



**LEGAL NOTICE:**

© **Copyright 2008 - 2022 NVM Express®, Inc. ALL RIGHTS RESERVED.**

This Technical Proposal is proprietary to the NVM Express, Inc. (also referred to as “Company”) and/or its successors and assigns.

**NOTICE TO USERS WHO ARE NVM EXPRESS, INC. MEMBERS:** Members of NVM Express, Inc. have the right to use and implement this Technical Proposal subject, however, to the Member’s continued compliance with the Company’s Intellectual Property Policy and Bylaws and the Member’s Participation Agreement.

**NOTICE TO NON-MEMBERS OF NVM EXPRESS, INC.:** If you are not a Member of NVM Express, Inc. and you have obtained a copy of this document, you only have a right to review this document or make reference to or cite this document. Any such references or citations to this document must acknowledge NVM Express, Inc. copyright ownership of this document. The proper copyright citation or reference is as follows: “© 2008 - 2022 NVM Express, Inc. ALL RIGHTS RESERVED.” When making any such citations or references to this document you are not permitted to revise, alter, modify, make any derivatives of, or otherwise amend the referenced portion of this document in any way without the prior express written permission of NVM Express, Inc. Nothing contained in this document shall be deemed as granting you any kind of license to implement or use this document or the specification described therein, or any of its contents, either expressly or impliedly, or to any intellectual property owned or controlled by NVM Express, Inc., including, without limitation, any trademarks of NVM Express, Inc.

**LEGAL DISCLAIMER:**

THIS DOCUMENT AND THE INFORMATION CONTAINED HEREIN IS PROVIDED ON AN “AS IS” BASIS. TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, NVM EXPRESS, INC. (ALONG WITH THE CONTRIBUTORS TO THIS DOCUMENT) HEREBY DISCLAIM ALL REPRESENTATIONS, WARRANTIES AND/OR COVENANTS, EITHER EXPRESS OR IMPLIED, STATUTORY OR AT COMMON LAW, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, TITLE, VALIDITY, AND/OR NONINFRINGEMENT.

All product names, trademarks, registered trademarks, and/or servicemarks may be claimed as the property of their respective owners.

The NVM Express® design mark is a registered trademark of NVM Express, Inc.

NVM Express Workgroup  
c/o VTM, Inc.  
3855 SW 153<sup>rd</sup> Drive  
Beaverton, OR 97003  
USA  
[info@nvmexpress.org](mailto:info@nvmexpress.org)

## NVM Express™ Technical Proposal (TP)

|                                   |   |
|-----------------------------------|---|
| <b>Technical Proposal ID</b>      | <b>4116 Optimal Read Size and Granularity</b>           |
| <b>Change Date</b>                | <b>2022-03-02</b>                                       |
| <b>Builds on Specification(s)</b> | <b>NVM Express Command Set Specification 1.0a</b>       |
| <b>References</b>                 | <b>NVMe – TP 4090 Enhanced Deallocation Granularity</b> |

## Technical Proposal Author(s)

| <b>Name</b>     | <b>Company</b> |
|-----------------|----------------|
| Andres Baez     | Intel          |
| Mike Allison    | Intel          |
| Jonathan Hughes | Intel          |

## Technical Proposal Overview

This proposal intends to do the following:

- Define a mechanism for specifying the optimal size and granularity for reading user data.

## Revision History

| <b>Revision Date</b> | <b>Change Description</b>   |
|----------------------|---|
| 2021-04-20           | Initial version   |
| 2021-05-14           | Addressed changes from feedback. <ul style="list-style-type: none"> <li>• Revised number of bytes reserved for new fields to accommodate all values.</li> <li>• Clarified Copy Command behavior for read and writes.</li> <li>• Clarified NOIOB usage for read and writes.</li> <li>• Converted NSFEAT bits list into a table.</li> </ul> |
| 2021-05-17           | Merged with ratified TP4090 that touched some of the same sections.   |
| 2021-05-20           | Updated Read and Write Commands with new definitions from ECN-001 2021.05.20  |
| 2021-05-27           | <ul style="list-style-type: none"> <li>• Explicitly listed all read and write commands in section 5.8.2</li> <li>• Updated the name of the OPTPERF field to clarify it applies to write commands only.</li> <li>• Updated content with newest available NVM Command Set Specification.</li> </ul>   |
| 2021-06-01           | <ul style="list-style-type: none"> <li>• Corrected minor phase 3-related wording.</li> <li>• Added sentence that specifies the relation between write granularity and granularity and read optimization.</li> <li>• Fixed content to match latest NVM Command Set Specification.</li> </ul>   |

| Revision Date | Change Description  |
|---------------|---|
| 2021-06-11    | <ul style="list-style-type: none"> <li>Consolidated list of commands that apply to write and read attributes in one section.</li> <li>Fixed NOIOB field command list to point to the performance section.</li> <li>Clarified the relation between NOIOB and the optional read/write performance attributes</li> <li>Separated the NORS field into 2. One that is write size dependent and one that is not write size dependent.</li> <li>Separated improved I/O performance section into write and read sub-sections</li> </ul> |
| 2021-06-17    | <ul style="list-style-type: none"> <li>Fixed NORDS and NORSI names.</li> </ul>  |
| 2021-06-18    | <ul style="list-style-type: none"> <li>Accepted all changes from the last version.</li> <li>Updated content with the final version of TP4090.</li> <li>Reverted separating NORS field into NORSD and NORSI into a single NORS field.</li> <li>Added sentence to specify how read performance may depending on how the host writes.</li> </ul>   |
| 2021-06-22    | <ul style="list-style-type: none"> <li>Updated content with reworked version of TP4090.</li> <li>Updated statement about read sizes smaller than NORS for clarity in informative section</li> <li>Reworded sentence that specifies how read performance may depending on writes as suggested.</li> </ul>  |
| 2021-06-24    | <ul style="list-style-type: none"> <li>Reworded NPRG usage without NORS adherence sentence for clarity.</li> <li>Accepted all comments</li> </ul>   |
| 2021-09-24    | <ul style="list-style-type: none"> <li>Increased the size of NPRG, NPRA and NORS fields to 32 bits.</li> <li>Updated content with the latest specification</li> </ul>   |
| 2021-09-30    | <ul style="list-style-type: none"> <li>Minor rewording changes for clarity per feedback</li> <li>Updated NPWA and NPRA figures to differentiate compliant vs non-compliant I/O</li> </ul>   |
| 2021-10-08    | <ul style="list-style-type: none"> <li>Fixed NORS wording in informative when a read does not conform to that field.</li> <li>Updated NPWA/NPRA figures using original diagrams</li> </ul>  |
| 2021-10-14    | <ul style="list-style-type: none"> <li>Accepted all changes and removed comments</li> </ul>   |
| 2021-11-03    | <ul style="list-style-type: none"> <li>Updated figures to better show alignment boundary in the lower layer.</li> <li>Addressed file format to comply with the latest template.</li> <li>Updated content with the latest ratified TP4074 and specification</li> </ul>   |
| 2021-11-12    | <ul style="list-style-type: none"> <li>Removed read conformant example in diagram that showed a read smaller than NPRG</li> </ul>   |
| 2021-11-18    | <ul style="list-style-type: none"> <li>Accepted all previous changes and comments.</li> <li>Added more details to the description of changes section.</li> </ul>  |
| 2021-12-02    | <ul style="list-style-type: none"> <li>Edited the description of changes section to remove unnecessary text and section numbers.</li> </ul>   |
| 2022-01-18    | <ul style="list-style-type: none"> <li>Integrated</li> </ul>  |
| 2022-02-27    | <ul style="list-style-type: none"> <li>Replaced NPWA/NPWG and NPRA/NPRG PNG images with Visio objects</li> </ul>  |
| 2022-02-28    | <ul style="list-style-type: none"> <li>Replaced <sup>TM</sup> with ® throughout document for NVM Express</li> </ul>   |
| 2022-03-02    | <ul style="list-style-type: none"> <li>Moved image 143 and 143a below text within section</li> <li>Addressed reviewer comments and editorial changes</li> </ul>   |

## Description for Changes Document for NVM Express Command Set 1.0a

New Features/Feature Enhancements/Required Changes:

- Optimal Read Size and Granularity (Optional)
  - Defined read performance attributes NPRG, NPRA, NORS.
  - Added new sub section in Improving Performance through I/O Size and Alignment Adherence separating section that describes how to use the new attributes to achieve improved read performance.
  - Restructured Improving Performance through I/O Size and Alignment Adherence separating section write and read performance into new sections.
  - References
    - Technical Proposal 4116

### Markup Conventions:

Black: Unchanged (however, hot links are removed)

~~Red Strikethrough~~: Deleted

Blue: New

Blue Highlighted: TBD values, anchors, and links to be inserted in new text.

<Green Bracketed>: Notes to editor

## Description of Specification Changes for NVM Command Set Specification 1.0a (version dated 2021.07.26)

### 4.1.5.1 NVM Command Set Identify Namespace data structure (CNS 00h)

Figure 97: Identify – Identify Namespace Data Structure, NVM Command Set

| Bytes | O/M<br>1        | Description   | Capability<br>Field |                 |      |      |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |
|-------|-----------------|---|---------------------|-----------------|------|------|--|--|--|------|------|------|-------|------|------|-----|----|----|----|----|----|----|-----|-----|-----|-----|----|-----|-----|-----|-----|-----|----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|----|
| 24    | M               | <p><b>Namespace Features (NSFEAT):</b> This field defines features of the namespace.</p> <p><del>Bits 7:6 are reserved.</del></p> <p><del>Bit 5:4 (<b>OPTPERF</b>) indicate support of alignment and granularity attributes of this namespace, as described in the following table:</del></p> <table><tr><th rowspan="2">Value</th><th colspan="6">Field Supported</th></tr><tr><th>NPWG</th><th>NPWA</th><th>NPDG</th><th>NPDGL</th><th>NPDA</th><th>NOWS</th></tr><tr><td>00b</td><td>No</td><td>No</td><td>No</td><td>No</td><td>No</td><td>No</td></tr><tr><td>01b</td><td>Yes</td><td>Yes</td><td>Yes</td><td>No</td><td>Yes</td><td>Yes</td></tr><tr><td>10b</td><td>Yes</td><td>Yes</td><td>No</td><td>Yes</td><td>Yes</td><td>Yes</td></tr><tr><td>11b</td><td>Yes</td><td>Yes</td><td>Yes</td><td>Yes</td><td>Yes</td><td>Yes</td></tr></table> <p><del>The use of these fields by the host for I/O optimization is described in section 5.10.2.</del></p> | Value               | Field Supported |      |      |  |  |  | NPWG | NPWA | NPDG | NPDGL | NPDA | NOWS | 00b | No | No | No | No | No | No | 01b | Yes | Yes | Yes | No | Yes | Yes | 10b | Yes | Yes | No | Yes | Yes | Yes | 11b | Yes | Yes | Yes | Yes | Yes | Yes | No |
| Value | Field Supported |   |                     |                 |      |      |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |
|       | NPWG            | NPWA  | NPDG                | NPDGL           | NPDA | NOWS |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |
| 00b   | No              | No  | No                  | No              | No   | No   |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |
| 01b   | Yes             | Yes   | Yes                 | No              | Yes  | Yes  |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |
| 10b   | Yes             | Yes   | No                  | Yes             | Yes  | Yes  |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |
| 11b   | Yes             | Yes   | Yes                 | Yes             | Yes  | Yes  |  |  |  |      |      |      |       |      |      |     |    |    |    |    |    |    |     |     |     |     |    |     |     |     |     |     |    |     |     |     |     |     |     |     |     |     |     |    |

**Figure 97: Identify – Identify Namespace Data Structure, NVM Command Set**

| Bytes | O/M<br>1  | Description   | Capability<br>Field |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
|-------|---|---|---------------------|-------------|----|---|----|----------|-------|--|----|---|----|--|----|---|----|--|--|
|       |   | <p>Bit 3 (<b>UIDREUSE</b>) This bit is as defined in the UIDREUSE bit in the I/O Command Set Independent Identify Namespace data structure (refer to the I/O Command Set Independent Identify Namespace data structure section in the NVMe Base Specification).</p> <p>Bit 2 (<b>DAE</b>) if set to '1' indicates that the controller supports the Deallocated or Unwritten Logical Block error for this namespace. If cleared to '0', then the controller does not support the Deallocated or Unwritten Logical Block error for this namespace. Refer to section 3.2.3.2.1.</p> <p>Bit 1 (<b>NSABP</b>) if set to '1' indicates that the fields NAWUN, NAWUPF, and NACWU are defined for this namespace and should be used by the host for this namespace instead of the AWUN, AWUPF, and ACWU fields in the Identify Controller data structure. If cleared to '0', then the controller does not support the fields NAWUN, NAWUPF, and NACWU for this namespace. In this case, the host should use the AWUN, AWUPF, and ACWU fields defined in the Identify Controller data structure in the NVMe Base Specification. Refer to section 2.1.4.</p> <p>Bit 0 (<b>THINP</b>) if set to '1' indicates that the namespace supports thin provisioning. If cleared to '0' indicates that thin provisioning is not supported. Refer to section 2.1.1 for details on the usage of this bit.</p> <table><tr><th>Bits</th><th>Description</th></tr><tr><td>07</td><td><b>Optional Read Performance (OPTRPERF):</b> If set to '1' indicates that the NPRG, NPRA, and NORS fields are defined for this namespace and should be used by the host for I/O optimization. If cleared to '0', then the controller does not support the NPRG, NPRA, and NORS fields for this namespace.</td></tr><tr><td>06</td><td>Reserved</td></tr><tr><td>05:04</td><td><b>Optional Write Performance (OPTPERF):</b> Indicates support of alignment and granularity attributes of this namespace, as described in the following table: <b>Figure TBD</b>.</td></tr><tr><td>03</td><td><b>UID Reuse (UIDREUSE):</b> This bit is as defined in the UIDREUSE bit in the I/O Command Set Independent Identify Namespace data structure (refer to the I/O Command Set Independent Identify Namespace data structure section in the NVMe Base Specification).</td></tr><tr><td>02</td><td><b>Deallocated Error (DAE):</b> If set to '1' indicates that the controller supports the Deallocated or Unwritten Logical Block error for this namespace. If cleared to '0', then the controller does not support the Deallocated or Unwritten Logical Block error for this namespace. Refer to section 3.2.3.2.1.</td></tr><tr><td>01</td><td><b>Namespace Supported Atomic Boundary &amp; Power (NSABP):</b> If set to '1' indicates that the fields NAWUN, NAWUPF, and NACWU are defined for this namespace and should be used by the host for this namespace instead of the AWUN, AWUPF, and ACWU fields in the Identify Controller data structure. If cleared to '0', then the controller does not support the fields NAWUN, NAWUPF, and NACWU for this namespace. In this case, the host should use the AWUN, AWUPF, and ACWU fields defined in the Identify Controller data structure in the NVMe Base Specification. Refer to section 2.1.4.</td></tr><tr><td>00</td><td><b>Thin Provisioning (THINP):</b> If set to '1' indicates that the namespace supports thin provisioning. If cleared to '0' indicates that thin provisioning is not supported. Refer to section 2.1.1 for details on the usage of this bit.</td></tr></table> | Bits                | Description | 07 | <b>Optional Read Performance (OPTRPERF):</b> If set to '1' indicates that the NPRG, NPRA, and NORS fields are defined for this namespace and should be used by the host for I/O optimization. If cleared to '0', then the controller does not support the NPRG, NPRA, and NORS fields for this namespace. | 06 | Reserved | 05:04 | <b>Optional Write Performance (OPTPERF):</b> Indicates support of alignment and granularity attributes of this namespace, as described in the following table: <b>Figure TBD</b> . | 03 | <b>UID Reuse (UIDREUSE):</b> This bit is as defined in the UIDREUSE bit in the I/O Command Set Independent Identify Namespace data structure (refer to the I/O Command Set Independent Identify Namespace data structure section in the NVMe Base Specification). | 02 | <b>Deallocated Error (DAE):</b> If set to '1' indicates that the controller supports the Deallocated or Unwritten Logical Block error for this namespace. If cleared to '0', then the controller does not support the Deallocated or Unwritten Logical Block error for this namespace. Refer to section 3.2.3.2.1. | 01 | <b>Namespace Supported Atomic Boundary &amp; Power (NSABP):</b> If set to '1' indicates that the fields NAWUN, NAWUPF, and NACWU are defined for this namespace and should be used by the host for this namespace instead of the AWUN, AWUPF, and ACWU fields in the Identify Controller data structure. If cleared to '0', then the controller does not support the fields NAWUN, NAWUPF, and NACWU for this namespace. In this case, the host should use the AWUN, AWUPF, and ACWU fields defined in the Identify Controller data structure in the NVMe Base Specification. Refer to section 2.1.4. | 00 | <b>Thin Provisioning (THINP):</b> If set to '1' indicates that the namespace supports thin provisioning. If cleared to '0' indicates that thin provisioning is not supported. Refer to section 2.1.1 for details on the usage of this bit. |  |
| Bits  | Description   |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 07    | <b>Optional Read Performance (OPTRPERF):</b> If set to '1' indicates that the NPRG, NPRA, and NORS fields are defined for this namespace and should be used by the host for I/O optimization. If cleared to '0', then the controller does not support the NPRG, NPRA, and NORS fields for this namespace.   |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 06    | Reserved  |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 05:04 | <b>Optional Write Performance (OPTPERF):</b> Indicates support of alignment and granularity attributes of this namespace, as described in the following table: <b>Figure TBD</b> .  |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 03    | <b>UID Reuse (UIDREUSE):</b> This bit is as defined in the UIDREUSE bit in the I/O Command Set Independent Identify Namespace data structure (refer to the I/O Command Set Independent Identify Namespace data structure section in the NVMe Base Specification).   |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 02    | <b>Deallocated Error (DAE):</b> If set to '1' indicates that the controller supports the Deallocated or Unwritten Logical Block error for this namespace. If cleared to '0', then the controller does not support the Deallocated or Unwritten Logical Block error for this namespace. Refer to section 3.2.3.2.1.  |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 01    | <b>Namespace Supported Atomic Boundary &amp; Power (NSABP):</b> If set to '1' indicates that the fields NAWUN, NAWUPF, and NACWU are defined for this namespace and should be used by the host for this namespace instead of the AWUN, AWUPF, and ACWU fields in the Identify Controller data structure. If cleared to '0', then the controller does not support the fields NAWUN, NAWUPF, and NACWU for this namespace. In this case, the host should use the AWUN, AWUPF, and ACWU fields defined in the Identify Controller data structure in the NVMe Base Specification. Refer to section 2.1.4. |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |
| 00    | <b>Thin Provisioning (THINP):</b> If set to '1' indicates that the namespace supports thin provisioning. If cleared to '0' indicates that thin provisioning is not supported. Refer to section 2.1.1 for details on the usage of this bit.  |   |                     |             |    |   |    |          |       |  |    |   |    |  |    |   |    |  |  |

|       |   |   |    |
|-------|---|---|----|
| 47:46 | O | <p><b>Namespace Optimal I/O Boundary (NOIOB):</b> This field indicates the optimal I/O boundary for this namespace. This field is specified in logical blocks. The host should construct <del>R</del>ead and <del>W</del>rite commands that do not cross the I/O boundary to achieve optimal performance. A value of 0h indicates that no optimal I/O boundary is reported.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance</p>   | No |
| 63:48 | O | <p><b>NVM Capacity (NVMCAP):</b> This field indicates the total size of the NVM allocated to this namespace. The value is in bytes. This field shall be supported if the Namespace Management capability (refer to section 5.3) is supported.</p> <p>Note: This field may not correspond to the logical block size multiplied by the Namespace Size field. Due to thin provisioning or other settings (e.g., endurance), this field may be larger or smaller than the product of the logical block size and the Namespace Size reported.</p> <p>If the controller supports Asymmetric Namespace Access Reporting (refer to the CMIC field), and the relationship between the controller and the namespace is in the ANA Inaccessible state (refer to the ANA Inaccessible state section in the NVMe Base Specification) or the ANA Persistent Loss state (refer to the ANA Persistent Loss state section in the NVMe Base Specification), then this field shall be cleared to 0h.</p> | No |
| 65:64 | O | <p><b>Namespace Preferred Write Granularity (NPWG):</b> This field indicates the smallest recommended write granularity in logical blocks for this namespace. This is a 0's based value. If this field is not supported as <del>defined</del>indicated by the OPTPERF field, then this field is reserved.</p> <p>The size indicated by this field should be less than or equal to the size indicated by the Maximum Data Transfer Size (MDTS) field (refer to the NVMe Base Specification), which is specified in units of minimum memory page size. The value of this field may change if the namespace is reformatted. The size should be a multiple of the Namespace Preferred Write Alignment (NPWA) field.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance</p>   | No |
| 67:66 | O | <p><b>Namespace Preferred Write Alignment (NPWA):</b> This field indicates the recommended write alignment in logical blocks for this namespace. This is a 0's based value. If this field is not supported as <del>defined</del>indicated by the OPTPERF field, then this field is reserved.</p> <p>The value of this field may change if the namespace is reformatted.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance</p>   | No |
| 69:68 | O | <p><b>Namespace Preferred Deallocate Granularity (NPDG):</b> This field indicates the recommended granularity in logical blocks for the Dataset Management command with the Attribute – Deallocate bit set to '1' in Dword 11. This is a 0's based value. If this field is not supported as <del>defined</del>indicated by the OPTPERF field, then this field is reserved.</p> <p>The value of this field may change if the namespace is reformatted. The size should be a multiple of the Namespace Preferred Deallocate Alignment (NPDA) field.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance.</p>  | No |
| 71:70 | O | <p><b>Namespace Preferred Deallocate Alignment (NPDA):</b> This field indicates the recommended alignment in logical blocks for the Dataset Management command with the Attribute – Deallocate bit set to '1' in Dword 11. This is a 0's based value. If this field is not supported as <del>defined</del>indicated by the OPTPERF field, then this field is reserved.</p> <p>The value of this field may change if the namespace is reformatted.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance.</p>  | No |

|       |   |  |    |
|-------|---|--|----|
| 73:72 | O | <p><b>Namespace Optimal Write Size (NOWS):</b> This field indicates the size in logical blocks for optimal write performance for this namespace. This is a 0's based value. If this field is not supported as <b>defined</b> indicated by the OPTPERF field, then this field is reserved.</p> <p>If this namespace is associated with an NVM Set and:</p> <ul style="list-style-type: none"> <li>a) this field is supported as defined by the OPTPERF field, then this field shall equal the value indicated by the Optimal Write Size field in the NVM Set Attributes Entry (refer to the Namespace Identification Descriptor in the NVMe Base Specification) for that NVM Set; or</li> <li>b) this field is not supported as defined by the OPTPERF field, then the host should use the Optimal Write Size field in the NVM Set Attributes Entry for that NVM Set for I/O optimization (refer to section 5.8.2).</li> </ul> <p>The size indicated should be less than or equal to Maximum Data Transfer Size (MDTS) that is specified in units of minimum memory page size. The value of this field may change if the namespace is reformatted. The value of this field should be a multiple of the Namespace Preferred Write Granularity (NPWG) field.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance.</p> | No |
|-------|---|--|----|

**Figure TBD: Namespace Alignment and Granularity Attributes**

| Optimal Write Performance Value | Field Supported |      |      |       |      |      |
|---------------------------------|-----------------|------|------|-------|------|------|
|                                 | NPWG            | NPWA | NPDG | NPDGL | NPDA | NOWS |
| 00b                             | No              | No   | No   | No    | No   | No   |
| 01b                             | Yes             | Yes  | Yes  | No    | Yes  | Yes  |
| 10b                             | Yes             | Yes  | No   | Yes   | Yes  | Yes  |
| 11b                             | Yes             | Yes  | Yes  | Yes   | Yes  | Yes  |

The use of these fields by the host for I/O optimization is described in section 5.8.2.

#### 4.1.5.3 I/O Command Set Specific Identify Namespace Data Structure (CNS 05h)

Figure 100 defines the I/O Command Set specific Identify Namespace data structure for the NVM Command Set.

**Figure 100: NVM Command Set I/O Command Set Specific Identify Namespace Data Structure (CSI 00h)**

| Bytes   | O/M 1 | Description   | Capabilities Field |
|---------|-------|---|--------------------|
| 275:272 | O     | <p><b>Namespace Preferred Read Granularity (NPRG):</b> This field is the smallest recommended read granularity in logical blocks for this namespace. This is a 0's based value. If this field is not supported as indicated by the OPTPERF field, then this field is reserved.</p> <p>The size indicated by this field should be less than or equal to the size indicated by the Maximum Data Transfer Size (MDTS) field (refer to the NVMe Base Specification) which is specified in units of minimum memory page size. The value of this field may change if the namespace is reformatted. The size should be a multiple of the Namespace Preferred Read Alignment (NPRA).</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance.</p> | No                 |

**Figure 100: NVM Command Set I/O Command Set Specific Identify Namespace Data Structure (CSI 00h)**

| Bytes  | O/M<br>1 | Description   | Capabilities<br>Field |
|--|----------|---|-----------------------|
| 279:276  | O        | <p><b>Namespace Preferred Read Alignment (NPRA):</b> This field indicates the recommended read alignment in logical blocks for this namespace (refer to section 5.8.2.TBD2).</p> <p>This is a 0's based value. If this field is not supported as indicated by the OPTRPERF field, then this field is reserved.</p> <p>The value of this field may change if the namespace is reformatted.</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance.</p>  | No                    |
| 283:280  | O        | <p><b>Namespace Optimal Read Size (NORS):</b> This field indicates the size in logical blocks for optimal read performance for this namespace. This is a 0's based value. If this field is not supported as indicated by the OPTRPERF field, then this field is reserved.</p> <p>The size indicated should be less than or equal to Maximum Data Transfer Size (MDTS) that is specified in units of minimum memory page size. The value of this field may change if the namespace is reformatted. The value of this field should be a multiple of Namespace Preferred Read Granularity (NPRG).</p> <p>Refer to section 5.8.2 for how this field is utilized to improve performance and endurance.</p> | No                    |
| 4095:268284  | O        | Vendor Specific   | No                    |
| <p>NOTES:<br/>O/M definition: O = Optional, M = Mandatory.</p> |          |   |                       |

## 5.8.2 Improving Performance through I/O Size and Alignment Adherence

NVMe controllers may require constrained I/O sizes and alignments to achieve the full performance potential. There are ~~a number of~~several optional attributes that the controller uses to indicate these recommendations. If hosts do not follow these constraints, then the controller shall function correctly, but performance may be limited.

For best write performance, the host should issue ~~Copy~~, Write command, Write Uncorrectable command, ~~and~~ Write Zeroes command, and the write portion of the Copy commands that specify:

- a) a number of logical blocks that is:
    - a. a multiple of the Namespace Preferred Write Granularity (NPWG) field (refer to Figure 97), if the NPWG field is defined; and
    - b. a multiple of the number of logical blocks indicated by the Stream Write Size (SWS) field (refer to the Streams Directive – Return Parameters Data Structure figure in the NVMe Base Specification), if the Streams Directive is enabled;
- and
- b) a Starting LBA (SLBA) field that is aligned to the Namespace Preferred Write Alignment (NPWA) field (refer to Figure 97), if the NPWA field is defined.

Resolving conflicts between namespace attributes and Streams attributes is described in section 5.8.2.1.



The namespace preferred deallocate granularity is a number of logical blocks that is indicated by both the NPDG field (refer to Figure 97) and the NPDGL field. The NPDGL field is able to represent larger values than the NPDG field (refer to Figure 100). Support for these fields is indicated by the OPTPERF field (refer to Figure 97). If the NPDG field and the NPDGL field are both supported and indicate different values of namespace preferred deallocate granularity, then the host should use the value indicated by the NPDGL field.

The namespace preferred deallocate alignment is a number of logical blocks that is indicated by the NPDA field (refer to Figure 97).

For best performance, the host should issue Dataset Management commands with the Attribute – Deallocate (AD) bit set to ‘1’ that specify:

- a) a Length in Logical Blocks field that is a multiple of the namespace preferred deallocate granularity, if the namespace preferred deallocate granularity is defined; and
- b) a Starting LBA field that is:
  - a. a multiple of the namespace preferred deallocate alignment, if the namespace preferred deallocate alignment is defined; and
  - b. a multiple of the Stream Granularity Size (SGS) field (refer to the Streams Directive – Return Parameters Data Structure figure in the NVMe Base Specification) if the Streams Directive is enabled.

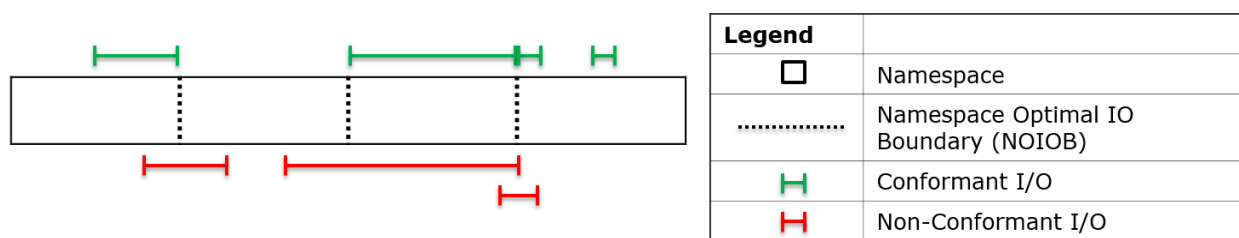
For best read performance, the host should issue Compare command, Read command, Verify command and the read portion of the Copy command that specify:

- a) a number of logical blocks that is a multiple of the Namespace Preferred Read Granularity (NPRG) field (refer to Figure 100), if the NPRG field is supported; and
- b) a Starting LBA (SLBA) field that is aligned to the Namespace Preferred Read Alignment (NPRA) field (refer to Figure 100), if the NPRA field is supported.

### 5.8.2.1 Improved I/O examples (Informative)

It is recommended that the host utilize the I/O attributes as reported by the controller to receive optimal performance from the NVM subsystem. This section summarizes performance related attributes from namespaces, streams, NVM Sets and the NVM command set. The I/O commands discussed throughout this section include those that interact with non-volatile storage in either a Read, Compare, Copy, Verify, Write, Write Uncorrectable, Write Zeroes operation, or Dataset Management operation with the Attribute - Deallocate bit set to ‘1’. The I/O command properties of length and alignment are discussed throughout this section.

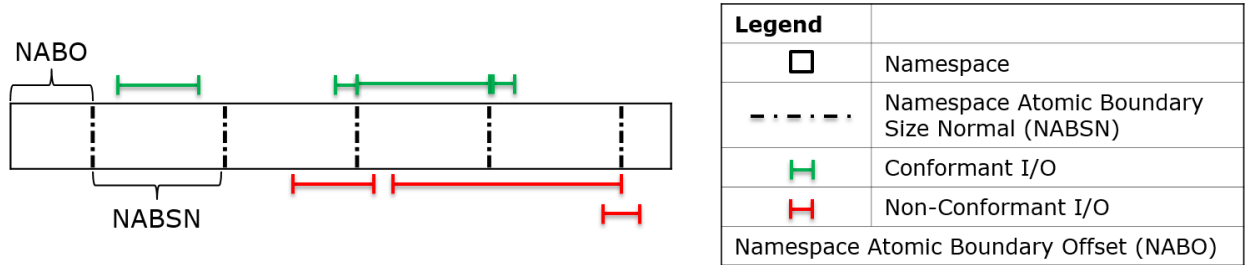
**Figure 141: An example namespace with four NOIOBs**



In Figure 141 an example namespace is diagrammed with three Namespace I/O Boundaries (NOIOB) (refer to Figure 97). The NOIOB attribute should be applied to ~~Read, Compare, Copy, Verify, Write, Write Uncorrectable, and Write Zeroes I/O commands~~ read and write commands as specified in section 5.8.2. An I/O command may see its performance limited if it does not conform to the NOIOB attribute considerations described in this section. The four green lines are example I/O commands from the host that

adhere to the recommendations of NOIOB settings for this namespace. None of the four I/O commands shown in green on the top of Figure 141 cross an NOIOB. The three I/O commands shown in red on the bottom of Figure 141 violate the recommendations for improved performance. The longest I/O command shown in red crosses one NOIOB and ends aligned with a different NOIOB. The remaining two I/O commands shown in red also cross an NOIOB. All three of these example I/O commands shown in red could be split into two I/O commands that adhere to the recommendations provided by the namespace for NOIOB.

Figure 142: Example namespace illustrating a potential NABO and NABSN



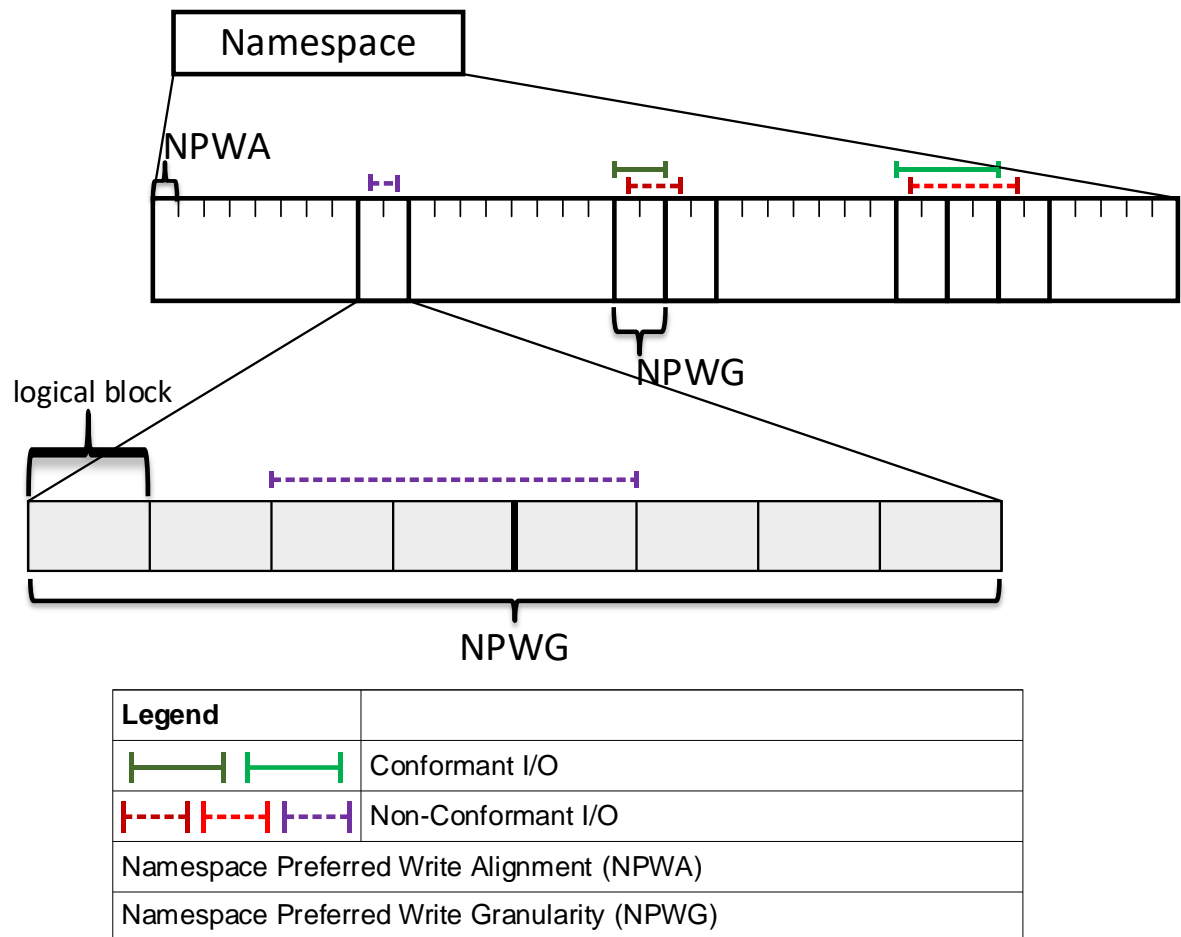
Continuing with the same namespace example from Figure 141, an illustration of Namespace Atomic Boundary Offset (NABO) (refer to Figure 97) and Namespace Atomic Boundary Size Normal (NABSN) (refer to Figure 97) is shown in Figure 142. NABSN and NABO attributes apply to Write, Write Uncorrectable, and Write Zeroes commands. NABSN and NOIOB may not be related to each other, and there may be an offset of NABO to locate the first NABSN starting. The NOIOBs are not shown in Figure 142. The I/O commands shown in green on the top of Figure 142 illustrate I/O commands that adhere to the namespace’s guidance for optimal performance. The I/O commands shown in red on the bottom illustrate I/O commands that do not follow the optimal performance guidelines.

The I/O command examples shown in red in Figure 141 and Figure 142 both illustrate commands that could be restructured to conform to the namespace attributes for Optimal I/O relative to NOIOB, NABO, and NABSN. Each of these example I/O commands shown in red in Figure 141 and Figure 142 could be split into two different I/O commands that adhere to the recommendations. While this increases the number of commands sent to the controller, the expectation is that adherence to the boundary recommendations improves the performance for the controller. Avoiding host traffic that demands non-optimal I/O commands is the most recommendable solution for a host.

5.8.2.TBD1      Alignment and Write Performance

NPWG and NPWA are namespace internal constructs, and they are illustrated in Figure 143. The box at the top of Figure 143 is the namespace. The series of boxes in the middle layer indicate many namespace optimal write units described by NPWA (refer to Figure 97) and NPWG (refer to Figure 97), and the bottom layer is a series of eight logical blocks that in aggregate form the NPWG for this example. Sometimes NPWG are useful because several sequential logical blocks (refer to Figure 97) may be placed and tracked together on the media, or the NPWG may be related to NVM subsystem data reliability implementation constraints.

Figure 143: Example namespace broken down to illustrate potential NPWA and NPWG settings

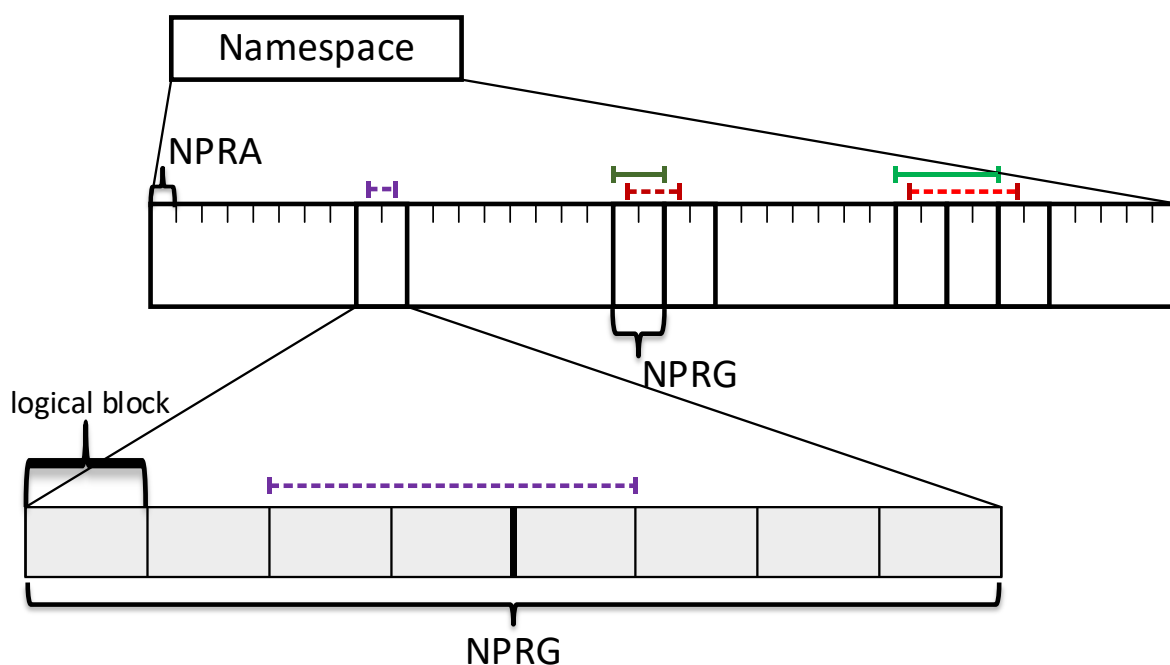


For a list of commands that apply to NPWG and NPWA attributes ~~apply to Copy, Write, Write Uncorrectable, and Write-Zeroes commands.~~ refer to section 5.8.2.

### 5.8.2.TBD2 Alignment and Read Performance

NPWG and NPRA are namespace internal constructs, and they are illustrated in Figure 143a. The box at the top of Figure 143a is the namespace. The series of boxes in the middle layer indicate many namespace optimal read units described by NPRA (refer to Figure 100) and NPWG (refer to Figure 100). In the figure, the read alignment is 4 logical blocks. The bottom layer is a series of eight logical blocks that in aggregate form the NPWG for this example. Sometimes NPWG are useful because several sequential logical blocks (refer to Figure 100) may be placed and tracked together on the media, or the NPWG may be related to NVM subsystem data reliability implementation constraints. Host reads that are of a length less than NPRA may see their performance impacted if they violate read alignment as described in Figure 100.

**Figure 143a:** Example namespace broken down to illustrate potential NPRA and NPRG settings



The host should issue reads that meet the recommendations of NPRG and NPRA and may achieve optimal read performance by issuing reads that meet the recommendation of NORS. If NORS is greater than NPRG, reads that are a multiple of NPRG and not equal to NORS may see improved performance; however, read performance may not be optimal.

Non-adherence to write-related performance attributes (i.e. NPWG, NPWA, NPDG, NPDGL, NPDA, and NOWS), across all the namespaces in:

- the same NVM Set;
- the same Endurance Group when NVM Sets are not supported; or
- the NVM subsystem when Endurance Groups are not supported,

may affect the level of read optimization achievable through the usage of NORS as described in this section.

For a list of commands that apply to NPRG and NPRA attributes refer to [section 5.8.2](#).