# Zoned Storage Models

Version 1.0

ABSTRACT: This SNIA document defines recommended behavior for hardware and software that supports Zoned Storage.

This document has been released and approved by SNIA. SNIA believes that the ideas, methodologies and technologies described in this document accurately represent SNIA goals and are appropriate for widespread distribution. Suggestions for revisions should be directed to https://www.snia.org/feedback/.

## SNIA Standard

July 2, 2023

# USAGE

# DISCLAIMER

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

## Revision History

| Major | Revision | Editor | Date | Comments |
|---|---|---|---|---|
| 0.1 | R0 | Yoni Shternhell | 18-Jul-2022 | Initial draft. |
| 0.1 | R1 | Yoni Shternhell | 19-Jul-2022 | Modification after TWG call |
| 0.1 | R2 | Yoni Shternhell | 25-Jul-22 | Add to the Architecture section each of the attribute used in the 'Characteristics' in each model. |
| 01 | R4 | Yoni Shternhell | 18-Aug-22 | Added text and figured in the Theory of Operation section for block storage devices |
| 0.1 | R7 | Yoni Shternhell | 20-Sep-22 | Integrated Paul S. comments |
| 0.1 | R8 | Yoni Shternhell | 20-Sep-22 | Incorporated comments after the group call |
| 0.1 | R8 | Yoni Shternhell | 26-Sep-22 | Clean version |
| 0.1 | R10 | Yoni Shternhell | 10-Oct-22 | Modified figure 2<br>Added reference to TP4115 and TP4076a |
| 0.1 | R11 | Yoni Shternhell | 26-Oct-22 | Modifications from the SNIA F2F meeting |
| 0.9 | R0 | Yoni Shternhell | 31-Oct-22 | Clean version for Public and Members review |
| 0.9 | R1 | Yoni Shternhell | 16-Nov-22 | Changes from 2022-11-15 meeting of the ZSTWG, to resolve RFC ballot comments.<br>Changes are tracked and comments resolved.<br>Updated date on title page. |
| 0.9 | R2 | Yoni Shternhell | 16-Nov-22 | Clean version |
| 0.9 | R3 | Yoni Shternhell | 07-Feb-23 | Javier comments on sections 4.3.1 (NVM Express Reliability Requirement) 4.3.2 (Offline Zone(s)) |
| 0.9 | R4 | Matias Bjørling | 10-Feb-23 | - Text modifications in section 4.3 (Common Requirements) that eliminates gaps due to write errors.<br>- Delete section 4.3.2 (Offline Zones) as the offline zones are no longer possible. |
| 0.9 | R5 | Yoni Shternhell | 14-Feb-23 | Modified sections 4.3 and 4.3.2 |
| 0.9 | R6 | Yoni Shternhell | 21-Feb-23 | Clean version |
| 0.9 | R7 | Yoni Shternhell | 26-Feb-23 | Incorporated agreed text from the Zone_Size_Proposal_Notes document into section 4 |
| 0.9 | R8 | Mike Allison | 28-Feb-23 | Reformatted. |
| 0.9 | R9 | Yoni Shternhell | 28-Feb-23 | Clean version after WG call for RFC ballot. |
| 0.9 | R10 | Yoni Shternhell | 14-Mar-23 | Incorporated RFC ballot comments |
| 0.9 | R11 | Yoni Shternhell | 21-Mar-23 | Continue Incorporate RFC ballot and Dave L. comments |
| 0.9 | R12 | Yoni Shternhell | 28-Mar-23 | Continue Incorporate RFC ballot and Dave L. comments. We made the following changes:<br>• The device does not support thin provisioning<br>• The device does not support zones in the Offline state. |
| 0.9 | R13 | Paul Suhler | 11-Apr-23 | Corrected bullet items for variations 2 and 3.<br>Removed all comments.<br>Updated revision number, date, table of contents, and table of figures.<br>Accepted all comments and stopped tracking.<br>Ready for forwarding to the TC for public review and thence to the BoD for a membership vote. |

# Table of Contents

# Table of Figures

Zone size of Z LBAs

| LBA X | LBA X+1 | LBA X+2 | ● ● ● | LBA X+Z-3 | LBA X+Z-2 | LBA X+Z-1 |

Zone capacity of Z LBAs that are writeable by the host

Zone size of Z LBAs

| LBA X | LBA X+1 | ● ● ● | LBA X+Y-2 | LBA X+Y-1 | LBA X+Y | LBA X+Y+1 | ● ● ● | LBA X+Z-2 | LBA X+Z-1 |

Zone capacity of Y LBAs that are writeable by the host

Non-writeable but readable LBAs

# FOREWORD

The SNIA Zoned Storage TWG was formed to facilitate a common industry understanding of Zoned Storage use cases, device architectures and programming model, providing a framework to enable the development of a robust Zoned Storage software and hardware ecosystem.

This SNIA standard outlines the architecture and use case models for Zoned Storage devices. As this standard is developed, requirements in interface standards and specific APIs may be proposed as separate documents and developed in the appropriate organizations.

# 1  Scope

This standard defines the requirements and use case models that apply to Zoned Storage devices.

A Zoned Storage device has several aspects:

- **Common Characteristics**. The properties of a Zoned Storage device that host software expects.
- **Security.** The security requirements for a Zoned Storage device.
- **Models**. Describes two Zoned Storage device models that cover known uses cases.

# 2  References

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

For references available from ANSI, contact ANSI Customer Service Department at (212) 642-4980 (phone), (212) 302-1286 (fax) or via the World Wide Web at https://www.ansi.org.

NVM Express® Base Specification
Published specification, available from https://nvmexpress.org

NVM Express® Zoned Namespace Command Set Specification
Published specification, available from https://nvmexpress.org

NVM Express® TP4115 Namespace Management Zoned Namespace Enhancement
Published ratified technical proposal, available from https://nvmexpress.org

NVM Express® TP4076b Zoned Random Write Area
Published ratified technical proposal, available from https://nvmexpress.org

INCITS 537-2016 Information Technology – Zoned Device ATA Command Set (ZAC)
Approved standard, available from https://webstore.ansi.org

INCITS 536-2016 Information Technology – Zoned Block Commands (ZBC)
Approved standard, available from https://webstore.ansi.org

# 3 Definitions, abbreviations, and conventions

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

## 3.1 Definitions

### 3.1.1 zone

A contiguous range of logical block addresses that are managed as a single unit.

### 3.1.2 zoned namespace

An NVMe namespace that is divided into zones, as described in the Zoned Namespace Command Set Specification.

## 3.2 Keywords

In the remainder of this standard, the following keywords are used to indicate text related to compliance:

### 3.2.1 mandatory

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

### 3.2.2 may

a keyword that indicates flexibility of choice with no implied preference; "may" is equivalent to "may or may not"

### 3.2.3 may not

keywords that indicate flexibility of choice with no implied preference; "may not" is equivalent to "may or may not"

### 3.2.4 need not

keywords indicating a feature that is not required to be implemented; "need not" is equivalent to "is not required to"

### 3.2.5 optional

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

### 3.2.6 shall

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

### 3.2.7 should

a keyword indicating flexibility of choice with a strongly preferred alternative

## 3.3 Abbreviations

HDD       Hard Disk Drive

LBA       Logical Block Address

LUN       logical unit number

NVM       non-volatile memory

UBER      uncorrectable bit error rate

SMR       shingled magnetic recording

SSD       Solid State Drive

WAL       Write Ahead Log

ZNS       Zoned Namespace

ZRWA      Zoned Random Write Area

# 4 Theory of Operation

A Zoned Storage device is a block storage device that has its LBA space divided into zones, as shown in Figure 1. A zone is of a certain type, which defines the rules for accessing its LBAs.

For example, one such type is Sequential Write Required, which requires that LBAs within a zone are written in sequential order, but can be read in any order. Each write must begin at the Write Pointer (WP), as shown below.



Figure 1 – Zone Abstraction

Typically, in the standards which define zoned storage, a zone is of a certain type, which defines the rules for accesses to its LBAs. For example, the most common zone type is Sequential Write Required, which requires that LBAs within a zone are written in sequential order, but can be read in any order. In a Sequential Write Required zone, the only way to change the contents of a logical block already written to a zone is to reset the WP (i.e., deleting all the data in the zone and restarting writing from the beginning of the zone). Reading data has no restrictions and the data can be read in the same manner as on traditional storage devices.

In SSDs, the zone abstraction allows a host to align its writes to the required properties of the Zoned Storage device, and thereby optimizes data placement on the device's media. Note that the management of media reliability continues to be the sole responsibility of the Zoned Storage SSD and is managed the same way as for a conventional device.

**Conventional SSD Device** — Flash — Conventional SSD: Device controls data placement

**Zoned SSD Device** — Flash — Zone Zone Zone — Zoned SSD: Applications controls data placement in zones

Figure 2 – Conventional device and Zoned Storage device internal data placement

Shingled magnetic recording (SMR) technology was introduced in Hard Disk Drives (HDDs) to enable increased areal density and larger capacities, and to improve the cost-effectiveness of HDDs. In SMR, unlike conventional recording, tracks are written in an overlapping manner.  This allows tracks to be more tightly packed and hence to achieve a higher recording density. However, once the tracks are overlapped, a logical block within a zone cannot be written independently. To manage the recording, the disk surface is divided into zones with a gap left between zones.  This allows each zone to be written and erased independently. Multiple approaches are possible to manage the recording restriction.

A conventional device handles the recording constraint internally and exposes a conventional interface (i.e., an interface allowing random writing) to the host.  However, in large-scale systems, where performance and space utilization must be carefully managed, the host cannot rely on the device to manage recording.  Therefore, the host is required to manage the performance and space utilization of large storage systems.

Zoned Storage Models

**Conventional**

**Zoned**

Zone

Figure 3 – Internal data placement in conventional and zoned HDDs

The Zoned Storage Device Model is standardized for storage devices as described in:

- ZBC: Zoned Block Commands in T10/SCSI
- ZAC: Zoned Device Command Set in T13/ATA
- ZNS: Zoned Namespace Command Set Specification in NVM Express

## 4.1 Overview

This section provides an overview of the Zoned Storage Model.

A Zoned Storage Model consists of a set of base requirements that applies to all SNIA Zoned Storage Models, followed by an additional set of requirements for a specific Zoned Storage Model. A generic architecture description of the Zoned Storage Model is illustrated in Figure 4.

Figure 4 – An Architectural view of the Zoned Storage Model

The Zoned Storage Model allows the host storage stack to always assume that a Zoned Storage device will support certain properties.

The ATA Command Set and the SCSI Primary Commands standards define interfaces for host software to communicate with storage devices. The Zoned Device ATA Command Set standard and the Zoned Block Commands standard, respectively, define additional functions for Zoned Storage devices.

This standard defines comprehensive requirements that apply to all SNIA Zoned Storage Models.

## 4.2 Characteristics

A Zoned Storage device has several key characteristics that are used by the Zoned Storage Models.

### 4.2.1 Zoned Device Protocol

The Zoned Device Protocol characteristic defines the protocol used by a Zoned Storage device (i.e., NVMe ZNS, T13 ZAC, or T10 ZBC).

### 4.2.2 Zone Type

The Zone Type characteristic defines the rules for reading from and writing to a zone (e.g., a zone type of Sequential Write Required).

### 4.2.3   Zone Capacity

The Zone Capacity characteristic defines the writeable capacity of a zone.

### 4.2.4   Zone Active Resources Available

The Zone Active Resources Available characteristic defines the total number of active resources (allocated and unallocated active resources). This characteristic may only apply to certain Zone Device Protocols.

### 4.2.5   Zone Open Resources Available

The Zone Open Resources Available characteristic defines the total number of open resources (allocated and unallocated open resources). This characteristic may only apply to certain Zoned Device Protocols.

### 4.2.6   Performance

The Performance characteristic describes host access patterns and their effects on performance.

### 4.2.7   Mandatory I/O Access Commands

The Mandatory I/O Access Commands characteristic defines the mandatory commands for a Zoned Storage device (e.g., Read command, Write command, etc.).

### 4.2.8   Mandatory Access Command

The Mandatory Access Command characteristic defines the mandatory commands for an NVMe Zoned Storage device. This characteristic applies only to ZNS SSDs.

### 4.2.9   ZRWA

The optional ZRWA (Zoned Random Write Area) feature defines an area with a set of assigned LBAs which start at the write pointer for a given zone in which the logical blocks that are mapped to that area may be written in random order as well as overwritten. Data flushed from ZRWA to a zone is written sequentially to the zone at the write pointer. The ZRWA feature is defined by the NVM Express Technical Proposal TP4076.

## 4.3 Common Requirements

This section defines the properties of a Zoned Storage device that host software expects.

### 4.3.1 General Requirements

The device shall manage media reliability issues caused by:

- Write errors, if correctable by the device.

- Prematurely worn-out flash blocks associated with a zone (i.e., flash block(s) associated with a zone must not be fixed and should be wear-leveled across zones).
- Read/program disturbs caused by open zones, excessive reads, or similar events.

The number of active and open resources should be equal.

The device does not support thin provisioning.

If the controller is not able to successfully write to all logical blocks specified by a command that initiates a write operation, then the write pointer shall be set to one greater than the last LBA that was written successfully. In the event of internal persistent controller errors that prohibits further writes (i.e., Write failure) that cannot be corrected by the device (i.e., unrecoverable error), the device shall transition the device to a read-only condition. The zoned storage device shall not transition zones to either the Readonly state or the Offline state.

### 4.3.2 NVM Express Zoned Namespace Specific Requirements

The controller shall not exhibit Zone Active Excursions related to Active Zones (i.e., the controller shall not transition open zones to the Full State due to one or more vendor-specific excursion events). Refer to the Zone Active Excursions section in the NVM Express Zoned Namespace Command Set Specification.

The controller shall maintain a fixed number of writeable LBAs (i.e., fixed zone capacity) within a zone from the time a zoned namespace is formatted or created and until the zone namespace is reformatted or deleted (i.e., the controller is not able to change the writeable capacity of a zone between resets). Refer to the Variable Zone Capacity bit in the NVM Express Zoned Namespace Command Set Specification for further information.

### 4.3.3 HDD Specific Requirements

This standard defines no requirements specific to HDDs.

### 4.3.4 Security

This standard defines no security requirements for Zoned Storage devices. Security requirements for general storage devices apply to Zoned Storage devices.

## 4.4 Zone Size and Zone Capacity Considerations

Zone size and zone capacity are properties that are orthogonal to the existing models. To keep the number of models to a minimum, this section describes the trade-offs around the different possibilities.

### 4.4.1 Variable Zone Capacity vs. Fixed Zone Capacity

One of the trade-offs present around zone size and zone capacity is the fact that the capacity of a zone might be allowed to shrink. The Variable Zone Capacity bit determines if the device is allowed to change the capacity of any zone when a zone is reset (i.e., at runtime).

Changing the capacity of a zone by shrinking it is a mechanism to facilitate increasing the lifetime of the device by utilizing portions of the available NAND once it begins to degrade.

This feature comes at the cost of two sources of complexity in the host:

1. The host must deal with variable zone capacities instead of simply using calculations. This requires tracking the capacities of all used zones and implementing the capability to deal with variable application objects. Otherwise, mapping a single object to several zones due to a capacity decrease will incur a host-side write amplification; and
2. The host must deal with failed writes.

Even if the host deals with the excursions, these are non-blocking.

### 4.4.2 Initial Zone Capacity = Zone Size vs. Zone Capacity < Zone Size

The other trade-off related to zone size and zone capacity is whether the initial capacity of a zone must be equal to the zone size.

This is a consequence of some operating systems that support rotational zoned devices (e.g., SMR HDDs) traditionally assuming that zone sizes are powers of two. Note that this is orthogonal to the previous trade-off on variable zone capacity.

If the capacity of a zone is less than the size of that zone, then the host must deal with the following sources of complexity:

1. When the zone capacity is initially less than the zone size, the logical block address space presented to the host contains gaps (refer to Figure 6). These gaps inflate the last valid addresses (e.g., for a namespace with a size of 3TB (i.e., NSZE), there may only be 2TB of writable capacity (i.e., NCAP) for that namespace). This is a departure from how existing storage devices present themselves to the host, which might produce unexpected behavior in existing storage stacks;
2. When these gaps are present, the host read path must be aware of the zone geometry. This requires extra logic in the host read path to deal with valid and invalid addresses as well as with potential I/O splits when a read spans several zones. Such logic may translate into performance degradation;
3. At namespace creation, a host is able to specify the Zone size (using the NSZE field), but does not know the resulting Zone Capacity until the namespace is created; and
4. The inclusion of these gaps is a departure from how existing zoned-aware applications operate.

This creates fragmentation across zoned devices as a function of the media they use and forces application changes. Note that zoned applications have normally used the zone size to iterate across the address space; the benefits of specific alignment for zone sizes emanate from this. However, when the zone capacity is less than the zone size, the host needs to use this new value for the calculations.

### 4.4.3 Zone Size and Zone Capacity Requirements

This section describes the requirements based on the zone size and zone capacity.

#### 4.4.3.1 General requirements

When host software manages zoned block devices, the host software's minimum requirements for writing and reading are as described below:

Zoned devices complying with one of the SNIA Zoned Storage Models handle media reliability issues without involvement of the host. For example, under normal operation zones in the Offline state and write errors are not exposed to the host.

Zone devices handle media reliability issues entirely in the device. This means that under normal operation the host does not have to deal with zones in the Offline state or write errors. When writing to a zoned block device, the host software shall take the following into account:

1. The state of each zone (I.e., Empty, Open/Closed, Full, Readonly/Offline):

    a. To manage free space, host-side GC, and other specifics, the host shall know the state of each zone.

2. Sequential Writes

    a. Hosts shall write sequentially within a zone following the write pointer.

3. Write errors to writeable LBAs may fail due to a partial write failure.

When reading from a zoned block device, the host software shall split reads across zone boundaries into multiple Read commands when read across zone boundaries is not supported.

#### 4.4.3.2 Device Specific Requirements

The zone models are allowed to implement three variations with regards to the relationship between zone size and zone capacity.

**Variation 1:** If the device zone size is a power of two LBAs and the zone capacity is equal to the zone size, then the device:

1. There are no gaps in the LBA address space.

2. The device capacity is managed as in a conventional storage device.

3. If reading across zone boundaries is supported by the device, the host read path does not need to be aware of zone states and reads can be issued to the zoned device as in a conventional storage device.

4. If reading across zone boundaries is not supported by the device, the host read path shall use the zone size to split reads into multiple Read commands so that no Read command crosses zone boundaries.

Note #1: This variation is followed by SMR HDDs.

- It is supported by the mainline zoned storage block device in the Linux Kernel.

- It aligns with existing zoned applications based on SMR HDDs.

**Variation 2:** If the device zone size is not a power of two, and the zone capacity is equal to the zone size, then:

1. The host has no gaps in writing the LBA address space as illustrated in Figure 5.

2. If reading across zone boundaries is supported by the device, the host read path does not need to be aware of zone states and reads can be issued to the zoned device as in a conventional storage device.

3. If reading across zone boundaries is not supported by the device, the host read path shall use the zone size to split reads into multiple Read commands so that no Read command crosses zone boundaries.

Zone size of Z LBAs

| LBA X | LBA X+1 | LBA X+2 | ● ● ● | LBA X+Z-3 | LBA X+Z-2 | LBA X+Z-1 |

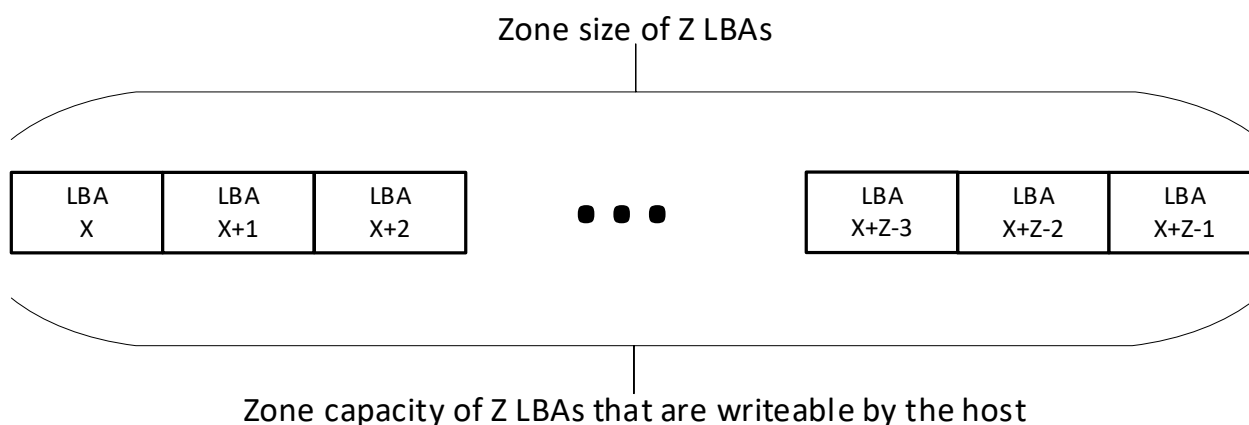Zone capacity of Z LBAs that are writeable by the host

Figure 5 – Zone Capacity is equal to the Zone Size

Note #1: This variation differs from SMR HDDs.

- It is not supported by the mainline zoned storage block device in the Linux Kernel.

- This variation requires changes to existing zoned applications based on SMR HDDs. The changes are related to the fact that zone size is not a power-of-two.

**Variation 3:** If the zone capacity is less than the zone size, then:

1. There are gaps in the LBA address space within the zones that are not writable as illustrated in Figure 6. While the gaps are not writeable, they are readable by the host as deallocated LBAs. The gaps do not contain user data.

2. The host has to manage these gaps.

   a. The host read path in the host is required to be aware of zones as it has to deal with gaps.

   b. This differs from how conventional storage devices manage and report usable capacity to the host. See Note #1 for an example.

3. If reading across zone boundaries is supported by the device, the host read path does not need to be aware of zone states and reads can be issued to the zoned device as in a conventional storage device.

4. If reading across zone boundaries is not supported by the device, the host read path shall use the zone capacity to split reads into multiple Read commands so that no Read command crosses zone boundaries.

Zone size of Z LBAs

| LBA X | LBA X+1 | ••• | LBA X+Y-2 | LBA X+Y-1 | LBA X+Y | LBA X+Y+1 | ••• | LBA X+Z-2 | LBA X+Z-1 |

Zone capacity of Y LBAs that are writeable by the host

Non-writeable but readable LBAs

Figure 6 – Zone Capacity is less than the Zone Size

Note #1: Assume 1TB device with 1GB zone size and 800MB zone capacity. The LBA gaps will consume 256MB worth of capacity per zone. This means that the last LBA will point to a capacity of 1280GB (1TB of usable capacity + 256GB of gaps).

Note #2: Calculations to identify the zone starting LBA uses the zone size; in the case that zone size is a power-of-2, the host can leverage shift operations to make this calculation. The host uses the zone capacity to identify when a zone is full. Calculations to identify the starting LBA of a zone use the zone size; even if the zone size is a power-of-2, requires the host to use division

operations. The host needs to take into account the remaining zone capacity when writing to a zone to avoid a "Full Zone" write error.

Note #3: This variation departs from SMR HDDs.

- It is supported by the mainline zoned storage block device in the Linux Kernel.

- It requires changes to existing zoned applications based on SMR HDDs. The changes are related to the management of zone gaps.

### 4.4.3.3　　Writable Storage Capacity

Because the zoned storage model:

- does not support zones in the Offline state;
- does not support thin provisioning; and
- the zone capacity is able to be less than the zone size,

the device writable storage capacity often needs to be calculated by the host software.

The writeable storage capacity can be derived by the following:

Writable storage capacity = (#TotalZones) * (Zone Capacity)

### 4.4.4　　Eco-System Support

| Type | Zone Configuration | Eco-System Support |
|------|--------------------|--------------------|
| A | Zone Size = Zone Capacity | **Pow2: Supported** in the Linux kernel eco-system since v4.10 (Feb 2017) <br><br> **Non-pow2: Not yet supported** in the Linux kernel eco-system |
| B | Zone Capacity < Zone Size | **Supported** in the Linux kernel eco-system since v5.9 (Oct 2020) |

# 5 Models

## 5.1 Model A

### 5.1.1 Overview

This model describes the requirements for a Zoned Storage device that is a good all-round device.

#### 5.1.1.1 Applicable Use Cases

This Zoned Device Model minimizes the changes required to host software to support zoned block devices.

Host software must respect the sequential write requirement of the zone type, and must reset a zone to rewrite a zone.

Common use cases include, but are not limited to:

- Streaming applications (e.g., sequential writes and random reads).
- Database applications (e.g., the Write Ahead Log (WAL) and log-structured writes).
- Storage arrays (e.g., strong data protection and high performance).

#### 5.1.1.2 Characteristics

| Characteristic | Value | Note(s) | Reference |
|---|---|---|---|
| Zone Type | Sequential Write Required | | See section 4.2.2. |
| Zone Active Resources Available | 12 or more recommended. Recommend that the number of active and open resources are equal. | Does not apply to ZBC/ZAC devices (i.e., SMR HDDs). | See section 4.2.4 and the NVM Express Zoned Namespace Command Set Specification. |
| Zone Open Resources Available | 12 or more recommended. | | See section 4.2.5 and the NVM Express Zoned Namespace Command Set. |
| Performance | Accessing 1 to 4 zones concurrently should achieve the maximum throughput of the associated media to the namespace and/or device. | | See section 4.2.6 |

| Characteristic | Value | Note(s) | Reference |
|---|---|---|---|
| Mandatory I/O Access Commands | Read and Write commands | | See section 4.2.7. |
| Mandatory Access Command (ZNS SSD only) | Zone Append | | See section 4.2.8. |

Figure 7 – Model A Characteristics

## 5.2 Model B

### 5.2.1 Overview

This model describes the requirements for a Zoned Storage device that provides high performance.

#### 5.2.1.1 Applicable Use cases

This Zoned Device Model minimizes the changes required to host software to support zoned block devices but requires high host I/O parallelism to achieve the full media bandwidth of a given device. The host software must:

- respect the sequential write requirement of the zone type, and must reset a zone to rewrite a zone;
- must access multiple zones in parallel to achieve the full bandwidth of the media; and
- must perform adequate parity protection to account for lower device UBER.

Common use cases include, but are not limited to:

- Archival storage (e.g., storage with host-defined erasure encoding).

#### 5.2.1.2 Characteristics

| Characteristic | Value | Note(s) | Reference |
|---|---|---|---|
| Zone Type | Sequential Write Required | | See section 4.2.2. |
| Zone Active Resources Available | Depends on device. Recommend that the number of active and open resources are equal. | Does not apply to ZBC/ZAC devices (i.e., SMR HDDs). | See section 4.2.4 and the NVM Express Zoned Namespace Command Set Specification. |
| Zone Open Resources Available | Depends on device | | See section 4.2.5 and the NVM Express Zoned Namespace Command Set Specification. |

| Characteristic | Value | Note(s) | Reference |
|---|---|---|---|
| Performance | Depends on device | Host must access at least the minimum number of zones concurrently, as defined by the device, to achieve the maximum throughput of the associated media to the namespace and/or device. | See section 4.2.6. |

Figure 8 – Model B Characteristics