no error occurs, the written data has been  
2366 recorded in the persistence domain. Note that this postcondition is equivalent to the  
2367 guarantee offered by the *fdatasync* system call in the original example.

2368 **See also:**

- 2369 • Microsoft Corp, FlushFileBuffers function (Windows),  
2370 <http://msdn.microsoft.com/en-us/library/windows/desktop/aa364439.aspx>
- 2371 • Oracle Corp, Synchronized I/O section in the Programming Interfaces Guide,  
2372 available from  
2373 <http://docs.oracle.com/cd/E19683-01/816-5042/chap7rt-57/index.html>
- 2374 • The Open Group, “The Open Group Base Specification Issue 6”, section 3.373  
2375 “Synchronized Input and Output”, available from  
2376 [http://pubs.opengroup.org/onlinepubs/007904975/basedefs/xbd\\_chap03.html#ta](http://pubs.opengroup.org/onlinepubs/007904975/basedefs/xbd_chap03.html#tag_03_373)  
2377 [g\\_03\\_373](http://pubs.opengroup.org/onlinepubs/007904975/basedefs/xbd_chap03.html#tag_03_373)

2378 **10.4.5 Persistent Memory Transaction Logging**

2379 **Purpose/Triggers:**

2380 An application developer wishes to implement a transaction log that maintains data  
2381 integrity through system crashes, system resets, and power failures. The underlying  
2382 storage is byte-granular persistent memory.

2383 **Scope/Context:**

2384 NVM.PM.VOLUME and NVM.PM.FILE

2385 For notational convenience, this use case will use the term “file” to apply to either a file  
2386 in the conventional sense which is accessed through the NVM.PM.FILE interface, or a  
2387 specific subset of memory ranges residing on an NVM device which are accessed  
2388 through the NVM.BLOCK interface.

2389 **Inputs:**

- 2390 • A set of changes to the persistent state to be applied as a single transaction.
- 2391 • The data and log files.

2392 **Outputs:**

- 2393 • An indication of transaction commit or abort.

2394 **Postconditions:**  
2395     • If an abort indication was returned, the data was not committed and the previous  
2396         contents have not been modified.  
2397     • If a commit indication was returned, the data has been entirely committed.  
2398     • After a system crash, reset, or power failure followed by system restart and  
2399         execution of the application transaction recovery process, the data has either  
2400         been entirely committed or the previous contents have not been modified.

2401 **Success Scenario:**

2402 The application transaction logic uses a log file in combination with its data file to  
2403 atomically update the persistent state of the application. The log may implement a  
2404 before-image log or a write-ahead log. The application transaction logic should  
2405 configure itself to handle torn or interrupted writes to the log or data files.

2406 Since persistent memory may be byte-granular, torn writes may occur at any point  
2407 during a series of stores. The application should be prepared to detect a torn write of  
2408 the record and either discard or recover such a torn record during the recovery process.  
2409 One common way of detecting such a torn write is for the application to compute a hash  
2410 of the record and record the hash in the record. Upon reading the record, the application  
2411 re-computes the hash and compares it with the recorded hash; if they do not match, the  
2412 record has been torn.

2413 **10.4.5.1 NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMICITY is true**

2414 If the NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMICITY is true, then writes which  
2415 are interrupted by a system crash, system reset, or power failure occur atomically. In  
2416 other words, upon restart the contents of persistent memory reflect either the state  
2417 before the store or the state after the completed store.

2418 In this case, the application need not handle interrupted writes to the log or data files.

2419 **10.4.5.2 NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMICITY is false**

2420 NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMICITY is false, then writes which are  
2421 interrupted by a system crash, system reset, or power failure do not occur atomically. In  
2422 other words, upon restart the contents of persistent memory may be such that  
2423 subsequent loads may create exceptions depending on the value of the  
2424 FUNDAMENTAL\_ERROR\_RANGE attribute.

2425 In this case, the application should be prepared to handle an interrupted write to the log  
2426 or data files.

2427 *10.4.5.2.1 NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE > 0*

2428 If the NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE is greater than zero, the  
2429 application should align the log or data records with the  
2430 NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE and pad the record size to be an  
2431 integral multiple of NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE. This prevents



2432 more than one record from residing in the same fundamental error range. The  
2433 application should be prepared to discard or recover the record if a load returns an  
2434 exception when subsequently reading the record during the recovery process. (See also  
2435 SQLite.org, *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>.)

2436 *10.4.5.2.2 NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE = 0*

2437 If the NVM.PM.FILE.FUNDAMENTAL\_ERROR\_RANGE is zero, the application lacks  
2438 sufficient information to handle interrupted writes to the log or data files.

2439 **Failure Scenarios:**

2440 Consider the recovery of an error resulting from an interrupted write on a persistent  
2441 memory volume or file system where the  
2442 NVM.PM.FILE.INTERRUPTED\_STORE\_ATOMICITY is false. This error may be  
2443 persistent and may be returned whenever the affected fundamental error range is read.  
2444 To repair this error, the application should be prepared to overwrite such a range.

2445 One common way of ensuring that the application will overwrite a range is by assigning  
2446 it to the set of internal free space managed by the application, which is never read and  
2447 is available to be allocated and overwritten at some point in the future. For example, the  
2448 range may be part of a circular log. If the range is marked as free, the transaction log  
2449 logic will eventually allocate and overwrite that range as records are written to the log.

2450 Another common way is to record either a before-image or after-image of a data range  
2451 in a log. During recovery after a system crash, system reset, or power failure, the  
2452 application replays the records in the log and overwrites the data range with either the  
2453 before-image contents or the after-image contents.

2454 **See also:**

- 2455 • SQLite.org, *Atomic Commit in SQLite*, <http://www.sqlite.org/atomiccommit.html>
- 2456 • SQLite.org, *Powersafe Overwrite*, <http://www.sqlite.org/psow.html>
- 2457 • SQLite.org, *Write-Ahead Logging*, <http://www.sqlite.org/wal.html>

## 2458 **Annex A (Informative) NVM pointers**

2459 Pointers are data types that hold virtual addresses of data in memory. When  
2460 applications use pointers with volatile memory, the value of the pointer must be re-  
2461 assigned each time the program is run (a consequence of the memory being volatile).  
2462 When applications map a file (or a portion of a file) residing in persistent memory to  
2463 virtual addresses, it may or may not be assigned the same virtual address. If not, then  
2464 pointers to values in that mapped memory will not reference the same data. There are  
2465 several possible solutions to this problem:

- 2466 1) Relative pointers
- 2467 2) Regions are mapped at fixed addresses
- 2468 3) Pointers are relocated when region is remapped

2469 All three approaches are problematic, and involve different challenges that have not  
2470 been fully addressed.

2471 None, except perhaps the third one, handles C++ vtable pointers inside persistent  
2472 memory, or pointers to string constants, where the string physically resides in the  
2473 executable, and not the memory-mapped file. Both of those issues are common.

2474 Option (1) implies that no existing pointer-containing library data structures can be  
2475 stored in NVM, since pointer representations change. Option (2) requires careful  
2476 management of virtual addresses to ensure that memory-mapped files that may need to  
2477 be accessed simultaneously are not assigned to the same address. It may also limit  
2478 address space layout randomization. Option (3) presents challenges in, for example, a  
2479 C language environment in which pointers may not be unambiguously identifiable, and  
2480 where they may serve as hash table indices or the like. Pointer relocation would  
2481 invalidate such hash tables. It may be significantly easier in the context of a Java-like  
2482 language.

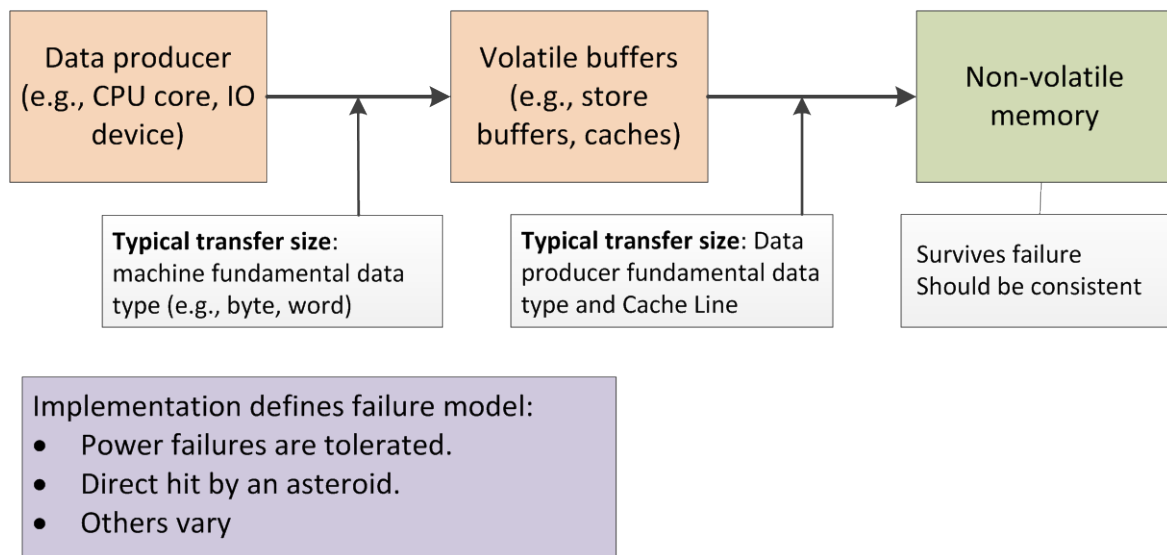
2483 **Annex B (Informative) Consistency**

2484 Persistent memory as defined in the NVM.PM.VOLUME and NVM.PM.FILE modes of  
2485 the SNIA NVM Programming Model must exhibit certain data consistency properties in  
2486 order to enable applications to operate correctly. Directly mapped load/store accessible  
2487 NVM must have the following properties.

- 2488 • Usable as ordinary (not just Durable) memory
- 2489 • Consistent and Durable even after a failure

2490 The context for providing these properties is illustrated in Figure 15 Consistency  
2491 overview.

2492 **Figure 15 Consistency overview**



2493

2494 For the purpose of this annex, CPU's and I/O devices are data producers whose write  
2495 data passes through volatile buffers on the way to non-volatile memory. Generally  
2496 transmission involves a CPU memory bus where typical transfer sizes are determined  
2497 by either a fundamental data type of the machine or the CPU's cache line size. All  
2498 references to cache in this annex refer to the CPU's cache. The Nonvolatile memory  
2499 block illustrated here is referred to in the NVM Programming Model specification as a  
2500 "Persistence Domain".

2501 Implementations vary as to the level of fault tolerance although in general power failures  
2502 must be tolerated but site scale catastrophes are not addressed at this level of the  
2503 system.

2504 When persistence behavior is ignored, memory mapped NVM must appear to operate  
2505 like ordinary memory. Normally compiled code without durability expectations should  
2506 continue to run correctly. This includes the following.

- 2507 • Accessible through load/store/atomic read/modify/write.
- 2508 • Subject to existing processor cache coherency and "uncacheable" models

- 2509 • Load/store/atomic read/modify/write retain their current semantics.
- 2510     ○ Even when accessed from multiple threads.
- 2511     ○ Even if locks or lock-protected data lives in NVM.
- 2512 • Able to use existing code (e.g., sort function) on NVM data.
- 2513 • Holds for all data producers: CPU and, where relevant, I/O.
- 2514 • “Execute In Place” capability
- 2515 • Supports pointers to NVM data structures

2516 At the implementation level, Fences in thread libraries etc. must have usual semantics,  
 2517 and usual thread visibility rules must be obeyed.

- 2518 In order to be consistent and durable even after failure, two properties are mandatory.
- 2519 • Atomicity. Some stores are all-or-nothing: They can’t be partly visible even after a  
 2520 failure.
  - 2521 • Strict Write ordering relative to NVM.PM.FILE.SYNC ().

2522 For example, consider the following code segment where msync implements  
 2523 NVM.PM.FILE.SYNC:

```

2524     // a, a_end in NVM
2525     a[0] = foo();
2526     msync(&(a[0]), ...);
2527     a_end = 0;
2528     msync(&a_end, ...);
2529     . . .
2530     n = a_end + 1;
2531     a[n] = foo();
2532     msync(&(a[n]), ...);
2533     a_end = n;
2534     msync(&a_end, ...);
  
```

2535 For correctness of this code the following assertions must apply:

- 2536 • a[0 .. a\_end] always contains valid data, even after a failure in this code.
- 2537 • a\_end must be written atomically to NVM, so that the second store to a\_end occurs  
 2538 no earlier than the store to a[n].

2539 To achieve failure atomicity, Aligned operations on fundamental data types reach NVM  
 2540 atomically. After a failure (allowed by the failure model), each such store is fully  
 2541 reflected in the resulting NVM state or not at all. Formally “aligned operations on  
 2542 fundamental data types” is implementation defined. These are usually exactly the same  
 2543 operations that under normal operation become visible to other threads/data producers  
 2544 atomically. They are already well-defined for most settings:

- 2545 • Instruction Set Architectures already define them.
- 2546     ○ E.g., for x86, MOV instructions with naturally aligned operands of at most 64  
 2547 bits qualify.

- 2548 • They're generated by known high-level language constructs, e.g.:
- 2549     ○ C++11 lock-free atomic<T>, C11 `_Atomic(T)`, Java & C# volatile, OpenMP  
2550     atomic directives.
- 2551 The fundamental data types that enable atomicity generally fit within CPU cache lines.
- 2552 At least two facilities may be useful to achieve strict ordering relative to `msync()`:
- 2553 • `msync`: Wait for all writes in a range to complete.
- 2554 • optimization using an intra-cache-line ordering guarantee.
- 2555 To elaborate on these, `msync(address_range)` must ensure that if any effects from code  
2556 following the call are visible, then so are all NVM stores to the `address_range` preceding  
2557 it. While this low-level least-common-denominator primitive can be used to implement  
2558 logging, etc., high level code often does not know what address range needs to be  
2559 flushed.
- 2560 Intra-cache-line ordering requires that “thread-ordered” stores to a single cache line  
2561 become visible in NVM in the order in which they are issued. The term “thread-ordered”  
2562 refers to certain stores that are already known in today’s implementations to reach  
2563 coherent cache in order, such as the following.
- 2564 • E.g., x86 MOV
- 2565 • some C11, C++11 atomic stores
- 2566 • Java & C# volatile stores.
- 2567 The CPU core and compiler may not reorder these. Within a single cache line, this order  
2568 is normally preserved when the lines are evicted to NVM. This last point is a critical  
2569 consideration as the preservation of thread-ordered stores during eviction to NVM is  
2570 sometimes not guaranteed for posted I/O.
- 2571 Posted I/O (or Store) refers to any write I/O (or Store) that returns a commitment to the  
2572 application before the data has reached a persistence domain. Posted I/O does not  
2573 provide any ordering guarantee in the face of failures. Write-back caching is analogous  
2574 to posted I/O. On some architectures write-through caching may preserve thread  
2575 ordering during eviction if constrained to a single path to a single persistence domain. If  
2576 better performance than write-through caching is desired or if consistency is mandatory  
2577 over multiple paths and/or multiple persistence domains, then additional mechanisms  
2578 such as synchronous snapshots or write-ahead logging must be used.
- 2579 PCIe has the synchronization primitives for software to determine that memory mapped  
2580 writes make it to a persistence domain. Specifically, Reads and Writes cannot pass  
2581 (other) Writes
- 2582 • within the same logical channel
- 2583 • ...and if relaxed ordering is not enabled

2584 The following pseudo-code illustrates a means of creating non-posted store behavior  
2585 with write-back cacheable persistent memory on PCIe.

```
2586     For a list of sequentially executed writes which don't individually  
2587     cross a cache line {  
2588         Execute the first write in the list  
2589         for the remaining list of writes {  
2590             if the write is not the same cache line as the previous write{  
2591                 flush previous cache line  
2592             }  
2593             Execute the new write  
2594         }  
2595         Flush previous cache line;  
2596         If commit is desired {  
2597             read an uncacheable location from the PCIe device  
2598         }  
2599     }
```

2600 This example depends on the strict ordering of writes within a cache line, which is  
2601 characteristic of current processor architectures. It is not clear whether this is viewed as  
2602 a hard requirement by processor vendors. Other approaches to guaranteed ordering in  
2603 persistence domains exist, some of which are specific to hardware implementations  
2604 other than PCIe.

## 2605 **Annex C (Informative) PM error handling**

2606 Persistent memory error handling for NVM.PM.FILE.MAP ranges is unique because  
2607 unlike I/O, there is no acknowledgement to a load or store instruction. Instead  
2608 processors equipped to detect memory access failures respond with machine checks. In  
2609 some cases these can be routed to user threads as asynchronous events.

2610 This annex only describes the handling of errors resulting from load instructions that  
2611 access memory. As will be described later in this annex, no application level recovery is  
2612 enabled at the point of a store error. These errors should appear in CPU event logs,  
2613 leading to an administrative response that is outside the scope of this annex.

2614 Applications needing timely assurance that recently stored data is recoverable should  
2615 use the NVM.PM.FILE.OPTIMIZED\_FLUSH\_AND\_VERIFY (see 10.2.7) action to read  
2616 data back from NVM after it is flushed. Errors during verify are handled in the manner  
2617 described in this annex.

2618 There are several scenarios that can arise in the handling of machine checks related to  
2619 persistent memory errors while reading data from memory locations such as can occur  
2620 during “load” instructions. Concepts are introduced here in an attempt to advance the  
2621 state of the art in persistent memory error handling. The goal is to provide error  
2622 reporting and recovery capability to applications that is equivalent to the current practice  
2623 for I/O.

2624 We need several definitions to assist in reasoning about asynchronous events.

- 2625 - Machine check: an interrupt. In this case interrupts that result from memory errors  
2626 are of specific interest.
- 2627 - Precise machine check – an interrupt that allows an application to resume at the  
2628 interrupted instruction
- 2629 - Error containment – this is an indication of how well the system can determine the  
2630 extent of an error. This enables a range of memory affected by an error that caused  
2631 an interrupt to be returned to the application.
- 2632 - Real time error recovery – This refers to scenarios in which the application can  
2633 continue execution after an error as opposed to being restarted.
- 2634 - Asynchronous event handler – This refers to code provided by an application that  
2635 runs in response to an asynchronous event, in this case an event that indicates a  
2636 memory error. An application’s event handler uses information about the error to  
2637 determine whether execution can safely continue from within the application or  
2638 whether a partial or full restart of the application is required to recover from the error.

2639 The ability to handle persistent memory errors depends on the capability of the  
2640 processor and memory system. It is useful to categorize error handling capability into  
2641 three levels:

- 2642 - No memory error detection – the lowest end systems have little or no memory error  
2643 detection or correction capability such as ECC, CRC or parity.

- 2644 - Non-precise or uncontained memory error detection – these systems detect memory  
2645 errors but they do not provide information about the location of the error and/or fail to  
2646 offer enough information to resume execution from the interrupted instruction.  
2647 - Precise, contained memory error detection – these systems detect memory errors  
2648 and report their locations in real time. These systems are also able to contain many  
2649 errors more effectively. This increases the range of errors that allowing applications  
2650 to continue execution rather than resetting the application or the whole system. This  
2651 capability is common when using higher RAS processors.

2652 Only the last category of systems can, with appropriate operating system software  
2653 enhancement, meet the error reporting goal stated above. The other two categories of  
2654 systems risk scenarios where persistent memory errors are forced to repeatedly reset  
2655 threads or processors, rendering them incapable of doing work. Unrecovered persistent  
2656 memory errors are more problematic than volatile memory errors because they are less  
2657 likely to disappear during a processor reset or application restart.

2658 Systems with precise memory error detection capability can experience a range of  
2659 scenarios depending on the nature of the error. These can be categorized into three  
2660 types.

- 2661 - Platform can't capture error  
2662     • Perhaps application or operating system dies  
2663     • Perhaps hardware product include diagnostic utilities  
2664 - Platform can capture error, considered fatal  
2665     • Operating system crashes  
2666     • Address info potentially stored by operating system or hardware/firmware  
2667     • Application could use info on restart  
2668 - Platform can capture error & deliver to application  
2669     • Reported to application using asynchronous "event"  
2670     • Example: SIGBUS on UNIX w/address info

2671 If the platform can't capture the error then no real time recovery is possible. The system  
2672 may function intermittently or not at all until diagnostics can expose the problem. The  
2673 same thing happens whether the platform lacks memory error detection capability or the  
2674 platform has the capability but was unable to use it due to a low probability error  
2675 scenario.

2676 If the platform can capture the error but it is fatal then real time recovery is not possible,  
2677 however then the system may make information about the error available after system  
2678 or application restart. For this scenario, actions are proposed below to obtain error  
2679 descriptions.

2680 If the platform can deliver the error to the application then real time recovery may be  
2681 possible. An action is proposed below to represent the means that the application uses  
2682 to obtain error information immediately after the failure.

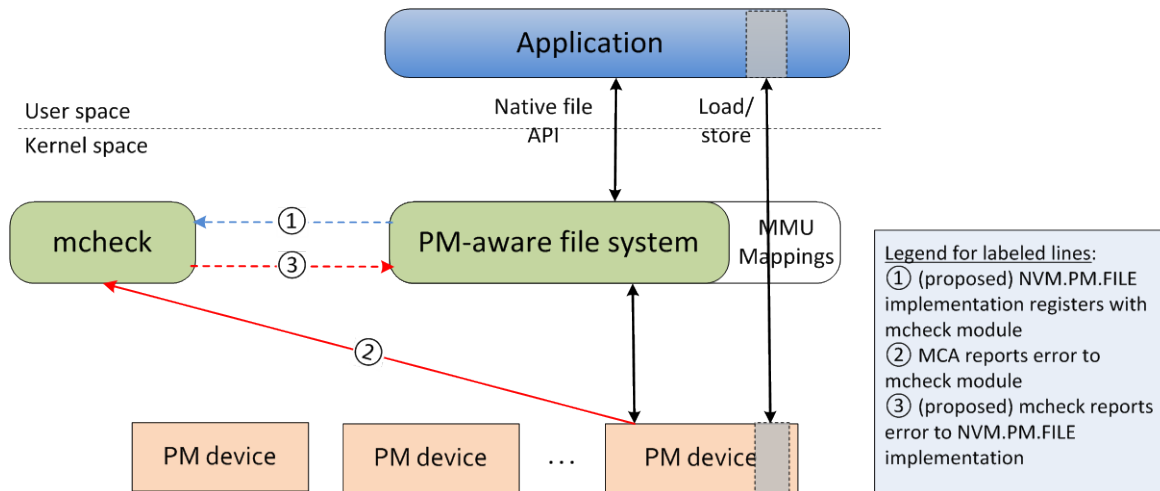


2683 As stated at the beginning of this annex, only errors during load are addressed by this  
 2684 annex. As with other storage media, little or no error checking occurs during store  
 2685 instructions (aka writes). In addition, memory write pipelines within CPU's effectively  
 2686 preclude error handling during memory accesses that result from store instructions. For  
 2687 example, errors detected during the process of flushing the CPU's write pipeline are  
 2688 more likely to be associated with that pipeline than the NVM itself. Errors that arise  
 2689 within the CPU's write pipeline are generally not contained so no application level  
 2690 recovery is enabled at the point of the error.

2691 Continuing to analyze the real time error delivery scenario, the handling of errors on  
 2692 load instructions is sufficient in today's high RAS systems to avoid the consumption of  
 2693 erroneous data by the application. Several enhancements are required to meet the goal  
 2694 of I/O-like application recoverability.

2695 Using Linux running on the Intel architecture as an example, memory errors are  
 2696 reported using Intel's Machine Check Architecture (MCA). When the operating system  
 2697 enables this feature, the error flow on an uncorrectable error is shown by the solid red  
 2698 arrow (labeled ②) in Figure 16 Linux Machine Check error flow with proposed new  
 2699 interface, which depicts the mcheck component getting notified when the bad location in  
 2700 PM is accessed.

2701 **Figure 16 Linux Machine Check error flow with proposed new interface**



2702 As mentioned above, sending the application a SIGBUS (a type of asynchronous event)  
 2703 allows the application to decide what to do. However, in this case, remember that the  
 2704 NVM.PM.FILE manages the PM and that the location being accessed is part of a file on  
 2705 that file system. So even if the application gets a signal preventing it from using  
 2706 corrupted data, a method for recovering from this situation must be provided. A system  
 2707 administrator may try to back up rest of the data in the file system before replacing the  
 2708 faulty PM, but with the error mechanism we've described so far, the backup application  
 2709 would be sent a SIGBUS every time it touched the bad location. What is needed in this  
 2710 case is a way for the NVM.PM.FILE implementation to be notified of the error so it can  
 2711 isolate the affected PM locations and then continue to provide access to the rest of the  
 2712

2713 PM file system. The dashed arrows in the figure above show the necessary modification  
2714 to the machine check code in Linux. On start-up, the NVM.PM.FILE implementation  
2715 registers with the machine code to show it has responsibility for certain ranges of PM.  
2716 Later, when the error occurs, NVM.PM.FILE gets called back by the mcheck component  
2717 and has a chance to handle the error.

2718 This suggested machine check flow change enables the file system to participate in  
2719 recovery while not eliminating the ability to signal the error to the application. The  
2720 application view of errors not corrected by the file system depends on whether the error  
2721 handling was precise and contained. Imprecise error handling precludes resumption of  
2722 the application, in which case the one recovery method available besides restart is a  
2723 non-local go-to. This resumes execution at an application error handling routine which,  
2724 depending on the design of the application, may be able to recover from the error  
2725 without resuming from the point in the code that was interrupted.

2726 Taking all of this into account, the proposed application view of persistent memory  
2727 errors is as described by the NVM.PM.FILE.MAP action (section 10.2.3) and the  
2728 NVM.PM.FILE.GET\_ERROR\_INFO action (section 10.2.6).

2729 The following actions have been proposed to provide the application with the means  
2730 necessary to obtain error information after a fatal error.

- 2731 • PM.FILE.ERROR\_CHECK(file, offset, length): Discover if range has any outstanding  
2732 errors. Returns a list of errors referenced by file and offset.
- 2733 • PM.FILE.ERROR\_CLEAR(file, offset, length): Reset error state (and data) for a  
2734 range: may not succeed

2735 The following attributes have been proposed to enable application to discover the error  
2736 reporting capabilities of the implementation.

- 2737 • NVM.PM.FILE.ERROR\_CHECK\_CAPABLE - System supports asking if range is in  
2738 error state

2739 **Annex D (Informative) Deferred behavior**  
2740 This annex lists some behaviors that are being considered for future specifications.

2741 **D.1 Remote sharing of NVM**  
2742 This version of the specification talks about the relationship between DMA and  
2743 persistent memory (see 6.6 Interaction with I/O devices) which should enable a network  
2744 device to access NVM devices. But no comprehensive approach to remote share of  
2745 NVM is addressed in this version of the specification.

2746 **D.2 MAP\_CACHED OPTION FOR NVM.PM.FILE.MAP**  
2747 This would enable memory mapped ranges to be either cached or uncached by the  
2748 CPU.

2749 **D.3 NVM.PM.FILE.DURABLE.STORE**  
2750 This might imply that through this action things become durable and visible at the same  
2751 time, or not visible until it is durable. Is there a special case for atomic write that, by the  
2752 time the operation completes, it is both visible and durable? The prospective use case is  
2753 an opportunity for someone with a hardware implementation that does not require  
2754 separation of store and sync. This is not envisioned as the same as a file system write.  
2755 It still implies a size of the store. The use case for NVM.FILE.DURABLE.STORE is to  
2756 force access to the persistence domain.

2757 **D.4 Enhanced NVM.PM.FILE.WRITE**  
2758 Add an NVM.PM.FILE.WRITE action where the only content describes error handling.

2759 **D.5 Management-only behavior**  
2760 Several management-only behaviors have been discussed, but deferred to a future  
2761 revision; including:

- 2762 • Secure Erase
- 2763 • Behavior enabling management application to discover PM devices (and  
2764 behavior to fill gaps in the discovery of block NVM attributes)
- 2765 • Attribute exposing flash erase block size for management of disk partitions

2766 **D.6 Access hints**  
2767 Allow applications to suggest how data is placed on storage

2768 **D.7 Multi-device atomic multi-write action**  
2769 Perform an atomic write to multiple extents in different devices.

2770 **D.8 NVM.BLOCK.DISCARD\_IF\_YOU\_MUST action**  
2771 The text below was partially developed, before being deferred to a future revision.

2772 10.4.6 **NVM.BLOCK.DISCARD\_IF\_YOU\_MUST**

2773 Proposed new name MARK\_DISCARDABLE

2774 Purpose - discard blocks to prevent write amplification

2775 This action notifies the NVM device that some or all of the blocks which constitute a  
2776 volume are no longer needed by the application, but the NVM device should defer  
2777 changes to the blocks as long as possible. This action is a hint to the device.

2778 If the data has been retained, a subsequent read shall return “success” along with the  
2779 data. Otherwise, it shall return an error indicating the data does not exist (and the data  
2780 buffer area for that block is undefined).

2781 Inputs: a range of blocks (starting LBA and length in logical blocks)

2782 Status: Success indicates the request is accepted but not necessarily acted upon.

2783 Existing implementations of TRIM may work this way.

2784 10.4.7 **DISCARD\_IF\_YOU\_MUST use case**

2785 **Purpose/triggers:**

2786 An NVM device may allocate blocks of storage from a common pool of storage. The  
2787 device may also allocate storage through a thin provisioning mechanism. In each of  
2788 these cases, it is useful to provide a mechanism which allows an application or NVM  
2789 user to notify the NVM storage system that some or all of the blocks which constitute  
2790 the volume are no longer needed by the application. This allows the NVM device to  
2791 return the memory allocated for the unused blocks to the free memory pool and make  
2792 the unused blocks available for other consumers to use.

2793 DISCARD\_IF\_YOU\_MUST operation informs the NVM device that that the specified  
2794 blocks are no longer required. DISCARD\_IF\_YOU\_MUST instructs the NVM device to  
2795 release previously allocated blocks to the NVM device’s free memory pool. The NVM  
2796 device releases the used memory to the free storage pool based on the specific  
2797 implementation of that device. If the device cannot release the specified blocks, the  
2798 DISCARD\_IF\_YOU\_MUST operation returns an error.

2799 **Scope/context:**

2800 This use case describes the capabilities of an NVM device that the NVM consumer can  
2801 determine.

2802 **Inputs:**

2803 The range to be freed.

2804 **Success scenario:**

2805 The operation succeeds unless an invalid region is specified or the NVM device is  
2806 unable to free the specified region.

2807 **Outputs:**  
2808 The completion status.

2809 **Postconditions:**  
2810 The specified region is erased and released to the free storage pool.

2811 **See also:**  
2812 DISCARD\_IF\_YOU\_CAN

2813 EXISTS

## 2814 **D.9 Atomic write action with Isolation**

2815 Offer alternatives to ATOMIC\_WRITE and ATOMIC\_MULTIWRITE that also include  
2816 isolation with respect to other atomic write actions. Issues to consider include whether  
2817 order is required, whether isolation applies across multiple paths, and how isolation  
2818 applies to file mapped I/O.

## 2819 **D.10 Atomic Sync/Flush action for PM**

2820 The goal is a mechanism analogous to atomic writes for persistent memory. Since  
2821 stored memory may be implicitly flushed by a file system, defining this mechanism may  
2822 be more complex than simply defining an action.

## 2823 **D.11 Hardware-assisted verify**

2824 Future PM device implementations may provide a capability to perform the verify step of  
2825 OPTIMIZED\_FLUSH\_AND\_VERIFY without requiring an explicit load instruction. This  
2826 capability may require the addition of actions and attributes in NVM.PM.VOLUME mode;  
2827 this change is deferred until we have examples of this type of device.